

**National Institute for Health and
Care Excellence**

Advanced breast cancer: diagnosis and management

NICE guideline CG81

Methods

March 2026

Draft for consultation

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 All NICE guidance is subject to regular review and may be updated or withdrawn.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 2026. All rights reserved. Subject to [Notice of rights](#).

ISBN:

Contents

Development of the guideline	4
Remit	4
Methods	5
Developing the review questions and outcomes.....	5
Reviewing research evidence	6
Review protocols	6
Searching for evidence	6
Selecting studies for inclusion	6
Incorporating published evidence syntheses.....	6
Methods of combining evidence.....	9
Data synthesis for intervention studies.....	9
Pairwise meta-analysis.....	9
Data synthesis for diagnostic accuracy data.....	10
Appraising the quality of evidence.....	12
Intervention studies (relative effect estimates).....	12
Diagnostic accuracy studies	15
Reviewing economic evidence.....	19
References.....	24

1 **Development of the guideline**

2 **Remit**

3 The methods outlined in this document relate to the update of [Advanced breast](#)
4 [cancer: diagnosis and management \(CG81\)](#) published in 2026.

5 This methods document does not cover recommendations published as part of the
6 original guideline in 2009. Methods for these recommendations are available as part
7 of the original [evidence review](#).

8 To see “What this guideline covers” and “What this guideline does not cover”, see the
9 guideline scope for [Advanced breast cancer: diagnosis and management](#).

1 **Methods**

2 This evidence review was developed using the methods and process described in
3 [Developing NICE guidelines: the manual](#).

4 Declarations of interest were recorded according to [NICE's conflicts of interest policy](#).

5 **Developing the review questions and outcomes**

6
7 The 3 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
9 development team and refined and validated by the guideline committee.

10 The review questions were based on the following frameworks:

- 11 • Population, Intervention, Comparator and Outcome and Study type (PICOS) for
12 reviews of interventions
- 13 • Population, index test(s), reference standard and outcome for reviews of
14 diagnostic and predictive accuracy

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

1 **Reviewing research evidence**

2 **Review protocols**

3 Review protocols were developed with the guideline committee to outline the
4 inclusion and exclusion criteria used to select studies for each evidence review.

5 **Searching for evidence**

6 Evidence was searched for each review question using the methods specified in
7 [Developing NICE guidelines: the manual](#).

8 **Selecting studies for inclusion**

9 All references identified by the literature searches and from other sources (for
10 example, previous versions of the guideline or studies identified by committee
11 members) were uploaded into EPPI reviewer software (version 5) and de-duplicated.
12 Titles and abstracts were assessed for possible inclusion using the criteria specified
13 in the review protocol. At least 10% of the abstracts were reviewed by two reviewers,
14 with any disagreements resolved by discussion or, if necessary, a third independent
15 reviewer.

16 Priority screening was not used for any reviews in this guideline and therefore the full
17 database was screened for each review.

18 The full text of potentially eligible studies was retrieved and assessed according to
19 the criteria specified in the review protocol. A standardised form was used to extract
20 data from included studies. There were no instances that it was thought necessary to
21 contact study investigators for missing data.

22 **Incorporating published evidence syntheses**

23 If published evidence syntheses were identified sufficiently early in the review
24 process (for example, from the surveillance review or early in the database search),
25 they were considered for use as the primary source of data, rather than extracting
26 information from primary studies. Syntheses considered for inclusion in this way were
27 quality assessed to assess their suitability using the appropriate checklist, as outlined
28 in [Table 1](#). Note that this quality assessment was solely used to assess the quality of

29 the synthesis in order to decide whether it could be used as a source of data, as
30 outlined in [Table 2](#), not the quality of evidence contained within it, which was
31 assessed in the usual way as outlined in the section on ‘Appraising the quality of
32 evidence’.

33 **Table 1: Checklists for published evidence syntheses**

Type of synthesis	Checklist for quality appraisal
Systematic review of quantitative evidence	ROBIS

34

35 Each published evidence synthesis was classified into one of the following three
36 groups:

- 37 • High quality – It is unlikely that additional relevant and important data would be
38 identified from primary studies compared to that reported in the review, and
39 unlikely that any relevant and important studies have been missed by the review.
- 40 • Moderate quality – It is possible that additional relevant and important data would
41 be identified from primary studies compared to that reported in the review, but
42 unlikely that any relevant and important studies have been missed by the review.
- 43 • Low quality – It is possible that relevant and important studies have been missed
44 by the review.

