

Draft

Neonatal infection: antibiotics for prevention and treatment

NICE guideline: methods

NICE guideline <number>

Methods

December 2020

Draft for Consultation

*Evidence reviews were developed by the
NICE guideline updates team*

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 All NICE guidance is subject to regular review and may be updated or withdrawn.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, 2019. All rights reserved. Subject to [Notice of rights](#).

ISBN:

Contents

Development of the guideline	6
What this guideline covers.....	6
What this guideline does not cover.....	6
Methods	7
Developing the review questions and outcomes.....	7
Reviewing research evidence.....	7
Review protocols.....	7
Searching for evidence.....	7
Selecting studies for inclusion.....	7
Incorporating published systematic reviews.....	8
Quality assessment.....	9
Evidence synthesis and meta-analyses of pair-wise data.....	10
Evidence of effectiveness of interventions.....	10
Quality assessment.....	10
Methods for combining intervention evidence.....	11
Minimal clinically important differences (MIDs) for intervention reviews.....	12
Modified GRADE for pairwise meta-analyses of interventional evidence.....	12
Publication bias.....	13
Methods for combining direct and indirect evidence (network meta-analysis) for interventions.....	14
Modified GRADE for network meta-analyses.....	14
Diagnostic test accuracy evidence.....	15
Quality assessment.....	16
Methods for combining diagnostic test accuracy evidence.....	17
GRADE for diagnostic test accuracy evidence.....	17
Publication bias.....	19
Predictive (prognostic) accuracy evidence.....	19
Quality assessment.....	20
Methods for combining predictive accuracy evidence.....	20
Modified GRADE for predictive accuracy evidence.....	21
Publication bias.....	22
Assessing c-statistics.....	22
Association studies.....	24
Quality assessment.....	24
Methods for combining association studies.....	25
Modified GRADE for association studies.....	25
Qualitative evidence.....	26
Quality assessment.....	26

Methods for combining qualitative evidence	26
CERQual for qualitative studies	26
Health economics	27

1 **Development of the guideline**

2 **What this guideline covers**

3 This method document describes the evidence reviews that were part of the 2021
4 update of this guideline. These evidence reviews included:

5 Evidence review A: Info and support

6 Evidence review B: Intrapartum antibiotics

7 Evidence review C: PPRM

8 Evidence review D: Risk factors for early onset

9 Evidence review E: Risk factors for late onset

10 Evidence review F: Intravascular catheters

11 Evidence review G: Investigations

12 Evidence review H: Antibiotics

13 Evidence review I: Antifungals

14 **What this guideline does not cover**

15 This method does not cover the parts of the 2012 guideline that were not updated in
16 2021. The methods used to develop these sections of the guideline are described in
17 the full version of the guideline.

1 **Methods**

2 This guideline was developed using the methods described in the NICE guidelines
3 manual.

4 Declarations of interest were recorded according to the NICE conflicts of interest
5 policy.

6 **Developing the review questions and outcomes**

7 The 13 review questions developed for this guideline were based on the key areas
8 identified in the guideline [scope](#). They were drafted by the NICE guideline updates
9 team and refined and validated by the guideline committee.

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

12 **Reviewing research evidence**

13 **Review protocols**

14 Review protocols were developed with the guideline committee to outline the
15 inclusion and exclusion criteria used to select studies for each evidence review.
16 Where possible, review protocols were prospectively registered in the [PROSPERO](#)
17 [register of systematic reviews](#).

18 **Searching for evidence**

19 Evidence was searched for each review question using the methods specified in the
20 NICE guidelines manual.

21 **Selecting studies for inclusion**

22 All references identified by the literature searches and from other sources (for
23 example, a previous version of the guideline or studies identified by committee
24 members) were uploaded into EPPI reviewer software and de-duplicated. Titles and
25 abstracts were assessed for possible inclusion using the criteria specified in the
26 review protocol. 10% of the abstracts were reviewed by two reviewers, with any
27 disagreements resolved by discussion or, if necessary, a third independent reviewer.

28 The reviews undertaken for this guideline all made use of the priority screening
29 functionality with the EPPI-reviewer systematic reviewing software. This uses a
30 machine learning algorithm (specifically, an SGD classifier) to take information on
31 features (1, 2 and 3 word blocks) in the titles and abstract of papers marked as being
32 'includes' or 'excludes' during the title and abstract screening process, and re-orders
33 the remaining records from most likely to least likely to be an include, based on that
34 algorithm. This re-ordering of the remaining records occurs every time 25 additional
35 records have been screened.

36 For most reviews in this guideline, priority screening was used to highlight the most
37 relevant records earlier in the search but was not used as a method to stop abstract

1 screening early. Consequently, the whole abstract database was searched for most
2 review questions.

3 For two reviews (Review G – Investigations and Review H – antibiotics), priority
4 screening was used to prioritise the articles that were most relevant to the review and
5 to allow screening to be stopped early.

6 Research is currently ongoing as to what are the appropriate thresholds where
7 reviewing of abstract can be stopped, assuming a defined threshold for the
8 proportion of relevant papers it is acceptable to miss on primary screening. As a
9 conservative approach until that research has been completed, the following rules
10 were adopted during the production of this guideline:

- 11 • At least 50% of the identified abstract (or 1,000 records, if that is a greater
12 number) were always screened.
- 13 • After this point, screening was only terminated if at least 500 additional
14 abstracts were screened without a single new include being identified.
- 15 • A random 10% sample of the studies remaining in the database when the
16 threshold was reached were additionally screened, to check if a substantial
17 number of relevant studies were not being correctly classified by the
18 algorithm, with the full database being screened if concerns were identified.

19 For review G (Investigations), 2796 abstracts (64% of the database) were screened
20 and for review H (antibiotics), 2949 abstracts (60% of the database) were screened
21 before the stopping criteria was met. As an additional check to ensure this approach
22 did not miss relevant studies, the included studies lists of systematic reviews were
23 searched to identify any papers not identified through the primary search.

24 **Incorporating published systematic reviews**

25 For all review questions where a literature search was undertaken looking for a
26 particular study design, systematic reviews containing studies of that design were
27 also included. All included studies from those systematic reviews were screened to
28 identify any additional relevant primary studies not found as part of the initial search.

29
30 If published evidence syntheses were identified sufficiently early in the review
31 process (for example, from the surveillance review or early in the database search),
32 they were considered for use as the primary source of data, rather than extracting
33 information from primary studies. Syntheses considered for inclusion in this way were
34 quality assessed to assess their suitability using ROBIS checklist. Note that this
35 quality assessment was solely used to assess the quality of the synthesis in order to
36 decide whether it could be used as a source of data, as outlined in Table 1: Criteria
37 for using published evidence syntheses as a source of data

38

1 Quality assessment

2 Individual systematic reviews that were considered as a direct source of data were
3 quality assessed using the ROBIS tool, with each classified into one of the following
4 three groups:

- 5 • High quality – It is unlikely that additional relevant and important data would be
6 identified from primary studies compared to that reported in the review, and
7 unlikely that any relevant and important studies have been missed by the review.
- 8 • Moderate quality – It is possible that additional relevant and important data would
9 be identified from primary studies compared to that reported in the review, but
10 unlikely that any relevant and important studies have been missed by the review.
- 11 • Low quality – It is possible that relevant and important studies have been missed
12 by the review.

