

Otitis media with effusion in under 12s

NICE guideline: methods

NICE guideline NG233

Supplement 1: Methods

August 2023

Final

*These supplements were developed by
NICE*

Disclaimer

The recommendations in this guideline represent the view of NICE, arrived at after careful consideration of the evidence available. When exercising their judgement, professionals are expected to take this guideline fully into account, alongside the individual needs, preferences and values of their patients or service users. The recommendations in this guideline are not mandatory and the guideline does not override the responsibility of healthcare professionals to make decisions appropriate to the circumstances of the individual patient, in consultation with the patient and/or their carer or guardian.

Local commissioners and/or providers have a responsibility to enable the guideline to be applied when individual health professionals and their patients or service users wish to use it. They should do so in the context of local and national priorities for funding and developing services, and in light of their duties to have due regard to the need to eliminate unlawful discrimination, to advance equality of opportunity and to reduce health inequalities. Nothing in this guideline should be interpreted in a way that would be inconsistent with compliance with those duties.

NICE guidelines cover health and care in England. Decisions on how they apply in other UK countries are made by ministers in the [Welsh Government](#), [Scottish Government](#), and [Northern Ireland Executive](#). All NICE guidance is subject to regular review and may be updated or withdrawn.

Copyright

© NICE, 2023. All rights reserved. Subject to Notice of rights.

ISBN: 978-1-4731-5346-2

Contents

Development of the guideline	5
Remit.....	5
What this guideline covers	5
What this guideline does not cover	5
Methods	6
Developing the review questions and outcomes.....	6
Searching for evidence	8
Scoping search	8
Systematic literature search.....	8
Economic systematic literature search	9
Reviewing research evidence	9
Systematic review process	9
Type of studies and inclusion/exclusion criteria.....	10
Methods of combining evidence.....	11
Data synthesis for intervention studies	11
Data synthesis for epidemiological reviews	13
Data synthesis for diagnostic test accuracy reviews.....	13
Data synthesis for prognostic reviews	14
Data synthesis for qualitative reviews.....	14
Appraising the quality of evidence	14
Intervention studies.....	14
Epidemiological studies	20
Prognostic studies	21
Diagnostic studies.....	23
Qualitative studies	25
Reviewing economic evidence.....	29
Appraising the quality of economic evidence.....	29
Economic modelling.....	29
Cost effectiveness criteria.....	30
Developing recommendations.....	30
Guideline recommendations	30
Research recommendations	31
Validation process.....	31
Updating the guideline	31
References	32

Development of the guideline

Remit

The National Institute for Health and Care Excellence (NICE) updated the following clinical guideline:

Title: Otitis media with effusion in under 12s: surgery.

Surgery was removed from the title of the guideline, because the scope of the update covers other topics as well as surgery.

What this guideline covers

Groups that will be covered:

All children under 12 years with suspected or confirmed otitis media with effusion (OME).

Settings that will be covered:

All settings where NHS-commissioned care is provided.

Key areas that will be covered:

1. Risk factors for OME.
2. Recognition of OME (to help identify when to refer for further investigation).
3. Natural history of OME (to help identify when intervention and follow-up is needed).
4. Interventions for children with OME.
5. Care during and after surgery.
6. Information for children, parents and carers.

What this guideline does not cover

1. Diagnosing or managing acute otitis media.
 - This is a different condition and is covered by the NICE guideline on antimicrobial prescribing for acute otitis media.
2. Specific methods of assessing hearing in children.
 - This is not specific to OME.

Methods

This guideline was developed using the methods described in the 2018 NICE guidelines manual. Declarations of interest were recorded according to the NICE conflicts of interest policy.

Developing the review questions and outcomes

The review questions developed for this guideline were based on the key areas identified in the guideline [scope](#). They were drafted by the NICE technical team, and refined and validated by the guideline committee.

The review questions were based on the following frameworks:

- population, intervention, comparator and outcome (PICO) for reviews of interventions and epidemiological reviews
- diagnostic reviews – using population, diagnostic test (index test), reference standard and target condition (PIRT)
- prognostic reviews – using population, presence or absence of a prognostic, risk or predictive factor and outcome (PPO)
- qualitative reviews – using population, phenomenon of interest and context (PICo)

Full literature searches, critical appraisals and evidence reviews were completed for all review questions.

The review questions and evidence reviews corresponding to each question (or group of questions) are summarised below.

Table 1: Summary of review questions and index to evidence reviews

Evidence review	Review question	Type of review
[A] Modifiable risk factors for developing OME in children	What are the modifiable risk factors for developing OME in children under 12 years?	Prognostic
[B] Presenting features associated with OME in children	What presenting features are associated with OME in children under 12 years?	Diagnostic
[C] Natural history of OME without hearing loss	What is the progression, resolution and recurrence (natural history) of OME without hearing loss at presentation in children under 12 years?	Epidemiological
[D] Natural history of OME-related hearing loss	What is the progression, resolution and recurrence (natural history) of OME-related hearing loss at presentation in children under 12 years?	Epidemiological

Evidence review	Review question	Type of review
[E] Ventilation tubes for children with OME	What is the effectiveness of ventilation tubes for managing otitis media with effusion (OME) with associated hearing loss in children under 12 years?	Intervention ¹
[F] Adenoidectomy for children with OME	What is the effectiveness of adenoidectomy (with or without ventilation tubes) for managing OME with associated hearing loss in children under 12 years?	Intervention
[G] Antibiotics for children with OME	What is the effectiveness of antibiotics for managing OME in children under 12 years?	Intervention
[H] Non-antimicrobial pharmacological interventions for children with OME	What is the effectiveness of non-antimicrobial pharmacological interventions (such as steroids, antihistamines, leukotriene receptor antagonists, mucolytics and decongestants) for managing OME in children under 12 years?	Intervention
[I] Auto-inflation for children with OME	What is the effectiveness of auto-inflation for managing OME with associated hearing loss in children under 12 years?	Intervention
[J] Hearing aids/devices for hearing loss associated with OME in children under 12 years ¹	What is the effectiveness of air conduction and bone conduction hearing aids/devices for hearing loss associated with OME in children under 12 years?	Intervention
[K] Intraoperative or postoperative interventions for preventing otorrhoea after surgery for hearing loss associated with OME in children	What intraoperative or postoperative interventions are effective at preventing otorrhoea (ear discharge) after surgery for otitis media with effusion (OME)-related hearing loss in children under 12 years?	Intervention
[L] Interventions for treating otorrhoea after surgery for hearing loss	What interventions are effective for treating otorrhoea (ear discharge) after surgery for otitis media with effusion (OME)-related hearing loss in children under 12 years?	Intervention

Evidence review	Review question	Type of review
associated with OME in children		
[M] Follow-up strategy after surgical treatment for OME-related hearing loss	What should the follow-up strategy be after surgical treatment for OME-related hearing loss in children under 12 years?	Intervention
[N] Information for suspected or confirmed OME	What information is valued by children under 12 years with suspected or confirmed otitis media with effusion (OME) and their parents and carers?	Qualitative

¹Original health economic analysis conducted
 OME: otitis media with effusion

The COMET database was searched for core outcome sets relevant to this guideline and a core outcome set was identified (Liu 2020) which informed the protocols.

Additional information related to development of the guideline is contained in:

- NGA staff list
- Supplement 3 (Previous guideline evidence related to Other non-surgical interventions).

Searching for evidence

Scoping search

During the scoping phase, searches were conducted for previous guidelines, economic evaluations, health technology assessments, systematic reviews, randomised controlled trials, observational studies and qualitative research.

Systematic literature search

Systematic literature searches were undertaken to identify published evidence relevant to each review question.

Databases were searched using subject headings, free-text terms and, where appropriate, study type filters. Where possible, searches were limited to retrieve studies published in English. Limits to exclude animal studies, letters, editorials, news and conferences were applied where possible. All the searches were conducted in the following databases: Medline, Cochrane Central Register of Controlled Trials (CENTRAL), Cochrane Database of Systematic Reviews (CDSR), Embase, Epistemonikos and International Network of Agencies for Health Technology Assessments (INAHTA) .

