

Heavy menstrual bleeding (update)

NICE guideline: methods

NICE guideline 88

Methods

March 2018

Final

*Evidence reviews were developed by
National Guideline Alliance, hosted by the
Royal College of Obstetricians and
Gynaecologists*

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 All NICE guidance is subject to regular review and may be updated or withdrawn.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 2018. All rights reserved. Subject to [Notice of rights](#).

ISBN: 978-1-4731-2777-7

Contents

Development of the guideline	5
Remit.....	5
What this guideline covers.....	5
Groups that are covered.....	5
Clinical areas that are covered.....	5
What this guideline does not cover.....	5
Groups that are not covered.....	5
Clinical areas that are not covered.....	6
Methods	7
Developing the review questions and outcomes.....	7
Searching for evidence.....	8
Clinical search literature.....	8
Health economics search literature.....	9
Call for evidence.....	9
Reviewing clinical evidence.....	9
Systematic review process.....	9
Type of studies and inclusion/ exclusion criteria.....	10
Methods of combining evidence.....	10
Appraising the quality of evidence.....	14
Evidence statements.....	20
Reviewing economic evidence.....	21
Inclusion and exclusion of economic studies.....	21
Health economic modelling.....	21
Cost effectiveness criteria.....	22
Developing recommendations.....	22
Guideline recommendations.....	22
Research recommendations.....	22
Validation process.....	22
Updating the guideline.....	23
Funding.....	23
References	24

1 Development of the guideline

2 Remit

3 The National Institute for Health and Care Excellence (NICE) commissioned the
4 National Guideline Alliance (NGA) to produce the update for this guideline.

5 The remit for this guideline update is to revise the NICE clinical guideline on the
6 structural diagnosis of causes of heavy menstrual bleeding and the management of
7 heavy menstrual bleeding.

8 What this guideline covers

9 Groups that are covered

10 The guideline update covers women with heavy menstrual bleeding, including:

- 11 • women with suspected or confirmed fibroids
- 12 • women with suspected or confirmed adenomyosis
- 13 • women with no identified pathology.

14 Women who wish to preserve their fertility have been identified as a subgroup
15 needing specific consideration.

16 Clinical areas that are covered

17 The guideline update covers the following clinical issues:

- 18 • clinical and cost-effectiveness of hysteroscopy and pelvic ultrasound scan to
19 detect causes of heavy menstrual bleeding
- 20 • clinical and cost-effectiveness of diagnostic imaging techniques to detect
21 adenomyosis in women presenting with heavy menstrual bleeding
- 22 • clinical and cost-effectiveness of pharmacological and surgical management of
23 heavy menstrual bleeding.

24 Note that guideline recommendations will normally fall within licensed indications.
25 Exceptionally, and only if clearly supported by evidence, use outside a licensed
26 indication may be recommended. This guideline will assume that prescribers will use
27 a drug's summary of product characteristics to inform decisions made with individual
28 patients.

29 For further details please refer to the scope on the NICE website
30 (<https://www.nice.org.uk/guidance/gid-ng10012/documents/final-scope>).

31 What this guideline does not cover

32 Groups that are not covered

33 The guideline does not cover the following groups:

- 34 • women without heavy menstrual bleeding who have other gynaecological
35 bleeding, for example, intermenstrual bleeding or post-coital bleeding

- 1 • women with gynaecological conditions in which heavy menstrual bleeding is not
2 the main problem, for example, women with endometriosis.

3 **Clinical areas that are not covered**

4 The following areas in the published guideline were not updated:

- 5 • definition of heavy menstrual bleeding
6 • education and information provision
7 • competencies:
8 ○ training
9 ○ maintenance
10 ○ governance
11 • the clinical and cost-effectiveness of treatment with progesterone receptor
12 modulators for fibroids of 3 cm or more in diameter (this topic was reviewed by the
13 NICE standing committee, and an addendum to the NICE guideline on Heavy
14 menstrual bleeding [CG44] was published in August 2016).

15 Recommendations in areas that were not updated were edited to ensure that they
16 meet the current editorial standard, and reflect the current policy and practice
17 context.

18 Management of endometriosis associated with heavy menstrual bleeding is not
19 covered by this guideline but is covered in the [NICE guideline on Endometriosis:
20 diagnosis and management](#) published in September 2017.
21

1 Methods

2 This chapter sets out in detail the methods used to review the evidence and to
3 generate recommendations in the guideline. This guideline was developed using the
4 methods described in the [2014 NICE guidelines manual \(NICE 2014\)](#).

5 Declarations of interest were recorded according to the 2014 NICE conflicts of
6 interest policy.

7 Developing the review questions and outcomes

8 The 3 review questions developed for this guideline were based on the key areas
9 identified in the guideline update scope (see [https://www.nice.org.uk/guidance/gid-
10 ng10012/documents/final-scope](https://www.nice.org.uk/guidance/gid-ng10012/documents/final-scope)). They were drafted by the NGA and refined and
11 validated by the guideline committee (see Table 1).

12 The review questions were based on the following frameworks:

- 13 • intervention review: population, intervention, comparator and outcome (PICO);
- 14 • diagnostic test accuracy review: population, index tests, reference standard and
15 target condition.

16 These frameworks guided the development of the review protocols, the literature
17 searching process, the critical appraisal and synthesis of evidence and facilitated the
18 development of recommendations by the guideline committee.

