# General Methods[a]

Version 4.1 of 28 November 2013

---

**Address of publisher:**

Institute for Quality and Efficiency in Health Care
Im Mediapark 8 (KölnTurm)
50670 Cologne
Germany

Tel.: +49 (0)221 – 35685-0
Fax: +49 (0)221 – 35685-1
E-mail: methoden@iqwig.de
Internet: www.iqwig.de

# Preamble

The Institute for Quality and Efficiency in Health Care (IQWiG[b]) is an establishment of the Foundation for Quality and Efficiency in Health Care. IQWiG is a professionally independent scientific institute. Information on the structure and organization of the Foundation and the Institute is available on the website www.iqwig.de.

The *General Methods* explain the legal and scientific basis of the Institute. Its tasks are described in this document, as are the scientific tools applied in the preparation of its products. The corresponding methods for the assessment of the cost-benefit relation of statutory health insurance (SHI) services are presented in the *General Methods for the Assessment of the Relation of Benefits to Costs*. Hence the Institute's methods papers provide an important contribution towards transparency in the Institute's mode of operation.

The *General Methods* are primarily directed at researchers. In order to make the information on the Institute's mode of operation accessible to as many interested persons as possible, the authors have aimed to produce a comprehensible document. However, as with any scientific text, a certain level of prior knowledge on the topic is assumed.

The *General Methods* aim to describe the Institute's procedures in a general manner. What specific individual steps the Institute undertakes in the assessment of specific medical interventions depend, among other things, on the research question posed and the available scientific evidence. The *General Methods* should therefore be regarded as a kind of framework. How the assessment process is designed in individual cases is presented in detail for each specific project.

The Institute's methods are usually reviewed annually with regard to any necessary revisions, unless errors in the document or relevant developments necessitate prior updating. Project-specific methods are defined on the basis of the methods version valid at that time. If changes are made to the general methodological procedures during the course of a project, then it will be assessed whether project-specific procedures need to be modified accordingly. In order to continuously further develop and improve its mode of operation, the Institute presents its *General Methods* for public discussion. This applies to the currently valid version, as well as to drafts of future versions.

---

[b]Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen

# What is new?

In comparison with Version 4.0 of the Institute's *General Methods* of 23 September 2011, in Version 4.1 minor errors were corrected and a few editorial changes made. The following changes to content were made:

- description of the external review for preliminary reports as an optional step in Sections 2.1.1 and 2.2.3

- division of the previous Section 3.1.4 into the new Sections 3.1.4 and 3.1.5 and specification of the requirements for the evidence base to formulate conclusions on benefit with different certainties of conclusions

- operationalization of the determination of the extent of added benefit, as well as the corresponding rationale, in Section 3.3.3 and in the new Appendix

- use of prediction intervals for meta-analyses with random effects in Section 7.3.8

# Table of contents

**List of tables**

**List of figures**

**List of abbreviations**

| Abbreviation | Definition |
| --- | --- |
| AGREE | Appraisal of Guidelines Research and Evaluation in Europe |
| AHP | analytic hierarchy process |
| AMSTAR | Assessment of Multiple Systematic Reviews |
| ANV | Arzneimittel-Nutzenbewertungsverordnung, AM-NutzenV (Regulation for Early Benefit Assessment of New Pharmaceuticals) |
| AWMF | Arbeitsgemeinschaft der Wissenschaftlichen Medizinischen Fachgesellschaften (Association of the Scientific Medical Professional Societies) |
| CA | conjoint analysis |
| CDSR | Cochrane Database of Systematic Reviews |
| CONSORT | Consolidated Standards of Reporting Trials |
| CPG | clinical practice guideline |
| DARE | Database of Reviews of Effects |
| DELB instrument | Deutsches Leitlinien-Bewertungs-(Instrument) (German Instrument for Methodological Guideline Appraisal) |
| DIPEx | Database of Personal Experiences of Health and Illness |
| DMP | disease management programme |
| EBM | evidence-based medicine |
| EMA | European Medicines Agency |
| EU | European Union |
| G-BA | Gemeinsamer Bundesausschuss (Federal Joint Committee) |
| GoR | grade(s) of recommendation |
| GRADE | Grading of Recommendations, Assessment, Development and Evaluation |
| HON | Health On the Net |
| HTA | health technology assessment |
| INAHTA | HTA Database of the International Network of Agencies for Health Technology Assessment |
| IPD | individual patient data |
| IQWiG | Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen (Institute for Quality and Efficiency in Health Care) |
| ITT | intention to treat |
| LoE | level(s) of evidence |
| MORE | McMaster Online Rating of Evidence |
| MTC | mixed treatment comparison |
| PRO | patient-reported outcome |

| Abbreviation | Definition |
|---|---|
| QALY | quality-adjusted life year |
| RCT | randomized controlled trial |
| ROC | receiver operating characteristic |
| SGB | Sozialgesetzbuch (Social Code Book) |
| SHI | statutory health insurance |
| STARD | Standards for Reporting of Diagnostic Accuracy |
| STE | surrogate threshold effect |
| WHO | World Health Organization |

*A chief cause of poverty in science is mostly imaginary wealth. The aim of science is not to open a door to infinite wisdom but to set a limit to infinite error.*

Bertolt Brecht. Life of Galileo. Frankfurt: Suhrkamp. World premiere, first version, Zurich theatre, 1943.

# 1 The Institute for Quality and Efficiency in Health Care

## 1.1 Legal responsibilities

The Institute was founded within the framework of the German Health Care Reform of 2004 [119] as an establishment of the Foundation for Quality and Efficiency in Health Care. The legal basis and responsibilities of the Institute have been anchored in Social Code Book V (SGB V[3]) [2] and adapted and extended several times in the course of further health care reforms. More information on the Institute's structure and organization is available on the website www.iqwig.de.

The Institute addresses issues of fundamental relevance for the quality and efficiency of SHI services. Its specific responsibilities are outlined in detail in § 139a SGB V:

- search for, assessment and presentation of current scientific evidence on diagnostic and therapeutic procedures for selected diseases

- preparation of scientific reports, expert opinions, and comments on quality and efficiency issues of SHI services, taking age, gender, and personal circumstances into account

- appraisal of evidence-based clinical practice guidelines (CPGs) on the most relevant diseases from an epidemiological point of view

- issue of recommendations on disease management programmes (DMPs)

- assessment of the benefit and cost of drugs

- provision of easily understandable information for all patients and consumers on the quality and efficiency of health care services, as well as on the diagnosis and treatment of diseases of substantial epidemiological relevance

The modalities of the commissioning and performance of tasks are specified in §139b SGB V. According to this law, only the Federal Joint Committee (G-BA[4]) or the Federal Ministry of

---

[3] Sozialgesetzbuch: regulates the statutory health care services.

[4] Gemeinsamer Bundesausschuss: The G-BA is the decision-making body of the self-government of the German health care system. More information on the Committee's responsibilities is provided at http://www.english.g-ba.de/.

Health (BMG[5]) may commission the Institute. In the case of commissioning by the Ministry, the Institute can reject a commission as unfounded, unless the Ministry funds the project.

The Institute must ensure that external experts are involved in the work on commissions. In order to ensure the Institute's scientific independence, these experts are required to disclose all connections to associations and contract organizations, particularly in the pharmaceutical and medical devices industries, including details on the type and amount of any remuneration received (see Section 2.2.2).

The Institute submits the results of the work on commissions awarded by the G-BA to this body in the form of recommendations. According to the law, the G-BA must consider these recommendations in its decision-making processes.

The Institute is largely funded by contributions of SHI members. For this purpose, a levy is determined by the G-BA in accordance with §139c SGB V. This levy is paid by all German medical practices and hospitals treating SHI-insured patients.

Within the framework of the Act on the Reform of the Market for Medicinal Products (AMNOG[6]), at the beginning of 2011, the Institute's responsibilities were extended to the assessment of the benefit of drugs with new active ingredients shortly after market entry [120]. For this purpose manufacturers must submit dossiers summarizing the results of studies. The G-BA is responsible for this "early benefit assessment"; however, it may commission the Institute or third parties to examine and assess the dossiers.

The new regulations in §35a SGB V are the basis for these assessments. They are supplemented by a legal decree of the Ministry of Health [70], which has also been effective since the beginning of 2011, and the G-BA's Code of Procedure [198].

In connection with a benefit assessment the G-BA can also commission the Institute to conduct a health economic evaluation. In this context, the benefit of a medical technology is related to the cost, with the aim of determining an appropriate price. The framework of these health economic evaluations is specified in §35b SGB V and §139a SGB V. For example, in its commission the G-BA must designate which comparator treatments should be considered and which patient groups the assessment is to be conducted for.

The health economic evaluation itself is based on a comparison with other drug or non-drug interventions. In particular, the following criteria to determine the benefit for patients are named in the law: increase in life expectancy, improvement in health status and quality of life (QoL), and reduction in disease duration and adverse effects. The definition of a "patient-

---

[5] Bundesministerium für Gesundheit
[6] Arzneimittelmarktneuordnungsgesetz

relevant benefit" valid for the Institute is inferred from the above specifications in the law (see Section 3.1).

Depending on the commission, the Institute determines the methods and criteria for the preparation of assessments on the basis of the international standards of evidence-based medicine (EBM) and health economics recognized by the relevant experts. The term "evidence-based medicine", its development and the underlying concept are described in detail in Section 1.2.

During the preparation of its reports, the Institute ensures the high transparency of procedures and appropriate involvement of third parties. In all important phases of report preparation the law obliges the Institute to provide the opportunity of comment to experts, manufacturers and relevant organizations representing the interests of patients and self-help groups of chronically ill and disabled persons, as well as to the Federal Government Commissioner for Patients' Affairs. The Institute goes beyond this obligation by allowing all interested persons and institutions the opportunity to submit comments on its reports, and considers these comments in its assessments.

The implementation of these regulations is described in Section 2.1.1 in connection with the production of report plans (protocols) and preliminary reports.

In addition, the Institute publishes the results of its work and supplementary information on its publicly accessible website. Those interested can also subscribe to the Institute's e-mail service (info service), where subscribers themselves can specify what type of information they would like to receive from the Institute.

## 1.2   Evidence-based medicine

EBM refers to patient health care that is not only based on opinions and consensus, but considers "evidence" – i.e. proof (e.g. of the benefit of a medical intervention) determined with the most objective scientific methods possible. EBM comprises tools and strategies designed to safeguard against false decisions and false expectations. In this context, a false decision can mean that beneficial interventions are not implemented in health care (or implemented with delay), or that useless or even harmful interventions are widely applied [17,164,218,222].

However, tools designed to prevent subjective (and therefore often biased) assessments (see also Chapter 7) were not first invented with the introduction of the term "EBM", but originated decades ago. In Germany, as early as 1932 Paul Martini described the main elements of a fair assessment of drug effectiveness in his monograph *Methodology of Therapeutic Studies* [358]. In the early 1960s, the method of randomly allocating study participants to comparator groups (randomization) in order to assess the effectiveness and safety of medical interventions became the internationally accepted standard [251]. Starting in the United States, in this period this type of study became the precondition for the approval of

drugs and (in some cases) medical devices regulated by authorities, legislation and other regulations [31]. About 20 years later, clinical epidemiologists attempted to establish this methodology in clinical practice [170]. Accompanied at times by serious controversy, this was not actually achieved until the 1990s, at the same time as the concept was defined as "EBM". Since this time, clinical studies and the systematic search for and assessment of these studies (systematic reviews) have formed the basis of the international scientific standard for health technology assessments (HTAs) [30].

The Institute is legally obliged to conduct an "assessment of the medical benefit [of interventions] according to internationally recognized standards of EBM" (§139a [4] SGB V). It is the task of the Institute's methods paper to describe the methods and strategies that define these international standards. EBM is not a rigid concept: which standard tool is to be applied, and when, depends on the question to be answered and the decision to be made. Despite the application of standards, decisions for which no international specifications are (as yet) available have to be made repeatedly in the search for, and the processing and assessment of studies. EBM also includes the freedom to define one's own specifications in such situations. However, this freedom is linked to the obligation to define such specifications preferably a priori, and to explain assessments in a transparent manner, so that the rationale is comprehensible. This chapter explains that in the implementation of EBM and the definition of specifications, an institution such as IQWiG is in a different situation from clinicians seeking support for a treatment decision.

### 1.2.1 Practical evidence-based medicine

The EBM concept is a strategy for physicians who, from a range of possible interventions, seek the most promising alternatives suited best to the needs of their patients, and who aim to offer prospects of success in an objective manner. This implementation of EBM in daily clinical practice for "individual patients" was defined by David Sackett et al. [445] as follows: "EBM is the conscientious, explicit and judicious use of current best evidence in making decisions about the care of the individual patient. It means integrating individual clinical expertise with the best available external clinical evidence from systematic research" (1996).

However, the "best available evidence" is often incomplete or unreliable. EBM has developed instruments to assess uncertainty; evidence levels are often used for illustration. In this way, EBM helps physicians and patients to recognize the type and degree of uncertainty; they can then discuss how to deal with this. Especially in uncertain situations, personal preferences are important and determine what option patients choose. Apart from being based on evidence, decisions are also ideally based on the clinical condition and circumstances of the individual patient, as well as on his or her preferences and actions [239]. At the same time, the description of the identified gaps in knowledge creates the precondition for medical research targeted towards patients' needs.

EBM is based on a critical approach [306]. The importance of scepticism is underlined by the fact that over the past few decades, several insufficiently tested but widely applied therapies have been assessed with EBM methods; these assessments have shown that a hasty, overoptimistic approach to a new intervention can have dangerous consequences for patients [140,429]. It is the Institute's task to assess objectively with what certainty the benefit of medical interventions has been demonstrated, in order to counter inappropriate judgements.

### 1.2.2   The relevance of evidence-based medicine for the Institute

The Institute's main task is to provide the most reliable answer possible to the question specified by the contracting agency as to whether evidence is available of the benefits or harms from an intervention. The aim is to present sufficiently reliable proof that "Treatment A" is better for patients than "Alternative B" for a specific disease. In short: What is the benefit of A compared with B?

The Institute's remit is therefore intentionally not aimed towards treating individual patients with their potential specific characteristics, but towards determining for which patient groups proof of a benefit of an intervention is available. In its decisions, the G-BA then considers aspects of patient care that are beyond the scope of a benefit assessment [198].

### 1.2.3   Evidence-based decision-making in health care

The reports prepared by the Institute are to serve the G-BA as a basis for decisions that in principle apply to all SHI members. Other Institute products are, for example, to serve as information for the G-BA. The type of decisions made by institutions like the G-BA has an effect on the application of EBM methods.

### 1.2.4   Strategies of evidence-based medicine

A characteristic standard element of EBM is the structured and systematic approach to the search for a response to a medical question:

1) The medical question must be worded precisely. Medicine (nearly) always deals with the choice between at least 2 alternatives. This can refer to treatments, diagnostic tests or complex changes in life style. From this, the following question is always inferred: Is Option A better than Option B? In this context, the decision not to undergo treatment can also be an option that should be thoroughly reviewed. However, it should be stressed that such an option (e.g. "watchful waiting") is not the same as "doing nothing".

2) It must be defined how the benefit of treatment (or diagnosis or lifestyle change) should be measured. A standard element of EBM is the question about relevant consequences for patients: Can life expectancy be increased? Can symptoms and quality of life be improved?

3) In EBM it is explicitly noted that in medicine, only probability statements or only conclusions about groups of patients are usually possible with regard to the benefit of

treatment, diagnostic procedures, or lifestyle changes. Benefit is demonstrated by showing that an intervention increases the probability of a beneficial outcome and/or reduces the risk of a non-beneficial outcome. In order to prove the benefit of an intervention, studies in sufficiently large groups of suitable patients are required. International researchers have developed a range of rules and tools for the planning, conduct, and analysis of such studies. The most important aim is to minimize (or, if this is impossible, at least document) factors that can distort the results of a comparison. The effects of such confounding factors are referred to as "bias". The rules and tools that are internationally accepted as the prevailing standard, and are under continuous development, are the methodological basis of EBM and the Institute's work.

4) A further key EBM strategy is to identify all "appropriate" studies (i.e. whose design and conduct are of appropriate quality) on a question and, in this way, to summarize the reliable evidence available. In this context, if large differences are shown between the results of individual studies (heterogeneity), an attempt should be made to explain them. The findings of these summaries and assessments are referred to as systematic reviews; the statistical analyses are referred to as meta-analyses.

## 1.2.5  The relevance of certainty of results

A specific characteristic of EBM is that it allows assessment as to what extent the available evidence is reliable. Decisions made by the G-BA must be based on highly reliable scientific evidence, as they have far-reaching consequences for all SHI members (e.g. exclusion of services from reimbursement).

The assessment of the certainty of results therefore plays a key role in the Institute's reports. Numerous details on how studies are planned, conducted, analysed, and published have an impact on how reliable the available results are. It is an international EBM standard to test and assess these aspects critically. However, how the certainty of results needed to answer a question can be achieved also depends on the disease and on the effect size of an intervention: If 2 athletes pass the finishing line of a fair race with a great distance between them, no stopwatch is needed to identify the winner. For example, the benefit of a new therapy that results in the cure of a previously always fatal disease can be proven by a relatively small number of surviving patients. In this case, the judgement is also ultimately based on a comparison, but in interventions with such dramatic effects, the comparison between historical and current patients may already provide sufficient certainty. However, therapies that show such dramatic benefits are very rare in modern medicine.

In chronically ill patients in particular, differences between 2 therapy alternatives are mostly smaller and may be easily confounded by a fluctuant course of disease. In these cases, precise methods and appropriate study designs are required in order to be able to recognize therapy effects under such fluctuations.

It can be assumed that the Institute will be specifically commissioned to compare such interventions where it is not immediately recognizable which alternative will be more beneficial. However, the smaller the expected differences between 2 alternatives are, the more reliable the studies must be in order to be sufficiently certain that an observed effect is not caused by chance or measurement errors (a world record over 100 metres can no longer be measured with an hourglass). In the event of small differences, their clinical relevance must also be judged.

The following requirements for precision and reliability determine the Institute's mode of operation:

1) For every question investigated, it is an international EBM standard to specify the study type (measuring tool) that minimizes the risk of unjustifiably discriminating against one of the alternatives.

2) The Institute's assessments on the benefits and harms of interventions are therefore normally based only on studies with sufficient certainty of results. This ensures that the decisions made by the G-BA, which are based on the Institute's recommendations, are supported by a sound scientific foundation. Moreover, an assessment that includes a literature search for studies with insufficient certainty of results would be costly and time consuming.

3) If it emerges that studies of the required quality and precision are generally lacking, it is the core task of the Institute to describe the circumstances and conclude that on the basis of the "currently best available" evidence, it is not possible to make reliable recommendations.

4) It is the G-BA's responsibility to take this uncertainty into account in its decision-making processes. In addition to considering scientific evidence, the G-BA also considers other aspects in its decisions, such as the efficiency of interventions as well as the needs and values of people [209]. In an uncertain scientific situation, such aspects become more important. In addition, the G-BA also has the option to call for studies in order to close the evidence gaps identified.

### 1.2.6  The connection between certainty of results and proximity to everyday conditions

The great value placed on the assessment of the certainty of results is often criticized. One argument is that studies with a high certainty of results (especially randomized controlled trials, RCTs) may have high internal validity, but often do not represent patient care under everyday conditions, and are therefore not transferable, i.e. have only low external validity. In this context it must be examined how well the patient population investigated in the studies, the interventions applied, and the outcome criteria analysed are in accordance with everyday conditions in health care. This criticism is then often connected to the call to include other study types without randomization, in order to better consider everyday conditions.

However, this criticism conflates levels of arguments that should be clearly separated. The following aspects should be taken into account:

1) The basis of a benefit assessment is the demonstration of causality. An indispensable precondition for such a demonstration is a comparative experiment, which has to be designed in such a way that a difference between intervention groups – an effect – can be ascribed to a single determining factor – the intervention tested. This goal requires considerable efforts in clinical trials, as there are numerous confounding factors that feign or mask effects (bias). The strongest of these distorting influences are unequal baseline conditions between comparator groups. Randomization (together with careful concealment) is currently the best available tool to minimize this type of bias. Random allocation of participants to groups ensures that there are no systematic differences between groups, neither regarding known factors (e.g. age, gender, disease severity), nor unknown factors. For this reason, RCTs provide a basic precondition for the demonstration of causality. However, randomization alone does not guarantee high certainty of results. To achieve this, the unbiased assessment, summarization and publication of results, for example, are also required.

2) Study types other than RCTs are usually not suited to demonstrate causality. In non-randomized comparative studies, as a matter of principle structural equality of groups cannot be assumed. They therefore always provide a potentially biased result and mostly cannot answer with sufficient certainty the relevant question as to whether a difference observed is caused by the intervention tested. The use of non-randomized studies as proof of the causality of an intervention therefore requires particular justification or specific preconditions and special demands on quality.

3) It is correct that many randomized studies do not reflect aspects of everyday patient care, for example, by excluding patients with accompanying diseases that are common in everyday life. However, this is not a consequence of the randomization technique, but of other factors (e.g. definition of narrow inclusion and exclusion criteria for the study, choice of interventions or outcome criteria). In addition, patients in randomized studies are often cared for differently (more intensively and more closely) than in everyday practice. However, these are intentional decisions made by those persons who wish to answer a specific question in a study. Dispensing with randomization does not change these decisions. There is also a selection of participants in non-randomized studies through inclusion and exclusion criteria and other potential design characteristics, so that external validity is not given per se in this study type any more than in RCTs.

4) Even if patient groups in an RCT differ from everyday health care, this does not mean the external validity of study results must be questioned. The decisive issue is in fact whether it is to be expected that a therapy effect determined in a population varies in a different population.

5) It depends on the individual case how the intensity of care provided in a study influences outcomes. For example, it is conceivable that a benefit of an intervention actually exists

---

only if patients are cared for by specially qualified physicians, as under everyday conditions too many complications may otherwise occur. However, it is also possible that intensified care of patients is more likely to reduce differences between groups. For example, differences in treatment adherence may be smaller in studies where, as a matter of principle, patients are cared for intensively.

6) However, the initiators of a clinical trial are responsible for the specification of study conditions. They can define research questions and outcomes rated as so relevant that they should be investigated in a study. If, for example, a drug manufacturer regards treatment adherence to be an important aspect of the benefit of a product, the obvious consequence would be to initiate studies that can measure this aspect with the greatest possible certainty of results and proximity to everyday conditions, and at the same time demonstrate its relevance for patients.

The above remarks show that certainty of results and proximity to everyday conditions (or internal and external validity) have no fixed relationship. High certainty of results and proximity to everyday conditions do not exclude one another, but only require the appropriate combination of study type, design and conduct.

Even if criticism of the lack of proximity to everyday practice may actually be justified for many studies, nothing would be gained by dispensing with high certainty of results in favour of greater proximity to everyday practice, because one would thereby be attempting to compensate one deficit by accepting another, more serious, one [237].

Studies that combine proximity to everyday conditions and high certainty of results are both desirable and feasible. RCTs are indeed feasible that neither place demands on patients beyond everyday health care nor specify fixed study visits. Such studies are being discussed at an international level as "real world trials", "practical trials" or "pragmatic trials" [185,187,206,357,518]. However, such pragmatic trials may themselves also lead to interpretation problems. For example, if very broad inclusion criteria are chosen, the question arises as to whether the (overall) study results can be applied to the overall study population [550], which, at least to some extent, would ultimately have to be answered by means of appropriate subgroup analyses.

### 1.2.7 Benefit in individual cases

The aim of a benefit assessment is to make robust predictions for future patients using results of studies suited to demonstrate causal effects. The conclusions drawn always apply to groups of patients with certain characteristics. Conclusions on the benefit of an intervention in terms of predictions of success for individual cases are, as a matter of principle, not possible. Vice versa, experiences based on individual cases (except for specific situations, e.g. dramatic effects) are unsuitable for a benefit assessment, as it is not possible to ascribe the results of an individual case (i.e. without a comparison) to the effect of an intervention.

For certain research questions (therapy optimization in individual patients) so-called (randomized) single patient trials (or "n-of-1" trials) can be conducted [219,224,290,461]. However, these are usually not suited to assess the benefit of a treatment method for future patients.

## 2   The Institute's products

According to its legal remit, the Institute prepares a variety of products in the form of scientific reports and easily understandable health information for consumers and patients. This chapter describes procedures and general methods applied in the preparation of the Institute's products. At first the individual products are named and product-specific procedures presented (Section 2.1). The next section outlines further aspects independent of products (Section 2.2).

### 2.1   Product-specific procedures

The Institute's products include

- report

- rapid report

- dossier assessment

- addendum

- health information

- working paper

The preparation of reports and rapid reports is conducted on the basis of the award of individual commissions through the G-BA or Ministry of Health. The basis of this are the Institute's responsibilities described in §139a SGB V (see also Section 1.1). Accordingly, reports and rapid reports can be prepared on the benefit assessment of drug and non-drug interventions, on health economic evaluations, and on the appraisal of CPGs. The main difference between reports and rapid reports is that commenting procedures (hearings) are only conducted for reports, but not for rapid reports. Accordingly, rapid reports are particularly intended for recommendations at short notice, for which, from the point of view of the contracting agency, no hearings by the Institute are required.

Dossier assessments are commissioned by the G-BA. The foundation for this is §35a SGB V, which regulates the assessment of the benefit of new active ingredients on the basis of a dossier by the pharmaceutical company (see also Section 3.3.3). No hearing by the Institute is intended for dossier assessments according to §35a SGB V; this is conducted in the further procedure by the G-BA.

Addenda can be commissioned by the G-BA or Ministry of Health in cases where, after the completion of a report, rapid report or dossier assessment, the need for additional work on the commission arises during the course of consultations.

Health information can be prepared on the basis of an individual commission; it can also be the consequence of a commission in other areas of the Institute's work (easily understandable

version of other products of the Institute, e.g. a report) or be prepared within the framework of the general legal remit to provide health information.

Working papers are prepared under the Institute's own responsibility; specific commissioning by the G-BA or Ministry of Health is not required. This takes place either on the basis of the general commission (see Section 2.1.6), with the aim of providing information on relevant developments in health care, or within the framework of the legal remit to develop the Institute's methods. The Institute's *General Methods* are not to be understood as a working paper in this sense, and are subjected to a separate preparation and updating procedure, which is outlined in the preamble of this document.

An overview of the Institute's various products is shown in Table 1 below. Product-specific procedures are described in the subsequent Sections 2.1.1 to 2.1.6.

Table 1: Overview of the Institute's products

| Product | Objective | Procedure | Commissioned by |
|---|---|---|---|
| Report | Recommendations on tasks described in §139a SGB V, including hearing | Described in Section 2.1.1 | G-BA, Ministry of Health |
| Rapid report | Recommendations on tasks described in §139a SGB V, insofar as no hearing on interim products is required; in particular provision of information at short notice on current topics | Described in Section 2.1.2 | G-BA, Ministry of Health |
| Dossier assessment | Assessment of the benefit of drugs with new ingredients according to §35a SGB V | Described in Section 2.1.3 | G-BA |
| Addendum | Supplementary information provided at short notice by the Institute on issues that have arisen during the consultation on its completed products | Described in Section 2.1.4 | G-BA, Ministry of Health |
| Health information | Easily understandable information for consumers and patients; wide scope of topics | Described in Section 2.1.5 | G-BA, Ministry of Health/own initiative of the Institute |
| Working paper | Information on relevant developments in health care or methodological aspects | Described in Section 2.1.6 | Own initiative of the Institute |
| G-BA: Gemeinsamer Bundesausschuss (Federal Joint Committee); SGB: Sozialgesetzbuch (Social Code Book) | | | |

### 2.1.1 Report

**A) Procedure for report production**

The procedure for report production is presented in Figure 1. All working steps are performed under the Institute's responsibility and regularly involve external experts (see Section 2.2.1). If necessary, the Institute's Scientific Advisory Board is also involved. The internal quality assurance process is not outlined in this flow chart.

After **commissioning** by the G-BA or Ministry of Health, the research question is formulated. Depending on the topic, the determination of outcome criteria is also required (e.g. in benefit assessments). As a rule, relevant patient organizations are involved, especially in the definition of patient-relevant outcomes; the opinion of individual affected patients may also be heard. Subsequently, the report plan (protocol) is prepared.

The **report plan** contains the precise scientific research question, including the outcome criteria (e.g. patient-relevant outcomes), as well as the inclusion and exclusion criteria of the information to be used in the assessment. This plan also includes a description of the project-specific methodology applied in the retrieval and assessment of information. The preliminary version of the report plan is first forwarded to the contracting agency as well as to the Foundation's Board of Directors, the Foundation Council and the Board of Trustees. It is normally published on the Institute's website 5 working days later.

For a period of at least 4 weeks, the public is given the **opportunity to submit written comments (hearing)**. This opportunity particularly refers to the project-specific methodological approach applied to answer the research question. The research question itself is usually specified by the commission, and is not an object of the commenting procedure. Optionally, an oral scientific debate including persons submitting comments may be held. This debate serves the potentially necessary clarification of aspects of the written comments and aims at improving the scientific quality of the report plan.

| Commissioning |
| By Federal Joint Committee (G-BA) or Ministry of Health |

Formulation of research question(s)

Depending on topic: determination of outcome criteria (e.g. patient-relevant outcomes, with involvement of individual patients/patient representatives)

**Report plan (preliminary version)**

**Hearing\***

**Report plan**

**Potential amendment to report plan**

Information retrieval and scientific evaluation

**Preliminary report**

**Hearing\***

External review (optional)

Compilation and appraisal of comments and external review: Update of information

**Final report**

\* The hearing is conducted by inviting written comments. In addition, an optional oral scientific debate may be held to discuss any unclear aspects of the written comments.

Figure 1: Procedure for the production of a report

After the analysis of the comments and (if appropriate) the conduct of the oral debate, the revised report plan, together with the documentation of the hearing on the report plan, are first forwarded to the contracting agency, the Foundation's Board of Directors, the Foundation Council and Board of Trustees. This document is usually published on the Institute's website 5 working days later. The revised report plan is the basis for the preparation of the preliminary report. If further relevant methodological changes are required in the course of the preparation of the preliminary report, these are usually presented in one or more amendments to the report plan. An opportunity to submit comments is usually also provided after publication of an amendment, following the conditions outlined above.

The results of the information retrieval and the scientific assessment are presented in the **preliminary report**. In order to avoid undue delay in the Institute's work, the retrieval and assessment of information already start before completion of the hearing on the report plan on the basis of the criteria formulated in the preliminary report plan. However, the result of the hearing is explicitly not anticipated, as these criteria may be modified on grounds of the hearing on the preliminary version of the report plan. This may also lead to supplementation and/or modification of the retrieval and assessment of information.

The preliminary report includes the preliminary recommendation to the G-BA. After completion it is first forwarded to the contracting agency as well as to the Foundation's Board of Directors, the Foundation Council and the Board of Trustees. The preliminary report is usually published on the Institute's website 5 working days after it is sent to the contracting agency.

For a period of at least 4 weeks, the public is then given the **opportunity to submit written comments (hearing)**. The results of the retrieval and assessment of information presented in the preliminary report are in particular the subject of the commenting procedure. Optionally, an oral scientific debate with those submitting comments may be held. This debate serves the potentially necessary clarification of aspects of the written comments and aims at improving the scientific quality of the final report.

The **final report**, which is based upon the preliminary report and contains the assessment of the scientific findings (considering the results of the hearing on the preliminary report), represents the concluding product of the work on the commission. The final report and the documentation of the hearing on the preliminary report are first forwarded to the contracting agency, as well as to the Foundation's Board of Directors and Foundation Council, and subsequently (usually 4 weeks later) forwarded to the Foundation's Board of Trustees. These documents are then published on the Institute's website (usually a further 4 weeks later). If comments are received on final reports that contain substantial evidence not considered, or if the Institute receives information on such evidence from other sources, the contracting agency will be sent well-founded information on whether, in the Institute's opinion, a new commission on the topic is necessary (if appropriate, a report update) or not. The contracting agency then decides on the commissioning of the Institute. Such an update is conducted

according to the general methodological and procedural requirements for the Institute's products.

## B) General remarks on the commenting procedure (hearing)

### Organizations entitled to submit comments

In accordance with §139a (5) SGB V, the Institute must ensure that the following parties are given the opportunity to submit comments in all important phases of the assessment procedure: medical, pharmaceutical, and health economic experts (from research and practice), drug manufacturers, relevant organizations representing the interests of patients and self-help groups for the chronically ill and disabled, as well as the Federal Government Commissioner for Patients' Affairs. Their comments must be considered in the assessment. These requirements are taken into account by the fact that hearings on the report plan and preliminary report are conducted and that the circle of people entitled to submit comments is not restricted. Moreover, all the Institute's products, in accordance with §139a SGB V, are sent to the Board of Trustees before publication. Patient organizations, the Federal Government Commissioner for Patients' Affairs, organizations of service providers and social partners, as well as the G-BA's self-government bodies are represented in the Board of Trustees.

### Formal requirements

In order to avoid undue delay in the Institute's work, the comments must fulfil certain formal requirements. Further information on the commenting procedure, including the conditions for participation in a scientific debate, can be found in a guideline published on the Institute's website.

### Publication of comments

Comments that fulfil the formal requirements are published in a separate document on the Institute's website (*Documentation and appraisal of the hearing*). In order to ensure transparency, documents that are submitted together with the comments and are not publicly accessible (e.g. manuscripts) are also published.

### Submission of documents within the framework of the hearing

During both the hearing on the report plan and the one on the preliminary report, the opportunity is provided to submit any document of appropriate quality, which, according to the person submitting comments, is suited to answer the research question of the report. If the search strategy defined in the report plan is restricted to RCTs, for example, non-randomized studies may nevertheless be submitted within the framework of the commenting procedure. However, in such cases, appropriate justification of the validity of the causal interpretation of the effects described in these studies is also required.

### 2.1.2 Rapid report

The procedure for the production of a **rapid report** is presented in Figure 2. All working steps are performed under the responsibility of the Institute, involving external experts where appropriate (see Section 2.2.1). If necessary, the Institute's Scientific Advisory Board is also involved. The internal quality assurance process is not presented in this flow chart.



Figure 2: Procedure for the production of a rapid report

Rapid reports are primarily produced with the aim of providing information at short notice on relevant developments in health care (e.g. new technologies, publication of milestone studies). A shorter production period is usually required here. Interim products are therefore not published and are not the subject of a hearing.

After **commissioning** by the G-BA or Ministry of Health, the research question is formulated. Depending on the topic, the determination of outcome criteria is also required (e.g. in benefit assessments). In this context, patient organizations may be involved or the opinion of individual affected patients sought, especially for the definition of patient-relevant outcomes. Subsequently, the project outline is prepared.

The **project outline** summarizes the main steps of the information retrieval and scientific assessment. It forms the basis for the production of the rapid report. The project outline is not published.

The **rapid report** presents the results of the information retrieval and scientific assessment. Before completion, as a further quality assurance step, optionally a draft of the rapid report may be reviewed by one or more external reviewers (see Section 2.2.3) with proven methodological and/or topic-related competence. After completion the rapid report is then sent to the contracting agency, the Foundation's Board of Directors and Foundation Council, as well as (usually a week later) to the Board of Trustees. The rapid report is usually published on the Institute's website 4 weeks after it is sent to the contracting agency and Board of Directors. If comments on rapid reports are received that contain substantial evidence not considered, or if the Institute receives such evidence from other sources, the contracting agency will be provided with well-founded information on whether, in the Institute's opinion, a new commission on the topic is necessary (if appropriate, a rapid report update) or not. The contracting agency then decides on the commissioning of the Institute. Such an update is conducted according to the general methodological and procedural requirements for the Institute's products.

### 2.1.3  Dossier assessment

The procedure for the production of a dossier assessment is presented in Figure 3. All working steps are performed under the Institute's responsibility and regularly involve external expertise (see Section 2.2.1). If necessary, the Institute's Scientific Advisory Board is also involved. The internal quality assurance process is not outlined in this flow chart.

Figure 3: Procedure for the production of a dossier assessment

After the **forwarding of the dossier** by the G-BA, the assessment of the dossier content is conducted under the responsibility of the Institute. In this context, medical expertise and the patient perspective are as a rule involved via external experts and patients/patient organizations respectively.

**Medical expertise is primarily involved** on the basis of a questionnaire sent to external experts at the beginning of the assessment. In its assessment the Institute considers the external experts' feedback. In addition, external experts may if necessary be drawn upon to clarify specific questions arising during the course of the assessment. External experts are identified via the Institute's own database for external experts (see Section 2.2.1).

The **patient perspective is considered** on the basis of a questionnaire sent to patients/patient organizations at the beginning of the assessment. In its assessment the Institute considers the information provided in this questionnaire, e.g. on relevant outcomes and important subgroups. Patients/patient organizations are identified via the relevant organizations named in §140f SGB V.

The basis of the assessment is the dossier submitted to the G-BA by the pharmaceutical company and then forwarded to the Institute. The Institute may optionally perform its **own literature search** to support the assessment.

The preparation of the **dossier assessment** is the final step in the process. In accordance with §35a SGB V, the assessment must be completed no later than 3 months after the relevant date for the submission of the dossier. After its completion the dossier assessment is delivered to the G-BA. Shortly afterwards it is subsequently forwarded to the Foundation's Board of Directors, the Foundation Council and the Foundation's Board of Trustees and then published on the Institute's website.

## 2.1.4  Addendum

The procedure for the production of an addendum is presented in Figure 4. All working steps are performed under the responsibility of the Institute, involving the Institute's Scientific Advisory Board where appropriate. The internal quality assurance process is not outlined in this flow chart.

Figure 4: Procedure for the production of an addendum

An addendum can be commissioned if the need for additional work on the commission arises during the consultations on products completed by the Institute. Depending on the type and extent of the research question it may be meaningful to prepare a **project outline** in which the main steps of the information retrieval and scientific assessment are summarized. The project outline is not published.

In the work on the **addendum**, depending on the type and extent of the research question, it may be meaningful to involve those external experts who were involved in preparing the underlying product of the Institute. After completion, the addendum is first sent to the G-BA, as well as to the Foundation Council and the Board of Directors. It is usually forwarded to the Foundation's Board of Trustees 1 week later and published on the Institute's website a further 3 weeks later.

### 2.1.5   Health information

The Institute produces **health information** for patients and the general public in various formats, which are presented in more detail in Chapter 5. The Institute's information products primarily include

- **feature articles:** comprehensive article which forms the basis of a set of related products on an important health issue

- **fact sheets:** short, easily-understandable information

- **research summaries:** summaries of systematic reviews, HTAs or larger studies, including summaries of the other products of the Institute, insofar as they are relevant for patients

In addition to these information products, supplementary items are also produced. These aim to make the key messages of the health information more understandable and interesting.

The production process for health information is presented in Figure 5. When information products are produced, external expertise is involved at various stages. The tasks fulfilled by these external people, depending on the type of product, are described in more detail in Chapter 5. The internal review process is also not outlined in the flow chart below. The process for the selection of health information topics on the Institute's own initiative is described in Section 5.3.1. The Institute's health information for patients and the general public is produced

- in response to commissions received from the G-BA or Ministry of Health

- to summarize other products published by the Institute and as accompanying information for these products

- to fulfil its legislative responsibility to provide consumers with health information, as well as on its own initiative within the framework of the G-BA's general commission

The Institute's general commission (see Section 2.1.6) was amended in July 2006 and in March 2008 as regards the production of health information to specifically include informing the general public. The process of evidence scanning, which the Institute undertakes to develop potential topics for this information, is described in Section 5.3.1. Chapter 5 also describes the methodology of literature searches, as well as patient involvement.

```
┌─────────────────────────────────┐      ┌─────────────────────────────────┐
│ Topic chosen on the Institute's │      │ IQWiG report, rapid report,     │
│ initiative or commission        │◄─────│ working paper or evidence       │
│ awarded by the Federal Joint    │      │ evaluated within the            │
│ Committee or Federal Ministry   │      │ framework of evidence scanning  │
│ of Health                       │      │                                 │
└─────────────────────────────────┘      └─────────────────────────────────┘
        │                  │                           │
   ┌────────────┐   ┌────────────┐           ┌──────────────────┐
   │  Feature   │   │ Fact sheet │           │ Research summary │
   │  article   │   │            │           │                  │
   └────────────┘   └────────────┘           └──────────────────┘
        │                  │                           │
┌──────────────┐  ┌──────────────────┐                │
│ Formation of │  │ Formation of the │                │
│ the project  │  │ project group    │                │
│ group        │  │ (for commissions)│                │
└──────────────┘  └──────────────────┘                │
        │                  │                           │
   ┌──────────────────────────────┐                   │
   │ Information retrieval and     │                   │
   │ scoping                       │                   │
   └──────────────────────────────┘                   │
                    │                                  │
   ┌──────────────────────────────┐      ┌──────────────────┐
   │ Scientific evaluation         │─────▶│ Text preparation │
   └──────────────────────────────┘      └──────────────────┘
                                                   │
                                      ┌──────────────────────────┐
                                      │ External review          │
                                      │ (excluding research      │
                                      │ summaries on IQWiG       │
                                      │ products)                │
                                      └──────────────────────────┘
                                                   │
                                      ┌──────────────────────────┐
                                      │ Product draft            │
                                      └──────────────────────────┘
                                                   │
                                      ┌──────────────────────────┐
                                      │ Presentation to the      │
                                      │ contracting agency /     │
                                      │ Board of Trustees /      │
                                      │ Scientific Advisory Board│
                                      └──────────────────────────┘
                                                   │
                                      ┌──────────────────────────┐
                                      │ Compilation and          │
                                      │ evaluation of comments   │
                                      │ Preparation of           │
                                      │ supplementary items      │
                                      └──────────────────────────┘
                                           │                │
                              ┌──────────────────┐  ┌──────────────────────┐
                              │ Final report     │─▶│ Publication on       │
                              │ (for commissions)│  │ informedhealthonline.│
                              │                  │  │ org                  │
                              └──────────────────┘  └──────────────────────┘
```

Figure 5: Procedure for the production of health information

After **commissioning** by the G-BA or Ministry of Health, an internal project group is formed. A project group includes at least one staff member of the Institute who does not belong to the Health Information Department. A project group is also formed for each feature article.

After the **text** has been prepared and a departmental quality assurance process performed, the drafts are sent out for **external review**. The drafts of feature articles and fact sheets are sent to at least one external reviewer. Research summaries are sent to the author of the research that was summarized. The research summaries of the Institute's products (accompanying information) are usually only reviewed internally, but can also be sent to external reviewers.

The supplementary items are subjected to the same internal review processes as the corresponding feature articles, fact sheets and research summaries. If necessary, they may also be reviewed externally. Patient stories (interviews, see Section 5.4.4) are only published if written consent of the patient involved has been given.

The final draft of a health information product is sent to the contracting agency and to the bodies of the Institute, including the Board of Trustees, for limited **submission of comments** within a 1-month consultation period. The Board of Trustees includes organizations of service providers and of employers and employees, and representatives of relevant organizations responsible for representing the interests of patients and self-help groups of chronically ill and disabled persons, as well as the Federal Government Commissioner for Patients' Affairs. Unlike preliminary reports, patient information drafts are not published on the Institute's website. All feature articles, fact sheets and research summaries – generally at the same time as the commenting procedure – undergo external user testing. In user testing, a group of patients or potential users comment on the texts regarding their content and understandability.

The comments submitted during the consultation period are summarized and reviewed. The summaries and the comprehensive versions are also provided to the project group (if applicable) and to the Institute Management.

With commissioned patient information, the contracting agency is also provided with a report on the comments received during consultation. Final reports on commissioned patient information are subjected to the same publication procedures as other final reports. They are sent initially to the contracting agency, the Board of Directors and the Foundation Board and then forwarded to the Board of Trustees, usually 4 weeks later. They are then published on the Institute's website www.iqwig.de (usually a further 4 weeks later). The corresponding health information is subsequently published on the Institute's website for consumers and patients at www.gesundheitsinformation.de/www.informedhealthonline.org. The patient information explaining directives by the G-BA are published only after publication of the directives themselves on www.gesundheitsinformation.de/www.informedhealthonline.org.

Corrections, improvements, and updates of the Institute's health information are made internally. If more extensive updates of content are made, external experts may be involved. A more detailed description of the Institute's updating mechanisms is provided in Chapter 5.

### 2.1.6  Working paper

The procedure for the production of a **working paper** is presented in Figure 6. All working steps are performed under the responsibility of the Institute, involving external experts or the Institute's Scientific Advisory Board, where appropriate. The internal quality assurance process is not presented in this flow chart.

Figure 6: Procedure for the production of a working paper

The production of working papers is conducted (among other things) within the framework of the **general commission** awarded by the G-BA on 21 December 2004. This commission was further specified and adapted in July 2006 and March 2008 with regard to the production of health information. According to the general commission, the Institute was commissioned "by means of documenting and analysing the relevant literature, continuously to study and assess medical developments of fundamental importance and their effects on the quality and efficiency of health care in Germany, and to relay its findings to the G-BA on a regular basis. In this context, the G-BA assumes that, within the framework of the tasks assigned in accordance with §139a (3) SGB V, the Institute will work not only on individual commissions awarded by the G-BA, but will also take on scientific projects on its own responsibility, and relay essential information on relevant health care developments to the G-BA so that it can fulfil its legal obligations. Against the background of this information, the Institute will also develop concrete proposals for individual commissions that it considers relevant."

The need to conduct independent scientific projects therefore results from the Institute's legal remit and the general commission. This also includes projects on the further development of methods, which can also be published as working papers.

The **topic selection** takes place within the Institute, particularly on the basis of the criteria defined in the general commission. The formulation of the research question may take place

by involving patient organizations or seeking the opinion of individual affected patients, especially for the definition of patient-relevant outcomes. The project outline is then prepared.

The **project outline** summarizes the main steps in the information retrieval and scientific assessment. It forms the basis for the preparation of the working paper. The project outline is not published.

The working paper presents the results of the information retrieval and scientific assessments. The quality assurance process can (optionally) include an external review. After completion the working paper is first sent to the G-BA as well as to the Foundation's Board of Directors and Foundation Council. It is then forwarded to the Foundation's Board of Trustees (usually a week later) and after 3 further weeks published on the IQWiG website. If comments on working papers are received that contain substantial unconsidered evidence, or if the Institute receives such evidence from other sources, the Institute assesses whether it considers it necessary to update the document or not. The general methodological and procedural requirements for the Institute's products apply to such an update.

## 2.2   General aspects in the preparation of products

The following procedures and aspects that are valid for all products are presented in this chapter:

- selection of external experts for collaboration in the preparation of products
- guarantee of scientific independence in the preparation of products
- review of products
- publication of products

### 2.2.1   Selection of external experts

In accordance with its legal remit, the Institute involves external experts in its work. External experts are persons who are awarded research commissions within the framework of the preparation of the Institute's products or their review or who advise the Institute on medical or other topic-related research questions. The Institute awards these commissions within the framework of a procedure that follows general procurement principles.

Current commissions to be awarded by the Institute in accordance with §139b (3) SGB V are published on the Institute's website. Commissions with a volume above the current threshold value of the procurement regulations of the European Union (EU) are advertised throughout the EU. The specific requirements regarding the suitability of applicants are published in the corresponding announcements or tendering documents.

The commissioning of external experts for benefit assessments in accordance with §35a SGB V is conducted on the basis of information provided by interested persons in a database for external experts. For inclusion in the database for external experts, the Institute's

website offers an access point via which interested experts can enter their profile, including details of their specialty and professional expertise. For the projects to be awarded, in each case the most suitable applicant of the relevant specialty is selected from this expert database by means of a criteria list and then commissioned. Further information on the selection procedure is published on the Institute's website.

### 2.2.2 Guarantee of professional independence

The scientific and professional independence of the Institute and of the products it is responsible for and publishes have their legal basis in §139a SGB V, as well as in the *Charter* of the Foundation.

### A) Guarantee of internal professional independence

The Institute's scientific staff are prohibited from performing paid external assignments that could in principle query their professional independence. All external assignments must be approved by the Institute's Management. External assignments in the broadest sense also include unpaid honorary positions such as positions on boards in organizations and societies.

### B) Guarantee of the independence of external experts

Before a contract is signed between the Institute and an external expert or external institution with regard to the preparation of a product, in accordance with §139b SGB V, "all connections to associations and contract organizations, particularly in the pharmaceutical and medical devices industries, including details on the type and amount of any remuneration received" must be disclosed to the Institute.

Following the usual practice in research to disclose such connections as potential conflicts of interest [339,345], within the framework of the selection of external experts, the Institute interprets this regulation as a responsibility to assess these disclosures with regard to the professional independence and impartiality of applicants. In this context, the Institute assesses whether a conflict of interest for the specific topic of a commission exists due to the financial connections reported. If this is the case, in a second step it is assessed whether this conflict of interest leads to serious concerns with regard to appropriate collaboration on the commission. If this is the case, collaboration on the topic of this commission is usually not possible or only possible under specific provisions. As this assessment is performed in relation to a specific commission, collaboration on topics of other commissions is indeed possible. The further process for the selection of external experts is outlined in Section 2.2.1.

The main basis of the assessment of conflicts of interest is self-disclosure using the *Form for disclosure of potential conflicts of interest*, which is published on the Institute's website. Self-disclosure refers to the following 6 types of financial connections:

- dependent occupation/employment
- advisory activities

- payments, e.g. for presentations, comments, as well as organization and/or participation in conferences and seminars

- financial support for research activities, other scientific services or patent registrations

- other financial or other cash-value support (e.g. for equipment, staff or travel expenses, without providing scientific services in return)

- shares, equity warrants or other shares in a business

The Institute reserves the right to draw upon additional information and verify the completeness and correctness of the reported information.

The names of external experts involved in the preparation of the Institute's products are usually published in these products. As a matter of principle, these publications are freely accessible via the Institute's website. The information on conflicts of interest is only published in a summarized form. In this context, for the types of connections covered by the *Form for disclosure of potential conflicts of interest*, it is only stated whether this type of connection existed or not. Specific details, for example, concerning business partners or the amount of any remuneration received, are not published.

If specifically requested by external experts or the contracting agency, or if other relevant reasons exist, it is possible to withhold publication of the external experts' names in order to ensure their independence and prevent interest-driven attempts to influence them.

### 2.2.3 Review of the Institute's products

The review of the Institute's products aims in particular at ensuring their high scientific quality. Moreover, other aims may be relevant for individual products, such as comprehensibility for the general public.

All products (including interim ones) are subjected to a comprehensive multi-stage internal quality assurance process. In addition, during the preparation of products, an external review procedure may be performed as an optional further quality assurance step. The choice of internal and external reviewers is primarily made on the basis of their methodological and/or professional expertise.

External reviewers can be identified by a literature search, the expertise of the project group, by contacting scientific societies, or by application during the tendering procedure for work on a commission, etc. In each case, external reviewers must also disclose potential conflicts of interest.

