

Title: Dynamic methods for ongoing assessment of site-level risk in risk-based monitoring of clinical trials: a scoping review

Running Head: Dynamic methods for site-level risk assessment

Version + Date: Version 1.0, 29th February 2020

Word count: Abstract: 409
Text (excluding abstract, references, appendix, tables and figures): 3984

Tables/Figures: 4 tables, 2 figures

Supplement: 3 files

Authors:	Name	Affiliation
	William J Cragg ORCID: 0000-0002-1274-8521	1: MRC Clinical Trials Unit at UCL, London, UK; 2: Clinical Trials Research Unit, Leeds Institute of Clinical Trials Research, University of Leeds, Leeds, UK
	Caroline Hurley ORCID: 0000-0002-3793-3408	Health Research Board-Trials Methodology Research Network (HRB-TMRN), National University of Ireland, Galway, Ireland
	Victoria Yorke-Edwards ORCID: 0000-0002-6465-7330	MRC Clinical Trials Unit at UCL, London, UK
	Sally P Stenning ORCID: 0000-0001-7009-5488	MRC Clinical Trials Unit at UCL, London, UK

Correspondence: William Cragg
w.cragg@leeds.ac.uk
Tel: 0113 343 8398
Clinical Trials Research Unit

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

Leeds Institute of Clinical Trials Research
University of Leeds
Leeds
LS2 9JT,
United Kingdom

Research support: Medical Research Council (MC_UU_12023/24)
The TEMPER study was funded by a grant from Cancer Research UK (C1495/A13305)

Summary of figures and tables

Table 1: general characteristics of included studies

Table 2: full listing of all included reports

Table 3: types of assessments and evidence presented by reports that included some assessments of their methods' effectiveness

Table 4: best reported information on methods' classification ability, where available or deductible

Figure 1: PRISMA flow diagram

Figure 2: publications by year and type

Summary of supplementary information

1. Completed PRISMA-ScR Checklist
2. Search strategy from Medline
3. Annotated list of data collection variables

1 **Abstract**

2 **Background/Aims**

3 It is increasingly recognised that reliance on frequent site visits for monitoring clinical trials
4 is inefficient. Regulators and trialists have in recent years encouraged more risk-based
5 monitoring. Risk assessment should take place before a trial begins in order to define the
6 overarching monitoring strategy. It can also be done on an ongoing basis, in order to target
7 sites for monitoring activity. Various methods have been proposed for such prioritisation,
8 often using terms like 'central statistical monitoring', 'triggered monitoring' or, as in ICH
9 Good Clinical Practice guidance, 'targeted on-site monitoring'. We conducted a scoping
10 review to identify such methods, to establish if any published methods were supported by
11 adequate evidence to allow wider implementation, and to point the way to future
12 developments in this field of research.

13 **Methods**

14 We used 7 publication databases, 2 sets of methodological conference abstracts and an
15 internet search engine to look for methods for using centrally held trial data to assess site
16 conduct during a trial. We included only reports in English, and excluded reports published
17 before 1996 and reports not directly relevant to our research question. We used reference
18 and citation searches to find additional relevant reports. We extracted data using a pre-
19 defined template. We contacted authors to request additional information about included
20 reports and to check whether reports might be eligible.

21 **Results**

22 We included 30 reports in our final dataset, of which 21 were peer-reviewed publications.
23 20 reports described central statistical monitoring methods (of which 7 focussed on

1 detection of fraud or misconduct) and 9 described triggered monitoring methods. 21 reports
2 included some assessment of their methods' effectiveness. Most commonly this involved
3 exploring the methods' characteristics using real trial data with no known integrity issues. Of
4 the 21 with some effectiveness assessment, most presented limited or no information about
5 whether or not concerns identified through central monitoring constituted meaningful
6 problems. Some reports commented on cost savings from reduced on-site monitoring, but
7 none gave detailed costings for the development and maintenance of central monitoring
8 methods themselves.

9 **Conclusions**

10 Our review identified various proposed methods, some of which could be combined within
11 the same trial. The apparent emphasis on fraud detection may not be proportionate in all
12 trial settings. Although some methods have self-justifying benefits for data cleaning activity,
13 many have limitations that may currently prevent their routine use for targeting trial
14 monitoring activity. The implementation costs, or uncertainty about these, may also be a
15 barrier. We make recommendations for how the evidence-base supporting these methods
16 could be improved.

17

18 **Keywords:** trial monitoring, risk-based monitoring, triggered monitoring, central statistical
19 monitoring, Good Clinical Practice, research misconduct, data fabrication

20

1 Introduction

2

3 Monitoring, a major component of assuring the quality of clinical trials, has traditionally
4 relied on frequent on-site monitoring visits,¹ particularly to facilitate sometimes extensive
5 source data verification (SDV).² However, it is increasingly recognised that this model may
6 be inefficient and unnecessary in many cases,^{3,4} with trialists questioning the value of 100%
7 SDV.⁵⁻⁷ In recent years, regulators⁸⁻¹⁰ and trialists^{1,11} have proposed a risk-based approach
8 to monitoring, whereby monitoring methods, including the frequency and nature of on-site
9 visits, vary across trials depending on the risks specific to each one.