Searches were run once for all reviews during development. Searches for the following questions were updated in November 2022, 5 weeks in advance of the final committee meeting.

- [A] Modifiable risk factors for developing OME in children

-
- [J] Hearing aids/devices for hearing loss associated with OME in children under 12 years
 - [K] Intraoperative or postoperative interventions for preventing otorrhoea after surgery for hearing loss associated with OME in children
 - [L] Interventions for treating otorrhoea after surgery for hearing loss associated with OME in children
 - [M] Follow-up strategy after surgical treatment for OME-related hearing loss
 - [N] Information for suspected or confirmed OME

Details of the search strategies, including the study-design filters used and databases searched, are provided in Appendix B of each evidence review.

Economic systematic literature search

Systematic literature searches were also undertaken to identify published economic evidence. Databases were searched using subject headings, free-text terms and, where appropriate, an economic evaluations search filter.

A single search, using the population search terms used in the evidence reviews, was conducted to identify economic evidence in the NHS Economic Evaluation Database (NHS EED) and the International Network of Agencies for Health Technology Assessments (INAHTA) database. Another single search, using the population search terms used in the evidence reviews combined with an economic evaluations search filter, was conducted in Medline, Cochrane Central Register of Controlled Trials (CENTRAL), and Embase. Where possible, searches were limited to studies published in English. Limits to exclude animal studies, letters, editorials, news were applied where possible.

As with the general literature searches, the economic literature searches were updated in November 2022, 5 weeks in advance of the final committee meeting before consultation on the draft guideline.

Details of the search strategies, including the study-design filters used and databases searched, are provided in Appendix B of each evidence review.

Quality assurance

Search strategies were quality assured by cross-checking reference lists of relevant studies, analysing search strategies from published systematic reviews and asking members of the committee to highlight key studies. The principal search strategies for each search were also quality assured by a second information scientist using an adaptation of the PRESS 2015 Guideline Evidence-Based Checklist (McGowan 2016). In addition, all publications highlighted by stakeholders at the time of the consultation on the draft scope were considered for inclusion.

Reviewing research evidence

Systematic review process

The evidence was reviewed in accordance with the following approach.

-
- Potentially relevant articles were identified from the search results for each review question by screening titles and abstracts. Full-text copies of the articles were then obtained.
 - Full-text articles were reviewed against pre-specified inclusion and exclusion criteria in the review protocol (see Appendix A of each evidence review).
 - Key information was extracted from each article on study methods and results, in accordance with factors specified in the review protocol. The information was presented in a summary table in the corresponding evidence review and in a more detailed evidence table (see Appendix D of each evidence review).
 - Included studies were critically appraised using an appropriate checklist as specified in [Developing NICE guidelines: the manual](#). Further detail on appraisal of the evidence is provided below.
 - Summaries of quantitative evidence by outcome and qualitative evidence by theme were presented in the corresponding evidence review and discussed by the committee.

Internal quality assurance processes included consideration of the outcomes of screening, study selection and data extraction and the committee reviewed the results of study selection and data extraction. The review protocol for each question specifies whether dual screening and study selection was undertaken for that particular question. Drafts of all evidence reviews were quality assured by a senior reviewer.

Type of studies and inclusion/exclusion criteria

Inclusion and exclusion of studies was based on criteria specified in the corresponding review protocol. If at least 75% of the population of a study matched the population specified in the protocol, the evidence was considered to be from a directly relevant population. If less than 50% of the population of a study matched the population specified in the protocol and data was not presented separately for the relevant population, the study was excluded.

Systematic reviews with meta-analyses or meta-syntheses were considered to be the highest quality evidence that could be selected for inclusion.

For intervention reviews, systematic reviews of randomised controlled trials (RCTs) and primary RCTs were prioritised for inclusion because they are considered to be the most robust type of study design that could produce an unbiased estimate of intervention effects. Where there was insufficient evidence from RCTs to inform guideline decision making, non-randomised studies (NRS) were considered for inclusion. Sufficiency was judged taking into account the number, quality and sample size of RCTs, as well as outcomes reported and availability of data from subgroups of interest. When NRS were considered for inclusion, priority was given to controlled studies, with separate control groups that were not allocated on the basis of the outcome, that adjusted for relevant confounders or matched participants on important confounding domains.

For epidemiological reviews, observational studies (non-comparative studies or untreated control arms from comparative studies) were prioritised for inclusion because they are less likely than experimental studies to have strict eligibility criteria that could restrict the population of interest. Where there was insufficient evidence from observational studies to inform guideline decision making, untreated control

arms from experimental studies were considered for inclusion; where there was insufficient evidence from observational and experimental studies, case series were considered for inclusion. Sufficiency was judged taking into account outcomes reported and availability of data from subgroups of interest.

For diagnostic reviews, cross-sectional diagnostic accuracy studies were prioritised for inclusion. Studies that used single-gate designs were prioritised.

For prognostic reviews, prospective and retrospective cohort studies were considered for inclusion. Studies that included multivariable analysis were prioritised.

For qualitative reviews, systematic reviews of qualitative studies and primary qualitative studies using focus groups, structured interviews or semi-structured interviews, observations or surveys with open-ended questions were considered for inclusion. Where qualitative evidence was sought, data from surveys or other types of questionnaires were considered for inclusion only if they provided data from open-ended questions, but not if they reported only quantitative data.

The committee was consulted about any uncertainty regarding inclusion or exclusion of studies. A list of excluded studies for each review question, including reasons for exclusion is presented in Appendix J of the corresponding evidence review.

Narrative reviews, posters, letters, editorials, comment articles, unpublished studies and studies published in languages other than English were only included in the reviews done by Cochrane. Conference abstracts were not considered for inclusion because conference abstracts typically do not have sufficient information to allow for full critical appraisal.

Methods of combining evidence

When planning reviews (through preparation of protocols), the following approaches for data synthesis were discussed and agreed with the committee.

Data synthesis for intervention studies

Pairwise meta-analysis

Meta-analysis to pool results from comparative intervention studies was conducted where possible using Cochrane Review Manager (RevMan5) software.

For dichotomous outcomes, such as mortality, the Mantel–Haenszel method with a fixed effect model was used to calculate risk ratios (RRs). For all outcomes with zero events in both arms the risk difference was presented. For outcomes in which the majority of studies had low event rates (<1%), Peto odds ratios (ORs) were calculated as this method performs well when events are rare (Bradburn 2007).

For continuous outcomes, measures of central tendency (mean) and variation (standard deviation; SD) are required for meta-analysis. Data for continuous outcomes, such as quality of life, were meta-analysed using an inverse-variance method for pooling weighted mean differences (WMDs). Where SDs were not reported for each intervention group, the standard error (SE) of the mean difference was calculated from other reported statistics (p values or 95% confidence intervals; CIs) and then meta-analysis was conducted as described above.

If a study reported only the summary statistic and 95% CI the generic-inverse variance method was used to enter data into RevMan5. If the control event rate was reported this was used to generate the absolute risk difference in GRADEpro. Where a study reported multiple adjusted estimates for the same outcome, the one that minimised the risk of bias due to confounding was chosen.

When evidence was based on studies that reported descriptive data or medians with interquartile ranges or p values, this information was included in the corresponding GRADE tables (see below) without calculating relative or absolute effects. Consequently, certain aspects of quality assessment such as imprecision of the effect estimate could not be assessed as per standard methods for this type of evidence and subjective ratings or ratings based on sample size cut-offs were considered instead.

For some reviews, evidence was either stratified from the outset or separated into subgroups when heterogeneity was encountered. The stratifications and potential subgroups were pre-defined at the protocol stage (see the protocols for each review for further detail). Where evidence was stratified or subgrouped the committee considered on a case by case basis if separate recommendations should be made for distinct groups. Separate recommendations may be made where there is evidence of a differential effect of interventions in distinct groups. If there is a lack of evidence in one group, the committee considered, based on their experience, whether it was reasonable to extrapolate and assume the interventions will have similar effects in that group compared with others.