External reviewers are selected by the Institute and their number is not limited. The external reviews are assessed with regard to their relevance for the particular product; they are not published. The names of the external reviewers of final reports and rapid reports are usually

published in these documents, including a presentation of their potential conflicts of interests, in analogy to the procedure for external experts.

In addition to the external quality assurance processes described above with the involvement of reviewers selected and commissioned by the Institute, an open and independent reviewing process is guaranteed by the publication of the Institute's products and the associated opportunity to submit comments.

### 2.2.4  Publication of the Institute's products

One of the Institute's main tasks is to determine the available evidence on a topic by performing a careful assessment of the information available, and to publish the results of this assessment. It is legally specified that the Institute "must at regular intervals publicly report on its working processes and results, including the bases for decision-making" (§139a [4] SGB V).

To maintain the Institute's independence, it must be ruled out that the contracting agencies or any other interested third parties can exert any influence on the content of the reports. This could lead to conflation of scientific findings with political and/or economic aspects and/or interests. At the same time, it must be avoided that the Institute itself withholds certain findings. All the results obtained by the Institute within the framework of its legal responsibilities are therefore published as soon as possible. In the case of reports, this also includes the report plan. Product-specific features are noted in those sections in which the procedures are described. In justified exceptional cases, timelines may deviate from the stipulated norm (period between completion and publication of a document).

Unless otherwise agreed, all copyright is held by the Institute.

### 3   Benefit assessment of medical interventions

### 3.1   Patient-relevant medical benefit and harm

### 3.1.1   Definition of patient-relevant medical benefit and harm

The term **"benefit"** refers to positive causal effects, and the term **"harm"** refers to negative causal effects of a medical intervention on patient-relevant outcomes (see below). In this context, "causal" means that it is sufficiently certain that the observed effects can be ascribed solely to the intervention to be tested [549]. The terms "benefit" and "harm" refer to a comparison with a placebo (or another type of sham intervention) or no treatment.

In the case of a comparison between the medical intervention to be assessed and a clearly defined alternative medical intervention, the following terms are used in the comparative assessment of beneficial or harmful aspects (the terms are always described from the point of view of the intervention to be assessed):

- Beneficial aspects:
    - In the case of a greater benefit, the term "added benefit" is used.
    - In the case of lesser or comparable benefit, the terms "lesser" or "comparable" benefit are used.

- Harmful aspects:
    - The terms "greater", "comparable" and "lesser" harm are used.

The assessment of the evidence should preferably come to a clear conclusion that either there is proof of a(n) (added) benefit or harm of an intervention, or there is proof of a lack of a(n) (added) benefit or harm, or there is no proof of a(n) (added) benefit or harm or the lack thereof, and it is therefore unclear whether the intervention results in a(n) (added) benefit or harm.

In addition, in the case of (added) benefit or harm that is not clearly proven, it may be meaningful to perform a further categorization as to whether at least "indications" or even only "hints" of an (added) benefit or harm are available (see Section 3.1.4).

As the benefit of an intervention should be related to the patient, this assessment is based on the results of studies investigating the effects of an intervention on patient-relevant outcomes. In this connection, "patient-relevant" refers to how a patient feels, functions or survives [44]. Consideration is given here to both the intentional and unintentional effects of the intervention that in particular allow an assessment of the impact on the following patient-relevant outcomes to determine the changes related to disease and treatment:

1) mortality

2) morbidity (symptoms and complications)

3) health-related quality of life

As supplementary information, consideration can be given to the time and effort invested in relation to the disease and the intervention, as well as treatment satisfaction of patients. However, a benefit or added benefit cannot be determined on the basis of these outcomes alone. For all listed outcomes it may be necessary that an assessment is made in relation to information on how other outcomes are affected by the intervention. In the event of particularly serious or even life-threatening diseases, for example, it is usually not sufficient only to demonstrate an improvement in quality of life by application of the intervention to be assessed, if at the same time it cannot be excluded with sufficient certainty that serious morbidity or even mortality are adversely affected to an extent no longer acceptable. This is in principle consistent with the ruling by the highest German judiciary that certain (beneficial) aspects must be assessed only if therapeutic effectiveness has been sufficiently proven [71]. On the other hand, in many areas (particularly in palliative care) an impact on mortality cannot be adequately assessed without knowledge of accompanying (possibly adverse) effects on quality of life.

In accordance with §35b SGB V, the following outcomes related to patient benefit are to be given appropriate consideration: increase in life expectancy, improvement in health status and quality of life, as well as reduction in disease duration and adverse effects. These dimensions of benefit are represented by the outcomes listed above; for example, the improvement in health status and the reduction in disease duration are aspects of direct disease-related morbidity; the reduction in adverse effects is an aspect of therapy-related morbidity.

Those outcomes reliably and directly representing specific changes in health status are primarily considered. In this context, individual affected persons as well as organizations of patient representatives and/or consumers are especially involved in the topic-related definition of patient-relevant outcomes. In the assessment of quality of life and patient satisfaction, only instruments should be used that are suited for application in clinical trials and have been evaluated accordingly [160]. In addition, valid surrogate endpoints can be considered in the benefit assessment (see Section 3.1.2).

Both beneficial and harmful aspects can have different relevance for the persons affected; these aspects may become apparent through qualitative surveys or the Institute's consultations with affected persons and organizations of patient representatives and/or consumers in connection with the definition of patient-relevant outcomes (examples of corresponding methods are listed at the end of Section 3.1.4). In such a situation it may be meaningful to establish a hierarchy of outcomes. General conclusions on benefit and harm are then primarily based on proof regarding higher-weighted outcomes. Planned subgroup and sensitivity analyses are then primarily conducted for higher-weighted outcomes, whereas such analyses are not routinely conducted for the remaining ones.

Diagnostic tests can be of indirect benefit by being a precondition for therapeutic interventions through which it is possible to achieve an effect on the patient-relevant outcomes mentioned above. The precondition for the benefit of such tests is therefore the existence and the proven benefit of a treatment for patients, depending on the test result.

Interventions can also have consequences for those indirectly affected, for example, relatives and carers. If appropriate, these consequences can also be considered within the framework of the Institute's reports.

The term **"benefit assessment"** refers to the whole process of the assessment of medical interventions with regard to their positive and negative causal effects compared with a clearly defined alternative treatment, a placebo (or a different type of sham intervention), or no treatment. In this context, beneficial and harmful aspects of an intervention are initially assessed on an outcome-specific basis and then presented. In addition, a combined evaluation of outcome-related beneficial and harmful aspects is possible (see Section 3.1.4) so that, for example, when the effects on all other outcomes have been analysed, the outcome-specific "lesser harm" from an intervention (in terms of a reduction in adverse effects) can lead to the balanced conclusion of an "added benefit".

### 3.1.2 Surrogates of patient-relevant outcomes

Surrogate endpoints are frequently used in medical research as a substitute for patient-relevant outcomes, mostly to arrive at conclusions on patient-relevant (added) benefits earlier and more simply [15,180,419]. Most surrogate endpoints are, however, unreliable in this regard and can be misleading when used in a benefit assessment [207,215]. Surrogate endpoints are therefore normally considered in the Institute's benefit assessments only if they have been validated beforehand by means of appropriate statistical methods within a sufficiently restricted patient population and within comparable interventions (e.g. drugs with a comparable mode of action). A surrogate endpoint can be regarded as valid if the effect of an intervention on the patient-relevant outcome to be substituted is explained to a sufficient degree by the effect on the surrogate endpoint [27,541]. The necessity to evaluate surrogate endpoints may have particular relevance within the framework of the early benefit assessment of drugs (see Section 3.3.3), as regulatory approval procedures primarily investigate the efficacy of a drug, but not always its patient-relevant benefit or added benefit.

There is neither a standard procedure for surrogate endpoint validation nor a general best estimation method nor a generally accepted criterion which, if fulfilled, would demonstrate validity [356]. However, the current methodological literature frequently discusses correlation-based procedures for surrogate validation, with estimation of correlation measures at a study level and individual level [268]. The Institute's benefit assessments therefore give preference to validations on the basis of such procedures. These procedures usually require a meta-analysis of several randomized studies, in which both the effects on the surrogate endpoint and those on the patient-relevant outcome of interest are investigated [77,374]. Alternative methods [541] are only considered in justified exceptional cases.

For correlation-based procedures the following conditions are normally required to demonstrate validity: on the one hand, a high correlation between the surrogate and the patient-relevant outcome at the individual level, and on the other hand, a high correlation between effects on the surrogate and effects on the patient-relevant outcome at a study level [77,78]. As in the Institute's benefit assessments, conclusions related to groups of patients are drawn, the assessment of the validity of a surrogate endpoint is primarily based on the degree of correlation at the level of treatment effects, i.e. the study level. In addition to the degree of correlation, for the assessment of validity of a surrogate endpoint the reliability of results of the validation process is considered. For this purpose, various criteria are drawn upon [268]. For example, associations observed between a surrogate endpoint and the corresponding patient-relevant outcome for an intervention with a specific mode of action are not necessarily applicable to other interventions used to treat the same disease, but with a different mode of action [179,207,215,356]. The studies on which the validation was based must therefore have been conducted in patient populations and interventions that allow conclusions on the therapeutic indication investigated in the benefit assessment as well as on the test intervention and comparator intervention. In order to assess transferability, in validation studies including various disease entities or interventions, analyses on heterogeneity should at least be available.

In the event that a surrogate endpoint cannot be validated conclusively (e.g. if correlation is not high enough), it is also possible to apply the "surrogate threshold effect (STE) concept" [76,268]. For this purpose, the effect on the surrogate resulting from the studies included in the benefit assessment is related to the STE [78,374].

For the Institute's benefit assessments, conclusions on patient-relevant outcomes can be drawn from the effects on the surrogate, depending on verification of the validity of the surrogate or the evaluation of the STE. The decisive factor for the first point is the degree of correlation of the effects on the surrogate and the patient-relevant outcome and the reliability of validation in the validation studies. In the evaluation of an STE, the decisive criterion is the size of the effect on the surrogate in the studies included in the benefit assessment compared with the STE. In the case of a statistically significant effect on the surrogate endpoints, all gradations of conclusions on the (added) benefit with regard to the corresponding patient-relevant outcome according to Section 3.1.4 are possible, depending on the constellation.

Surrogate endpoints that are not valid or for which no adequate validation procedure was conducted can nevertheless be presented in the Institute's reports. However, independent of the observed effects, such endpoints are not suited to provide proof of verification of an (added) benefit of an intervention.

Depending on the proximity to a corresponding patient-relevant outcome, the literature uses various other terms to describe surrogate endpoints (e.g. intermediate endpoint). However, we dispense with such a distinction here, as the issue of the necessary validity remains unaffected by this. In addition it should be considered that an endpoint can at the same time represent a

patient-relevant outcome and, beyond this, can also be regarded as a surrogate (i.e. a substitute) for a different patient-relevant outcome.

### 3.1.3  Assessment of the harm of medical interventions

The use of any type of medical intervention (drug, non-drug, surgical, diagnostic, preventive, etc.) carries per se the risk of adverse effects. In this context, the term "adverse effects" refers to all events and effects representing individually perceived or objectively detectable physical or mental harm that may to a greater or lesser extent cause a short- or long-term reduction in life expectancy, an increase in morbidity, or impairment in quality of life. It should be noted that if the term "adverse effects" is used, a causal relationship to the intervention is assumed, whereas the issue of causality still remains open with the term "adverse events" [92].

The term "harm" describes the occurrence of adverse effects when using a medical intervention. The description of harm is an essential and equal component in the benefit assessment of an intervention. It ensures the informed, population-related, but also individual weighing of benefit and harm [557]. A prerequisite for this is that the effect sizes of a medical intervention can be described by means of the data available, both for its desired as well as its adverse effects, and compared with therapy alternatives, for example.

However, in a systematic review, the analysis, assessment, and reporting of the harm of a medical intervention are often far more difficult than those of the (added) benefit. This applies in particular to unexpected adverse events [92]. Studies are typically designed to measure the effect of a medical intervention on a few predefined outcomes. In most cases, these are outcomes representing effectiveness, while adverse effects are concomitantly recorded as adverse events. The results for adverse events depend heavily on the underlying methods for data collection. For example, explicit queries on defined adverse events normally result in the determination of higher event rates than do general queries [41,279]. To detect unexpected adverse events in particular, general queries about the well-being of patients are however required. In addition, studies designed to specifically detect rare, serious adverse effects (including the description of a causal relationship to the medical intervention) are considerably underrepresented in medical research [48,149,278]. Moreover, reporting of adverse events in individual studies is of poor quality, which has also led to amendment of the CONSORT[7] statement for RCTs [277]. Finally, the systematic assessment of the adverse effects of an intervention is also made more difficult by the fact that the corresponding coding in bibliographic databases is insufficient, so that the specific search for relevant scientific literature often produces an incomplete picture [112].

The obstacles noted above often make the investigation of harm more difficult. In cases where complete clinical study reports are available for the assessment, at least sufficient data transparency is also given for adverse events. However, it is still necessary to find a

---

[7] Consolidated Standards of Reporting Trials

---

meaningful balance between the completeness of the evaluation of aspects of harm and the resources invested. Consequently, it is necessary to limit the evaluation and reporting to relevant adverse effects. In particular, those adverse effects can be defined as relevant that may

- completely or almost completely offset the benefit of an intervention

- substantially vary between 2 or more otherwise equivalent treatment options

- occur predominantly with treatment options that may be particularly effective

- have a dose-effect relationship

- be regarded by patients as especially important

- be accompanied by serious morbidity or even increased mortality, or be associated with substantial impairment in quality of life

The Institute observes the following principles when evaluating and reporting adverse effects: In the benefit assessment, the initial aim is to compile a selection of potentially relevant adverse effects that are essential in deciding for or against the use of the intervention to be assessed. In this context, the selection of adverse effects and events is made in accordance with the criteria outlined above. This compilation is made within the framework of the preliminary literature search for the particular research question, especially on the basis of data from controlled intervention studies in which the benefit of the intervention was specifically investigated. In addition, and if appropriate, the compilation is made on the basis of available epidemiological data (e.g. from cohort or case-control studies), as well as pharmacovigilance and regulatory data, etc. In individual cases, data obtained from animal trials and experiments to test pathophysiological constructs may be useful. The compilation of potentially relevant adverse effects described above forms the foundation for assessment of harm on the basis of the studies included in the benefit assessment. In this context, if possible and meaningful, pooled analyses (e.g. overall rates of serious adverse events) may also be drawn upon.

### 3.1.4  Outcome-related assessment

The benefit assessment and the estimation of the extent of the (un)certainty of results generally follow international EBM standards as developed, for example, by the GRADE[8] group [23].

Medical interventions are compared with other interventions, sham interventions (e.g. placebo), or no intervention in respect of their effects on defined patient-relevant outcomes, and their (added) benefit and harm are described in summary. For this purpose, on the basis of the analysis of the scientific data available, for each predefined patient-relevant outcome

---

[8] Grading of Recommendations, Assessment, Development and Evaluation

separately a conclusion on the evidence base of the (added) benefit and harm is drawn in 4 levels with regard to the respective certainty of the conclusion: The data provide either "proof" (highest certainty of conclusions), an "indication" (medium certainty of conclusions), a "hint" (weakest certainty of conclusions) in respect of the benefit or harm of an intervention, or none of these 3 situations applies. The latter is the case if no data are available or the data available do not allow any of the other 3 conclusions to be drawn.

Depending on the research question, the conclusions refer to the presence or lack of a(n) (added) benefit or harm. The prerequisite for conclusions on the lack of a(n) (added) benefit or harm are well-founded definitions of irrelevance ranges (see Section 7.3.6).

The certainty of results is an important criterion for the inference of conclusions on the evidence base. In principle, every result from an empirical study or systematic review of empirical studies is potentially uncertain and therefore the certainty of results must be examined. In this context, one distinguishes between qualitative and quantitative certainty of results. The qualitative certainty of results is impaired by systematic errors (bias; see Section 7.3.11) such as information errors, selection errors and confounding. The quantitative certainty of results is influenced by random errors caused by sampling (statistical uncertainty).

The qualitative certainty of results is thus determined by the study design, from which evidence levels can be inferred (see Section 7.1.3). It is also determined by (outcome-related) measures for further prevention or minimization of potential bias, which must be assessed depending on the study design (see Section 7.1.4). Such measures include, for example, the blinded assessment of outcomes, an analysis based on all included patients (potentially supported by the application of adequate imputation methods for missing values), and, if appropriate, the use of valid measurement instruments.

The quantitative certainty of results is directly connected to the sample size (i.e. the number of patients investigated in a study or the number of [primary] studies included in a systematic review), as well as to the variability observed within and between studies. If the underlying data allow for this, the statistical uncertainty can be quantified and assessed as the standard error or confidence interval of parameter estimates (precision of the estimate).

The Institute uses the following 3 categories to grade the degree of qualitative certainty at the individual study level and outcome level:

- **high qualitative certainty of results:** results on an outcome from a randomized study with a low risk of bias

- **moderate qualitative certainty of results**: results on an outcome from a randomized study with a high risk of bias

- **low qualitative certainty of results:** results on an outcome from a non-randomized comparative study

In the inference of the evidence base for an outcome, the number of available studies, their qualitative certainties of results, as well as the effects found in the studies are of crucial importance. If at least 2 studies are available, it is first distinguished whether, due to existing heterogeneity within a meta-analysis (see Section 7.3.8), a common effect estimate can be meaningfully formed or not. In the case of homogenous results that can be meaningfully pooled, the common effect estimate must be statistically significant to infer proof, an indication or a hint according to the existing certainty of results. If the estimated results are too heterogeneous to meaningfully form a pooled common effect estimate, one distinguishes between effects that are "not in the same direction", "moderately in the same direction" and "clearly in the same direction". These are defined as follows:

Effects in the same direction are present if the prediction interval for displaying heterogeneity in a meta-analysis with random effects (see Section 7.3.8) is presented and does not cover the zero effect. In other cases (no presentation of the prediction interval or this interval covers the zero effect) effects in the same direction are present in the following situation:

The effect estimates of 2 or more studies point in the same direction. For these "directed" studies, all of the following conditions apply:

- The overall weight of these studies is ≥ 80%.

- At least 2 of these studies show statistically significant results.

- At least 50% of the weight of these studies is based on statistically significant results.

In this context, the weights of these studies generally come from a meta-analysis with random effects (see Section 7.3.8). If no meta-analysis is meaningful, the relative sample size corresponds to the weight.

If effects in the same direction are moderately or clearly in the same direction, if possible, a decision is made on the basis of the location of the prediction interval. As the prediction interval is generally only presented if at least 4 studies are available (see Section 7.3.8), the classification into effects that are moderately or clearly in the same direction depends on the number of studies.

- 2 studies: Effects in the same direction are always clearly in the same direction.

- 3 studies:

  - All studies show statistically significant results. The effects in same direction are clearly in the same direction.

  - Not all of the 3 studies show statistically significant results. The effects in the same direction are moderately in the same direction.

- 4 or more studies:

▫ The prediction interval does not cover the zero effect: The effects in the same direction are clearly in the same direction.

▫ The prediction interval covers the zero effect: The effects in the same direction are moderately in the same direction.

For the case that the available studies show the same qualitative certainty of results or only one study is available, with these definitions the regular requirements for the evidence base to infer conclusions with different certainties of conclusions can be specified. As described above, the Institute distinguishes between 3 different certainties of conclusions: "proof", "indication" and "hint".

A conclusion on proof generally requires that a meta-analysis of studies with a high qualitative certainty of results shows a corresponding statistically significant effect. If a meta-analysis cannot be conducted, at least 2 studies conducted independently of each other and showing a high qualitative certainty of results and a statistically significant effect should be present, the results of which are not called into question by further comparable studies with a high certainty of results (consistency of results). These 2 studies do not need to have an exactly identical design. Which deviations in design between studies are still acceptable depends on the research question. Accordingly, a meta-analysis of studies with a moderate qualitative certainty of results or a single study with a high qualitative certainty of results can generally provide only a hint, despite statistically significant effects.

On the basis of only one study, in exceptional cases proof can be inferred for a specific (sub)population with regard to an outcome. This requires the availability of a clinical study report according to the International Conference on Harmonization (ICH) guidelines and the fulfilment of the other requirements stipulated for proof. In addition, the study must fulfil the following specific requirements:

▪ The study is a multi-centre study with at least 10 centres.

▪ The effect estimate observed has a very small corresponding p-value ($p < 0.001$).

▪ The result is consistent within the study. For the (sub)population of interest, analyses of different further subpopulations are available (particularly subpopulations of study centres), which in each case provide evaluable and sufficiently homogeneous effect estimates. This assessment of consistency is only possible for binary data if a certain minimum number of events has occurred.

▪ The analyses for the subpopulations addressed above are available for all relevant outcomes, i.e. these analyses are not restricted to individual selected outcomes.

It is possible that in the case of the existence of only one study, which alone provides only an indication or a hint, the evidence base may be changed by additional indirect comparisons. However, high methodological demands must be placed on indirect comparisons (see Section 7.3.9). In addition, in the case of a homogeneous data situation, it is possible that by adding

indirect comparisons the precision of the effect estimate increases, which plays an important role when determining the extent of added benefit (see Section 3.3.3).

A meta-analysis of studies with a low qualitative certainty of results or an individual study with a moderate qualitative certainty of results (both with a statistically significant effect) generally only provides a hint.

An overview of the regular operationalization is shown in Table 2. In justified cases further factors influence these evaluations. The assessment of surrogate endpoints (see Section 3.1.2), the presence of serious deficiencies in study design or justified doubts about the transferability to the treatment situations in Germany may, for example, lead to a reduction in the certainty of conclusions. On the other hand, great effects or a clear direction of an existing risk of bias, for example, can justify an increase in certainty.

Table 2: Certainty of conclusions regularly inferred for different evidence situations if studies with the same qualitative certainty of results are available

| | | **Number of studies** | | | | |
|---|---|---|---|---|---|---|
| | | 1 (with statistically significant effect) | $\geq 2$ | | | |
| | | | Homogeneous | Heterogeneous | | |
| | | | Meta-analysis statistically significant | Effects in the same direction[a] | | |
| | | | | Clear | Moderate | No |
| **Qualitative certainty of results** | High | Indication | Proof | Proof | Indication | – |
| | Moderate | Hint | Indication | Indication | Hint | – |
| | Low | – | Hint | Hint | – | – |
| a: See text for explanation of term. | | | | | | |

If several studies with a different qualitative certainty of results are available, then first only the studies with the higher-quality certainty of results are examined, and conclusions on the evidence base are inferred on this basis according to Table 2. In the inference of conclusions on the evidence base for the whole study pool the following principles then apply:

- The conclusions on the evidence base, when restricted to higher-quality studies, are not weakened by the addition of the other studies, but at best upgraded.

- The confirmation (replication) of a statistically significant result of a study with a high qualitative certainty of results, which is required to infer proof, can be provided by one or more results of moderate (but not low) qualitative certainty of results. In this context the weight of the study with a high qualitative certainty of results should have an appropriate size (between 25 and 75%).

- If the meta-analytical result for the higher-quality studies is not statistically significant or if no effects in the same direction are shown in these studies, then conclusions on the evidence base are to be inferred on the basis of results of the whole study pool, whereby the certainty of conclusions is determined by the minimum qualitative certainty of results of all studies included.

According to these definitions and principles, a corresponding conclusion on benefit is inferred for each outcome separately. Considerations on the assessment across outcomes are presented in the following section (see Section 3.1.5).

### 3.1.5 Summarizing assessment

These conclusions, drawn separately for each patient-relevant outcome within the framework of the deduction of conclusions on the evidence base, are then summarized (as far as possible) in one evaluating conclusion in the form of a weighing of benefits and harms. If proof of a(n) (added) benefit and/or harm exists with regard to Outcomes 1 to 3 of Section 3.1.1, the Institute presents (insofar as is possible on the basis of the data available)

1) the benefit

2) the harm

3) (if appropriate) the weighing of benefit and harm

In this context, characteristics related to age, gender, and personal circumstances are considered.

One option in the conjoint evaluation of benefit and harm is to compare the outcome-related beneficial and harmful aspects of an intervention. In this context, the effects on all outcomes (qualitative or semi-quantitative) are weighed against each other, with the aim of drawing a conclusion across outcomes with regard to the benefit or added benefit of an intervention. A further option in the conjoint evaluation is to aggregate the various patient-relevant outcomes into a single measure. In this case the conclusions made by the Institute would be weighted for each individual patient-relevant outcome, for example, by being considered in a summarizing score. The conjoint evaluation of benefit and harm is specified depending on the topic of interest and should, if this is prospectively possible, be described in the report plan (protocol), or otherwise in the preliminary report. A quantitative weighting using summarizing scores or indices should be done prospectively at the time the outcomes to be investigated are selected.

So-called "utility values" are often used to determine a person's state of health; such values are supposed to express positive and negative aspects perceived by respondents in an index score. If the duration of the corresponding states of health is incorporated, these utility values can, for example, be transformed into quality-adjusted life years (QALYs). The recording and

calculation of utility values is, for example, presented in: [141,354,464]. In health economic evaluations the Institute can draw on QALYs as an overall measure of benefit [267].

There is no resumption here to the scientific debate about the ethical and methodological problems of the QALY concept itself and their solution or a linked willingness-to-pay threshold in health economic evaluations, as well as the use of the QALY for pure weighing of benefit and harm. In this context we refer to the following selected publications: [121,137,195,232,343,364,394,406,420,530].

If a measure of overall benefit for the comparison of interventions is to be determined, in addition to the utility value QALY, procedures for multi-criteria decision-making or determining preferences such as the analytic hierarchy process (AHP) and the conjoint analysis (CA) can be applied.

For the AHP [135,136] a problem in decision-making is broken down into various criteria. These are then arranged in a hierarchy. For example, a drug can be assessed by means of the criteria "effect", "adverse effects", and "quality of life". The criterion "effect" can be broken down into further subcriteria that correspond to outcomes [260]. Participants in the AHP then respond to questions about the criteria in a binary way, i.e. on a specified scale they choose how much more a certain criterion means to them than another. By means of a procedure for matrix multiplication [440-442] the weights for the criteria and subcriteria can be determined via a so-called "right eigenvector"; these weights must add up to 1.

The CA belongs to the group of stated-preference techniques [59]. A decision is also broken down into attributes. The respondents are then confronted with a set of (theoretical) scenarios. In each scenario all attributes and their range of effect sizes are compared with each other. From the choice of scenarios a weighting factor for each attribute is then calculated by means of a regression model. These weights can in turn be standardized to 1.

## 3.2 Special aspects of the benefit assessment

### 3.2.1 Impact of unpublished study results on conclusions

An essential prerequisite for the validity of a benefit assessment is the complete availability of the results of the studies conducted on a topic. An assessment based on incomplete data or possibly even selectively compiled data may produce biased results [165,271] (see also Section 7.3.11).

The distortion of published evidence through publication bias and outcome reporting bias has been described comprehensively in the literature [142,363,486]. In order to minimize the consequences of such bias, the Institute has extended information retrieval beyond a search in bibliographic databases, for example, by screening trial registries. In addition, at the beginning of an assessment the Institute normally contacts the manufacturers of the drugs or medical devices to be assessed, and requests the transfer of complete information on studies investigating these interventions (see also Section 6.1.5).

This transfer of information by manufacturers can only solve the problem of bias caused by unpublished evidence if the transfer is itself not selective but complete. An incomplete transfer of information carries a risk of bias for the result of the benefit assessment. This risk should be considered by the Institute in the conclusions of a benefit assessment.

Table 3 below describes what constellations carry a risk of bias for assessment results, and what consequences arise for the conclusions of a benefit assessment.

If the data transfer was complete and no evidence is available that a relevant amount of data is missing, bias seems improbable (Scenario 1). The inferences drawn from the assessment of data can therefore be adopted without limitation in the conclusions of the benefit assessment.

Table 3: Scenarios for data transfer by third parties and consequences for the conclusions of a benefit assessment

| Scenario | Data transfer by third parties (e.g. manufacturer data) | Evidence that a relevant amount of data is missing | Bias | Assessment/Impact on the conclusions |
|---|---|---|---|---|
| 1 | Complete | No | Improbable | No limitation of the conclusions of the benefit assessment |
| 2 | Incomplete | No | Possible | Conclusions are made with reservations |
| 3 | Incomplete | Yes | Probable | Description of the available and missing data; no proof (or indication or hint) of benefit or harm |
| 4 | Complete | Yes (e.g. other manufacturers, investigator-initiated trials) | Possible | Conclusions are drawn with reservations |

If the data transfer is incomplete, the consequences for the conclusions depend on whether additional search steps demonstrate that a relevant amount of data is missing. If this is not the case (Scenario 2), bias may still be possible, as data transfer may have been selective and further unpublished data may exist that were not identified by the search steps. In such cases the conclusions are therefore drawn with reservations. If it was demonstrated that a relevant amount of data is missing (Scenario 3), it can be assumed that the data transfer was selective. In this situation, further analysis of the available limited data and any conclusions inferred from them with regard to benefit or harm are probably seriously biased and therefore do not form a valid decision-making basis for the G-BA. Consequently, no proof (nor indication nor hint) of a benefit or harm of the intervention to be assessed can be determined in this situation, independently of whether the available data show an effect of the intervention or not.

If the manufacturer completely transfers data and additional literature searches demonstrate that a relevant amount of data from studies inaccessible to the manufacturer is missing (Scenario 4), then no selective data transfer by the manufacturer is evident. In this situation, bias caused by missing data is still possible. The conclusions are therefore drawn with reservation.

### 3.2.2  Dramatic effect

If the course of a disease is certainly or almost certainly predictable, and no treatment options are available to influence this course, then proof of a benefit of a medical intervention can also be provided by the observation of a reversal of the (more or less) deterministic course of the disease in well-documented case series of patients. If, for example, it is known that it is highly probable that a disease leads to death within a short time after diagnosis, and it is described in a case series that, after application of a specific intervention, most of those affected survive for a longer period of time, then this "dramatic effect" may be sufficient to provide proof of a benefit. An example of such an effect is the substitution of vital hormones in diseases with a failure of hormone production (e.g. insulin therapy in patients with diabetes mellitus type 1). An essential prerequisite for classification as a "dramatic effect" is sufficiently reliable documentation of the fateful course of the disease in the literature and of its diagnosis in the patients included in the study to be assessed. In this context, possible harms of the intervention should also be taken into account. Glasziou et al. [202] have attempted to operationalize the classification of an intervention as a "dramatic effect". In a first approach they propose to regard an observed effect as not explicable solely by the impact of confounding factors if it was significant at a level of 1% and, expressed as the relative risk, exceeded the value of 10 [202]. This magnitude serves as orientation for the Institute and does not represent a rigid threshold. Glasziou et al. [202] made their recommendation on the basis of results of simulation studies, according to which an observed relative risk of 5 to 10 can no longer be plausibly explained only by confounding factors. This illustrates that a corresponding threshold also depends on the attendant circumstances (among other things, the quality of studies used to determine the existence of a dramatic effect). This dependence is also reflected in the recommendations of other working groups (e.g. the GRADE group) [321].

If, in the run-up to the work on a specific research question, sufficient information is available indicating that a dramatic effect caused by the intervention to be assessed can be expected (e.g. because of a preliminary literature search), then information retrieval will also include a search for studies that show a higher uncertainty of results due to their design.

### 3.2.3  Study duration

Study duration is an essential criterion in the selection of studies relevant to the benefit assessment. In the assessment of a therapeutic intervention for acute diseases where the primary objective is, for example, to shorten disease duration and alleviate acute symptoms, it is not usually meaningful to call for long-term studies, unless late complications are to be

expected. On the other hand, in the assessment of therapeutic interventions for chronic diseases, short-term studies are not usually suitable to achieve a complete benefit assessment of the intervention. This especially applies if treatment is required for several years, or even lifelong. In such cases, studies covering a treatment period of several years are particularly meaningful and desirable. As both benefits and harms can be distributed differently over time, in long-term interventions the meaningful comparison of the benefits and harms of an intervention is only feasible with sufficient certainty if studies of sufficient duration are available. However, individual aspects of the benefits and harms may quite well be investigated in short-term studies.

With regard to the selection criterion "minimum study duration", the Institute primarily follows standards for demonstrating the effectiveness of an intervention. In the assessment of drugs, the Institute will in particular resort to information provided in guidelines specific to therapeutic indications, which are published by regulatory authorities (e.g. [162]). As the benefit assessment of an intervention also includes aspects of harm, the generally accepted standards in this respect are also relevant when determining the minimum study duration. Moreover, for long-term interventions as described above, the Institute will resort to the relevant guidelines for the criterion "long-term treatment" [263]. In individual cases, the Institute may deviate from this approach (and will justify this deviation), for example, if a topic requires longer follow-up, or if specific (sub)questions apply to a shorter period. Such deviations may also be indicated if short-term effects are a subject of the assessment (e.g. in the assessment of newly available/approved interventions and/or technologies where no appropriate treatment alternative exists).

### 3.2.4   Patient-reported outcomes

The patient-relevant dimensions of benefit outlined in Section 3.1.1 can also include patient-reported outcomes (PROs). In addition to health-related quality of life and treatment satisfaction, PROs can also cover other dimensions of benefit, for example, disease symptoms. As in the assessment of quality of life and treatment satisfaction, instruments are required that are suitable for use in clinical trials [160]. In the selection of evidence (especially study types) to be considered for the demonstration of an effect, the same principles as with other outcomes usually apply [183]. This means that also for PROs (including health-related quality of life and treatment satisfaction), RCTs are best suited to demonstrate an effect.

As information on PROs is subjective due to their nature, open studies in this area are of limited validity. The size of the effect observed is an important decision criterion for the question as to whether an indication of a benefit of an intervention with regard to PROs can be inferred from open studies. Empirical evidence shows a high risk of bias for subjective outcomes in open studies [555]. This should be considered in their interpretation (see also Sections 7.1.4 and 7.3.4). However, situations are conceivable where blinding of physicians and patients is not possible. In such situations, if possible, other efforts are required to

minimize and assess bias (e.g. blinded documentation and assessment of outcomes). Further aspects on the quality assessment of studies investigating PROs are outlined in [183].

### 3.2.5   Benefits and harms in small populations

In small populations (e.g. patients with rare diseases or special subgroups of patients with common diseases), there is no convincing argument to deviate in principle from the hierarchy of evidence levels. In this connection, it is problematical that no international standard definition exists as to what is to be understood under a "rare" disease [552]. Independent of this, patients with rare diseases also have the right to the most reliable information possible on treatment options [157]. Non-randomized studies require larger sample sizes than randomized ones because of the need of adjustment for confounding factors. However, due to the rarity of a disease it may sometimes be impossible to include enough patients to provide the study with sufficient statistical power. A meta-analytical summary of smaller studies may be particularly meaningful in such cases. Smaller samples generally result in lower precision in an effect estimate, accompanied by wider confidence intervals. Because of the relevance of the assumed effect of an intervention, its size, the availability of treatment alternatives, and the frequency and severity of potential therapy-related harms, for small sample sizes it may be meaningful to accept a higher p-value than 5% (e.g. 10%) to demonstrate statistical significance, thus increasing quantitative uncertainty. Similar recommendations have been made for other problematical constellations [159]. Such an approach must, however, be specified a priori and well justified. Likewise, for small sample sizes it may be more likely that is necessary to substitute a patient-relevant outcome that occurs too rarely with surrogate endpoints. However, these surrogates must also be valid for small sample sizes [161].

In the case of extremely rare diseases or very specific disease constellations, the demand for (parallel) comparative studies may be inappropriate [552]. Nevertheless, in such cases it is also possible at least to document and assess the course of disease in such patients appropriately, including the expected course without applying the intervention to be assessed (e.g. using historical patient data) [72]. The fact that a situation is being assessed involving an extremely rare disease or a very specific disease constellation is specified and explicitly highlighted in the report plan.

### 3.3   Benefit assessment of drugs

One main objective of the benefit assessment reports on drugs is to support the G-BA's decisions on directives concerning the reimbursement of drugs by the SHI. For this purpose, it is necessary to describe whether a drug's benefit has been demonstrated (or whether, when compared with a drug or non-drug alternative, a higher benefit [added benefit] has been demonstrated).

The G-BA's decisions on directives do not usually consider particular cases, but the general one. Consequently, the Institute's reports do not usually refer to decisions on particular cases.

Because of the objective of the Institute's benefit assessments, these assessments only include studies with an evidence level principally suited to demonstrate a benefit of an intervention. Thus, studies that can only generate hypotheses are generally not relevant for the benefit assessment. The question as to whether a study can demonstrate a benefit mainly depends on the certainty of results of the data analysed.

### 3.3.1   Relevance of the drug approval status

The commissioning of the Institute by the G-BA to assess the benefit of drugs usually takes place within the framework of the approval status of the drug to be investigated (therapeutic indication, dosage, contra-indications, concomitant treatment, etc.). For this reason, the Institute's recommendations to the G-BA, which are formulated in the conclusions of the benefit assessment report, usually refer to the use of the assessed drug within the framework of the current approval status.

It is clarified on a project-by-project basis how to deal with studies (and the evidence inferred from them) that were not conducted according to the use of a drug as outlined in the approval documents. In principle, it is conceivable that studies in which a drug was used outside the scope of the approval status described in the Summary of Product Characteristics ("off-label use"), over- or underestimated a drug's benefit and/or harm. This may lead to a misjudgement of the benefit and/or harm in patients treated within the framework of the drug's approval status. However, if it is sufficiently plausible or has even been demonstrated that the results obtained in these studies are applicable to patients treated according to the drug's approval status, these results can be considered in the benefit assessment.

Therefore, for studies excluded from the assessment only because they were off-label studies (or because it was unclear whether they fulfilled the requirements of the approval status), each case is assessed to establish to what extent the study results are applicable to patients treated according to the approval requirements.

Results from off-label studies are regarded as "applicable" if it is sufficiently plausible or has been demonstrated that the effect estimates for patient-relevant outcomes are not greatly affected by the relevant characteristic of the drug approval status (e.g. pretreatment required). As a rule, the equivalence of effects should be proven with appropriate scientific studies. These studies should be targeted towards the demonstration of equivalence of the effect between the group with and without the characteristic. Results applicable to patients treated according to a drug's approval status can be considered in the conclusion of the assessment.

Results from studies are regarded as "not applicable" if their applicability has not been demonstrated and if plausible reasons against the transferability of results exist. As a rule, study results are regarded to be "not applicable" if, for example, the age range or disease severity treated lay outside the approved range or severity, if off-label combinations including other active ingredients were used, or if studies were conducted in patients with contra-

indications for the intervention investigated. The results of these studies are not presented in the reports, as they cannot be considered in the assessment

If results from off-label studies are regarded as applicable, this is specified in the report plan. As a rule the results of studies showing the following characteristics are discussed, independently of the applicability of study results to the use specified in the approval of the drug:

- They refer to patients with the disease specified in the commission.
- They refer to patients treated with the drug to be assessed.
- They are of particular relevance due to factors such as sample size, study duration, or outcomes investigated.

### 3.3.2  Studies on the benefit assessment of drugs

The results of the Institute's benefit assessment of drugs may have an impact on patient health care in Germany. For this reason, high standards are required regarding the certainty of results of studies included in the benefit assessment.

The certainty of results is defined as the certainty with which an effect (or the lack of an effect) can be inferred from a study. This refers to both "positive" aspects (benefit) as well as "negative" aspects (harm). The certainty of results of an individual study is essentially influenced by 3 components:

- the study design
- the internal validity (which is design-specific and determined by the specific way the study was conducted)
- the size of an expected or observed effect

In the benefit assessment of drugs, not only individual studies are assessed, but the results of these studies are incorporated into a systematic review. The certainty of results of a systematic review is in turn based on the certainty of results of the studies included. In addition, it is determined in particular by the following factor:

- the consistency of the results of several studies

The study design has considerable influence on the certainty of results insofar as a causal association between intervention and effect cannot usually be shown with prospective or retrospective observational studies, whereas controlled intervention studies are in principle suited for this purpose [214]. This particularly applies if other factors influencing results are completely or almost completely eliminated. For this reason, an RCT represents the gold standard in the assessment of drug and non-interventions [322].

In the assessment of drugs, RCTs are usually possible and practically feasible. As a rule, the Institute therefore considers RCTs in the benefit assessment of drugs and only uses non-randomized intervention studies or observational studies in justified exceptional cases. Reasons for exception are, on the one hand, the non-feasibility of an RCT (e.g. if the therapist and/or patient have a strong preference for a specific therapy alternative) or, on the other, the fact that other study types may also provide sufficient certainty of results for the research question posed. For diseases that would be fatal within a short period of time without intervention, several consistent case reports may provide sufficient certainty of results that a particular intervention prevents this otherwise inevitable course [338] (dramatic effect, see also Section 3.2.2). The special obligation to justify a non-randomized design when testing drugs can also be found within the framework of drug approval legislation in the directives on the testing of medicinal products (Directive 2001/83/EC, Section 5.2.5 [311]).

In the preparation of the report plan (see also Section 2.1.1), the Institute therefore determines beforehand which study types can be regarded as feasible on the basis of the research question posed, and provide sufficient certainty of results (with high internal validity). Studies not complying with these minimum quality standards (see also Section 7.1.4) are not given primary consideration in the assessment process.

Sections 3.1.4 and 7.1 present information on the assessment of the internal validity of studies, as well as on further factors influencing certainty of results, such as the consistency of the results of several studies and the relevance of the size of the effect to be expected.

In addition to characterizing the certainty of results of the studies considered, it is necessary to describe whether – and if yes, to what extent – the study results are transferable to local settings (e.g. population, health care sector etc.), or what local study characteristics had (or could have had) an effect on the results or their interpretation. From this perspective, studies are especially relevant in which the actual German health care setting is represented as far as possible. However, the criteria for certainty of results outlined above must not be ignored. Finally, the transferability of study results (generalizability or external validity) must be assessed in a separate process initially independent of the study design and quality.

### 3.3.3   Benefit assessment of drugs according to §35a SGB V

A benefit assessment of a drug according to §35a SGB V is based on a dossier of the pharmaceutical company in which the company provides the following information:

1) approved therapeutic indications

2) medical benefit

3) added medical benefit compared with an appropriate comparator therapy

4) number of patients and patient groups for whom a therapeutically relevant added benefit exists

5)  cost of treatment for the SHI

6)  requirements for quality-assured usage of the drug

The requirements for form and content of the dossier are outlined in dossier templates, which are a component of the G-BA's Code of Procedure [198]. In the dossier, specifying the validity of the evidence, the pharmaceutical company must describe the likelihood and the extent of added benefit of the drug to be assessed compared with an appropriate comparator therapy. The information provided must be related both to the number of patients and to the extent of added benefit. The costs for the drug to be assessed and the appropriate comparator therapy must be declared (based on the pharmacy sales price und taking the Summary of Product Characteristics and package information leaflet into account).

The probability of the added benefit describes the certainty of conclusions on the added benefit. In the dossier, the extent of added benefit should be described according to the categories of the Regulation for Early Benefit Assessment of New Pharmaceuticals (ANV[9]) (major, considerable, minor, non-quantifiable added benefit; no added benefit proven; benefit of the drug to be assessed smaller than benefit of the appropriate comparator therapy) [70].

In the benefit assessment the validity and completeness of the information in the dossier are examined. It is also examined whether the comparator therapy selected by the pharmaceutical company can be regarded as appropriate in terms of §35a SGB V and the ANV. In addition, the Institute assesses the effects described in the documents presented, taking the certainty of results into account. In this assessment, the qualitative and quantitative certainty of results within the evidence presented, as well as the size of observed effects and their consistency, are appraised. The benefit and cost assessments are conducted on the basis of the standards of evidence-based medicine described in this methods paper and those of health economic standards, respectively. As a result of the assessment, the Institute presents its own conclusions, which may confirm or deviate from those arrived at by the pharmaceutical company (providing a justification in the event of deviation).

The operationalization for determining the extent of added benefit comprises 3 steps:

1)  In the first step the probability of the existence of an effect is examined for each outcome separately (qualitative conclusion). For this purpose, the criteria for inferring conclusions on the evidence base are applied (see Section 3.1.4). Depending on the quality of the evidence, the probability is classified as a hint, an indication or proof.

2)  In the second step, for those outcomes where at least a hint of the existence of an effect was determined in the first step, the extent of the effect size is determined for each outcome separately (quantitative conclusion). The following quantitative conclusions are possible: major, considerable, minor, and non-quantifiable.

[9]Arzneimittel-Nutzenbewertungsverordnung, AM-NutzenV

3) In the third and last step, the overall conclusion on the added benefit according to the 6 specified categories is determined on the basis of all outcomes, taking into account the probability and extent at outcome level within the overall picture. These 6 categories are as follows: major, considerable, minor, and non-quantifiable added benefit; no added benefit proven; the benefit of the drug under assessment is less than the benefit of the appropriate comparator therapy.

The quality of the outcome, as well as the effect size, are essential in determining the extent at outcome level in the second step. The rationale for this operationalization is presented in the Appendix A *Rationale of the methodological approach for determining the extent of added benefit*. The basic approach aims to derive thresholds for confidence intervals for relative effect measures depending on the effects to be achieved, which in turn depend on the quality of the outcomes and the extent categories.

It will not always be possible to quantify the extent at outcome level. For instance, if a statistically significant effect on a sufficiently valid surrogate is present, but no reliable estimate of this effect on a patient-relevant outcome is possible, then the (patient-relevant) effect cannot be quantified. In such and similar situations, an effect of a non-quantifiable extent is concluded, with a corresponding explanation.

On the basis of the case of a quantifiable effect, the further approach depends on the scale of the outcome. One distinguishes between the following scales:

- binary (analyses of 2x2 tables)

- time to event (survival time analyses)

- continuous or quasi-continuous, in each case with available responder analyses (analyses of mean values and standard deviations)

- other (e.g. analyses of nominal data)

In the following text, first the approach for binary outcomes is described. The other scales are subsequently based on this approach.

On the basis of the effect measure "relative risk", denominator and numerator are always chosen in such a way that the effect (if present) is realized as a value $< 1$, i.e. the lower the value, the stronger the effect.

**A) Binary outcomes**

To determine the extent of the effect in the case of binary outcomes, the two-sided 95% confidence interval for the relative risk is used; if appropriate, this is calculated by the Institute itself. If several studies are pooled quantitatively, the meta-analytical result for the relative risk is used.

Depending on the quality of the outcome, the confidence interval must lie completely below a certain threshold for the extent to be regarded as minor, considerable or major. It is thus decisive that the upper limit of the confidence interval is smaller than the respective threshold.

The following 3 categories for the quality of the outcome are formed:

- all-cause mortality
- serious (or severe) symptoms (or late complications) and adverse events, as well as health-related quality of life
- non-serious (or non-severe) symptoms (or late complications) and adverse events

The thresholds are specified separately for each category. The more serious the event, the bigger the thresholds (in terms of lying closer to 1). The greater the extent, the smaller the thresholds (in terms of lying further away from 1). For the 3 extent categories (minor, considerable, major), the following Table 4 shows the thresholds to be undercut for each of the 3 categories of quality of the outcomes.

Table 4: Thresholds for determining the extent of an effect

<table>
<tr><td colspan="2" rowspan="2"></td><td colspan="3"><strong>Outcome category</strong></td></tr>
<tr><td>All-cause mortality</td><td>Serious (or severe) symptoms (or late complications) and adverse events, as well as health-related quality of life[a]</td><td>Non-serious (or non-severe) symptoms (or late complications) and adverse events</td></tr>
<tr><td rowspan="3"><strong>Extent category</strong></td><td>Major</td><td>0.85</td><td>0.75<br>and risk $\geq 5\%$[b]</td><td>Not applicable</td></tr>
<tr><td>Considerable</td><td>0.95</td><td>0.90</td><td>0.80</td></tr>
<tr><td>Minor</td><td>1.00</td><td>1.00</td><td>0.90</td></tr>
<tr><td colspan="5">a: Precondition (as for all patient-reported outcomes): use of a validated or established instrument, as well as a validated or established response criterion.<br>b: Risk must be at least 5% for at least 1 of the 2 groups compared.</td></tr>
</table>

The relative risk can generally be calculated in 2 ways, depending on whether the risk refers to events or counter-events (e.g. survival vs. death, response vs. non-response). This is irrelevant for the statement on significance specified in Step 1 of the approach (conventional, non-shifted hypotheses), as in such a case the p-value of a single study is invariant and plays a subordinate role in meta-analysis. However, this does not apply to the distance of the confidence interval limits to the zero effect. To determine the extent of effect for each binary outcome (by means of content criteria under consideration of the type of outcome and underlying disease), it must therefore be decided what type of risk is to be assessed, that of an event or counter-event.

**B) Time to event**

The two-sided 95% confidence interval for the hazard ratio is required to determine the extent of the effect in the case of outcomes representing a "time to event". If several studies are pooled quantitatively, the meta-analytical result for the hazard ratio is used. If the confidence interval for the hazard ratio is not available, it is approximated on the basis of the available information, if possible [512]. The same limits as for the relative risk are set for determining the extent (see Table 4).

If a hazard ratio is neither available nor calculable, or if the available hazard ratio cannot be interpreted meaningfully (e.g. due to relevant violation of the proportional hazard assumption), it should be examined whether a relative risk (referring to a meaningful time point) can be calculated. It should also be examined whether this operationalization is adequate in the case of transient outcomes for which the outcome "time to event" was chosen. If appropriate, the calculation of a relative risk at a time point is also indicated here.

**C) Continuous or quasi-continuous outcomes, in each case with available responder analyses**

Responder analyses are used to determine the extent of added benefit in the case of continuous or quasi-continuous outcomes. For this purpose, a validated or established response criterion or cut-off value is required. On the basis of the responder analyses (2x2 tables) the relative risks are calculated directly from them. According to Table 4 the extent of the effect is then determined.

**D) Other outcomes**

In the case of other outcomes where no responder analyses with inferable relative risks are available either, it should be examined in the individual case whether relative risks can be approximated [101] to set the corresponding thresholds for determining the extent. Otherwise the extent is to be classified as "non-quantifiable".

For the third step of the operationalization of the overall conclusion on the extent of added benefit, when all outcomes are examined together, a strict formalization is not possible, as no sufficient abstraction is currently known for the value judgements to be made in this regard. In its benefit assessment the Institute will compare the conclusions on probability and on the extent of the effects and provide a justified proposal for an overall conclusion.

## 3.4   Non-drug therapeutic interventions

Even if the regulatory preconditions for the market access of drugs and non-drug therapeutic interventions differ, there is nevertheless no reason to apply a principally different standard concerning the certainty of results in the assessment of the benefits and harms of an intervention. For example, the G-BA's Code of Procedure [198] envisages, as far as possible, the preferential consideration of RCTs, independent of the type (drug/non-drug) of the medical intervention to be assessed. For medical devices, this is weakened by the conformity

evaluation in the EN ISO Norm 14155-2 (Section 4.7 [122]), where RCTs are not presented as the design of choice; however, the choice of design must be justified.