13 Each individual systematic review was also classified into one of three groups for its
14 applicability as a source of data, based on how closely the review matches the
15 specified review protocol in the guideline. Studies were rated as follows:

- 16 • Fully applicable – The identified review fully covers the review protocol in the
17 guideline.
- 18 • Partially applicable – The identified review fully covers a discrete subsection of the
19 review protocol in the guideline (for example, some of the factors in the protocol
20 only).
- 21 • Not applicable – The identified review, despite including studies relevant to the
22 review question, does not fully cover any discrete subsection of the review
23 protocol in the guideline.

24
25

Table 1: Criteria for using published evidence syntheses as a source of data

Quality	Applicability	Use of published evidence synthesis
High	Fully applicable	Data from the published evidence synthesis were used instead of undertaking a new literature search or data analysis. Searches were only done to cover the period of time since the search date of the review. If the review was considered up to date (following discussion with the guideline committee and NICE lead for quality assurance), no additional search was conducted.
High	Partially applicable	Data from the published evidence synthesis were used instead of undertaking a new literature search and data analysis for the relevant subsection of the protocol. For this section, searches were only done to cover the period of time since the search date of the review. If the review was considered up to date (following discussion with the guideline committee and NICE lead for quality assurance), no additional search was conducted. For other sections not covered by the evidence synthesis, searches were undertaken as normal.
Moderate	Fully applicable	Details of included studies were used instead of undertaking a new literature search. Full-text papers of included studies were still retrieved for the purposes of data analysis. Searches were

Quality	Applicability	Use of published evidence synthesis
		only done to cover the period of time since the search date of the review.
Moderate	Partially applicable	Details of included studies were used instead of undertaking a new literature search for the relevant subsection of the protocol. For this section, searches were only done to cover the period of time since the search date of the review. For other sections not covered by the evidence synthesis, searches were undertaken as normal.
Low	Any	The published evidence synthesis was not used as a source of data and a full literature review was completed.
Any	Not applicable	

1 For most reviews, no additional studies were identified from systematic reviews. For
2 one review (review D - early onset risk factors), 2 additional studies which examined
3 the effectiveness of prognostic models for early-onset infection were found from
4 systematic reviews. The primary studies were reviewed and included in this review.
5 For review I (antifungals), most results were taken directly from 2, high quality, fully
6 applicable, systematic reviews which examined the effectiveness of antifungals as
7 prophylaxis in very low birthweight and preterm babies.

8 Evidence synthesis and meta-analyses of pair-wise data

9 Where possible, meta-analyses were conducted to combine the results of
10 quantitative studies for each outcome. For continuous outcomes analysed as mean
11 differences, where change from baseline data were reported in the trials and were
12 accompanied by a measure of spread (for example standard deviation), these were
13 extracted and used in the meta-analysis. Where measures of spread for change from
14 baseline values were not reported, the corresponding values at study end were used
15 and were combined with change from baseline values to produce summary estimates
16 of effect. These studies were assessed to ensure that baseline values were balanced
17 across the treatment groups; if there were significant differences at baseline these
18 studies were not included in any meta-analysis and were reported separately. Where
19 there were differences in populations or interventions, meta-analyses were not
20 performed, and instead the results of individual studies were presented.

21 Evidence of effectiveness of interventions

22 Quality assessment

23 Individual RCTs and quasi-randomised controlled trials were quality assessed using
24 the Cochrane Risk of Bias 2.0 Tool. Other studies were quality assessed using the
25 ROBINS-I tool. Each individual study was classified into one of the following three
26 groups:

- 27 • Low risk of bias – The true effect size for the study is likely to be close to the
28 estimated effect size.
- 29 • Moderate risk of bias – There is a possibility the true effect size for the study is
30 substantially different to the estimated effect size.
- 31 • High risk of bias – It is likely the true effect size for the study is substantially
32 different to the estimated effect size.

1 Each individual study was also classified into one of three groups for directness,
2 based on if there were concerns about the population, intervention, comparator
3 and/or outcomes in the study and how directly these variables could address the
4 specified review question. Studies were rated as follows:

- 5 • Direct – No important deviations from the protocol in population, intervention,
6 comparator and/or outcomes.
- 7 • Partially indirect – Important deviations from the protocol in one of the population,
8 intervention, comparator and/or outcomes.
- 9 • Indirect – Important deviations from the protocol in at least two of the following
10 areas: population, intervention, comparator and/or outcomes.

11 **Methods for combining intervention evidence**

12 Meta-analyses of interventional data were conducted with reference to the Cochrane
13 Handbook for Systematic Reviews of Interventions (Higgins et al. 2011).

14 Where outcomes measured the same underlying construct but used different
15 instruments/metrics, data were analysed using standardised mean differences
16 (Hedges' g).

17 A pooled relative risk was calculated for dichotomous outcomes (using the Mantel–
18 Haenszel method) reporting numbers of people having an event, and a pooled
19 incidence rate ratio was calculated for dichotomous outcomes reporting total
20 numbers of events. Both relative and absolute risks were presented, with absolute
21 risks calculated by applying the relative risk to the risk in the comparator arm of the
22 meta-analysis (calculated as the total number events in the comparator arms of
23 studies in the meta-analysis divided by the total number of participants in the
24 comparator arms of studies in the meta-analysis).

25 Fixed- and random-effects models (der Simonian and Laird) were fitted for all
26 syntheses, with the presented analysis dependent on the degree of heterogeneity in
27 the assembled evidence. Fixed-effects models were the preferred choice to report,
28 but in situations where the assumption of a shared mean for fixed-effects model were
29 clearly not met, even after appropriate pre-specified subgroup analyses were
30 conducted, random-effects results are presented. Fixed-effects models were deemed
31 to be inappropriate if one or both of the following conditions was met:

- 32 • Significant between study heterogeneity in methodology, population,
33 intervention or comparator was identified by the reviewer in advance of data analysis.
34 This decision was made and recorded before any data analysis was undertaken.
- 35 • The presence of significant statistical heterogeneity in the meta-analysis,
36 defined as $I^2 \geq 50\%$.

37 In situations where subgroup analyses were conducted, pooled results and results for
38 the individual subgroups are reported when there was evidence of between group
39 heterogeneity, defined as a statistically significant test for subgroup interactions (at
40 the 95% confidence level). Where no such evidence as identified, only pooled results
41 are presented.