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

21 **Table 1: Description of review questions**

Chapter or section	Type of review	Review question	Outcomes
Evidence reviews for diagnostic test accuracy in investigation for women presenting with heavy menstrual bleeding	Diagnostic	What is the diagnostic accuracy of ultrasound and hysteroscopy for investigation of women presenting with heavy menstrual bleeding?	<ul style="list-style-type: none"> • Sensitivity • Specificity • Positive likelihood ratio (LR+) • Negative likelihood ratio (LR-) • Area under the curve (AUC) if meta-analysis can be conducted • Patient satisfaction and acceptability of the test, including pain score
Evidence reviews for diagnostic test accuracy in investigation for women presenting with heavy menstrual bleeding	Diagnostic	What is the most clinically effective imaging strategy for diagnosing adenomyosis in women with heavy menstrual bleeding?	<ul style="list-style-type: none"> • Sensitivity • Specificity • LR+ • LR-

Chapter or section	Type of review	Review question	Outcomes
			<ul style="list-style-type: none"> • AUC if meta-analysis can be conducted • Patient satisfaction and acceptability of the test, including pain score
Evidence reviews for management of heavy menstrual bleeding	Intervention	<p>What is the most clinically and cost-effective treatment (pharmacological/surgical) for heavy menstrual bleeding in women with:</p> <ul style="list-style-type: none"> • suspected or diagnosed fibroids • suspected or diagnosed adenomyosis • no identified pathology? 	<ul style="list-style-type: none"> • Reduction in blood loss • Quality of life • Patient satisfaction • Adverse events

1 AUC: area under the curve; LR+: positive likelihood ratio; LR-: negative likelihood ratio

2 Searching for evidence

3 Clinical search literature

4 Systematic literature searches were undertaken to identify all published clinical
5 evidence relevant to the review questions on 13th October 2016 (Diagnosis question)
6 and 23rd November 2016 (Management question).

7 Databases were searched using relevant medical subject headings, free-text terms
8 and study type filters where appropriate. Studies published in languages other than
9 English were not reviewed. Where possible, searches were restricted to retrieve only
10 articles published in English. All searches were conducted in MEDLINE, Embase and
11 The Cochrane Library.

12 Any studies added to the databases after this date (even those published prior to this
13 date) were not included unless specifically stated in the text.

14 Search strategies were quality assured by cross-checking reference lists of highly
15 relevant papers, analysing search strategies in other systematic reviews and asking
16 the group members to highlight any additional studies. The questions, the study
17 types applied, the databases searched and the years covered can be found in
18 Appendix E in each evidence review chapter.

19 Searching for grey literature or unpublished literature was not undertaken. Searches
20 for electronic, ahead-of-print publications were not routinely undertaken unless
21 indicated by the guideline committee. All references suggested by stakeholders at
22 the scoping consultation were initially considered.

1 Health economics search literature

2 A global search of economic evidence was undertaken in December 2016. The
3 following databases were searched:

- 4 • MEDLINE (Ovid)
- 5 • EMBASE (Ovid)
- 6 • Cochrane Central Register of Controlled Trials (CCTR)
- 7 • HTA database (HTA)
- 8 • NHS Economic Evaluations Database (NHS EED).

9 Further to the database searches, the committee was contacted with a request for
10 details of relevant published and unpublished studies of which they may have
11 knowledge; reference lists of key identified studies were also reviewed for any
12 potentially relevant studies. Finally, the NICE website was searched for any recently
13 published guidance relating to heavy menstrual bleeding that had not been already
14 identified via the database searches.

15 The search strategy for existing economic evaluations combined terms capturing the
16 target condition (heavy menstrual bleeding) and, for searches undertaken in
17 MEDLINE, EMBASE and CCTR, terms to capture economic evaluations. No
18 restrictions on language or setting were applied to any of the searches, but a
19 standard exclusions filter was applied (letters, animals, etc). Conference abstracts
20 were considered for inclusion from 1st January 2014, as high-quality studies reported
21 in abstract form before 2014 were expected to have been published in a peer-
22 reviewed journal. Full details of the search strategies are presented in Appendix E of
23 each evidence review chapter.

24 Call for evidence

25 No call for evidence was made.

26 Reviewing clinical evidence

27 Systematic review process

28 The evidence was reviewed following these steps.

- 29 • Potentially relevant studies were identified for each review question from the
30 relevant search results by reviewing titles and abstracts. Full papers were then
31 obtained.
- 32 • Full papers were reviewed against pre-specified inclusion and exclusion criteria in
33 the review protocols (in Appendix A of each evidence review chapter).
- 34 • Key information was extracted on the study's methods, according to the factors
35 specified in the protocols and results. These were presented in summary tables (in
36 each review chapter) and evidence tables (in Appendix F of each evidence review
37 chapter).
- 38 • Relevant studies were critically appraised using the appropriate checklist as
39 specified in the [NICE guidelines manual \(NICE 2014\)](#).
- 40 • Summaries of evidence were generated by outcome (included in the relevant
41 review chapters) and were presented in committee meetings.

- 1 • Randomised studies: meta-analysis was carried out where appropriate and results
2 were reported in GRADE profiles (for intervention reviews).
- 3 • Diagnostic studies: data were presented individually by study as measures of
4 diagnostic test accuracy (sensitivity and specificity, positive and negative
5 likelihood ratios) and were presented in modified GRADE profiles.
- 6 To assure quality of the study identification, a 10% sample of all the titles and
7 abstracts for each review question were assessed for possible inclusion by a second
8 independent reviewer. Possible discrepancies were resolved by discussion between
9 the two reviewers.
- 10 All drafts of reviews were checked by a second reviewer. Any discrepancies were
11 resolved by discussion between the 2 reviewers.

12 **Type of studies and inclusion/ exclusion criteria**

13 Systematic reviews (SRs) with meta-analyses were considered the highest quality
14 evidence to be selected for inclusion.