Compared with studies on drug interventions, studies on non-drug interventions are often associated with specific challenges and difficulties [362]. For example, the blinding of the staff performing the intervention will often be impossible, and the blinding of patients will either be difficult or also impossible. In addition, it can be assumed that therapists' and patients' preferences for certain treatment options will make the feasibility of studies in these areas particularly problematic. In addition, it may be necessary especially in the assessment of complex interventions to consider the possibility of contamination effects. It may also be necessary to consider the distinction between the effects caused by the procedure or (medical) device to be assessed on the one hand, and those caused by the expertise and skills of those applying the intervention on the other. Moreover, depending on the time of assessment, learning effects need to be taken into account.

In order to give consideration to the aspects outlined above, studies of particularly good quality are required in order to achieve sufficient certainty of results. Paradoxically, the opposite has rather been the case in the past; i.e. sound randomized studies are often lacking, particularly in the area of non-drug interventions (e.g. in surgery [362]). In order to enable any conclusions at all to be drawn on the relevance of a specific non-drug therapeutic intervention, it may therefore also be necessary to consider non-randomized studies in the assessment. Nonetheless, quality standards also apply in these studies, in particular regarding measures taken to ensure structural equality. However, such studies will usually at best be able to provide hints of a(n) (added) benefit or harm of an intervention due to their inherently lower certainty of results. The inclusion of studies with lower evidence levels is consistent with the corresponding regulation in the G-BA's Code of Procedure [198]. However, the specific obligation to provide a justification is emphasized. In this regulation it is noted: "However, in order to protect patients, recognition of a method's medical benefit on the basis of documents with lower evidence levels requires all the more justification the greater the deviation from evidence level 1 (in each case, the medical necessity of the method must also be considered). For this purpose, the method's potential benefit for patients is in particular to be weighed against the risks associated with the demonstration of effectiveness based on studies of lower evidential value" [198]. This means that the non-availability of studies of the highest evidence level alone cannot generally be viewed as sufficient justification for a benefit assessment based on studies with lower evidence levels.

In the assessment of non-drug therapeutic interventions, it may also be necessary to consider the marketability or CE marking (according to the German Medical Devices Act) and the approval status of drugs (according to the German Pharmaceutical Act), insofar as the test interventions or comparator interventions comprise the use of medical devices or drugs (see Section 3.3.1). The corresponding consequences must subsequently be specified in the report plan (see Section 2.1.1).

### 3.5   Diagnostic tests

In general, the evaluation process for diagnostic tests can be categorized into different hierarchy phases or levels, analogously to the evaluation of drugs [190,307]. Phase 4 prospective, controlled diagnostic studies according to Köbberling et al. [307], or Level 5 studies according to Fryback and Thornbury [190] have an (ideally random) allocation of patients to a strategy with or without application of the diagnostic test to be assessed or to a group with or without disclosure of the (diagnostic) test results. These studies can be seen as corresponding to Phase 3 (drug) approval trials ("efficacy trials"). Accordingly, they are allocated to the highest evidence level (see, for example, the G-BA's Code of Procedure [198]). The US Food and Drug Administration also recommends such studies for specific indications in the approval of drugs and biological products developed in connection with diagnostic imaging techniques [182]. Examples show that they can be conducted with comparatively moderate effort [16,525].

The Institute follows this logic and primarily conducts benefit assessments of diagnostic tests on the basis of studies designed as described above that investigate patient-relevant outcomes. The main features of the assessment comply with the explanations presented in Sections 3.1 and 3.4. In this context, patient-relevant outcomes refer to the same benefit categories as in the assessment of therapeutic interventions, namely mortality, morbidity, and health-related quality of life. The impact of diagnostic tests on these outcomes can be achieved by the avoidance of high(er) risk interventions or by the (more) targeted use of interventions. If the diagnostic test itself is associated with a high(er) risk, a lower-risk diagnostic test may have patient-relevant advantages, namely if (in the case of comparable test quality) the conduct of the test itself causes lower mortality and morbidity rates, or fewer restrictions in quality of life.

Studies in which the interaction between the diagnostic information and the therapeutic benefit is investigated also have a high evidence level and are to be given preference in the benefit assessment of diagnostic tests [455].

If such studies are not available or are of insufficient quantity or quality, an assessment of the diagnostic chain can be performed [366]. In this context, the accuracy of the diagnostic test is assessed by means of generally applied test quality criteria determined in studies showing sufficient certainty of results (usually Phase 3 according to Köbberling et al. [307]). It is also reviewed to what extent it is proven that the consequences resulting from the test results are associated with a benefit. In the case of therapeutic consequences (which is mostly assumed), such proof can be inferred from randomized intervention studies (with patient-relevant outcomes) in which a specific (test) result of the diagnostic test to be assessed was defined as an inclusion criterion. Such studies alone may possibly be drawn upon to provide evidence of a benefit of a diagnostic test, even without a specific assessment of diagnostic accuracy [340], if (with sufficient certainty) conclusions can be inferred from them regarding the interaction between the diagnostic information and the (mostly therapeutic) consequences.

In the assessment of the certainty of results of studies on diagnostic accuracy, the Institute primarily follows the QUADAS[10] criteria [546], which, however, may be adapted for the specific project. The STARD[11] criteria [51,52] are applied in order to decide on the inclusion or exclusion of studies not published in full text on a case-by-case basis.

Level 3 and 4 studies according to Fryback and Thornbury [190] investigate the effect of the (diagnostic) test to be assessed on considerations regarding (differential) diagnosis and/or subsequent therapeutic (or other management) decisions, i.e. it is investigated whether the result of a diagnostic test actually leads to any changes in decisions. However, such studies or study concepts have the major disadvantage that they are not sharply defined, and are therefore of rather theoretical nature. A principal (quality) characteristic of these studies is that it was clearly planned to question the physicians involved regarding the probability of the existence of the disease (and their further diagnostic and/or therapeutic approach) *before* the conduct of the diagnostic test to be assessed or the disclosure of results. This is done in order to determine the change in attitude caused by the test result. In contrast, retrospective appraisals and theoretical estimates are susceptible to bias [190,225]. The relevance of such ultimately uncontrolled studies within the framework of the benefit assessment of diagnostic tests must be regarded as largely unclear. Information on management changes alone cannot therefore be drawn upon to provide evidence of a benefit, as long as no information on the patient-relevant consequences of such changes is available.

It will not always be necessary to reinvestigate the whole diagnostic chain regarding modifications of diagnostic tests already available and for which a patient-relevant benefit has been demonstrated or can be postulated with sufficient plausibility. In such cases it can, for example, be sufficient only to verify equivalent or improved intra-test variability. In a comparison between 2 or more diagnostic tests regarding specific test characteristics, studies with the highest certainty of results are those with a random allocation of the sequence of the test performance (conducted independently of one another and preferably blinded) in the same patients, or with random allocation of the test to different patients. These studies are therefore given primary consideration in the Institute's reports.

It is also conceivable that a new diagnostic test is incorporated in an already existing diagnostic strategy; for example, if a new test precedes (triage test) or follows (add-on test) an established test in order to reduce the frequency of application of the established test or new test, respectively [50]. However, against the background of the subsequent therapeutic (or other types of) consequences, it should be considered that through such a combination of tests, the patient populations ensuing from the respective combined test results differ from those ensuing from the individual test results. This difference could in turn influence subsequent therapeutic (or other types of) consequences and their effectiveness. If such an

---

[10] Quality Assessment of Diagnostic Accuracy Studies
[11] Standards for Reporting of Diagnostic Accuracy

influence cannot be excluded with sufficient certainty, comparative studies on diagnostic strategies including and excluding the new test may be required [182,346].

In the assessment of diagnostic tests, it may also be necessary to consider the result of the conformity assessment procedure for CE marking and the approval status of drugs used in diagnostics (see Section 3.3.1). The corresponding consequences must subsequently be specified in the report plan (see Section 2.1.1).

## 3.6   Early diagnosis and screening

Screening programmes are composed of different modules, which can be examined either in part or as a whole [104,479]. The assessment of a screening test generally follows internationally accepted standards and criteria, for example, those of the UK National Screening Committee (UK NSC [521]), the US Preventive Services Task Force (US PSTF [233,410,458]), or the New Zealand National Health Committee (NHC) [384].

According to the criteria outlined above, the Institute primarily assesses the benefit of screening tests by means of prospective comparative intervention studies on the whole screening chain, which include the (ideally random) allocation of participants to a strategy with or without application of the screening test (or to different screening strategies) and which investigate patient-relevant outcomes. In this context, the main features of the assessment comply with the explanations outlined in Sections 3.1 to 3.4.

If such studies are not available or are of insufficient quantity or quality, an assessment of the single components of the screening chain can be performed. In this context, the accuracy of the diagnostic test is assessed by means of generally applied test quality criteria, determined in studies showing sufficient certainty of results (usually Phase 3 according to Köbberling et al. [307]), (see Section 3.5), and it is reviewed to what extent it is proven that the consequences resulting from the test outcomes are associated with a benefit. In the case of therapeutic consequences (which are mostly assumed), proof can be inferred from randomized intervention studies in which an early (earlier) intervention was compared with a late(r) one. The benefit of an early (earlier) vs. a late(r) intervention may also be assessed by means of intervention studies in which the interaction between the earliness of the start of the intervention and the intervention's effect can be investigated. This can be performed either directly within a study or indirectly by comparing studies with different starting points for the intervention, but with otherwise comparable study designs. Here too, the main features of the assessment comply with the explanations outlined in Sections 3.1 to 3.4.

## 3.7   Prevention

Prevention is directed at avoiding, reducing the probability of, or delaying health impairment [537]. Whereas primary prevention comprises all measures employed before the occurrence of detectable biological impairment in order to avoid the triggering of contributory causes, secondary prevention comprises measures to detect clinically asymptomatic early stages of

diseases, as well as their successful early therapy (see also Section 3.6). Primary and secondary prevention measures are characterized by the fact that, in contrast to curative measures, whole population groups are often the focus of the intervention. Tertiary prevention in the narrowest sense describes specific interventions to avoid permanent (especially social) functional deficits occurring after the onset of disease [238]. This is not the focus of this section, but is addressed in the sections on the benefit assessment of drug and non-drug interventions (see Sections 3.3 and 3.4).

The Institute also primarily performs benefit assessments of prevention programmes (other than screening programmes) by means of prospective, comparative intervention studies that have an (ideally random) allocation of participants to a strategy with or without application of the prevention measure, and that investigate patient-relevant outcomes. Alternatively, due to potential "contamination" between the intervention and control group, studies in which clusters were allocated to the study arms may also be eligible [513].

In individual cases, it needs to be assessed to what extent the consideration of other study designs is meaningful [283]. For example, mass-media campaigns are often evaluated within the framework of "interrupted time-series analyses" (e.g. in [529]), and the use of this study design is also advocated for community intervention research [43]. In the quality assessment of these studies, the Institute uses for orientation the criteria developed by the Cochrane Effective Practice and Organisation of Care Review Group [90].

For the benefit on the population level, not only the effectiveness of the programme is decisive, but also the participation rate. In addition, the question is relevant as to which persons are reached by prevention programmes; research indicates that population groups with an increased risk of disease participate less often in such programmes [323]. Special focus is therefore placed on both of these aspects in the Institute's assessments.

## 3.8 Prognosis studies

An essential basis for the assessment of prognosis studies is the precise formulation of a research question, as studies conducted to evaluate prognostic characteristics have different objectives (e.g. evaluation of risk factors, score development or validation). The discrimination between prognosis studies and diagnostic and/or screening studies can be difficult. Depending on the study objective, in the quality assessment of prognosis studies, different assessment principles are applied.

A prognostic characteristic provides information that should not be an end in itself, but should have a consequence that constitutes a verifiable benefit for the patient. In this context, the (general) requirements applying to a prognostic procedure are similar to those applying to a diagnostic test. If a prognostic characteristic is to be applied in the sense of a screening or prevention programme, then the principles formulated in Section 3.6 need to be considered in the assessment.

No generally accepted quality criteria exist for the assessment of prognosis studies [11,236,481]. Simon and Altman [481] describe guidelines for the planning and conduct of prognosis studies in oncology. Laupacis et al. [328] suggest a general framework for the quality assessment of prognosis studies. Hayden et al. [236] developed guidelines for the quality assessment of prognosis studies with regard to potential sources of bias. A good overview of the relevance of prognosis studies and the different approaches to the development, validation and application of prognostic models for clinical practice is provided by a series in the *British Medical Journal* [14,376,377,437]. The development and validation of prognostic models in the event of missing data on predictors are described by Vergouwe et al. [528]. As systematic reviews of prognosis studies are often limited by deficits in the planning, analysis and reporting of these studies, Hemingway et al. [241] have made proposals for improving prognosis research.

In the assessment of prognosis studies the following points, which result from the underlying data source as well as the data analysis applied, should always be considered:

- Clear formulation of a research question and the corresponding planning of the study. This includes sample size planning, which can, for example, be orientated towards the desired precision of the estimate (width of the confidence interval), and requires an estimate of both the prevalence and incidence of the exposure regarding the outcome variable concerned.

- Clear description of the target and sample population (e.g. population-, register- or general practitioner-based) and justification of their selection.

- Clear description of the selection of study participants and the recruitment procedure.

- Homogeneity of the population investigated. If the population is heterogeneous, it needs to be considered that a prognostic statement can be made as constantly as possible across the subgroups causing heterogeneity (e.g. existence of different baseline risks for the outcome variable of interest).

- Clear definition of the outcome variable(s) towards which the prognostic significance should be orientated.

- Clear definition of the prognostic characteristics, including the statistical handling (e.g. dichotomization or assessment of terziles or quartiles, etc., for a quantitative characteristic), and justification of the procedure selected.

- Clear specification and definition of potential confounders and interactions, including the statistical handling.

- Clear description of the development of the statistical model.

- For cohort studies: completeness of follow-up or measures to achieve as complete a follow-up as possible. Estimation of possible selection effects if follow-up is incomplete.

- Clear description of the handling of missing data.

- In the assessment of prognostic scores: distinction between score development and score validation, e.g. score development within a "training sample" and validation in a test sample.

As multifactorial regression models often play a central role in prognosis studies [377], Section 7.3.7 should also be taken into account. Typical study designs for the evaluation of prognostic characteristics in terms of risk factors include cohort studies and possibly also case-control studies [377]. In exceptional cases (e.g. when investigating constant characteristics), cross-sectional studies may also play a role. The basic principles for the assessment of such studies beyond the aspects mentioned above are described in Section 7.1.4.

The literature search for the evaluation of prognostic characteristics (within the framework of a systematic review) is for example more difficult than for therapeutic studies, and no generally accepted optimum search strategy exists (yet). Furthermore, it is assumed that this research field is especially susceptible to publication bias [11,377,481]. The methodological quality of studies (or their publications) on prognostic characteristics is frequently insufficient [241,415], so that the extraction of the required data is difficult or even impossible. Meta-analyses of prognosis studies (not, however, systematic reviews per se) are therefore often inappropriate and their findings should be utilized with reservation [11]. Some important problems with meta-analyses of prognosis studies can be avoided if individual patient data (IPD) are available [11].

Besides using the results of studies investigating single or (mainly) multiple prognostic characteristics, risk charts (also called risk engines) are being increasingly used to assess the individual risk of patients (or clinically healthy persons) of experiencing an adverse event. Multi-factorial estimates for the concurrence of several risk factors are made in these charts (e.g. the Sheffield Table [536] or the Joint British Chart [60]). The basis for these risk charts are mainly formed by multi-factorial regression models, whose results, for easier handling, are presented in tables or points systems [498]. It should be noted that risks derived for such charts are not "personal" estimates for specific individuals, but statistical estimates of the average risks of a population with a specific risk profile for a defined period (e.g. 10 years). The following factors should be considered when assessing such instruments:

- which population the estimated risks apply to

- what type of study the underlying data originate from

- whether the variables included were analysed together in these studies

- whether, and if so, how a multi-factorial statistical analysis was conducted in these underlying studies

- whether these instruments were ever validated in subsequent studies (test samples)

## 4   Clinical practice guidelines and health care analysis

### 4.1   Background

CPGs are systematically developed decision aids for service providers and patients enabling an appropriate approach to specific health problems. Their aim is to improve patient care. Their recommendations are informed by a systematic review of the evidence and an assessment of the benefits and harms of alternative treatment options [177,208]. CPGs can normatively describe standards in all areas of the health care chain, i.e. in diagnosis, treatment, rehabilitation or after-care. These health care standards contain essential information on the quality of care aimed for in a health care system. Determining a health care standard is a key precondition for drawing conclusions on the quality of care in a health care system.

The identification and description of health care standards by means of high-quality CPGs serve as the basis for different scientific analyses, for example, as a starting point for the development and update of DMPs (see Section 4.3). Likewise, by comparing these standards with specific health care structures, processes and outcomes, gaps in health care and potential for improvement can be detected (see Section 4.4). In the following text, this is described as a "health care analysis". Such an analysis enables conclusions on quality and efficiency issues of services provided within the framework of SHI (see §139a SGB V [3] No. 2).

The focus is on providing an overview of the whole picture of a disease. In addition, individual procedures or technologies may be examined, for example as a basis for further assessment in systematic reviews.

The aim is to present current (or to document lacking) health care standards for decision makers and other players in the health care system and, depending on the research question, to compare them with the specific health care situation in order to enable well-founded decisions to improve the quality of care in the health care system.

### 4.2   Identification of health care standards by means of clinical practice guidelines

### 4.2.1   Health care standards in clinical practice guidelines

Medical standard is defined by medical practice that, according to medical and scientific evidence and/or clinical experience, is accepted in the profession [234]. A CPG is a means of establishing a medical standard scientifically and institutionally.

Evidence-based CPGs are normally drawn upon in our department's reports to answer questions on health care standards. Evidence-based CPGs refer to CPGs whose recommendations are based on a systematic literature search, and are linked as a matter of principle to a level of evidence (LoE) and/or grade of recommendation (GoR), as well as to citations of the underlying primary and/or secondary literature (modified according to

AGREE[12] [4]). An evidence-based CPG does not assume that each individual recommendation included is linked to a high LoE. In general, CPGs that were prepared systematically and transparently, and are therefore evidence-based, also include recommendations founded on a weak evidence base [515].

### 4.2.2 Methodological appraisal of clinical practice guidelines

Information retrieval is conducted according to the procedures described in Chapter 6.

On an international level different instruments are used for the methodological appraisal of CPGs [533]. The AGREE instrument [4,352] and its revised version AGREE II [5,65,66] were developed and validated by a network of researchers and health policy makers and are the most widespread tools internationally. The German-language DELB[13] instrument of the Association of the Scientific Medical Professional Societies (AWMF[14]) and the Agency for Quality in Medicine (ÄZQ[15]) is also based on the appraisal tool of the AGREE Collaboration. To simplify any potential future comparison between the results of a CPG appraisal by the Institute and CPG appraisals published in other studies, AGREE is used as a rule in the Institute's methodological appraisal of CPGs. The Institute is actively involved in the further development of the DELB instrument.

When preparing the report plan, the Institute specifies a priori whether, on the grounds of a research question, a methodological appraisal of CPGs should be performed with the AGREE instrument [4]. This tool consists of 23 key items assessed by means of a scale and organized in 6 domains. Each domain covers a separate dimension of CPG quality:

- Domain 1: scope and purpose

- Domain 2: stakeholder involvement

- Domain 3: rigour of development

- Domain 4: clarity and presentation

- Domain 5: applicability

- Domain 6: editorial independence

Each CPG appraisal is performed by 2 reviewers independently of each other.

### A) Standardized domain scores

The domain scores are independent of each other, which is why for each CPG sum scores are calculated separately for the individual domains. As specified in the AGREE instrument, a

---

[12]Appraisal of Guidelines Research and Evaluation in Europe
[13]Deutsches Leitlinien-Bewertungs-(Instrument) (German Instrument for Methodological Guideline Appraisal)
[14]Arbeitsgemeinschaft der Wissenschaftlichen Medizinischen Fachgesellschaften
[15]Ärztliches Zentrum für Qualität in der Medizin

standardization of the calculated domain scores is performed for better comparability between domains. These are presented in tables in the reports.

The instrument does not allow thresholds to be set for the assessment of domains. However, the individual standardized domain scores can be used for the comparison of CPGs.

**B) Overall appraisal of methodological quality of clinical practice guidelines**

In addition to calculation of the standardized domain scores, according to the procedures specified in the report plan, an overall assessment of CPG quality can be performed with the AGREE instrument [3]. As in the proposal by the AGREE Collaboration, 3 categories are distinguished: "strongly recommend", "recommend (with provisos or alterations)", and "would not recommend".

### 4.2.3  Structured processing and evaluation of recommendations

**A) Clinical practice guideline recommendations, levels of evidence and grades of recommendation**

A recommendation is defined as a proposal for action for clinical or health system decisions. The recipients are generally professionals. In principle, in CPGs those statements are identified as recommendations that are formally indicated as such by the CPG authors. In addition, depending on the research question, recommendations not formally indicated may be identified by a linguistic label (e.g. "is recommended, must, should, can, could, be considered", including negations or negative recommendations).

The developers of evidence-based CPGs use different systems to classify the LoE underlying recommendations and grade the strength of recommendations [24,143,221,320,468]. LoE should inform the reader or user of a CPG in brief about the strength (quality and quantity) of the evidence underlying the recommendation. "Evidence" is understood here to be the primary and secondary literature systematically searched for and evaluated by CPG developers. LoE with regard to the (benefit) assessment of medical interventions are generally based on a hierarchy of evidence and study types.

GoR provide the reader or user of a CPG with information on the strength of a CPG recommendation. They go beyond LoE, as they consider not only the strength of the underlying evidence, but in general also include the balancing of medical, social, patient-specific and/or economic benefits and corresponding risks of a recommendation [24,221,468]. They may also refer to the specific health care situation in a health care system.

**B) Structured processing of recommendations**

Within the framework of structured processing, the recommendations from the CPGs included are first listed in tables in their original language and separately for the health care aspects "prevention", "diagnosis", "therapy", and/or "rehabilitation or after-care". In addition, the LoE and/or GoR for a recommendation are listed in the extraction tables, insofar as they have

been awarded by the CPG developers. Depending on the research question, further information may be presented.

As there is to date no internationally consented standardization of grading systems for evidence and recommendations, the LoE and GoR used by the individual CPG developers are generally noted and the corresponding grading systems documented. In order to better compare the systems of different CPG developers, if possible or planned, comparable LoE and GoR from different developers are summarized in comprehensive evidence and recommendation categories.

## C) Evaluation of the recommendations extracted

### *Summarization of the recommendations*

The evaluation of the extraction tables initially consists of a summarization of the content of the extracted recommendations on the individual aspects of care. In this context, if noted in the CPGs, both the LoE and GoR of the corresponding recommendations are presented.

### *Synthesis of key statements*

If specified in the report plan, the information on content can be summarized from the recommendations of different CPGs on the same aspect of health care in a comprehensive "key statement". Key statements are presented in tables together with the information as to which CPG supports the particular statement with what evidence and/or recommendation category.

## D) Identification of gaps or discrepancies in the presentation of health care standards

The structured processing and evaluation of CPG recommendations enable the identification of gaps, deviations or consistencies in the presentation of existing health care standards.

Depending on their objective, CPGs address certain health care areas up to the complete care chain of a disease. If recommendations are lacking in individual CPGs on the addressed areas of the care chain (e.g. on rehabilitation or treatment), a gap exists in the presentation of the health care standard. This may have various causes. The specifics of the evidence base are a crucial factor (e.g. missing, deficient, insufficiently transferable evidence, etc.). Aspects of health care are also of importance, such as the approval and reimbursement status in a health care system or opportunities available in the corresponding context.

Differences in CPG recommendations or in the allocation to LoE/GoR constitute deviations in the presentation of health care standards. Such deviations may, for example, occur if the content of recommendations differs or if nearly identical recommendations or key statements on the same health care aspect are supported by very different LoE/GoR. The detection of deviations may, among other things, indicate an uncertain evidence base or state of consensus for a distinct aspect or the influence of context-specific factors.

## 4.3 Recommendations on disease management programmes

The health care standards identified by means of the procedure described in Section 4.2 can serve as a basis for the preparation of DMPs. In addition, comparison of health care standards with existing DMP recommendations can determine a potential need for revision of the DMP. Recommendations that are consistent in content and were allocated a high GoR in the underlying CPGs are particularly suitable as a basis for the preparation or revision of DMPs. If GoR are lacking, a high LoE is taken.

## 4.4 Health care analysis

### 4.4.1 Background

#### A) Health care

Health care is defined as the medical and psychosocial care of sick people, as well as measures for prevention and health promotion offered by medical and non-medical providers of health care services. Medical care comprises diagnosis, treatment, nursing care, rehabilitation and after-care. The care provided offers all measures within the health care system that are directly or indirectly targeted towards improving or sustaining the health status (mortality, morbidity, and quality of life) of certain individuals or populations [20].

#### B) Health care standard

Medical standard is defined by medical practice that, according to medical and scientific evidence and/or clinical experience, is accepted in the profession [234]. This standard is referred to as the health care standard, which may be specified by laws, regulations and directives, or identified in CPGs (see also Section 4.2). The reference values of quality indicators can also be interpreted as health care standards [176].

#### C) Quality of health care

For the assessment of quality of care the actual health care situation referring to structures, processes and outcomes is compared with the particular health care standard specified through norms, directives and CPGs [235,273]. By comparing the target status with the actual status, conclusions on the current quality of health care become possible; in this connection the current health care situation represents the "actual status" and the current health care standard represents the "target status", whereby the latter describes the goals to be achieved in health care, i.e. "ideal" health care. This analysis/evaluation is conducted for an area of health care defined by the research question of the commission. The precondition for determining the quality of health care is the availability of health care data that were systematically collected and analysed and that a corresponding health care standard can be determined. In this context, "systematic" is understood to be a planned data collection with uniform documentation instructions (coding instructions, e.g. International Classification of Diseases [ICD] version), standardized data collection forms, as well as a complete, and, if possible, comprehensive collection of data (depending on the research question).

### 4.4.2  Content aspects of a health care analysis

The health care analysis comprises the current and systematic description, analysis and assessment of health care aspects of a defined population group with regard to a specific medical or system-related research question (see §139a SGB V [3] Nos. 1 and 2). How detailed the analysis is depends on the type of commission.

The analysis usually examines the German health care situation, potentially supplemented by international comparison. The health care analysis allows the examination of complex interventions referring to both patient-relevant outcomes and outcomes related to the health care system. For the health care analysis, different individual medical as well as population and health system-related data and studies can be compiled in a modular system. In health sciences the term "individual medicine" is used for "classical" medicine involving the patient, in order to make a distinction from the term "population medicine"; the latter term is a component of public health.

The health care analysis can describe and assess different levels and/or several health care aspects. Basically one distinguishes between 2 areas: an epidemiological area and one comprising the social organization of health care. The first area describes the distribution and frequency of diseases in the population. If one examines a health care problem, this area is important for estimation of how many and what type of patients are affected and whether, in the attempt to solve the health care problem, certain subgroups need to be focused on, e.g. elderly people or socially disadvantaged persons. The second area addresses, for example, issues of health care-related structures and processes.

The health care analysis can examine different resources of the health care system (input), structures and processes (throughput), health care services (output), and/or results (outcomes) [411]. In order to assess the quality of health care, the health care situation is compared with a normative standard, the health care standard, insofar as such a standard exists.

### 4.4.3  Aims of a health care analysis

The superordinate aim of a health care analysis is to assess the quality of care.

The following points can be subgoals of the health care analysis:

- examination of the implementation of standards within health care and identification of possible potential for improvement

- investigation of the effects of health care models or measures of quality assurance on the population or on patient groups/population groups

- provision of (background) information for the development of quality indicators or for the prioritization of research questions

- presentation of references to a potential over-, under- or inappropriate provision of health care [443] and, if applicable, formulation of suggestions for improvement in terms of the optimized use of available resources

- identification of a potential need for research (e.g. clinical research, HTA, health care system research)

For feasibility reasons, the focus within the framework of a project is usually on one or a small number of the aims described above with regard to a certain disease.

### 4.4.4   Research question of a health care analysis

The precondition for the systematic description, investigation, and assessment of health care areas is the formulation of a specific research question. The definition of the research question comprises the specification of the following points:

- population (age; gender; disease; if relevant, subgroup or severity of disease)

- the interventions to be investigated (e.g. care of diabetic patients in general practice)

- outcome measures/patient-relevant outcomes (e.g. structural characteristics or health-related quality of life)

- health care setting (e.g. outpatient care, acute inpatient care or cross-sector care)

When formulating the research question it needs to be specified from which perspective (e.g. patients, society, cost carriers, etc.) health care is to be described and assessed, as the focus of the investigation and the selection of outcomes may change depending on the perspective. In this context, specific attention may be paid to the interests of vulnerable groups.

Regional variations (disparities), international comparisons, as well as temporal developments (trends) may also be addressed according to the research question.

### 4.4.5   Potential health care parameters

Different parameters can be used within the framework of a health care analysis. Health care parameters are, for example, epidemiological indices or indicators that help describe various areas of the health care system (see Table 5).

Table 5: Examples of potential health care parameters

| Examples of potential health care parameters | |
| --- | --- |
| **Indicators** | **Health care parameters** |
| Incidence, prevalence, morbidity | Disease burden |
| Case fatality rate | Disease severity |
| Impairments and disabilities according to International Classification of Functioning (ICF), early retirements, mortality | Consequences of disease |
| Number of doctors per 1000 inhabitants, number of service providers per spatial unit, number of hospital beds per 1000 inhabitants etc. | Structure of the health care system (e.g. in Germany) |
| Utilization of services or service provision | Volume of services |
| Quality indicators for inpatient/outpatient sector, e.g. for patient safety, guideline-compliant care of patients | Quality of health care |
| E.g. neonatal and/or maternal mortality, vaccination rates, length of hospital stays | Structures, processes, and outcomes of health care in an international comparison |

Epidemiological indices, for example, prevalence of a disease, can be drawn upon to obtain an overview of the extent of a health care problem. They provide information on the frequency of disease [326]. Disease severity can be estimated by means of the case fatality rate [242]. The consequences of a disease can be assessed by means of data according to the International Classification of Functioning (ICF) and pension fund data (e.g. on invalidity pensions) [117,490]. Health care studies, as well as data from cost carriers or service providers (health insurance funds, associations of SHI physicians, etc.), can identify patients' utilization of health care services. They thus provide information on how often such services are requested, provided or made use of. Quality indicators for the structural, process and outcome quality of inpatient and/or outpatient care may supplement the data pool. They serve quality assurance purposes and may indicate specific health care problems related to structural characteristics, process steps or individual outcomes. In addition, patient safety data from hospital quality reports and registries, as well as clinical and qualitative studies (as far as available), may be incorporated into a health care analysis. For example, they may disclose avoidable adverse events. Evaluation reports on model projects according to §63 SGB V may indicate potential new health care paths. At a system level, further parameters can be used to describe the health care situation and compared with international data. Examples are vaccination rates, disease-specific life expectancy, the number of hospital beds per 1000 inhabitants, and the proportion of expenditure on health care services in relation to the gross domestic product [297,303,542].

Depending on the research question, the above-mentioned parameters (and possibly others) can be combined and thus enable a comprehensive overview of individual health care areas. The health care standards allocated to these areas are identified as described in Section 4.4.8.

### 4.4.6  Procedure for a health care analysis

An example of a procedure for a health care analysis is presented in Figure 7.
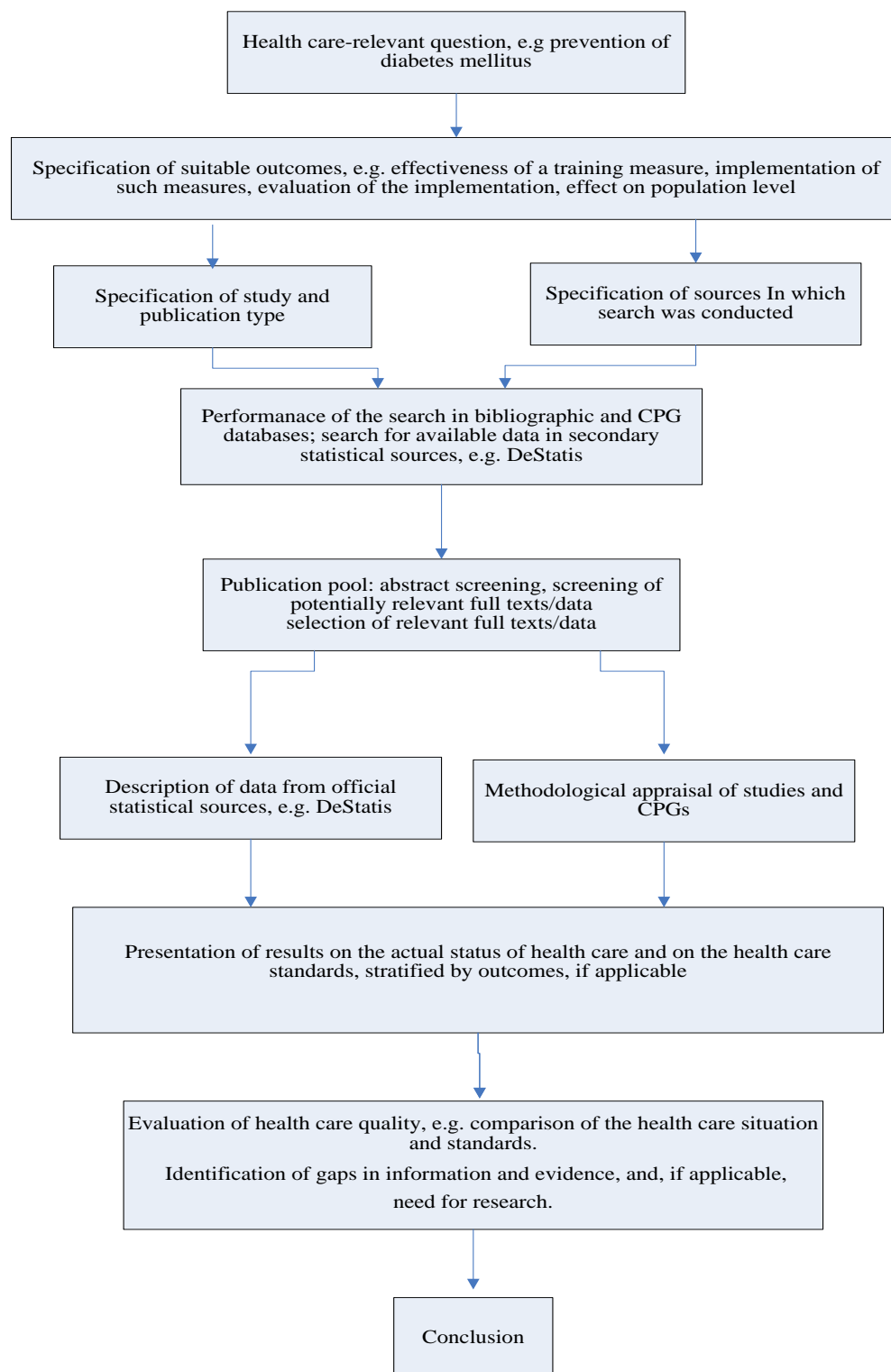
Figure 7: Example of a procedure for a health care analysis

### 4.4.7   Levels of a health care analysis

Health care can be described by means of the above-mentioned parameters relating to 3 different levels: that of individual medicine, of population medicine and of the health care system.

The first level refers to individual patients or patient groups in a clinical setting. Typical outcome measures on this level are patient-relevant outcomes such as mortality, morbidity and health-related quality of life.

The second level refers to population-based studies in the sense of evidence-based health care (population medicine) [245,326]. Outcome measures on this level are also patient-relevant outcomes such as mortality and morbidity, referring, however, to the general population [326]. Beyond this further outcome measures can be investigated, for example, rates of or reasons for participation in screening or vaccination programmes or in health care models such as DMPs.

The third level is the health care system [21,110]. Outcome measures on this level can be the utilization of health care services or the provision of services in different settings (inpatient/outpatient care) or by different professions/providers.

After a description of areas in German health care an international comparison may be meaningful. Depending on the research question, the description of health care in a modular system can refer to information from all 3 levels (individual medicine, population medicine, health care system). In addition, at all levels, temporal developments and regional variations (disparities) can be investigated [115]; for this purpose, geographic information systems can be used, amongst other things.

### 4.4.8   Methodological features of a health care analysis

With regard to the complexity of the health care system and the above-mentioned levels (see Sections 4.4.5 and 4.4.7), different study and publication types may be considered within the framework of a health care analysis.

In addition, it may be necessary to examine different research questions on health care with different quantitative and qualitative methods (pluralism of methods). Moreover, data from several sources are drawn upon (see Section 4.4.9) and processed with different methods. As far as possible, the methodological assessment is performed with suitable instruments (see Section 4.4.9).

In addition, the consideration of sociocultural and ethical aspects may be necessary in the assessment of quality of health care in certain groups of patients, for example, access to health care.

### 4.4.9 Information retrieval

Depending on the research question, different sources may be searched. The search is developed according to the requirements of the source. Both the literature search and the search for CPGs are conducted according to the Institute's *General Methods* (see Section 6.1).

**A) Determination of the health care standard**

The type of health care standard is inferred from the research question for the health care analysis. The first preference is to identify health care standards via evidence-based CPGs. The systematic approach to identify health care standards via CPGs is described in Section 4.2. Laws, regulations and directives define the legally binding framework of health care/medical care.

Structures and processes are mostly assessed by means of quality indicators. High-quality CPGs designate quality indicators, among other things. These refer to measures that indirectly represent the quality of health care. They can be applied to the quality of structures, processes and outcomes. The reference range of the quality indicator specifies the health care goal, i.e. the standard. An indicator always only refers to one health care area, therefore it is meaningful to combine several indicators in order to assess quality [10]. Table 6 provides an overview of potential sources for identifying health care standards.

Table 6: Information sources for identifying German health care standards

| Information on | Examples of data providers |
|---|---|
| Health care/medical standards (CPGs) | Association of the Scientific Medical Professional Societies (AWMF) Guidelines International Network (G-I-N) National Guideline Clearinghouse (NGC) |
| Laws (Social Code Book, SGB) and regulations | Federal Ministry of Justice and Consumer Protection (BMJ) Federal Ministry of Health (BMG) |
| Directives | Federal Joint Committee (G-BA German Medical Association (BÄK) |
| Indicators for the quality of structures, processes and outcomes | National Association of Statutory Health Insurance Physicians (KBV), e.g. Ambulatory Quality Indicators and Key Measures (AQUIK) Federal Office for Quality Assurance (BQS) Institute for Applied Quality Promotion and Research in Health Care (AQUA) |
| AQUA: Institut für angewandte Qualitätsförderung und Forschung im Gesundheitswesen; AQUIK: Ambulante Qualitätsindikatoren und Kennzahlen; AWMF: Arbeitsgemeinschaft der Wissenschaftlichen Medizinischen Fachgesellschaften; BÄK: Bundesärztekammer; BMJ: Bundesministerium der Justiz und für Verbraucherschutz; BMG: Bundesministerium für Gesundheit; BQS: Institut für Qualität und Patientensicherheit (previously Bundesgeschäftsstelle Qualitätssicherung); CPG: clinical practice guideline; G-BA: Gemeinsamer Bundesausschuss; SGB: Sozialgesetzbuch | |

**B) Identification of data sources for health care data**

All relevant data sources for the particular research question should be identified and, as far as possible, used to describe the provision of health care. Following the general principles of topic-related information retrieval (see Section 6.1), selection of data sources is specified in the report plan and is binding (sources: e.g. bibliographic databases, databases of organizations providing official statistics, morbidity registries, handsearch in selected professional journals, contacts with experts, patient organizations, and, if applicable, industry). Potential data sources for identifying health care data are named below (see Table 7).

Table 7: Data sources for identifying health care data

| Information on | Examples of sources |
|---|---|
| Morbidity and mortality, e.g. incidence and prevalence rates (population level) | ▪ Health report of federal and state organizations (e.g. child and youth health survey of the Robert Koch Institute)<br>▪ Report of the Federal Statistical Office (e.g. hospital discharge diagnoses, statistics on causes of death).<br>▪ Morbidity registries (e.g. epidemiological cancer registries)<br>▪ Routine data, e.g. of health care funds and Associations of Statutory Health Insurance Physicians |
| Health care needs (e.g. regional needs analyses) | ▪ Health care studies |
| Utilization and prescription behaviour | ▪ Drug prescription report (Research Institute of the Local Health Care Fund, WidO)<br>▪ Hospital report (WidO)<br>▪ Remedy report (WidO)<br>▪ ICD-10 key codes according to specialty groups (Central Institute for Health Care provided by Statutory Health Insurance Physicians)<br>▪ Routine data, e.g. of health care funds or Associations of Statutory Health Insurance Physicians |
| Patient safety | ▪ Arbitration boards of the Regional Medical Associations<br>▪ Quality indicators of the Organization for Economic Co-operation and Development (OECD)<br>▪ Further publications of the statutory health insurance funds |
| Measurement of health care quality with indicators<br>▪ Quality of health care at a system level<br>▪ Quality of outpatient medical care<br><br>▪ Quality of inpatient care<br><br><br><br><br>▪ Quality of nursing care | <br>▪ OECD (e.g. access to health care)<br>▪ Quality reports of the Associations of Statutory Health Insurance Physicians<br>▪ Hospital quality reports according to §137<br>▪ Publications of the Federal Office for Quality Assurance (BQS)/Institute for Applied Quality Promotion and Research in Health Care (AQUA)<br>▪ Nursing care reports of the Medical Review Board of the Statutory Health Insurance Funds (MDK) |
| ▪ DMP | ▪ Evaluation reports of DMPs |
| Health care system/Comparison of systems | ▪ e.g. WHO publications (e.g. World Health Report) |
| AQUA: Institut für angewandte Qualitätsförderung und Forschung im Gesundheitswesen; BQS: Institut für Qualität und Patientensicherheit (previously Bundesgeschäftsstelle Qualitätssicherung); DMP: disease management programme; ICD: International Classification of Diseases; MDK: Medizinischer Dienst der Krankenversicherung; WidO: Wissenschaftliches Institut der Allgemeinen Ortskrankenkasse; WHO: World Health Organization || |

**C) Search procedure for health care data**

The search procedure follows the usual approach in the Institute. Regarding the search for CPGs and the bibliographic search, this is described in a different section (see Section 6.1). Health care data, e.g. from official statistical sources or morbidity registries, are specifically searched for. The search and search result are documented. For relevant data published exclusively on the Internet, the search strategy to be used is determined by the structure and options of the relevant websites.

Depending upon the specific research question, different data sources/study types are drawn upon to describe health care.

### 4.4.10  Assessment of the data identified

The following aspects need to be considered in the assessment of the data identified:

- Assessment of the study and publication quality of the studies included: Studies are assessed following the Institute's *General Methods*. Supplementations, e.g. regarding evaluation studies or qualitative studies [201], need to be justified.

- Assessment of studies with constructs as endpoints: For patient-relevant outcomes that are constructs, e.g. health-related quality of life, the validity of the survey instrument is assessed. Non-validated instruments are not suitable for comparison.

- Assessment of official statistics: Such data, e.g. from DeStatis, are not assessed methodologically, as it is often impossible to assess such statistics (e.g. on mortality) in this way. In addition, they are already subjected to strict quality criteria by the issuing organizations [68,403]. Publication of these data through third parties, e.g. in journal articles, are assessed according to the Institute's *General Methods*.

- Assessment of registry data: If endpoints are presented by means of registry data, the validity of the registry should be addressed (data quality [i.e. completeness and plausibility of data sets], completeness, currentness) [403,520].

- Conclusive assessment of study and publication quality: The evaluation of the potential risk of bias in the studies/publications to be assessed is conducted following the Institute's *General Methods*.

- Assessment of the methodological quality of the CPGs: see Section 4.2.3.

### 4.4.11 Information synthesis and analysis

The synthesis and analysis of information is conducted as follows: At first the available literature is checked for relevant information on the outcomes specified in the report plan, assessed according to the Institute's *General Methods*, and then described. The results are subsequently summarized. On the basis of the results of the health care analysis an assessment of health care quality is conducted.

**4.4.12 Description and assessment of health care quality**

The assessment of health care quality comprises 3 steps:

**Step 1:** description of the actual status

The actual status of health care is described as specified in the report plan. In this context, the following questions need to be considered:

▪ Are data, studies or publications available on the goals or health care aspects of the investigation?

▪ How reliable are the results found?

**Step 2:** description of the target status

In a further step, health care standards are identified and described. Here too, the availability and the methodological quality of standards are checked.

▪ Does a health care standard exist for the goals/health care aspects stated in the report plan?

▪ How reliable are the results found?

**Step 3:** comparison between actual and target status

Then the actual health care status is compared to the health care standards. Taking the following questions into account, this leads to the assessment of health care quality:

▪ Is the health care standard implemented in everyday health care?

▪ How great are the deviations between the actual and the target status? In which direction does the actual status deviate from the target status (over- or underprovision of health care)?

▪ What conclusions can be drawn from the above comparison?

A final evaluation is made in the conclusions of the report. The evaluation enables us to judge whether gaps in information and/or evidence exist, whether there is a need for research, and/or whether potential for improvement exists.

**4.5   Validity of clinical practice guideline recommendations**

**4.5.1   Background**

Even though the methodology for CPG development is increasingly being further developed [19,220]. CPGs still differ considerably in their development process, presentation, methodological quality, and not least in their content [74,75,80,253,262,359,361,365]. In addition, high methodological quality of CPGs does not necessarily correlate with the content quality of the recommendations included in them [539]. For many research questions the examination of the methodological quality of a CPG is therefore insufficient to assess the value of individual recommendations [220]. This results in the necessity to analyse and

review the contents of CPGs, particularly with regard to the validity of individual recommendations.

### 4.5.2 Definitions of internal and external validity

One distinguishes between the "internal" and "external" validity of CPG recommendations. These are defined as follows:

**Internal validity:** ensures the minimization of potential bias in the development of CPG recommendations.

**External validity:** describes the applicability of a recommendation under the conditions of the health care setting described in the CPG. This can comprise both the clinical setting as well as the use of a CPG at a system level.

The appraisal of the internal validity of CPG recommendations is understood to be the appraisal of the handling of the literature underlying the recommendation, as well as the appraisal of the consensus process. The appraisal of external validity involves consideration of context aspects (e.g. availability, patient preferences, and ethical aspects) or costs in the generation and formulation of the recommendation.

External validity is distinguished from the term "transferability", which describes to what extent a recommendation is transferable to a different context. This can refer to transferability between different health care systems, as well as within a system (e.g. different setting, different target groups of patients).

### 4.5.3 Aim of the analysis and appraisal of clinical practice guideline recommendations

The aim of the methodological approach is to appraise the internal validity of CPG recommendations. Aspects of external validity are only considered if they are helpful in the appraisal of internal validity. For example, the designated context for the CPGs or the user target group may need to be taken into account when assessing the underlying evidence. Closer examination of external validity going beyond this is not conducted.

### 4.5.4 Approach to the analysis and appraisal of internal validity

The appraisal of the internal validity of individual CPG recommendations is conducted by means of

1) Identification and documentation of potentially biasing factors that might impair the internal validity of a CPG recommendation. Such factors are identified and documented at the level of the characteristics and structure of the CPG or its recommendation, the evidence base of the recommendation, and the consensus process.

2) Identification of the need for adjustment of the CPG recommendation: This results from the potential risk of bias identified under Point 1. In this context, depending on the

severity of the determined deficiencies a distinction must be made between a potential and mandatory need for adjustment. A recommendation is classified as "not valid" if a mandatory need for adjustment is identified for this recommendation.

### 4.5.5 Potential research questions

The methods for analysis and appraisal of the internal validity of CPG recommendations are applicable to various research questions. Potential research questions are

- appraisal of individual recommendations from several CPGs on a disease or on a health care aspect comprising several interventions: e.g. "Appraisal of the internal validity of CPG recommendations from evidence-based CPGs on preoperative diagnostics"

- appraisal of recommendations from several CPGs on an intervention: e.g. "Appraisal of the internal validity of CPG recommendations from evidence-based CPGs on the treatment of diabetes mellitus type 2 with long-acting insulin analogues"

- appraisal of recommendations from a specific CPG

Moreover, the analysis and appraisal of the internal validity of CPG recommendations can also contribute to the appraisal of evidence-based CPGs for diseases with the greatest epidemiological relevance (see §139a SGB V [3] No. 3).

## 5  Evidence-based health information for consumers

### 5.1  Background and goals

The Institute has a legislative responsibility for providing health information to consumers, but not direct advice to individuals. The Institute's goal is to improve health and patient autonomy through the provision of health information that aims to advance health and scientific literacy [34,88,127,300,315]. The goals of the Institute's health information are therefore to

- support active and informed decision-making about health issues

- promote the critical use of health care services

- improve understanding of physical, mental and emotional health

- improve understanding of medical and scientific information, including the concept of evidence-based medicine

- enable support of patients by family and friends

To achieve these goals, the Institute needs to be a reliable, trusted and patient-centred information provider. The integration of patient values in decision-making is core to the concept of evidence-based medicine [446], and thus to evidence-based health information. There are several definitions of evidence-based patient information [96,123,155,453,492]. Each definition requires the information to include evidence on benefits, harms and uncertainties of interventions, but the requirements of evidence-based health information go beyond that [123,492]. The Institute defines evidence-based health information as information where

- the content is based on clear scientific evidence, particularly systematic reviews

- the information is developed following systematic methods which aim to minimise bias and maintain neutrality

- evidence-based communication techniques are used to meet the goals of informing, supporting and empowering users

- uncertainties as well as the potential for benefit and harm are discussed.

- language and framing are neutral and non-directive, so that people can make their decisions in accordance with their own values

- the information is updated so that it remains evidence-based

The Institute's primary medium for communication is the Internet, through which it aims to reach consumers, patients and those who advise or provide information to patients (key communicators). Internet-based and offline computer-based health information can positively affect consumers' and patients' knowledge, choices, health and wellbeing [42,293,334,383,387,487,538,554]. However, information interventions can also be

ineffective or harmful [146,148,212,387]. Evidence-based information is also unfamiliar to many people, which brings a set of communication challenges [156,205,493].

## 5.2 Patient-centred communication

### 5.2.1 Communication standards

A key challenge for evidence-based health information is to communicate in a way that is widely understandable while remaining scientifically accurate and objective. The objective and easy-to-use methods of measuring the level of readability in English, French and Spanish have only limited applicability in Germany, and no similarly validated local tool exists [199]. Literacy levels also vary widely among people [300].