10 Risk-based monitoring methods can be applied at different stages of a trial. Pre-trial risk
11 assessments can help define the overarching strategies appropriate to the trial's risks. In
12 some models,^{12,13} this is predominantly a one-off assessment during trial setup. However, it
13 is also possible to modify the monitoring strategy, or build in flexibility and responsiveness,
14 based on changing or emerging risks during the course of the trial.¹⁴

15 Risk-based monitoring is often associated with fewer on-site visits than 'traditional'
16 monitoring.¹⁵ Although effective central monitoring methods alone could, in some respects,
17 provide adequate trial monitoring in place of visits, on-site visits offer particular benefits
18 over central monitoring, for example in any activities that require access to site-held source
19 data (such as patients' medical notes). On-site visits may also be necessary, for example, to
20 investigate potential fraudulent activity. In a risk-based monitoring framework, visits to sites
21 may not be routine, but can be based on assessed risk; we therefore need methods to
22 assess site-level risk on an ongoing basis. We can interpret these methods as assessing the
23 risk of *not going to site now*. If the risk seems too high, a visit – or some other corrective
24 action – is triggered. Methods of this kind have been referred to using various terms,

1 including ‘triggered monitoring’¹⁴ or, as in ICH GCP guidance, ‘targeted monitoring’,¹⁶ and
2 may employ data-driven approaches from methods known collectively as ‘central statistical
3 monitoring’,¹⁷ or more subjective assessments.^{14,18,19}

4 A recent systematic review has established the breadth of tools available to assess overall
5 trial risk (and to use this assessment to define the monitoring strategy) in the setup stage,²⁰
6 but so far there has been no such exercise for methods to assess ongoing site-level risk once
7 a trial has started. We conducted a scoping review²¹ to identify and characterise available
8 methods.

9 Our aims were 1) for trialists, to establish if any published methods were supported by
10 adequate evidence to support implementation in routine practice and 2) for researchers in
11 this area, to consolidate the existing evidence and point towards future developments in
12 this growing field.

13

1 **Methods**

2 We conducted a scoping review to identify methods for using centrally held clinical trial data
3 to assess site-level risk of deviations from Good Clinical Practice (GCP) or the trial protocol,
4 or research misconduct, and thereby to target sites for further monitoring activity. We
5 chose scoping review methodology as we anticipated finding a variety of methods and study
6 types, and we wanted to characterise the extent, range and nature of research activity.²²

7 There is no published protocol for this scoping review.

8

9 **Eligibility criteria**

10 We defined our eligibility criteria before beginning any searches, with minor refinements
11 (mainly to the exclusion criteria) after having piloted the search strategy.

12 We included reports:

- 13 - Describing methods for using centrally held data (i.e. at the clinical trials unit or
14 other trial coordinating centre) to assess, in ongoing trials, site-level risk of protocol
15 or GCP deviation, risk of data fabrication or research misconduct, or to target sites in
16 some other way for corrective action based on assessed risk (regardless of whether
17 the corrective action involved an on-site monitoring visit or not);
- 18 - With methods described in enough detail that we considered them – subjectively –
19 reproducible;
- 20 - Either published in peer-reviewed journals or available as grey literature;
- 21 - About clinical trials, not limited to trials of Investigational Medicinal Products;
- 22 - In English.

1

2 We excluded reports:

- 3 - Published before 1996 (the year of the first version of ICH GCP Guidance, E6[R1]²³);
- 4 - About quality assurance only in the context of intervention fidelity²⁴ or ‘rater
- 5 differences’²⁵ for subjective trial outcome measures;
- 6 - About ‘data monitoring’ in general, for example data monitoring committees, unless
- 7 including methods for assessing site conduct;
- 8 - About ‘monitoring’ in any sense other than the GCP sense, e.g. clinical monitoring;
- 9 - Focusing only on trial recruitment;
- 10 - About more efficient alternatives to standard on-site activity, for example remote
- 11 source data verification;
- 12 - About site selection during trial setup;
- 13 - Featuring only opinions, or lacking enough detail allowing theoretical reproduction
- 14 of methods;
- 15 - Reviewing methods presented elsewhere without any new, original methods.