Data from RCTs and NRS, or from NRS with substantially different designs (i.e., cohort studies and case-control studies), that were theoretically possible to pool were entered into RevMan5 as subgroups based on study design. This was to take into account the likelihood of increased heterogeneity from studies with different design features and different approaches to appraising the quality of evidence based on study design (see appraising the quality of evidence: intervention studies below).

When meta-analysis was undertaken, the results were presented visually using forest plots generated using RevMan5 (see Appendix E of relevant evidence reviews).Included Cochrane Reviews

During the development of this guideline, 5 registered Cochrane protocols were identified which matched the committee's intended review questions:

- [E] Ventilation tubes for children with OME
- [F] Adenoidectomy for children with OME
- [G] Antibiotics for children with OME
- [H] Non-antimicrobial pharmacological interventions for children with OME
- [I] Auto-inflation for children with OME.

The Cochrane review team completed 5 reviews investigating the effectiveness of ventilation tubes (MacKeith 2023a), adenoidectomy (MacKeith 2023b), antibiotics (Mulvaney 2023a), steroids (Mulvaney 2023b), and auto-inflation (Webster 2023) for children with OME during guideline development and presented their results to the guideline committee, who used them to make recommendations.

Cochrane's methods are closely aligned to standard NICE methods, minor deviations (the use of GRADE only on main outcomes, summary of findings tables instead of full

GRADE tables, defining primary and secondary outcomes as opposed to critical and important, assessing the risk of bias in primary studies using version 1 (as opposed to version 2) of the Cochrane Risk of Bias tool (Higgins 2011), how clinically important differences are determined, and including countries from a broader range of income categories than the majority of the other reviews in the guideline) relevant to the topic area were highlighted to the committee and taken into account in discussions of the evidence. Full details of the Cochrane review, including methods, are available in the above reviews.

Data synthesis for epidemiological reviews

Proportion data

Meta-analysis to pool proportion data was conducted where possible using the metafor package in R (Viechtbauer 2010), which allows meta-analysing of data from single group studies.

Evidence was stratified based on certain parameters at the outset and further separated into subgroups when heterogeneity was encountered. The stratifications and potential subgroups were pre-defined at the protocol stage (see the protocols for each review for further detail). Where evidence was stratified or subgrouped the committee considered on a case by case basis if separate recommendations should be made for distinct groups. Separate recommendations may be made where there is evidence of differences in natural history in distinct groups. If there is a lack of evidence in one group, the committee considered, based on their experience, whether it was reasonable to extrapolate and assume the natural history will be similar in that group compared with others.

When meta-analysis was undertaken, the results were presented visually using forest plots generated using R (see Appendix E of relevant evidence reviews).

Time-to-event data

The intention was to pool time-to-event data and present the results as summary survival curves using the metaSurvival package in R (Pandey 2020). However, there was insufficient time-to-event data available to generate summary survival curves. Therefore, such data were converted to proportion data to allow for direct comparison, and where applicable pooling, with the proportion data.

Data synthesis for diagnostic test accuracy reviews

When diagnostic test accuracy was measured dichotomously, sensitivity, specificity, and positive and negative predictive values were used as outcomes. Where possible, diagnostic accuracy parameters were calculated by the NICE technical team using data from 2x2 tables reported in the articles; alternatively, parameters were obtained directly from results reported in the source articles. For sensitivity and specificity, 95% CIs were reported.

No meta-analyses of diagnostic test accuracy data were possible due to insufficient similarities in index tests, populations and reference standards across studies.

Data synthesis for prognostic reviews

ORs or RRs with 95% CIs reported in published studies were extracted or calculated by the NICE technical team to examine relationships between risk factors and outcomes of interest. Ideally analyses would have adjusted for key confounders (such as age or severity of hearing loss) to be considered for inclusion, but studies with unadjusted analyses were included if data from adjusted analyses was insufficient for guideline decision making.

No meta-analyses of prognostic data were possible due to variation across studies in terms of risk factors and outcomes.

Data synthesis for qualitative reviews

Where possible, a meta-synthesis was conducted to combine evidence from more than one study into a theme or sub-theme. Whenever studies identified a qualitative theme relevant to the protocol, this was extracted and the main characteristics were summarised. When all themes had been extracted from studies, common concepts were categorised and tabulated. This included information on how many studies had contributed to each theme identified by the NICE technical team.

The technical team were guided in their data extraction, synthesis and formulation of review findings, or themes, by a framework of phenomena developed by the guideline committee. This framework consisted of the themes that the committee anticipated would be covered by the included studies and these were set out a priori in the corresponding review protocol. Themes identified from the included studies, which were not set out in the protocol but which were considered relevant to answering the review question, were also extracted

Themes from individual studies were integrated into a wider context and, when possible, overarching categories of themes with sub-themes were identified. Themes were derived from data presented in individual studies and sub-theme and theme names were assigned by the NICE technical team.

Emerging themes were placed into a thematic map representing the relationship between themes and overarching categories. The purpose of such a map is to show relationships between overarching categories and associated themes.

Appraising the quality of evidence

Intervention studies

Pairwise meta-analysis

GRADE methodology for intervention reviews

For intervention reviews, the evidence for outcomes from included RCTs and comparative non-randomised studies was evaluated and presented using the Grading of Recommendations Assessment, Development and Evaluation (GRADE) methodology developed by the international GRADE working group.

When GRADE was applied, software developed by the GRADE working group (GRADEpro) was used to assess the quality of each outcome, taking account of

individual study quality factors and any meta-analysis results. Results were presented in GRADE profiles (GRADE tables).

The selection of outcomes for each review question was agreed during development of the associated review protocol in discussion with the committee. The evidence for each outcome was examined separately for the quality elements summarised in Table 2. Criteria considered in the rating of these elements are discussed below. Each element was graded using the quality ratings summarised in Table 3. Footnotes to GRADE tables were used to record reasons for grading a particular quality element as having a 'serious' or 'very serious' quality issue. The ratings for each component were combined to obtain an overall assessment of quality for each outcome as described in Table 4.

The initial quality rating was based on the study design: RCTs and NRS assessed by ROBINS-I start as 'high' quality evidence, other non-randomised studies start as 'low' quality evidence. The rating was then modified according to the assessment of each quality element (Table 2). Each quality element considered to have a 'serious' or 'very serious' quality issue was downgraded by 1 or 2 levels respectively (for example, evidence starting as 'high' quality was downgraded to 'moderate' or 'low' quality). In addition, there was a possibility to upgrade evidence from non-randomised studies (provided the evidence for that outcome had not previously been downgraded) if there was a large magnitude of effect, a dose–response gradient, or if all plausible confounding would reduce a demonstrated effect or suggest a spurious effect when results showed no effect.

Table 2: Summary of quality elements in GRADE for intervention reviews

Quality element	Description
Risk of bias ('Study limitations')	This refers to limitations in study design or implementation that reduce the internal validity of the evidence
Inconsistency	This refers to unexplained heterogeneity in the results
Indirectness	This refers to differences in study populations, interventions, comparators or outcomes between the available evidence and inclusion criteria specified in the review protocol
Imprecision	This occurs when a study has few participants or few events of interest, resulting in wide confidence intervals that cross minimally important thresholds
Publication bias	This refers to systematic under- or over-estimation of the underlying benefit or harm resulting from selective publication of study results

Table 3: GRADE quality ratings (by quality element)

Quality issues	Description
None or not serious	No serious issues with the evidence for the quality element under consideration
Serious	Issues with the evidence sufficient to downgrade by 1 level for the quality element under consideration
Very serious	Issues with the evidence sufficient to downgrade by 2 levels for the quality element under consideration

Table 4: Overall quality of the evidence in GRADE (by outcome)

Overall quality grading	Description
High	Further research is very unlikely to change the level of confidence in the estimate of effect
Moderate	Further research is likely to have an important impact on the level of confidence in the estimate of effect and may change the estimate
Low	Further research is very likely to have an important impact on the level of confidence in the estimate of effect and is likely to change the estimate
Very low	The estimate of effect is very uncertain

Assessing risk of bias in intervention reviews

Bias is a systematic error, or consistent deviation from the truth in results obtained. When a risk of bias is present the true effect can be either under- or over-estimated.