42 Meta-analyses were performed in Cochrane Review Manager V5.3.

1 Minimal clinically important differences (MIDs) for intervention reviews

2 The Core Outcome Measures in Effectiveness Trials (COMET) database was
3 searched to identify published minimal clinically important difference thresholds
4 relevant to this guideline. Identified MIDs were assessed to ensure they had been
5 developed and validated in a methodologically rigorous way, and were applicable to
6 the populations, interventions and outcomes specified in this guideline. In addition,
7 the Guideline Committee were asked to prospectively specify any outcomes where
8 they felt a consensus MID could be defined from their experience. MIDs identified
9 through this process were intended to be used to inform discussions on the clinical
10 importance of effects and the precision of effect estimates. No published MIDs were
11 found through this process and the committee did not wish to pre specify consensus
12 MIDs for any outcome. The clinical importance of effects was judged by the
13 committee taking into account evidence across all outcomes and absolute effect
14 estimates. These discussions are documented in the committee discussion section of
15 each evidence review.

16

17 Modified GRADE for pairwise meta-analyses of interventional evidence

18 GRADE was used to assess the quality of evidence for the selected outcomes as
19 specified in 'Developing NICE guidelines: the manual'. Data from randomised
20 controlled trials, non-randomised controlled trials and cohort studies (which were
21 quality assessed using the Cochrane risk of bias tool or ROBINS-I) were initially
22 rated as high quality while data from other study types were initially rated as low
23 quality. Ratings were subsequently downgraded once for each serious source of
24 uncertainty and twice of each very serious source of uncertainty, as outlined in the
25 table below.

26 This evidence review for this guideline was conducted using a modified version of the
27 GRADE approach to rate the certainty of evidence in systematic reviews. This is part
28 of a pilot project being undertaken by NICE, to examine the assessment of certainty
29 of evidence in systematic reviews. Instead of using predefined MIDs to assess
30 imprecision in GRADE tables, imprecision was assessed qualitatively during
31 committee discussions. These discussions involved consideration of published MIDs
32 where they exist, but the committee were also encouraged to make judgements of
33 imprecision based on the 95% confidence intervals and sample sizes reported in the
34 GRADE tables. The committee were not aware of any published MIDs for any of the
35 outcomes in the intervention reviews and so the discussions were based on the width
36 of confidence intervals and whether they crossed the line of no effect. This should
37 enable judgements of clinical importance to be made in the context of wider decision
38 making, taking into account evidence across all outcomes and analyses, including
39 health economic analyses.

40 Committee discussions regarding the clinical importance of effects was recorded in
41 the 'imprecision and clinical importance of effects' section of the evidence review. In
42 particular, this included consideration of whether the whole effect of a treatment
43 (which may be felt across multiple independent outcome domains) would be likely to
44 be clinically meaningful, rather than simply whether each individual sub outcome
45 might be meaningful in isolation. The impact of imprecision on the recommendations
46 was presented in the 'quality of the evidence' section of the committee discussion in
47 the evidence review.

1
2**Table 2: Rationale for downgrading quality of evidence for intervention studies**

GRADE criteria	Reasons for downgrading quality
Risk of bias	<p>Not serious: If less than 33.3% of the weight in a meta-analysis came from studies at moderate or high risk of bias, the overall outcome was not downgraded.</p> <p>Serious: If greater than 33.3% of the weight in a meta-analysis came from studies at moderate or high risk of bias, the outcome was downgraded one level.</p> <p>Very serious: If greater than 33.3% of the weight in a meta-analysis came from studies at high risk of bias, the outcome was downgraded two levels.</p>
Indirectness	<p>Not serious: If less than 33.3% of the weight in a meta-analysis came from partially indirect or indirect studies, the overall outcome was not downgraded.</p> <p>Serious: If greater than 33.3% of the weight in a meta-analysis came from partially indirect or indirect studies, the outcome was downgraded one level.</p> <p>Very serious: If greater than 33.3% of the weight in a meta-analysis came from indirect studies, the outcome was downgraded two levels.</p>
Inconsistency	<p>Concerns about inconsistency of effects across studies, occurring when there is unexplained variability in the treatment effect demonstrated across studies (heterogeneity), after appropriate pre-specified subgroup analyses have been conducted. This was assessed using the I^2 statistic.</p> <p>N/A: Inconsistency was marked as not applicable if data on the outcome was only available from one study.</p> <p>Not serious: If the I^2 was less than 33.3%, the outcome was not downgraded.</p> <p>Serious: If the I^2 was between 33.3% and 66.7%, the outcome was downgraded one level.</p> <p>Very serious: If the I^2 was greater than 66.7%, the outcome was downgraded two levels.</p>
Imprecision	<p>This was not included in the GRADE table, but was considered during committee discussions of the evidence, taking into account 95% confidence intervals around the point estimate of the effect, any relevant MIDs, committee expertise and the effect of a single intervention based on multiple outcomes.</p>

3 The quality of evidence for each outcome was upgraded if any of the following three
4 conditions were met:

- 5 • Data from non-randomised studies showing an effect size sufficiently large that it
6 cannot be explained by confounding alone.
- 7 • Data showing a dose-response gradient.
- 8 • Data where all plausible residual confounding is likely to increase our confidence
9 in the effect estimate.

10 Publication bias

11 Where 10 or more studies were included as part of a single meta-analysis, a funnel
12 plot was produced to graphically assess the potential for publication bias. When a
13 funnel plot showed convincing evidence of publication bias, or the review team
14 became aware of other evidence of publication bias (for example, evidence of
15 unpublished trials where there was evidence that the effect estimate differed in
16 published and unpublished data), the outcome was downgraded once. If no
17 evidence of publication bias was found for any outcomes in a review (as was often
18 the case), this domain was excluded from GRADE profiles to improve readability.

1

2 **Methods for combining direct and indirect evidence** 3 **(network meta-analysis) for interventions**

4 Hierarchical Bayesian Network Meta-Analysis (NMA) was performed using WinBUGS
5 version 1.4.3. The models used reflected the recommendations of the NICE Decision
6 Support Unit's Technical Support Documents (TSDs) on evidence synthesis,
7 particularly TSD 2 ('A generalised linear modelling framework for pairwise and
8 network meta-analysis of randomised controlled trials'; see
9 <http://www.nicedsu.org.uk>). The WinBUGS code provided in the appendices of the
10 TSDs was used without substantive alteration to specify synthesis models.

11 Three separate chains with different initial values were used. Results were assessed
12 for convergence to determine the length of 'burn-in' period required by examining the
13 'bgdiag' and 'history' plots. Results were reported summarising at least 10,000
14 samples from the posterior distribution of each model, having run and discarded the
15 'burn-in' iterations. The MC error was assessed to check that it was sufficiently small
16 (less than 5% of the standard deviation of the posterior distribution for each
17 parameter) and additional samples were summarised if this was the case.