16

17

18 **Information sources and search strategies**

19 **1. Database searches**

20 We designed search strategies for the following databases:

- 21 - PubMed

- 1 - Embase (Ovid)
- 2 - Medline (Ovid)
- 3 - Web of Science (Clarivate Analytics)
- 4 - Cinahl
- 5 - Cochrane Central
- 6 - Scopus

7

8 Full database searches took place on 23 October 2017 (run and extracted by WC). The
9 search strategy for Medline is given in Supplementary Information. We developed our
10 search strategy following review of systematic reviews in this area^{1,26} to identify relevant
11 search terms. We developed the final list through an iterative process in September and
12 October 2017. The final search term combined searches around two concepts: clinical trials
13 (using terms based on those used in a previous systematic review of monitoring methods²⁶)
14 and targeted or risk-based clinical trial monitoring. Search terms were modified as required
15 to suit the database being used. No database filters were applied.

16 Both reviewers (WC, CH) imported results from the seven databases into reference
17 management software and used in-built tools to remove duplicate entries. Both reviewers
18 carried out initial title and abstract screening for relevant reports, producing an initial
19 shortlist of potential papers. These were reduced through review and discussion to a final
20 list of relevant reports for full-text review. We reviewed and discussed these, using full-text
21 reports where possible, to agree a final list of relevant reports. Throughout the process, SS
22 acted as third reviewer where required.

1 In order to ensure our results were current, this element of the search strategy was
2 repeated on 28 August 2018. WC ran the searches and conducted the title and abstract
3 screening. A shortlist of potentially relevant reports was shared with SS and CH; SS and WC
4 agreed a final list of additional relevant reports from this repeated search.

5

6 **2. Conference abstracts**

7 We searched for relevant conference abstracts from the first four International Clinical Trials
8 Methodology Conferences (occurring between 2011 and 2017) and all annual meetings of
9 the Society for Clinical Trials since 1996 (initial searches completed on 8 December 2017).

10 Our methods to search conference abstracts varied depending on the year and the
11 conference: we used website search functions where available and manual PDF searches
12 otherwise. Keywords used for the conference abstract search were based on the key
13 database search strategy terms. These were: 'monitor', 'supervision', 'oversight', 'risk',
14 'performance', 'metric', 'quality', 'fraud', 'fabrication', and 'error'.

15 Both WC and CH performed the abstract searches. This produced an initial shortlist of
16 potentially relevant abstracts. A final list was agreed through discussion, with SS acting as
17 third reviewer where required.

18

19 **3. Internet searches**

20 We conducted structured searches through Google internet search engine (searches carried
21 out 15-19 December 2017).

1 Google searches were performed without limitations or use of quotes. Search terms were
2 based on the main database search: ‘Risk based monitoring’, ‘Risk adapted monitoring’,
3 ‘Central monitoring’, ‘Central statistical monitoring’, ‘Triggered monitoring’, ‘Targeted
4 monitoring’, ‘Performance metric’, ‘Site metric’, ‘Key risk indicator’, ‘Site performance’,
5 ‘Centre performance’, ‘Detect fraud’ and ‘Detect fabrication’. We reviewed results on the
6 first 20 pages; if there were no relevant results on any three consecutive pages before this,
7 we stopped reviewing.
8 WC and CH conducted the searches. Any potential additions to the included list of reports
9 were discussed and agreed, with SS acting as third reviewer where required.

10

11 **4. References, citations and author contact**

12 To identify other relevant reports, we reviewed references (manually) and citations (using
13 Web of Science) of all papers included or considered for inclusion in the final results, and of
14 review articles relevant to the topic. Whenever required, we contacted report authors to
15 help ascertain if given reports should be included, and to ask about the availability of full-
16 text articles.

17

18

19 **Data collection**

20 We extracted data from full journal articles, where available. We recorded data into an
21 Excel-based tool. WC carried out the final data collection used for this report, with SS
22 double-checking all data for inclusion; consensus was reached on any areas of

1 disagreement. Article authors were contacted (two attempts maximum) for missing
2 descriptive data and further clarifications.

3 Our data collection template was designed and agreed prior to any data collection, with
4 minor refinement after a first review of all relevant papers (a list of data collection variables
5 is available as Supplementary Information). We collected descriptive data about each of the
6 included reports, including year of publication, type of report and details of the trial(s) it was
7 embedded in. We also looked for any information on cost implications of the proposed
8 methods.

9 When designing this study, although we predicted we would find a range of methods, we
10 agreed that most of them would in essence address a classification problem, i.e. methods to
11 assign sites a status as ‘concerning’ or ‘not-concerning’, with a ‘true’ deviation status – i.e.
12 confirmed existence of meaningful problems – that could be uncovered by further review.
13 The ‘gold standard’ reference test required to assess true status might be study-specific, but
14 could be on-site monitoring or, if the true status was created through simulation, prior
15 knowledge.