To support the differing levels of understanding throughout the community, the Institute produces information at a variety of comprehension levels ranging from fairly simple fact sheets to more complex feature articles, and including multi-media elements (see below). As the Institute produces its information in English as well as German, it can also make use of English readability tools in assessing the level of readability of its information. The Institute aims at a readability below university level. The primary means of quality assessment for understandability, however, will remain reviews of drafts by test readers and reader ratings of understandability (see below).

Explaining evidence and remaining objective in communicating about health-related information raise additional challenges [156,300,467,514]. To be objective and non-directive, the Institute's health information should, on the one hand, not exaggerate what is scientifically known and, on the other, not tell people what they "should" do. This is achieved by not making recommendations and by using neutral language.

There is increasingly more evidence on effective communication with patients and the public [492,514]. Important conclusions for the production of health information can be drawn from this evidence. There are some restrictions, however [145]; for example the fact that research mainly considers cognitive outcomes. Conclusions regarding the effects on mental, behavioural or health-relevant outcomes can often not be deduced in the same way. These studies are often done with groups of people of a certain level of education, which makes it difficult to generalize for the entire population [145]. In addition, the studies often compare very different interventions with one another, which in themselves are very complex, so that it remains unclear what part of the intervention in the end is decisive for the effect on a certain outcome.

Drawing on the evidence that has accumulated on communicating research findings, the Institute aims to

- present information in consistent formats to aid understanding, supplemented with additional elements that could enhance the understandability of medical words and numerical information

- communicate the degree of uncertainty associated with evidence

- considering its transferability, indicate to which groups of people the evidence is applicable

- distinguish very clearly and carefully between "absence of evidence" and "evidence of no effect"

- if there are data that are reliable and relevant for decision-making, give risks as absolute risk and supplement it with other information such as relative risk, if suitable

- avoid biasing information in favour or to the detriment of the products of any particular company, by using generic names of products whenever possible and only using brand names when it is essential for understanding and/or all products on the market in Germany can be named

There is evidence from screenings that the presentation of risk estimates can help patients in their personal orientation when making decisions [148]. If there are reliable data that can help people assess their risk the Institute will present this information. One method to help people individually weigh benefit and harm is offering decision aids. Although uncertainties remain about formal decision aids [397], especially for individual use on the Internet [85,153], the Institute may develop decision aids for some topics, particularly by drawing on the experience of specific decision aids which have been shown to be effective in randomized trials. The Institute develops its decision aids in accordance with the International Patient Decision Aid Standards (IPDAS) [153,254].

Giving people information is not the only purpose of health communication. An important role of health information is to provide users with emotional support [173], and it can also play a role in enhancing patient autonomy. Health communication needs to be patient-centred if it is to be empowering and emotionally supportive. According to the definition of the World Health Organization (WHO), empowerment in health includes the ability to make choices and take actions in accordance with your own goals [396]. These abilities enable consumers to think and act autonomously. Empowering health communication addresses what consumers want to know, shows interest and respect for what patients think and respects their competence [129,299,532].

Historically, patient information tended to be paternalistic [129], assuming ignorance on the part of the patient, together with a need for them to be protected from uncertainty and distressing information: they were to be told what to do. Although it is clear today that this is not helpful [417], it remains, to a significant extent, a feature of patient information and discussions about its role today [129].

Patients in Germany receive less information than they would like [127], and more patients in Germany report that they do not have enough information in comparison with people in similar countries [459,523]. German users of the Internet are possibly even more interested in

hearing about the latest research results than they are in hearing the advice of experts [416]. Nevertheless, the effects of making people more aware of scientific uncertainty are to a large extent unknown. As an argument in favour of information on uncertainties too it has been put forward that coping with uncertainty is an integral part of coping and self-realization for adults [184]. More research is needed on the question of how to measure effect of different health information for the public.

As well as seeking to be understandable, objective and accurate in its information, the Institute aims to

- demonstrate sensitivity and respect for user knowledge, values and concerns, autonomy, cultural differences as well as gender, age and disability-related interests [191]

- maintain a patient-centred, non-judgmental, non-directive and neutral style of language

- respect readers' time

A neutral style of language has to make sure that information reaches both women and men and that both genders feel addressed to in the same way. Continuously referring to people in the masculine form (generic masculine) leads to a mental underrepresentation of women, which is to be seen as linguistic discrimination of women [280]. In the texts of Gesundheitsinformation.de/Informed Health Online, the Institute therefore uses a gender-neutral style of language, which, if possible, does not use generic masculine forms at all. Wherever possible, both genders are explicitly named if both are meant, or gender-neutral expressions are chosen.

On the basis of available evidence and the experience of other groups, the Institute has developed a style guide for its products as well as awareness-raising and other measures for its editorial staff [386]. It will continue to develop its communication standards in the light of the monitoring and evaluation of its products, as well as emerging evidence in the area of evidence-based communication.

### 5.2.2 Method of multidimensional patient pathways

Patient-centred health information is oriented towards users' questions, takes patients' experience and views into account and acknowledges their competency. Patient-centred health information not only aims at answering medical-scientific questions and making an informed decision possible, but also at offering emotional support. To do this, it is necessary on the one hand to know the questions users might be interested in. On the other hand, it is important for the authors to approach the patients' or their family members' perspectives and to develop an understanding of what it means to be living with a certain illness. To do this the Institute employs a method that traces the possible paths patients with a certain illness can follow. This method is called "method of multidimensional pathways" (in German: "Methode der Patientenwege") in the following.

Multidimensional patient pathways summarize and illustrate the different social, emotional, cognitive and clinical dimensions that can be associated with an illness. The method follows the medical-sociological "illness trajectory" model [93] and "patient career" model [200,319] as well as different models of "patient's journey" [325].

Medical sociology started early to look into the effects of illnesses on the patients' lives. In this context the term "patient career" (in German: "Patientenkarriere") was coined in Germany. Some of the contributions to be mentioned here are the ones made by Goffman, Gerhardt and Dörner [200,319]. Another approach is the "illness trajectory" as described by Corbin and Strauss [93].

The Institute derived its method of multidimensional patient pathways from these different approaches. Even though patients become experts in living with a certain illness, and therefore do, in a certain sense, pursue a kind of "career", the Institute prefers the term "multidimensional patient pathways" (in German: "Patientenwege"). This term comprises as many potential courses of patients as possible. The aim is to find out what the different pathways are in a certain illness and what different challenges and decisions the patient will face.

One of the aims of developing multidimensional patient pathways is to set the framework for the contents of the Institute's health information. To do this, the following questions should be answered:

- Who might read the information?

- What content-related questions might the readers have?

- What might be the emotional state of the reader?

- Which information might be used at what point during the course of the disease?

- What decisions are patients faced with, and when will they have to make these decisions?

- What effects might health information on this topic have?

This method mainly aims to help the authors of the Institute's health information systematically develop a good understanding of patients and their relatives as well as of their interaction with information. The orientation towards the dimensions given in Table 8 supports this aim.

Table 8: Different dimensions of a patient pathway

| **(Everyday) life** | Effects of the disease on social relations and roles: family and relationship, job, quality of life, "performance", etc. |
|---|---|
| **Doing/coping** | Anything that is done with regards to the illness, such as visiting a doctor, taking medications, looking for information, self-help, etc. |
| **Feeling** | Feelings that come up during the course of the disease and the treatment, such as grief, fears, worries, etc. |
| **Knowledge** | What do consumers know already? What information might they need? |
| **Decisions** | What decisions must the person affected make in each phase? |
| **Clinic** | Description of the medical phases, such as risk factors, symptoms, diagnosis, treatment, rehabilitation, etc. |
| **Contact point in the health care system** | Who in the health care or social welfare system can be contacted in each phase, for example, doctors, nurses, physiotherapists, psychotherapists, social workers, counselling centres, insurance funds? |

Multidimensional patient pathways can be mapped for the more comprehensive products (fact sheets and feature articles). The effects an illness may have on the life of patients are to be identified systematically. The method should be comprehensible and reproducible, and reflect everyday health care. The sources used include rapid reviews of qualitative studies, personal accounts from the Database of Personal Experiences of Health and Illness (DIPEx) database [128], literature on factors influencing adherence, literature on patients' informational needs systematic reviews on communication and information on health care issues.

(Clinical) care pathways can help identify important diagnostic and therapeutic steps and other "milestones" on a patient's pathway. Care pathways are multidisciplinary care and treatment plans. They describe how basic diagnostic and therapeutic steps in typical patients with a certain diagnosis or illness are optimally coordinated and organized. They aim to connect evidence and practice, and to detect patients' expectations and preferences in order to eventually facilitate optimal care [379,431].

### 5.2.3 Consumer involvement

There is some evidence that getting people affected involved in the development of health information can increase its relevance [393]. One of the requirements of evidence-based health information is its orientation towards the consumers' perspective and their informational needs [123]. This is a key element for the Institute when producing health information [558]. The different measures used include the following: When selecting a topic, topics proposed by website users, experiences from consultations with self-help groups and results of online polls on www.gesundheitsinformation.de/www.informedhealthonline.org are taken into account (see Section 5.3.1). By presenting the individual stories of patients as well

as those close to them, the Institute would like to enable patients and other interested people to find out about the different aspects of living with a condition and nursing care. This is intended to complement the other health information (see Section 5.4.4). As part of the quality assurance, patients or patient representatives are asked to review certain text drafts. Another procedure to include the public is the limited commenting procedure. The Board of Trustees and others are given the opportunity to comment on all feature articles, fact sheets and research summary drafts. The Board of Trustees also includes representatives of patients' interests and representatives of self-help organizations of chronically ill and disabled people. Usually all feature articles, fact sheets and research summaries also undergo external user testing at the same time as they are submitted for comments. In user testing, a group of people affected by the given condition or disease or potential users comment on the texts regarding their content and understandability. In addition, users of the website www.gesundheitsinformation.de/www.informedhealthonline.org have different ways to contact the Institute with their feedback. They can evaluate and comment on the individual information products. There is also an ongoing user survey on the website (see Sections 5.2.4 and 5.5.2).

### 5.2.4   Visual communication and multi-media

Text alone may not be as understandable and memorable as information where explanations are supported by pictures [156,258,312,342,514]. Spoken text may also enhance understanding [258,454]. Explanations where text, pictures and sound are combined may be the most understandable form of communication, especially for people of lower literacy [258]. Where appropriate, the Institute supports its texts with visuals and sound to enhance the effectiveness of its information and reach a wider audience. These include anatomical diagrams and short animated films on key topics that combine visuals, text and sound (videos). Graphics and pictograms also help many people understand numerical data and other aspects of scientific evidence [147,342,467]. Visual and multi-media elements should not replace text, but enhance the material covered in the texts. This ensures that the information is also accessible to people who are visually or hearing impaired.

The Internet enables health information to be presented in multi-media formats. As the technology of the Internet constantly improves and access to the Internet is no longer limited only to computers, communicating effectively with vision and sound on websites is becoming increasingly feasible for more users. The Internet also enables interactivity with users, so that communication need not flow only towards them. Showing an interest in what is important to patients is a critical element in patient-centred and empowering communication [129,299,532]. While the Institute cannot provide individual health advice, there are nevertheless multiple ways in which the Institute offers its users the opportunity to share their views and concerns, including

- online rating of articles
- topic suggestion and general online contact form

- an ongoing survey of the website's usability

- occasional online polls on specific health topics [304]

### 5.2.5  Accessibility

The Internet has both particular advantages and disadvantages in terms of accessibility [32]. For example, its availability 24 hours a day to those with access to the Internet makes it a highly accessible medium [127,462]. Access to the Internet continues to increase: More than half the people in Germany use the Internet for health information, and the number continues to grow [329].

The rate of use of the Internet by people with chronic diseases may be particularly high. Studies in patients attending orthopaedic clinics in Germany, for example, found that up to 70% were using the Internet for information about their condition [407,424]. Over one-third (38%) of patients had visited the Internet on the subject of their consultation before they arrived at the clinic [407]. However, it is not necessary for people to have direct personal access to the Internet for them to benefit from health information on the Internet: Relatives or friends will often search the Internet on their behalf, and key communicators such as doctors, self-help groups and journalists routinely use the Internet. Health information is often shared widely among family members [469]. In the early years of the World Wide Web there were clear gender differences in terms of access to the Internet, and the same is true regarding age and formal education. But these differences appear to be getting smaller [32,266].

For Internet use, there are several accessibility issues, including

- disabilities, particularly (but not only) visual and hearing impairment

- poor reading skills

- insufficient computer skills

- technological capacity (affecting speed and access to multi-media)

- language (the user's mother tongue)

The Institute ensures that its website meets internationally accepted disability accessibility criteria [540], and the German Barrier-Free Information Technology Regulation (BITV[16]) [69]. It will continue to evaluate and optimize the usability of the website, and develop tools that assist with understanding the Institute's information.

Publishing press releases helps health information content reach people who do not use the Internet to look for information about health topics through other media.

---

[16]Barrierefreie Informationstechnik-Verordnung

Close to 10% of people living in Germany have another nationality and close to a further 10% of Germans have a migrant family background [491]. The largest cultural group among these is people from Turkey [491]. People from non-German backgrounds as a group may have greater needs for health information [105]. For this reason it is important that the fact sheets are written in easily-understandable language. Ideally, culturally and gender-appropriate health information would be available in the language people understand best.

The Institute publishes health information in both German and English. The best possible quality assurance requires broad international involvement. Because the Institute's health information is published in English it can profit from feedback from international scientists and particularly from the reviews of the authors of systematic reviews. The availability of an English version broadens the opportunities for translation into further languages.

Translating health information requires high standards, but the quality of these translations is often inadequate [194]. It is difficult to assess the quality of translations according to objective criteria. It is therefore necessary that the people who do the translations and proofread them have linguistic and specialist qualifications. In principle, there are various different ways to translate texts: Translations can be done word for word, for example, or they can capture the intent of the original text into the target language [381]. The translations of the Institute's health information aim to produce a functional informational offer appropriate for the target group. To do this, the contents of the source texts are transferred into the target language. All translations of the Institute's health information are double-checked by a second person with proficiency in both languages.

## 5.3   Topic selection, research and evaluation of evidence

### 5.3.1   Topic selection

The Institute's health information is produced

- in response to commissions received from the G-BA or Ministry of Health

- to summarize other products published by the Institute and as accompanying information for these products

- to fulfil its legislative responsibility to provide consumers with health information, as well as on its own initiative within the framework of the G-BA's general commission

Health information is potentially limitless in scope, and informing everyone about everything is not possible. As with other health care priority-setting decisions, deciding on priorities for health information involves simultaneous analysis of multiple sources of information [28,29].

§139a of the German Social Code Book V (SGB V) sets the following task for the Institute: "provision of easily-understandable general information to citizens on the quality and efficiency of health care services, as well as on the diagnosis and therapy of diseases of substantial epidemiological relevance". The Institute's general commission was amended in

July 2006. According to this it has to "continuously monitor and evaluate fundamentally important developments in medicine" and report on these. This general commission was adapted to the Institute's health information in 2008 [197].

The Institute uses the evidence-scanning system to implement the continuous monitoring and evaluation of medical developments. This is described in more detail in Section 5.3.4. It was not possible to come up with a broadly acceptable definition or a clear list of "diseases of substantial epidemiological relevance". The epidemiological relevance of a disease with practical impact can only be determined using factors for which burden of disease data can be identified. Epidemiologically relevant factors could include

- mortality

- frequency (prevalence/incidence)

- frequency of utilization of health care services

- treatment costs

- absence from work due to illness

- impairment of quality of life and other consequences that have a significant impact on the lives of those affected by the condition

The Institute uses a variety of sources when setting priorities for topics, including causes of death, diagnoses for hospital admission, absence from work due to illness, common diagnoses and medication prescriptions in Germany, as well as the list of diseases chosen within the framework of the morbidity-oriented risk structure compensation scheme in Germany [73].

The legislative responsibility of the Institute to provide consumers with information is consumer-oriented. This includes healthy people as well as patients. For this reason the spectrum of information should also include topics that not only cover the perspective of people directly affected by an illness, but also of those who are not. To meet its goals, the Institute needs to offer information that helps users make choices and take actions to realize their own health goals [396], supports self-management, and addresses what people want to know [129,299,532]. When selecting a topic, the Institute also has to take into account what consumers might be interested in and what they will find in the realities of the health care system.

The Institute uses a number of sources to find out what consumers, healthy or ill, would like to know:

- surveys, primary qualitative research and reviews of qualitative research on people's information needs

- enquiries made to call centres of the SHI funds

- the experiences of other information providers, patient advice services and self-help groups

- enquiries made to the Federal Government Commissioner for Patients' Affairs [463]

- topics that are entered into the search engine of the IQWiG website www.gesundheitsinformation.de/www.informedhealthonline.org as well as other data concerning Internet use

- topics suggested by the website users

- results of the Institute's own online polls about information needs and interests

The Institute also considers reviews of effective information interventions in specific health and illness-related topics to help it determine which health information might be valuable.

The Institute's choice of topics is not solely based on predefined (interest-driven) issues. In its monitoring it gives priority to questions for which evidence-based answers exist. The aspects considered when setting priorities for topics are shown in Table 9. Different sources are used to assess these aspects. After a topic is selected the Institute Management is given an internal product outline of a health information product, which ought to contain information on the following aspects, if possible:

- the product and messages planned

- results of the evaluation of evidence

- statements on epidemiological relevance and burden of disease

- statements on relevance from the McMaster Online Rating of Evidence (MORE) [238]

- statement on relevance on the basis of German guidelines

The methods of topic selection will be further developed in the future, and evaluation of procedure will also play a role. Updating health information offers the opportunity to evaluate the procedure of topic selection. In the process of updating, the number of visitors and their feedback are taken into account (see Section 5.3.1), so that over the years, topics that receive less response can be archived, freeing up capacities for topics of greater user interest. Another important indicator of successful topic selection will be in how far the chosen topic spectrum on the website www.gesundheitsinformation.de/www.informedhealthonline.org covers the most frequent enquires to call centres of SHI funds.

Table 9: Aspects of prioritizing in topic selection

| Evidence | Editorial considerations | Patient/user interest |
|---|---|---|
| Systematic reviews of the benefit of health-related interventions | Balanced range of topics | Patient/user interest |
| Evidence on the effect of information on the topic | Up-to-dateness of the topic | Information looked for by users |
| | Possible adverse effects of the health information | Topic arousing interest in the reader/user |
| | Priorities of contracting agencies | Unfulfilled information needs |
| | Workload and resources | Burden of disease |
| | | Information needs from an expert's point of view |

### 5.3.2 Literature search

The Institute relies predominantly on systematic reviews and qualitative research to develop its health information. When researching a topic in depth, the Institute generally looks for the following information to help identify issues of interest and concern for patients and consumers:

- rapid appraisals of primary qualitative studies as well as reviews of qualitative studies (see Section 7.4)

- reviews of the effects of communication

- reviews of adherence studies

- freely accessible patient information on the Internet as well as self-help group websites

The Internet and other sources are also searched to identify the interventions being used by, or offered to, consumers.

The results of this primary assessment of patients' and information needs help form the Institute's picture of the different stations a person affected with a certain health problem has to go through, the psychological and emotional problems that can occur in relation to that topic, and at what points decisions need to be made. Patient representatives are also generally interviewed to identify further issues and discuss the relevance in Germany of the Institute's findings from research.

The Institute also searches for systematic reviews of causes, prognoses, diagnosis, treatments and adverse effects. This usually covers the entire disease, with a scoping exercise conducted later internally in the department to focus on areas that the health information will cover.

When articles are produced the scoping decisions are discussed with the project group and the prerequisites for information retrieval are taken into consideration (see Chapter 6). The literature search includes, but is not limited to, the Database of Reviews of Effects (DARE) [55,390], the Cochrane Database of Systematic Reviews (CDSR), the HTA Database of the International Network of Agencies for Health Technology Assessment (INAHTA) and MEDLINE. Only reviews with searches undertaken in the last 5 years are considered. Reviews are generally considered to be up-to-date if the search was carried out within the last 3 years [476,477].

The Institute sometimes considers doing a search for trials. For example, if there is no more recent review on an important subject, an update search is considered if there is a high quality review with a search conducted more than 3 years ago [476]. In some instances, updates become relevant if the time span is less than 3 years. This depends on the strength of the evidence in the review and the extent of research activity in that field. An update search for trials, to test how up-to-date a review is, is generally conducted using the Cochrane Controlled Trials Register, MEDLINE and EMBASE [238]. Other databases may be used to complement this search. Searches are also done for qualitative primary literature. Screening of search results is done by 2 people.

### 5.3.3   Evaluation of evidence

The health information produced by the Institute is mainly based on systematic reviews (see Section 7.2). The Institute only uses systematic reviews on the effects of an intervention for their health information if they fulfil certain minimum requirements, which means that they are only allowed to have few methodological flaws according to the Oxman and Guyatt Index [284,400,402]. To be the subject of a research summary suggesting treatment benefit, a review should include at least one trial judged to be of adequate quality by the review's authors, and include data on at least one patient-relevant outcome. The Institute also takes into consideration the relevance and applicability of the evidence to the reader, particularly in terms of gender and age (see Section 7.4).

When more than one systematic review of adequate methodological quality addresses a particular subject or outcome, a further quality assessment is carried out. The aim is to determine whether there are qualitative differences between the different reviews or whether individual reviews are less suitable. The aspects compared include the following:

- main content of the review, especially as regards its relevance for patient information

- extensiveness and up-to-dateness of search

- sensitivity analyses and handling of heterogeneity

- addressing and dealing with any bias potential

- statement about update periods/provision of regular updates (e.g. Cochrane or Agency for Healthcare Research and Quality [AHRQ])

The results of the highest quality review for a particular topic are the source of numerical data used in the Institute's health information. Where reviews come to contradictory conclusions, the possible reasons for this are explored [285].

The methods of the GRADE working group may be used to formally assess the strength of primary evidence in a particular systematic review [23]. The GRADE system explicitly assesses the quality of the evidence and describes how trustworthy the estimates of specific treatment effects are, such as the estimated mortality rate associated with the treatment in question.

For issues concerning aspects like the aetiology or prognosis of a condition, or qualitative descriptions of patients' experiences, other types of primary studies are suitable for inclusion in systematic reviews [204]. When assessing such systematic reviews, the Institute uses the criteria of the Oxford Centre for Evidence-Based Medicine and the McMaster University evidence-rating system [83,238]. The Institute's methods for assessing qualitative research are described in Section 7.4.

### 5.3.4   Updating

A critical part of evidence-based health information is making sure that its conclusions are not out-of-date. Regular updating is one of the quality criteria determined by the EU for health-related websites [91] and which the German position paper *Gute Praxis Gesundheitsinformation* (Good Practice Health Information) [123] describes. Evidence is growing exponentially. This is the case for both trials [33,516] and systematic reviews [33,371]. New evidence can render existing reviews obsolete or out-of-date [189,451,477,534], although new evidence often leads to no change or a strengthening of the original conclusions [276,418,502].

A study of guideline recommendations concluded that after 3 years, over 90% of recommendations may still be current, while after 6 years, about 50% of the recommendations in guidelines may be obsolete [476]. For some topics, for example where the evidence is very strong, the half-life of evidence can be much longer, and in other areas it can be less than 3 years [477]. However, as evidence continues its exponential growth, the half-life of information is likely to shorten: that is, information will become out-of-date more quickly. The Institute sees 3 years as the usual time after which its information requires review. The half-life of the Institute's health information is monitored to inform future updating methods.

Updating can be very resource-intensive [287,372]. It has been estimated that a full update of a guideline, for example, can take almost as long as developing a new guideline [144]. Traditional mechanisms of updating are to schedule a review for a set date. However, this is only sustainable for providers of multiple pieces of evidence-based information in the long term if there is also a continuous increase in resources: An updating workload will continue to grow exponentially over time. The Cochrane Collaboration, for example, has the goal of having the searches for at least 80% of its reviews updated every 2 years [309]. However, this

has not been possible and the reviews are instead probably becoming ever more out-of-date [33]. Using this standard approach of a 2-yearly update, the updating workload for the Institute's health information would already have exceeded the Institute's capacity to both stay up-to-date and keep producing new information.

The Institute uses the following model to keep its information relatively up-to-date:

### A) Preparing regular updates

Before a feature article, a fact sheet or a research summary is published, the most important conclusions from the production process are recorded in an editorial meeting, which are to be considered in a future update. A date is specified when the information is to be updated at the latest. This period will usually be 3 years.

### B) The evidence-scanning system

Evidence scanning means continually identifying all systematic reviews in German or English that might concern information products already published or in the process of being produced. In order to do this, 2 people regularly screen the following sources: CDSR, DARE, INAHTA, MORE, and PubMed. Regulatory authorities in Germany, the United Kingdom, and the United States as well as the European Medicines Agency (EMA) are also continually monitored for the publication of safety alerts. Regularly updated evidence-based information offered to physicians, including Clinical Evidence and EBM Guidelines, are also taken into account.

Each potentially relevant review, study or communication that might concern one or several of the Institute's health information products is considered for the following process and graded according to relevance and editorial interest. This grading can trigger an immediate need to update a health information product, confirm the original scheduling or lead to postponement of the planned update, for example.

### C) Triggering an update

The following factors can serve as triggers of an update:

- letters to the editor or readers' comments
- internal or external criticism
- significant new information from the evidence-scanning system
- arriving at the date scheduled for updating

### D) Planned updating

Six months before the date scheduled for updating, the information products are assessed as to the amount of resources needed for updating, and the updating process is planned. The following aspects are taken into account:

- the summary from the editorial meeting of the original publication

- results from the evidence scanning system

- the authors' planned updates of the systematic reviews the information product is based on (enquiry to the authors)

- changes in methods or editorial standards implemented in the meantime

- user online-rating

- number of visitors

- the information's ranking in Google searches

- necessity of user testing, external review or new commenting procedure

Additional searches for systematic reviews or sometimes studies may also be regarded as necessary. Decisions concerning revision and further searches up to and including the withdrawal of a health information product are usually made in an editorial meeting.

After each update of a health information product, the necessity of making adjustments to the next scheduled date for an update will be considered.

### E) Quality assurance of updating

The extent of individual quality assurance measures depends on the needs for revision of the health information product in question. Updating a product with minor changes due to an updated Cochrane review, for example, will have only internal quality assurance. Generally, the possible measures range from minor corrections in the chief editor's responsibility to a new review. All newsletter subscribers are informed after the updated health information product has been published. Readers' feedback can lead to another update.

### F) Grading of updating

The scope of each update of a feature article, a fact sheet or a research summary is put into 1 of 3 categories. These categories are published on the website history: minor, normal and major (see Table 10).

Table 10: Update categories

| Minor | Normal | Major |
|---|---|---|
| Minor corrections or language improvements of the text that do not affect the message of the health information product. | Revision of a systematic review the product is based on that has led to no or only very few new data. The message of the health information product has only changed slightly or not at all. Exchange of a systematic review the product is based on without substantially changing the message of the health information product. | Update of a systematic review or exchange of a review the product is based on with substantially changing the message of the health information product (e.g. changing the direction of an effect). Withdrawal of a systematic review the product is based on, or of the health information product. |

## 5.4 Information products

### 5.4.1 Feature articles, fact sheets and research summaries

The Institute produces health information for consumers in different formats. This is intended to help meet the needs of audiences who have differing information needs, differing reading levels, and varying time for reading.

The information products of the Institute include

- feature articles: comprehensive articles which form the basis of a set of related products on an important health issue

- fact sheets: short, easily-understandable information

- research summaries: summaries of individual systematic reviews or larger studies including those produced by other departments in the Institute

Together, these products and supplementary items constitute an evidence-based health encyclopaedia. Studies show that the greatest interest most people looking for information have is in information about treatments and what they can actively do themselves [173,298,373,439,478]. In addition, German patients feel particularly under-informed about diagnostic tests and their results [523]. This reflects the fact that evidence-based information should mainly focus on the effects of treatments, diagnostic tests and self-management strategies. For this reason, the majority of the Institute's information products are fact sheets and research summaries.

The feature articles and fact sheets are similar in format to conventional patient information, while the research summaries are more similar to newspaper reports. Feature articles are information products that are usually more than 20 pages in length, directed at people who are

interested in more detailed information on a topic. The level of readability of these extended articles reflects the more interested and motivated readers likely to read them. These readers may feel patronized by overly simple language [129].

The contents of the individual feature articles may vary from topic to topic. For each topic, information taken from the following areas is screened and assessed regarding its relevance for the individual information product:

- Descriptions of the health condition or symptoms, including

  - anatomy

  - physiology

  - different forms of the condition

  - causes of the condition

  - recognizing symptoms

  - natural course of the condition

  - prognosis

  - possible complications

  - recovery/rehabilitation

  - possible recurrence of the condition (relapse)

  - recognizing recurrences

  - risk groups (including family members)

- Preventive and health promotion measures, including

  - diet

  - physical activity

  - screening methods

  - information

- Diagnostic options, including complementary diagnostic procedures

- Treatment options, including

  - medication

  - surgical interventions

  - other non-pharmacological options

- Rehabilitation

- Other health services

- Psychosocial aspects, including the personal experiences of patients who have the condition in question, as well as of other people affected by the condition, such as carers, family members and friends

Fact sheets are written at a reading level that should be more widely understandable, again reflecting the expected use of this information. The fact sheets can be used in patient-doctor consultations and are written for people who want a quick overview of the information.

Research summaries can be thought of as research-based frequently asked questions (FAQs). German Internet users are very interested in information on the latest research, which may interest them even more than the opinions of experts [416]. The research summaries offer the opportunity to make the results of high quality scientific evidence more widely accessible in Germany.

Section 2.1.5 describes how health information is produced. Information on health research should have a similar level of quality assurance as the research report itself [432]. Assuring quality and relevance of patient information could also be done using patient interviews [393]. This is why feature articles and fact sheets are subjected to external review, sometimes also by patients. Quality assurance of all feature articles, fact sheets and research summary also generally includes giving the individual authors of the systematic reviews drawn upon the opportunity to comment on the respective patient information product or products. In a limited commenting procedure the feature articles, fact sheets and research summary drafts are given to the Institute's Board of Trustees, among others. This also ensures that the patient representatives in this body will also have the opportunity to comment on these drafts. In addition, external user testing is done on all feature articles, fact sheets and research summaries and in some cases also supplementary items. The patients who have reported on their condition and recovery are also invited to comment on the patient information drafts that correspond to their respective accounts (see Section 5.4.4).

### 5.4.2 Supplementary items

In addition to the feature articles, fact sheets and research summaries, additional products (supplementary items) are produced. These aim to make the key messages of the health information more understandable and interesting. The inclusion of pictures, sound and animated films may increase the understandability of the website, especially for people with lower literacy levels (see Section 5.2.4). The animated films are the easiest to understand of all of the Institute's health information products.

People may also prefer and trust websites which are more attractive and which include multi-media elements [96,298,478]. Indeed, high quality content can be rejected solely because of poor design [478].

The supplementary items include

- graphics, photos and other images

- short animated films with text and sound

- interactive quizzes

- an online dictionary, which can be switched on or off

- short explanatory texts on subjects such as recognizing the signs of a disease

- texts and interactive tools explaining evidence-based medicine, to improve understanding of research and numbers

- calculators

- online polls and questionnaires

- patient stories

The goals of the supplementary items are to

- promote general understanding of health and medical issues

- help users to understand and weigh up the potential benefits and harms of medical interventions

- support self-management strategies

Interactive supplementary items are tested for usability in-house, which will also be a critical focus of any user evaluation of the website. Accessibility is a particular focus.

### 5.4.3  Press releases

The Institute issues press releases at the same time as they publish selected health information products. Press releases are an important way to establish the website www.gesundheitsinformation.de/www.informedhealthonline.org as a source of reliable patient-centred health information and to reach the Institute's goals described in Section 5.1 [213]. The coverage following the press releases transports the key messages of the health information, raises interest in the Institute's information offer and makes the website www.gesundheitsinformation.de/www.informedhealthonline.org more widely known. Coverage in different media also addresses people who are not reached by the website alone. The contents of the press releases are produced following the same communication standards as the Institute's health information. The Department of Health Information evaluates the resulting coverage in terms of quantity and analysis of content. The factors influencing which health information products are accompanied by press releases include

- which illnesses, diagnostic procedures or treatments affect many people

- what topics are currently being discussed in the media and the general public

- what consumers might want to know

- whether the research results in question are already publicly known

- whether public perception contradicts the latest research results

- what topics are suitable to promote general health and scientific literacy

### 5.4.4 Patient stories

Patients may trust health websites more if they include the experiences of people affected by the respective condition [478].

Many patients would like to hear or read about the experiences of people affected by the same health condition as them [243,503]. Patient stories are commonly used to impart information in both the fields of journalism and patient information. They represent one means of conveying scientific evidence and making it accessible to the general public [205]. The importance of patient stories in medical practice and in health care is increasingly recognized [210,494,556].

Patient stories provide the following functions [503]:

- They offer the opportunity to compare people's own experiences with those of others.

- Reading about the feelings of others might "allow" acceptance of similar emotions.

- They can show people who are affected that they are not alone with their experiences.

By presenting the individual stories of patients as well as those close to them, the Institute would like to enable patients and other interested people to find out about the different aspects of living with a condition and nursing care. This is intended as a complementary source of health information, in addition to the other products. The content of the patient stories should not contradict the evidence-based health information.

One example of patient stories associated with evidence-based health information is DIPEx, an evaluated multimedia website [128], which is available free of charge on the Internet [243,244,556]. The Institute's methods for gathering, editing and publishing patient stories are based on DIPEx's established approach.

The Institute prepares patient stories using the following process:

1) Interview partners are found, most often via self-help organisations.

2) Informed consent is sought regarding the interview procedure and how the story will be used.

3) The interviews are carried out.

4) The interviews are documented and edited, and the interview partners give their informed consent regarding the publication of the final version.

5) The patient story is published on the website with the permission of the interview partner.

Particular importance is placed on extensively briefing the interview partners before the interview, on the fact that they can withdraw their informed consent to publish the story at any time, on preparing the interviews well, on carrying out the interviews based on predefined criteria, as well as on the anonymity of the interviews. If possible, every feature article should be accompanied by at least 2 patient stories.

## 5.4.5 Website

The primary dissemination vehicle for the Institute's health information is the bi-lingual website, www.gesundheitsinformation.de/www.informedhealthonline.org. The Institute aims to maintain high website standards in

- usability and accessibility [264,318,392] (see Section 5.2.5)
- privacy and data protection [269]
- transparency
- search engine visibility [509]
- attractiveness to users
- user interactivity

The Institute aims to achieve usability and user interaction through a variety of means, including

- navigation through a graphic of the human body
- linkage of related topics to each other
- online rating of individual information items
- help and website tour functions

The website also includes a free electronic newsletter, with the choice of biweekly or monthly subscription. The newsletter contains information on what is new on the website, including when information is updated. In addition, the Institute maintains a version of the website for handheld computers (personal digital assistants) and provides RSS feeds to enable individuals to subscribe by RSS. This also allows the contents of the website to be automatically integrated into other websites.

People's statements about what they trust in a website often focus on factors such as credibility and being clearly non-commercial. User behaviour suggests that, in practice, good design and attractiveness also play a large role in user trust of websites [35,97,298].

The Institute's website is certified by the "Health On the Net" (HON) Foundation and fulfils the 8 requirements of the HON Code of Conduct (HONcode) for medical and health-related

websites [240]. Based in Switzerland, this is an accreditation programme for health website standards in content, transparency and privacy. The Institute has chosen HONcode because it is internationally recognized, covers multiple quality dimensions, and because HON regularly reviews its accredited websites to ensure that they continue to meet the HONcode.

## 5.5 Monitoring and evaluation

### 5.5.1 Routine monitoring

The Institute routinely monitors and analyses the use of its health information website, in particular

- website usage, including comparison with similar websites

- user ratings and feedback, including responses to its ongoing online user survey

- the information's position in Google searches and the website's Google PageRank

- the technical performance of the website

- newsletter subscriptions and retention of subscribers

- adoption of the Institute's information by those who advise or provide information to patients (key communicators)

Commonly used metrics for website use, such as number of hits, provide an inflated impression of the use of websites: Numbers in the hundreds of thousands can in fact represent very small numbers of people actually reading the information. Terms such as "hits" are not good indicators of website readership as they measure technical aspects of delivering Internet information rather than actual readers or reach [264,391]. The Institute differentiates between several main categories of website metrics:

- measurement of website traffic (the number of people who have looked for the website or come across it by chance)

- determination of the "source" of visitors (search engines and links from other websites)

- measurement of the number of pages of information viewed

- readership of and interaction with the website, including searches

In order to be able to compare user traffic with that of other websites, the Institute routinely gathers and analyses data on [264,265]

- the number of individual pages opened by users (page impressions or page views)

- the total number of individual website viewing sessions (visits)

Page impressions and visits by Internet robots (crawlers) are excluded, as are the use of the website by the Institute itself and its website development team. Care is taken not to gather

data in forms that can identify the user. The Institute's privacy and data protection policy is described in detail on the website [269].

The traffic of the website is the total number of people who enter the website, not the number of people who actually read items on it [264]. The Institute therefore monitors and analyses more critically indicators of the number of people who are apparently actually reading information, for example

- searching for information

- navigating through articles

- clicking on glossary terms and visiting related information

- viewing animated films or using quizzes

- downloading PDFs

- visiting the website to view new information after receiving the newsletter

- rating the health information

In addition, the Institute estimates the extent to which key communicators have adopted the Institute's health information. This involves some analysis of how many other websites link to www.gesundheitsinformation.de or www.informedhealthonline.org, in particular the websites of the German SHI funds. In addition, the number of printed versions of the Institute's information, particularly those produced by SHI funds, is monitored.

### 5.5.2   Feedback, corrections and improvements

As well as doing prepublication quality assurance, the Institute's website continually encourages postpublication comments for improvements by readers. Although the Institute does not give medical advice, answer individual health-related questions or distribute any information other than that on its website or other information published by the Institute, reader feedback is an important element of the Institute's efforts to offer patient-friendly and useful health information.

Amendments to health information are classified as minor, normal or major. Minor changes only include linguistic improvements, while major changes are content-related. Standard procedures cover the allocation of these categories and the level of quality assurance required for corrections and improvements associated with them. Each article, fact sheet and research summary is accompanied by a document history online, showing the date and level of change that was made, and each version is archived.

### 5.5.3  Evaluation

As well as routine monitoring, the extent to which the website is meeting its goals is assessed by

- user testing

- online user surveys

- consultation with users and patient representatives

- independent evaluation

User testing of individual pieces of information is part of the prepublication process for some interactive products. These tests are undertaken by members of the Department of Health Information. More formal user testing of the website is undertaken by people outside of the Department.

From the launch of the website (2006) the online survey used by the University of Bielefeld to evaluate online pilot projects in accordance with §65b of SGB V was implemented with permission of the developers on the Institute's website [460]. This was done to enable benchmarking between the website www.gesundheitsinformation.de and the more formally evaluated pilot projects (15 websites). Analysis of 2561 completed surveys on the website www.gesundheitsinformation.de showed that on many criteria, such as understandability, www.gesundheitsinformation.de's rating fell in the middle range of the 15 other websites.

Methods for evaluating websites are an ongoing area of methodological work for the Department of Health Information. There are many instruments and guidelines for the evaluation of health information and health websites. Many evaluate process but not content, or content but not process [97,99,313]. There is no instrument that has been shown to be a reliable indicator of the quality of health information or health websites [166,193,286,378]. There are major omissions in commonly used instruments: DISCERN, for example, does not address the quality of content [97,313].

Surveys of patients, including patients in Germany [167,523], indicate that some of the issues suggested as important in evaluating health information may not in fact be important to most patients. Some recommendations common in such instruments may actually reduce the quality of health information. One example is expecting links to other information and to self-help groups: Only a minority of users may rate this as of significance [416], and poor quality in linked information may result in misinformation. Regularly checking links to see if they remain of high quality is a very resource-intensive task that few websites could realistically maintain.

External evaluation and particularly qualitative evaluation by potential users are important for the ongoing development of the Institute's information and website [106]. The Institute

commissions external experts to evaluate the content of individual information products and information packets by potential users [18,252,405].

External user testing is a continuous key element for quality assurance of prepublication health information, which is done generally at the same time as the commenting procedure. In user testing, potential users, including people affected, comment the text drafts regarding their content and understandability so that already at this stage numerous improvements can be included in the final products. If possible, all people interviewed for patient stories should also be included in the evaluation of the information drafts they contributed to.

## 6   Information retrieval

Various types of information form the basis of the Institute's reports (e.g. results from scientific studies, CPGs, registry data and other data collections, documents from regulatory authorities, and dossiers from pharmaceutical companies). This chapter describes the process of a topic-related search for scientific literature.

In the following Section 6.1 the Institute's approach to conducting its own information retrieval is described. The approach to examining information retrieval conducted by others is presented in Section 6.2.

If data are submitted to the Institute that are not allowed to be published, then these data cannot be considered in the Institute's assessments, as this would contradict the principle of transparency.

### 6.1   Information retrieval conducted by the Institute itself

A systematic literature search aims to identify all publications **relevant** to the particular research question (i.e. publications that contribute to a gain in knowledge on the topic). The search for primary literature is normally orientated towards the aim of achieving high sensitivity.

If a benefit assessment is based on systematic reviews, completeness in terms of complete consideration of all available primary studies is not aimed for. If the completeness of the pool of primary studies used in a systematic review is in doubt and the robustness of results is no longer ensured, a benefit assessment is conducted on the basis of primary studies. In this context, robustness is understood to be sufficient certainty that the result would not be considerably changed by the inclusion of additional information or studies.

The following aspects have to be defined a priori in the systematic literature search:

- The inclusion criteria of the report plan or project outline with regard to
  - medical criteria (e.g. target population, intervention)
  - the study design or type of guideline
  - formal characteristics of the publication (e.g. abstract publication, language, etc.)
- The data sources to be included (e.g. bibliographic databases, guideline databases, handsearching in selected scientific journals, contacts with experts/industry/patient organizations, etc.)

Studies and examples on this topic are provided by numerous publications [152,174,255,256,292,353,370,412,434,449,484]. The relevance of the above criteria varies, depending on the different research questions. The type of product to be prepared (e.g. report,

rapid report, working paper) and the resulting timeframe also have an impact on the approach to information retrieval.

### 6.1.1  Search procedure

The search in bibliographic databases, trial registries as well as guideline databases and websites of guideline providers consists of the following steps:

1) if necessary, specification of the research question posed

2) modification of the research question to a searchable research question

3) formulation of a search concept

4) selection of databases

5) identification of search terms

6) formulation of the search strategies

7) quality assurance (in the case of a bibliographic search)

8) performance of the search

9) storage of the search results in text files and import into a reference management software programme (if a standardized export is possible)

10) documentation of the search

Relevant publications identified in the preliminary search are usually drawn upon to identify search terms and formulate the search strategy for bibliographic databases. As a quality assurance step, it is tested whether the search strategy developed in this way identifies known relevant primary publications (test set) with sufficient certainty. The test set is generated by using previous publications by other working groups (systematic reviews on the topic of interest). In addition, a formal internal quality assurance is performed taking the review by Sampson into account [450].

### 6.1.2  Bibliographic databases

**A) Search for primary literature**

The selection of databases for each product is generally based on the focus (i.e. regarding content, methods, and region) of the bibliographic databases. At least 2 large biomedical databases (e.g. MEDLINE and EMBASE) are always selected. For the preparation of health information a search for qualitative studies is additionally conducted in CINAHL and PsycInfo.

**B) Search for systematic reviews**

In the search for systematic reviews, some different sources from those used in the search for primary literature need consideration. As a rule, databases are searched that exclusively or

largely contain systematic reviews. In addition, a selection of biomedical databases also (but not primarily) containing systematic reviews is searched (e.g. MEDLINE and EMBASE).

Depending on the topic investigated, it is decided what databases or other sources (e.g. websites of individual HTA agencies) are also relevant and should be searched. HTAs that are not free of charge are considered in exceptional cases, if it is assumed, for example, that additional relevant information can be retrieved from them, or if no information is otherwise available.

### 6.1.3   Search in trial registries

The systematic search should identify not only published but also unpublished studies. In this context, "unpublished" means that the studies (or individual data sets) have not been published (or only partly) in scientific journals. Study publications are generally identified by means of a search in bibliographic databases such as MEDLINE or EMBASE. Trial registries can be drawn upon in the search for unpublished studies or data [332].

As a rule, the Institute's benefit assessments involve a search in large general trial registries, as well as meta-registries thereof. In particular these include the trial registry ClinicalTrials.gov of the US National Institute of Health (NIH) as well as the WHO's meta-registry "International Clinical Trials Registry Platform Search Portal" (ICTRP). In the benefit assessment of drugs, trial registries of the pharmaceutical industry (individual companies and meta-registries) are additionally screened. Searches in disease-specific trial registries are only performed in exceptional cases. Due to a lack of functions [548] and incomplete data set, the trial registry "clinicaltrialsregister.eu", which is located at EMA, is not a source that can usually be considered in IQWiG's benefit assessments in its current form.

In addition to information on the existence of a study, some registries are also increasingly including study results. This applies, for example, to ClinicalTrials.gov and trial registries of the pharmaceutical industry. Providing the study in question is in principle relevant to the assessment, results from trial registries can be considered in the Institute's reports.

### 6.1.4   Clinical practice guideline databases and providers

If the aim of the search is to identify CPGs, it is primarily conducted in guideline databases (e.g. Guidelines International Network [G-I-N], the AWMF, or the National Guideline Clearinghouse [NGC]), and may be followed by a search on the websites of providers of specialist and multi-disciplinary guidelines. Whether a supplementary search for guidelines is performed in bibliographic databases depends on the type of report to be prepared.

For the search in guideline databases or websites of guideline providers, the search strategy to be applied is targeted towards the structure and options of the particular websites. Only a few websites allow a search with key words, so that generally the complete list of a website's published guidelines is screened. In addition, for the search in guideline databases or websites

of guideline providers, a standardized export is often not possible. For this reason, the search and number of hits are documented in a standardized search protocol. The potentially relevant hits are documented in a literature management programme. Otherwise, the procedure is followed as described in Section 6.1.1.

Within a benefit assessment, guidelines are not categorically excluded as a source of information. However, a systematic search for guidelines is not usually conducted.

### 6.1.5  Requests to manufacturers

Within the framework of the Institute's benefit assessments, the manufacturers of the technologies to be assessed are usually asked to provide previously unpublished information. The aim of this request is to identify all studies and other information relevant to the benefit assessment, independent of their publication status. For drug assessments this request is usually made in 2 steps. In the first, the Institute asks the manufacturer to supply a complete overview of all studies conducted by the manufacturer on the drug to be assessed. If appropriate, the Institute defines the project-specific inclusion criteria for this overview. In the second, the Institute identifies studies relevant to the benefit assessment from the overview, and requests detailed information on these studies. This may refer to a request for unpublished studies, or for supplementary, previously unpublished information on published studies. Previously unpublished information considered in the benefit assessment will also be published in the Institute's reports in order to ensure transparency. The basis for the incorporation of previously unpublished information into the benefit assessment is the conclusion of an agreement on the transfer and publication of study information. This agreement is made between the Institute and the manufacturer involved before the submission of data (see sample contract [270]). It specifies the procedure, the requirements for the documents to be submitted, as well as their confidential and non-confidential components. If the manufacturer concerned does not agree to this contract and therefore does not agree in particular to the complete transfer of all information requested by the Institute, or does not completely transfer the information requested despite conclusion of the agreement, no further requests to the manufacturer will be made. This is to prevent biased results due to selective provision of information.

### 6.1.6  Other data sources for the search

**A) Proceedings of abstracts and selected scientific journals**

Besides bibliographical database searches, it can be useful (depending on the research question) to conduct a handsearch in selected scientific journals and proceedings of abstracts from scientific meetings. This is decided on a case-by-case basis.

**B) Publicly accessible documents from regulatory authorities**

In the case of drug assessments, but also of assessments of specific (non-drug) medicinal products, publicly accessible drug approval databases or correspondence with regulatory authorities are further potential sources of information.

**C) Information from authors of individual publications**

Within the framework of guideline appraisals or benefit assessments it may be meaningful in individual cases to contact authors of publications or guidelines. For example, the requests may refer to specific details on individual guidelines or to unpublished information on journal publications.

**D) Documents transferred by the G-BA or Ministry of Health**

If documents are provided by the contracting agency (G-BA or Ministry of Health), they are regarded as a component of information retrieval. In the subsequent procedure, these documents are handled following the other principles of the search for and assessment of information.

### 6.1.7    Selection of relevant publications

Due to the primarily sensitive approach, the literature search in bibliographic databases results in a large number of citations that are not relevant to the assessment. The selection of relevant publications is made in several steps:

- Exclusion of definitely irrelevant publications (i.e. publications not fulfilling the inclusion or exclusion criteria of the report plan or project outline) through perusal of the titles, and, if available, the abstracts. This step can be divided into 2 in order to distinguish completely irrelevant publications from topic-related ones which, however, do not fulfil the inclusion or exclusion criteria. "Topic-related" refers, for example, to studies investigating the topic of interest but with a different study design or duration from that specified in the report plan or project outline.

- The full texts of the remaining potentially relevant publications are obtained. The decision on the inclusion of the study in the assessment concerned is then made on the basis of these documents.

- Depending on the research question, a supplementary third step is performed for the search in (clinical practice) guideline databases and on websites of guideline providers, where it is examined whether a methodological approach was adopted in the development and formulation of the guideline. This usually refers to the evidence base of the guideline (see Section 4.2). When preparing the report plan the Institute specifies a priori whether on the basis of the research question only evidence-based guidelines are to be considered in the particular report.

All selection steps are performed by 2 persons independently of each other. Discrepancies are resolved by discussion. In the first selection step, if doubts exist as to the relevance of a study, the corresponding full text is obtained and assessed. In this step, completely irrelevant publications may also be distinguished from topic-related ones.

The languages of publication are usually restricted to those of Western Europe. However, other foreign-language publications may also be included if the available information on these

publications indicates that additional and relevant information for answering the research question is to be expected.

In the search for guidelines, the steps for the full-text screening (from the second screening onwards) are performed by 2 persons independently of each other. The quality assurance of the first screening step is conducted with the help of a standardized search protocol.

### 6.1.8  Documentation of the search

All steps in the search in bibliographic databases are documented. This especially includes

- the search strategy for the databases selected

- the search date

- the user interface

- the number of hits

- after perusal of all hits: documentation of the publications judged relevant to the research question posed (citations)

- after perusal of the full texts: documentation of the citations not judged relevant; alternatively, documentation of the topic-related publications that were, however, irrelevant for the report (in each case providing a reason for exclusion)

All other steps in the information retrieval procedure are also documented (e.g. correspondence with authors, queries to manufacturers, etc.).