45 Each published evidence synthesis was also classified into one of three groups for its
46 applicability as a source of data, based on how closely the review matches the
47 specified review protocol in the guideline. Studies were rated as follows:

- 48 • Fully applicable – The identified review fully covers the review protocol in the
49 guideline.
- 50 • Partially applicable – The identified review fully covers a discrete subsection of the
51 review protocol in the guideline (for example, some of the factors in the protocol
52 only).
- 53 • Not applicable – The identified review, despite including studies relevant to the
54 review question, does not fully cover any discrete subsection of the review
55 protocol in the guideline.

56 The way that a published evidence synthesis was used in the evidence review
 57 depended on its quality and applicability, as defined in [Table 2](#). When published
 58 evidence syntheses were used as a source of primary data, data from these were
 59 quality assessed and presented in GRADE/CERQual tables in the same way as if
 60 data had been extracted from primary studies. In questions where data was extracted
 61 from both systematic reviews and primary studies, these were checked to ensure
 62 none of the data had been double counted through this process.

63 **Table 2: 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 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.

64

1 **Methods of combining evidence**

2 **Data synthesis for intervention studies**

3 Where possible, meta-analyses were conducted to combine the results of
4 quantitative studies for each outcome. When there were 2 treatment alternatives,
5 pairwise meta-analysis was used to compare interventions.

6 **Pairwise meta-analysis**

7 Pairwise meta-analyses were performed in Cochrane RevMan (web version). A
8 pooled relative risk was calculated for dichotomous outcomes (using the Mantel–
9 Haenszel method) reporting numbers of people having an event, and a pooled
10 incidence rate ratio was calculated for dichotomous outcomes reporting total
11 numbers of events. Both relative and absolute risks were presented, with absolute
12 risks calculated by applying the relative risk to the risk in the comparator arm of the
13 meta-analysis (calculated as the total number events in the comparator arms of
14 studies in the meta-analysis divided by the total number of participants in the
15 comparator arms of studies in the meta-analysis). If a study reported only the
16 summary statistic and 95% CI the generic-inverse variance method was used to
17 enter data into the web version of Cochrane Review Manager.

18 No evidence presented as continuous outcomes was identified as part of this update.

19 Random effects models were fitted when significant between-study heterogeneity in
20 methodology, population, intervention or comparator was identified by the reviewer in
21 advance of data analysis. This decision was made and recorded before any data
22 analysis was undertaken. For all other syntheses, fixed- and random-effects models
23 were fitted, with the presented analysis dependent on the degree of heterogeneity in
24 the assembled evidence. Fixed-effects models were the preferred choice to report,
25 but in situations where the assumption of a shared mean for fixed-effects model were
26 clearly not met, even after appropriate pre-specified subgroup analyses were
27 conducted, random-effects results are presented. Fixed-effects models were deemed
28 to be inappropriate if there was significant statistical heterogeneity in the meta-
29 analysis, defined as I^2 more than or equal to 40%.

1 Results from subgroup analyses (rather than the full analysis) were assessed using
2 GRADE only when statistically significant subgroup differences were identified (p
3 <0.05). In these cases, where the full analysis displayed significant statistical
4 heterogeneity, but the pre-specified subgroup analyses were less heterogeneous
5 (with $I^2 <40\%$), the results from these subgroups were reported using fixed effects
6 models.

7 **Data synthesis for diagnostic accuracy data**

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 Where multiple observer interpretations of the same images were reported in a study,
17 all interpretations were included.

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

- 20 • **Sensitivity** is the probability that the feature will be positive in a person with the
21 condition.
 - 22 ○ $\text{sensitivity} = \text{TP}/(\text{TP}+\text{FN})$
- 23 • **Specificity** is the probability that the feature will be negative in a person without
24 the condition.
 - 25 ○ $\text{specificity} = \text{TN}/(\text{FP}+\text{TN})$
- 26 • **Positive likelihood ratios** describe how many times more likely positive features
27 are in people with the condition compared to people without the condition. Values
28 greater than 1 indicate that a positive result makes the condition more likely.
 - 29 ○ $\text{LR}^+ = (\text{TP}/[\text{TP}+\text{FN}])/(\text{FP}/[\text{FP}+\text{TN}])$

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 Where studies did not provide 2x2 data, a Microsoft excel calculator was used to
6 convert diagnostic accuracy measures (sensitivity, specificity, diagnostic accuracy,
7 predictive values or likelihood ratios) to the raw 2x2 data.