15 For the review on the management of heavy menstrual bleeding, randomised
16 controlled trials (RCTs) were prioritised for inclusion because they are considered the
17 most robust study design to estimate the true effect of the interventions. RCTs with
18 less than 10 participants in any intervention arm were excluded.

19 For the diagnostic test accuracy reviews, studies were included in which the index
20 test and the reference standard were compared in the same individual and in which
21 2x2 tables could be constructed. The study designs considered for inclusion included
22 test and treat RCTs, cross-sectional studies, and prospective cohort studies. Case-
23 control studies were excluded.

24 Conference abstracts, posters, letters, editorials, comment articles, unpublished
25 studies and studies not in the English language were excluded. Narrative reviews
26 were also excluded, but individual references were checked for inclusion.

27 The inclusion and exclusion of studies was based on the review protocols, which can
28 be found in Appendix A of each evidence review chapter. Excluded studies and the
29 reasons for their exclusion are listed in Appendix I of each evidence review chapter.
30 In addition, the guideline committee was consulted about any uncertainty regarding
31 inclusion or exclusion.

32 **Methods of combining evidence**

33 **Data synthesis for intervention review**

34 ***Pairwise meta-analysis***

35 Pairwise meta-analysis was conducted whenever it could be robustly performed to
36 combine the results of studies using Review Manager 5 (RevMan 5) software.

37 For binary outcomes, such as occurrence of adverse events, the Mantel-Haenszel
38 method of statistical analysis was used to calculate risk ratios (relative risks, RRs)
39 with 95% confidence intervals (CIs).

40 For continuous outcomes, measures of central tendency (mean) and variation
41 (standard deviation) are required for meta-analysis. Data for continuous outcomes

1 (such as health-related quality of life score or length of hospital stay) were analysed
2 using an inverse variance method for pooling weighted mean differences. When the
3 only evidence was based on studies summarising results by presenting medians
4 (and interquartile ranges) or only p values were given, this information was assessed
5 in terms of the study's sample size, and was included in the GRADE tables without
6 calculating the relative or absolute effects. Consequently, aspects of quality
7 assessment, such as imprecision of effect, could not be assessed for evidence of this
8 type.

9 Forest plots were generated to visually present the results.

10 Statistical heterogeneity was assessed by visually examining the forest plots (please
11 see Appendix H of each evidence review chapter) and by considering the chi-
12 squared test for significance at $p < 0.1$ or an I^2 squared inconsistency statistic (with an
13 I^2 value of more than 50% indicating considerable heterogeneity). Where
14 considerable heterogeneity was present, predefined subgroup analyses were
15 performed.

16 **Network meta-analysis**

17 As is the case for ordinary pairwise meta-analysis, network meta-analysis (NMA)
18 may be conducted using either fixed or random effect models. A fixed effect model
19 typically assumes that there is no variation in relative effects across trials for a
20 particular pairwise comparison and any observed differences are solely due to
21 chance. For a random effects model, it is assumed that the relative effects are
22 different in each trial but that they are from a single common distribution. The
23 variance reflecting heterogeneity is often assumed to be constant across trials.

24 For continuous outcomes, where standard errors (SEs) could not be calculated from
25 the data, we imputed them from other studies that reported measures of
26 uncertainty/variance, using the median standard deviation (SD) of other study arms
27 in the analysis that used the same treatment.

28 In a Bayesian analysis, for each parameter the evidence distribution is weighted by a
29 distribution of prior beliefs. The Markov Chain Monte Carlo (MCMC) algorithm was
30 used to generate a sequence of samples from a joint posterior distribution of 2 or
31 more random variables and is particularly well adapted to sampling the treatment
32 effects (known as a posterior distribution) of a Bayesian network. A non-informative
33 prior distribution was used to maximise the weighting given to the data and to
34 generate the posterior distribution for each log odds ratio (OR), log mean ratio (MR)
35 or mean difference (MD) of interest in the networks. We used the median of the
36 distribution as our point estimate and the centiles provided the 95% Credible
37 Intervals (CrIs).

38 Non-informative priors were used that were normally distributed with a mean of 0 and
39 SD of 100. However, for discontinuation due to adverse events, as there was sparse
40 data on a number of treatments, we investigated whether the use of informative
41 priors generated from empirical data would give a more stable between-study
42 variance (Turner 2012).

43 For the analyses, a series of 40,000 burn-in simulations were run to allow the
44 posterior distributions to convergence and then a further 100,000 simulations were
45 run to produce the outputs. Convergence was assessed by examining the history,
46 autocorrelation and Brooks-Gelman-Rubin plots.

1 Goodness-of-fit of the model was also estimated by using the posterior mean of the
2 sum of the deviance contributions for each item by calculating the residual deviance
3 and deviance information criteria (DIC). If the residual deviance was close to the
4 number of unconstrained data points (the number of trial arms in the analysis) then
5 the model was explaining the data at a satisfactory level. The choice of a fixed effect
6 or random effects model can be made by comparing their goodness-of-fit to the data.

7 Incoherence in NMA between direct and indirect evidence can be assessed in closed
8 treatment loops within the network. These closed treatment loops are regions within
9 a network where direct evidence is available on at least 3 different treatments that
10 form a closed 'circuit' of treatment comparisons (for example, A versus B, B versus
11 C, C versus A). If closed treatment loops existed then discrepancies between direct
12 and indirect evidence was assessed for each loop using node-splitting (van
13 Valkenhoef 2016). The outputs of the NMA were as follows.