18 Non-informative prior distributions were used in all models. Unless otherwise
19 specified, trial-specific baselines and treatment effects were assigned Normal (0,
20 10000) priors, and the between-trial standard deviations used in random-effects
21 models for dichotomous outcomes were given Uniform (0, 5) priors. These are
22 consistent with the recommendations in TSD 2 for dichotomous outcomes. For the 1
23 continuous outcome (length of stay) we used Uniform (0, 10) priors, which is
24 substantially wider than the expected variance for this outcome.

25 Fixed-effect and random-effects models were explored for each outcome, with the
26 final choice of model based on the total residual deviance and deviance information
27 criterion (DIC): if DIC was at least 3 points lower for the random-effects model, it was
28 preferred; otherwise, the fixed-effect model was considered to provide an equivalent
29 fit to the data in a more parsimonious analysis, and was preferred.

30 Inconsistency between direct and indirect evidence was assessed by inspecting NMA
31 outputs compared with pairwise direct results. We also fitted 'inconsistency models'
32 and compared the posterior mean deviance contribution for each data point to the
33 analogous value from the main ('consistency') models to identify areas of
34 inconsistency (see TSD 4 'Inconsistency in networks of evidence based on
35 randomised controlled trials'; <http://www.nicedsu.org.uk>). Datapoints with an absolute
36 difference in deviance of 0.5 or greater between the 2 models were deemed worthy
37 of additional investigation, to see if any reason for underlying heterogeneity could be
38 established.

39 **Modified GRADE for network meta-analyses**

40 A modified version of the standard GRADE approach for pairwise interventions was
41 used to assess the quality of evidence across the network meta-analyses
42 undertaken. While most criteria for pairwise meta-analyses still apply, it is important
43 to adapt some of the criteria to take into consideration additional factors, such as how
44 each 'link' or pairwise comparison within the network applies to the others. As a

1 result, the following was used when modifying the GRADE framework to a network
 2 meta-analysis. It is designed to provide a single overall quality rating for an NMA,
 3 which can then be combined with pairwise quality ratings for individual comparisons
 4 (if appropriate), to judge the overall strength of evidence for each comparison.
 5

6 **Table 3: Rationale for downgrading quality of evidence for intervention studies**

GRADE criteria	Reasons for downgrading quality
Risk of bias	Not serious: If fewer than 33.3% of the studies in the network meta-analysis were at moderate or high risk of bias, the overall network was not downgraded. Serious: If greater than 33.3% of the studies in the network meta-analysis were at moderate or high risk of bias, the network was downgraded one level. Very serious: If greater than 33.3% of the studies in the network meta-analysis were at high risk of bias, the network was downgraded two levels.
Indirectness	Not serious: If fewer than 33.3% of the studies in the network meta-analysis were partially indirect or indirect, the overall network was not downgraded. Serious: If greater than 33.3% of the studies in the network meta-analysis were partially indirect or indirect, the network was downgraded one level. Very serious: If greater than 33.3% of the studies in the network meta-analysis were indirect, the network was downgraded two levels.
Inconsistency	N/A: Inconsistency was marked as not applicable if there were no links in the network where data from multiple studies (either direct or indirect) were synthesised. For network meta-analyses conducted under a Bayesian framework, the network was downgraded one level if the DIC for an inconsistency model was more than 3 points higher than the corresponding consistency model.
Imprecision	This was not included in the GRADE table, but was considered during committee discussions of the evidence, taking into account 95% confidence intervals around the point estimate of the effect, any relevant MIDs, committee expertise and the effect of a single intervention based on multiple outcomes..

7 Diagnostic test accuracy evidence

8 In this guideline, diagnostic test accuracy (DTA) data are classified as any data in
 9 which a feature – be it a symptom, a risk factor, a test result or the output of some
 10 algorithm that combines many such features – is observed in some people who have
 11 the condition of interest at the time of the test and some people who do not. Such
 12 data either explicitly provide, or can be manipulated to generate, a 2x2 classification
 13 of true positives and false negatives (in people who, according to the reference
 14 standard, truly have the condition) and false positives and true negatives (in people
 15 who, according to the reference standard, do not).

16 The 'raw' 2x2 data can be summarised in a variety of ways. Those that were used for
 17 decision making in this guideline are as follows:

- 18 • **Positive likelihood ratios** describe how many times more likely positive features
 19 are in people with the condition compared to people without the condition. Values
 20 greater than 1 indicate that a positive result makes the condition more likely.

21 ○ $LR^+ = (TP/[TP+FN]) / (FP/[FP+TN])$

1 • **Negative likelihood ratios** describe how many times less likely negative features
2 are in people with the condition compared to people without the condition. Values
3 less than 1 indicate that a negative result makes the condition less likely.

4 ○ $LR^- = (FN/[TP+FN]) / (TN/[FP+TN])$

5 • **Sensitivity** is the probability that the feature will be positive in a person with the
6 condition.

7 ○ $sensitivity = TP / (TP + FN)$

8 • **Specificity** is the probability that the feature will be negative in a person without
9 the condition.

10 ○ $specificity = TN / (FP + TN)$

11 The following schema, adapted from the suggestions of Jaeschke et al. (1994), was
12 used to interpret the likelihood ratio findings from diagnostic test accuracy reviews.

13
14 **Table 4: Interpretation of likelihood ratios**

Value of likelihood ratio	Interpretation
$LR \leq 0.1$	Very large decrease in probability of disease
$0.1 < LR \leq 0.2$	Large decrease in probability of disease
$0.2 < LR \leq 0.5$	Moderate decrease in probability of disease
$0.5 < LR \leq 1.0$	Slight decrease in probability of disease
$1.0 < LR < 2.0$	Slight increase in probability of disease
$2.0 \leq LR < 5.0$	Moderate increase in probability of disease
$5.0 \leq LR < 10.0$	Large increase in probability of disease
$LR \geq 10.0$	Very large increase in probability of disease

15 The schema above has the effect of setting a clinical decision threshold for positive
16 likelihoods ratio at 2, and a corresponding clinical decision threshold for negative
17 likelihood ratios at 0.5. Likelihood ratios (whether positive or negative) falling
18 between these thresholds were judged to indicate no meaningful change in the
19 probability of disease.

20 Quality assessment

21 Individual studies were quality assessed using the QUADAS-2 tool, which contains
22 four domains: patient selection, index test, reference standard, and flow and timing.
23 Each individual study was classified into one of the following three groups:

- 24 • Low risk of bias – The true effect size for the study is likely to be close to the
25 estimated effect size.
- 26 • Moderate risk of bias – There is a possibility the true effect size for the study is
27 substantially different to the estimated effect size.
- 28 • High risk of bias – It is likely the true effect size for the study is substantially
29 different to the estimated effect size.