8 Meta-analysis of diagnostic accuracy data was conducted with reference to the
9 Cochrane Handbook for Systematic Reviews of Diagnostic Test Accuracy Version
10 2.0 (Deeks et al. 2023).

11 Where three or more studies were available for all included strata, a bivariate model
12 was fitted using the using the metaDTA app (<https://crsu.shinyapps.io/MetaDTA/>),
13 which accounts for the correlations between positive and negative likelihood ratios,
14 and between sensitivities and specificities (alternative approaches were required
15 where convergence wasn't achieved).

16 MetaDTA is based on the glmer function in the lme-4 package in R which runs a
17 bivariate meta-analysis using a generalised linear mixed model (GLMM), as
18 recommended in the Cochrane Handbook for Systematic Reviews of Diagnostic Test
19 Accuracy.

20 Where an analysis only contained 2 studies, regardless of whether the data
21 converged in MetaDTA, meta-analysis was not completed because a minimum of 3
22 studies are needed to estimate the parameters needed for bivariate meta-analysis.
23 Sensitivity and specificity forest plots were obtained from Cochrane RevMan (version
24 8.16.0), and likelihood ratios were obtained using a Microsoft Excel calculator. The
25 range of point estimates were presented in the GRADE table as a pooled effect
26 estimate could not be obtained.

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 Version 2.0 (Deeks et al. 2023).

1 **Appraising the quality of evidence**

2 **Intervention studies (relative effect estimates)**

3 RCTs were quality assessed using the Cochrane Risk of Bias Tool. Risk of bias for
4 single studies were conducted once for objective outcomes, once for subjective
5 outcomes, and once for adverse events.

6 Where systematic reviews were used as a source of evidence for RCTs but they use
7 the Cochrane Risk of Bias Tool 1 for risk of bias, the judgements were taken from
8 that review and converted to Cochrane risk of bias Tool 2 judgements so that all
9 RCTs were assessed in the same way. Descriptions of the approach taken are
10 written in the methods specific to the review.

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

- 15 • Direct – No important deviations from the protocol in population, intervention,
16 comparator and/or outcomes.
- 17 • Partially indirect – Important deviations from the protocol in one of the following
18 areas: population, intervention, comparator and/or outcomes.
- 19 • Indirect – Important deviations from the protocol in at least two of the following
20 areas: population, intervention, comparator and/or outcomes.

21 **Minimally important differences (MIDs) and clinical decision thresholds**

22 The Core Outcome Measures in Effectiveness Trials (COMET) database was
23 searched to identify published minimal clinically important difference thresholds
24 relevant to this guideline that might aid the committee in identifying clinical decision
25 thresholds for the purpose of GRADE. Identified MIDs were assessed to ensure they
26 had been developed and validated in a methodologically rigorous way, and were
27 applicable to the populations, interventions and outcomes specified in this guideline.
28 In addition, the Guideline Committee were asked to prospectively specify any
29 outcomes where they felt a consensus clinical decision threshold could be defined

1 from their experience. In particular, any questions looking to evaluate non-inferiority
 2 (that one treatment is not meaningfully worse than another) required a clinical
 3 decision threshold to be defined to act as a non-inferiority margin.

4 Clinical decision thresholds were used to assess imprecision using GRADE and aid
 5 interpretation of the size of effects for different outcomes. Clinical decision thresholds
 6 that were used in the guideline are given in [Table 3](#).