16 We considered a key measure of the reported methods’ effectiveness to be a demonstrated
17 ability, ideally in a real-life setting, not only to detect ‘true’ sites of concern, but also to
18 show with confidence that sites apparently not of concern are performing well. We
19 therefore aimed to summarise the available information on classification, i.e. any or all of
20 specificity, sensitivity, positive and negative predictive value. We gathered the best reported
21 classification statistics for each method, or, if this was not reported, used available statistics,
22 e.g. number of true and false positives, to calculate these. These calculations were verified

1 by an independent statistician at the Medical Research Council Clinical Trials Unit at

2 University College London.

3 We did not formally assess the quality of the studies. However, review of the QUADAS-2

4 tool for quality assessment in diagnostic accuracy studies²⁷ informed development of our

5 data collection template.

6

7 **Synthesis of results**

8 We did not combine results for individual studies, as it was clear through preliminary review

9 of relevant papers that we would have a variety of study types. Instead, the evidence is

10 summarised descriptively.

11

1 Results

2

3 **Figure 1** gives a PRISMA flow diagram²⁸ showing the different stages of the review. From the
4 various data sources, we ultimately included 30 reports in our final dataset. 21 of these are
5 peer-reviewed publications. The results are characterised in **Table 1** and listed in full in
6 **Table 2. Figure 2** shows reports by year of publication.

7 When specific trials were mentioned, they involved various health conditions and were in
8 various geographical settings. Information on trial intervention was not often available, but
9 where it was, methods had most often been used in phase III trials of investigational
10 medicinal products (IMP). The IMP risk category,²⁹ when known, was either ‘licensed and
11 used within its licensed indication’, or ‘licensed and used outside its licensed indication’ (i.e.
12 we found no reports involving trials of unlicensed IMPs).

13 We classified 20/30 of our results as central statistical monitoring methods, of which 7
14 focussed on detection of investigator fraud or research misconduct. We classified 9,
15 including one of the 20 that used central statistical monitoring, as ‘triggered monitoring’, i.e.
16 review of each trial site against pre-set thresholds in key performance metrics, usually
17 without any statistical testing. A final 2 we could not fit into either of these categories; these
18 involved using measured site metrics to directly compare sites against one another.^{30,31}

19 21/30 reports included some assessment of the effectiveness of the methods; these are
20 summarised in **Table 3**. The most common experimental designs were to explore the
21 methods’ characteristics using real trial data with no known integrity issues (n=9), and
22 simulating data integrity problems at sites within real trial datasets then using the method
23 to try to identify the problem sites (n=6).

1 Of the 21 reports, 9 had no information about sites' 'true' status, i.e. whether the problems
2 identified through central monitoring constitute meaningful problems (either recorded
3 through on-site monitoring or audit activity, or known because statuses were created
4 through simulation). One report³² only contained case studies, i.e. partial and selective
5 reporting. Seven^{14,33-38} had partial information, e.g. some of sites' true statuses were
6 reported, but not all. Two explored classification ability through extensive simulation,^{39,40}
7 and 2 had detailed information from a limited set of scenarios on the number of true and
8 false positives and negatives.^{19,41} The best reported or deducible classification ability for the
9 11 papers with at least some information on sites' 'true' status is shown in **Table 4**; in many
10 cases this is limited by considerable amounts of missing or unclear data.

11 Some papers report on actual or theoretical cost savings from reduced on-site
12 monitoring,^{33,42,43} and others comment on the risk of incurring costs if their proposed
13 central monitoring method identifies sites that do not in fact have meaningful problems (i.e.
14 false positives).^{19,37} However, no papers give detailed costings for the development,
15 implementation and maintenance of the central monitoring methods themselves.

16

1 Discussion

2

3 We conducted a scoping review to identify published methods for assessing the risk of not
4 taking corrective action at trial sites at a given time. Although our search looked for reports
5 from any time after 1995, over half of our results are from after 2013, highlighting the
6 recent growth of risk-based monitoring concepts. Most trials for these methods were most
7 often phase III trials of licensed investigational medicinal products, i.e. somewhat lower risk
8 compared to earlier phase trials of newer treatments.²⁹ Around a third of our results were
9 not full, peer-reviewed reports, reflecting a wider problem with evidence supporting trial
10 conduct methods, i.e. that researchers can feasibly produce posters and abstracts for
11 conferences, but may not have time or incentive to produce comprehensive reports.⁴⁴

12 Identified methods were mainly in two broad categories (with some overlap). Most were
13 about central statistical monitoring, which uses statistical testing of all or a subset of trial
14 data items to compare sites and identify atypical trial centres. A minority described
15 triggered monitoring techniques, whereby sites are assessed against pre-specified site
16 metric threshold rules (usually binary), with sites meeting the greatest number of ‘triggers’
17 being considered the most concerning. Several authors note that central statistical
18 monitoring needs sufficient overall and per-site sample sizes for adequate statistical
19 power.^{17,19,37} Triggered monitoring, however, can be used at any stage of a trial’s
20 recruitment (especially with trigger rules based on single instances of a given protocol
21 violation, for example). We therefore suggest use of the techniques is not mutually
22 exclusive.