Risk of bias in RCTs was assessed using the Cochrane risk of bias 2.0 tool (see Appendix H in Developing NICE guidelines: the manual).

The Cochrane risk of bias tool assesses the following possible sources of bias:

- randomisation process
- deviations from the intended interventions
- missing outcome data
- measurement of the outcome
- selection of the reported result.

A study with a poor methodological design does not automatically imply high risk of bias; the bias is considered individually for each outcome and it is assessed whether the chosen design and methodology will impact on the estimation of the intervention effect.

More details about the Cochrane risk of bias 2.0 tool can be found in Section 8 of the Cochrane Handbook for Systematic Reviews of Interventions (Higgins 2011).

For systematic reviews the ROBIS checklist was used (see Appendix H in Developing NICE guidelines: the manual).

For non-randomised controlled studies, cohort studies or historical controlled studies the ROBINS-I checklist was used (see Appendix H in Developing NICE guidelines: the manual). *Assessing inconsistency in intervention reviews*

Inconsistency refers to unexplained heterogeneity in results of meta-analysis. When estimates of treatment effect vary widely across studies (that is, there is heterogeneity or variability in results), this suggests true differences in underlying effects. Inconsistency is, thus, only truly applicable when statistical meta-analysis is conducted (that is, results from different studies are pooled). When outcomes were derived from a single study the rating 'no serious inconsistency' was used when assessing this domain, as per GRADE methodology (Santesso 2016).

Inconsistency was assessed visually by inspecting forest plots and observing whether there was considerable heterogeneity in the results of the meta-analysis (for

example if the point estimates of the individual studies consistently showed benefits or harms). This was supported by calculating the I-squared statistic for the meta-analysis with an I-squared value of more than 50% indicating serious heterogeneity, and more than 80% indicating very serious heterogeneity. When serious or very serious heterogeneity was observed, possible reasons were explored and subgroup analyses were performed as pre-specified in the review protocol where possible.

When no plausible explanation for the serious or very serious heterogeneity could be found, the quality of the evidence was downgraded in GRADE for inconsistency and the meta-analysis was re-run using the Der-Simonian and Laird method with a random effects model and this was used for the final analysis.

Assessing indirectness in intervention reviews

Directness refers to the extent to which populations, interventions, comparisons and outcomes reported in the evidence are similar to those defined in the inclusion criteria for the review and was assessed by comparing the PICO elements in the studies to the PICO defined in the review protocol. Indirectness is important when such differences are expected to contribute to a difference in effect size, or may affect the balance of benefits and harms considered for an intervention. Evidence was considered seriously indirect if 25% to 50% of the evidence for an outcome differed from one of the PICO elements of the protocol and very seriously indirect if more than 50% of the evidence for an outcome differed on one of the PICO elements from the protocol or if 25% to 50% of the evidence for an outcome differed on two or more of the PICO elements from the protocol.

Assessing imprecision and importance in intervention reviews

Imprecision in GRADE methodology refers to uncertainty around the effect estimate and whether or not there is an important difference between interventions (that is, whether the evidence clearly supports a particular recommendation or appears to be consistent with several candidate recommendations). Therefore, imprecision differs from other aspects of evidence quality because it is not concerned with whether the point estimate is accurate or correct (has internal or external validity). Instead, it is concerned with uncertainty about what the point estimate actually represents. This uncertainty is reflected in the width of the CI.

The 95% CI is defined as the range of values within which the population value will fall on 95% of repeated samples, were the procedure to be repeated. The larger the study, the smaller the 95% CI will be and the more certain the effect estimate.

Imprecision was assessed in the guideline evidence reviews by considering whether the width of the 95% CI of the effect estimate was relevant to decision making, considering each outcome independently. This is illustrated in Figure 1, which considers a positive outcome for the comparison of two treatments. Three decision-making zones can be differentiated, bounded by the thresholds for minimal importance (minimally important differences; MIDs) for benefit and harm.

When the CI of the effect estimate is wholly contained in 1 of the 3 zones there is no uncertainty about the size and direction of effect, therefore, the effect estimate is considered precise; that is, there is no imprecision.

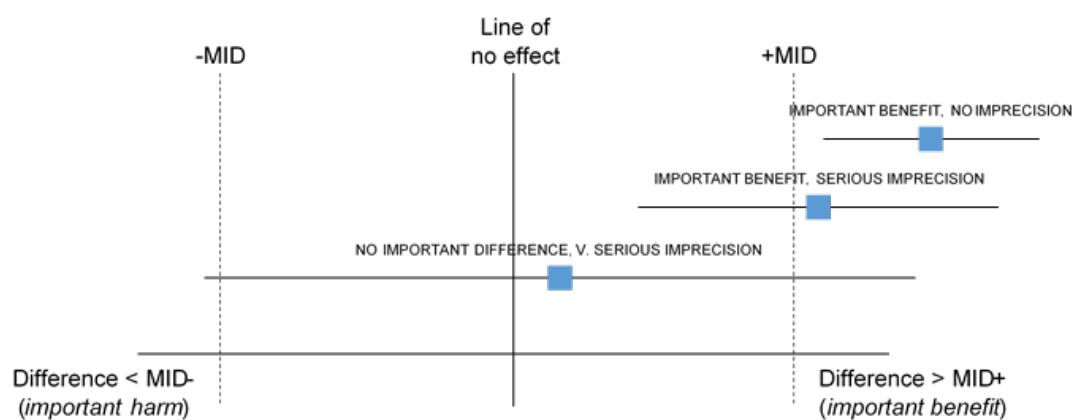
When the CI crosses 2 zones, it is uncertain in which zone the true value of the effect estimate lies and therefore there is uncertainty over which decision to make. The CI

is consistent with 2 possible decisions, therefore, the effect estimate is considered to be imprecise in the GRADE analysis and the evidence is downgraded by 1 level ('serious imprecision').

When the CI crosses all 3 zones, the effect estimate is considered to be very imprecise because the CI is consistent with 3 possible decisions and there is therefore a considerable lack of confidence in the results. The evidence is therefore downgraded by 2 levels in the GRADE analysis ('very serious imprecision').

Implicitly, assessing whether a CI is in, or partially in, an important zone, requires the guideline committee to estimate an MID or to say whether they would make different decisions for the 2 confidence limits.

Figure 1: Assessment of imprecision and importance in intervention reviews using GRADE



MID, minimally important difference

Defining minimally important differences for intervention reviews

The committee was asked whether there were any recognised or acceptable MIDs in the published literature and community relevant to the review questions under consideration. The committee was not aware of any MIDs that could be used for the guideline.

In the absence of published or accepted MIDs, the committee agreed to use the GRADE default MIDs to assess imprecision. For dichotomous outcomes minimally important thresholds for a RR of 0.8 and 1.25 respectively were used as default MIDs in the guideline. The committee also chose to use 0.8 and 1.25 as the MIDs for ORs in the absence of published or accepted MIDs. ORs were predominantly used in the guideline when Peto OR were indicated due to low event rates, at low event rates OR are mathematically similar to RR making the extrapolation appropriate.

If risk difference was used for meta-analysis, for example if the majority of studies had zero events in either arm, imprecision was assessed based on sample size using 200 and 400 as cut-offs for very serious and serious imprecision respectively. The committee used these numbers based on commonly used optimal information size thresholds.