### 6.1.9  Benefit assessments based on systematic reviews – supplementary search

In most cases a supplementary search for current primary studies is required; this search covers the period between the last date of the search conducted in the systematic review and the date of the search conducted by IQWiG in the preparation of its report. In benefit assessments based on systematic reviews, a supplementary update search for primary literature can only be dispensed with in justified exceptional cases. This applies, for example, if it is sufficiently certain that the result of the assessment would not be considerably changed by the inclusion of additional information or studies (robustness).

In addition, it may be necessary to conduct supplementary searches for primary literature for specific research questions not addressed in the systematic review.

### 6.2  Evaluation of the information retrieval for dossiers

In its dossier assessments, the Institute does not primarily conduct its own information retrieval; instead, the information retrieval presented in the dossiers is evaluated.

A search in bibliographic databases is not always required in the preparation of a dossier. In contrast, a search in publicly accessible trial registries must be conducted as a matter of

principle by the pharmaceutical company; more details are provided in the G-BA's Code of Procedure [198].

On the one hand, the Institute conducts an evaluation of form and content of information retrieval for the dossier assessment: This refers to the search in bibliographic databases and trial registries and is based on the quality assurance procedures described in Section 6.1.1, as well as on the document templates for the preparation of dossiers included in the requirements of the G-BA's Code of Procedure [198].

On the other hand, depending on the results arising from the evaluation of form and content, the Institute subsequently conducts its own search and/or study selection in order to evaluate the completeness of information retrieval. For this purpose, various strategies are available, for example, random checks of the literature citations included in the dossier, the application of specific functions of literature databases (e.g. "related articles" feature in PubMed [452], as well as the conduct of a complete literature search). The result of the evaluation of information retrieval in the dossier and the description of the approach in this regard are part of the dossier assessment.

## 7   Assessment of information

In research the term "bias" means a systematic deviation between research results and the "truth" [444]. For example, this may refer to an erroneously too high (or too low) estimation of a treatment effect.

A main objective in the benefit assessment of medical services is to estimate the actual effect of therapies and interventions as reliably and unbiasedly as possible. In order to minimize bias in the benefit assessment of medical services, different approaches are adopted internationally; these include using scientifically robust methods, ensuring wide participation in the relevant studies, as well as avoiding conflicts of interest [89]. All these methods also form the legal basis of the Institute's work.

### 7.1   Quality assessment of individual studies

### 7.1.1   Criteria for study inclusion

The problem often arises that studies relevant to a benefit assessment do not completely fulfil the inclusion criteria for the patient population and/or the test and comparator intervention defined in the systematic review. In this case the Institute usually proceeds according to the following criteria:

For the inclusion criterion with regard to the study population, it suffices if at least 80% of the patients included in the study fulfil this criterion. Corresponding subgroup analyses are drawn upon if they are available in such studies. Studies in which the inclusion criterion for the study population is fulfilled in less than 80% of the patients included in the study are only included in the analysis if corresponding subgroup analyses are available, or if it has been demonstrated with sufficient plausibility or has been proven that the findings obtained from this study are applicable to the target population of the systematic review (see Section 3.3.1 for applicability).

Studies are also included in which at least 80% of patients fulfil the inclusion criterion regarding the test intervention (intervention group of the study) and at least 80% fulfil the inclusion criterion regarding the comparator intervention (comparator group of the study). If 1 of the 2 criteria is violated in a study, it is excluded from the benefit assessment.

### 7.1.2   Relationship between study type and research question

Only the most relevant study designs that play a role in benefit assessments in medical research (depending on the research question posed) are summarized here.

It is primarily the inclusion of a control group that is called for in the benefit assessment of interventions. In a design with dependent samples without a control group, proof of the effect of an intervention cannot usually be inferred from a pure "before-after" comparison. Exceptions include diseases with a deterministic (or practically deterministic) course (e.g. ketoacidotic diabetic coma; see Section 3.2.2). Randomization and blinding are quality

criteria that increase the evidential value of controlled studies. Parallel group studies [414], cross-over studies [289], and cluster randomized studies [138] are common designs for clinical trials. If interim analyses are planned, the use of appropriate sequential designs must be considered [544].

Case reports or case series often provide initial information on a topic. These are susceptible to all kinds of bias, so that, depending on the research question, only limited reliable evidence can be inferred from this type of study. The prevalence of diseases can be estimated from population-based cross-sectional studies. Other fundamental and classical study types in epidemiology are case-control studies [57] to investigate the association between exposures and the occurrence of rare diseases, as well as cohort studies [58] to investigate the effect of an exposure over time. Cohort studies designed for this purpose are prospective, although retrospective cohort studies are also conducted in which past exposure is recorded (this type of study is frequently found in occupational or pharmacological epidemiology). In principle, prospective designs are preferable to retrospective designs. However, case-control studies, for example, are frequently the only feasible way of obtaining information on associations between exposures and rare diseases. Newer study designs in modern epidemiology contain elements of both case-control and cohort studies and can no longer be clearly classified as retrospective or prospective [294].

Diagnostic and screening studies may have very different aims, so that the assessment depends on the choice of an appropriate design (see Sections 3.5 and 3.6).

### 7.1.3  Ranking of different study types/evidence levels

Different approaches exist within the framework of systematic reviews or guideline development for allocating specific evidence levels to particular study types [223,228]. These levels can be used to create a ranking with regard to the validity of evidence from different study types. However, no evidence assessment system currently exists that is generally accepted and universally applicable to all systematic reviews [295]. Due to the complexity of the appraisal of studies, no conclusive judgement on quality can be inferred from the hierarchy of evidence. In general, the Institute follows the rough hierarchy of study types described below, which is widely accepted and is also largely consistent with the evidence classification of the G-BA [198], and has been incorporated in the regulation on the benefit assessment of drugs according to §35a SGB V [70]. The highest evidence level is allocated to RCTs and systematic reviews of RCTs, at least within the framework of therapeutic studies. In some classifications, individual RCTs are further graded into those of higher or lower quality. In this context, the conflation of the quality of concept and the quality of results has been criticized by some authors [553]. The next levels include non-randomized intervention studies, prospective observational studies, retrospective observational studies, non-experimental studies (case reports and case series) and, at the lowest evidence level, expert opinions not based on scientific rationale. The Institute will adapt this rough grading system

to the particular situation and research question and, if necessary, present it in more detail [228].

### 7.1.4  Aspects of the assessment of the risk of bias

One main aspect of the interpretation of study results is the assessment of the risk of bias (see qualitative uncertainty of results, Section 3.1.4). In this context, the research question, the study type and design, and the conduct of the study play a role, as well as the availability of information. The risk of bias is substantially affected by the study quality; however, its assessment is not equivalent to the quality assessment of a study. For example, individual outcomes may also be considerably biased in a high-quality study. Other studies, however, may provide high certainty of results for specific outcomes in individual cases, despite being of low quality. As a rule, the Institute will therefore estimate the extent of the risk of bias in a problem-orientated manner for all relevant results (both for the study and the specific outcomes).

In principle, a recognized standardized concept should be followed in a study; from planning to conduct, data analysis, and reporting. This includes a study protocol describing all the important methods and procedures. For (randomized) clinical trials, the usual standards are defined by the basic principles of good clinical practice (GCP) [275,310]; for epidemiological studies, they are defined by guidelines and recommendations to ensure good epidemiological practice (GEP) [116]. In this context, a key criterion to avoid bias is whether the study was actually analysed in the way planned. This cannot usually be reliably concluded from the relevant publications. However, a section on sample size planning may at least provide indications in this regard. In addition, a comparison with the study protocol (possibly previously published) or with the corresponding publication on the study design is useful.

The following important documents were developed to improve the quality of publications:

- the CONSORT statement on RCTs [465] and the corresponding explanatory document [368]

- a proposal for an extension of the CONSORT statement for randomized studies on non-drug interventions [54] and the corresponding explanatory document [53]

- the CONSORT statement on cluster-randomized trials [81]

- the CONSORT statement on the documentation of adverse events [277]

- the CONSORT statement on non-inferiority and equivalence studies [413]

- the CONSORT statement on pragmatic studies [559]

- the PRISMA[17] statement on meta-analyses of randomized trials [369] and the corresponding explanatory document [337]

- the TREND[18] statement on non-randomized intervention trials [113]

- the STROBE[19] statement for observational studies in epidemiology [535] and the corresponding explanatory document [527]

- the MOOSE[20] checklist for meta-analysis of observational studies in epidemiology [497]

- the STARD statement on diagnostic studies [51] and the corresponding explanatory document [52]

If a publication fails to conform to these standards, this may be an indicator of an increased risk of bias of the results of the relevant study. Additional key publications on this issue describe fundamental aspects concerning the risk-of-bias assessment [150,222,248].

Key aspects of the Institute's risk-of-bias assessment of the results of RCTs comprise

- adequate concealment, i.e. the unforeseeability and concealment of allocation to groups (e.g. by external randomization in trials that cannot be blinded)

- blinded outcome assessment in trials where blinding of physicians and patients is not possible

- appropriate application of the "intention-to-treat" (ITT) principle

There must be a more cautious interpretation of the results of unblinded trials, or of trials where unblinding (possibly) occurred, compared with the interpretation of blinded studies. Randomization and the choice of appropriate outcome variables are important instruments to prevent bias in studies where a blinding of the intervention was not possible. In studies that cannot be blinded, it is crucial to ensure adequate concealment of the allocation of patients to the groups to be compared. It is also necessary that the outcome variable is independent of the (non-blinded) treating staff or assessed in a blinded manner independent of the treating staff (blinded assessment of outcomes). If a blinded assessment of outcome measures is not possible, a preferably objective outcome should be chosen which can be influenced as little as possible (with regard to its dimension and the stringency of its recording) by the (non-blinded) person assessing it.

In the production of reports standardized assessment forms are generally used to assess the risk of bias of study results. As a rule, for controlled studies on the benefit assessment of interventions the following items across and specific to outcomes are considered in particular:

---

[17] Preferred Reporting Items for Systematic Reviews and Meta-Analyses
[18] Transparent Reporting of Evaluations with Non-randomized Designs
[19] Strengthening the Reporting of Observational Studies in Epidemiology
[20] Meta-analysis of Observational Studies in Epidemiology

Items across outcomes:

- appropriate generation of a randomization sequence (in randomized studies)

- allocation concealment (in randomized studies)

- temporal parallelism of the intervention groups (in non-randomized studies)

- comparability of intervention groups and appropriate consideration of prognostically relevant factors (in non-randomized studies)

- blinding of patients and treating staff/staff responsible for follow-up treatment

- reporting of all relevant outcomes independent of results

Outcome-specific items:

- blinding of outcome assessors

- appropriate implementation of the ITT principle

- reporting of individual outcomes independent of results

On the basis of these aspects, in randomized studies the risk of bias is summarized and classified as "high" or "low". A low risk of bias is present if it can be excluded with great probability that the results are relevantly biased. Relevant bias is understood to be a change in the basic message of the results if the bias were to be corrected.

In the assessment of an outcome, the risk of bias across outcomes is initially classified as "high" or "low". If classified as "high", the risk of bias for the outcome is also usually classified as "high". Apart from that, the outcome-specific items are taken into account.

The classification as "high" of the risk of bias of the result for an outcome does not lead to exclusion from the benefit assessment. This classification rather serves the discussion of heterogeneous study results and affects the certainty of the conclusion.

No summarizing risk-of-bias assessment is usually performed for non-randomized comparative studies, as their results generally carry a high risk of bias due to the lack of randomization.

If a project of the Institute involves the assessment of older studies that do not satisfy current quality standards because they were planned and conducted at a time when these standards did not exist, then the Institute will present the disadvantages and deficiencies of these studies and discuss possible consequences. A different handling of these older studies compared with the handling of newer studies that have similar quality deficits is however only necessary if this is clearly justifiable from the research question posed or other circumstances of the assessment.

The assessment of formal criteria provides essential information on the risk of bias of the results of studies. However, the Institute always conducts a risk-of-bias assessment that goes beyond purely formal aspects in order, for example, to present errors and inconsistencies in publications, and to assess their relevance in the interpretation of results.

### 7.1.5 Interpretation of composite outcomes

A "composite outcome" comprises a group of events defined by the investigators (e.g. myocardial infarctions, strokes, cardiovascular deaths). In this context the individual events in this group often differ in their severity and relevance for patients and physicians (e.g. hospital admissions and cardiovascular deaths). Therefore, when interpreting composite outcomes one needs to be aware of the consequences thereby involved [94,175,188]. The following explanations describe the aspects to be considered in the interpretation of results. However, they specifically do not refer to a (possibly conclusive) assessment of benefit and harm by means of composite outcomes, if, for example, the potential harm from an intervention (e.g. increase in severe bleeding events) is included in an outcome together with the benefit (e.g. decrease in the rate of myocardial infarctions).

A precondition for consideration of a composite outcome is that the individual components of the composite outcome all represent patient-relevant outcomes defined in the report plan. In this context surrogate endpoints can be only included if they are specifically accepted by the Institute as valid (see Section 3.1.2). The results for every individual event included in a composite outcome should also be reported separately. The components should be of similar severity; this does not mean that they must be of identical relevance. For example, the outcome "mortality" can be combined with "myocardial infarction" or "stroke", but not with "silent myocardial infarction" or "hospital admission".

If a composite outcome fulfils the preconditions stated above, then the following aspects need to be considered in the interpretation of conclusions on benefit and harm:

- Does the effect of the intervention on the individual components of the composite outcome usually take the same direction?

- Was a relevant outcome suited to be included in the composite outcome not included, or excluded, without a comprehensible and acceptable justification?

- Was the composite outcome defined a priori or introduced post hoc?

Insofar as the available data and data structures allow, sensitivity analyses may be performed by comparing the exclusion versus the inclusion of individual components.

If the relevant preconditions are fulfilled, individual outcomes may be determined and calculated from a composite outcome within the framework of a benefit assessment.

### 7.1.6   Interpretation of subgroup analyses

In the methodological literature, subgroup analyses are a matter of controversy [22,401]. The interpretation of results of subgroup analyses at a study level is complicated mainly by 3 factors:

- No characteristic of proof: Subgroup analyses are rarely planned a priori and are rarely a component of the study protocol (or its amendments). If subgroup analyses with regard to more or less arbitrary subgroup-forming characteristics are conducted post hoc, the results cannot be regarded as a methodologically correct testing of a hypothesis.

- Multiple testing: If several subgroups are analysed, results in a subgroup may well reach statistical significance, despite actually being random.

- Lack of power: The sample size of a subgroup is often too small to enable the detection of moderate differences (by means of inferential statistics), so that even if effects actually exist, significant results cannot be expected. The situation is different if an adequate power for the subgroup analysis was already considered in the sample size calculation and a correspondingly larger sample size was planned [64].

The results of subgroup analyses should be considered in the assessment, taking the above 3 issues into account and not dominating the result of the primary analysis, even more so if the primary study objective was not achieved. An exception from this rule may apply if social law implications (see below) necessitate such analyses. Moreover, subgroup analyses are not interpretable if the subgroup-forming characteristic was defined after initiation of treatment (after randomization), e.g. in responder analyses. These aspects also play a role in the conduct and interpretation of subgroup analyses within the framework of meta-analyses (see Section 7.3.8).

The statistical demonstration of different effects between various subgroups should be conducted by means of an appropriate homogeneity or interaction test. The finding that a statistically significant effect was observed in one subgroup, but not in another, cannot be interpreted (by means of inferential statistics) as the existence of a subgroup effect.

Analyses of subgroups defined a priori represent the gold standard for subgroup analyses, where stratified randomization by means of subgroups and appropriate statistical methods for data analysis (homogeneity test, interaction test) are applied [98].

Despite the limitations specified above, for some research questions subgroup analyses may represent the best scientific evidence available in the foreseeable future in order to assess effects in subgroups [186], since factors such as ethical considerations may argue against the reproduction of findings of subgroup analyses in a validation study. Rothwell [430] presents an overview of reasons for conducting subgroup analyses. Sun et al. [500] identified criteria to assess the credibility of subgroup analyses.

Possible heterogeneity of an effect in different, clearly distinguishable patient populations is an important reason for conducting subgroup analyses [314,430]. If a priori information is available on a possible effect modifier (e.g. age, pathology), it is in fact essential to investigate possible heterogeneity in advance with regard to the effect in the various patient groups. If such heterogeneity exists, then the estimated total effect across all patients cannot be interpreted meaningfully [314]. It is therefore important that information on a possible heterogeneity of patient groups is considered appropriately in the study design. It may even be necessary to conduct several studies [216]. Within the framework of systematic reviews, the analysis of heterogeneity between individual studies (and therefore, if applicable, subgroup analyses) is a scientific necessity (see Section 7.3.8), but also a necessity from the perspective of social law, as according to §139a (2) SGB V, the Institute is obliged to consider characteristics specific to age, gender, and life circumstances. In addition, according to the official rationale for the SHI Modernization Act (GMG[21]), the Institute is to elaborate in which patient groups a new drug is expected to lead to a relevant improvement in treatment success, with the aim of providing these patients with access to this new drug [118]. A corresponding objective can also be found in §35a SGB V regarding the assessment of the benefit of drugs with new active ingredients [120]. In this assessment, patient groups should be identified in whom these drugs show a therapeutically relevant added benefit. According to social law, a further necessity for subgroup analyses may arise due to the approval status of drugs. On the one hand, this may be the consequence of the decision by regulatory authorities that, after balancing the efficacy and risks of a drug, may determine that it will only be approved for part of the patient population investigated in the approval studies. These considerations may also be based on subgroup analyses conducted post hoc. On the other hand, studies conducted after approval may include patient groups for whom the drug is not approved in Germany; the stronger approvals differ on an international level, the more this applies. In such cases, subgroup analyses reflecting the approval status of a drug may need to be used, independently of whether these analyses were planned a priori or conducted post hoc.

### 7.1.7   Assessment of data consistency

To assess the evidential value of study results, the Institute will review the consistency of data with regard to their plausibility and completeness. Implausible data are not only produced by incorrect reporting of results (typing, formatting, or calculation errors), but also by the insufficient or incorrect description of the methodology, or even by forged or invented data [9]. Inconsistencies may exist within a publication, and also between publications on the same study.

One problem with many publications is the reporting of incomplete information in the methods and results sections. In particular, the reporting of lost-to-follow-up patients, withdrawals, etc., as well as the way these patients were considered in the analyses, are often not transparent.

---

[21]Gesetzliche Krankenversicherung-(GKV-)Modernisierungsgesetz

It is therefore necessary to expose potential inconsistencies in the data. For this purpose, the Institute reviews, for example, calculation steps taken, and compares data presented in text, tables, and graphs. In practice, a common problem in survival-time analyses arises from inconsistencies between the data on lost-to-follow-up patients and those on patients at risk in the survival curve graphs. For certain outcomes (e.g. total mortality), the number of lost-to-follow-up patients can be calculated if the Kaplan-Meier estimates are compared with the patients at risk at a point in time before the minimum follow-up time. Statistical techniques may be useful in exposing forged and invented data [9].

If relevant inconsistencies are found in the reporting of results, the Institute's aim is to clarify these inconsistencies and/or obtain any missing information by contacting authors, for example, or requesting the complete clinical study report and further study documentation. However, it should be considered that firstly, enquiries to authors often remain unanswered, especially concerning older publications, and that secondly, authors' responses may produce further inconsistencies. In the individual case, a weighing-up of the effort involved and the benefit of such enquiries is therefore meaningful and necessary. If inconsistencies cannot be resolved, the potential impact of these inconsistencies on effect sizes (magnitude of bias), uncertainty of results (increase in error probability), and precision (width of the confidence intervals) will be assessed by the Institute. For this purpose, sensitivity analyses may be conducted. If it is possible that inconsistencies may have a relevant impact on the results, this will be stated and the results will be interpreted very cautiously.

## 7.2 Consideration of systematic reviews

Systematic reviews are publications that summarize and assess the results of primary studies in a systematic, reproducible, and transparent way. This also applies to HTA reports, which normally aim to answer a clinical and/or patient-relevant question. HTA reports also often seek to answer additional questions of interest to contracting agencies and health policy decision makers [139,333,408]. There is no need to differentiate between systematic reviews and HTA reports for the purposes of this section. Therefore, the term "systematic review" also includes HTA reports.

### 7.2.1 Classification of systematic reviews

Relying on individual scientific studies can be misleading. Looking at one or only a few studies in isolation from other similar studies on the same question can make treatments appear more or less useful than they actually are [1]. High quality systematic reviews aim to overcome this form of bias by identifying, assessing and summarizing the evidence systematically rather than selectively [139,150,204,408].

Systematic reviews identify, assess and summarize the evidence from one or several study types that can provide the best answer to a specific and clearly formulated question. Systematic and explicit methods are used to identify, select and critically assess the relevant studies for the question of interest. If studies are identified, these data are systematically

extracted and analysed. Systematic reviews are non-experimental studies whose methodology must aim to minimize systematic errors (bias) on every level of the review process [1,150,248].

For systematic reviews of the effects of medical interventions, RCTs provide the most reliable answers. However, for other questions such as aetiology, prognosis or the qualitative description of patients' experiences, the appropriate evidence base for a systematic review will consist of other primary study types [204]. Systematic reviews of diagnostic and screening tests also show some methodological differences compared with reviews of treatment interventions [107].

In the production of the Institute's reports, systematic reviews are primarily used to identify potentially relevant (primary) studies. However, an IQWiG report can be based partially or even solely on systematic reviews (see Section 7.2.2). Health information produced by the Institute for patients and consumers is to a large part based on systematic reviews. This includes systematic reviews of treatments, and reviews addressing other questions such as aetiology, adverse effects and syntheses of qualitative research (see Section 5.3.3).

The minimal prerequisite for a systematic review on the effects of treatments to be used by the Institute is that it has only minimal methodological flaws according to the Oxman and Guyatt index [284,400,402] or the AMSTAR[22] instrument [473-475]. In addition to considering the strength of evidence investigated in systematic reviews, the Institute will also consider the relevance and applicability of the evidence. This includes investigating the question as to whether the results have been consistent among different populations and subgroups as well as in different healthcare contexts. The following factors are usually considered: the population of the participants in the included studies (including gender and baseline disease risk); the healthcare context (including the healthcare settings and the medical service providers); and the applicability and likely acceptance of the intervention in the form in which it was assessed [47,103].

### 7.2.2 Benefit assessment on the basis of systematic reviews

A benefit assessment on the basis of systematic reviews can provide a resource-saving and reliable evidence base for recommendations to the G-BA or the Ministry of Health, provided that specific preconditions have been fulfilled [95,330]. In order to use systematic reviews in a benefit assessment these reviews must be of sufficiently high quality, that is, they must

- show only a minimum risk of bias

- present the evidence base in a complete, transparent, and reproducible manner

---

[22] Assessment of Multiple Systematic Reviews

and thus allow clear conclusions to be drawn [23,400,547]. In addition, it is an essential prerequisite that the searches conducted in the systematic reviews do not contradict the Institute's methodology and that it is possible to transfer the results to the research question of the Institute's report, taking the defined inclusion and exclusion criteria into account.

The methodology applied must provide sufficient certainty that a new benefit assessment based on primary literature would not reach different conclusions from one based on systematic reviews. For example, this is usually not the case if a relevant amount of previously unpublished data is to be expected.

### A) Research questions

In principle, this method is suited for all research questions insofar as the criteria named above have been fulfilled. The following points should be given particular consideration in the development of the research question:

- definition of the population of interest

- definition of the test intervention and comparator intervention of interest

- definition of all relevant outcomes

- if appropriate, specification of the health care setting or region affected (e.g. Germany, Europe)

The research question defined in this way also forms the basis for the specification of the inclusion and exclusion criteria to be applied in the benefit assessment, and subsequently for the specification of the relevance of the content and methods of the publications identified. On the basis of the research question, it is also decided which type of primary study the systematic reviews must be based on. Depending on the research question, it is possible that questions concerning certain parts of a commission are answered by means of systematic reviews, whereas primary studies are considered for other parts.

### B) Minimum number of relevant systematic reviews

All systematic reviews that are of sufficient quality and relevant to the topic are considered. In order to be able to assess the consistency of results, at least 2 high-quality publications (produced independently of each other) should as a rule be available as the foundation of a report based on systematic reviews. If only one high-quality publication is available and can be considered, then it is necessary to justify the conduct of an assessment based only on this one systematic review.

### C) Quality assessment of publications, including minimum requirements

The assessment of the general quality of systematic reviews is performed with Oxman and Guyatt's validated quality index for systematic reviews [399,400,402] or with the AMSTAR Instrument [473-475]. According to Oxman and Guyatt's index, systematic reviews are regarded to be of sufficient quality if they have been awarded at least 5 of 7 possible points in

the overall assessment, which is performed by 2 reviewers independently of one another. No such threshold is defined for the AMSTAR Instrument and therefore should, if appropriate, be defined beforehand. In addition, as a rule, the sponsors of systematic reviews, as well as authors' conflicts of interests, are documented and discussed. Depending on the requirements of the project, the particular index criteria can be supplemented by additional items (e.g. completeness of the search, search for unpublished studies, for example in registries, or additional aspects regarding systematic reviews of diagnostic studies).

## D) Results

For each research question, the results of a benefit assessment based on systematic reviews are summarized in tables, where possible. If inconsistent results on the same outcome are evident in several publications, possible explanations for this heterogeneity are described [285].

If the compilation of systematic reviews on a topic indicates that a new benefit assessment on the basis of primary studies could produce different results, then such an assessment will be performed.

## E) Conclusion/recommendations

Reports based on systematic reviews summarize the results of the underlying systematic reviews and, if necessary, they are supplemented by a summary of up-to-date primary studies (or primary studies on questions not covered by the systematic reviews). Independent conclusions are then drawn from these materials.

The recommendations made on the basis of systematic reviews are not founded on a summary of the recommendations or conclusions of the underlying systematic reviews. In HTA reports, they are often formulated against the background of the specific socio-political and economic setting of a particular health care system, and are therefore rarely transferable to the health care setting in Germany.

### 7.2.3   Consideration of published meta-analyses

Following international EBM standards, the Institute's assessments are normally based on a systematic search for relevant primary studies, which is specific to the research question posed. If it is indicated and possible, results from individual studies identified are summarized and evaluated by means of meta-analyses. However, the Institute usually has access only to aggregated data from primary studies, which are extracted from the corresponding publication or the clinical study report provided. Situations exist where meta-analyses conducted on the basis of IPD from relevant studies have a higher value (see Section 7.3.8). This is especially the case if, in addition to the effect caused solely by the intervention, the evaluation of other factors possibly influencing the intervention effect is also of interest (interaction between intervention effect and covariables). In this context, meta-analyses including IPD generally provide greater certainty of results, i.e. more precise results not affected by ecological bias,

when compared with meta-regressions based on aggregated data [480]. In individual cases, these analyses may lead to more precise conclusions, particularly if heterogeneous results exist that can possibly be ascribed to different patient characteristics. However, one can only assume a higher validity of meta-analyses based on IPD if such analyses are actually targeted towards the research question of the Institute's assessment and also show a high certainty of results. The prerequisite for the assessment of the certainty of results of such analyses is maximum transparency; this refers both to the planning and to the conduct of analyses. Generally valid aspects that are relevant for the conduct of meta-analyses are outlined, for example, in a document published by EMA [158]. In its benefit assessments, the Institute considers published meta-analyses based on IPD if they address (sub)questions in the Institute's reports that cannot be answered with sufficient certainty by meta-analyses based on aggregated data. In addition, high certainty of results for the particular analysis is required.

## 7.3   Specific statistical aspects

### 7.3.1   Description of effects and risks

The description of intervention or exposure effects needs to be clearly linked to an explicit outcome variable. Consideration of an alternative outcome variable also alters the description and size of a possible effect. The choice of an appropriate effect measure depends in principle on the measurement scale of the outcome variable in question. For continuous variables, effects can usually be described using mean values and differences in mean values (if appropriate, after appropriate weighting). For categorical outcome variables, the usual effect and risk measures of 2x2 tables apply [36]. Chapter 9 of the *Cochrane Handbook for Systematic Reviews of Interventions* [109] provides a well-structured summary of the advantages and disadvantages of typical effect measures. Agresti [6,7] describes the specific aspects to be considered for ordinal data.

It is essential to describe the degree of statistical uncertainty for every effect estimate. For this purpose, the calculation of the standard error and the presentation of a confidence interval are methods frequently applied. Whenever possible, the Institute will state appropriate confidence intervals for effect estimates, including information on whether one- or two-sided confidence limits apply, and on the confidence level chosen. In medical research, the two-sided 95% confidence level is typically applied; in some situations, 90% or 99% levels are used. Altman et al. [13] give an overview of the most common calculation methods for confidence intervals.

In order to comply with the confidence level, the application of exact methods for the interval estimation of effects and risks should be considered, depending on the particular data situation (e.g. very small samples) and the research question posed. Agresti [8] provides an up-to-date discussion on exact methods.

### 7.3.2   Evaluation of statistical significance

With the help of statistical significance tests it is possible to test hypotheses formulated a priori with control for type 1 error probability. The convention of speaking of a "statistically

significant result" when the p-value is lower than the significance level of 0.05 (p<0.05) may often be meaningful. Depending on the research question posed and hypothesis formulated, a lower significance level may be required. Conversely, there are situations where a higher significance level is acceptable. The Institute will always explicitly justify such exceptions.

A range of aspects should be considered when interpreting p-values. It must be absolutely clear which research question and data situation the significance level refers to, and how the statistical hypothesis is formulated. In particular, it should be evident whether a one- or two-sided hypothesis applies [45] and whether the hypothesis tested is to be regarded as part of a multiple hypothesis testing problem [517]. Both aspects, whether a one- or two-sided hypothesis is to be formulated, and whether adjustments for multiple testing need to be made, are a matter of repeated controversy in scientific literature [172,305].

Regarding the hypothesis formulation, a two-sided test problem is traditionally assumed. Exceptions include non-inferiority studies. The formulation of a one-sided hypothesis problem is in principle always possible, but requires precise justification. In the case of a one-sided hypothesis formulation, the application of one-sided significance tests and the calculation of one-sided confidence limits are appropriate. For better comparability with two-sided statistical methods, some guidelines for clinical trials require that the typical significance level should be halved from 5% to 2.5% [274]. The Institute generally follows this approach. The Institute furthermore follows the central principle that the hypothesis formulation (one- or two-sided) and the significance level must be specified clearly a priori. In addition, the Institute will justify deviations from the usual specifications (one-sided instead of two-sided hypothesis formulation; significance level unequal to 5%, etc.) or consider the relevant explanations in the primary literature.

If the hypothesis investigated clearly forms part of a multiple hypothesis problem, appropriate adjustment for multiple testing is required if the type I error is to be controlled for the whole multiple hypothesis problem [40]. The problem of multiplicity cannot be solved completely in systematic reviews, but should at least be considered in the interpretation of results [37]. If meaningful and possible, the Institute will apply methods to adjust for multiple testing. In its benefit assessments (see Section 3.1). The Institute attempts to control type I errors separately for the conclusions on every single benefit outcome. A summarizing evaluation is not usually conducted in a quantitative manner, so that formal methods for adjustment for multiple testing cannot be applied here either.

The Institute does not evaluate a statistically non-significant finding as evidence of the absence of an effect (absence or equivalence) [12]. For the demonstration of equivalence, the Institute will apply appropriate methods for equivalence hypotheses.

In principle, Bayesian methods may be regarded as an alternative to statistical significance tests [488,489]. Depending on the research question posed, the Institute will, where necessary, also apply Bayesian methods (e.g. for indirect comparisons, see Section 7.3.9).

### 7.3.3   Evaluation of clinical relevance

The term "clinical relevance" refers to different concepts in the literature. On the one hand, at a group level, it may address the question as to whether a difference between 2 treatment alternatives for a patient-relevant outcome (e.g. serious adverse events) is large enough to recommend the general use of the better alternative. On the other hand, clinical relevance is understood to be the question as to whether a change (e.g. the observed difference of 1 point on a symptom scale) is relevant for individual patients. Insofar as the second concept leads to the inspection of group differences in the sense of a responder definition and corresponding responder analyses, both concepts are relevant for the Institute's assessments.

In general, the evaluation of the clinical relevance of group differences plays a particular role within the framework of systematic reviews and meta-analyses, as they often achieve the power to "statistically detect" the most minor effects [526]. In this context, in principle, the clinical relevance of an effect or risk cannot be derived from a p-value. Statistical significance is a statement of probability, which is not only influenced by the size of a possible effect but also by data variability and sample size. When interpreting the relevance of p-values, particularly the sample size of the underlying study needs to be taken into account [433]. In a small study, a very small p-value can only be expected if the effect is marked, whereas in a large study, highly significant results are not uncommon, even if the effect is extremely small [171,261]. Consequently, the clinical relevance of a study result can by no means be derived from a p-value.

Widely accepted methodological procedures for evaluating the clinical relevance of study results do not yet exist, regardless of which of the above-mentioned concepts are being addressed. For example, only a few guidelines contain information on the definition of relevant or irrelevant differences between groups [324,505]. Methodological manuals on the preparation of systematic reviews also generally provide no guidance or no clear guidance on the evaluation of clinical relevance at a system or individual level (e.g. the *Cochrane Handbook* [248]). However, various approaches exist for evaluating the clinical relevance of study results. For example, the observed difference (effect estimate and the corresponding confidence interval) can be assessed solely on the basis of medical expertise without using predefined thresholds. Alternatively, it can be required as a formal relevance criterion that the confidence interval must lie above a certain "irrelevance threshold" to exclude a clearly irrelevant effect with sufficient certainty. This then corresponds to the application of a statistical test with a shifting of the null hypothesis in order to statistically demonstrate clinically relevant effects [551]. A further proposal plans to evaluate relevance solely on the basis of the effect estimate (compared to a "relevance threshold"), provided that there is a statistically significant difference between the intervention groups [301]. In contrast to the use of a statistical test with a shifting of the null hypothesis, the probability of a type 1 error cannot be controlled thorough the evaluation of relevance by means of the effect estimate. Moreover, this approach may be less efficient. Finally, a further option in the evaluation of relevance is to formulate a relevance criterion individually, e.g. in terms of a responder

definition [302]. In this context there are also approaches in which the response criterion within a study differs between the investigated participants by defining individual therapy goals a priori [426].

In the assessment of patient-relevant outcomes that have been operationalized by using (complex) scales, in addition to evaluating the statistical significance of effects, it is particularly important to evaluate the relevance of the observed effects of the interventions under investigation. This is required because the complexity of the scales often makes a meaningful interpretation of minor differences difficult. It therefore concerns the issue as to whether the observed difference between 2 groups is at all tangible to patients. This evaluation of relevance can be made on the basis of differences in mean values as well as responder analyses [466]. A main problem in the evaluation of relevance is the fact that scale-specific relevance criteria are not defined or that appropriate analyses on the basis of such relevance criteria (e.g. responder analyses) are lacking [375]. Which approach can be chosen in the Institute's assessments depends on the availability of data from the primary studies.

In order to do justice to characteristics specific to scales and therapeutic indications, the Institute as a rule uses the following hierarchy for the evaluation of relevance, the corresponding steps being determined by the presence of different relevance criteria.

1) If a justified irrelevance threshold for the group difference (mean difference) is available or deducible for the corresponding scale, this threshold is used for the evaluation of relevance. If the corresponding confidence interval for the observed effect lies completely above this irrelevance threshold, it is statistically ensured that the effect size does not lie within a range that is certainly irrelevant. The Institute judges this to be sufficient for demonstration of a relevant effect, as in this case the effects observed are normally realized clearly above the irrelevance threshold (and at least close to the relevance threshold). On the one hand, a validated or established irrelevance threshold is suitable for this criterion. On the other hand, an irrelevance threshold can be deduced from a validated, established or otherwise well-justified relevance threshold (e.g. from sample size estimations). One option is to determine the lower limit of the confidence interval as the irrelevance threshold; this threshold arises from a study sufficiently powered for the classical null hypothesis if the estimated effect corresponds exactly to the relevance threshold.

2) If scale-specific justified irrelevance criteria are not available or deducible, responder analyses may be considered. It is required here that a validated or established response criterion was used in these analyses (e.g. in terms of an individual minimally important difference [MID]) [423]. If a statistically significant difference is shown in such an analysis in the proportions of responders between groups, this is seen as demonstrating a relevant effect (unless specific reasons contradict this), as the responder definition already includes a threshold of relevance.

3) If neither scale-specific irrelevance thresholds nor responder analyses are available, a general statistical measure for evaluating relevance is drawn upon in the form of standardized mean differences (SMD expressed as Hedges' g). An irrelevance threshold of 0.2 is then used: If the confidence interval corresponding to the effect estimate lies completely above this irrelevance threshold, it is assumed that the effect size does not lie within a range that is certainly irrelevant. This is to ensure that the effect can be regarded at least as "small" with sufficient certainty [169].

### 7.3.4 Evaluation of subjective outcomes in open-label study designs

Various empirical studies have shown that in non-blinded RCTs investigating subjective outcomes, effects are biased on average in favour of the test intervention. These subjective outcomes include, for example, PROs, as well as outcomes for which the documentation and assessment strongly depend on the treating staff or outcome assessors. Wood et al. provide a summary of these studies [555]. According to this such results show a potential high risk of bias. A generally accepted approach to this problem within the framework of systematic reviews does not exist. In this situation the Institute will normally infer neither proof of benefit nor harm from statistically significant results.

One possibility to take the high risk of bias for subjective outcomes in open-label studies into account is the definition of an adjusted decision threshold. Only if the confidence interval of the group difference of interest shows a certain distance to the zero effect is the intervention effect regarded as so large that it cannot only be explained by bias. The usual procedure for applying an adjusted decision threshold is to test a shifted null hypothesis. This procedure has been applied for decades; among other things, it is required in the testing of equivalence and non-inferiority hypotheses [159]. The prospective determination of a specific threshold value is required in the application of adjusted decision thresholds. If applied, the Institute will justify the selection of a threshold value on a project-specific basis by means of empirical data, as provided, for example, by Wood et al. [555].

### 7.3.5 Demonstration of a difference

Various aspects need to be considered in the empirical demonstration that certain groups differ with regard to a certain characteristic. It should first be noted that the "demonstration" (of a difference) should not be understood as "proof" in a mathematical sense. With the help of empirical study data, statements can only be made by allowing for certain probabilities of error. By applying statistical methods, these probabilities of error can, however, be specifically controlled and minimized in order to "statistically demonstrate" a hypothesis. A typical method for such a statistical demonstration in medical research is the application of significance tests. This level of argumentation should be distinguished from the evaluation of the clinical relevance of a difference. In practice, the combination of both arguments provides an adequate description of a difference based on empirical data.

When applying a significance test to demonstrate a difference, the research question should be specified a priori, and the outcome variable, the effect measure, and the statistical hypothesis formulation should also be specified on the basis of this question. It is necessary to calculate the sample size required before the start of the study, so that the study is large enough for a difference to be detected. In simple situations, in addition to the above information, a statement on the clinically relevant difference should be provided, as well as an estimate of the variability of the outcome measure. For more complex designs or research questions, further details are required (e.g. correlation structure, recruitment scheme, estimate of drop-out numbers, etc.) [46,114].

Finally, the reporting of results should include the following details: the significance level for a statement; a confidence interval for the effect measure chosen (calculated with appropriate methods); descriptive information on further effect measures to explain different aspects of the results; as well as a discussion on the clinical relevance of the results, which should be based on the evaluation of patient-relevant outcomes.

### 7.3.6  Demonstration of equivalence

One of the most common serious errors in the interpretation of medical data is to rate the non-significant result of a traditional significance test as evidence that the null hypothesis is true [12]. To demonstrate "equivalence", methods to test equivalence hypotheses need to be applied [288]. In this context, it is important to understand that demonstrating exact "equivalence" (e.g. that the difference in mean values between 2 groups is exactly zero) is not possible by means of statistical methods. In practice, it is not demonstration of exact equivalence that is required, but rather demonstration of a difference between 2 groups that is "at most irrelevant". To achieve this objective, it must, of course, first be defined what an irrelevant difference is, i.e. an equivalence range must be specified.

To draw meaningful conclusions on equivalence, the research question and the resulting outcome variable, effect measure, and statistical hypothesis formulation need to be specified a priori (similar to the demonstration of a difference). In addition, in equivalence studies the equivalence range must be clearly defined. This range can be two-sided, resulting in an equivalence interval, or one-sided in terms of an "at most irrelevant difference" or "at most irrelevant inferiority". The latter is referred to as a "non-inferiority hypothesis" [100,274,428].

As in superiority studies, it is also necessary to calculate the required sample size in equivalence studies before the start of the study. The appropriate method depends on the precise hypothesis, as well as on the analytical method chosen [427].

Specifically developed methods should be applied to analyse data from equivalence studies. The confidence interval approach is a frequently used technique. If the confidence interval calculated lies completely within the equivalence range defined a priori, then this will be classified as the demonstration of equivalence. To maintain the level of $\alpha = 0.05$, it is

sufficient to calculate a 90% confidence interval [288]. However, following the international approach, the Institute generally uses 95% confidence intervals.

Compared with superiority studies, equivalence studies show specific methodological problems. On the one hand, it is often difficult to provide meaningful definitions of equivalence ranges [324]; on the other hand, the usual study design criteria, such as randomization and blinding, no longer sufficiently protect from bias [470]. Even without knowledge of the treatment group, it is possible, for example, to shift the treatment differences to zero and hence in the direction of the desired alternative hypothesis. Moreover, the ITT principle should be applied carefully, as its inappropriate use may falsely indicate equivalence [288]. For this reason, particular caution is necessary in the evaluation of equivalence studies.

### 7.3.7 Adjustment principles and multi-factorial methods

Primarily in non-randomized studies, multi-factorial methods that enable confounder effects to be compensated play a key role [296]. Studies investigating several interventions are a further important field of application for these methods [360]. In the medical literature, the reporting of results obtained with multi-factorial methods is unfortunately often insufficient [38,380]. To be able to assess the quality of such an analysis, the description of essential aspects of the statistical model formation is necessary [231,435], as well as information on the quality of the model chosen (goodness of fit) [257]. The most relevant information for this purpose is usually

- a clear description and a priori specification of the outcome variables and all potential explanatory variables

- information on the measurement scale and on the coding of all variables

- information on the selection of variables and on any interactions

- information on how the assumptions of the model were verified

- information on the goodness of fit of the model

- inclusion of a table with the most relevant results (parameter estimate, standard error, confidence interval) for all explanatory variables

Depending on the research question posed, this information is of varying relevance. If it concerns a good prediction of the outcome variable within the framework of a prognosis model, a high-quality model is more important than in a comparison of groups, where an adjustment for important confounders must be made.

Inadequate reporting of the results obtained with multi-factorial methods is especially critical if the (inadequately described) statistical modelling leads to a shift of effects to the "desired" range, which is not recognizable with mono-factorial methods. Detailed comments on the

requirements for the use of multi-factorial methods can be found in various reviews and guidelines [26,39,296].

The Institute uses modern methods in its own regression analysis calculations [230]. In this context, results of multi-factorial models that were obtained from a selection process of variables should be interpreted with great caution. When choosing a model, if such selection processes cannot be avoided, a type of backward elimination will be used, as this procedure is preferable to the procedure of forward selection [230,499]. A well-informed and careful preselection of the candidate predictor variable is essential in this regard [111]. If required, modern methods such as the lasso technique will also be applied [511]. For the modelling of continuous covariates, the Institute will, if necessary, draw upon flexible modelling approaches (e.g. regression using fractional polynomials [436,457]) to enable the appropriate description of non-monotonous associations.

### 7.3.8 Meta-analyses

### A) General comments

Terms used in the literature, such as "literature review", "systematic review", "meta-analysis", "pooled analysis", or "research synthesis", are often defined differently and not clearly distinguished [150]. The Institute uses the following terms and definitions:

- A "non-systematic review" is the assessment and reporting of study results on a defined topic, without a sufficiently systematic and reproducible method for identifying relevant research results on this topic. A quantitative summary of data from several studies is referred to as a "pooled analysis". Due to the lack of a systematic approach and the inherent subjective component, reviews and analyses not based on a systematic literature search are extremely prone to bias.

- A "systematic review" is based on a comprehensive, systematic approach and assessment of studies, which is applied to minimize potential sources of bias. A systematic review may, but does not necessarily have to, contain a quantitative summary of study results.

- A "meta-analysis" is a statistical summary of the results of several studies within the framework of a systematic review. In most cases this analysis is based on aggregated study data from publications. An overall effect is calculated from the effect sizes measured in individual studies, taking sample sizes and variances into account.

- More efficient analysis procedures are possible if IPD are available from the studies considered. An "IPD meta-analysis" is the analysis of data on the patient level within the framework of a general statistical model of fixed or random effects, in which the study is considered as an effect and not as an observation unit.

- The Institute sees a "prospective meta-analysis" as a statistical summary (planned a priori) of the results of several prospective studies that were jointly planned. However, if other studies are available on the particular research question, these must also be considered in the analysis in order to preserve the character of a systematic review.

The usual presentation of the results of a meta-analysis is made by means of forest plots in which the effect estimates of individual studies and the overall effect (including confidence intervals) are presented graphically [335]. On the one hand, models with a fixed effect are applied, which provide weighted mean values of the effect sizes (e.g. weighting by inversing the variance). On the other hand, random-effects models are frequently chosen in which an estimate of the variance between individual studies (heterogeneity) is considered. The question as to which model should be applied in which situation has long been a matter of controversy [154,471,531]. If information is available that the effects of the individual studies are homogeneous, a meta-analysis assuming a fixed effect is sufficient. However, such information will often not be available, so that in order to evaluate studies in their totality, an assumption of random effects is useful [472]. Moreover, it should be noted that the confidence intervals calculated from a fixed-effect model may show a substantially lower coverage probability with regard to the expected overall effect, even if minor heterogeneity exists when compared with confidence intervals from a random-effects model [61]. The Institute therefore primarily uses random-effects models and only switches to models with a fixed effect in well-founded exceptional cases. In this context, if the data situation is homogeneous, it should be noted that meta-analytical results from models with random and fixed effects at best show marginal differences. As described in the following text, the Institute will only perform a meta-analytical summary of strongly heterogeneous study results if the reasons for this heterogeneity are plausible and still justify such a summary.

## B) Heterogeneity

Before a meta-analysis is conducted, it must first be considered whether the pooling of the studies investigated is in fact meaningful, as the studies must be comparable with regard to the research question posed. In addition, even in the case of comparability, the studies to be summarized will often show heterogeneous effects [250]. In this situation it is necessary to assess the heterogeneity of study results [203]. The existence of heterogeneity can be statistically tested; however, these tests usually show very low power. Consequently, it is recommended that a significance level between 0.1 and 0.2 is chosen for these tests [282,308]. However, it is also important to quantify the extent of heterogeneity. For this purpose, specific new statistical methods are available, such as the $I^2$ measure [249]. Studies exist for this measure that allow a rough classification of heterogeneity, for example, into the categories "might not be important" (0 to 40%), "moderate" (30 to 60%), "substantial" (50 to 90%) and "considerable" (75 to 100%) heterogeneity [109]. If the heterogeneity of the studies is too large, the statistical pooling of the study results may not be meaningful [109]. The specification as to when heterogeneity is "too large" depends on the context. A pooling of data is usually dispensed with if the heterogeneity test yields a p-value of less than 0.2. In this context, the location of the effects also plays a role. If the individual studies show a clear effect in the same direction, then pooling heterogeneous results by means of a random effects model can also lead to a conclusion on the benefit of an intervention. However, in this situation a positive conclusion on the benefit of an intervention may possibly be drawn without the quantitative pooling of data (see Section 3.1.4). In the other situations the Institute

will not conduct a meta-analysis. However, not only statistical measures, but also reasons of content should be considered when making such a decision, which must be presented in a comprehensible way. In this context, the choice of the effect measure also plays a role. The choice of a certain measure may lead to great study heterogeneity, yet another measure may not. For binary data, relative effect measures are frequently more stable than absolute ones, as they do not depend so heavily on the baseline risk [192]. In such cases, the data analysis should be conducted with a relative effect measure, but for the descriptive presentation of data, absolute measures for the specific baseline risks may possibly be inferred from relative ones.

In the case of great heterogeneity of the studies, it is necessary to investigate potential causes. Factors that could explain the heterogeneity of effect sizes may possibly be detected by means of meta-regression [506,524]. In a meta-regression, the statistical association between the effect sizes of individual studies and the study characteristics is investigated, so that study characteristics can possibly be identified that explain the different effect sizes, i.e. the heterogeneity. However, when interpreting results, it is important that the limitations of such analyses are taken into account. Even if a meta-regression is based on randomized studies, only evidence of an observed association can be inferred from this analysis, not a causal relationship [506]. Meta-regressions that attempt to show an association between the different effect sizes and the average patient characteristics in individual studies are especially difficult to interpret. These analyses are subject to the same limitations as the results of ecological studies in epidemiology [211]. Due to the high risk of bias, which in analyses based on aggregate data cannot be balanced by adjustment, definite conclusions are only possible on the basis of IPD [480,506] (see also Section 7.2.3).

The Institute uses prediction intervals to display heterogeneity within the framework of a meta-analysis with random effects [217,246,425]. In contrast to the confidence interval, which quantifies the precision of an estimated effect, the 95% prediction interval covers the true effect of a single (new) study with a probability of 95%. In this context it is important to note that a prediction interval cannot be used to assess the statistical significance of an effect. The Institute follows the proposal by Guddat et al. [217] to insert the prediction interval – clearly distinguishable from the confidence interval – in the form of a rectangle in a forest plot. The use of meta-analyses with random effects and related prediction intervals in the event of very few studies (e.g. less than 5) is critically discussed in the literature, as potential heterogeneity can only be estimated very imprecisely [246]. The Institute generally presents prediction intervals in forest plots of meta-analyses with random effects if at least 4 studies are available and if the graphic display of heterogeneity is important. This is particularly the case if, due to great heterogeneity, no pooled effect is presented.

Prediction intervals are therefore particularly used in forest plots if no overall effect can be estimated and displayed due to great heterogeneity. In these heterogeneous situations, the prediction interval is a valuable aid in evaluating whether the study effects are in the same

direction or not or whether in the former case these effects are moderately or clearly in the same direction (see Section 3.1.4).

**C) Subgroup analyses within the framework of meta-analyses**

In addition to the general aspects requiring consideration in the interpretation of subgroup analyses (see Section 7.1.6), there are specific aspects that play a role in subgroup analyses within the framework of meta-analyses. Whereas in general subgroup analyses conducted post hoc on a study level should be viewed critically, in a systematic review one still depends on the use of the results of such analyses on a study level if the review is supposed to investigate precisely these subgroups. In analogy to the approach of not pooling studies with too great heterogeneity by means of meta-analyses, results of subgroups should not be summarized to a common effect estimate if the subgroups differ too strongly from each other. Within the framework of meta-analyses, the Institute usually interprets the results of a heterogeneity or interaction test regarding important subgroups as follows: A significant result at the level of $\alpha = 0.05$ is classified as proof of different effects in the subgroups; a significant result at the level of $\alpha = 0.20$ is classified as an indication of different effects. If the data provide at least an indication of different effects in the subgroups, then the individual subgroup results are reported in addition to the overall effect. If the data provide proof of different effects in the subgroups, then the results for all subgroups are not pooled to a common effect estimate. In the case of more than 2 subgroups, pairwise statistical tests are conducted, if possible, to detect whether subgroup effects exist. Pairs that are not statistically significant at the level of $\alpha = 0.20$ are then summarized in a group. The results of the remaining groups are reported separately and separate conclusions on the benefit of the intervention for these groups are inferred [482].