1 Although the reports on central statistical monitoring described a range of methods, there
2 were some commonalities. Previous papers have reviewed these methods in more
3 detail.^{17,41,45-47}

4 Nearly half of the central statistical monitoring reports focused on identification of data
5 fabrication. The possibility of fraud is a serious concern to trialists and a threat to wider
6 trust in science.⁴⁸ It was possibly an important factor in establishing 100% on-site
7 verification of trial data as a common monitoring approach.^{49,50} This may help explain the
8 prevalence of reports about fraud detection, as some may see the priority in risk-based
9 monitoring to be establishing its fraud detection ability compared with 100% source data
10 verification. However, although the incidence of data fraud is difficult to quantify, cases of
11 extensive data fabrication appear rare enough to have individual notoriety.⁵¹ Further,
12 methods to detect fraud are necessarily rather selective, and therefore may not alone be
13 suitable for trialists looking to detect more common, lower-level data integrity issues.

14 We collected data on how the proposed methods we identified had been evaluated. A
15 number of reports only presented proposed, untested methods, or only selected case
16 studies to demonstrate the methods' performance. Of those that presented more detailed
17 evaluation, a common limitation was that the 'true' status both of identified problem sites
18 and sites apparently not of concern was often not available, or only partially available. It
19 was therefore difficult to know if the 'concern' status of sites in central monitoring results
20 represented meaningful problems or not. In addition, a number of studies use simulation to
21 create 'true' sites of concern; these raise the additional question of whether these
22 simulations reflect real-life issues, though the involvement of clinicians (i.e. those who
23 would provide real-life trial data) in the simulation process of some reports^{19,36} is reassuring.

1 Clearly, a balance must be struck between minimising cost and adequately mitigating key
2 risks. Many false positives could be expensive (e.g. additional on-site monitoring visits for no
3 gain), but many false negatives could be catastrophic (e.g. missing serious risks to trial
4 participants' safety). Some authors do acknowledge this, and even suggest that a high false
5 positive rate could be acceptable in a method used to broadly target additional scrutiny at
6 particular sites.³⁸

7 It is important to recognise the limitations of the available 'gold-standards' in the
8 classification of sites. When methods are tested using simulated or real-but-adjusted data, it
9 may be difficult to know how well these accurately recreate real-life situations. When
10 central monitoring methods are tested in real, ongoing trials, on-site monitoring may be an
11 imperfect reference test, in that it may not be able to identify all problems. By contrast, it is
12 clear that central monitoring, with its enhanced inter- and intra-site review, can identify
13 issues that a single team at one site for a limited time might not.⁴⁷

14 It could be argued that at least some of the methods we have identified do not need
15 extensive evaluation because they prove their own worth. For instance, they help identify
16 outliers that in some cases are self-evidently meaningful problems to resolve. We
17 acknowledge that some central monitoring activities identify 'known' problems (e.g.
18 identifying weekend visit dates, which in most cases are unlikely to be correct) and are
19 valuable for data cleaning purposes. However, we were specifically interested in the more
20 nuanced use of these methods to identify sites of 'concern', at which monitoring activity
21 may be targeted, and consequently sites 'not of concern', monitoring of which may be
22 reduced or omitted. This element does need adequate proof before wider adoption and, if it
23 is shown to be effective, could have significant benefits. In light of the limitations we have

1 described here, we do not believe any methods have yet fully demonstrated that they
2 should be adopted more widely as a means to limit on-site monitoring of sites deemed,
3 based on review of centrally-held data, not to be of concern.

4 Aside from some comments on the potential cost of investigating false positive central
5 monitoring results,^{19,52} the reports we identified contained limited information on the cost
6 of developing and implementing their methods. As well as uncertainty about how to
7 develop relevant methods, uncertainty or concern about costs involved is a substantial
8 barrier to adoption of risk-based monitoring.⁵³

9 Through our various searches, we identified 24 commercial companies which, we were at
10 least reasonably certain, had a method relevant to our search (data not shown). Only a small
11 proportion of these had published details of their methods in peer-reviewed journals,
12 although just over half had some detail in grey literature sources. This highlights another
13 difficulty in disseminating new trial conduct methods, i.e. that they may be commercially
14 sensitive and, unlike evidence about treatments being tested in trials, there is no
15 compulsion to publish new evidence. It is somewhat contradictory that while risk-based
16 monitoring has come about partly to reduce costs of trials, some risk-based monitoring
17 methods are available only for a fee.

18 Further work is needed to fully demonstrate the effectiveness of these dynamic site risk
19 assessment methods which, alongside pre-trial risk assessments, form the core of risk-based
20 monitoring. We therefore recommend the following:

21 **1. Coordinate research efforts.** From the scoping review and contact with report
22 authors, it was clear that various small research projects relevant to this topic

1 were ongoing, but mostly in isolation. Researchers in this area should take stock
2 of existing research, and set clear priorities to ensure research time is well-spent.