7 **Table 3: Identified clinical decision thresholds**

Outcome	Clinical decision threshold	Source
Quality of life FACT-G total	3 to 7 points	Eton DT, Cella D, Yost KJ, Yount SE, Peterman AH, Neuberg DS, Sledge GW, Wood WC. A combination of distribution- and anchor-based approaches determined minimally important differences (MIDs) for four endpoints in a breast cancer scale. J Clin Epidemiol. 2004 Sep;57(9):898-910. doi: 10.1016/j.jclinepi.2004.01.012. PMID: 15504633.
Quality of life FACT-B total	7 to 8 points	Eton DT, Cella D, Yost KJ, Yount SE, Peterman AH, Neuberg DS, Sledge GW, Wood WC. A combination of distribution- and anchor-based approaches determined minimally important differences (MIDs) for four endpoints in a breast cancer scale. J Clin Epidemiol. 2004 Sep;57(9):898-910. doi: 10.1016/j.jclinepi.2004.01.012. PMID: 15504633.
Quality of life TOI (trial outcome index) of FACT-B	5 to 6 points	Eton DT, Cella D, Yost KJ, Yount SE, Peterman AH, Neuberg DS, Sledge GW, Wood WC. A combination of distribution- and anchor-based approaches determined minimally important differences (MIDs) for four endpoints in a breast cancer scale. J Clin Epidemiol. 2004 Sep;57(9):898-910. doi: 10.1016/j.jclinepi.2004.01.012. PMID: 15504633.
Quality of life BCS of FACT-B	2 to 3 points	Eton DT, Cella D, Yost KJ, Yount SE, Peterman AH, Neuberg DS, Sledge GW, Wood WC. A combination of distribution- and anchor-based approaches determined minimally important differences (MIDs) for four endpoints in a breast cancer scale. J Clin Epidemiol. 2004 Sep;57(9):898-910. doi: 10.1016/j.jclinepi.2004.01.012. PMID: 15504633.
Quality of life EORTC QLQ-C30	11 points for improvement Minus 8 points for deterioration	Kawahara, T., Taira, N., Shirowa, T. et al. Minimal important differences of EORTC QLQ-C30 for metastatic breast cancer patients: Results from a randomized clinical trial. Qual Life Res 31, 1829–1836 (2022). https://doi.org/10.1007/s11136-021-03074-y
Quality of life	1 point	Den Oudsten, B.L., Zijlstra, W.P. & De Vries, J. The minimal clinical important difference in the World Health

Advanced breast cancer methods DRAFT FOR CONSULTATION (March 2026)

WHOQOL-100	Organization Quality of Life instrument—100 . Support Care Cancer 21, 1295–1301 (2013). https://doi.org/10.1007/s00520-012-1664-8
------------	--

1

2 **GRADE for intervention studies analysed using pairwise analysis**

3 GRADE was used to assess the quality of evidence for the outcomes specified in the
 4 review protocol. Data from randomised controlled trials was initially rated as high
 5 quality. The quality of the evidence for each outcome was downgraded or not from
 6 this initial point, based on the criteria given in [Table 4](#). These criteria were used to
 7 apply preliminary ratings, but were overridden in cases where, in the view of the
 8 analyst or committee the uncertainty identified was unlikely to have a meaningful
 9 impact on decision making.

10 **Table 4: Rationale for downgrading quality of evidence for intervention studies**

GRADE criteria	Reasons for downgrading quality
Risk of bias	<ul style="list-style-type: none"> • Not serious (don't downgrade): less than 50% overall weighting some concerns/high risk of bias • Serious (downgrade 1 level): more than 50% some concerns/high risk of bias • Very serious (downgrade 2 levels): more than 50% high risk of bias.
Indirectness	<ul style="list-style-type: none"> • Not serious (don't downgrade): less than 50% of overall weighting partially direct or indirect. • Serious (downgrade 1 level): more than 50% of overall weighting partially direct or indirect. • Very serious (downgrade 2 levels): more than 50% of overall weighting indirect
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 I² statistic.</p> <ul style="list-style-type: none"> • Not serious (don't downgrade) I² = less than 40%; • Serious (downgrade 1 level) I² = 40-60%; • Very serious (downgrade 2 levels) I² = more than 60%.
Imprecision	<p>If an MID other than the line of no effect was defined for the outcome, the outcome was downgraded once if the 95% confidence interval for the effect size crossed one line of the MID, and twice if it crossed both lines of the MID.</p> <p>If the line of no effect was defined as an MID for the outcome, it was downgraded once if the 95% confidence interval for the effect size crossed</p>

Advanced breast cancer methods DRAFT FOR CONSULTATION (March 2026)