**D) Small number of events**

A common problem of meta-analyses using binary data is the existence of so-called "zero cells", i.e. cases where not a single event was observed in an intervention group of a study. the Institute follows the usual approach here; i.e. in the event of zero cells, the correction value of 0.5 is added to each cell frequency of the corresponding fourfold table [109]. This approach is appropriate as long as not too many zero cells occur. In the case of a low overall number of events, it may be necessary to use other methods. In the case of very rare events the Peto odds-ratio method can be applied; this does not require a correction term in the case of zero cells [56,109].

If studies do exist in which no event is observed in either study arm (so-called "double-zero studies") then in practice these studies are often excluded from the meta-analytic calculation. This procedure should be avoided if too many double-zero studies exist. Several methods are available to avoid the exclusion of double-zero studies. The absolute risk difference may possibly be used as an effect measure which, especially in the case of very rare events, often does not lead to the heterogeneities that otherwise usually occur. A logistic regression with random effects represents an approach so far rarely applied in practice [519]. Newer methods such as exact methods [510] or the application of the arcsine difference [438] represent

interesting alternatives, but have not yet been investigated sufficiently. Depending on the particular data situation, the Institute will select an appropriate method and, if applicable, examine the robustness of results by means of sensitivity analyses.

**E) Meta-analyses of diagnostic studies**

The results of studies on diagnostic accuracy can also be statistically pooled by means of meta-analytic techniques [124,281]. However, as explained in Section 3.5, studies investigating only diagnostic accuracy are mostly of subordinate relevance in the evaluation of diagnostic tests, so that meta-analyses of studies on diagnostic accuracy are likewise of limited relevance.

The same basic principles apply to a meta-analysis of studies on diagnostic accuracy as to meta-analyses of therapy studies [124,421]. Here too, it is necessary to conduct a systematic review of the literature, assess the methodological quality of the primary studies, conduct sensitivity analyses, and examine the potential influence of publication bias.

In practice, in most cases heterogeneity can be expected in meta-analyses of diagnostic studies; therefore it is usually advisable here to apply random-effects models [124]. Such a meta-analytical pooling of studies on diagnostic accuracy can be performed by means of separate models for sensitivity and specificity. However, if a summarizing receiver operating characteristic (ROC) curve and/or a two-dimensional estimate for sensitivity and specificity are of interest, newer bivariate meta-analyses with random effects show advantages [227,422]. These methods also enable consideration of explanatory variables [226]. Results are presented graphically either via the separate display of sensitivities and specificities in the form of modified forest plots or via a two-dimensional illustration of estimates for sensitivity and specificity. In analogy to the confidence and prediction intervals in meta-analyses of therapy studies, confidence and prediction regions can be presented in the ROC area in bivariate meta-analyses of diagnostic studies.

**F) Cumulative meta-analyses**

For some time it has been increasingly discussed whether, in the case of repeated updates of systematic reviews, one should calculate and present meta-analyses included in these reviews as cumulative meta-analyses with correction for multiple testing [49,62,63,395,507,543]. As a standard the Institute applies the usual type of meta-analyses and normally does not draw upon methods for cumulative meta-analyses.

However, if the conceivable case arises that the Institute is commissioned with the regular update of a systematic review to be updated until a decision can be made on the basis of a statistically significant result, the Institute will consider applying methods for cumulative meta-analyses with correction for multiple testing.

### 7.3.9  Indirect comparisons

Methods are currently being developed that enable the combination of evidence from direct and indirect comparisons. These techniques are called "mixed treatment comparison (MTC) meta-analysis" [347-349], "multiple treatment meta-analysis" (MTM) [79], or "network meta-analysis" [350,447]. These methods represent an important further development of the usual meta-analytic techniques. However, there are still several unsolved methodological problems, so that currently the routine application of these methods within the framework of benefit assessments is not advisable [196,448,485,501]. For this reason, in its benefit assessments of interventions, the Institute primarily uses direct comparative studies (placebo-controlled studies as well as head-to-head comparisons); this means that conclusions for benefit assessments are usually inferred only from the results of direct comparative studies.

In certain situations, as, for example, in assessments of the benefit of drugs with new active ingredients [120], as well as in health economic evaluations (see below), it can however be necessary to consider indirect comparisons and infer conclusions from them for the benefit assessment, taking a lower certainty of results into account.

For the health economic evaluation of interventions, conjoint quantitative comparisons of multiple (i.e. more than 2) interventions are usually required. Limiting the study pool to direct head-to-head comparisons would mean limiting the health economic evaluation to a single pairwise comparison or even making it totally impossible. In order to enable a health economic evaluation of multiple interventions, the Institute can also consider indirect comparisons to assess cost-benefit relations [267], taking into account the lower certainty of results (compared with the approach of a pure benefit assessment).

However, appropriate methods for indirect comparisons need to be applied. The Institute disapproves the use of non-adjusted indirect comparisons (i.e. the naive use of single study arms); it accepts solely adjusted indirect comparisons. These particularly include the approach by Bucher et al. [67], as well as the MTC meta-analysis methods mentioned above. Besides the assumptions of pairwise meta-analyses, which must also be fulfilled here, in indirect comparisons sufficient consistency is also required in the effects estimated in the individual studies. The latter is a critical point, as indirect comparisons provide valid results only if assumptions on consistency are fulfilled. Even though techniques to examine inconsistencies are being developed [126,348], many open methodological questions in this area still exist. It is therefore necessary to describe completely the model applied, together with any remaining unclear issues [501]. In addition, an essential condition for consideration of an indirect comparison is that it is targeted towards the overall research question of interest and not only towards selective components such as individual outcomes.

### 7.3.10 Handling of unpublished or partially published data

In the quality assessment of publications, the problem frequently arises in practice that essential data or information is partially or entirely missing. This mainly concerns "grey

literature" and abstracts, but also full-text publications. Moreover, it is possible that studies have not (yet) been published at the time of the Institute's technology assessment.

It is the Institute's aim to conduct an assessment on the basis of a data set that is as complete as possible. If relevant information is missing, the Institute therefore tries to complete the missing data, among other things by contacting the authors of publications or the study sponsors (see Sections 3.2.1 and 6.1.5). However, depending on the type of product prepared, requests for unpublished information may be restricted due to time limits.

A common problem is that important data required for the conduct of a meta-analysis (e.g. variances of effect estimates) are lacking. However, in many cases, missing data can be calculated or at least estimated from the data available [125,259,404]. If possible, the Institute will apply such procedures.

If data are only partly available or if estimated values are used, the robustness of results will be analysed and discussed, if appropriate with the support of sensitivity analyses (e.g. by presenting best-case and worst-case scenarios). However, a worst-case scenario can only be used here as proof of the robustness of a detected effect. From a worst-case scenario not confirming a previously found effect it cannot be concluded that this effect is not demonstrated. In cases where relevant information is largely or completely lacking, it may occur that a publication cannot be assessed. In such cases, it will merely be noted that further data exist on a particular topic, but are not available for assessment.

### 7.3.11 Description of types of bias

Bias is the systematic deviation of the effect estimate (inferred from study data) from the true effect. Bias may be produced by a wide range of possible causes [86]. The following text describes only the most important types; a detailed overview of various types of bias in different situations is presented by Feinstein [170].

"Selection bias" is caused by a violation of the random principles for sampling procedures, i.e. in the allocation of patients to intervention groups. Particularly in the comparison of 2 groups, selection bias can lead to systematic differences between groups. If this leads to an unequal distribution of important confounders between groups, the results of a comparison are usually no longer interpretable. When comparing groups, randomization is the best method to avoid selection bias [247], as the groups formed do not differ systematically with regard to known as well as unknown confounders. However, structural equality can only be ensured if the sample sizes are sufficiently large. In small studies, despite randomization, relevant differences between groups can occur at random. When comparing groups with structural inequality, the effect of known confounders can be taken into account by applying multi-factorial methods. However, the problem remains of a systematic difference between the groups due to unknown or insufficiently investigated confounders.

Besides the comparability of groups with regard to potential prognostic factors, equality of treatment and equality of observation for all participants play a decisive role. "Performance bias" is bias caused by different types of care provided (apart from the intervention to be investigated). A violation of the equality of observation can lead to detection bias. Blinding is an effective protection against both performance and detection bias [291], which are summarized as "information bias" in epidemiology.

If not taken into account, protocol violations and study withdrawals can cause a systematic bias of study results, called "attrition bias". To reduce the risk of attrition bias, in studies that aim to show superiority, the ITT principle can be applied, where all randomized study participants are analysed within the group to which they were randomly assigned, independently of protocol violations [291,317].

Missing values due to other causes present a similar problem. Missing values not due to a random mechanism can also cause bias in a result [344]. The possible causes and effects of missing values should therefore be discussed on a case-by-case basis and, if necessary, statistical methods should be applied to account or compensate for bias. In this context, replacement methods (imputation methods) for missing values are only one class of various methods available, of which none are regarded to be generally accepted. For example, EMA recommends comparison of various methods for handling missing values in sensitivity analyses [163].

When assessing screening programmes, it needs to be considered that earlier diagnosis of a disease often results only in an apparent increase in survival times, due to non-comparable starting points ("lead time bias"). Increased survival times may also appear to be indicated if a screening test preferably detects mild or slowly progressing early stages of a disease ("length bias"). The conduct of a randomized trial to assess the effectiveness of a screening test can protect against these bias mechanisms [181].

"Reporting bias" is caused by the selective reporting of only part of all relevant data and may lead to an overestimation of the benefit of an intervention in systematic reviews. If, depending on the study results, some analyses or outcomes are not reported or reported in less detail within a publication, or reported in a way deviating from the way originally planned, then "selective" or "outcome reporting bias" is present [84,142,247]. In contrast, "publication bias" describes the fact that studies finding a statistically significant negative difference or no statistically significant difference between the test intervention and control group are not published at all or published later than studies with positive and statistically significant results [496]. The pooling of published results can therefore result in a systematic bias of the common effect estimate. Graphic methods such as the funnel plot [151] and statistical methods such as meta-regression can be used to identify and consider publication bias. These methods can neither certainly confirm nor exclude the existence of publication bias, which underlines the importance of also searching for unpublished data. For example, unpublished

information can be identified and obtained by means of trial registries or requests to manufacturers [327,351,409,495,496].

In studies conducted to determine the accuracy of a diagnostic strategy (index test), results may be biased if the reference test does not correctly distinguish between healthy and sick participants ("misclassification bias"). If the reference test is only conducted in a non-random sample of participants receiving the index test ("partial verification bias") or if the reference test applied depends on the result of the index test ("differential verification bias"), this may lead to biased estimates of diagnostic accuracy. Cases in which the index test itself is a component of the reference test may lead to overestimates of diagnostic accuracy ("incorporation bias") [331].

"Spectrum bias" is a further type of bias mentioned in the international literature. This plays a role in studies where the sample for validation of a diagnostic test consists of persons who are already known to be sick and healthy volunteers as a control group [341]. The validation of a test in such studies often leads to estimates for sensitivity and specificity that are higher than they would be in a clinical situation where patients with a suspected disease are investigated [545]. However, the use of the term "bias" (in the sense of a systematic impairment of internal validity) in this connection is unfortunate, as the results of such studies may well be internally valid if the study is conducted appropriately [545]. Nonetheless, studies of the design described above may have features (particularly regarding the composition of samples) due to which they are not informative for clinical questions in terms of external validity.

As in intervention studies, in diagnostic studies it is necessary to completely consider all study participants (including those with unclear test results) in order to avoid systematic bias of results [331]. While numerous investigations are available on the relevance and handling of publication bias in connection with intervention studies, this problem has been far less researched for diagnostic accuracy studies [331].

A general problem in the estimation of effects is bias caused by measurement errors in the study data collected [82,87]. In practice, measurement errors can hardly be avoided and it is known that non-differential measurement errors can also lead to a biased effect estimate. In the case of a simple linear regression model with a classical measurement error in the explanatory variable, "dilution bias" occurs, i.e. a biased estimate in the direction of the zero effect. However, in other models and more complex situations, bias in all directions is possible. Depending on the research question, the strength of potential measurement errors should be discussed, and, if required, methods applied to adjust for bias caused by measurement errors.

## 7.4  Qualitative methods

### 7.4.1  Qualitative studies

Qualitative research methods are applied to explore and understand subjective experiences, individual actions, and the social world [130,229,355,382]. They can enable access to opinions and experiences of patients, relatives, and medical staff with respect to a certain disease or intervention.

The instruments of qualitative research include focus groups conducted with participants of a randomized controlled trial, for example. Qualitative data can also be collected by means of interviews, observations, and written documents, such as diaries.

An analysis follows collection of data, which mainly aims to identify and analyse overlapping topics and concepts in the data collected. Among other things, qualitative methods can be used as an independent research method, in the preparation of or as a supplement to quantitative studies, within the framework of the triangulation or mixed-method approach, or after the conduct of quantitative studies, in order to explain processes or results. Qualitative research is seen as a method to promote the connection between evidence and practice [132].

Systematic synthesis of various qualitative studies investigating a common research question is also possible [25,316,367,508]. However, no generally accepted approach exists for the synthesis of qualitative studies and the combination of qualitative and quantitative data [132,133].

**A) Qualitative studies in the production of health information**

In the development of health information the Institute uses available qualitative research findings to identify (potential) information needs, as well as to investigate experiences with a certain disease or an intervention.

Relevant publications are then selected by means of prespecified inclusion and exclusion criteria, and the study quality is assessed by means of criteria defined beforehand. The results of the studies considered are extracted, organized by topic, and summarized in a descriptive manner for use in the development of health information. The Institute may also take this approach in the production of reports.

In recent years various instruments for evaluating the quality of qualitative studies have been developed [102]. The main task of the Institute in the assessment of qualitative studies is to determine whether the study design, study quality, and reliability are appropriate for the research question investigated. There is a weaker general consensus with regard to the validity of criteria for the conduct, assessment, and synthesis of qualitative studies when compared with other research areas [130,133,229,382].

**B) Qualitative studies in the production of reports**

Different sources of information can support the integration of systematic reviews [131,336,504]. One possible source are research results from qualitative studies [229,336,384,504]. Qualitative studies seem to be establishing themselves in systematic reviews on the benefit assessment of medical services [130,131,384].

Qualitative research can provide information on the acceptability and suitability of interventions in clinical practice [25,130]. The results of qualitative research can be helpful in the interpretation of a systematic review [504] and may be used in the context of primary studies or systematic reviews on determining patient-relevant outcomes [130,132,316,382,384].

The Institute can use qualitative research findings to identify patient-relevant outcomes, and to present background information on patients' experiences and on the patient relevance of the intervention to be assessed. The Institute can also use these findings in the discussion and interpretation of results of a systematic review.

### 7.4.2 Consultation techniques

The processing of research questions and tasks commissioned to the Institute often requires the consultation of patients, patient representatives, and national and international experts. To do this the Institute uses various consultation techniques.

In the production of reports, the Institute uses these techniques to identify patient-relevant outcomes and to involve national and international experts, and also uses them in the Institute's formal consultation procedure. In the development of health information, consultation techniques serve to involve patients and patient representatives in the identification of information needs, the evaluation of health information, and during consultation.

The Institute uses the following consultation techniques:

- key informant interviews [522], e.g. interviews with patient representatives to identify patient-relevant outcomes

- group meetings and consultations [385,388,389], e.g. within the framework of scientific debates on the Institute's products

- group interviews and focus groups [130,522], e.g. with patients with respect to the evaluation of health information

- surveys and polling (including online polling and feedback mechanisms), e.g. to identify information needs of readers of
www.gesundheitsinformation.de/www.informedhealthonline.org

If a deeper understanding of experiences and opinions is necessary, then the Institute should use the scientific findings obtained from qualitative research. The use of consultation techniques and the involvement of experts are associated with an additional use of resources. However, the involvement of patients in research processes enables the consideration of patient issues and needs as well as the orientation of research towards these issues and needs [398].

**Appendix A – Rationale of the methodological approach for determining the extent of added benefit**

This appendix describes the rationale of the methodological approach for determining the extent of added benefit according to the ANV.

According to §5 (4) Sentence 1 of ANV, the dossier must present and consequently also assess "the extent to which there is added benefit". For this purpose, §5 (7) ANV contains a classification into 6 categories: (1) major added benefit, (2) considerable added benefit, (3) minor added benefit, (4) non-quantifiable added benefit, (5) no added benefit proven, (6) less benefit. For the Categories 1 to 3, §5 (7) ANV also provides a definition, as well as examples of criteria for particular consideration, as orientation for the presentation and assessment. These criteria describe qualitative characteristics (type of outcome) and also explicitly quantitative characteristics (e.g. "major" vs. "moderate" increase in survival time). In addition, a hierarchical ranking of outcomes is obviously intended, as sometimes the same modifier (e.g. "relevant") results in a different extent of added benefit for different outcomes. The corresponding details of the primarily relevant extent categories of added benefit (minor, considerable, major) are shown in Table 11. On the basis of these requirements, it was IQWiG's responsibility to operationalize the extent of added benefit for the benefit assessment.

The criteria provided in §5 (7) ANV for the extent of added benefit designate (legal) terms. Some of these terms are clearly defined (e.g. "survival time", "serious adverse events") and some are not (e.g. "alleviation of serious symptoms"). In addition, the criteria listed are not allocated to all categories. For instance, examples of "survival time" are given only for the categories "considerable" and "major" added benefit.

By using the wording "in particular" in §5 (7) with regard to the Categories 1 to 3, the legislator makes it clear that the criteria allocated to the categories are not to be regarded as conclusive. For instance, even if an increase in survival time is classified as less than "moderate", it cannot be assumed that the legislator would not at least acknowledge a "minor" added benefit. Furthermore, the outcome "(health-related) quality of life", which is explicitly defined as a criterion of benefit in §2 (3) ANV, is not mentioned at all in the list of criteria for the extent of added benefit.

Table 11: Determination of extent of added benefit – Criteria according to the ANV

| | | | | | |
|---|---|---|---|---|---|
| **Extent category** | **Major** <br> **sustained and great improvement** in the therapy-relevant benefit, which has not previously been achieved versus the appropriate comparator therapy | Cure | Major increase in survival time | Long-term freedom from serious symptoms | Extensive avoidance of serious adverse events |
| | **Considerable** <br> **marked improvement** in the therapy-relevant benefit, which has not previously been achieved versus the appropriate comparator therapy | Perceptible alleviation of the disease | Moderate increase in survival time | Alleviation of serious symptoms | Relevant avoidance of serious adverse events <br><br> Important avoidance of other adverse events |
| | **Minor** <br> **moderate and not only marginal improvement** in the therapy-relevant benefit, which has not previously been achieved versus the appropriate comparator therapy | | | Reduction in non-serious symptoms | Relevant avoidance of adverse events |
| ANV: Arzneimittel-Nutzenbewertungsverordnung (Regulation for Early Benefit Assessment of New Pharmaceuticals) | | | | | |

In a first step it is thus reasonable to extend the list of criteria by means of criteria that are qualitatively and quantitatively comparable. These amendments to the ANV requirements are shown in Table 12. In this context, the criteria "cure" and "perceptible alleviation of the disease" were not explicitly considered. The former generally requires operationalization. This should in principle be based on criteria referring to the outcomes "mortality" and "morbidity" (e.g. survival over a defined minimum period in patients with oncological diseases). As the ANV links "cure" solely to a major added benefit, the respective specific operationalization, on the basis of the outcomes used, must be examined with regard to whether this equals a relevant improvement in mortality or serious events. In this sense, a reduction in the duration of symptoms, for instance, in patients with simple infections, is not regarded as a "cure".

On the basis of the above amendments the outcome categories are restructured to illustrate the ranking of outcomes intended in the ANV and to consider disease severity according to §5 (7) ANV. For this purpose, the outcomes are grouped as follows, according to their relevance (see Table 13):

1) all-cause mortality

2) serious (or severe) symptoms (or late complications); serious (or severe) adverse events; health-related quality of life

3) non-serious (or non-severe) symptoms (or late complications); non-serious (or non-severe) adverse events

Health-related quality of life is regarded to be of equal importance as serious (or severe symptoms), late complications and adverse events. The potential categories of extent of added benefit for non-serious outcomes are restricted to "minor" and "considerable".

The requirements of the ANV make it clear that to determine the extent of added benefit, first the effect sizes must be described at outcome level. For each outcome separately the effect size – independent of its direction – is classified into 1 of the 3 extent categories (minor, considerable, major). Within the overall weighing of benefits and harms, these individual outcomes are then summarized into a global conclusion on the extent of added benefit. This step-by-step approach is described in Section 3.3.3.

Table 12: Determination of extent of added benefit – Criteria according to the ANV plus amendments[a]

| | | Outcome category | | | |
|---|---|---|---|---|---|
| | | *All-cause mortality* | *Symptoms (morbidity)* | *Health-related quality of life* | *Adverse events* |
| **Extent category** | **Major** **sustained and great improvement** in the therapy-relevant benefit, which has not previously been achieved versus the appropriate comparator therapy | Major increase in survival time | Long-term freedom from serious (*or severe*) symptoms (*or late complications*) | *Major improvement in quality of life* | Extensive avoidance of serious (*or severe*) adverse events |
| | **Considerable** **marked improvem**ent in the therapy-relevant benefit, which has not previously been achieved versus the appropriate comparator therapy | Moderate increase in survival time | Alleviation of serious (*or severe*) symptoms (*or late complications*) *Important reduction in non-serious (or non-severe) symptoms (or late complications)* | *Important improvement in quality of life* | Relevant avoidance of serious (*or severe*) adverse events Important avoidance of other (*non-serious or non-severe*) adverse events |
| | **Minor** **moderate and not only marginal improvement** in the therapy-relevant benefit, which has not previously been achieved versus the appropriate comparator therapy | *Any increase in survival time* | *Any reduction in serious (or severe) symptoms (or late complications)* Reduction in non-serious (*or non-severe*) symptoms *(or late complications)* | *Relevant improvement in quality of life* | *Any statistically significant reduction in serious (or severe) adverse events* Relevant avoidance of *(other, non-serious or non-severe)* adverse events |
| a: Amendments to the ANV in *italics.* | | | | | |
| ANV: Arzneimittel-Nutzenbewertungsverordnung (Regulation for Early Benefit Assessment of New Pharmaceuticals) | | | | | |

Table 13: Determination of extent of added benefit – Ranked criteria according to the ANV plus amendments[a]

| | | Outcome category | | | |
|---|---|---|---|---|---|
| | | **All-cause mortality** | **Serious** *(or severe)* **symptoms** *(or late complications)* **and adverse events** | *Health-related quality of life* | **Non-serious** *(or non-severe)* **symptoms** *(or late complications)* **and adverse events** |
| **Extent category** | **Major** <br> **sustained and great improvement** in the therapy-relevant benefit, which has not previously been achieved versus the appropriate comparator therapy | Major increase in survival time | Long-term freedom or extensive avoidance | *Major improvement* | *Not applicable* |
| | **Considerable** <br> **marked improvement** in the therapy-relevant benefit, which has not previously been achieved versus the appropriate comparator therapy | Moderate increase in survival time | Alleviation or relevant avoidance | *Important improvement* | Important avoidance |
| | **Minor** <br> **moderate and not only marginal improvement** in the therapy-relevant benefit, which has not previously been achieved versus the appropriate comparator therapy | *Any increase in survival time* | *Any reduction* | *Relevant improvement* | Relevant avoidance |
| a: Amendments to the ANV in *italics*. <br> ANV=Arzneimittel-Nutzenbewertungsverordnung (Regulation for Early Benefit Assessment of New Pharmaceuticals) | | | | | |

In accordance with §2 (3) ANV, the term "benefit" is defined as an "effect" and in §2 (4) ANV the term "added benefit" is defined as such an effect compared with the appropriate comparator therapy. It can be inferred from these definitions that the extent of added benefit must be determined by taking into account both the hierarchy of outcomes and effect sizes.

The ANV provides no details on the questions as to which effect sizes for the individual outcomes result in which extent category, or which effect measures should be chosen for the assessment. In principle, these questions can only be partly answered from a methodological point of view. Nevertheless, IQWiG is required to assess the extent of added benefit presented in the dossiers (§7 (2) ANV) and to draw its own conclusions on the extent. To restrict to a minimum at this stage the value judgements that will necessarily be made in the further deliberation process and to reveal them, the following measures are required:

- explicit operationalization to ensure a transparent approach

- abstract operationalization to achieve the best possible consistency between benefit assessments

Against this background a suitable effect measure must first be chosen. The initial focus is on the situation with binary data (analysis of 2x2 tables). In this context, relative effect measures – these mainly comprise the relative risk (RR) and the odds ratio (OR) – show the following advantages over absolute measures such as the risk difference (RD):

- The risk difference does not describe the effectiveness of therapy as such, as this difference strongly depends on the baseline risk in the control group. However, the baseline risk varies between regions, populations and over the course of time, as well as particularly between control groups receiving different comparator therapies. A risk difference should thus be interpreted as a descriptive measure of a specific study, not as a fixed measure of a specific treatment procedure; this is also and primarily a problem in meta-analyses [483]. This great susceptibility to external conditions calls into question the transferability of absolute effect measures from clinical studies to the daily healthcare setting. It is therefore common practice preferably to express effects shown in clinical studies as relative risks, odds ratios or hazard (or incidence) ratios [108].

- The degree of the risk difference is limited by the degree of the baseline risk (absolute risk in the control group). If this baseline risk is 1%, then the risk difference can never exceed 0.01 (or if it is 10%, the risk difference can never exceed 0.1 etc.). The risk difference could only reach the optimum value of 1 if the baseline risk was 100%. For instance, if an absolute risk reduction of at least 20% was defined as a substantial therapeutic improvement, then, for this example of a requirement, in diseases with (long-term) survival rates of > 80%, generally a major added benefit (for the corresponding outcome) would no longer be presentable.

▪ A further disadvantage of the use of the absolute risk reduction as an effect measure to operationalize the determination of the extent of added benefit is that an exact time point must be defined at which this absolute risk reduction is determined (e.g. after 1, 2, 5 or 10 years), if no generally accepted definitions are available (e.g. 30-day mortality for myocardial infarction).

In summary, absolute risk reductions may have more of an impact in a situation of individual decision making, but relative effect measures are more suitable for general conclusions in terms of an assessment of the added benefit of a drug.

Relative measures have in common that the zero effect (no group difference) is 1. In the following text we address effects below 1, from which effects above 1 can be calculated by using the reciprocal. For the result to be classified as a minor, considerable or major added benefit, the approach stipulates that the (two-sided) 95% confidence interval of the effect undercuts the respective threshold in terms of a shift in the hypothesis boundary. In comparison with the examination of point estimates, such an inferential statistical approach has 2 main advantages: (i) The precision of the estimate is considered in the assessment; and accordingly, (ii) the probability of statistical errors can be limited to the usual small values (e.g. 5%).

The thresholds vary with regard to the 2 dimensions "outcome category" and "extent category (of the effect)" displayed in Table 13. The greater the relevance ascribed to the outcome, the closer the thresholds should lie to 1 (below 1). This takes into account the ANV's requirement to consider disease severity. In contrast, the greater the determined extent of the effect, the further the thresholds should lie from 1 (below 1).

Following the explicit and abstract operationalization above, a division of the thresholds in step sizes of 0.05 is planned [272]. The further development of the methodological approach leading to these thresholds is briefly explained in the following text. The further deliberations will show that the choice of 0.05 is applicable in practice and leads to reasonable conclusions.

The starting point was formed by the question as to how large the actual effects have to be in order to be classified, for instance, as effects of a "major" extent. For this purpose, a relative risk of 0.50 – proposed by Djulbegovic et al. [134] as a requirement for a "breakthrough" – was defined as an effect of a major extent for the outcome "all-cause mortality" [272].

For this actual effect (0.5) the question arises as to how the threshold should be chosen to really achieve the extent "major" with adequate power. Details of the corresponding considerations can be found in the first dossier assessment conducted by the Institute [272], but are also addressed again at the end of this appendix. Following these considerations, the simultaneous requirement for feasibility and stringency can be regarded as fulfilled for a threshold of 0.85.

In a next step, for the matrix of the extent, the other actual effects are specified and the corresponding thresholds determined. In this context it should be considered that, on the basis of the outcome category "mortality", the requirements should increase for less serious outcomes, and on the basis of the extent category "major", should decrease for lower extent categories. In this context, a division into sixths for the actual effects was shown to be a pragmatical solution. The thresholds for the respective extent categories are described in the following text.

**1. All-cause mortality**

With the usual significance level of 5%, any statistically significant increase in survival time is at least classified as "minor added benefit", since for all-cause mortality the requirement that an effect should be "more than marginal" is regarded to be fulfilled by the outcome itself. The threshold referring to the 95% confidence interval is thus 1 here. An increase in survival time is classified as a "considerable" effect if a threshold of 0.95 is undercut. An increase in survival time is classified as being "major" if the threshold of 0.85 is undercut by the upper limit of the 95% confidence interval.

**2. Serious (or severe) symptoms (or late complications), serious or (severe) adverse events, health-related quality of life**

For serious (or severe) symptoms (or late complications) and serious (or severe) adverse events, any statistically significant reduction also represents at least a "minor" effect, as the requirement of "more than marginal" is already fulfilled by the quality of the outcome itself. In contrast to the desired effects on all-cause mortality, a "considerable" effect requires that a threshold of 0.90 must be undercut and a "major" effect requires that a threshold of 0.75 is undercut. To derive a major effect from these outcomes also requires that the risk of the examined event should be at least 5% in at least one of the groups compared. This additional criterion supports the relevance of the event at population level and allows for the special requirements for this category of added benefit.

The precondition for determining the extent of added benefit for outcomes on health-related quality of life (as for all PROs) is that both the instruments applied and the response criteria must be validated or at least generally established. If these results are dichotomous in terms of responders and non-responders, the above criteria for serious symptoms apply (risk for the category "major" should be at least 5%).

**3. Non-serious (or non-severe) symptoms (or late complications), non-serious (or non-severe) adverse events**

The specification of thresholds for the non-serious (or non-severe) symptoms (or late complications) and the non-serious (or non-severe) adverse events takes into account the lower severity compared with Categories 1 and 2. As a matter of principle, the effect for non-serious outcomes should not be classified as "major". To classify an effect as "considerable" or "minor" the thresholds of 0.80 or 0.90 respectively must be undercut. In the latter case, this

is based on the requirement for minor added benefit specified in §5 (7) ANV that there must be a moderate, and not only marginal, improvement. The procedure thus implies that effects (also statistically significant ones) only assessed as "marginal" lead to classification into the category "no added benefit".

The corresponding thresholds for all extent categories and outcome categories are presented in the following Table 14.

Table 14: Inferential statistical thresholds (hypotheses boundaries) for relative effect measures

<table>
<tr><td colspan="2" rowspan="2"></td><td colspan="3"><strong>Outcome category</strong></td></tr>
<tr><td>All-cause mortality</td><td>Serious (or severe) symptoms (or late complications) and adverse events, as well as quality of life[a]</td><td>Non-serious (or non-severe) symptoms (or late complications) and adverse events</td></tr>
<tr><td rowspan="3"><strong>Extent category</strong></td><td>Major</td><td>0.85</td><td>0.75 and risk ≥ 5%[b]</td><td>Not applicable</td></tr>
<tr><td>Considerable</td><td>0.95</td><td>0.90</td><td>0.80</td></tr>
<tr><td>Minor</td><td>1.00</td><td>1.00</td><td>0.90</td></tr>
<tr><td colspan="5">a: Precondition (as for all patient-reported outcomes): use of a validated or established instrument, as well as a validated or established response criterion. Values apply to non-response.<br>b: Risk must be at least 5% for at least 1 of the 2 groups compared.</td></tr>
</table>

**Detailed methodological rationale for determination of thresholds**

The starting point is the planning of a (fictional) study to test the conventional hypotheses

$$H_0: RR \geq RR_0 \quad vs. \quad H_1: RR < RR_0$$

on the basis of the relative risk $RR_0 = 1$. The required sample size is calculated by specifying the significance level, the power, the risk in the control group, and the actual effect ($RR_1$).

For all hypothesis boundaries shifted from 1 ($RR_0 < 1$) a study of this sort has reduced power. In order to maintain the same power for the shifted hypothesis boundary of interest (the thresholds named above) as specified for the testing of the conventional (non-shifted) hypotheses, the sample size must be increased – either within the study or through a combination of several studies. Assuming the normal case of 2 (e.g. pivotal) studies, it can be assumed that the sample size is twice as large. The hypothesis boundary for the shifted hypotheses is then precisely selected so that the power for the conventional hypotheses of the 2 individual studies corresponds to the power for the shifted hypotheses of the combined (pooled) analysis. This hypothesis boundary serves as the threshold for the upper limit of the two-sided 95% confidence interval for the relative risk. For instance, the specification of a significance level of 5% (two-sided) and a power of 90% (both for the conventional and for

the shifted hypothesis boundary), as well as a doubling of the sample size for the shifted hypothesis boundary resulted in a threshold of (rounded) 0.85 for the actual effect of 0.5 postulated for the outcome "mortality" and the extent category "major".

The formula included in Appendix A of the benefit assessment on ticagrelor [272] for the relationship between the actual effect and the threshold is independent of the other requirements and is based on the algorithm used in the "power" procedure of the software SAS. The corresponding documentation for this algorithm [456] refers to the work by Fleiss et al. [178], A query to Mr Röhmel (former Speaker of the Working Group "Pharmaceutical Research" of the German Region of the International Biometric Society), as well as directly to the Technical Support Section of SAS, showed that documentation of the validity of this algorithm has evidently not been published. The question arises as to which actual effects are required in more precise calculations to reach the respective extent category with high probability.

The actual effects were thus determined by means of Monte Carlo simulations as follows:

1.  The significance level for the above hypothesis is 2.5% and the power is 90%. The parameter $RR_1$ runs through all values between 0.2 and 0.95 at step sizes of 0.01. The risk in the control group $p_C$ runs through all values between 0.05 and 0.95 at step sizes of 0.05. For each of these tuples $(RR_1, p_C)$ the required sample size $n$ is calculated using $RR_0 = 1$ according to the formula by Farrington and Manning [168] and then doubled ($m := 2n$).

2.  For each triple $(RR_1, p_C, m)$ a threshold $T$ runs through all values between 1 and 0 in a descending order with a step size of -0.005. For each $T$ the power for the above hypothesis is approximated with $RR_0 = T$. The significance level is 2.5%. For this purpose 50 000 2x2 tables are simulated with a random generator, the upper confidence interval limit for the relative risk is calculated by means of the normal distribution approximation and the delta method for estimation of variance. Subsequently, the proportion of simulation cycles is determined for which the upper confidence interval limit is smaller than $T$. The $T$ cycle is stopped as soon as an approximated power is smaller than 90%. The corresponding triple $(RR_1, p_C, T)$ is documented in a list.

3.  After the cycle of all parameters in Steps 1 and 2, all triples are chosen from the list for which the threshold $T$ deviates less than 0.01 from one of the values 0.75, 0.80, 0.85, 0.90 or 0.95.

Figure 8 shows the resulting (more precise) actual effects, depending on the risk in the control group for all thresholds specified above (points approximated by smoothed curves).
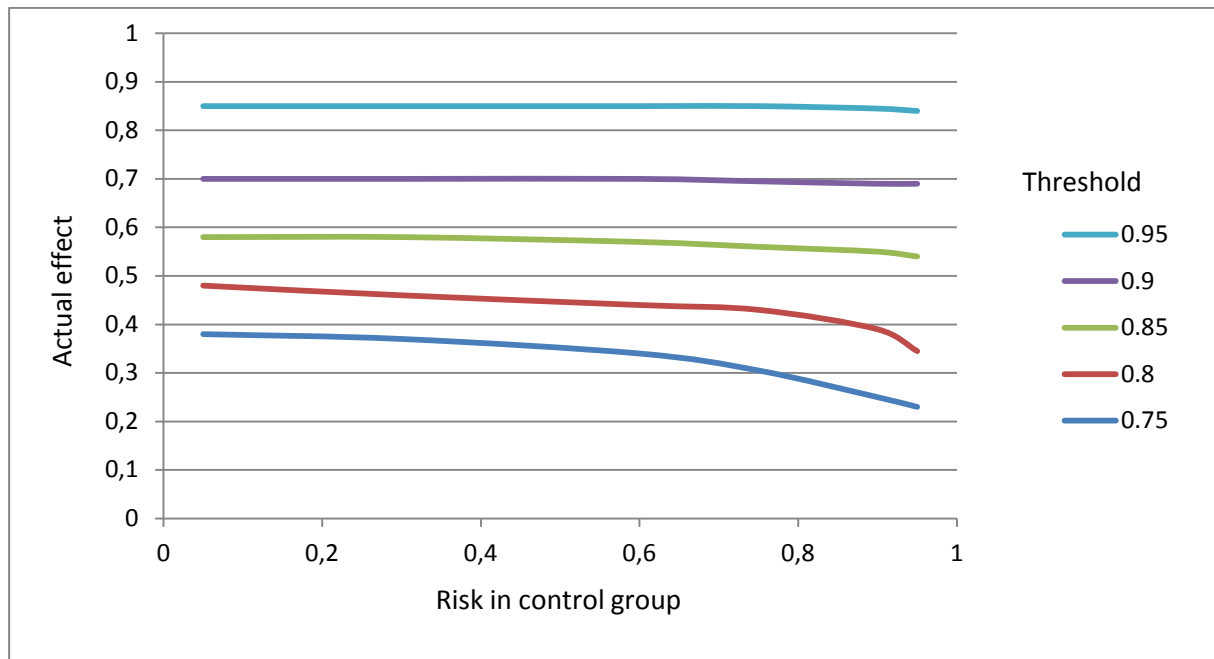
Figure 8: Actual effects depending on the baseline risk

Table 15 again contains the ranges (depending on the risk of the control group) in which the actual effects are realized, per outcome category and extent category.

Table 15: Actual effects for the relative risk

| | | Outcome category | | |
|---|---|---|---|---|
| | | All-cause mortality | Serious (or severe) symptoms (or late complications) and adverse events as well as quality of life | Non-serious (or non-severe) symptoms (or late complications) and adverse events |
| **Extent category** | Major | 0.53–0.58 | 0.24–0.38 | Not applicable |
| | Considerable | 0.84–0.85 | 0.69–0.71 | 0.34–0.48 |
| | Minor | Not applicable | Not applicable | 0.69–0.71 |

In relation to all-cause mortality, actual relative risks of about 0.55 – i.e. still corresponding to about a halving of the risk – are to be specified for the extent "major". For the extent "considerable" the actual effect must lie at about 0.85. For serious symptoms and comparable outcomes, to be classified as a "major" extent, an actual reduction in risk to about a quarter to a third of the risk is required. Compared with the originally specified actual effects [272] good consistency is provided for thresholds lying close to 1. For the thresholds lying further away from 1, the simulation results show slightly more moderate requirements for the strength of the actual effects. The division of the thresholds as defined in Table 14 seems reasonable and practicable.

## 8   References

1. Editorial commentary: avoiding biased comparisons [online]. In: James Lind Library. 2007 [accessed: 19 April 2013]. URL: http://www.jameslindlibrary.org/essays/bias/avoiding-biased-comparisons.html.

2. SGB V Handbuch Sozialgesetzbuch V: Krankenversicherung. Altötting: KKF-Verlag; 2011.

3. AGREE Collaboration. Appraisal of guidelines for research & evaluation: AGREE instrument; training manual [online]. January 2003 [accessed: 8 August 2011]. URL: http://www.agreecollaboration.org/pdf/aitraining.pdf.

4. AGREE Collaboration. Appraisal of guidelines for research & evaluation (AGREE) instrument. London: AGREE Research Trust; 2006. URL: http://www.agreetrust.org/?o=1085.

5. AGREE Next Steps Consortium. Appraisal of guidelines for research and evaluation II: AGREE II instrument [online]. May 2009 [accessed: 23 April 2013]. URL: http://www.agreetrust.org/index.aspx?o=1397.

6. Agresti A. Modelling ordered categorical data: recent advances and future challenges. Stat Med 1999; 18(18): 2191-2207.

7. Agresti A (Ed). Categorical data analysis. Hoboken: Wiley; 2002.

8. Agresti A. Dealing with discreteness: making 'exact' confidence intervals for proportions, differences of proportions, and odds ratios more exact. Stat Methods Med Res 2003; 12(1): 3-21.

9. Al-Marzouki S, Evans S, Marshall T, Roberts I. Are these data real? Statistical methods for the detection of data fabrication in clinical trials. BMJ 2005; 331(7511): 267-270.

10. Altenhofen L, Blumenstock G, Diel F, Döbler K, Geraedts M, Jäckel WH et al. Qualitätsindikatoren: Manual für Autoren. Neukirchen: Make a Book; 2009. (ÄZQ-Schriftenreihe; Volume 36). URL: http://www.aezq.de/mdb/edocs/pdf/schriftenreihe/schriftenreihe36.pdf.

11. Altman DG. Systematic reviews of evaluations of prognostic variables. In: Egger M, Davey Smith G, Altman DG (Ed). Systematic reviews in health care: meta-analysis in context. London: BMJ Publishing Group; 2001. p. 228-247.

12. Altman DG, Bland JM. Statistic notes: absence of evidence is not evidence of absence. BMJ 1995; 311(7003): 485.

13. Altman DG, Machin D, Bryant TM, Gardner MJ. Statistics with confidence: confidence intervals and statistical guidelines. London: BMJ Publishing Group; 2000.

14. Altman DG, Vergouwe Y, Royston P, Moons KGM. Prognosis and prognostic research: validating a prognostic model. BMJ 2009; 338(7708): 1432-1435.

15. American Society of Clinical Oncology. Outcomes of cancer treatment for technology assessment and cancer treatment guidelines. J Clin Oncol 1996; 14(2): 671-679.

16. Amir E, Seruga B, Martinez-Lopez J, Kwong R, Pandiella A, Tannock IF et al. Oncogenic targets, magnitude of benefit, and market pricing of antineoplastic drugs. J Clin Oncol 2011; 29(18): 2543-2549.

17. Antman EM, Lau J, Kupelnick B, Mosteller F, Chalmers TC. A comparison of results of meta-analyses of randomized control trials and recommendations of clinical experts: treatments for myocardial infarction. JAMA 1992; 268(2): 240-248.

18. Arbeits- und Koordinierungsstelle Gesundheitsversorgungsforschung. Evaluation der Gesundheitsinformationsseite "Spezial: Wechseljahre": Interviews mit Nutzerinnen zum Informationspaket „Wechseljahre" des Instituts für Qualität und Wirtschaftlichkeit im Gesundheitswesen"; Abschlussbericht [online]. 15 August 2008 [accessed: 22 April 2013]. URL: http://www.akg.uni-bremen.de/pages/download.php?ID=27&SPRACHE=de&TABLE=AP&TYPE=PDF.

19. Arbeitsgemeinschaft der Wissenschaftlichen Medizinischen Fachgesellschaften, Ärztliche Zentralstelle Qualitätssicherung. Das Leitlinien-Manual von AWMF und ÄZQ. Z Arztl Fortbild Qualitatssich 2001; 95(Suppl 1): 5-84.

20. Arbeitskreis Versorgungsforschung beim Wissenschaftlichen Beirat. Definition und Abgrenzung der Versorgungsforschung [online]. 8 September 2004 [accessed: 23 April 2013]. URL: http://www.bundesaerztekammer.de/downloads/Definition.pdf.

21. Arnold M. Gesundheitssystemforschung. In: Hurrelmann K, Laaser U (Ed). Gesundheitswissenschaften: Handbuch für Lehre, Forschung und Praxis. Weinheim: Beltz; 1993. p. 423-437.

22. Assmann SF, Pocock SJ, Enos LE, Kasten LE. Subgroup analysis and other (mis)uses of baseline data in clinical trials. Lancet 2000; 355(9209): 1064-1069.

23. Atkins D, Best D, Briss PA, Eccles MP, Falck-Ytter Y, Flottorp S et al. Grading quality of evidence and strength of recommendations. BMJ 2004; 328(7454): 1490.

24. Atkins D, Eccles MP, Flottorp S, Guyatt GH, Henry D, Hill S et al. Systems for grading the quality of evidence and the strength of recommendations; I: critical appraisal of existing approaches. BMC Health Serv Res 2004; 4: 38.

25. Atkins S, Lewin S, Smith H, Engel M, Fretheim A, Volmink J. Conducting a meta-ethnography of qualitative literature: lessons learnt. BMC Med Res Methodol 2008; 8: 21.

26. Bagley SC, White H, Golomb BA. Logistic regression in the medical literature: standards for use and reporting, with particular attention to one medical domain. J Clin Epidemiol 2001; 54(10): 979-985.

27. Baker SG. Surrogate endpoints: wishful thinking or reality? J Natl Cancer Inst 2006; 98(8): 502-503.

28. Baltussen R, Niessen L. Priority setting of health interventions: the need for multi-criteria decision analysis. Cost Eff Resour Alloc 2006; 4: 14.

29. Baltussen R, Stolk E, Chisholm D, Aikins M. Towards a multi-criteria approach for priority setting: an application to Ghana. Health Econ 2006; 15(7): 689-696.

30. Banta D. The development of health technology assessment. Health Policy 2003; 63(2): 121-132.

31. Barron BA, Bukantz SC. The evaluation of new drugs: current Food and Drug Administration regulations and statistical aspects of clinical trials. Arch Intern Med 1967; 119(6): 547-556.

32. Bastian H. Health information on the internet. In: Heggenhougen HK, Quah SR (Ed). International encyclopedia of public health: volume 3; G-I. Amsterdam: Academic Press; 2008. p. 678-682.

33. Bastian H, Glasziou P, Chalmers I. Seventy-five trials and eleven systematic reviews a day: how will we ever keep up? PLoS Med 2010; 7(9): e1000326.

34. Bastian H, Kaiser T, Matschewsky S. Förderung allgemeiner Gesundheits- und Wissenschaftskenntnisse mittels Bürger- und Patienteninformationen: die Rolle des IQWiG. Z Arztl Fortbild Qualitatssich 2005; 99(6): 379-385.

35. Bates BR, Romina S, Ahmed R, Hopson D. The effect of source credibility on consumers' perceptions of the quality of health information on the internet. Med Inform Internet Med 2006; 31(1): 45-52.

36. Bender R. Interpretation von Effizienzmaßen der Vierfeldertafel für Diagnostik und Behandlung. Med Klin 2001; 96(2): 116-121.

37. Bender R, Bunce C, Clarke M, Gates S, Lange S, Pace NL et al. Attention should be given to multiplicity issues in systematic reviews. J Clin Epidemiol 2008; 61(9): 857-865.

38. Bender R, Grouven U. Logistic regression models used in medical research are poorly presented. BMJ 1996; 313(7057): 628.

39. Bender R, Grouven U. Ordinal logistic regression in medical research. J R Coll Physicians Lond 1997; 31(5): 546-551.

40. Bender R, Lange S. Adjusting for multiple testing: when and how? J Clin Epidemiol 2001; 54(4): 343-349.

41. Bent S, Padula A, Avins AL. Brief communication: better ways to question patients about adverse medical events; a randomized, controlled trial. Ann Intern Med 2006; 144(4): 257-261.

42. Bessell TL, McDonald S, Silagy CA, Anderson JN, Hiller JE, Sansom LN. Do internet interventions for consumers cause more harm than good? A systematic review. Health Expect 2002; 5(1): 28-37.

43. Biglan A, Ary D, Wagenaar AC. The value of interrupted time-series experiments for community intervention research. Prev Sci 2000; 1(1): 31-49.

44. Biomarkers Definitions Working Group. Biomarkers and surrogate endpoints: preferred definitions and conceptual framework. Clin Pharmacol Ther 2001; 69(3): 89-95.

45. Bland JM, Altman DG. One and two sided tests of significance. BMJ 1994; 309(6949): 248.

46. Bock J, Toutenburg H. Sample size determination in clinical research. In: Rao CR, Chakraborty R (Ed). Statistical methods in biological and medical sciences. Amsterdam: Elsevier; 1991. p. 515-538. (Handbook of Statistics; Volume 8).

47. Bonell C, Oakley A, Hargreaves J, Strange V, Rees R. Assessment of generalisability in trials of health interventions: suggested framework and systematic review. BMJ 2006; 333(7563): 346-349.

48. Bonhoeffer J, Zumbrunn B, Heininger U. Reporting of vaccine safety data in publications: systematic review. Pharmacoepidemiol Drug Saf 2005; 14(2): 101-106.

49. Borm GF, Donders AR. Updating meta-analyses leads to larger type I errors than publication bias. J Clin Epidemiol 2009; 62(8). 825-830, 830.e1-830.e10.

50. Bossuyt PM, Irwig LM, Craig J, Glasziou PP. Comparative accuracy: assessing new tests against existing diagnostic pathways. BMJ 2006; 332(7549): 1089-1092.

51. Bossuyt PM, Reitsma JB, Bruns DE, Gatsonis CA, Glasziou PP, Irwig LM et al. Towards complete and accurate reporting of studies of diagnostic accuracy: the STARD initiative. Ann Intern Med 2003; 138(1): 40-44.

52. Bossuyt PM, Reitsma JB, Bruns DE, Gatsonis CA, Glasziou PP, Irwig LM et al. The STARD statement for reporting studies of diagnostic accuracy: explanation and elaboration. Ann Intern Med 2003; 138(1): W1-W12.

53. Boutron I, Moher D, Altman DG, Schulz KF, Ravaud P. Extending the CONSORT statement to randomized trials of nonpharmacologic treatment: explanation and elaboration. Ann Intern Med 2008; 148(4): 295-309.

54. Boutron I, Moher D, Altman DG, Schulz KF, Ravaud P. Methods and processes of the CONSORT group: example of an extension for trials assessing nonpharmacologic treatments. Ann Intern Med 2008; 148(4): W60-W66.

55. Boynton J, Glanville J, McDaid D, Lefebvre C. Identifying systematic reviews in MEDLINE: developing an objective approach to search strategy design. Journal of Information Science 1998; 24(3): 137-155.

56. Bradburn MJ, Deeks JJ, Berlin JA, Russell Localio A. Much ado about nothing: a comparison of the performance of meta-analytical methods with rare events. Stat Med 2007; 26(1): 53-77.