The same thresholds were used as default MIDs in the guideline for all dichotomous outcomes considered in intervention evidence reviews. For continuous outcomes default MIDs are equal to half the SD of the control group at baseline (or at follow-up if the SD is not available a baseline).

MIDs, the line of no effect, and both 95% and 90% confidence intervals (CIs) were used to assess whether there were important differences in outcomes between groups. Outcomes were considered to have an important benefit/ harm, possible important benefit/ harm, no evidence of an important difference, or no important difference using the following approach:

- Where the point estimate (PE) is greater than the upper MID and the 95% CI do not cross line of no effect, an intervention was described as having an important benefit
- Where the PE is greater than the upper MID and the 95% CI do cross the line of no effect, but the 90% CI do not, an intervention was described as having a possible important benefit
- Where the PE is greater than the upper MID **or** lower than the lower MID, and the 90% CI cross the line of no effect, the result was described as no evidence of an important difference
- Where the PE is between two MIDs, the result was described as no important difference
- Where the PE is lower than the lower MID and the 95% CI do cross the line of no effect, but the 90% CI do not, an intervention is described as having a possible important harm
- Where the PE is lower than the lower MID and the 95% CI do not cross line of no effect, an intervention was described as having an important harm.

This approach was used for all evidence reviews which informed decision making on the guideline, including when interpreting results from evidence reviews conducted by the Cochrane Collaboration. Please note that the above descriptions are based on positive outcomes (where high values indicate better outcomes or events are positive). If the outcomes were negative (where high values indicate worse outcomes or events are negative) then whether an intervention is considered to have an important benefit or important harm would be switched (for example, where the PE is greater than the upper MID and the 95% CI do not cross line of no effect, an intervention would be described as having an important harm; where the PE is lower than the lower MID and the 95% CI do not cross line of no effect, an intervention would be described as having an important benefit).

90% CI are reported in the summary of the evidence section of the evidence reviews only when they were used to determine a possible importance difference (that is, when interventions had a possible important benefit/ harm).

Assessing publication bias in intervention reviews

Where 10 or more studies were included as part of a single meta-analysis, a funnel plot was produced to graphically assess the potential for publication bias. Where fewer than 10 studies were included for an outcome, the committee subjectively assessed the likelihood of publication bias based on factors such as the proportion of trials funded by industry and the propensity for publication bias in the topic area.

Epidemiological studies

Adapted GRADE methodology for prevalence reviews

For prevalence reviews with evidence from comparative studies an adapted GRADE approach was used. As noted above, GRADE methodology is designed for intervention reviews, but the quality assessment elements were adapted for epidemiological reviews.

The evidence for each outcome in the epidemiological reviews was examined separately for the quality elements listed and defined in Table 6. The criteria considered in the rating of these elements are discussed below. Each element was graded using the quality levels summarised in Table 3. Footnotes to GRADE tables were used to record reasons for grading a particular quality element as having 'serious' or 'very serious' quality issues. The ratings for each component were combined to obtain an overall assessment of quality for each outcome as described in Table 4.

Table 5: Adaptation of GRADE quality elements for epidemiological reviews

Quality element	Description
Risk of bias ('Study limitations')	Limitations in study design and implementation may bias estimates and interpretation of the effect of the prognostic/risk factor. High risk of bias for the majority of the evidence reduces confidence in the estimated effect. Epidemiological studies are not usually randomised and therefore would not be downgraded for study design from the outset (they start as high quality)
Inconsistency	This refers to unexplained heterogeneity between studies looking at the same outcome, resulting in wide variability in the epidemiological estimates, with little or no overlap in confidence intervals
Indirectness	This refers to any departure from inclusion criteria listed in the review protocol (such as differences in study populations or measurement of outcomes), that may affect the generalisability of results
Imprecision	This occurs when a study has relatively few participants or events of interest, resulting in wide confidence intervals

Assessing risk of bias in prevalence reviews

The Joanna Briggs Institute (JBI) checklist for prevalence studies developed by Munn 2015 was used to assess risk of bias in studies included in epidemiological reviews (see Appendix H in the Developing NICE guidelines: the manual). The risk of bias in each study was determined by assessing the following domains:

- appropriate sample frame
- appropriate participant sampling
- adequate sample size
- detailed description of study subjects and setting
- sufficient coverage of identified sample
- valid methods for identifying the condition
- standardised and reliable measurement of the condition

-
- appropriate statistical analysis
 - adequate response rate.

Assessing inconsistency in prevalence reviews

Where multiple results were deemed appropriate to meta-analyse (that is, there was sufficient similarity between populations, outcome under investigations and unit of analysis) inconsistency was assessed by visually inspecting forest plots and observing whether there was considerable heterogeneity in the results of the meta-analysis. This was assessed by calculating the I-squared statistic for the meta-analysis with an I-squared value of more than 50% indicating serious heterogeneity, and more than 80% indicating very serious heterogeneity. When serious or very serious heterogeneity was observed, possible reasons were explored and subgroup analyses were performed as pre-specified in the review protocol where possible.

When no plausible explanation for the heterogeneity could be found, the quality of the evidence was downgraded in GRADE for inconsistency.

Assessing indirectness in prevalence reviews

Indirectness in prevalence reviews was assessed by comparing the populations, prognostic factors and outcomes in the evidence to those defined in the review protocol. Evidence was considered seriously indirect if 25% to 50% of the evidence for an outcome differed from one of the elements of the protocol and very seriously indirect if more than 50% of the evidence for an outcome differed on one of the elements from the protocol or if 25% to 50% of the evidence for an outcome differed on two or more of the elements from the protocol.

Assessing imprecision and importance in prevalence reviews

Minimally important differences are not appropriate for non-comparative data. Therefore, imprecision was judged based on optimal information size criteria. Evidence was considered seriously imprecise if there were less than 300 events, based on the rule-of-thumb specified in version 3.2 of the GRADE handbook (Schünemann 2009), and very seriously imprecise if there were less than 150 events. The threshold for very serious imprecision was a pragmatic decision, in the absence of a rule-of-thumb being available, based on the fact that this is half the number required for serious imprecision, which would be consistent with approach suggested for continuous outcomes.

The importance of outcomes was assessed qualitatively during committee discussions and documented in the committee's discussion and interpretation of the evidence.

Prognostic studies

Adapted GRADE methodology for prognostic reviews

For prognostic reviews with evidence from comparative studies an adapted GRADE approach was used. As noted above, GRADE methodology is designed for intervention reviews but the quality assessment elements were adapted for prognostic reviews.

The evidence for each outcome in the prognostic reviews was examined separately for the quality elements listed and defined in Table 6. The criteria considered in the rating of these elements are discussed below. Each element was graded using the quality levels summarised in Table 3. Footnotes to GRADE tables were used to record reasons for grading a particular quality element as having ‘serious’ or ‘very serious’ quality issues. The ratings for each component were combined to obtain an overall assessment of quality for each outcome as described in Table 4.

Table 6: Adaptation of GRADE quality elements for prognostic reviews

Quality element	Description
Risk of bias (‘Study limitations’)	Limitations in study design and implementation may bias estimates and interpretation of the effect of the prognostic/risk factor. High risk of bias for the majority of the evidence reduces confidence in the estimated effect.
Inconsistency	This refers to unexplained heterogeneity between studies looking at the same prognostic/risk factor, resulting in wide variability in estimates of association (such as RRs or ORs), with little or no overlap in confidence intervals
Indirectness	This refers to any departure from inclusion criteria listed in the review protocol (such as differences in study populations or prognostic/risk factors), that may affect the generalisability of results
Imprecision	This occurs when a study has relatively few participants or events of interest, resulting in wide confidence intervals

RR, relative risk; OR, odds ratio

Assessing risk of bias in prognostic reviews

The Quality in Prognosis Studies (QUIPS) tool developed by Hayden 2013 was used to assess risk of bias in studies included in prognostic reviews (see Appendix H in the Developing NICE guidelines: the manual). The risk of bias in each study was determined by assessing the following domains:

- selection bias
- attrition bias
- prognostic factor bias
- outcome measurement bias
- control for confounders
- appropriate statistical analysis.