- 14 • Treatment specific log ORs, log MRs and MDs with their 95% CrIs were generated
15 for every possible pair of comparisons by combining direct and indirect evidence
16 in each network.
- 17 • The probability that each treatment is ranked within the best 3 or worst 3
18 treatments, based on the proportion of Markov chain iterations in which the
19 treatment effect for an intervention is ranked best, second best and so forth. This
20 was calculated by taking the treatment effect of each drug compared to placebo
21 and counting the proportion of simulations of the Markov chain in which each
22 intervention had the highest treatment effect.
- 23 • The ranking of treatments compared to the reference treatment (typically placebo
24 or levonorgestrel-releasing intrauterine system (LNG-IUS) presented as median
25 rank and its 95% CrI.

26 One of the main advantages of the Bayesian approach is that the method leads to a
27 decision framework that supports decision making. The Bayesian approach also
28 allows the probability that each intervention is best for achieving a particular
29 outcome, as well as its ranking, to be calculated.

30 We adapted a random effects model template for continuous and dichotomous data
31 available from NICE Decision Support Unit (DSU) technical support document
32 number 2: [http://www.nicedsu.org.uk/Evidence-Synthesis-TSD-series-\(2391675\).htm](http://www.nicedsu.org.uk/Evidence-Synthesis-TSD-series-(2391675).htm).
33 This model accounts for the within-study correlation between treatment effects
34 induced by multi-arm trials.

35 For further description of outcomes and the specific results of the NMA please see
36 the evidence review chapter for the management of heavy menstrual bleeding.

37 **Data synthesis for diagnostic test accuracy reviews**

38 ***Diagnostic data and outcomes***

39 Sensitivity, specificity, positive and negative likelihood ratios, and area under the
40 curve (AUC) were used as outcomes for diagnostic test accuracy reviews in this
41 guideline. These diagnostic accuracy parameters (with 95% CIs) were obtained from
42 the studies or calculated by the technical team using data from the studies (see
43 Table 2).

44 Sensitivity and specificity are measures of the ability of a test to correctly classify a
45 person as having a condition or not having a condition. When sensitivity is high, a

1 negative test result rules out the target condition. When specificity is high, a positive
2 test result rules in the target condition. An ideal test would be both highly sensitive
3 and highly specific, but this is frequently not possible and typically there is a trade-off.

4 The following cut-offs were used when summarising the levels of sensitivity or
5 specificity for the guideline committee:

- 6 • high: more than 90%
- 7 • moderate: 75% to 90%
- 8 • low: less than 75%.

9 Positive and negative likelihood ratios are measures of the association between a
10 test result and the target condition. A positive likelihood ratio (LR+) greater than 1
11 indicates a positive test result and is associated with having the disorder, whilst a
12 negative likelihood ratio (LR-) less than 1 indicates a negative test result and is
13 associated with not having the disorder. A high LR+ would indicate that the test is
14 useful in ruling in the condition whereas a low LR- would indicate that the test is
15 useful in ruling out the condition.

16 The following cut-offs were used when summarising the likelihood ratios for the
17 guideline committee:

- 18 • very useful test: LR+ higher than 10.0, LR- lower than 0.1
- 19 • moderately useful test: LR+ 5.0 to 10.0, LR- 0.1 to 0.2
- 20 • not a useful test: LR+ lower than 5.0, LR- higher than 0.2.

21 **Table 2: 2x2 table for calculating diagnostic test accuracy parameters**

	Condition present (according to reference standard)	No condition (according to reference standard)	Total
Index test positive	True positive (TP)	False positive (FP)	TP+FP = total number of subjects positive index test result
Index test negative	False negative (FN)	True negative (TN)	FN+TN = total number of subjects with negative index test result
Total	TP+FN = total number of subjects with condition	FP+TN = total number of subject without condition	TP+FP+FN+TN = Total number of subjects in study
Calculations for diagnostic test accuracy parameters:			
Sensitivity = TP/(TP+FN)		LR+ = sensitivity/(1-specificity)	
Specificity = TN/(TN+FP)		LR- = (1-sensitivity)/specificity	

22 *FN: false negative; FP: false positive; LR+: positive likelihood ratio; LR-: negative likelihood ratio; TP:*
23 *true negative; TP: true positive*

24 **Diagnostic meta-analysis**

25 When data from 5 or more studies were available, a diagnostic meta-analysis was
26 carried out by using statistical software STATA with metandi package (Harbord and
27 Whiting 2009; Harbord 2008). The metandi package performs bivariate meta-analysis
28 of sensitivity and specificity using a generalised linear mixed model approach.

1 Forest plots and hierarchical summary receiver operating characteristic (HSROC)
2 plots were created to visually present the results.

3 **Appraising the quality of evidence**

4 **Intervention reviews**

5 ***Pairwise analysis***

6 **GRADE methodology (The Grading of Recommendations Assessment, 7 Development and Evaluation)**

8 For intervention reviews, the evidence for outcomes from the included RCTs were
9 evaluated and presented using GRADE, which was developed by the international
10 GRADE working group.

11 The software developed by the GRADE working group (GRADEpro) was used to
12 assess the quality of each outcome, taking into account individual study quality
13 factors and the meta-analysis results. The clinical/economic evidence profile tables
14 include details of the quality assessment and pooled outcome data, where
15 appropriate, an absolute measure of intervention effect and the summary of quality of
16 evidence for that outcome. In this table, the columns for intervention and control
17 indicate summary measures of effect and measures of dispersion (such as mean and
18 SD or median and range) for continuous outcomes and frequency of events (n/N; the
19 sum across studies of the number of patients with events divided by sum of the
20 number of completers) for binary outcomes. Reporting or publication bias was taken
21 into consideration in the quality assessment and reported in the clinical evidence
22 profile tables if it was apparent.

23 The selection of outcomes for each review question was decided when each review
24 protocol was discussed with the guideline committee, and was informed by
25 committee discussion and by key papers, for example, previous NMAs. The
26 systematic review by Herman (2016) describing the outcomes used in published
27 systematic reviews and RCTs was also used to ensure all the main primary and
28 secondary outcomes reported in trials were considered.