57. Breslow NE, Day NE. Statistical methods in cancer research; volume I: the analysis of case-control studies. Lyon: International Agency for Research on Cancer; 1980. (IARC Scientific Publications; Volume 32). URL: http://www.iarc.fr/en/publications/pdfs-online/stat/sp32/SP32.pdf.

58. Breslow NE, Day NE. Statistical methods in cancer research; volume II: the design and analysis of cohort studies. Lyon: International Agency for Research on Cancer; 1987. (IARC Scientific Publications; Volume 82). URL: http://www.iarc.fr/en/publications/pdfs-online/stat/sp82/SP82.pdf.

59. Bridges JF, Kinter ET, Kidane L, Heinzen RR, McCormick C. Things are looking up since we started listening to patients: trends in the application of conjoint analysis in health 1982-2007. Patient 2008; 1(4): 273-282.

60. British Cardiac Society, British Hyperlipidaemia Association, British Hypertension Society. Joint British recommendations on prevention of coronary heart disease in clinical practice. Heart 1998; 80(Suppl 2): 1-29.

61. Brockwell SE, Gordon IR. A comparison of statistical methods for meta-analysis. Stat Med 2001; 20(6): 825-840.

62. Brok J, Thorlund K, Gluud C, Wetterslev J. Trial sequential analysis reveals insufficient information size and potentially false positive results in many meta-analyses. J Clin Epidemiol 2008; 61(8): 763-769.

63. Brok J, Thorlund K, Wetterslev J, Gluud C. Apparently conclusive meta-analyses may be inconclusive: trial sequential analysis adjustment of random error risk due to repetitive testing of accumulating data in apparently conclusive neonatal meta-analyses. Int J Epidemiol 2009; 38(1): 287-298.

64. Brookes ST, Whitely E, Egger M, Davey Smith G, Mulheran PA, Peters TJ. Subgroup analyses in randomized trials: risks of subgroup-specific analyses; power and sample size for the interaction test. J Clin Epidemiol 2004; 57(3): 229-236.

65. Brouwers MC, Kho ME, Browman GP, Burgers JS, Cluzeau F, Feder G et al. Development of the AGREE II; part 1: performance, usefulness and areas for improvement. CMAJ 2010; 182(10): 1045-1052.

66. Brouwers MC, Kho ME, Browman GP, Burgers JS, Cluzeau F, Feder G et al. Development of the AGREE II; part 2: assessment of validity of items and tools to support application. CMAJ 2010; 182(10): E472-E478.

67. Bucher HC, Guyatt GH, Griffith LE, Walter SD. The results of direct and indirect treatment comparisons in meta-analysis of randomized controlled trials. J Clin Epidemiol 1997; 50(6): 683-691.

68. Bundesministerium der Justiz. Gesetz über die Statistik für Bundeszwecke (Bundesstatistikgesetz - BStatG) [online]. 7 September 2007 [accessed: 6 May 2013]. URL: http://www.gesetze-im-internet.de/bundesrecht/bstatg_1987/gesamt.pdf.

69. Bundesministerium der Justiz. Verordnung zur Schaffung barrierefreier Informationstechnik nach dem Behindertengleichstellungsgesetz (Barrierefreie-Informationstechnik-Verordnung - BITV 2.0) [online]. 12 September 2011 [accessed: 6 May 2013]. URL: http://www.gesetze-im-internet.de/bitv_2_0/BJNR184300011.html.

70. Bundesministerium für Gesundheit. Verordnung über die Nutzenbewertung von Arzneimitteln nach § 35a Absatz 1 SGB V für Erstattungsvereinbarungen nach § 130b SGB V (Arzneimittel-Nutzenbewertungsverordnung – AM-NutzenV). Bundesgesetzblatt Teil 1 2010; (68): 2324-2328.

71. Bundessozialgericht. Urteil: Aktenzeichen B 6 A 1/08 R [online]. 6 May 2009 [accessed: 19 April 2013]. URL: http://juris.bundessozialgericht.de/cgi-bin/rechtsprechung/document.py?Gericht=bsg&Art=en&sid=965bc60820d25990f7f287c0fa2b4c2c&nr=11110&pos=0&anz=1.

72. Bundesverfassungsgericht. Leitsatz zum Beschluss des Ersten Senats: Aktenzeichen 1 BvR 347/98 [online]. 6 December 2005 [accessed: 23 April 2013]. URL: http://www.bverfg.de/entscheidungen/rs20051206_1bvr034798.html.

73. Bundesversicherungsamt. Festlegungen [online]. [Accessed: 23 April 2013]. URL: http://www.bundesversicherungsamt.de/cln_115/nn_1046668/DE/Risikostrukturausgleich/Festlegungen/festlegungen__node.html?__nnn=true.

74. Burgers JS. Guideline quality and guideline content: are they related? Clin Chem 2006; 52(1): 3-4.

75. Burgers JS, Bailey JV, Klazinga NS, Van der Bij AK, Grol R, Feder G. Inside guidelines: comparative analysis of recommendations and evidence in diabetes guidelines from 13 countries. Diabetes Care 2002; 25(11): 1933-1939.

76. Burzykowski T, Buyse M. Surrogate threshold effect: an alternative measure for meta-analytic surrogate endpoint validation. Pharm Stat 2006; 5(3): 173-186.

77. Burzykowski T, Molenberghs G, Buyse M (Ed). The evaluation of surrogate endpoints. New York: Springer; 2005.

78. Buyse M, Molenberghs G, Burzykowski T, Renard D, Geys H. The validation of surrogate endpoints in meta-analyses of randomized experiments. Biostatistics 2000; 1(1): 49-67.

79. Caldwell DM, Ades AE, Higgins JP. Simultaneous comparison of multiple treatments: combining direct and indirect evidence. BMJ 2005; 331(7521): 897-900.

80. Campbell F, Dickinson HO, Cook JV, Beyer FR, Eccles M, Mason JM. Methods underpinning national clinical guidelines for hypertension: describing the evidence shortfall. BMC Health Serv Res 2006; 6: 47.

81. Campbell MK, Elbourne DR, Altman DG. CONSORT statement: extension to cluster randomised trials. BMJ 2004; 328(7441): 702-708.

82. Carroll RJ, Ruppert D, Stefanski LA, Crainiceanu CM. Measurement error in nonlinear models: a modern perspective. London: Chapman & Hall; 2006.

83. Centre for Evidence-based Medicine. Levels of evidence (March 2009) [online]. March 2009 [accessed: 19 April 2013]. URL: http://www.cebm.net/index.aspx?o=1025.

84. Chan AW, Hróbjartsson A, Haahr MT, Gøtzsche PC, Altman DG. Empirical evidence for selective reporting of outcomes in randomized trials: comparison of protocols to published articles. JAMA 2004; 291(20): 2457-2465.

85. Charles C, Gafni A, Whelan T, O'Brien MA. Treatment decision aids: conceptual issues and future directions. Health Expect 2005; 8(2): 114-125.

86. Chavalarias D, Ioannidis JP. Science mapping analysis characterizes 235 biases in biomedical research. J Clin Epidemiol 2010; 63(11): 1205-1215.

87. Cheng CL, Van Ness JW. Statistical regression with measurement error. London: Arnold; 1999.

88. Clement S, Ibrahim S, Crichton N, Wolf M, Rowlands G. Complex interventions to improve the health of people with limited literacy: a systematic review. Patient Educ Couns 2009; 75(3): 340-351.

89. Cochrane Collaboration. Our principles [online]. 22 June 2012 [accessed: 23 April 2013]. URL: http://www.cochrane.org./about-us/our-principles.

90. Cochrane Effective Practice and Organisation of Care Review Group. The data collection checklist [online]. June 2002 [accessed: 23 April 2013]. URL: http://epoc.cochrane.org/sites/epoc.cochrane.org/files/uploads/datacollectionchecklist.pdf.

91. Commission of the European Communities. eEurope 2002: quality criteria for health related websites [online]. 29 November 2002 [accessed: 4 March 2011]. URL: http://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=COM:2002:0667:FIN:EN:PDF.

92. Committee of Experts on Management of Safety and Quality in Health Care, Expert Group on Safe Medication Practices. Glossary of terms related to patient and medication safety [online]. 20 October 2005 [accessed: 19 April 2013]. URL: http://www.who.int/patientsafety/highlights/COE_patient_and_medication_safety_gl.pdf.

93. Corbin JM, Strauss AL. Weiterleben lernen: Verlauf und Bewältigung chronischer Krankheit. Bern: Huber; 2003.

94. Cordoba G, Schwartz L, Woloshin S, Bae H, Gøtzsche PC. Definition, reporting, and interpretation of composite outcomes in clinical trials: systematic review. BMJ 2010; 341: c3920.

95. Cornell JE, Laine C. The science and art of deduction: complex systematic overviews. Ann Intern Med 2008; 148(10): 786-788.

96. Coulter A. Evidence based patient information is important, so there needs to be a national strategy to ensure it. BMJ 1998; 317(7153): 225-226.

97. Coulter A, Ellins J, Swain D, Clarke A, Heron P, Rasul F et al. Assessing the quality of information to support people in making decisions about their health and healthcare. Oxford: Picker Institute Europe; 2006. URL: http://www.pickereurope.org/assets/content/pdf/Project_Reports/Health-information-quality-web-version-FINAL.pdf.

98. Cui L, Hung HMJ, Wang SJ, Tsong Y. Issues related to subgroup analysis in clinical trials. J Biopharm Stat 2002; 12(3): 347-358.

99. Cullen R. Health information on the internet: a study of providers, quality, and users. Westport: Praeger; 2006.

100. D'Agostino RB Sr, Massaro JM, Sullivan LM. Non-inferiority trials: design concepts and issues; the encounters of academic consultants in statistics. Stat Med 2003; 22(2): 169-186.

101. Da Costa BR, Rutjes AWS, Johnston BC, Reichenbach S, Nüesch E, Tonia T et al. Methods to convert continuous outcomes into odds ratios of treatment response and numbers needed to treat: meta-epidemiological study. Int J Epidemiol 2012; 41(5): 1445-1459.

102. Daly J, Willis K, Small R, Green J, Welch N, Kealy M et al. A hierarchy of evidence for assessing qualitative health research. J Clin Epidemiol 2007; 60(1): 43-49.

103. Dans AL, Dans LF, Guyatt GH, Richardson S. Users' guides to the medical literature; XIV: how to decide on the applicability of clinical trial results to your patient. JAMA 1998; 279(7): 545-549.

104. Dans LF, Silvestre MA, Dans AL. Trade-off between benefit and harm is crucial in health screening recommendations; part I: general principles. J Clin Epidemiol 2011; 64(3): 231-239.

105. David M, Borde T, Kentenich H. Knowledge among German and Turkish women about specifically female bodily functions, contraception, preventative medical examinations and menopause. Ethn Health 2000; 5(2): 101-112.

106. De Jonchere K, Gartlehner G, Goologly L, Mustajoki P, Permanand G. Gesundheitsinformationen für Patienten und die Öffentlichkeit: zusammengestellt vom Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen; Gutachten der Weltgesundheitsorganisation 2008/2009 [online]. 2010 [accessed: 23 April 2013]. URL: http://www.euro.who.int/__data/assets/pdf_file/0004/94990/E93735g.pdf.

107. Deeks JJ. Systematic reviews in health care: systematic reviews of evaluations of diagnostic and screening tests. BMJ 2001; 323(7305): 157-162.

108. Deeks JJ. Issues in the selection of a summary statistic for meta-analysis of clinical trials with binary outcomes. Stat Med 2002; 21(11): 1575-1600.

109. Deeks JJ, Higgins JPT, Altman DG. Analysing data and undertaking meta-analyses. In: Higgins JPT, Green S (Ed). Cochrane handbook for systematic reviews of interventions. Chichester: Wiley; 2008. p. 243-296.

110. Delnoij DM, Groenewegen PP. Health services and systems research in Europe: overview of the literature 1995-2005. Eur J Public Health 2007; 17(Suppl 1): 10-13.

111. Derksen S, Keselman HJ. Backward, forward, and stepwise automated subset selection algorithms: frequency of obtaining authentic and noise variables. Br J Math Stat Psychol 1992; 45(2): 265-282.

112. Derry S, Loke YK, Aronson JK. Incomplete evidence: the inadequacy of databases in tracing published adverse drug reactions in clinical trials. BMC Med Res Methodol 2001; 1: 7.

113. Des Jarlais DC, Lyles C, Crepaz N. Improving the reporting quality of nonrandomized evaluations of behavioral and public health interventions: the TREND statement. Am J Public Health 2004; 94(3): 361-366.

114. Desu MM, Raghavarao D. Sample size methodology. Boston: Academic Press; 1990.

115. Detmer DE. Building the national health information infrastructure for personal health, health care services, public health, and research. BMC Med Inform Decis Mak 2003; 3: 1.

116. Deutsche Gesellschaft für Epidemiologie. Leitlinien und Empfehlungen zur Sicherung von guter epidemiologischer Praxis (GEP): Langversion [online]. March 2008 [accessed: 23 April 2013]. URL: http://www.gmds.de/pdf/publikationen/stellungnahmen/stell_gep_ergaenzung.pdf.

117. Deutsche Rentenversicherung Bund (Ed). Rentenversicherung in Zeitreihen: Ausgabe 2012. Berlin: DRV; 2008. (DRV-Schriften; Volume 22). URL: http://www.deutsche-rentenversicherung.de/cae/servlet/contentblob/238700/publicationFile/50912/rv_in_zeitreihen.pdf.

118. Deutscher Bundestag. Gesetzentwurf der Fraktionen SPD, CDU/CSU und BÜNDNIS 90/DIE GRÜNEN: Entwurf eines Gesetzes zur Modernisierung der gesetzlichen Krankenversicherung (GKV-Modernisierungsgesetz - GMG); BT Drucksache 15/1525 [online]. 8 September 2003 [accessed: 16 January 2013]. URL: http://dipbt.bundestag.de/doc/btd/15/015/1501525.pdf.

119. Deutscher Bundestag. Gesetz zur Modernisierung der gesetzlichen Krankenversicherung (GKV-Modernisierungsgesetz - GMG). Bundesgesetzblatt Teil 1 2003; (55): 2190-2258.

120. Deutscher Bundestag. Gesetz zur Neuordnung des Arzneimittelmarktes in der gesetzlichen Krankenversicherung (Arzneimittelmarktneuordnungsgesetz – AMNOG) vom 22. Dezember 2010. Bundesgesetzblatt Teil 1 2010; (67): 2262-2277.

121. Deutscher Ethikrat (Ed). Nutzen und Kosten im Gesundheitswesen: zur normativen Funktion ihrer Bewertung; Stellungnahme. Berlin: Deutscher Ethikrat; 2011. URL: http://www.ethikrat.org/dateien/pdf/stellungnahme-nutzen-und-kosten-im-gesundheitswesen.pdf.

122. Deutsches Institut für Normung. Klinische Prüfung von Medizinprodukten an Menschen: gute klinische Praxis (ISO 14155:2011 + Cor. 1:2011); deutsche Fassung EN ISO 14155:2011 + AC:2011. Berlin: Beuth; 2012.

123. Deutsches Netzwerk Evidenzbasierte Medizin. Die "Gute Praxis Gesundheitsinformation". Z Evid Fortbild Qual Gesundhwes 2010; 104(1): 66-68.

124. Devillé WL, Buntinx F, Bouter LM, Montori VM, De Vet HCW, Van der Windt DAWM et al. Conducting systematic reviews of diagnostic studies: didactic guidelines. BMC Med Res Methodol 2002; 2: 9.

125. Di Pietrantonj C. Four-fold table cell frequencies imputation in meta analysis. Stat Med 2006; 25(13): 2299-2322.

126. Dias S, Welton NJ, Caldwell DM, Ades AE. Checking consistency in mixed treatment comparison meta-analysis. Stat Med 2010; 29(7-8): 932-944.

127. Dierks ML, Seidel G, Horch K, Schwartz FW. Bürger- und Patientenorientierung im Gesundheitswesen. Berlin: Robert Koch-Institut; 2006. (Gesundheitsberichterstattung des Bundes; Volume 32). URL: http://www.rki.de/DE/Content/Gesundheitsmonitoring/Gesundheitsberichterstattung/GBEDownloadsT/buergerorientierung.pdf?__blob=publicationFile.

128. DIPEx. Healthtalkonline [online]. [Accessed: 19 April 2013]. URL: http://www.healthtalkonline.org.

129. Dixon-Woods M. Writing wrongs? An analysis of published discourses about the use of patient information leaflets. Soc Sci Med 2001; 52(9): 1417-1432.

130. Dixon-Woods M, Agarwal S, Young B, Jones D, Sutton A. Integrative approaches to qualitative and quantitative evidence. London: Health Development Agency; 2004. URL: http://www.nice.org.uk/niceMedia/pdf/Integrative_approaches_evidence.pdf.

131. Dixon-Woods M, Fitzpatrick R. Qualitative research in systematic reviews: has established a place for itself. BMJ 2001; 323(7316): 765-766.

132. Dixon-Woods M, Fitzpatrick R, Roberts K. Including qualitative research in systematic reviews: opportunities and problems. J Eval Clin Pract 2001; 7(2): 125-133.

133. Dixon-Woods M, Sutton A, Shaw R, Miller T, Smith J, Young B et al. Appraising qualitative research for inclusion in systematic reviews: a quantitative and qualitative comparison of three methods. J Health Serv Res Policy 2007; 12(1): 42-47.

134. Djulbegovic B, Kumar A, Soares HP, Hozo I, Bepler G, Clarke M et al. Treatment success in cancer: new cancer treatment successes identified in phase 3 randomized controlled trials conducted by the National Cancer Institute-sponsored cooperative oncology groups, 1955 to 2006. Arch Intern Med 2008; 168(6): 632-642.

135. Dolan JG. Shared decision-making: transferring research into practice; the Analytic Hierarchy Process (AHP). Patient Educ Couns 2008; 73(3): 418-425.

136. Dolan JG, Isselhardt BJ Jr, Cappuccio JD. The Analytic Hierarchy Process in medical decision making: a tutorial. Med Decis Making 1989; 9(1): 40-50.

137. Dolan P, Shaw R, Tsuchiya A, Williams A. QALY maximisation and people's preferences: a methodological review of the literature. Health Econ 2005; 14(2): 197-208.

138. Donner A, Klar J. Design and analysis of cluster randomization trials in health research. London: Arnold; 2000.

139. Draborg E, Gyrd-Hansen D, Poulsen PB, Horder M. International comparison of the definition and the practical application of health technology assessment. Int J Technol Assess Health Care 2005; 21(1): 89-95.

140. Drazen JM. COX-2 inhibitors: a lesson in unexpected problems. N Engl J Med 2005; 352(11): 1131-1132.

141. Drummond MF, Sculpher MJ, Torrance GW, O'Brian BJ, Stoddart GL. Methods for the economic evaluation of health care programmes. Oxford: Oxford University Press; 2005.

142. Dwan K, Altman DG, Arnaiz JA, Bloom J, Chan AW, Cronin E et al. Systematic review of the empirical evidence of study publication bias and outcome reporting bias. PLoS One 2008; 3(8): e3081.

143. Ebell MH, Siwek J, Weiss BD, Woolf SH, Susman J, Ewigman B et al. Strength of Recommendation Taxonomy (SORT): a patient-centered approach to grading evidence in the medical literature. Am Fam Physician 2004; 69(3): 548-556.

144. Eccles MP, Rousseau N, Freemantle N. Updating evidence-based clinical guidelines. J Health Serv Res Policy 2002; 7(2): 98-103.

145. Edwards A, Elwyn G (Ed). Shared decision-making in health care: achieving evidence-based patient choice. Oxford: Oxford University Press; 2009.

146. Edwards AGK, Bastian H. Risk communication: making evidence part of patient choices? In: Edwards AGK, Elwyn GJ (Ed). Evidence-based patient choice: inevitable or impossible? Oxford: Oxford University Press; 2001. p. 144-160.

147. Edwards AGK, Elwyn GJ, Mulley A. Explaining risks: turning numerical data into meaningful pictures. BMJ 2002; 324(7341): 827-830.

148. Edwards AGK, Evans R, Dundon J, Haigh S, Hood K, Elwyn GJ. Personalised risk communication for informed decision making about taking screening tests. Cochrane Database Syst Rev 2006; (4): CD001865.

149. Edwards JE, McQuay HJ, Moore RA, Collins SL. Reporting of adverse effects in clinical trials should be improved: lessons from acute postoperative pain. J Pain Symptom Manage 1999; 18(6): 427-437.

150. Egger M, Davey Smith G, Altman DG (Ed). Systematic reviews in health care: meta-analysis in context. London: BMJ Publishing Group; 2001.

151. Egger M, Davey Smith G, Schneider M, Minder C. Bias in meta-analysis detected by a simple, graphical test. BMJ 1997; 315(7109): 629-634.

152. Egger M, Jüni P, Bartlett C, Holenstein F, Sterne J. How important are comprehensive literature searches and the assessment of trial quality in systematic reviews? Empirical study. Health Technol Assess 2003; 7(1): 1-76.

153. Elwyn GJ, O'Connor A, Stacey D, Volk R, Edwards AGK, Coulter A et al. Developing a quality criteria framework for patient decision aids: online international Delphi consensus process. BMJ 2006; 333(7565): 417-424.

154. Engels EA, Schmid CH, Terrin N, Olkin I, Lau J. Heterogeneity and statistical significance in meta-analysis: an empirical study of 125 meta-analyses. Stat Med 2000; 19(13-35): 1707-1728.

155. Entwistle VA, Sheldon TA, Sowden A, Watt IS. Evidence-informed patient choice: practical issues of involving patients in decisions about health care technologies. Int J Technol Assess Health Care 1998; 14(2): 212-225.

156. Epstein RM, Alper BS, Quill TE. Communicating evidence for participatory decision making. JAMA 2004; 291(19): 2359-2366.

157. Europäisches Parlament, Rat der Europäischen Union. Verordnung (EG) Nr. 141/2000 des Europäischen Parlaments und des Rates vom 16. Dezember 1999 über Arzneimittel für seltene Leiden. Amtsblatt der Europäischen Gemeinschaften 2000; 43(L18): 1-5.

158. European Medicines Agency. Points to consider on application with: 1. meta-analyses; 2. one pivotal study [online]. 31 May 2001 [accessed: 22 September 2010]. URL: http://www.ema.europa.eu/docs/en_GB/document_library/Scientific_guideline/2009/09/WC5 00003657.pdf.

159. European Medicines Agency. Guideline on the choice of the non-inferiority margin [online]. 27 July 2005 [accessed: 23 April 2013]. URL: http://www.ema.europa.eu/docs/en_GB/document_library/Scientific_guideline/2009/09/WC5 00003636.pdf.

160. European Medicines Agency. Reflection paper on the regulatory guidance for the use of Health Related Quality of Life (HRQL) measures in the evaluation of medicinal products [online]. 27 July 2005 [accessed: 23 April 2013]. URL: http://www.ema.europa.eu/docs/en_GB/document_library/Scientific_guideline/2009/09/WC5 00003637.pdf.

161. European Medicines Agency. Guideline on clinical trials in small populations [online]. 27 July 2006 [accessed: 23 April 2013]. URL: http://www.ema.europa.eu/docs/en_GB/document_library/Scientific_guideline/2009/09/WC5 00003615.pdf.

162. European Medicines Agency. Guideline on clinical investigation of medicinal products in the treatment of diabetes mellitus: draft [online]. 20 January 2010 [accessed: 23 April 2013]. URL:
http://www.ema.europa.eu/docs/en_GB/document_library/Scientific_guideline/2010/02/WC5 00073570.pdf.

163. European Medicines Agency. Guideline on missing data in confirmatory clinical trials [online]. 2 July 2010 [accessed: 23 April 2013]. URL:
http://www.ema.europa.eu/docs/en_GB/document_library/Scientific_guideline/2010/09/WC5 00096793.pdf.

164. Evidence-Based Medicine Working Group. Evidence-based medicine: a new approach to teaching the practice of medicine. JAMA 1992; 268(17): 2420-2425.

165. Eyding D, Lelgemann M, Grouven U, Harter M, Kromp M, Kaiser T et al. Reboxetine for acute treatment of major depression: systematic review and meta-analysis of published and unpublished placebo and selective serotonin reuptake inhibitor controlled trials. BMJ 2010; 341: c4737.

166. Eysenbach G. Recent advances: consumer health informatics. BMJ 2000; 320(7251): 1713-1716.

167. Eysenbach G, Köhler C. How do consumers search for and appraise health information on the world wide web? Qualitative study using focus groups, usability tests, and in-depth interviews. BMJ 2002; 324(7337): 573-577.

168. Farrington CP, Manning G. Test statistics and sample size formulae for comparative binomial trials with null hypothesis of non-zero risk difference or non-unity relative risk. Stat Med 1990; 9(12): 1447-1454.

169. Fayers P, Machin D. Quality of life: the assessment, analysis and interpretation of patient-reported outcomes. Chichester: Wiley; 2007.

170. Feinstein AR. Clinical epidemiology: the architecture of clinical research. Philadelphia: Saunders; 1985.

171. Feinstein AR. Invidious comparisons and unmet clinical challenges. Am J Med 1992; 92(2): 117-120.

172. Feise RJ. Do multiple outcome measures require p-value adjustment? BMC Med Res Methodol 2002; 2: 8.

173. Feldman-Stewart D, Brennenstuhl S, Brundage MD. A purpose-based evaluation of information for patients: an approach to measuring effectiveness. Patient Educ Couns 2007; 65(3): 311-319.

174. Fergusson D, Laupacis A, Salmi LR, McAlister FA, Huet C. What should be included in meta-analyses? An exploration of methodological issues using the ISPOT meta-analyses. Int J Technol Assess Health Care 2000; 16(4): 1109-1119.

175. Ferreira-Gonzáles I, Busse JW, Heels-Ansdell D, Montori VM, Alk EA, Byrant DM et al. Problems with use of composite end points in cardiocascular trials: systematic review of randomized controlled trials. BMJ 2007; 334(7597): 786-792.

176. Fessler J, Fischer J, Franzen D, Geraedts M, Graf HJ, Kroegel C et al. Leitlinien-Clearingbericht "COPD": Leitlinien-Clearingverfahren von Bundesärztekammer und Kassenärztlicher Bundesvereinigung in Kooperation mit Deutscher Krankenhausgesellschaft, Spitzenverbänden der Krankenkassen und Gesetzlicher Rentenversicherung. Niebüll: Videel; 2003. (ÄZQ-Schriftenreihe; Volume 14). URL: http://www.leitlinien.de/mdb/edocs/pdf/schriftenreihe/schriftenreihe14.pdf.

177. Field MJ, Lohr KN (Ed). Clinical practice guidelines: directions for a new program. Washington: National Academy Press; 1990.

178. Fleiss JL, Tytun A, Ury HK. A simple approximation for calculating sample sizes for comparing independent proportions. Biometrics 1980; 36(2): 343-346.

179. Fleming TR. Surrogate endpoints and FDA's accelerated approval process. Health Aff (Millwood) 2005; 24(1): 67-78.

180. Fleming TR, DeMets DL. Surrogate end points in clinical trials: are we being misled? Ann Intern Med 1996; 125(7): 605-613.

181. Fletcher RH, Fletcher SW. Klinische Epidemiologie: Grundlagen und Anwendung. Bern: Huber; 2007.

182. Food and Drug Administration. Guidance for industry: developing medical imaging drug and biological products; part 2: clinical indications [online]. June 2004 [accessed: 23 April 2013]. URL: http://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/ucm071603.pdf.

183. Food and Drug Administration. Guidance for industry: patient-reported outcome measures; use in medical product development to support labeling claims [online]. December 2009 [accessed: 23 April 2013]. URL: http://www.fda.gov/downloads/Drugs/GuidanceComplianceRegulatoryInformation/Guidances/UCM193282.pdf.

184. Førde OH. Is imposing risk awareness cultural imperialism? Soc Sci Med 1998; 47(9): 1155-1159.

185. Fransen GA, Van Marrewijk CJ, Mujakovic S, Muris JW, Laheij RJ, Numans ME et al. Pragmatic trials in primary care: methodological challenges and solutions demonstrated by the DIAMOND-study. BMC Med Res Methodol 2007; 7: 16.

186. Freemantle N. Interpreting the results of secondary end points and subgroup analyses in clinical trials: should we lock the crazy aunt in the attic? BMJ 2001; 322(7292): 989-991.

187. Freemantle N, Blonde L, Bolinder B, Gerber RA, Hobbs FD, Martinez L et al. Real-world trials to answer real-world questions. Pharmacoeconomics 2005; 23(8): 747-754.

188. Freemantle N, Calvert M. Weighing the pros and cons for composite outcomes in clinical trials. J Clin Epidemiol 2007; 60(7): 658-659.

189. French SD, McDonald S, McKenzie JE, Green SE. Investing in updating: how do conclusions change when Cochrane systematic reviews are updated? BMC Med Res Methodol 2005; 5: 33.

190. Fryback DG, Thornbury JR. The efficacy of diagnostic imaging. Med Decis Making 1991; 11(2): 88-94.

191. Fujimori M, Uchitomi Y. Preferences of cancer patients regarding communication of bad news: a systematic literature review. Jpn J Clin Oncol 2009; 39(4): 201-216.

192. Furukawa TA, Guyatt GH, Griffith LE. Can we individualize the 'number needed to treat'? An empirical study of summary effect measures in meta-analyses. Int J Epidemiol 2002; 31(1): 72-76.

193. Gagliardi A, Jadad AR. Examination of instruments used to rate quality of health information on the internet: chronicle of a voyage with an unclear destination. BMJ 2002; 324(7337): 569-573.

194. Garcia-Castillo D, Fetters MD. Quality in medical translations: a review. J Health Care Poor Underserved 2007; 18(1): 74-84.

195. Garrison LP Jr, Towse A, Bresnahan BW. Assessing a structured, quantitative health outcomes approach to drug risk-benefit analysis. Health Aff (Millwood) 2007; 26(3): 684-695.

196. Gartlehner G, Moore CG. Direct versus indirect comparisons: a summary of the evidence. Int J Technol Assess Health Care 2008; 24(2): 170-177.

197. Gemeinsamer Bundesausschuss. Beschluss des Gemeinsamen Bundesausschusses über die Anpassung der Beauftragung des Instituts für Qualität und Wirtschaftlichkeit im Gesundheitswesen zur Erstellung von Patienteninformationen [online]. 13 March 2008 [accessed: 23 April 2013]. URL: http://www.g-ba.de/downloads/39-261-650/2008-03-13-IQWiG-Anpassung-Generalauftrag.pdf.

198. Gemeinsamer Bundesausschuss. Verfahrensordnung des Gemeinsamen Bundesausschusses [online]. 6 December 2012 [accessed: 15 April 2013]. URL: http://www.g-ba.de/downloads/62-492-667/VerfO_2012-12-06.pdf.

199. Gemoets D, Rosemblat G, Tse T, Logan R. Assessing readability of consumer health information: an exploratory study. Medinfo 2004; 11(Pt 2): 869-873.

200. Gerhardt U. Patientenkarrieren. Frankfurt am Main: Suhrkamp; 1986.

201. Gesellschaft für Evaluation. Standards für Evaluation. Mainz: DeGEval; 2008. URL: http://www.degeval.de/images/stories/Publikationen/DeGEval_-_Standards.pdf.

202. Glasziou PP, Chalmers I, Rawlins M, McCulloch P. When are randomised trials unnecessary? Picking signal from noise. BMJ 2007; 334(7589): 349-351.

203. Glasziou PP, Sanders SL. Investigating causes of heterogeneity in systematic reviews. Stat Med 2002; 21(11): 1503-1511.

204. Glasziou PP, Vandenbroucke JP, Chalmers I. Assessing the quality of research. BMJ 2004; 328(7430): 39-41.

205. Glenton C, Nilsen ES, Carlsen B. Lay perceptions of evidence-based information: a qualitative evaluation of a website for back pain sufferers. BMC Health Serv Res 2006; 6: 34.

206. Godwin M, Ruhland L, Casson I, MacDonald S, Delva D, Birtwhistle R et al. Pragmatic controlled clinical trials in primary care: the struggle between external and internal validity. BMC Med Res Methodol 2003; 3: 28.

207. Gøtzsche PC, Liberati A, Torri V, Rossetti L. Beware of surrogate outcome measures. Int J Technol Assess Health Care 1996; 12(2): 238-246.

208. Graham RM, Mancher M, Miller-Wolman D, Greenfield S, Steinberg E. Clinical practice guidelines we can trust. Washington: National Academies Press; 2011. URL: http://www.awmf.org/fileadmin/user_upload/Leitlinien/International/IOM_CPG_lang_2011.pdf.

209. Gray JAM. How to get better value healthcare. Oxford: Offox Press; 2007.

210. Greenhalgh T, Hurwitz B. Narrative based medicine: why study narrative? BMJ 1999; 318(7175): 48-50.

211. Greenland S, Morgenstern H. Ecological bias, confounding, and effect modification. Int J Epidemiol 1989; 18(1): 269-274.

212. Griffiths F, Lindenmeyer A, Powell J, Lowe P, Thorogood M. Why are health care interventions delivered over the internet? A systematic review of the published literature. J Med Internet Res 2006; 8(2): e10-e36.

213. Grilli R, Ramsay C, Minozzi S. Mass media interventions: effects on health services utilisation. Cochrane Database Syst Rev 2002; (1): CD000389.

214. Grimes DA, Schulz K. An overview of clinical research: the lay of the land. Lancet 2002; 359(9300): 57-61.

215. Grimes DA, Schulz KF. Surrogate end points in clinical research: hazardous to your health. Obstet Gynecol 2005; 105(5 Pt 1): 1114-1118.

216. Grouin JM, Coste M, Lewis J. Subgroup analyses in randomized clinical trials: statistical and regulatory issues. J Biopharm Stat 2005; 15(5): 869-882.

217. Guddat C, Grouven U, Bender R, Skipka G. A note on the graphical presentation of prediction intervals in random-effects meta-analyses. Systematic Reviews 2012; 1: 34.

218. Guyatt GH. Evidence-based medicine. ACP J Club 1991; 14(Suppl 2): A16.

219. Guyatt GH, Jaeschke R, Roberts R. N-of-1 randomized clinical trials in pharmacoepidemiology. In: Strom BL (Ed). Pharmacoepidemiology. Chichester: Wiley; 2005. p. 665-680.

220. Guyatt GH, Oxman AD, Kunz R, Falck-Ytter Y, Vist GE, Liberati A et al. Going from evidence to recommendations. BMJ 2008; 336(7652): 1049-1051.

221. Guyatt GH, Oxman AD, Vist GE, Kunz R, Falck-Ytter Y, Alonso-Coello P et al. GRADE: an emerging consensus on rating quality of evidence and strength of recommendations. BMJ 2008; 336(7650): 924-926.

222. Guyatt GH, Rennie D. Users' guides to the medical literature: a manual for evidence-based clinical practice. Chicago: American Medical Association; 2002.

223. Guyatt GH, Sackett DL, Sinclair JC, Hayward R, Cook DJ, Cook RJ. Users' guides to the medical literature; IX: a method for grading health care recommendations. JAMA 1995; 274(22): 1800-1804.

224. Guyatt GH, Sackett DL, Taylor DW, Chong J, Roberts R, Pugsley S. Determining optimal therapy: randomized trials in individual patients. N Engl J Med 1986; 314(14): 889-892.

225. Guyatt GH, Tugwell P, Feeny DH, Drummond MF, Haynes RB. The role of before-after studies of therapeutic impact in the evaluation of diagnostic technologies. J Chronic Dis 1986; 39(4): 295-304.

226. Hamza TH, Van Houwelingen HC, Heijenbrok-Kal MH, Stijnen T. Associating explanatory variables with summary receiver operating characteristic curves in diagnostic meta-analysis. J Clin Epidemiol 2009; 62(12): 1284-1291.

227. Harbord RM, Whiting P, Sterne JA, Egger M, Deeks JJ, Shang A et al. An empirical comparison of methods for meta-analysis of diagnostic accuracy showed hierarchical models are necessary. J Clin Epidemiol 2008; 61(11): 1095-1103.

228. Harbour RT, Miller J. A new system for grading recommendations in evidence based guidelines. BMJ 2001; 323(7308): 334-336.

229. Harden A, Garcia J, Oliver S, Rees R, Shepherd J, Brunton G et al. Applying systematic review methods to studies of people's views: an example from public health research. J Epidemiol Community Health 2004; 58(9): 794-800.

230. Harrell FE. Regression modeling strategies: with applications to linear models, logistic regression, and survival analysis. New York: Springer; 2001.

231. Harrell FE, Lee KL, Mark DB. Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. Stat Med 1996; 15(4): 361-387.

232. Harris J. QALYfying the value of life. J Med Ethics 1987; 13(3): 117-123.

233. Harris RP, Helfand M, Woolf SH, Lohr KN, Mulrow CD, Teutsch SM et al. Current methods of the US Preventive Services Task Force: a review of the process. Am J Prev Med 2001; 20(3 Suppl): 21-35.

234. Hart D (Ed). Ärztliche Leitlinien im Medizin- und Gesundheitsrecht: Recht und Empirie professioneller Normbildung. Baden-Baden: Nomos; 2005. (Gesundheitsrecht und Gesundheitswissenschaften; Volume 9).

235. Harteloh P. The meaning of quality in health care: a conceptual analysis. Health Care Anal 2003; 11(3): 259-267.

236. Hayden JA, Cote P, Bombardier C. Evaluation of the quality of prognosis studies in systematic reviews. Ann Intern Med 2006; 144(6): 427-437.

237. Haynes RB. Forming research questions. J Clin Epidemiol 2006; 59(9): 881-886.

238. Haynes RB, Cotoi C, Holland J, Walters L, Wilczynski N, Jedraszewski D et al. Second-order peer review of the medical literature for clinical practitioners. JAMA 2006; 295(15): 1801-1808.

239. Haynes RB, Devereaux PJ, Guyatt GH. Clinical expertise in the era of evidence-based medicine and patient choice. Evid Based Med 2002; 7(2): 36-38.

240. Health on the Net Foundation. Der HONcode in Kürze [online]. April 1997 [accessed: 19 April 2013]. URL: http://www.hon.ch/HONcode/German.

241. Hemingway H, Riley RD, Altman DG. Ten steps towards improving prognosis research. BMJ 2009; 339: b4184.

242. Hennekens CH, Buring JE. Epidemiology in medicine. Boston: Little, Brown and Company; 1987.

243. Herxheimer A, McPherson A, Miller R, Shepperd S, Yaphe J, Ziebland S. Database of Patients' Experiences (DIPEx): a multi-media approach to sharing experiences and information. Lancet 2000; 355(9214): 1540-1543.

244. Herxheimer A, Ziebland S. DIPEx: fresh insights for medical practice. J R Soc Med 2003; 96(5): 209-210.

245. Hicks NJ. Evidence-based health care. Bandolier 1997; 4(5): 8.

246. Higgins JP, Thompson SG, Spiegelhalter DJ. A re-evaluation of random-effects meta-analysis. J R Stat Soc A 2009; 172(1): 137-159.

247. Higgins JPT, Altman DG. Assessing risk of bias in included studies. In: Higgins JPT, Green S (Ed). Cochrane handbook for systematic reviews of interventions. Chichester: Wiley; 2008. p. 187-242.

248. Higgins JPT, Green S (Ed). Cochrane handbook for systematic reviews of interventions. Chichester: Wiley; 2008.

249. Higgins JPT, Thompson SG. Quantifying heterogeneity in a meta-analysis. Stat Med 2002; 21(11): 1539-1558.

250. Higgins JPT, Thompson SG, Deeks JJ, Altman DG. Measuring inconsistency in meta-analyses. BMJ 2003; 327(7414): 557-560.

251. Hill AB (Ed). Controlled clinical trials. Oxford: Blackwell; 1960.

252. Hirschberg I. Bewertung und Wirkung von evidenzbasierten Gesundheitsinformationen: die Perspektive der Nutzer; Public Health; Magisterarbeit. München: GRIN Verlag; 2010. (Patientenorientierung und Gesundheitskompetenz; Volume 1).

253. Hirsh J, Guyatt G. Clinical experts or methodologists to write clinical guidelines? Lancet 2009; 374(9686): 273-275.

254. Holmes-Rovner M. International Patient Decision Aid Standards (IPDAS): beyond decision aids to usual design of patient education materials. Health Expect 2007; 10(2): 103-107.

255. Hopewell S, Clarke M, Lefebvre C, Scherer R. Handsearching versus electronic searching to identify reports of randomized trials. Cochrane Database Syst Rev 2007; (2): MR000001.

256. Hopewell S, McDonald S, Clarke M, Egger M. Grey literature in meta-analyses of randomized trials of health care interventions. Cochrane Database Syst Rev 2007; (2): MR000010.

257. Hosmer DW, Taber S, Lemeshow S. The importance of assessing the fit of logistic regression models: a case study. Am J Public Health 1991; 81(12): 1630-1635.

258. Houts PS, Doak CC, Doak LG, Loscalzo MJ. The role of pictures in improving health communication: a review of research on attention, comprehension, recall, and adherence. Patient Educ Couns 2006; 61(2): 173-190.

259. Hozo SP, Djulbegovic B, Hozo I. Estimating the mean and variance from the median, range, and the size of a sample. BMC Med Res Methodol 2005; 5: 13.

260. Hummel JM, IJzerman MJ. The use oft the Analytic Hierarchy Process in health care decision making. Enschede: University of Twente; 2009.

261. Hung HM, O'Neill RT, Bauer P, Köhne K. The behavior of the P-value when the alternative hypothesis is true. Biometrics 1997; 53(1): 11-22.

262. Hussain T, Michel G, Shiffman RN. The Yale Guideline Recommendation Corpus: a representative sample of the knowledge content of guidelines. Int J Med Inf 2009; 78(5): 354-363.

263. ICH Expert Working Group. ICH harmonised tripartite guideline: the extent of population exposure to assess clinical safety for drugs intended for long-term treatment of non-life-threatening conditions; E1; current step 4 version [online]. 27 October 1994 [accessed: 19 April 2013]. URL: http://www.ich.org/fileadmin/Public_Web_Site/ICH_Products/Guidelines/Efficacy/E1/Step4/ E1_Guideline.pdf.

264. Inan H. Measuring the success of your website: a customer-centric approach to website management. Frenchs Forest: Pearson Education Australia; 2002.

265. Informationsgemeinschaft zur Feststellung der Verbreitung von Werbeträgern. Definitionen und technische Erläuterungen: Version 1.9; Anlage 1 zu den IVW-Richtlinien für Online-Angebote [online]. 14 November 2005 [accessed: 19 April 2013]. URL: http://daten.ivw.eu/download/pdf/Online_RichtlinienV1_9_Anlage1.pdf.

266. Initi@tive D[21]. (N)Onliner Atlas 2009: eine Topographie des digitalen Grabens durch Deutschland; Nutzung und Nichtnutzung des Internets, Strukturen und regionale Verteilung [online]. June 2009 [accessed: 23 April 2013]. URL: http://www.initiatived21.de/wp-content/uploads/2009/06/NONLINER2009.pdf.

267. Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen. Allgemeine Methoden zur Bewertung von Verhältnissen zwischen Nutzen und Kosten [online]. 12 October 2009 [accessed: 4 November 2013]. URL: https://www.iqwig.de/download/Methodik_fuer_die_Bewertung_von_Verhaeltnissen_zwisch en_Kosten_und_Nutzen.pdf.

268. Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen. Aussagekraft von Surrogatendpunkten in der Onkologie: Rapid Report; Auftrag A10-05; Version 1.1 [online]. 21 November 2011 [accessed: 23 April 2013]. (IQWiG-Berichte; Volume 80). URL: https://www.iqwig.de/download/A10-05_Rapid_Report_Version_1-1_Surrogatendpunkte_in_der_Onkologie.pdf.

269. Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen. Nutzungsbedingungen [online]. 13 September 2012 [accessed: 6 May 2013]. URL: http://www.gesundheitsinformation.de/nutzungsbedingungen.51.de.html.

270. Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen. Vereinbarung über die vertrauliche Behandlung von Unterlagen [online]. 19 August 2005 [accessed: 19 April 2013]. URL: http://www.iqwig.de/download/IQWiG-VFA-Mustervertrag.pdf.

271. Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen. Bupropion, Mirtazapin und Reboxetin bei der Behandlung von Depression: Abschlussbericht; Auftrag A05-20C [online]. 9 November 2009 [accessed: 23 April 2013]. (IQWiG-Berichte; Volume 68). URL: https://www.iqwig.de/download/A05-20C_Abschlussbericht_Bupropion_Mirtazapin_und_Reboxetin_bei_Depressionen.pdf.

272. Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen. Ticagrelor: Nutzenbewertung gemäß § 35a SGB V; Dossierbewertung; Auftrag A11-02 [online]. 29 September 2011 [accessed: 4 November 2013]. (IQWiG-Berichte; Volume 96). URL: https://www.iqwig.de/de/projekte_ergebnisse/projekte/arzneimittelbewertung/a11_02_ticagrelor_nutzenbewertung_gemass_35a_sgb_v_dossierbewertung.1425.html.

273. Institute of Medicine. Crossing the quality chasm: a new health system for the 21st century. Washington: National Academy Press; 2001. URL: http://books.nap.edu/openbook.php?record_id=10027.

274. International Conference on Harmonisation Expert Working Group. ICH harmonised tripartite guideline: statistical principles for clinical trials; an introductory note on an international guideline. Stat Med 1999; 18(15): 1905-1942.

275. International Conference on Harmonisation of Technical Requirements for Registration of Pharmaceuticals for Human Use. Website [online]. [Accessed: 7 May 2013]. URL: http://www.ich.org.

276. Ioannidis JPA. Contradicted and initially stronger effects in highly cited clinical research. JAMA 2005; 294(2): 218-228.

277. Ioannidis JPA, Evans S, Gøtzsche PC, O'Neill RT, Altman DG, Schulz K et al. Better reporting of harms in randomized trials: an extension of the CONSORT statement. Ann Intern Med 2004; 141(10): 781-788.

278. Ioannidis JPA, Lau J. Completeness of safety reporting in randomized trials: an evaluation of 7 medical areas. JAMA 2001; 285(4): 437-443.

279. Ioannidis JPA, Mulrow CD, Goodman SN. Adverse events: the more you search, the more you find. Ann Intern Med 2006; 144(4): 298-300.

280. Irmen L, Linner U. Die Repräsentation generisch maskuliner Personenbezeichnungen: eine theoretische Integration bisheriger Befunde. Z Psychol 2005; 213(3): 167-175.

281. Irwig L, Tosteson AN, Gatsonis C, Lau J, Colditz G, Chalmers TC et al. Guidelines for meta-analyses evaluating diagnostic tests. Ann Intern Med 1994; 120(8): 667-676.

282. Jackson D. The power of the standard test for the presence of heterogeneity in meta-analysis. Stat Med 2006; 25(15): 2688-2699.

283. Jackson N, Waters E. Criteria for the systematic review of health promotion and public health interventions. Health Promot Int 2005; 20(4): 367-374.

284. Jadad AR. Randomised controlled trials: a user's guide. London: BMJ Books; 1998.

285. Jadad AR, Cook DJ, Browman GP. A guide to interpreting discordant systematic reviews. CMAJ 1997; 156(10): 1411-1416.

286. Jadad AR, Gagliardi A. Rating health information on the internet: navigating to knowledge or to Babel? JAMA 1998; 279(8): 611-614.

287. Johnston ME, Brouwers MC, Browman GP. Keeping cancer guidelines current: results of a comprehensive prospective literature monitoring strategy for twenty clinical practice guidelines. Int J Technol Assess Health Care 2003; 19(4): 646-655.

288. Jones B, Jarvis P, Lewis J, Ebbutt AF. Trials to assess equivalence: the importance of rigorous methods. BMJ 1996; 313(7048): 36-39.

289. Jones B, Kenward MG. Design and analysis of cross-over trials. London: Chapman and Hall; 1989. (Monographs on Statistics and Applied Probability; Volume 34 ).

290. Jull A, Bennett D. Do n-of-1 trials really tailor treatment? Lancet 2005; 365(9476): 1992-1994.

291. Jüni P, Altman DG, Egger M. Systematic reviews in health care: assessing the quality of controlled clinical trials. BMJ 2001; 323(7303): 42-46.

292. Jüni P, Holenstein F, Sterne J, Bartlett C, Egger M. Direction and impact of language bias in meta-analyses of controlled trials: empirical study. Int J Epidemiol 2002; 31(1): 115-123.

293. Kaltenthaler E, Brazier J, De Nigris E, Tumur I, Ferriter M, Beverley C et al. Computerised cognitive behaviour therapy for depression and anxiety update: a systematic review and economic evaluation. Health Technol Assess 2006; 10(33): 1-168.

294. Kass PH, Gold EB. Modern epidemiologic study designs. In: Ahrens W, Pigeot I (Ed). Handbook of epidemiology. Berlin: Springer; 2005. p. 321-344.

295. Katrak P, Bialocerkowski AE, Massy-Westropp N, Kumar S, Grimmer KA. A systematic review of the content of critical appraisal tools. BMC Med Res Methodol 2004; 4: 22.

296. Katz MH. Multivariable analysis: a primer for readers of medical research. Ann Intern Med 2003; 138(8): 644-650.

297. Kelley E, Hurst J. Health Care Quality Indicators Project: conceptual framework paper [online]. 9 March 2006 [accessed: 19 April 2013]. (OECD Health Working Papers; Volume 23). URL: http://www.oecd.org/dataoecd/1/36/36262363.pdf.

298. Kerr CE, Murray E, Stevenson F, Gore C, Nazareth I. Internet interventions for long-term conditions: patient and caregiver quality criteria. J Med Internet Res 2006; 8(3): e13-e27.

299. Kettunen T, Liimatainen L, Villberg J, Perko U. Developing empowering health counseling measurement: preliminary results. Patient Educ Couns 2006; 64(1-3): 159-166.

300. Kickbusch IS. Health literacy: addressing the health and education divide. Health Promot Int 2001; 16(3): 289-297.

301. Kieser M. Assessment of clinical relevance by considering point estimates and associated confidence intervals. Pharm Stat 2005; 4(2): 101-107.

302. Kieser M, Röhmel J, Friede T. Power and sample size determination when assessing the clinical relevance of trial results by 'responder analyses'. Stat Med 2004; 23(21): 3287-3305.

303. Klusen N, Meusch M (Ed). Wettbewerb und Solidarität im europäischen Gesundheitsmarkt. Baden-Baden: Nomos Verlagsgesellschaft; 2006. (Beiträge zum Gesundheitsmanagement; Volume 16).

304. Knelangen M, Zschorlich B, Büchter R, Fechtelpeter D, Rhodes T, Bastian H. Online-Umfragen auf Gesundheitsinformation.de: Ermittlung potenzieller Informationsbedürfnisse für evidenzbasierte Gesundheitsinformationen. Z Evid Fortbild Qual Gesundhwes 2010; 104(8-9): 667-673.