3 **2. Standardise monitoring studies.** Core outcome sets⁵⁴ or other mechanisms to
4 standardise studies about monitoring would improve study quality and may
5 facilitate cross-study evidence synthesis.

6 **3. Share evidence.** Time, commercial sensitivity and perceived reputational risk
7 could all be barriers to publishing evidence about monitoring practices. However,
8 additional, publicly-available evidence to support best monitoring practice will
9 allow trialists in all settings to adopt new methods with confidence.

10 **4. Publish full papers.** Conference abstracts and posters are a useful way to
11 disseminate basic information about new ideas, but rarely have enough detail to
12 allow replication or robustly demonstrate effectiveness. As this emerging field
13 cannot be built on abstracts alone, we encourage researchers to publish full,
14 peer-reviewed papers about their monitoring methods.

15 **5. Look to combine complementary methods.** Although work has been done on a
16 number of distinct risk-based monitoring methods, an optimal monitoring plan
17 might involve a combination of these, including both central statistical
18 monitoring and triggered monitoring. A collaborative approach to combining
19 existing methods could help develop and test such an idea.

20

21 We acknowledge several limitations. Our database searches identified relevant material
22 from disparate locations, including abstracts in conferences in unrelated research fields. It is
23 possible that other abstract collections include relevant material, but it was not feasible to

1 find and hand-search all of these. Although the internet searches made little contribution to
2 the final list of included reports, they may have been limited by known reproducibility
3 problems.⁵⁵

4 Scoping review methodology advises that relevant experts in a field are surveyed to help
5 identify other relevant work.⁵⁶ We have not formally done this. We have, however,
6 contacted most authors of included reports for clarifications, and this has not highlighted
7 any additional relevant reports.

8 Although we pre-specified our aims and our eligibility criteria, we added some detail to our
9 exclusion criteria during the course of our review to help us decide about certain reports. It
10 is possible that other researchers repeating the same review might result in a slightly
11 different list, but we believe this might only affect the 'method-only' papers, which do not
12 form the key part of our conclusions. The comprehensive nature of our search strategy,
13 including review of reference and citation lists, gives us confidence that our report is a
14 sound overview of the state of the evidence in this research area.

15 We have not performed a formal quality assessment of reports we found, however, this is
16 considered by some to be unnecessary in scoping review methodology.²² There is also no
17 validated way to review the quality of risk-based monitoring studies, although we used the
18 QUADAS-2 tool, designed to assess the quality of diagnostic studies, to inform our data
19 collection template.

20 Finally, we acknowledge that some time has now passed since we first conducted our search
21 for relevant evidence. Conscious of this, we repeated the main database search in 2018
22 (albeit with only one author conducting title and abstract screening) and added three
23 relevant reports. We are not aware of any research published since then that might change

1 our overall conclusions. If evidence is now available that addresses the limitations we have
2 highlighted in the existing literature, we would certainly consider this a positive
3 development.

4 Our scoping review highlighted some promising evidence for risk-based monitoring in
5 ongoing trials. However, currently published methods may not yet have demonstrated their
6 efficacy or cost-effectiveness well enough for trialists to implement them with confidence as
7 a means to target or omit on-site visits. A more coordinated, collaborative and transparent
8 approach to developing and sharing evidence in this field, including industry and academic
9 partners, could help it grow beyond its current nascent state, and could contribute to risk-
10 based monitoring more quickly entering routine practice.

11

1 **Acknowledgements**

2 We would like to thank the authors of reports – some referenced in this manuscript and
3 some ultimately not – for their time in answering our questions and discussing their work
4 with us. We would also like to thank: systematic reviewers at Medical Research Council
5 Clinical Trials Unit at University College London (MRC CTU at UCL) for their advice; Meredith
6 Martyn for verifying the classification statistics, Sharon B Love for project support and
7 Matthew R Sydes for advice, project support and review of this manuscript.

8

1 **Declaration of Conflicting Interests**

2 SS and WC are part of the TEMPER study team, one of the studies included in the final

3 results of this paper. CH and VYE declare no conflict of interest.

4

1 **Funding**

- 2 This work was supported by the MRC London Hub for Trial Methodology Research
- 3 (MC_UU_12023/24). The idea for this work arose from the TEMPER study, which was
- 4 funded by a grant from Cancer Research UK (C1495/A13305).