	<p>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.</p> <p>Outcomes meeting the criteria for downgrading above were not downgraded if the confidence interval was sufficiently narrow that the upper and lower bounds would correspond to clinically equivalent scenarios.</p>
Publication bias	<p>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. When a funnel plot showed convincing evidence of publication bias, or the review team became aware of other evidence of publication bias (for example, evidence of unpublished trials where there was evidence that the effect estimate differed in published and unpublished data), the outcome was downgraded once. If no evidence of publication bias was found for any outcomes in a review (as was often the case), this domain was excluded from GRADE profiles to improve readability.</p>

1 **Diagnostic accuracy studies**

2 Individual diagnostic accuracy studies were quality assessed using the QUADAS-2
3 tool. Each individual study was classified into one of the following three groups:

- 4 • Low risk of bias – The true effect size for the study is likely to be close to the
5 estimated effect size.
- 6 • Moderate risk of bias – There is a possibility the true effect size for the study is
7 substantially different to the estimated effect size.
- 8 • High risk of bias – It is likely the true effect size for the study is substantially
9 different to the estimated effect size.

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

- 14 • Direct – No important deviations from the protocol in population, index feature
15 and/or reference standard.
- 16 • Partially indirect – Important deviations from the protocol in one of the population,
17 index feature and/or reference standard.
- 18 • Indirect – Important deviations from the protocol in at least two of the population,
19 index feature and/or reference standard.

- 1 • Where a meta-analysis could not be conducted because only 2 studies were
2 included in the analysis, individual studies reporting sensitivity and specificity had
3 equal weighting. Where $\geq 50\%$ of the studies had some concerns for risk of bias,
4 the evidence was downgraded by one level and where $\geq 50\%$ of the studies had
5 high risk of bias the evidence was downgraded by two levels.

6 **GRADE for diagnostic accuracy evidence**

7 Evidence from diagnostic accuracy studies was initially rated as high quality and then
8 downgraded according to the standard GRADE criteria (risk of bias, inconsistency,
9 imprecision and indirectness) as detailed in [Table 5](#) below.

10 The choice of primary outcome for decision making was determined by the
11 committee and GRADE assessments were undertaken based on these outcomes.

12 In all cases, the downstream effects of diagnostic accuracy on patient-important
13 outcomes were considered. This was done explicitly during committee deliberations
14 and reported as part of the discussion section of the review detailing the likely
15 consequences of true positive, true negative, false positive and false negative test
16 results. In reviews where a decision model is being carried (for example, as part of
17 an economic analysis), these consequences were incorporated here in addition.

18 **Using sensitivity and specificity as the primary outcomes**

19 GRADE assessments were only undertaken for sensitivity and specificity but results
20 for positive and negative likelihood ratios are also presented alongside those data.

21 The committee were consulted to set 2 clinical decision thresholds for each measure:
22 the value above which a test would be recommended, and a second below which a
23 test would be considered of no clinical use. These values were used to judge
24 imprecision (see below).

25 If studies could not be pooled in a meta-analysis, GRADE assessments were
26 undertaken on the body of evidence that was considered for meta-analysis, rather
27 than on individual estimates or median estimates.

28 Sensitivity and specificity were rated as high, moderate or low based on the
29 following:

Advanced breast cancer methods DRAFT FOR CONSULTATION (March 2026)

- 1 • High: Point estimate is greater than or equal to the upper clinical decision-making
- 2 threshold. Upper clinical decision-making thresholds were 90% for sensitivity and
- 3 80% for specificity
- 4 • Moderate: Point estimate greater than or equal to the lower clinical decision-
- 5 making threshold but lower than the upper clinical decision-making threshold.
- 6 • Low: Point estimate is less than the lower clinical decision-making threshold.
- 7 Lower clinical decision-making thresholds were 70% for sensitivity and 60% for
- 8 specificity.
- 9 • Sensitivity and specificity were interpreted using the percentage point estimate.
- 10 For sensitivity, the value of the point estimate determines the proportion that is
- 11 ruled out and for specificity, the value of the point estimate determines the
- 12 proportion that are ruled in as having the condition.

13 These criteria were used to apply preliminary ratings, but were overridden in cases
 14 where, in the view of the analyst or committee the uncertainty identified was unlikely
 15 to have a meaningful impact on decision making.