29 The evidence for each outcome in the intervention reviews was examined separately
30 for the quality elements listed and defined in Table 3. Each element was graded
31 using the quality levels listed in Table 4.

32 The main criteria considered in the rating of these elements are discussed below.
33 Footnotes were used to describe reasons for grading a quality element as having
34 serious or very serious limitations. The ratings for each component were summed to
35 obtain an overall assessment for each outcome (Table 5).

36 **Table 3: Description of quality elements in GRADE for intervention reviews**

Quality element	Description
Risk of bias (study limitations)	Limitations in the study design and implementation may bias the estimates of the treatment effect. High risk of bias for the majority of the evidence decreases confidence in the estimate of the effect.
Inconsistency	Inconsistency refers to an unexplained heterogeneity of results or findings.

Quality element	Description
Indirectness	Indirectness refers to differences in study population, intervention, comparator and outcomes between the available evidence and the review question, or recommendation made, such that the effect estimate is changed. This is also related to applicability or generalisability of findings.
Imprecision	Results are imprecise when studies include relatively few patients and few events and thus have wide confidence intervals around the estimate of the effect. Imprecision results if the confidence interval includes the clinically important threshold.
Publication bias	Publication bias is a systematic underestimate or an overestimate of the underlying beneficial or harmful effect due to the selective publication of studies.

1 Table 4: Levels of quality elements in GRADE

Levels of quality elements in GRADE	Description
None/no serious	There are no serious issues with the evidence.
Serious	The issues are serious enough to downgrade the outcome evidence by 1 level.
Very serious	The issues are serious enough to downgrade the outcome evidence by 2 levels.

2 Table 5: Levels of overall quality of outcome evidence in GRADE

Overall quality of outcome evidence in GRADE	Description
High	Further research is very unlikely to change our confidence in the estimate of effect.
Moderate	Further research is likely to have an important impact on our confidence in the estimate of effect and may change the estimate.
Low	Further research is very likely to have an important impact on our confidence in the estimate of effect and is likely to change the estimate.
Very low	Any estimate of effect is very uncertain.

3 Assessing risk of bias in intervention reviews

- 4 Bias is a systematic error, or a consistent deviation from the truth in the results.
- 5 When a risk of bias is present the true effect can be either under- or overestimated.
- 6 Risk of bias in intervention studies was assessed using the Cochrane Risk of Bias
- 7 Tool ((see Appendix H in the [NICE guidelines manual 2014](#)).
- 8 The possible sources of bias in RCTs in the Cochrane risk of bias tool fit with these 5
- 9 categories:
- 10 • selection bias
- 11 • performance bias
- 12 • attrition bias
- 13 • detection bias

- 1 • reporting bias

2 It should be noted that a study with a poor methodological design does not
3 automatically imply high risk of bias; the bias is considered individually for each
4 outcome and it is assessed whether this poor design will impact on the estimation of
5 the intervention effect.

6 More details about this tool can be found here:

7 http://cobe.paginas.ufsc.br/files/2014/10/Cochrane.RCT_.pdf

8 **Assessing inconsistency in intervention reviews**

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

16 Heterogeneity was assessed by calculating the I^2 statistic for the meta-analysis. An I^2
17 of more than 50% was considered to indicate high heterogeneity. When high
18 heterogeneity was observed, possible reasons for it were explored and subgroup
19 analyses were performed as pre-specified in the review protocol.

20 When no plausible explanation for the heterogeneity could be found, the quality of
21 the evidence was downgraded in GRADE by 1 or 2 levels for the domain of
22 inconsistency, depending on the extent of heterogeneity in the results.

23 **Assessing indirectness in intervention reviews**

24 Directness refers to the extent to which the populations, intervention, comparisons
25 and outcome measures are similar to those defined in the inclusion criteria for the
26 reviews. Indirectness is important when these differences are expected to contribute
27 to a difference in effect size, or may affect the balance of harms and benefits
28 considered for an intervention.

29 **Assessing imprecision and clinical significance in intervention reviews**

30 Imprecision in guidelines concerns whether the uncertainty (CI) around the effect
31 estimate means that it is not clear whether there is a clinically important difference
32 between interventions or not (that is, whether the evidence would clearly support one
33 recommendation or appear to be consistent with several different types of
34 recommendations). Therefore, imprecision differs from the other aspects of evidence
35 quality because it is not really concerned with whether the point estimate is accurate
36 or correct (has internal or external validity). Instead, it is concerned with the
37 uncertainty about what the point estimate actually is. This uncertainty is reflected in
38 the width of the CI.

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

42 Imprecision in the evidence reviews was assessed by considering whether the width
43 of the 95% CI of the effect estimate was relevant to decision-making, taking each

1 outcome in isolation. This is explained in Figure 1, which considers a positive
 2 outcome for the comparison of treatment A versus treatment B. Three decision-
 3 making zones can be identified, bounded by the thresholds for clinical importance
 4 (minimally important difference, MID) for benefit and for harm. The MID for harm for a
 5 positive outcome means the threshold at which drug A is less effective than drug B
 6 by an amount that is clinically important to patients (favours B).