Assessing inconsistency in prognostic reviews

Where multiple results were deemed appropriate to meta-analyse (that is, there was sufficient similarity between risk factor and outcome under investigation) inconsistency was assessed by visually inspecting forest plots and observing whether there was considerable heterogeneity in the results of the meta-analysis. This was assessed by calculating the I-squared statistic for the meta-analysis with an I-squared value of more than 50% indicating serious heterogeneity, and more than 80% indicating very serious heterogeneity. When serious or very serious heterogeneity was observed, possible reasons were explored and subgroup analyses were performed as pre-specified in the review protocol where possible.

When no plausible explanation for the heterogeneity could be found, the quality of the evidence was downgraded in GRADE for inconsistency.

Assessing indirectness in prognostic reviews

Indirectness in prognostic reviews was assessed by comparing the populations, prognostic factors and outcomes in the evidence to those defined in the review protocol. Evidence was considered seriously indirect if 25% to 50% of the evidence for an outcome differed from one of the elements of the protocol and very seriously indirect if more than 50% of the evidence for an outcome differed on one of the elements from the protocol or if 25% to 50% of the evidence for an outcome differed on two or more of the elements from the protocol.

Assessing imprecision and importance in prognostic reviews

Prognostic studies may have a variety of purposes, for example, establishing typical prognosis in a broad population, establishing the effect of patient characteristics on prognosis, and developing a prognostic model. While by convention MIDs relate to intervention effects, the committee agreed to use GRADE default MIDs for risk ratios as a starting point from which to assess whether the size of an outcome effect (association) in a prognostic study would be large enough to be meaningful in practice. Specifically, the committee agreed that these values would correspond to a moderate association between the prognostic factor and the outcome, with any statistically significant association being considered a small association, and risk ratios <0.5 and >2.00 being considered a strong association between the factor and the outcome. The latter threshold was selected for consistency with estimated effect sizes where it is possible to consider upgrading non-RCT evidence in GRADE. Imprecision was judged based on optimal information size criteria. Evidence was considered seriously imprecise if there were less than 300 events, based on the rule-of-thumb specified in version 3.2 of the GRADE handbook (Schünemann 2009), and very seriously imprecise if there were less than 150 events. The threshold for very serious imprecision was a pragmatic decision, in the absence of a rule-of-thumb being available, based on the fact that this is half the number required for serious imprecision, which would be consistent with the approach suggested for continuous outcomes.

MIDs, the line of no effect, and both 95% and 90% confidence intervals (CIs) were used to assess whether there were important differences in outcomes between groups, as described above.

Diagnostic studies

Adapted GRADE methodology for diagnostic reviews

For diagnostic reviews, an adapted GRADE approach was used. GRADE methodology is designed for intervention reviews but the quality assessment elements and outcome presentation were adapted by the NICE technical team for diagnostic test accuracy reviews and prediction models. For example, GRADE tables were modified to include diagnostic test accuracy measures (sensitivity, specificity and positive and negative predictive values).

The evidence for each outcome in the diagnostic reviews and prediction models was examined separately for the quality elements listed and defined in Table 7. The

criteria considered in the rating of these elements are discussed below. Each element was graded using the quality levels summarised in Table 3. Footnotes to GRADE tables were used to record reasons for grading a particular quality element as having a ‘serious’ or ‘very serious’ quality issue. The ratings for each component were combined to obtain an overall assessment of quality for each outcome as described in Table 4.

The initial quality rating was based on the study design: cross-sectional or cohort studies start as ‘high’ quality and case–control studies start as ‘low’ quality.

Table 7: Adaptation of GRADE quality elements for diagnostic reviews

Quality element	Description
Risk of bias (‘Study limitations’)	Limitations in study design and implementation may bias estimates of diagnostic accuracy. High risk of bias for the majority of the evidence reduces confidence in the estimated effect. Diagnostic accuracy studies are not usually randomised and therefore would not be downgraded for study design from the outset (they start as high quality)
Inconsistency	This refers to unexplained heterogeneity in test accuracy measures (such as sensitivity and specificity) between studies
Indirectness	This refers to differences in study populations, index tests, reference standards or outcomes between the available evidence and inclusion criteria specified in the review protocol
Imprecision	This occurs when a study has relatively few participants and the probability of a correct diagnosis is low. Accuracy measures would therefore have wide confidence intervals around the estimated effect

Assessing risk of bias in diagnostic reviews and prediction models

Risk of bias in diagnostic reviews and prediction models was assessed using the Quality Assessment of Diagnostic Accuracy Studies version 2 (QUADAS-2) checklist (see Appendix H in Developing NICE guidelines: the manual).

Risk of bias in primary diagnostic accuracy reviews or prediction models in QUADAS-2 consists of 4 domains:

- participant selection
- index test
- reference standard
- flow and timing.

More details about the QUADAS-2 tool can be found on the developer’s website.

Assessing inconsistency in diagnostic reviews and prediction models

Inconsistency refers to the unexplained heterogeneity of the results in meta-analysis. When estimates of diagnostic accuracy and prediction model parameters vary widely across studies (that is, there is heterogeneity or variability in results), this suggests true differences in underlying effects. Inconsistency is, thus, only truly applicable when statistical meta-analysis is conducted (that is, results from different studies are pooled).

Inconsistency for diagnostic reviews was assessed based on visual inspection of the point estimates and confidence intervals of the included studies. If these varied widely (for example, point estimates for some studies lying outside the CIs of other studies) the evidence was downgraded for inconsistency.

Assessing indirectness in diagnostic reviews and prediction models

Indirectness in diagnostic reviews and prediction models was assessed using the QUADAS-2 checklist by assessing the applicability of the studies in relation to the review question in the following domains:

- participant selection
- index test
- reference standard.

Evidence was considered seriously indirect if 25% to 50% of the evidence for an outcome differed from one of the domains and very seriously indirect if more than 50% of the evidence for an outcome differed on one of the domains or if 25% to 50% of the evidence for an outcome differed on two or more of the domains.

More details about the QUADAS-2 tool can be found on the developer's website.

Assessing imprecision and importance in diagnostic reviews and prediction models

The judgement of precision for diagnostic and prediction model evidence was based on the CIs of the sensitivity and specificity. The committee defined 2 decision thresholds for each measure, a value above which the test could be recommended and a value below which the test would be considered of no use. These thresholds were based on the committee's experience and consensus.

The thresholds were:

- sensitivity: low threshold 60%, high threshold 90%
- specificity: low threshold 60%, low threshold 90%.

Outcomes were downgraded for imprecision when their 95% CI crossed at least 1 threshold. If the CI crossed 1 threshold, the outcome was downgraded once for imprecision. If the CI crossed 2 thresholds, the outcome was downgraded twice for imprecision. These assessments were made on the meta-analysed outcomes where applicable or if outcomes were not meta-analysed, on the individual study results themselves.

Decision making thresholds for positive and negative predictive values have not been defined a priori. Imprecision in and importance of positive and negative predictive values will be assessed qualitatively during committee discussions and documented in the committee's discussion and interpretation of the evidence.

Qualitative studies

GRADE-CERQual methodology for qualitative reviews

For qualitative reviews an adapted GRADE Confidence in the Evidence from Reviews of Qualitative research (GRADE-CERQual) approach (Lewin 2018) was used. In this approach the confidence in the evidence is considered according to themes in the evidence. The themes may have been identified in the primary studies

or they may have been identified by considering the reports of a number of studies. Confidence elements assessed using GRADE-CERQual are listed and defined in Table 8. Each element was graded using the levels of concern summarised in Table 9.