16 **Table 5: Rationale for downgrading quality of evidence for diagnostic accuracy**
 17 **data**

GRADE criteria	Reasons for downgrading quality
Risk of bias	<ul style="list-style-type: none"> • Not serious (don't downgrade): less than 50% overall weighting some concerns/high risk of bias • Serious (downgrade 1 level): more than 50% some concerns/high risk of bias • Very serious (downgrade 2 levels): more than 50% high risk of bias. Weight of the study was equal to the study sample as a proportion of the total number of participants contributing to the outcome.
Indirectness	<ul style="list-style-type: none"> • Not serious (don't downgrade): less than 50% of overall weighting partially direct or indirect. • Serious (downgrade 1 level): more than 50% of overall weighting partially direct or indirect. • Very serious (downgrade 2 levels): more than 50% of overall weighting indirect Weight of the study was equal to the study sample as a proportion of the total number of participants contributing to the outcome.
Inconsistency	Concerns about inconsistency of effects across studies, occurring when there is unexplained variability in the treatment effect demonstrated across

GRADE criteria	Reasons for downgrading quality
	<p>studies (heterogeneity), after appropriate pre-specified subgroup analyses have been conducted.</p> <p>Ranges of test performance are determined by the 2 clinical decision thresholds set for each measure (sensitivity and specificity).</p> <ul style="list-style-type: none"> • Not serious (don't downgrade): $\geq 80\%$ of point estimates are within the same range of test performance. • Serious (downgrade 1 level): $< 80\%$ of point estimates are within the same range of test performance, and the point estimates are separated by 1 decision making threshold. • Very serious (downgrade 2 levels): $< 80\%$ of point estimates are within the same range of test performance, and the point estimates are separated by 2 decision making thresholds. <p>Where there are apparent differences in effect size due consideration was given to the appropriateness of pooling studies.</p>
Imprecision	<p>If the 95% confidence interval for the outcome crossed one of the clinical decision thresholds, the outcome was downgraded one level. If the 95% confidence interval spanned both thresholds, the outcome was downgraded twice.</p> <p>See the section on 'Using sensitivity and specificity as the primary outcome' for a description of how clinical decision thresholds were agreed.</p>
Publication bias	<p>If the review team became aware of evidence of publication bias (for example, evidence of unpublished trials where there was evidence that the effect estimate differed in published and unpublished data), the outcome was downgraded once.</p> <p>If no evidence of publication bias was found for any outcomes in a review (as was often the case), this domain was excluded from GRADE profiles to improve readability.</p> <p>Funnel plot analyses were not carried out as this method is not recommended by the Cochrane Collaboration (Macaskill et al., 2023) to detect publication bias for DTA studies. Macaskill et al. (2023) explained that statistical tests to detect funnel plot asymmetry are designed primarily for use on reviews of randomised trials and that applying such tests in systematic reviews of test accuracy is likely to result in publication bias being incorrectly indicated by the test far too often. They also recommended exploring heterogeneity in test accuracy as participant and study characteristics may be associated with study size as well as test accuracy. Heterogeneity have been considered already as part of the GRADE assessment within this review.</p>

1

1 **Reviewing economic evidence**

2 The committee is required to make decisions based on the best available evidence of
3 both clinical effectiveness and cost effectiveness. Guideline recommendations should
4 be based on the expected costs of the different options in relation to their expected
5 health benefits (that is, their 'cost effectiveness') rather than the total implementation
6 cost. However, the committee will also need to be increasingly confident in the cost
7 effectiveness of a recommendation as the cost of implementation increases.

8 Therefore, the committee may require more robust evidence on the effectiveness and
9 cost effectiveness of any recommendations that are expected to have a substantial
10 impact on resources; any uncertainties must be offset by a compelling argument in
11 favour of the recommendation. The cost impact or savings potential of a
12 recommendation should not be the sole reason for the committee's decision.³

13 Health economic evidence was sought relating to the key clinical issues being
14 addressed in the guideline. Health economists:

- 15 • Undertook a systematic review of the published economic literature.
- 16 • Undertook new cost-effectiveness analysis in priority areas.

17 **Literature review**

18 The health economists:

- 19 • Identified potentially relevant studies for each review question from the health
20 economic search results by reviewing titles and abstracts. Full papers were then
21 obtained.
- 22 • Reviewed full papers against prespecified inclusion and exclusion criteria to
23 identify relevant studies (see below for details).
- 24 • Critically appraised relevant studies using economic evaluations checklists as
25 specified in the NICE guidelines manual.
- 26 • Extracted key information about the studies' methods and results into health
27 economic evidence tables (which can be found in appendices to the relevant
28 evidence reports).