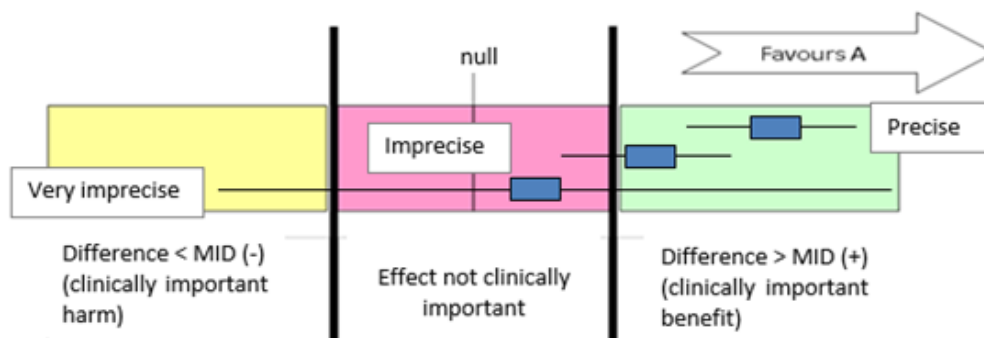
7 When the CI of the effect estimate is wholly contained in 1 of the 3 zones (for
 8 example, clinically important benefit), we are not uncertain about the size and
 9 direction of effect (whether there is a clinically important benefit, or the effect is not
 10 clinically important, or there is a clinically important harm), so there is no imprecision.

11 When a wide CI lies partly in each of 2 zones, it is uncertain in which zone the true
 12 value of effect estimate lies and therefore there is uncertainty over which decision to
 13 make (based on this outcome alone). The CI is consistent with 2 possible decisions
 14 and so this is considered to be imprecise in the GRADE analysis and the evidence is
 15 downgraded by 1 level ('serious imprecision').

16 If the CI of the effect estimate crosses into 3 zones, this is considered to be very
 17 imprecise evidence because the CI is consistent with 3 possible clinical decisions
 18 and there is therefore a considerable lack of confidence in the results. The evidence
 19 is therefore downgraded by 2 levels in the GRADE analysis ('very serious
 20 imprecision').

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

Figure 1: Illustration of precise, imprecise and very imprecise evidence based on the confidence interval of outcomes in forest plots



24 **Minimally important differences**

25 The literature was searched for established MIDs for the selected outcomes in the
 26 evidence reviews, such as blood loss or quality of life. In addition, the committee was
 27 asked whether they were aware of any acceptable MIDs in the clinical community.

28 If no published or acceptable MIDs were identified, the committee considered
 29 whether it was clinically acceptable to use the GRADE default MID to assess
 30 imprecision. For binary outcomes clinically important thresholds for a RR of 0.8 and
 31 1.25 respectively were used (due to the statistical distribution of this measure this
 32 means that this is not a symmetrical interval). This default MID was used for all the

1 binary outcomes in the intervention evidence reviews as a starting point and
 2 decisions on clinical importance were then considered based on the absolute risk
 3 difference. For continuous outcomes GRADE default MIDs were half of the median
 4 SD of the control group.

5 **Network meta-analysis**

6 For the NMAs, quality was assessed by looking at risk of bias across the included
 7 evidence (using the standard GRADE approach for this domain), as well as
 8 heterogeneity and incoherence.

9 The following limits of the upper 95% CrI for between-study standard deviation were
 10 used to assess heterogeneity for NMAs in which a random effects model was used:

- 11 • less than 0.3 – low heterogeneity
- 12 • 0.3 to 0.6 – moderate heterogeneity
- 13 • more than 0.6 to 0.9 – high heterogeneity
- 14 • more than 0.9 to 1.2 – very high heterogeneity.

15 Where significant incoherence was found it was considered to be serious when the
 16 direction of effect for both direct and indirect estimates was the same (for example,
 17 an OR of greater than 1 in both the direct and indirect estimates), and very serious
 18 when the direction of effect was different (for example, an OR of greater than 1 for
 19 the direct estimate but less than 1 for the indirect estimate).

20 For fixed-effect NMAs that did not model heterogeneity, or for networks in which
 21 incoherence could not be assessed as no closed treatment loops existed, these
 22 criteria were not considered to impact the quality of evidence.

23 **Diagnostic reviews**

24 **Adapted GRADE methodology**

25 The GRADE toolbox is designed for RCTs and observational studies, but we adapted
 26 the quality assessment elements and outcome presentation for diagnostic test
 27 accuracy reviews. For example, the GRADE clinical evidence tables were modified to
 28 include the most appropriate measures of diagnostic accuracy (sensitivity, specificity,
 29 and likelihood ratios).

30 The evidence for each outcome in the diagnostic test accuracy reviews was
 31 examined separately for the quality elements listed and defined in Table 6. Each
 32 element was graded using the quality levels listed in Table 4.

33 The main criteria considered in the rating of these elements are discussed below.
 34 Footnotes were used to describe reasons for grading a quality element as having
 35 serious or very serious limitations. The ratings for each component were summed to
 36 obtain an overall assessment for each outcome (Table 5).

37 **Table 6: Description of the elements in GRADE and how they are used to**
 38 **assess the quality for diagnostic accuracy reviews**

Quality element	Description
Risk of bias ('Study limitations')	Limitations in the study design and implementation may bias the estimates of the diagnostic accuracy. High risk of bias for the majority of the evidence decreases confidence in the estimate of

Quality element	Description
	the effect. Diagnostic accuracy studies are not usually randomised and therefore would not be downgraded for study design from the outset and start as high level evidence.
Inconsistency	Inconsistency refers to an unexplained heterogeneity of test accuracy measures, for example sensitivity or specificity, between studies.
Indirectness	Indirectness refers to differences in study population, index tests, reference standards and outcomes between the available evidence and the review question.
Imprecision	Results are considered imprecise when studies include relatively few patients and the confidence intervals were wide. Imprecision results if the confidence interval includes the clinically important threshold.

1 **Assessing risk of bias and indirectness in diagnostic test accuracy reviews**

2 Risk of bias in diagnostic test accuracy studies was assessed using the Quality
3 Assessment of Diagnostic Accuracy Studies version 2 (QUADAS-2) checklist (see
4 Appendix H in the [NICE guidelines manual 2014](#)).