1 Each individual study was also classified into one of three groups for directness,
2 based on if there were concerns about the population, index features and/or
3 reference standard in the study and how directly these variables could address the
4 specified review question. Studies were rated as follows:

- 5 • Direct – No important deviations from the protocol in population, index feature
6 and/or reference standard.
- 7 • Partially indirect – Important deviations from the protocol in one of the population,
8 index feature and/or reference standard.
- 9 • Indirect – Important deviations from the protocol in at least two of the population,
10 index feature and/or reference standard.

11 **Methods for combining diagnostic test accuracy evidence**

12 Meta-analysis of diagnostic test accuracy data was conducted with reference to the
13 Cochrane Handbook for Systematic Reviews of Diagnostic Test Accuracy (Deeks et
14 al. 2010).

15 Where applicable, diagnostic syntheses were stratified by:

- 16 • Presenting symptomatology (features shared by all participants in the study, but
17 not all people who could be considered for a diagnosis in clinical practice).
- 18 • The reference standard used for true diagnosis.

19 Where five or more studies were available for all included strata, a bivariate model
20 was fitted using the mada package in R v3.4.0, which accounts for the correlations
21 between positive and negative likelihood ratios, and between sensitivities and
22 specificities. Where sufficient data were not available (2-4 studies), separate
23 independent pooling was performed for positive likelihood ratios, negative likelihood
24 ratios, sensitivity and specificity, using Microsoft Excel. This approach is conservative
25 as it is likely to somewhat underestimate test accuracy, due to failing to account for
26 the correlation and trade-off between sensitivity and specificity (see Deeks 2010).

27 Random-effects models (der Simonian and Laird) were fitted for all syntheses, as
28 recommended in the Cochrane Handbook for Systematic Reviews of Diagnostic Test
29 Accuracy (Deeks et al. 2010).

30 **GRADE for diagnostic test accuracy evidence**

31 The choice of primary outcome for decision making was determined by the
32 committee and GRADE assessments were undertaken using the appropriate method
33 from those listed below.

34 In all cases, following completion of the GRADE table, the downstream effects of
35 these tests on patient- important outcomes were considered. This could be done
36 explicitly during committee deliberations and reported as part of the discussion
37 section of the review detailing the likely consequences of true positive, true negative,
38 false positive and false negative test results. Alternatively, in reviews where a
39 decision model is being carried (for example, as part of an economic analysis), these
40 consequences may be incorporated here instead.

1 Using likelihood ratios as the primary outcomes

2 GRADE assessments were only undertaken for positive and negative likelihood
3 ratios, as the clinical decision thresholds used to assess imprecision were based on
4 these outcomes, but results for sensitivity and specificity are also presented
5 alongside those data.

6 Evidence from diagnostic accuracy studies was initially rated as high-quality, and
7 then downgraded according to the standard GRADE criteria (risk of bias,
8 inconsistency, imprecision and indirectness) as detailed in Table 5: Rationale for
9 downgrading quality of evidence for diagnostic questions using likelihood ratio
10 measures.

11 The committee were consulted to set 2 clinical decision thresholds for each measure:
12 the likelihood ratio above (or below for negative likelihood ratios) which a test would
13 be recommended, and a second below (or above for negative likelihood ratios) which
14 a test would be considered of no clinical use. These were used to judge imprecision
15 (see below). If the committee were unsure which values to pick, then the default
16 values of 2 for LR+ and 0.5 for LR- were used based on

17 Table , with the line of no effect as the second clinical decision line in both cases.

18

19 **Table 5: Rationale for downgrading quality of evidence for diagnostic questions using**
20 **likelihood ratio measures.**

GRADE criteria	Reasons for downgrading quality
Risk of bias	<p>Not serious: If less than 33.3% of the weight in a meta-analysis came from studies at moderate or high risk of bias, the overall outcome was not downgraded.</p> <p>Serious: If greater than 33.3% of the weight in a meta-analysis came from studies at moderate or high risk of bias, the outcome was downgraded one level.</p> <p>Very serious: If greater than 33.3% of the weight in a meta-analysis came from studies at high risk of bias, the outcome was downgraded two levels.</p>
Indirectness	<p>Not serious: If less than 33.3% of the weight in a meta-analysis came from partially indirect or indirect studies, the overall outcome was not downgraded.</p> <p>Serious: If greater than 33.3% of the weight in a meta-analysis came from partially indirect or indirect studies, the outcome was downgraded one level.</p> <p>Very serious: If greater than 33.3% of the weight in a meta-analysis came from indirect studies, the outcome was downgraded two levels.</p>
Inconsistency	<p>Concerns about inconsistency of effects across studies, occurring when there is unexplained variability in the treatment effect demonstrated across studies (heterogeneity), after appropriate pre-specified subgroup analyses have been conducted. This was assessed using the I^2 statistic.</p> <p>N/A: Inconsistency was marked as not applicable if data on the outcome was only available from one study.</p> <p>Not serious: If the I^2 was less than 33.3%, the outcome was not downgraded.</p> <p>Serious: If the I^2 was between 33.3% and 66.7%, the outcome was downgraded one level.</p> <p>Very serious: If the I^2 was greater than 66.7%, the outcome was downgraded two levels.</p>

21 Upgrading of evidence was carried out with caution, and was considered when:

- 1 • Test accuracy was extremely high (and the associated confidence intervals also
2 only include extremely high accuracy estimates)
3 • All plausible confounding would act to reduce test accuracy

4 **Publication bias**

5 If the review team became aware of evidence of publication bias (for example,
6 evidence of unpublished trials where there was evidence that the effect estimate
7 differed in published and unpublished data), the outcome was downgraded once. If
8 no evidence of publication bias was found for any outcomes in a review (as was often
9 the case), this domain was excluded from GRADE profiles to improve readability.

10 **Predictive (prognostic) accuracy evidence**

11 In this guideline, predictive (or prognostic) accuracy data are classified as any data
12 in which a feature – be it a symptom, a risk factor, a test result or the output of some
13 algorithm that combines many such features – is observed in some people who go
14 on to develop the condition of interest and some people who do not. Such data either
15 explicitly provide, or can be manipulated to generate, a 2x2 classification of true
16 positives and false negatives (in people who, according to the reference standard,
17 truly develop the condition) and false positives and true negatives (in people who,
18 according to the reference standard, do not). This category would include studies
19 classed as prediction models under the TRIPOD statement, provided the data were
20 reported a 2x2 classification data.

21 The ‘raw’ 2x2 data can be summarised in a variety of ways. Those that were used for
22 decision making in this guideline are as follows:

- 23 • **Positive likelihood ratios** describe how many times more likely positive features
24 are in people who develop the condition compared to people who do not. Values
25 greater than 1 indicate that a positive result makes the condition more likely.
26 ○ $LR^+ = (TP/[TP+FN])/(FP/[FP+TN])$
- 27 • **Negative likelihood ratios** describe how many times less likely negative features
28 are in people who develop the condition compared to people who do not. Values
29 less than 1 indicate that a negative result makes the condition less likely.
30 ○ $LR^- = (FN/[TP+FN])/(TN/[FP+TN])$
- 31 • **Sensitivity** is the probability that the feature will be positive in a person who goes
32 on to develop the condition.
33 ○ $sensitivity = TP/(TP+FN)$
- 34 • **Specificity** is the probability that the feature will be negative in a person who does
35 not go on to develop the condition.
36 ○ $specificity = TN/(FP+TN)$

37 The following schema, adapted from the suggestions of Jaeschke et al. (1994), was
38 used to interpret the findings from predictive accuracy evidence.