- 1 • Generated summaries of the evidence in NICE health economic evidence profile
2 tables (included in the relevant evidence report for each review question) – see
3 below for details.

4 **Inclusion and exclusion criteria**

5 Full economic evaluations (studies comparing costs and health consequences of
6 alternative courses of action: cost–utility, cost-effectiveness, cost–benefit and cost–
7 consequences analyses) and comparative costing studies that addressed the review
8 question in the relevant population were considered potentially includable as health
9 economic evidence.

10 Studies that only reported cost per hospital (not per patient), or only reported average
11 cost effectiveness without disaggregated costs and effects were excluded. Literature
12 reviews, abstracts, posters, letters, editorials, comment articles, unpublished studies
13 and studies not in English were excluded. Studies published before 2010 and studies
14 from non-OECD countries or the USA were also excluded, on the basis that the
15 applicability of such studies to the present UK NHS context is likely to be too low for
16 them to be helpful for decision-making.

17 Remaining health economic studies were prioritised for inclusion based on their
18 relative applicability to the development of this guideline and the study limitations. For
19 example, if a high quality, directly applicable UK analysis was available, then other
20 less relevant studies may not have been included. Where exclusions occurred on this
21 basis, this is noted in the relevant evidence report.

22 For more details about the assessment of applicability and methodological quality
23 see [Table 6](#) below and the economic evaluation checklist ([appendix H of the NICE
24 guidelines manual 2014](#)) and the health economics review protocol, which can be
25 found in each of the evidence reports.

26 When no relevant health economic studies were found from the economic literature
27 review, relevant UK NHS unit costs related to the compared interventions were
28 presented to the committee to inform the possible economic implications of the
29 recommendations.

1 **NICE health economic evidence profiles**

2 NICE health economic evidence profile tables were used to summarise cost and
 3 cost-effectiveness estimates for the included health economic studies in each
 4 evidence review report. The health economic evidence profile shows an assessment
 5 of applicability and methodological quality for each economic study, with footnotes
 6 indicating the reasons for the assessment. These assessments were made by the
 7 health economist using the economic evaluation checklist from the NICE guidelines
 8 manual. It also shows the incremental costs, incremental effects (for example,
 9 quality-adjusted life years [QALYs]) and incremental cost-effectiveness ratio (ICER)
 10 for the base case analysis in the study, as well as information about the assessment
 11 of uncertainty in the analysis. See Table 6 for more details.

12 When a non-UK study was included in the profile, the results were converted into
 13 pounds sterling using the appropriate purchasing power parity.

14 **Table 6: Content of NICE health economic evidence profile**

Item	Description
Study	Surname of first author, date of study publication and country perspective with a reference to full information on the study.
Applicability	An assessment of applicability of the study to this guideline, the current NHS situation and NICE decision-making: ^(a) <ul style="list-style-type: none"> • Directly applicable – the study meets all applicability criteria, or fails to meet 1 or more applicability criteria but this is unlikely to change the conclusions about cost effectiveness. • Partially applicable – the study fails to meet 1 or more applicability criteria, and this could change the conclusions about cost effectiveness. • Not applicable – the study fails to meet 1 or more of the applicability criteria, and this is likely to change the conclusions about cost effectiveness. Such studies would usually be excluded from the review.
Limitations	An assessment of methodological quality of the study: ^(a) <ul style="list-style-type: none"> • Minor limitations – the study meets all quality criteria, or fails to meet 1 or more quality criteria, but this is unlikely to change the conclusions about cost effectiveness. • Potentially serious limitations – the study fails to meet 1 or more quality criteria, and this could change the conclusions about cost effectiveness. • Very serious limitations – the study fails to meet 1 or more quality criteria, and this is highly likely to change the conclusions about cost effectiveness. Such studies would usually be excluded from the review.