5 Risk of bias and applicability in primary diagnostic accuracy studies in QUADAS-2
6 consists of 4 domains:

- 7 • patient selection
- 8 • index test
- 9 • reference standard
- 10 • flow and timing.

11 More details about this tool can be found here: [http://www.bristol.ac.uk/social-](http://www.bristol.ac.uk/social-community-medicine/projects/quadas/quadas-2/)
12 [community-medicine/projects/quadas/quadas-2/](http://www.bristol.ac.uk/social-community-medicine/projects/quadas/quadas-2/)

13 **Assessing inconsistency in diagnostic test accuracy reviews**

14 Inconsistency refers to the unexplained heterogeneity of the results in meta-analysis.
15 When estimates of diagnostic accuracy parameters vary widely across studies (that
16 is, there is heterogeneity or variability in results), this suggests true differences in
17 underlying effects. Inconsistency is, thus, only applicable when statistical meta-
18 analysis was conducted (that is, results from different studies were pooled).
19 However, 'no inconsistency' is nevertheless used to describe this quality assessment
20 in the GRADE profiles for outcomes from single studies.

21 For the diagnostic test accuracy reviews, the heterogeneity of the pooled result was
22 assessed by visually inspecting the size of the 95% CI prediction region in the
23 HSROC plot. When considerable heterogeneity was observed, possible reasons for it
24 were explored and subgroup analyses were performed, when possible, according to
25 the pre-specified subgroups in the review protocol.

26 When no plausible explanation for the heterogeneity could be found, the quality of
27 the evidence was downgraded in GRADE by 1 or 2 levels for the domain of
28 inconsistency, depending on the extent of heterogeneity in the results.

1 **Assessing indirectness in diagnostic test accuracy reviews**

2 Indirectness in diagnostic test accuracy studies was assessed using the QUADAS-2
3 checklist by assessing the applicability of the studies in relation to the review
4 question in the following domains (see **Error! Reference source not found.**):

- 5 • patient selection
- 6 • index test
- 7 • reference standard.

8 **Assessing imprecision and clinical significance in diagnostic test accuracy** 9 **reviews**

10 In diagnostic accuracy measures, it was first considered whether sensitivity,
11 specificity, positive likelihood ratios or negative likelihood ratios would be given more
12 weight in the decision-making process. If one measure was given more importance
13 than the other, then imprecision was rated on this statistical measure using the
14 following MID thresholds:

- 15 • sensitivity and specificity
 - 16 ○ high: more than 90%
 - 17 ○ moderate: 75-90%
 - 18 ○ low: less than 75%
- 19 • positive likelihood ratio:
 - 20 ○ very useful test: more than 10
 - 21 ○ moderately useful test: 5-10
 - 22 ○ not a useful test: less than 5
- 23 • negative likelihood ratio:
 - 24 ○ very useful test: less than 0.1
 - 25 ○ moderately useful test: 0.1 to 0.2
 - 26 ○ not a useful test: more than 0.2.

27 **Evidence statements**

28 Evidence statements are summary statements that are presented after the GRADE
29 profiles, summarising the key features of the clinical evidence presented. The
30 wording of the evidence statements reflects the certainty or uncertainty in the
31 estimate of effect. The evidence statements are presented by outcome or theme and
32 encompass the following key features of the evidence:

- 33 • the quality of the evidence (GRADE rating)
- 34 • the number of studies and the number of participants for a particular outcome
- 35 • a brief description of the participants
- 36 • the clinical significance of the effect and an indication of its direction (for example,
37 if a treatment is clinically significant [beneficial or harmful] compared with another,
38 or whether there is no clinically significant difference between the tested
39 treatments).

1 Reviewing economic evidence

2 Inclusion and exclusion of economic studies

3 The titles and abstracts of papers identified through the searches were independently
4 assessed for inclusion using pre-defined eligibility criteria defined in Table 7.

5 **Table 7: Inclusion and exclusion criteria for the systematic reviews of**
6 **economic evaluations**

Inclusion criteria
Intervention or comparators according to the scope
Study population according to the scope
Full economic evaluations (cost-utility, cost-effectiveness, cost-benefit or cost-consequence analyses) that assess both the costs and outcomes associated with the interventions of interest
Exclusion criteria
Abstracts with insufficient methodological details
Cost-of-illness type studies
Conference papers pre January 2014

7 Once the screening of titles and abstracts was complete, full versions of the selected
8 papers were acquired for assessment. The Preferred Reporting Items for Systematic
9 Reviews and Meta-Analyses (PRISMA) for this search on economic evaluations is
10 presented in the Health Economics Chapter.

11 The quality of evidence was assessed using the economic evaluations checklist as
12 specified in the [NICE guidelines manual \(NICE 2014\)](#). Quality assessments of
13 included studies and data extraction tables are provided in Appendix B of the
14 evidence review chapters. The excluded economic studies list is presented in the
15 management evidence review chapter.

16 Health economic modelling

17 The aims of the health economic input to the guideline were to inform the guideline
18 committee of potential economic issues related to the diagnosis and management of
19 heavy menstrual bleeding to ensure that recommendations represented a cost-
20 effective use of healthcare resources. Health economic evaluations aim to integrate
21 data on healthcare benefits (ideally in terms of quality-adjusted life-years, QALYs)
22 with the costs of different care options. In addition, the health economic input aimed
23 to identify areas of high resource impact; recommendations which – while
24 nevertheless cost-effective – might have a large impact on Clinical Commissioning
25 Group or Trust finances and so need special attention.