39

40 **Table 6: Interpretation of likelihood ratios**

Value of likelihood ratio	Interpretation
$LR \leq 0.1$	Very large decrease in probability of disease
$0.1 < LR \leq 0.2$	Large decrease in probability of disease
$0.2 < LR \leq 0.5$	Moderate decrease in probability of disease
$0.5 < LR \leq 1.0$	Slight decrease in probability of disease
$1.0 < LR < 2.0$	Slight increase in probability of disease
$2.0 \leq LR < 5.0$	Moderate increase in probability of disease
$5.0 \leq LR < 10.0$	Large increase in probability of disease
$LR \geq 10.0$	Very large increase in probability of disease

1 The schema above has the effect of setting a clinical decision threshold for positive
2 likelihoods ratio at 2, and a corresponding clinical decision threshold for negative
3 likelihood ratios at 0.5. Likelihood ratios (whether positive or negative) falling
4 between these thresholds were judged to indicate no meaningful change to
5 probability of disease.

6 Quality assessment

7 Individual studies reporting clinical prediction models were quality assessed using the
8 PROBAST tool, which contains five domains: participant selection, predictors,
9 outcome, sample size and participant flow, analysis ([Wolff et al. 2018](#)). Cohort
10 studies reporting other predictive accuracy data were quality assessed using the
11 QUIPs checklist. Each individual study was classified into one of the following three
12 groups based on an assessment of the overall risk of bias:

- 13 • Low risk of bias – The true effect size for the study is likely to be close to the
14 estimated effect size.
- 15 • Moderate risk of bias – There is a possibility the true effect size for the study is
16 substantially different to the estimated effect size.
- 17 • High risk of bias – It is likely the true effect size for the study is substantially
18 different to the estimated effect size.

19 Each individual study was also classified into one of three groups for directness,
20 based on if there were concerns about the population, predictive features and/or
21 reference standard in the study and how directly these variables could address the
22 specified review question. Studies were rated as follows:

- 23 • Direct – No important deviations from the protocol in population, predictive feature
24 and/or reference standard.
- 25 • Partially indirect – Important deviations from the protocol in one of the population,
26 predictive feature and/or reference standard.
- 27 • Indirect – Important deviations from the protocol in at least two of the population,
28 predictive feature and/or reference standard.

29 Methods for combining predictive accuracy evidence

30 Where applicable, predictive accuracy syntheses were stratified by:

- 31 • Presenting symptomatology (features shared by all participants in the study, but
32 not all people in the full relevant clinical population).

- 1 • The length of time between the measurement of the predictive feature and the
2 final outcome.
- 3 • The reference standard used for categorising true positives.
- 4 Where five or more studies were available for all included strata, a bivariate model
5 was fitted using the mada package in R v3.4.0, which accounts for the correlations
6 between positive and negative likelihood ratios, and between sensitivities and
7 specificities. Where sufficient data were not available (2-4 studies), separate
8 independent pooling was performed for positive likelihood ratios, negative likelihood
9 ratios, sensitivity and specificity, using Microsoft Excel. This approach is likely to
10 somewhat underestimate test accuracy (see Deeks 2001).
- 11 Random-effects models (der Simonian and Laird) were fitted for all syntheses, due to
12 the expected level of between study heterogeneity in prognostic reviews.

13 **Modified GRADE for predictive accuracy evidence**

14 GRADE has not been developed for use with predictive accuracy studies; therefore a
15 modified approach was applied using the GRADE framework. GRADE assessments
16 were only undertaken for positive and negative likelihood ratios, as the clinical
17 decision thresholds used to assess imprecision were based on these outcomes.

18 Cohort studies were initially rated as high-quality evidence if well conducted, and
19 then downgraded according to the standard GRADE criteria (risk of bias,
20 inconsistency, imprecision and indirectness) as detailed in Table below.

21

22 **Table 7: Rationale for downgrading quality of evidence for evidence reporting**
23 **predictive accuracy data**

GRADE criteria	Reasons for downgrading quality
Risk of bias	<p>Not serious: If less than 33.3% of the weight in a meta-analysis came from studies at moderate or high risk of bias, the overall outcome was not downgraded.</p> <p>Serious: If greater than 33.3% of the weight in a meta-analysis came from studies at moderate or high risk of bias, the outcome was downgraded one level.</p> <p>Very serious: If greater than 33.3% of the weight in a meta-analysis came from studies at high risk of bias, the outcome was downgraded two levels.</p>
Indirectness	<p>Not serious: If less than 33.3% of the weight in a meta-analysis came from partially indirect or indirect studies, the overall outcome was not downgraded.</p> <p>Serious: If greater than 33.3% of the weight in a meta-analysis came from partially indirect or indirect studies, the outcome was downgraded one level.</p> <p>Very serious: If greater than 33.3% of the weight in a meta-analysis came from indirect studies, the outcome was downgraded two levels.</p>

GRADE criteria	Reasons for downgrading quality
Inconsistency	<p>Concerns about inconsistency of effects across studies, occurring when there is unexplained variability in the treatment effect demonstrated across studies (heterogeneity), after appropriate pre-specified subgroup analyses have been conducted. This was assessed using the I^2 statistic.</p> <p>N/A: Inconsistency was marked as not applicable if data on the outcome was only available from one study.</p> <p>Not serious: If the I^2 was less than 33.3%, the outcome was not downgraded.</p> <p>Serious: If the I^2 was between 33.3% and 66.7%, the outcome was downgraded one level.</p> <p>Very serious: If the I^2 was greater than 66.7%, the outcome was downgraded two levels.</p>
Imprecision	<p>If the 95% confidence interval for sensitivity crossed one of the clinical decision thresholds, the outcome was downgraded one level, as the data were deemed to be imprecise. If the 95% confidence interval spanned both thresholds, the outcome was downgraded twice, as suffering from very serious imprecision. Specificity was assessed for imprecision in a similar manner using the 2 previously defined clinical decision thresholds.</p>

1 Publication bias

2 If the review team became aware of evidence of publication bias (for example,
3 evidence of unpublished trials where there was evidence that the effect estimate
4 differed in published and unpublished data), the outcome was downgraded once. If
5 no evidence of publication bias was found for any outcomes in a review (as was often
6 the case), this domain was excluded from GRADE profiles to improve readability.