305. Knottnerus JA, Bouter LM. The ethics of sample size: two-sided testing and one-sided thinking. J Clin Epidemiol 2001; 54(2): 109-110.

306. Köbberling J. Der Zweifel als Triebkraft des Erkenntnisgewinns in der Medizin. In: Kunz R, Ollenschläger G, Raspe H, Jonitz G, Donner-Banzhoff N (Ed). Lehrbuch evidenzbasierte Medizin in Klinik und Praxis. Köln: Deutscher Ärzteverlag; 2007. p. 3-14.

307. Köbberling J, Trampisch HJ, Windeler J. Memorandum for the evaluation of diagnostic measures. J Clin Chem Clin Biochem 1990; 28(12): 873-879.

308. Koch A, Ziegler S. Metaanalyse als Werkzeug zum Erkenntnisgewinn. Med Klin 2000; 95(2): 109-116.

309. Koch G. No improvement: still less then half of the Cochrane reviews are up to date. In: XIV Cochrane Colloquium: programme and abstract book; 23.-26.10.2006; Dublin, Irland. 2006. p. 104.

310. Kolman J, Meng P, Scott G. Good clinical practice: standard operating procedures for clinical researchers. Chichester: Wiley; 1998.

311. Kommission der Europäischen Gemeinschaften. Richtlinie 2003/63/EG der Kommission vom 25. Juni 2003 zur Änderung der Richtlinie 2001/83/EG des Europäischen Parlaments und des Rates zur Schaffung eines Gemeinschaftskodexes für Humanarzneimittel. Amtsblatt der Europäischen Gemeinschaften 2003; 46(L159): 46-94.

312. Kools M, Van de Wiel MW, Ruiter RA, Kok G. Pictures and text in instructions for medical devices: effects on recall and actual performance. Patient Educ Couns 2006; 64(1-3): 104-111.

313. Köpke S, Berger B, Steckelberg A, Meyer G. In Deutschland gebräuchliche Bewertungsinstrumente für Patienteninformationen: eine kritische Analyse. Z Arztl Fortbild Qualitatssich 2005; 99(6): 353-357.

314. Kraemer HC, Frank E, Kupfer DJ. Moderators of treatment outcomes: clinical, research, and policy importance. JAMA 2006; 296(10): 1286-1289.

315. Kranich C. Patientenkompetenz: was müssen Patienten wissen und können? Bundesgesundheitsblatt Gesundheitsforschung Gesundheitsschutz 2004; 47(10): 950-956.

316. Kristensen FB, Sigmund H. Health technology assessment handbook. Kopenhagen: Danish Centre for Heath Technology Assessment; 2007. URL: http://www.sst.dk/~/media/Planlaegning%20og%20kvalitet/MTV%20metode/HTA_Handbook_net_final.ashx.

317. Kristman V, Manno M, Cote P. Loss to follow-up in cohort studies: how much is too much? Eur J Epidemiol 2004; 19(8): 751-760.

318. Krug S. Don't make me think! Web Usability; das intuitive Web. Heidelberg: mitp; 2006.

319. Kulbe A. Grundwissen Psychologie, Soziologie und Pädagogik: Lehrbuch für Pflegeberufe. Stuttgart: Kohlhammer; 2009.

320. Kunz R, Djulbegovic B, Schünemann HJ, Stanulla M, Muti P, Guyatt G. Misconceptions, challenges, uncertainty, and progress in guideline recommendations. Semin Hematol 2008; 45(3): 167-175.

321. Kunz R, Lelgemann M, Guyatt GH, Antes G, Falck-Ytter Y, Schünemann H. Von der Evidenz zur Empfehlung. In: Kunz R, Ollenschläger G, Raspe H, Jonitz G, Donner-Banzhoff N (Ed). Lehrbuch evidenzbasierte Medizin in Klinik und Praxis. Köln: Deutscher-Ärzte-Verlag; 2007. p. 231-247.

322. Kunz R, Vist GE, Oxman AD. Randomisation to protect against selection bias in healthcare trials. Cochrane Database Syst Rev 2007; (2): MR000012.

323. Laaser U, Hurrelmann K. Gesundheitsförderung und Krankheitsprävention. In: Hurrelmann K, Laaser U (Ed). Handbuch Gesundheitswissenschaften. Weinheim: Juventa Verlag; 1998. p. 395-424.

324. Lange S, Freitag G. Choice of delta: requirements and reality; results of a systematic review. Biom J 2005; 47(1): 12-27.

325. Lapsley P. The patient's journey: travelling through life with a chronic illness. BMJ 2004; 329(7466): 582-583.

326. Last JM, Spasoff RA, Harris SS, Thuriaux MC (Ed). A dictionary of epidemiology. Oxford: Oxford University Press; 2001.

327. Lau J, Ioannidis JPA, Terrin N, Schmid CH, Olkin I. The case of the misleading funnel plot. BMJ 2006; 333(7568): 597-600.

328. Laupacis A, Wells G, Richardson WS, Tugwell P. Users' guides to the medical literature; V: how to use an article about prognosis. JAMA 1994; 272(3): 234-237.

329. Lausen B, Potapov S, Prokosch HU. Gesundheitsbezogene Internetnutzung in Deutschland 2007. GMS Med Inform Biom Epidemiol 2007; 4(2): Doc06.

330. Lavis JN. How can we support the use of systematic reviews in policymaking? PLoS Med 2009; 6(11): e1000141.

331. Leeflang MM, Deeks JJ, Gatsonis C, Bossuyt PM. Systematic reviews of diagnostic test accuracy. Ann Intern Med 2008; 149(12): 889-897.

332. Lefebvre C, Manheimer E, Glanville J. Searching for studies. In: Higgins JPT, Green S (Ed). Cochrane handbook for systematic reviews of interventions. New York: Wiley; 2008. p. 95-150.

333. Lehoux P, Tailliez S, Denis JL, Hivon M. Redefining health technology assessment in Canada: diversification of products and contextualization of findings. Int J Technol Assess Health Care 2004; 20(3): 325-336.

334. Lewis D. Computer-based approaches to patient education: a review of the literature. J Am Med Inform Assoc 1999; 6(4): 272-282.

335. Lewis S, Clarke M. Forest plots: trying to see the wood and the trees. BMJ 2001; 322(7300): 1479-1480.

336. Leys M. Health care policy: qualitative evidence and health technology assessment. Health Policy 2003; 65(3): 217-226.

337. Liberati A, Altman DG, Tetzlaff J, Mulrow C, Gøtzsche PC, Ioannidis JPA et al. The PRISMA statement for reporting systematic reviews and meta-analyses of studies that evaluate healthcare interventions: explanation and elaboration. BMJ 2009; 339: b2700.

338. Liberati A, Sheldon TA, Banta HD. EUR-ASSESS project subgroup report on methodology: methodological guidance for the conduct of health technology assessment. Int J Technol Assess Health Care 1997; 13(2): 186-219.

339. Lieb K, Klemperer D, Koch K, Baethge C, Ollenschläger G, Ludwig WD. Interessenkonflikt in der Medizin: mit Transparenz Vertrauen stärken. Dtsch Arztebl 2011; 108(6): A256-A260.

340. Lijmer JG, Bossuyt PM. Various randomized designs can be used to evaluate medical tests. J Clin Epidemiol 2009; 62(4): 364-373.

341. Lijmer JG, Mol BW, Heisterkamp S, Bonsel GJ, Prins MH, Van der Meulen JH et al. Empirical evidence of design-related bias in studies of diagnostic tests. JAMA 1999; 282(11): 1061-1066.

342. Lipkus IM, Hollands JG. The visual communication of risk. J Natl Cancer Inst Monogr 1999; (25): 149-163.

343. Lipscomb J, Drummond M, Fryback D, Gold M, Revicki D. Retaining, and enhancing, the QALY. Value Health 2009; 12(Suppl 1): S18-S26.

344. Little RJA, Rubin DB. Statistical analysis with missing data. Hoboken: Wiley; 2002.

345. Lo B, Field MJ (Ed). Conflict of interest in medical research, education, and practice. Washington: National Academies Press; 2009.

346. Lord SJ, Irwig LM, Simes RJ. When is measuring sensitivity and specificity sufficient to evaluate a diagnostic test, and when do we need randomized trials? Ann Intern Med 2006; 144(11): 850-855.

347. Lu G, Ades AE. Combination of direct and indirect evidence in mixed treatment comparisons. Stat Med 2004; 23(20): 3105-3124.

348. Lu G, Ades AE. Assessing evidence inconsistency in mixed treatment comparisons. J Am Stat Assoc 2006; 101(474): 447-459.

349. Lu G, Ades AE, Sutton AJ, Cooper NJ, Briggs AH, Caldwell DM. Meta-analysis of mixed treatment comparisons at multiple follow-up times. Stat Med 2007; 26(20): 3681-3699.

350. Lumley T. Network meta-analysis for indirect treatment comparisons. Stat Med 2002; 21(16): 2313-2324.

351. Macaskill P, Walter SD, Irwig LM. A comparison of methods to detect publication bias in meta-analysis. Stat Med 2001; 20(4): 641-654.

352. MacDermid JC, Brooks D, Solway S, Switzer-McIntyre S, Brosseau L, Graham ID. Reliability and validity of the AGREE instrument used by physical therapists in assessment of clinical practice guidelines. BMC Health Serv Res 2005; 5: 18.

353. MacLean CH, Morton SC, Ofman JJ, Roth EA, Shekelle PG. How useful are unpublished data from the Food and Drug Administration in meta-analysis? J Clin Epidemiol 2003; 56(1): 44-51.

354. Maetzel A. Der Gebrauch von Nutzwerten im gesundheitsökonomischen Vergleich von Interventionen bei verschiedenen Krankheitsbildern: eine Einführung. Z Rheumatol 2004; 63(5): 380-384.

355. Malterud K. The art and science of clinical knowledge: evidence beyond measures and numbers. Lancet 2001; 358(9279): 397-400.

356. Mangiapane S, Velasco Garrido M. Surrogatendpunkte als Parameter der Nutzenbewertung [online]. 2009 [accessed: 24 April 2013]. (Schriftenreihe Health Technology Assessment; Volume 91). URL: http://portal.dimdi.de/de/hta/hta_berichte/hta250_bericht_de.pdf.

357. March JS, Silva SG, Compton S, Shapiro M, Califf R, Krishnan R. The case for practical clinical trials in psychiatry. Am J Psychiatry 2005; 162(5): 836-846.

358. Martini P. Methodenlehre der therapeutischen Untersuchung. Berlin: Springer; 1932.

359. Matthys J, De Meyere M, Van Driel ML, De Sutter A. Differences among international pharyngitis guidelines: not just academic. Ann Fam Med 2007; 5(5): 436-443.

360. McAlister FA, Straus SE, Sackett DL, Altman DG. Analysis and reporting of factorial trials: a systematic review. JAMA 2003; 289(19): 2545-2553.

361. McAlister FA, Van Diepen S, Padwal RS, Johnson JA, Majumdar SR. How evidence-based are the recommendations in evidence-based guidelines? PLoS Med 2007; 4(8): e250.

362. McCulloch P, Taylor I, Sasako M, Lovett B, Griffin D. Randomised trials in surgery: problems and possible solutions. BMJ 2002; 324(7351): 1448-1451.

363. McGauran N, Wieseler B, Kreis J, Schüler YB, Kölsch H, Kaiser T. Reporting bias in medical research: a narrative review. Trials 2010; 11(1): 37.

364. McGregor M, Caro JJ. QALYs: are they helpful to decision makers? Pharmacoeconomics 2006; 24(10): 947-952.

365. McMurray J, Swedberg K. Treatment of chronic heart failure: a comparison between the major guidelines. Eur Heart J 2006; 27(15): 1773-1777.

366. Medical Services Advisory Committee. Guidelines for the assessment of diagnostic technologies. Canberra: MSAC; 2005. URL: http://www.health.gov.au/internet/msac/publishing.nsf/Content/D81BE529B98B3DB6CA257 5AD0082FD1B/$File/Guidelines%20for%20the%20assessment%20of%20diagnostic%20tech nologies%20Sept%202005.pdf.

367. Mills E, Jadad AR, Ross C, Wilson K. Systematic review of qualitative studies exploring parental beliefs and attitudes toward childhood vaccination identifies common barriers to vaccination. J Clin Epidemiol 2005; 58(11): 1081-1088.

368. Moher D, Hopewell S, Schulz KF, Montori V, Gøtzsche PC, Devereaux PJ et al. CONSORT 2010 explanation and elaboration: updated guidelines for reporting parallel group randomised trials. BMJ 2010; 340: c869.

369. Moher D, Liberati A, Tetzlaff J, Altman DG. Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. BMJ 2009; 339: b2535.

370. Moher D, Pham B, Klassen TP, Schulz KF, Berlin JA, Jadad AR et al. What contributions do languages other than English make on the results of meta-analyses? J Clin Epidemiol 2000; 53(9): 964-972.

371. Moher D, Tetzlaff J, Tricco AC, Sampson M, Altman DG. Epidemiology and reporting characteristics of systematic reviews. PLoS Med 2007; 4(3): e78.

372. Moher D, Tsertsvadze A. Systematic reviews: when is an update an update? Lancet 2006; 367(9514): 881-883.

373. Molenaar S, Sprangers M, Oort F, Rutgers E, Luiten E, Mulder J et al. Exploring the black box of a decision aid: what information do patients select from an interactive CD-ROM on treatment options in breast cancer? Patient Educ Couns 2007; 65(1): 122-130.

374. Molenberghs G, Burzykowski T, Alonso A, Assam P, Tilahun A, Buyse M. A unified framework for the evaluation of surrogate endpoints in mental-health clinical trials. Stat Methods Med Res 2010; 19(3): 205-236.

375. Molnar FJ, Man-Son-Hing M, Fergusson D. Systematic review of measures of clinical significance employed in randomized controlled trials of drugs for dementia. J Am Geriatr Soc 2009; 57(3): 536-546.

376. Moons KG, Altman DG, Vergouwe Y, Royston P. Prognosis and prognostic research: application and impact of prognostic models in clinical practice. BMJ 2009; 338: b606.

377. Moons KG, Royston P, Vergouwe Y, Grobbee DE, Altman DG. Prognosis and prognostic research: what, why, and how? BMJ 2009; 338: b375.

378. Moult B, Franck LS, Brady H. Ensuring quality information for patients: development and preliminary validation of a new instrument to improve the quality of written health care information. Health Expect 2004; 7(2): 165-175.

379. Müller HP, Schmidt K, Conen D. Qualitätsmanagement: interne Leitlinien und Patientenpfade. Med Klin 2001; 96(11): 692-697.

380. Müllner M, Matthews H, Altman DG. Reporting on statistical methods to adjust for confounding: a cross-sectional survey. Ann Intern Med 2002; 136(2): 122-126.

381. Munday J. Introducing translation studies: theories and applications. London: Routledge; 2001.

382. Murphy E, Dingwall R, Greatbatch D, Parker S, Watson P. Qualitative research methods in health technology assessment: a review of the literature. Health Technol Assess 1998; 2(16): 1-274.

383. Murray E, Burns J, See TS, Lai R, Nazareth I. Interactive health communication applications for people with chronic disease. Cochrane Database Syst Rev 2005; (4): CD004274.

384. National Advisory Committee on Health and Disability. Screening to improve health in New Zealand: criteria to assess screening. Wellington: National Health Committee; 2003. URL: http://www.nsu.govt.nz/files/NSU/Screening_to_improve_health.pdf.

385. National Health and Medical Research Council. Statement on consumer and community participation in health and medical research. Canberra: Commonwealth of Australia; 2002. URL: http://www.nhmrc.gov.au/_files_nhmrc/file/publications/synopses/r22.pdf.

386. National Health and Medical Research Council. Cultural competency in health: a guide for policy, partnerships and participation. Canberra: Commonwealth of Australia; 2006. URL: http://www.nhmrc.gov.au/_files_nhmrc/file/publications/synopses/hp19.pdf.

387. National Institute of Clinical Studies. The impact of the internet on consumers health behaviour [online]. February 2003 [accessed: 24 April 2013]. URL: http://www.nhmrc.gov.au/_files_nhmrc/file/nics/material_resources/The%20Impact%20of%20the%20Internet%20on%20Consumers%20Health%20Behaviour.pdf.

388. National Resource Centre for Consumer Participation in Health. Methods and models of consumer participation [online]. 1 September 2008 [accessed: 6 May 2013]. (Information Series; Volume 2). URL: http://www.healthissuescentre.org.au/documents/items/2008/09/231154-upload-00001.pdf.

389. National Resource Centre for Consumer Participation in Health. Feedback, participation and consumer diversity: a literature review. Canberra: Commonwealth of Australia; 2000. URL: http://www.healthissuescentre.org.au/documents/items/2008/08/226293-upload-00001.pdf.

390. NHS Centre for Reviews and Dissemination. The Database of Abstracts of Reviews of Effects (DARE). Effectiveness Matters 2002; 6(2): 1-4.

391. Nielsen J. Designing web usability. Indianapolis: New Riders Publishing; 2000.

392. Nielsen J, Loranger H. Web Usability. München: Addison-Wesley; 2008.

393. Nilsen ES, Myrhaug HT, Johansen M, Oliver S, Oxman AD. Methods of consumer involvement in developing healthcare policy and research, clinical practice guidelines and patient information material. Cochrane Database Syst Rev 2006; (3): CD004563.

394. Nord E. Cost-value analysis in health care: making sense out of QALYs. Cambridge: Cambridge University Press; 1999.

395. Nüesch E, Jüni P. Commentary: which meta-analyses are conclusive? Int J Epidemiol 2009; 38(1): 298-303.

396. Nutbeam D. Health promotion glossary. Health Promot Int 1998; 13(4): 349-364.

397. O'Connor AM, Bennett CL, Stacey D, Barry M, Col NF, Eden KB et al. Decision aids for people facing health treatment or screening decisions. Cochrane Database Syst Rev 2009; (3): CD001431.

398. Oliver S, Clarke-Jones L, Rees R, Milne R, Buchanan P, Gabbay J et al. Involving consumers in research and development agenda setting for the NHS: developing an evidence-based approach. Health Technol Assess 2004; 8(15): 1-148.

399. Oxman AD, Guyatt GH. Guidelines for reading literature reviews. Can Med Assoc J 1988; 138(8): 697-703.

400. Oxman AD, Guyatt GH. Validation of an index of the quality of review articles. J Clin Epidemiol 1991; 44(11): 1271-1278.

401. Oxman AD, Guyatt GH. A consumer's guide to subgroup analyses. Ann Intern Med 1992; 116(1): 78-84.

402. Oxman AD, Guyatt GH, Singer J, Goldsmith CH, Hutchison BG, Milner RA et al. Agreement among reviewers of review articles. J Clin Epidemiol 1991; 44(1): 91-98.

403. Parkin DM, Chen VW, Ferlay J, Galceran J, Storm HH (Ed). Comparability and quality control in cancer registration. Lyon: International Agency for Research on Cancer; 1994. (IARC Technical Reports; Volume 19).

404. Parmar MK, Torri V, Stewart L. Extracting summary statistics to perform meta-analyses of the published literature for survival endpoints. Stat Med 1998; 17(24): 2815-2834.

405. Patientenuniversität an der Medizinischen Hochschule Hannover. Nutzertestung von Gesundheitsinformationen des Instituts für Qualität und Wirtschaftlichkeit (IQWiG): Abschlussbericht [online]. October 2010 [accessed: 22 April 2013]. URL: http://www.mh-hannover.de/fileadmin/institute/epidemiologie/public_health/downloads/KursbeschrMar2011/NutzertestungAbschlussberichtIQWIG_16_12.pdf.

406. Pearson SD, Rawlins MD. Quality, innovation, and value for money: NICE and the British National Health Service. JAMA 2005; 294(20): 2618-2622.

407. Pennekamp PH, Diedrich O, Schmitt O, Kraft CN. Prävalenz und Stellenwert der Internetnutzung orthopädischer Patienten. Zeitschrift für Orthopädie 2006; 144(5): 459-463.

408. Perleth M, Jakubowski E, Busse R. What is 'best practice' in health care? State of the art and perspectives in improving the effectiveness and efficiency of the European health care systems. Health Policy 2001; 56(3): 235-250.

409. Peters JL, Sutton A, Jones D, Abrams KR, Rushton L. Comparison of two methods to detect publication bias in meta-analysis. JAMA 2006; 295(6): 676-680.

410. Petitti DB, Teutsch SM, Barton MB, Sawaya GF, Ockene JK, DeWitt T. Update on the methods of the U.S. Preventive Services Task Force: insufficient evidence. Ann Intern Med 2009; 150(3): 199-205.

411. Pfaff H, Glaeske G, Neugebauer EA, Schrappe M. Memorandum III: Methoden für die Versorgungsforschung (Teil 1). Gesundheitswesen 2009; 71(8-9): 505-510.

412. Pham B, Platt R, McAuley L, Klassen TP, Moher D. Is there a "best" way to detect and minimize publication bias? An empirical evaluation. Eval Health Prof 2001; 24(2): 109-125.

413. Piaggio G, Elbourne DR, Altman DG, Pocock SJ, Evans SJ. Reporting of noninferiority and equivalence randomized trials: an extension of the CONSORT statement. JAMA 2006; 295(10): 1152-1160.

414. Pocock SJ. Clinical trials: a practical approach. Chichester: Wiley; 1983.

415. Pocock SJ, Collier TJ, Dandreo KJ, De Stavola BL, Goldman MB, Kalish LA et al. Issues in the reporting of epidemiological studies: a survey of recent practice. BMJ 2004; 329(7471): 883-887.

416. Poensgen A, Larsson S. Patients, physicians and the internet: myth, reality, and implications [online]. January 2001 [accessed: 23 April 2013]. URL: http://www.bostonconsulting.com.au/documents/file13629.pdf.

417. Pollock K. Concordance in medical consultations: a critical review. Oxford: Radcliffe Publishing; 2005.

418. Poynard T, Munteanu M, Ratziu V, Benhamou Y, Di Martino V, Taieb J et al. Truth survival in clinical research: an evidence-based requiem? Ann Intern Med 2002; 136(12): 888-895.

419. Prentice RL. Surrogate endpoints in clinical trials: definition and operational criteria. Stat Med 1989; 8(4): 431-440.

420. Puhan MA, Singh S, Weiss CO, Varadhan R, Boyd CM. A framework for organizing and selecting quantitative approaches for benefit-harm assessment. BMC Med Res Methodol 2012; 12: 173.

421. Raum E, Perleth M. Methoden der Metaanalyse von diagnostischen Genauigkeitsstudien. Köln: Deutsches Institut für Medizinische Dokumentation und Information; 2003. (Schriftenreihe Health Technology Assessment; Volume 2). URL: http://portal.dimdi.de/de/hta/hta_berichte/hta025_bericht_de.pdf.

422. Reitsma JB, Glas AS, Rutjes AW, Scholten RJ, Bossuyt PM, Zwinderman AH. Bivariate analysis of sensitivity and specificity produces informative summary measures in diagnostic reviews. J Clin Epidemiol 2005; 58(10): 982-990.

423. Revicki D, Hays RD, Cella D, Sloan J. Recommended methods for determining responsiveness and minimally important differences for patient-reported outcomes. J Clin Epidemiol 2008; 61(2): 102-109.

424. Richter JG, Becker A, Specker C, Monser R, Schneider M. Krankheitsbezogene Internetnutzung bei Patienten mit entzündlich-rheumatischen Systemerkrankungen. Z Rheumatol 2004; 63(3): 216-222.

425. Riley RD, Higgins JP, Deeks JJ. Interpretation of random effects meta-analyses. Br Med J 2011; 342: 964-967.

426. Rockwood K, Fay S, Song X, MacKnight C, Gorman M. Attainment of treatment goals by people with Alzheimer's disease receiving galantamine: a randomized controlled trial. Can Med Assoc J 2006; 174(8): 1099-1105.

427. Roebruck P, Elze M, Hauschke D, Leverkus F, Kieser M. Literaturübersicht zur Fallzahlplanung für Äquivalenzprobleme. Inform Biom Epidemiol Med Biol 1997; 28(2): 51-63.

428. Röhmel J, Hauschke D, Koch A, Pigeot I. Biometrische Verfahren zum Wirksamkeitsnachweis im Zulassungsverfahren: Nicht-Unterlegenheit in klinischen Studien. Bundesgesundheitsblatt Gesundheitsforschung Gesundheitsschutz 2005; 48(5): 562-571.

429. Rossouw JE. Estrogens for prevention of coronary heart disease: putting the brakes on the bandwagon. Circulation 1996; 94(11): 2982-2985.

430. Rothwell PM. Treating individuals 2: subgroup analysis in randomised controlled trials; importance, indications, and interpretation. Lancet 2005; 365(9454): 176-186.

431. Rotter T, Kinsman L, James E, Machotta A, Gothe H, Willis J et al. Clinical pathways: effects on professional practice, patient outcomes, length of stay and hospital costs. Cochrane Database Syst Rev 2010; (3): CD006632.

432. Royal Society. Science and the public interest: communicating the results of new scientific research to the public [online]. April 2006 [accessed: 23 April 2013]. URL: http://royalsociety.org/uploadedFiles/Royal_Society_Content/policy/publications/2006/8315. pdf.

433. Royall RM. The effect of sample size on the meaning of significance tests. Am Stat 1986; 40(4): 313-315.

434. Royle P, Waugh N. Literature searching for clinical and cost-effectiveness studies used in health technology assessment reports carried out for the National Institute for Clinical Excellence appraisal system. Health Technol Assess 2003; 7(34): 1-51.

435. Royston P. A strategy for modelling the effect of a continuous covariate in medicine and epidemiology. Stat Med 2000; 19(14): 1831-1847.

436. Royston P, Altman DG. Regression using fractional polynomials of continuous covariates: parsimonious parametric modelling. Appl Stat 1994; 43(3): 429-467.

437. Royston P, Moons KG, Altman DG, Vergouwe Y. Prognosis and prognostic research: developing a prognostic model. BMJ 2009; 338: b604.

438. Rücker G, Schwarzer G, Carpenter J, Olkin I. Why add anything to nothing? The arcsine difference as a measure of treatment effect in meta-analysis with zero cells. Stat Med 2009; 28(5): 721-738.

439. Rutten LJ, Arora NK, Bakos AD, Aziz N, Rowland J. Information needs and sources of information among cancer patients: a systematic review of research (1980-2003). Patient Educ Couns 2005; 57(3): 250-261.

440. Saaty TL. A scaling method for priorities in hierarchical structures. J Math Psychol 1977; 15(3): 234-281.

441. Saaty TL. Decision making with the Analytic Hierarchy Process. International Journal of Services Sciences 2008; 1(1): 83-98.

442. Saaty TL, Vargas LG. The Analytic Hierarchy Process: wash criteria should not be ignored. International Journal of Management and Decision Making 2006; 7(2/3): 180-188.

443. Sachverständigenrat für die Konzertierte Aktion im Gesundheitswesen. Bedarfsgerechtigkeit und Wirtschaftlichkeit; Band III: Über- Unter- und Fehlversorgung; Gutachten 2000/2001; ausführliche Zusammenfassung [online]. August 2001 [accessed: 6 May 2013]. URL: http://www.svr-gesundheit.de/fileadmin/user_upload/Gutachten/2000-2001/Kurzf-de-01.pdf.

444. Sackett DL. Bias in analytic research. J Chronic Dis 1979; 32(1-2): 51-63.

445. Sackett DL, Rosenberg WM, Gray JA, Haynes RB, Richardson WS. Evidence based medicine: what it is and what it isn't. BMJ 1996; 312(7023): 71-72.

446. Sackett DL, Straus SE, Richardson WS, Rosenberg W, Haynes RB. Evidence-based medicine: how to practice and teach EBM. Edinburgh: Churchill Livingstone; 2000.

447. Salanti G, Higgins JPT, Ades AE, Ioannidis JPA. Evaluation of networks of randomized trials. Stat Methods Med Res 2008; 17(3): 279-301.

448. Salanti G, Marinho V, Higgins JPT. A case study of multiple-treatments meta-analysis demonstrates that covariates should be considered. J Clin Epidemiol 2009; 62(8): 857-864.

449. Sampson M, Barrowman NJ, Moher D, Klassen TP, Pham B, Platt R et al. Should meta-analysts search Embase in addition to Medline? J Clin Epidemiol 2003; 56(10): 943-955.

450. Sampson M, McGowan J, Cogo E, Grimshaw J, Moher D, Lefebvre C. An evidence-based practice guideline for the peer review of electronic search strategies. J Clin Epidemiol 2009; 62(9): 944-952.

451. Sampson M, Shojania KG, McGowan J, Daniel R, Rader T, Iansavichene AE et al. Surveillance search techniques identified the need to update systematic reviews. J Clin Epidemiol 2008; 61(8): 755-762.

452. Sampson MJ. Updating searches for systematic reviews [Dissertation]. Aberystwyth: Universität; 2009.

453. Sänger S, Lang B, Klemperer D, Thomeczek C, Dierks ML. Manual Patienteninformation: Empfehlungen zur Erstellung evidenzbasierter Patienteninformationen. Berlin: Ärztliches Zentrum für Qualität in der Medizin; 2006. (ÄZQ-Schriftenreihe; Volume 25). URL: http://www.aezq.de/mdb/edocs/pdf/schriftenreihe/schriftenreihe25.pdf.

454. Santo A, Laizner AM, Shohet L. Exploring the value of audiotapes for health literacy: a systematic review. Patient Educ Couns 2005; 58(3): 235-243.

455. Sargent DJ, Conley BA, Allegra C, Collette L. Clinical trial designs for predictive marker validation in cancer treatment trials. J Clin Oncol 2005; 23(9): 2020-2027.

456. SAS Institute. SAS/STAT 9.2 user's guide: second edition [online]. 2009 [accessed: 4 November 2013]. URL: http://support.sas.com/documentation/cdl/en/statug/63033/PDF/default/statug.pdf.

457. Sauerbrei W, Royston P. Building multivariable prognostic and diagnostic models: transformation of the predictors by using fractional polynomials. J R Stat Soc Ser A 1999; 162(1): 71-94.

458. Sawaya GF, Guirguis-Blake J, LeFevre M, Harris R, Petitti D. Update on the methods of the U.S. Preventive Services Task Force: estimating certainty and magnitude of net benefit. Ann Intern Med 2007; 147(12): 871-875.

459. Sawicki PT. Qualität der Gesundheitsversorgung in Deutschland: ein randomisierter Sechs-Länder-Vergleich aus Patientensicht. Med Klin 2005; 100(11): 755-768.

460. Schaeffer D, Dierks ML, Hurrelmann K, Keller A, Krause H, Schmidt-Kaehler S et al. Evaluation der Modellprojekte zur Patienten- und Verbraucherberatung nach § 65B Sozialgesetzbuch V: Abschlussbericht der wissenschaftlichen Begleitforschung für die Spitzenverbände der GKV [online]. November 2004 [accessed: 6 May 2013]. URL: http://www.gkv-spitzenverband.de/media/dokumente/krankenversicherung_1/praevention__selbsthilfe__beratung/beratung/Abschlussbericht_UPD_1_Modellphase_Uni_Bielefeld.pdf.

461. Schluter PJ, Ware RS. Single patient (n-of-1) trials with binary treatment preference. Stat Med 2005; 24(17): 2625-2636.

462. Schmidt-Kaehler S. Patienteninformation Online: theoretische Grundlagen, Planung und Entwicklung eines Konzeptes für die Patientenschulung im Internet. Bern: Huber; 2004.

463. Schneider N, Dierks ML, Seidel G, Schwartz FW. The federal government commissioner for patient issues in Germany: initial analysis of the user inquiries. BMC Health Serv Res 2007; 7: 24.

464. Schöffski O, Graf von der Schulenburg JM (Ed). Gesundheitsökonomische Evaluationen. Berlin: Springer; 2012.

465. Schulz KF, Altman DG, Moher D. CONSORT 2010 statement: updated guidelines for reporting parallel group randomised trials. BMJ 2010; 340: c332.

466. Schünemann HJ, Akl EA, Guyatt GH. Interpreting the results of patient reported outcome measures in clinical trials: the clinician's perspective. Health Qual Life Outcomes 2006; 4: 62.

467. Schünemann HJ, Best D, Vist GE, Oxman AD. Letters, numbers, symbols and words: how to communicate grades of evidence and recommendations. Can Med Assoc J 2003; 169(7): 677-680.

468. Schünemann HJ, Fretheim A, Oxman AD. Improving the use of research evidence in guideline development; 9: grading evidence and recommendations. Health Res Policy Syst 2006; 4: 21.

469. Scott JT, Harmsen M, Prictor MJ, Entwistle VA, Sowden AJ, Watt I. Recordings or summaries of consultations for people with cancer. Cochrane Database Syst Rev 2008; (3): CD001539.

470. Senn SJ. Inherent difficulties with active control equivalence studies. Stat Med 1993; 12(24): 2367-2375.

471. Senn SJ. The many modes of meta. Drug Inf J 2000; 34(2): 535-549.

472. Senn SJ. Trying to be precise about vagueness. Stat Med 2007; 26(7): 1417-1430.

473. Shea BJ, Bouter LM, Peterson J, Boers M, Andersson N, Ortiz Z et al. External validation of a measurement tool to assess systematic reviews (AMSTAR). PLoS One 2007; 2(12): e1350.

474. Shea BJ, Grimshaw JM, Wells GA, Boers M, Andersson N, Hamel C et al. Development of AMSTAR: a measurement tool to assess the methodological quality of systematic reviews. BMC Med Res Methodol 2007; 7: 10.

475. Shea BJ, Hamel C, Wells GA, Bouter LM, Kristjansson E, Grimshaw J et al. AMSTAR is a reliable and valid measurement tool to assess the methodological quality of systematic reviews. J Clin Epidemiol 2009; 62(10): 1013-1020.

476. Shekelle PG, Ortiz E, Rhodes S, Morton SC, Eccles MP, Grimshaw JM et al. Validity of the agency for healthcare research and quality clinical practice guidelines: how quickly do guidelines become outdated? JAMA 2001; 286(12): 1461-1467.

477. Shojania KG, Sampson M, Ansari MT, Ji J, Doucette S, Moher D. How quickly do systematic reviews go out of date? A survival analysis. Ann Intern Med 2007; 147(4): 224-233.

478. Sillence E, Briggs P, Harris PR, Fishwick L. How do patients evaluate and make use of online health information? Soc Sci Med 2007; 64(9): 1853-1862.

479. Silvestre MAA, Dans LF, Dans AL. Trade-off between benefit and harm is crucial in health screening recommendations; part II: evidence summaries. J Clin Epidemiol 2011; 64(3): 240-249.

480. Simmonds MC, Higgins JPT. Covariate heterogeneity in meta-analysis: criteria for deciding between meta-regression and individual patient data. Stat Med 2007; 26(15): 2982-2999.

481. Simon R, Altman DG. Statistical aspects of prognostic factor studies in oncology. Br J Cancer 1994; 69(6): 979-985.

482. Skipka G, Bender R. Intervention effects in the case of heterogeneity between three subgroups: assessment within the framework of systematic reviews. Methods Inf Med 2010; 49(6): 613-617.

483. Smeeth L, Haines A, Ebrahim S. Numbers needed to treat derived from meta-analyses: sometimes informative, usually misleading. Br Med J 1999; 318(7197): 1548-1551.

484. Song F, Eastwood AJ, Gilbody S, Duley L, Sutton A. Publication and related biases. Health Technol Assess 2000; 4(10): 1-115.

485. Song F, Loke YK, Walsh T, Glenny AM, Eastwood AJ, Altman DG. Methodological problems in the use of indirect comparisons for evaluating healthcare interventions: survey of published systematic reviews. BMJ 2009; 338: b1147.

486. Song F, Parekh S, Hooper L, Loke YK, Ryder J, Sutton AJ et al. Dissemination and publication of research findings: an updated review of related biases. Health Technol Assess 2010; 14(8): 1-193.

487. Spek V, Cuijpers P, Nyklicek I, Riper H, Keyzer J, Pop V. Internet-based cognitive behaviour therapy for symptoms of depression and anxiety: a meta-analysis. Psychol Med 2007; 37(3): 319-328.

488. Spiegelhalter DJ, Freedman LS, Parmar MK. Bayesian approaches to randomized trials. J R Stat Soc Ser A 1994; 157(3): 357-416.

489. Spiegelhalter DJ, Myles JP, Jones D, Abrams KR. Methods in health service research: an introduction to Bayesian methods in health technology assessment. BMJ 1999; 319(7208): 508-512.

490. Statistisches Bundesamt. Statistik der schwerbehinderten Menschen 2007: Kurzbericht [online]. January 2009 [accessed: 22 April 2013]. URL: https://www.destatis.de/DE/Publikationen/Thematisch/Gesundheit/BehinderteMenschen/Sozi alSchwerbehinderteKB5227101079004.pdf?__blob=publicationFile.

491. Statistisches Bundesamt. Leben in Deutschland: Haushalte, Familien und Gesundheit; Ergebnisse des Mikrozensus 2005 [online]. June 2006 [accessed: 7 May 2013]. URL: https://www.destatis.de/DE/PresseService/Presse/Pressekonferenzen/2006/Mikrozensus/Press ebroschuere.pdf?__blob=publicationFile.

492. Steckelberg A, Berger B, Köpke S, Heesen C, Mühlhauser I. Kriterien für evidenzbasierte Patienteninformationen. Z Arztl Fortbild Qualitatssich 2005; 99(6): 343-351.

493. Steckelberg A, Kasper J, Redegeld M, Mühlhauser I. Risk information: barrier to informed choice? A focus group study. Soz Praventivmed 2004; 49(6): 375-380.

494. Steiner JF. The use of stories in clinical research and health policy. JAMA 2005; 294(22): 2901-2904.

495. Sterne J, Egger M, Smith GD. Systematic reviews in health care: investigating and dealing with publication and other biases in meta-analysis. BMJ 2001; 323(7304): 101-105.

496. Sterne JAC, Egger M, Moher D. Addressing reporting biases. In: Higgins JPT, Green S (Ed). Cochrane handbook for systematic reviews of interventions. Chichester: Wiley; 2008. p. 297-333.

497. Stroup DF, Berlin JA, Morton SC, Olkin I, Williamson GD, Rennie D et al. Meta-analysis of observational studies in epidemiology: a proposal for reporting. JAMA 2000; 283(15): 2008-2012.

498. Sullivan LM, Massaro JM, D'Agostino RB Sr. Presentation of multivariate data for clinical use: the Framingham study risk score functions. Stat Med 2004; 23(10): 1631-1660.

499. Sun GW, Shook TL, Kay GL. Inappropriate use of bivariable analysis to screen risk factors for use in multivariable analysis. J Clin Epidemiol 1996; 49(8): 907-916.

500. Sun X, Briel M, Walter SD, Guyatt GH. Is a subgroup effect believable? Updating criteria to evaluate the credibility of subgroup analyses. BMJ 2010; 340: c117.

501. Sutton A, Ades AE, Cooper N, Abrams K. Use of indirect and mixed treatment comparisons for technology assessment. Pharmacoeconomics 2008; 26(9): 753-767.

502. Sutton AJ, Donegan S, Takwoingi Y, Garner P, Gamble C, Donald A. An encouraging assessment of methods to inform priorities for updating systematic reviews. J Clin Epidemiol 2009; 62(3): 241-251.

503. Swift TL, Dieppe PA. Using expert patients' narratives as an educational resource. Patient Educ Couns 2005; 57(1): 115-121.

504. Thomas J, Harden A, Oakley A, Oliver S, Sutcliffe K, Rees R et al. Integrating qualitative research with trials in systematic reviews. BMJ 2004; 328(7446): 1010-1012.

505. Thomas S. Klinische Relevanz von Therapieeffekten: systematische Sichtung, Klassifizierung und Bewertung methodischer Konzepte [Dissertation]. Duisburg/Essen: Universität; 2009.

506. Thompson SG, Higgins JPT. How should meta-regression analyses be undertaken and interpreted? Stat Med 2002; 21(11): 1559-1573.

507. Thorlund K, Devereaux PJ, Wetterslev J, Guyatt G, Ioannidis JPA, Thabane L et al. Can trial sequential monitoring boundaries reduce spurious inferences from meta-analyses? Int J Epidemiol 2009; 38(1): 276-286.

508. Thorne S. The role of qualitative research within an evidence-based context: can metasynthesis be the answer? Int J Nurs Stud 2009; 46(4): 569-575.

509. Thurow S. Search engine visibility. Indianapolis: New Riders; 2003.

510. Tian L, Cai T, Pfeffer MA, Piankov N, Cremieux PY, Wei LJ. Exact and efficient inference procedure for meta-analysis and its application to the analysis of independent $2 \times 2$ tables with all available data but without artificial continuity correction. Biostatistics 2009; 10(2): 275-281.

511. Tibshirani R. Regression shrinkage and selection via the lasso. J R Stat Soc Series B Stat Methodol 1996; 58(1): 267-288.

512. Tierney JF, Stewart LA, Ghersi D, Burdett S, Sydes MR. Practical methods for incorporating summary time-to-event data into meta-analysis. Trials 2007; 8: 16.

513. Torgerson DJ. Contamination in trials: is cluster randomisation the answer? BMJ 2001; 322(7282): 355-357.

514. Trevena LJ, Davey HM, Barratt A, Butow P, Caldwell P. A systematic review on communicating with patients about evidence. J Eval Clin Pract 2006; 12(1): 13-23.

515. Tricoci P, Allen JM, Kramer JM, Califf RM, Smith SC Jr. Scientific evidence underlying the ACC/AHA clinical practice guidelines. JAMA 2009; 301(8): 831-841.

516. Tsay MY, Yang YH. Bibliometric analysis of the literature of randomized controlled trials. J Med Libr Assoc 2005; 93(4): 450-458.

517. Tukey JW. Some thoughts on clinical trials, especially problems of multiplicity. Science 1977; 198(4318): 679-684.

518. Tunis SR, Stryer DB, Clancy CM. Practical clinical trials: increasing the value of clinical research for decision making in clinical and health policy. JAMA 2003; 290(12): 1624-1632.

519. Turner RM, Omar RZ, Yang M, Goldstein H, Thompson SG. A multilevel model framework for meta-analysis of clinical trials with binary outcomes. Stat Med 2000; 19(24): 3417-3432.

520. Tyczynski JE, Demaret E, Parkin DM (Ed). Standards and guidelines for cancer registration in Europe: the ENCR recommendations vol.1. Lyon: IARC Press; 2003. (IARC Technical Publications; Volume 40).

521. UK National Screening Committee. Programme appraisal criteria: criteria for appraising the viability, effectiveness and appropriateness of a screening programme [online]. [Accessed: 18 April 2013]. URL: http://www.screening.nhs.uk/criteria.

522. USAID Center for Development Information and Evaluation. Conducting key informant interviews [online]. 1996 [accessed: 7 May 2013]. (Performance Monitoring and Evaluation TIPS; Volume 2). URL: http://pdf.usaid.gov/pdf_docs/PNABS541.pdf.

523. Van den Brink-Muinen A, Verhaak PF, Bensing JM, Bahrs O, Deveugele M, Gask L et al. Doctor-patient communication in different European health care systems: relevance and performance from the patients' perspective. Patient Educ Couns 2000; 39(1): 115-127.

524. Van Houwelingen HC, Arends LR, Stijnen T. Advanced methods in meta-analysis: multivariate approach and meta-regression. Stat Med 2002; 21(4): 589-624.

525. Van Tinteren H, Hoekstra OS, Smit EF, Van den Bergh JH, Schreurs AJ, Stallaert RA et al. Effectiveness of positron emission tomography in the preoperative assessment of patients with suspected non-small-cell lung cancer: the PLUS multicentre randomised trial. Lancet 2002; 359(9315): 1388-1393.

526. Van Tulder M, Furlan A, Bombardier C, Bouter L. Updated method guidelines for systematic reviews in the Cochrane Collaboration Back Review Group. Spine (Phila Pa 1976) 2003; 28(12): 1290-1299.

527. Vandenbroucke JP, Von Elm E, Altman DG, Gøtzsche PC, Mulrow CD, Pocock SJ et al. Strengthening the reporting of observational studies in epidemiology (STROBE): explanation and elaboration. Ann Intern Med 2007; 147(8): W163-W194.

528. Vergouwe Y, Royston P, Moons KG, Altman DG. Development and validation of a prediction model with missing predictor data: a practical approach. J Clin Epidemiol 2010; 63(2): 205-214.

529. Vidanapathirana J, Abramson MJ, Forbes A, Fairley C. Mass media interventions for promoting HIV testing. Cochrane Database Syst Rev 2005; (3): CD004775.

530. Vijan S. Should we abandon QALYs as a resource allocation tool? Pharmacoeconomics 2006; 24(10): 953-954.

531. Villar J, Mackey ME, Carroli G, Donner A. Meta-analyses in systematic reviews of randomized controlled trials in perinatal medicine: comparison of fixed and random effects models. Stat Med 2001; 20(23): 3635-3647.

532. Virtanen H, Leino-Kilpi H, Salantera S. Empowering discourse in patient education. Patient Educ Couns 2007; 66(2): 140-146.

533. Vlayen J, Aertgeerts B, Hannes K, Sermeus W, Ramaekers D. A systematic review of appraisal tools for clinical practice guidelines: multiple similarities and one common deficit. Int J Qual Health Care 2005; 17(3): 235-242.

534. Voisin CE, De la Varre C, Whitener L, Gartlehner G. Strategies in assessing the need for updating evidence-based guidelines for six clinical topics: an exploration of two search methodologies. Health Info Libr J 2008; 25(3): 198-207.

535. Von Elm E, Altman DG, Egger M, Pocock SJ, Gøtzsche PC, Vandenbroucke JP. The strengthening the reporting of observational studies in epidemiology (STROBE) statement: guidelines for reporting observational studies. Ann Intern Med 2007; 147(8): 573-577.

536. Wallis EJ, Ramsay LE, Ul Haq I, Ghahramani P, Jackson PR, Rowland-Yeo K et al. Coronary and cardiovascular risk estimation for primary prevention: validation of a new Sheffield table in the 1995 Scottish health survey population. BMJ 2000; 320(7236): 671-676.

537. Walter U, Schwartz FW. Prävention. In: Schwartz FW, Badura B, Busse R, Leidl R, Raspe H, Siegrist J et al (Ed). Das Public Health Buch: Gesundheit und Gesundheitswesen. München: Urban und Fischer; 2003. p. 189-214.

538. Wantland DJ, Portillo CJ, Holzemer WL, Slaughter R, McGhee EM. The effectiveness of web-based vs. non-web-based interventions: a meta-analysis of behavioral change outcomes. J Med Internet Res 2004; 6(4): e40.

539. Watine J, Friedberg B, Nagy E, Onody R, Oosterhuis W, Bunting PS et al. Conflict between guideline methodologic quality and recommendation validity: a potential problem for practitioners. Clin Chem 2006; 52(1): 65-72.

540. Web Accessibility Initiative. Website [online]. [Accessed: 22 April 2013]. URL: http://www.w3.org/WAI.

541. Weir CJ, Walley RJ. Statistical evaluation of biomarkers as surrogate endpoints: a literature review. Stat Med 2006; 25(2): 183-203.

542. Wendt C. Gesundheitssysteme im internationalen Vergleich. Gesundheitswesen 2006; 68(10): 593-599.

543. Wetterslev J, Thorlund K, Brok J, Gluud C. Trial sequential analysis may establish when firm evidence is reached in cumulative meta-analysis. J Clin Epidemiol 2008; 61(1): 64-75.

544. Whitehead J. The design and analysis of sequential clinical trials. Chichester: Horwood; 1983.

545. Whiting P, Rutjes AW, Reitsma JB, Glas AS, Bossuyt PM, Kleijnen J. Sources of variation and bias in studies of diagnostic accuracy: a systematic review. Ann Intern Med 2004; 140(3): 189-202.

546. Whiting P, Rutjes AWS, Reitsma JB, Bossuyt PM, Kleijnen J. The development of QUADAS: a tool for the quality assessment of studies of diagnostic accuracy included in systematic reviews. BMC Med Res Methodol 2003; 3: 25.

547. Whitlock EP, Lin JS, Chou R, Shekelle P, Robinson KA. Using existing systematic reviews in complex systematic reviews. Ann Intern Med 2008; 148(10): 776-782.

548. Wieseler B, McGauran N, Kaiser T. Still waiting for functional EU Clinical Trials Register. BMJ 2011; 342: d3834.

549. Windeler J. Bedeutung randomisierter klinischer Studien mit relevanten Endpunkten für die Nutzenbewertung [online]. In: Gesundheitsforschungsrat des Bundesministeriums für Bildung und Forschung (Ed). Diskussionsforum zur Nutzenbewertung im Gesundheitswesen: Begriffsdefinitionen und Einführung; Dokumentation des ersten gemeinsamen Workshops von GFR und IQWiG am 4. September 2007 in Berlin. December 2007. S. 26-31 [accessed: 19 April 2013]. URL: http://www.gesundheitsforschung-bmbf.de/_media/DLR_Nutzenbewert_07-11-22_Druckversion.pdf.

550. Windeler J. Externe Validität. Z Evid Fortbild Qual Gesundhwes 2008; 102(4): 253-259.

551. Windeler J, Conradt C. Wie können "Signifikanz" und "Relevanz" verbunden werden? Med Klin 1999; 94(11): 648-651.

552. Windeler J, Lange S. Nutzenbewertung in besonderen Situationen: seltene Erkrankungen. Z Evid Fortbild Qual Gesundhwes 2008; 102(1): 25-30.

553. Windeler J, Ziegler S. Evidenzklassifizierungen. Z Arztl Fortbild Qualitatssich 2003; 97(6): 513-514.

554. Wofford JL, Smith ED, Miller DP. The multimedia computer for office-based patient education: a systematic review. Patient Educ Couns 2005; 59(2): 148-157.

555. Wood L, Egger M, Gluud LL, Schulz KF, Jüni P, Altman DG et al. Empirical evidence of bias in treatment effect estimates in controlled trials with different interventions and outcomes: meta-epidemiological study. BMJ 2008; 336(7644): 601-605.

556. Ziebland S, McPherson A. Making sense of qualitative data analysis: an introduction with illustrations from DIPEx (personal experiences of health and illness). Med Educ 2006; 40(5): 405-414.

557. Ziegler DK, Mosier MC, Buenaver M, Okuyemi K. How much information about adverse effects of medication do patients want from physicians? Arch Intern Med 2001; 161(5): 706-713.

558. Zschorlich B, Knelangen M, Bastian H. Die Entwicklung von Gesundheitsinformationen unter Beteiligung von Bürgerinnen und Bürgern am Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen (IQWiG). Gesundheitswesen 2011; 73(7): 423-429.

559. Zwarenstein M, Treweek S, Gagnier JJ, Altman DG, Tunis S, Haynes B et al. Improving the reporting of pragmatic trials: an extension of the CONSORT statement. BMJ 2008; 337: a2390.