The ratings for each component were combined (as with other types of evidence) to obtain an overall assessment of confidence for each theme as described in Table 10. 'Confidence' in this context refers to the extent to which the review finding is a reasonable representation of the phenomenon of interest set out in the protocol. Similar to other types of evidence all review findings start off with 'high confidence' and are rated down by one or more levels if there are concerns about any of the individual CERQual components. In line with advice from the CERQual developers, the overall assessment does not involve numerical scoring for each component but in order to ensure consistency across and between guidelines, the NGA established some guiding principles for overall ratings. For example, a review finding would not be downgraded (and therefore would be assessed with 'high' confidence) if at least 2 of the individual components were rated as 'no or very minor'; and none of the components were rated as having moderate or serious concerns.

At the other extreme, a review finding would be downgraded 3 times (to 'very low') if at least 2 components had serious concerns or 3 had moderate concerns (as long as the 4th component was rated 'serious') or if all components had moderate concerns. A basic principle was that if any components had any serious concerns then overall confidence in the review finding would be downgraded at least twice, to low. Transparency about overall judgements is provided in the CERQual tables, with explanations for downgrading given in the individual domain cells.

Table 8: Adaptation of GRADE quality elements for qualitative reviews

Quality element	Description
Methodological limitations	Limitations in study design and implementation may bias interpretation of qualitative themes identified. High risk of bias for the majority of the evidence reduces our confidence that the review findings reflect the phenomena of interest. Qualitative studies are not usually randomised and therefore would not be downgraded for study design from the outset (they start as high confidence)
Relevance (or applicability) of evidence	This refers to the extent to which the context of the studies supporting the review findings is applicable to the context specified in the review question
Coherence of findings	This refers to the extent to which review findings are well grounded in data from the contributing primary studies and provide a credible explanation for patterns identified in the evidence. If the data from the underlying studies are ambiguous or contradict the review finding this would reduce our confidence in the finding.
Adequacy of data (theme saturation or sufficiency)	This corresponds to a similar concept in primary qualitative research, that is, whether a theoretical point of theme saturation was achieved, at which point no further citations or observations would provide more insight or suggest a different interpretation of the particular theme. Judgements are not based on the number of studies but do take account of the quantity and also richness of data underpinning a finding. The more complex the finding, the more detailed the supporting data need to be. For simple findings, relatively superficial data would be considered adequate to explain and explore the phenomenon being described.

Table 9: CERQual levels of concern (by quality element)

Level of concern	Definition
None or very minor concerns	Unlikely to reduce confidence in the review finding
Minor concerns	May reduce confidence in the review finding
Moderate concerns	Will probably reduce confidence in the review finding
Serious concerns	Very likely to reduce confidence in the review finding

Table 10: Overall confidence in the evidence in CERQual (by review finding)

Overall confidence level	Definition
High	It is highly likely that the review finding is a reasonable representation of the phenomenon of interest
Moderate	It is likely that the review finding is a reasonable representation of the phenomenon of interest
Low	It is possible that the review finding is a reasonable representation of the phenomenon of interest
Very low	It is unclear whether the review finding is a reasonable representation of the phenomenon of interest

Assessing methodological limitations in qualitative reviews

Methodological limitations in qualitative studies were assessed using the Critical Appraisal Skills Programme (CASP) checklist for qualitative studies (see appendix H in Developing NICE guidelines: the manual). Overall methodological limitations were derived by assessing the methodological limitations across the 6 domains summarised in Table 11.

Table 11: Methodological limitations in qualitative studies

Aim and appropriateness of qualitative evidence	This domain assesses whether the aims and relevance of the study were described clearly and whether qualitative research methods were appropriate for investigating the research question
Rigour in study design or validity of theoretical approach	This domain assesses whether the study approach was documented clearly and whether it was based on a theoretical framework (such as ethnography or grounded theory). This does not necessarily mean that the framework has to be stated explicitly, but a detailed description ensuring transparency and reproducibility should be provided
Sample selection	This domain assesses the background, the procedure and reasons for the method of

	selecting participants. The assessment should include consideration of any relationship between the researcher and the participants, and how this might have influenced the findings
Data collection	This domain assesses the documentation of the method of data collection (in-depth interviews, semi-structured interviews, focus groups or observations). It also assesses who conducted any interviews, how long they lasted and where they took place
Data analysis	This domain assesses whether sufficient detail was documented for the analytical process and whether it was in accordance with the theoretical approach. For example, if a thematic analysis was used, the assessment would focus on the description of the approach used to generate themes. Consideration of data saturation would also form part of this assessment (it could be reported directly or it might be inferred from the citations documented that more themes could be found)
Results	This domain assesses any reasoning accompanying reporting of results (for example, whether a theoretical proposal or framework is provided)

Assessing relevance of evidence in qualitative reviews

Relevance (applicability) of findings in qualitative research is the equivalent of indirectness for quantitative outcomes, and refers to how closely the aims and context of studies contributing to a theme reflect the objectives outlined in the guideline review protocol.

Assessing coherence of findings in qualitative reviews

For qualitative research, a similar concept to inconsistency is coherence, which refers to the way findings within themes are described and whether they make sense. This concept was used in the quality assessment across studies for individual themes. This does not mean that contradictory evidence was automatically downgraded, but that it was highlighted and presented, and that reasoning was provided. Provided the themes, or components of themes, from individual studies fit into a theoretical framework, they do not necessarily have to reflect the same perspective. It should, however, be possible to explain these by differences in context (for example, the views of health or social care professionals might not be the same as those of family members, but they could contribute to the same overarching themes).

Assessing adequacy of data in qualitative reviews

Adequacy of data (theme saturation or sufficiency) corresponds to a similar concept in primary qualitative research in which consideration is made of whether a

theoretical point of theme saturation was achieved, meaning that no further citations or observations would provide more insight or suggest a different interpretation of the theme concerned. As noted above, it is not equivalent to the number of studies contributing to a theme, but it does take account of the quantity of data supporting a review finding (for instance whether sufficient quotations or observations were provided to underpin the findings) and in particular the degree of 'richness' of supporting data. Concerns about richness arise when insufficient details are provided by the data to enable an understanding of the phenomenon being described. Generally, if a review finding is fairly simple then relatively superficial data will be needed to understand it. Data underpinning a more complex finding would need to offer greater detail, allowing for interpretation and exploration of the phenomenon being described. Therefore in assessing adequacy downgrading involved weighing up the complexity of the review finding against the explanatory contribution of the supporting data.

Reviewing economic evidence

Titles and abstracts of articles identified through the economic literature searches were independently assessed for inclusion using the predefined eligibility criteria listed in Table 12.

Table 12: Inclusion and exclusion criteria for systematic reviews of economic evaluations

Inclusion criteria
Intervention or comparators in accordance with the guideline scope
Study population in accordance with the guideline scope
Full economic evaluations (cost-utility, cost effectiveness, cost-benefit or cost-consequence analyses) assessing both costs and outcomes associated with interventions of interest
Exclusion criteria
Abstracts containing insufficient methodological details
Cost-of-illness type studies

Once the screening of titles and abstracts was completed, full-text copies of potentially relevant articles were requested for detailed assessment. Inclusion and exclusion criteria were applied to articles obtained as full-text copies.

Details of economic evidence study selection, lists of excluded studies and, economic evidence tables are presented in appendices G, H and J of the evidence report. The results of quality assessment of economic evidence (see below) and health economic profiles are provided in the main body of the evidence review.

Appraising the quality of economic evidence

The quality of economic evidence was assessed using the economic evaluations checklist specified in [Developing NICE guidelines: the manual](#).

Economic modelling

The aims of the economic input to the guideline were to inform the guideline committee of potential economic issues to ensure that recommendations represented

a cost effective use of healthcare resources. Economic evaluations aim to integrate data on healthcare benefits (ideally in terms of quality-adjusted life-years; QALYs) with the costs of different options. In addition, the economic input aimed to identify areas of high resource impact; these are recommendations which (while cost effective) might have a large impact on Clinical Commissioning Group or Trust finances and so need special attention.

The guideline committee prioritised the following review questions for economic modelling where it was thought that economic considerations would be particularly important in formulating recommendations.