Item	Description
Other comments	Information about the design of the study and particular issues that should be considered when interpreting it.
Incremental cost	The mean cost associated with one strategy minus the mean cost of a comparator strategy.
Incremental effects	The mean QALYs (or other selected measure of health outcome) associated with one strategy minus the mean QALYs of a comparator strategy.
Cost effectiveness	Incremental cost-effectiveness ratio (ICER): the incremental cost divided by the incremental effects (usually in £ per QALY gained).
Uncertainty	A summary of the extent of uncertainty about the ICER reflecting the results of deterministic or probabilistic sensitivity analyses, or stochastic analyses of trial data, as appropriate.

1 (a) *Applicability and limitations were assessed using the economic evaluation checklist in appendix H*
2 *of the NICE guidelines manual³*

3 **Undertaking new health economic analysis**

4 As well as reviewing the published health economic literature for each review
5 question, as described above, new health economic analysis was undertaken by the
6 health economist in all areas.

7 The following general principles were adhered to in developing the cost-effectiveness
8 analyses:

- 9 • Methods were consistent with the NICE reference case for interventions with
10 health outcomes in NHS settings.
- 11 • The committee was involved in the design of the model, selection of inputs and
12 interpretation of the results.
- 13 • Model inputs were based on the systematic review of the clinical literature
14 supplemented with other published data sources where possible.
- 15 • When published data were not available committee expert opinion was used to
16 populate the model.
- 17 • Model inputs and assumptions were reported fully and transparently.
- 18 • The results were subject to sensitivity analysis and limitations were discussed.
- 19 • The model was peer-reviewed by another health economist at NICE.

20 Full methods and results of the cost-effectiveness analysis are described in a
21 separate economic analysis report.

1 **Cost-effectiveness criteria**

2 NICE sets out the principles that committees should consider when judging whether
3 an intervention offers good value for money. In general, an intervention was
4 considered to be cost effective (given that the estimate was considered plausible) if
5 either of the following criteria applied:

- 6 • the intervention dominated other relevant strategies (that is, it was both less costly
7 in terms of resource use and more clinically effective compared with all the other
8 relevant alternative strategies), or
- 9 • the intervention cost less than £20,000 per QALY gained compared with the next
10 best strategy.

11 If the committee recommended an intervention that was estimated to cost more than
12 £20,000 per QALY gained, or did not recommend one that was estimated to cost less
13 than £20,000 per QALY gained, the reasons for this decision are discussed explicitly
14 in 'The committee's discussion of the evidence' section of the relevant evidence
15 report, with reference to issues regarding the plausibility of the estimate or to factors
16 set out in NICE methods manuals.

17 **In the absence of health economic evidence**

18 When making recommendations in areas not in the scope of the health economic
19 analysis and where no relevant published evidence was identified the committee
20 made a qualitative judgement about cost effectiveness by considering expected
21 differences in resource use between options and relevant UK NHS unit costs,
22 alongside the results of the review of clinical effectiveness evidence.

23 The UK NHS costs reported in the guideline are those that were presented to the
24 committee and were correct at the time recommendations were drafted. They may
25 have changed subsequently before the time of publication. However, we have no
26 reason to believe they have changed substantially.

1 **References**

- 2 Deeks JJ, Bossuyt PM, Leeflang MM, Takwoingi Y (editors). Cochrane Handbook for
3 Systematic Reviews of Diagnostic Test Accuracy. Version 2.0 (updated July 2023).
4 Cochrane, 2023. Available from [https://training.cochrane.org/handbook-diagnostic-](https://training.cochrane.org/handbook-diagnostic-test-accuracy/current)
5 [test-accuracy/current](https://training.cochrane.org/handbook-diagnostic-test-accuracy/current)
- 6 Follmann D, Elliott P, Suh I, Cutler J (1992) Variance imputation for overviews of
7 clinical trials with continuous response. *Journal of Clinical Epidemiology* 45:769–73
- 8 Fu R, Vandermeer BW, Shamliyan TA, et al. (2013) Handling Continuous Outcomes
9 in Quantitative Synthesis In: *Methods Guide for Effectiveness and Comparative*
10 *Effectiveness Reviews* [Internet]. Rockville (MD): Agency for Healthcare Research
11 and Quality (US); 2008-. Available from:
12 <http://www.ncbi.nlm.nih.gov/books/NBK154408/>
- 13 Norman G., Sloan JA., Wyrwich KW. (2003) Interpretation of changes in health-
14 related quality of life: the remarkable universality of half a standard deviation. *Med*
15 *Care* 41(5):582-92.