26 The guideline committee prioritised a single economic model on diagnosis and
27 management where it was thought that economic considerations would be
28 particularly important in formulating recommendations and a review of the health
29 economic literature was undertaken. This model covered multiple review questions,
30 as a complete health economic analysis of the treatment pathway required
31 consideration of all possible combinations of diagnostic strategy and treatment
32 strategy together.

1 Cost effectiveness criteria

- 2 NICE's report Social value judgements: principles for the development of NICE
3 guidance ([https://www.nice.org.uk/media/default/about/what-we-do/research-and-](https://www.nice.org.uk/media/default/about/what-we-do/research-and-development/social-value-judgements-principles-for-the-development-of-nice-guidance.pdf)
4 [development/social-value-judgements-principles-for-the-development-of-nice-](https://www.nice.org.uk/media/default/about/what-we-do/research-and-development/social-value-judgements-principles-for-the-development-of-nice-guidance.pdf)
5 [guidance.pdf](https://www.nice.org.uk/media/default/about/what-we-do/research-and-development/social-value-judgements-principles-for-the-development-of-nice-guidance.pdf)) sets out the principles that committees should consider when judging
6 whether an intervention offers good value for money. In general, an intervention was
7 considered to be cost effective if either of the following criteria applied (given that the
8 estimate was considered plausible):
- 9 • the intervention dominated other relevant strategies (that is, it was both less costly
10 in terms of resource use and more clinically effective compared with all the other
11 relevant alternative strategies), or
 - 12 • the intervention cost less than £20,000 per QALY gained compared with the next
13 best strategy, or
 - 14 • the intervention provided clinically significant benefits at an acceptable additional
15 cost when compared with the next best strategy.
- 16 The committee's considerations of cost-effectiveness are discussed explicitly under
17 the 'Consideration of economic benefits and harms' heading of the relevant sections.

18 Developing recommendations

19 Guideline recommendations

- 20 Recommendations were drafted on the basis of the committee's interpretation of the
21 available evidence, taking into account the balance of benefits, harms and costs
22 between different courses of action. When clinical and economic evidence was of
23 poor quality, conflicting or absent, the committee drafted recommendations based on
24 their expert opinion. The considerations for making consensus-based
25 recommendations include the balance between potential harms and benefits, the
26 economic costs or implications compared with the economic benefits, current
27 practices, recommendations made in other relevant guidelines, patient preferences
28 and equality issues.
- 29 The main considerations specific to each recommendation are outlined under the
30 'Recommendations and link to evidence' headings within each chapter.
- 31 For further details please refer to the [NICE guidelines manual \(NICE 2014\)](#).

32 Research recommendations

- 33 When areas were identified for which good evidence was lacking, the committee
34 considered making recommendations for future research. For further details please
35 refer to the [NICE guidelines manual \(NICE 2014\)](#).

36 Validation process

- 37 This guidance is subject to a 6-week public consultation and feedback as part of the
38 quality assurance and peer review of the document. All comments received from
39 registered stakeholders are responded to in turn and posted on the NICE website at
40 publication. For further details please refer to the [NICE guidelines manual \(NICE](#)
41 [2014\)](#).

1 **Updating the guideline**

- 2 Following publication, and in accordance with the NICE guidelines manual, NICE will
- 3 undertake a review of whether the evidence base has progressed significantly to alter
- 4 the guideline recommendations and warrant an update. For further details please
- 5 refer to the [NICE guidelines manual \(NICE 2014\)](#).

6 **Funding**

- 7 The NGA was commissioned by NICE to undertake the work on this guideline.

References

- 1 **Harbord 2008**
- 2 **Harbord 2008**
- 3 Harbord, R., metandi: Stata module for meta-analysis of diagnostic accuracy,
4 Statistical Software Components, Boston College Department of Economics, Revised
5 15 Apr 2008.
- 6 **Harbord and Whiting 2009**
- 7 Harbord, R. M., Whiting, P., metandi: Meta-analysis of diagnostic accuracy using
8 hierarchical logistic regression, *Stata Journal*, 9, 211-29, 2009
- 9 **Herman 2016**
- 10 Herman, M.C., Penninx, J.P., Geomini, P.M., Mol, B.W., Bongers, M.Y., Choice of
11 primary outcomes evaluating treatment for heavy menstrual bleeding, *BJOG: An
12 International Journal of Obstetrics & Gynaecology*, 1, 123(10), 1593-8, 2016
- 13 **NICE 2014**
- 14 National Institute for Health and Care Excellence (NICE), [NICE guidelines manual](#),
15 2014 [online, accessed 8th June 2017]
- 16 **Santesso 2016**
- 17 Santesso, N., Carrasco-Labra, A., Langendam, M., Brignardello-Petersen, R.,
18 Mustafa, R.A., Heus, P., Lasserson, T., Opiyo, N., Kunnamo, I., Sinclair, D. and
19 Garner, P., Improving GRADE evidence tables part 3: detailed guidance for
20 explanatory footnotes supports creating and understanding GRADE certainty in the
21 evidence judgments, *Journal of clinical epidemiology*, 74, 28-39, 2016.
- 22 **Turner 2012**
- 23 Turner, R. M., Davey, J., Clarke, M. J., Thompson, S. G., Higgins, J. P. T., Predicting
24 the extent of heterogeneity in meta-analysis, using empirical data from the Cochrane
25 Database of Systematic Reviews, *International Journal of Epidemiology*, 41, 818-
26 827, 2012
- 27 **Van Valkenhoef 2016**
- 28 Van Valkenhoef, G., Dias, S., Ades, A. E., Welton, N. J., Automated generation of
29 node-splitting models for assessment of inconsistency in network meta-analysis,
30 *Research Synthesis Methods*, 7, 80-93, 2016