- What is the effectiveness of ventilation tubes for OME in children under 12 years?
- What is the effectiveness of adenoidectomy (with or without ventilation tubes) for OME in children under 12 years?
- What is the effectiveness of air and bone conduction hearing aids for children with OME under 12 years?

The methods and results of the de novo economic analyses is reported in Appendix I of the Evidence review E. When new economic analysis was not prioritised, the committee made a qualitative judgement regarding cost effectiveness by considering expected differences in resource and cost use between options, alongside clinical effectiveness evidence identified from the clinical evidence review.

Cost effectiveness criteria

NICE sets out the [principles](#) that committees should consider when judging whether an intervention offers good value for money. In general, an intervention was considered to be cost effective if any of the following criteria applied (provided that the estimate was considered plausible):

- the intervention dominated other relevant strategies (that is, it was both less costly in terms of resource use and more effective compared with all the other relevant alternative strategies)
- the intervention cost less than £20,000 per QALY gained compared with the next best strategy
- the intervention provided important benefits at an acceptable additional cost when compared with the next best strategy.

The committee's considerations of cost effectiveness are discussed explicitly under the heading 'cost-effectiveness and resource use' in the relevant evidence reviews.

Developing recommendations

Guideline recommendations

Recommendations were drafted on the basis of the committee's interpretation of the available evidence, taking account of the balance of benefits, harms and costs between different courses of action. When effectiveness, qualitative and economic evidence was of poor quality, conflicting or absent, the committee drafted recommendations based on their expert opinion. The considerations for making consensus-based recommendations include the balance between potential benefits and harms, the economic costs or implications compared with the economic benefits,

current practices, recommendations made in other relevant guidelines, person's preferences and equality issues.

The main considerations specific to each recommendation are outlined under the heading 'The committee's discussion of the evidence' within each evidence review.

For further details refer to Developing NICE guidelines: the manual.

Research recommendations

When areas were identified for which evidence was lacking, the committee considered making recommendations for future research. For further details refer to Developing NICE guidelines: the manual and NICE's Research recommendations process and methods guide.

Validation process

This guideline was subject to a 6-week public consultation and feedback process. All comments received from registered stakeholders were responded to in writing and posted on the NICE website at publication. For further details refer to Developing NICE guidelines: the manual.

Updating the guideline

Following publication, NICE will undertake a surveillance review to determine whether the evidence base has progressed sufficiently to consider altering the guideline recommendations and warrant an update. For further details refer to Developing NICE guidelines: the manual.

References

Bradburn 2007

Bradburn, M. J., Deeks, J. J., Berlin, J. A., & Localio, A. R. Much ado about nothing: A comparison of the performance of meta-analytical methods with rare events. *Statistics in Medicine*, 26, 53–77, 2007.

Dixon-Woods 2005

Dixon-Woods M, Agarwal S, Jones D et al. (2005) Synthesising qualitative and quantitative evidence: a review of possible methods. *Journal of Health Services Research & Policy* 10(1), 45–53

Hayden 2013

Jill A. Hayden, Danielle A. van der Windt, Jennifer L. Cartwright, Pierre Côté, Claire Bombardier. Assessing Bias in Studies of Prognostic Factors. *Ann Intern Med*. 2013;158:280–286. doi: 10.7326/0003-4819-158-4-201302190-00009

Higgins 2011

Higgins JPT, Green S (editors) (2011) *Cochrane Handbook for Systematic Reviews of Interventions* Version 5.1.0 [updated 2019] The Cochrane Collaboration. Available from www.handbook.cochrane.org (accessed 27 March 2023)

Lewin 2018

Lewin S, Booth A, Glenton C, Munthe-Kaas H et al. (2018) Applying GRADE-CERQual to qualitative evidence synthesis findings: introduction to the series. *Implement Sci*. 2018 Jan 25;13 (Suppl1): 2

Liu 2020

Liu PZ, Ismail-Koch H, Stephenson K, Donne AJ et al. (2020) A core outcome set for research on the management of otitis media with effusion in otherwise-healthy children. *Int J Pediatr Otorhinolaryngol*. 2020 Jul 134;110029: doi: 10.1016/j.ijporl.2020.110029

MacKeith 2023a

MacKeith S, Mulvaney CA, Galbraith K, Webster KE, Connolly R, Paing A, Marom T, Daniel M, Venekamp RP, Rovers MM, Schilder AGM. Ventilation tubes (grommets) for otitis media with effusion (OME) in children. *Cochrane Database of Systematic Reviews* 2023. Art. No.: CD015215. DOI: 10.1002/14651858.CD015215.pub2

MacKeith 2023b

MacKeith S, Mulvaney CA, Galbraith K, Webster KE, Paing A, Connolly R, Marom T, Daniel M, Venekamp RP, Schilder AGM. Adenoidectomy for otitis media with effusion (OME) in children. *Cochrane Database of Systematic Reviews* 2023. Art. No.: CD015252. DOI: 10.1002/14651858.CD015252.pub2

Mulvaney 2023a

Mulvaney CA, Galbraith K, Webster KE, Rana M, Connolly R, Marom T, Daniel M, Venekamp RP, Schilder AGM, MacKeith S. Antibiotics for otitis media with effusion (OME) in children. Cochrane Database of Systematic Reviews 2023. Art. No.: CD015254. DOI: 10.1002/14651858.CD015254.pub2

Mulvaney 2023b

Mulvaney CA, Galbraith K, Webster KE, Rana M, Connolly R, Tudor-Green B, Marom T, Daniel M, Venekamp RP, Schilder AGM, MacKeith S. Topical and oral steroids for otitis media with effusion (OME) in children. Cochrane Database of Systematic Reviews 2023. Art. No.: CD015255. DOI: 10.1002/14651858.CD015255.pub2

Webster 2023

Webster KE, Mulvaney CA, Galbraith K, Rana M, Marom T, Daniel M, Venekamp RP, Schilder AGM, MacKeith S. Autoinflation for otitis media with effusion (OME) in children. Cochrane Database of Systematic Reviews 2023. Art. No.: CD015253. DOI: 10.1002/14651858.CD015253.pub2

McGowan 2016

McGowan J, Sampson M, Salzwedel DM et al. (2016) [PRESS Peer Review of Electronic Search Strategies: 2015 guideline statement](#). Journal of Clinical Epidemiology 75: 40–6

Munn 2015

Munn Z, Moola S, Lisy K, Riitano D, Tufanaru C, (2015) Methodological guidance for systematic reviews of observational epidemiological studies reporting prevalence and incidence data. Internal Journal of Evidence Based Healthcare, 13(3), 147-153

NICE 2018

National Institute for Health and Care Excellence (NICE) (2014) NICE Policy on conflicts of interest (updated 2017). Available from <https://www.nice.org.uk/Media/Default/About/Who-we-are/Policies-and-procedures/declaration-of-interests-policy.pdf> (accessed 27 March 2023)

Pandey 2020

Pandey S, (2022) metaSurvival: Meta-analysis of a single survival curve using the multivariate methodology of DerSimonian and Laird. R package version 0.1.0. Available from <https://cran.r-project.org/web/packages/metaSurvival/index.html> (accessed 27 March 2023)

Santesso 2016

Santesso N, Carrasco-Labra A, Langendam M et al. (2016) Improving GRADE evidence tables part 3: detailed guidance for explanatory footnotes supports creating and understanding GRADE certainty in the evidence judgments. Journal of clinical epidemiology 74, 28-39

Schünemann 2009

Schünemann H, Brožek J, Oxman A, (editors) (2009) GRADE handbook for grading quality of evidence and strength of recommendation. Version 3.2 [updated March 2009]

Viechtbauer 2010

Viechtbauer W, (2010) Conducting meta-analyses in R with the metafor package. *Journal of Statistical Software*, 36(3), 1-48. Available from <https://doi.org/10.18637/jss.v036.i0> (accessed 20 December 2022)