1 **Supplementary information**

- 2 1. Completed PRISMA-ScR Checklist
- 3 2. Search strategy from Medline
- 4 3. Annotated list of data collection variables
- 5

1 References

- 2
- 3 1 Macefield RC, Beswick AD, Blazeby JM, *et al.* A systematic review of on-site
4 monitoring methods for health-care randomised controlled trials. *Clin Trials*
5 2013;**10**:104–24.<http://ctj.sagepub.com/content/10/1/104.abstract>
- 6 2 Sheetz N, Wilson B, Benedict J, *et al.* Evaluating Source Data Verification as a Quality
7 Control Measure in Clinical Trials. *Ther Innov Regul Sci* 2014;**48**:671–
8 80.<http://dij.sagepub.com/content/48/6/671.abstract>
- 9 3 Reith C, Landray M, Devereaux PJ, *et al.* Randomized Clinical Trials — Removing
10 Unnecessary Obstacles. *N Engl J Med* 2013;**369**:1061–5. doi:10.1056/NEJMs1300760
- 11 4 Duley L, Antman K, Arena J, *et al.* Specific barriers to the conduct of randomized trials.
12 *Clin Trials* 2008;**5**:40–8.<http://ctj.sagepub.com/content/5/1/40.abstract>
- 13 5 Olsen R, Bihlet AR, Kalakou F, *et al.* The impact of clinical trial monitoring approaches
14 on data integrity and cost—a review of current literature. *Eur J Clin Pharmacol*
15 2016;**72**:399–412. doi:10.1007/s00228-015-2004-y
- 16 6 Tudur Smith C, Stocken DD, Dunn J, *et al.* The value of source data verification in a
17 cancer clinical trial. *PLoS One* 2012;**7**:12.
- 18 7 Buyse M, Squifflet P, Coart E, *et al.* The impact of data errors on the outcome of
19 randomized clinical trials. *Clin Trials* 2017;**14**:174077451771615.
20 doi:10.1177/1740774517716158
- 21 8 MRC/DH/MHRA Joint Project. Risk-adapted Approaches to the Management of
22 Clinical Trials of Investigational Medicinal Products.
23 2011.[https://www.gov.uk/government/uploads/system/uploads/attachment_data/file/343677/Risk-](https://www.gov.uk/government/uploads/system/uploads/attachment_data/file/343677/Risk-adapted_approaches_to_the_management_of_clinical_trials_of_investigational_medicinal_products.pdf)
24 [adapted_approaches_to_the_management_of_clinical_trials_of_investigational_medicinal_products.pdf](https://www.gov.uk/government/uploads/system/uploads/attachment_data/file/343677/Risk-adapted_approaches_to_the_management_of_clinical_trials_of_investigational_medicinal_products.pdf) (accessed 8 Sep 2017).
- 25
26
- 27 9 Food and Drug Administration. Guidance for Industry: Oversight of Clinical
28 Investigations — A Risk-Based Approach to Monitoring.
29 2013.<http://www.fda.gov/downloads/Drugs/.../Guidances/UCM269919.pdf>
30 (accessed 8 Sep 2017).
- 31 10 European Medicines Agency. Reflection paper on risk based quality management in
32 clinical trials.
33 2013.http://www.ema.europa.eu/docs/en_GB/document_library/Scientific_guideline/2013/11/WC500155491.pdf (accessed 16 Jan 2017).
- 34
- 35 11 Brosteanu O, Houben P, Ihrig K, *et al.* Risk analysis and risk adapted on-site
36 monitoring in noncommercial clinical trials. *Clin Trials* 2009;**6**:585–
37 96.<http://ctj.sagepub.com/content/6/6/585.abstract>
- 38 12 Brosteanu O, Schwarz G, Houben P, *et al.* Risk-adapted monitoring is not inferior to
39 extensive on-site monitoring: Results of the ADAMON cluster-randomised study. *Clin*
40 *Trials* 2017;**14**. doi:10.1177/1740774517724165
- 41 13 Journot V, Pignon JP, Gaultier C, *et al.* Validation of a risk-assessment scale and a risk-