7 Other prognostic evidence

8 Other prognostic studies were also included if they reported outcomes of c-statistics.
9 Hazard ratios were also included in the review protocols for some reviews, but no
10 data were found for these outcomes. These studies were also quality assessed using
11 the PROBAST checklist (in the case of studies reporting clinical prediction models) or
12 the QUIPs checklist (in the case of other prognostic studies), as in the predictive
13 accuracy section above.

14 Assessing c-statistics

15 C-statistics were assessed in a similar manner to likelihood ratios using the
16 categories in Table below.

17 **Table 8 Interpretation of c-statistics (Hosmer DW Jr, Lemeshow S, Sturdivant**
18 **RX. Applied logistic regression: John Wiley & Sons; 2013.)**

Value of c-statistic	Interpretation
c-statistic <0.6	Poor classification accuracy
0.6 ≤ c-statistic <0.7	Adequate classification accuracy
0.7 ≤ c-statistic <0.8	Good classification accuracy
0.8 ≤ c-statistic <0.9	Excellent classification accuracy
0.9 ≤ c-statistic < 1.0	Outstanding classification accuracy

1 Meta-analyses were carried out using the metamisc package in R v3.4.0, which
2 confines the analysis results to between 0 and 1 matching the limited range of values
3 that c-statistics can take. Random effects meta-analysis was used when the I^2 was
4 50% or greater.

5 In any meta-analyses where some (but not all) of the data came from studies at high
6 risk of bias, a sensitivity analysis was conducted, excluding those studies from the
7 analysis. Results from both the full and restricted meta-analyses are reported.
8 Similarly, in any meta-analyses where some (but not all) of the data came from
9 indirect studies, a sensitivity analysis was conducted, excluding those studies from
10 the analysis.

11 A modified version of GRADE was carried out to assess the quality of the meta-
12 analysed c-statistics as follows:

- 13 • Imprecision - the 95% CI boundaries were examined and if they crossed 2
14 categories of test classification accuracy as shown in Table 8 Interpretation of c-
15 statistics (Hosmer DW Jr, Lemeshow S, Sturdivant RX. Applied logistic
16 regression: John Wiley & Sons; 2013.)

17 then the study was downgraded once (imprecision rated as serious); if the
18 boundaries crossed 3 categories then the study was downgraded twice (very serious
19 imprecision).

- 20 • Inconsistency, indirectness and risk of bias were determined using the methods in
21 the section on GRADE for predictive accuracy evidence.

22 In cases where meta-analyses could not be carried out due to the large numbers of
23 studies without 95% CI, the following decision rules were used to assess risk of bias,
24 indirectness, imprecision and inconsistency for each outcome:

25 1. Risk of bias and indirectness were calculated as normal, but using the study
26 weight by population, rather than weight in the meta-analysis.

27 2. Imprecision

- 28 a. Single study with 95% CI: the 95% CI boundaries were examined and if
29 they crossed 2 categories of test classification accuracy then the study
30 was downgraded once (imprecision rated as serious); if the boundaries
31 crossed 3 categories then the study was downgraded twice (very serious
32 imprecision).
 - 33 i. Multiple studies with 95% CI: the individual studies were rated as
34 in a. and then if >33.3% of the studies by population weight were
35 rated serious then the analysis was downgraded once; if > 33.33%
36 were rated very serious the analysis was downgraded twice.
- 37 b. Single study or multiple studies without 95% CI: the mean sample size
38 was calculated and if this was < 250 then the analysis was downgraded
39 twice (very serious); if it was >250, but > 500 the analysis was
40 downgraded once (serious); if the mean was > 500 people/study then the
41 analysis was not downgraded (not serious).
- 42 c. Multiple studies with and without 95% CI: the studies without 95% CI were
43 analysed as in 2c; those with 95% CI were analysed as in 2b. The results
44 were averaged, but the number of studies in each group were also taken

- 1 into account with the result that if there were a lot more studies in one
2 group compared to the other then that group rating would be used. In
3 general, not serious and serious or not serious and very serious were
4 averaged to serious; serious and very serious resulted in a very serious
5 rating.
- 6 3. Inconsistency
- 7 a. Single study with or without 95% CI: N/A
- 8 b. Multiple studies with or without 95% CI: the highest and lowest point
9 estimates were examined. If they spanned < 2 categories of c-statistic
10 classification accuracy the analysis was rated as not serious for
11 inconsistency; if they spanned 2 categories this was rated as serious and
12 ≥ 3 categories was rated as very serious.
- 13

14 Association studies

15 In this guideline, association studies are defined those reporting data showing an
16 association of a predictor (either a single variable or a group of variables) and an
17 outcome variable, where the data are not reported in terms of outcome classification
18 (i.e. diagnostic/predictive accuracy). Data were reported as hazard ratios (if
19 measured over time) or odds ratios or risk ratios (if measured at a specific time-
20 point). Data reported in terms of model fit or predictive accuracy were not assessed
21 using this method.

22 Quality assessment

23 Individual cohort studies were quality assessed using the QUIPS checklist. Each
24 individual study was classified into one of the following three groups:

- 25 • Low risk of bias – The true effect size for the study is likely to be close to the
26 estimated effect size.
- 27 • Moderate risk of bias – There is a possibility the true effect size for the study is
28 substantially different to the estimated effect size.
- 29 • High risk of bias – It is likely the true effect size for the study is substantially
30 different to the estimated effect size.

31 Individual cross-sectional studies were quality assessed using the Joanna Briggs
32 Institute critical appraisal checklist for analytical cross sectional studies (2016), which
33 contains 8 questions covering: inclusion criteria, description of the sample, measures
34 of exposure, measures of outcomes, confounding factors, and statistical analysis.
35 Each individual study was classified into one of the following groups:

- 36 • Low risk of bias – The true effect size for the study is likely to be close to the
37 estimated effect size.
- 38 • Moderate risk of bias – There is a possibility the true effect size for the study is
39 substantially different to the estimated effect size.
- 40 • High risk of bias – It is likely the true effect size for the study is substantially
41 different to the estimated effect size.

1 Each individual study was also classified into one of three groups for directness,
2 based on if there were concerns about the population, predictors and/or outcomes in
3 the study and how directly these variables could address the specified review
4 question. Studies were rated as follows:

- 5 • Direct – No important deviations from the protocol in population, predictors and/or
6 outcomes.
- 7 • Partially indirect – Important deviations from the protocol in one of the population,
8 predictors and/or outcomes.
- 9 • Indirect – Important deviations from the protocol in at least two of the population,
10 predictors and/or outcomes.

11 **Methods for combining association studies**

12 Adjusted odds ratios, hazard ratios and risk ratios from multivariate models were only
13 considered for pooling if the same set of predictor variables were used across
14 multiple studies, the same thresholds to measure predictors were used across
15 studies, and the studies controlled for the same confounding factors. This was not
16 the case for any data in this guideline and so data was presented separately for
17 individual studies.

18 **Modified GRADE for association studies**

19 GRADE has not been developed for use with association studies; therefore a
20 modified approach was applied using the GRADE framework. Data from cohort
21 studies and cross sectional studies was initially rated as high quality, with the quality
22 of the evidence for each outcome then downgraded or not from this initial point.

23

24 **Table 9: Rationale for downgrading quality of evidence for association studies**

GRADE criteria	Reasons for downgrading quality
Risk of bias	Not serious: If the outcome was from a study judged at low risk of bias the outcome was not downgraded. Serious: If the outcome was from a study judged at moderate risk of bias the outcome was downgraded one level. Very serious: If the outcome was from a study that was judged at high risk of bias, the outcome was downgraded two levels.
Indirectness	Not serious: If the outcome was from a study judged directly applicable, the overall outcome was not downgraded. Serious: If the outcome was from a study judged partially applicable, the outcome was downgraded one level. Very serious: If the outcome was from a study judged indirectly applicable, the outcome was downgraded two levels.
Inconsistency	Results were not synthesised and were presented for individual studies. Inconsistency could therefore not be assessed as was rated as 'not applicable'.
Imprecision	The outcome was downgraded once if the 95% confidence interval for the effect size crossed the line of no effect (i.e. the outcome was not statistically significant), and twice if the sample size of the study was sufficiently small that it is not plausible any realistic effect size could have been detected.

1 Qualitative evidence

2 Quality assessment

3 Individual qualitative studies were quality assessed using the CASP qualitative
4 checklist. Each individual study was classified into one of the following three groups:

- 5 • Low risk of bias – The findings and themes identified in the study are likely to
6 accurately capture the true picture.
- 7 • Moderate risk of bias – There is a possibility the findings and themes identified in
8 the study are not a complete representation of the true picture.
- 9 • High risk of bias – It is likely the findings and themes identified in the study are not
10 a complete representation of the true picture

11 Each individual study was also classified into one of three groups for relevance,
12 based on if there were concerns about the perspective, population, phenomenon of
13 interest and/or setting in the included studies and how directly these variables could
14 address the specified review question. Studies were rated as follows:

- 15 • Highly relevant – No important deviations from the protocol in perspective,
16 population, phenomenon of interest and/or setting.
- 17 • Relevant – Important deviations from the protocol in one of the perspective,
18 population, phenomenon of interest and/or setting.
- 19 • Partially relevant – Important deviations from the protocol in at least two of the
20 perspective, population, phenomenon of interest and/or setting.

21 Methods for combining qualitative evidence

22 Only one qualitative study was included within the reviews. If multiple qualitative
23 studies had been identified for a single question, then information from the studies
24 would have been combined using a thematic synthesis. Instead, the main themes
25 were extracted from the single study and were evaluated to examine their relevance
26 to the review question. Each relevant theme was then presented to the committee.

27 CERQual for qualitative studies

28 CERQual was used to assess the confidence we have in the summary findings of
29 each of the identified themes. Evidence from all qualitative study designs (interviews,
30 focus groups etc.) was initially rated as high confidence and the confidence in the
31 evidence for each theme was then downgraded from this initial point as detailed in
32 Table101 Rationale for downgrading confidence in evidence for qualitative questions

33

34 Table101 Rationale for downgrading confidence in evidence for qualitative questions

CERQual criteria	Reasons for downgrading confidence
Methodological limitations	<p>Not serious: If the theme was identified in studies at low risk of bias, the outcome was not downgraded</p> <p>Serious: If the theme was identified only in studies at moderate or high risk of bias, the outcome was downgraded one level.</p> <p>Very serious: If the theme was identified only in studies at high risk of bias, the outcome was downgraded two levels.</p>

CERQual criteria	Reasons for downgrading confidence
Relevance	High: If the theme was identified in highly relevant studies, the outcome was not downgraded Moderate: If the theme was identified only in relevant and partially relevant studies, the outcome was downgraded one level. Low: If the theme was identified only in partially relevant studies, the outcome was downgraded two levels.
Coherence	Coherence was addressed based on two factors: Between study – does the theme consistently emerge from all relevant studies Theoretical – does the theme provide a convincing theoretical explanation for the patterns found in the data The outcome was downgraded once if there were concerns about one of these elements of coherence, and twice if there were concerns about both elements.
Adequacy of data	The outcome was downgraded if there was insufficient data to develop an understanding of the phenomenon of interest, either due to insufficient studies, participants or observations.

1 Health economics

2 The methods for the de novo models built for this guideline can be found in the
3 appendices of their respective evidence reviews. Literature reviews seeking to
4 identify published cost–utility analyses of relevance to the issues under consideration
5 were conducted for all questions. In each case, the search undertaken for the clinical
6 review was modified, retaining population and intervention descriptors, but removing
7 any study-design filter and adding a filter designed to identify relevant health
8 economic analyses. In assessing studies for inclusion, population, intervention and
9 comparator, criteria were always identical to those used in the parallel clinical search;
10 only cost–utility analyses were included. Economic evidence profiles, including
11 critical appraisal according to the Guidelines manual, were completed for included
12 studies.

13 Economic studies identified through a systematic search of the literature are
14 appraised using a methodology checklist designed for economic evaluations (NICE
15 guidelines manual; 2014). This checklist is not intended to judge the quality of a
16 study per se, but to determine whether an existing economic evaluation is useful to
17 inform the decision-making of the committee for a specific topic within the guideline.

18 There are 2 parts of the appraisal process. The first step is to assess applicability
19 (that is, the relevance of the study to the specific guideline topic and the NICE
20 reference case); evaluations are categorised according to the criteria in Table 10.

21

22

Table 11 Applicability criteria

Level	Explanation
Directly applicable	The study meets all applicability criteria, or fails to meet one or more applicability criteria but this is unlikely to change the conclusions about cost effectiveness
Partially applicable	The study fails to meet one or more applicability criteria, and this could change the conclusions about cost effectiveness
Not applicable	The study fails to meet one or more applicability criteria, and this is likely to change the conclusions about cost

Level	Explanation
	effectiveness. These studies are excluded from further consideration

1 In the second step, only those studies deemed directly or partially applicable are
 2 further assessed for limitations (that is, methodological quality); see categorisation
 3 criteria in Table 11 Applicability criteria

4

5 **Table 12 Methodological criteria**

Level	Explanation
Minor limitations	Meets all quality criteria, or fails to meet one or more quality criteria but this is unlikely to change the conclusions about cost effectiveness
Potentially serious limitations	Fails to meet one or more quality criteria and this could change the conclusions about cost effectiveness
Very serious limitations	Fails to meet one or more quality criteria and this is highly likely to change the conclusions about cost effectiveness. Such studies should usually be excluded from further consideration

6 Where relevant, a summary of the main findings from the systematic search, review
 7 and appraisal of economic evidence is presented in an economic evidence profile
 8 alongside the clinical evidence.