- 1 adapted monitoring plan for academic clinical research studies--the Pre-Optimon
2 study. *Contemp Clin Trials* 2011;**32**:16–24.
- 3 14 Stenning SP, Cragg WJ, Joffe N, *et al.* Triggered or routine site monitoring visits for
4 randomised controlled trials: results of TEMPER, a prospective, matched-pair study.
5 *Clin Trials* 2018;:174077451879337. doi:10.1177/1740774518793379
- 6 15 Sullivan LB. The Current Status of Risk-Based Monitoring.
7 <http://www.appliedclinicaltrialsonline.com/current-status-risk-based-monitoring>
8 (accessed 6 Jun 2018).
- 9 16 International Conference on Harmonisation of Technical Requirements for
10 Pharmaceuticals for Human Use (ICH). Integrated Addendum to ICH E6(R1): guideline
11 for good clinical practice E6(R2).
12 2016.https://www.ich.org/fileadmin/Public_Web_Site/ICH_Products/Guidelines/Efficacy/E6/E6_R2__Step_4_2016_1109.pdf (accessed 25 Feb 2017).
13
- 14 17 Kirkwood AA, Cox T, Hackshaw A. Application of methods for central statistical
15 monitoring in clinical trials. *Clin Trials* 2013;**10**:783–
16 806.<http://ctj.sagepub.com/content/10/5/783.abstract>
- 17 18 Tudur Smith C, Williamson P, Jones A, *et al.* Risk-proportionate clinical trial
18 monitoring: an example approach from a non-commercial trials unit. *Trials*
19 2014;**15**:1–10. doi:10.1186/1745-6215-15-127
- 20 19 Knepper D, Lindblad AS, Sharma G, *et al.* Statistical Monitoring in Clinical Trials: Best
21 Practices for Detecting Data Anomalies Suggestive of Fabrication or Misconduct. *Ther*
22 *Innov Regul Sci* 2016;**50**:144–54. doi:10.1177/2168479016630576
- 23 20 Hurley C, Shiely F, Power J, *et al.* Risk based monitoring (RBM) tools for clinical trials:
24 A systematic review. *Contemp Clin Trials* 2016;**51**:15–27.
25 doi:10.1016/j.cct.2016.09.003
- 26 21 Daudt HM, van Mossel C, Scott SJ. Enhancing the scoping study methodology: a large,
27 inter-professional team’s experience with Arksey and O’Malley’s framework. *BMC*
28 *Med Res Methodol* 2013;**13**:48. doi:10.1186/1471-2288-13-48
- 29 22 Pham MT, Rajić A, Greig JD, *et al.* A scoping review of scoping reviews: advancing the
30 approach and enhancing the consistency. *Res Synth Methods* 2014;**5**:371–85.
31 doi:10.1002/jrsm.1123
- 32 23 International Conference on Harmonisation of Technical Requirements for
33 Pharmaceuticals for Human Use (ICH). Guideline for good clinical practice E6(R1).
34 1996.http://www.ich.org/fileadmin/Public_Web_Site/ICH_Products/Guidelines/Efficacy/E6/E6_R1_Guideline.pdf (accessed 8 Sep 2017).
35
- 36 24 Devito Dabbs A, Song M-K, Hawkins R, *et al.* An intervention fidelity framework for
37 technology-based behavioral interventions. *Nurs Res* 2011;**60**:340–7.
38 doi:10.1097/NNR.0b013e31822cc87d
- 39 25 Opler MGA, Yavorsky C, Daniel DG. Positive and Negative Syndrome Scale (PANSS)
40 Training: Challenges, Solutions, and Future Directions. *Innov Clin Neurosci*
41 2017;**14**:77–81.<http://www.ncbi.nlm.nih.gov/pubmed/29410941> (accessed 4 Jun
42 2018).

- 1 26 Bakobaki J, Joffe N, Burdett S, *et al.* A systematic search for reports of site monitoring
2 technique comparisons in clinical trials. *Clin Trials* 2012;**9**:777–
3 80.<http://ctj.sagepub.com/content/9/6/777.abstract>
- 4 27 Whiting PF, Rutjes AWS, Westwood ME, *et al.* QUADAS-2: A Revised Tool for the
5 Quality Assessment of Diagnostic Accuracy Studies. *Ann Intern Med* 2011;**155**:529.
6 doi:10.7326/0003-4819-155-8-201110180-00009
- 7 28 Moher D, Liberati A, Tetzlaff J, *et al.* Preferred Reporting Items for Systematic Reviews
8 and Meta-Analyses: The PRISMA Statement. *PLoS Med* 2009;**6**:e1000097.
9 doi:10.1371/journal.pmed.1000097
- 10 29 Organisation for Economic Co-operation and Development. OECD Recommendation
11 on the Governance of Clinical Trials. 2013. [https://www.oecd.org/sti/sci-tech/oecd-](https://www.oecd.org/sti/sci-tech/oecd-recommendation-governance-of-clinical-trials.pdf)
12 [recommendation-governance-of-clinical-trials.pdf](https://www.oecd.org/sti/sci-tech/oecd-recommendation-governance-of-clinical-trials.pdf)
- 13 30 Smith A, Seltzer J, A. S. An In-Process Scaling Model: A Potential Framework for Data
14 Monitoring Committees and Clinical Trial Quality Improvement. *Drug Inf J* 2012;**46**:8–
15 12. doi:10.1177/0092861511427864
- 16 31 Djali S, Janssens S, Van Yper S, *et al.* How a Data-Driven Quality Management System
17 Can Manage Compliance Risk in Clinical Trials. *Drug Inf J* 2010;**44**:359–
18 73.<http://ovidsp.ovid.com/ovidweb.cgi?T=JS&PAGE=reference&D=emed12&NEWS=N>
19 [&AN=359239259](http://ovidsp.ovid.com/ovidweb.cgi?T=JS&PAGE=reference&D=emed12&NEWS=N)
- 20 32 Edwards P, Shakur H, Barnetson L, *et al.* Central and statistical data monitoring in the
21 Clinical Randomisation of an Antifibrinolytic in Significant Haemorrhage (CRASH-2)
22 trial. *Clin Trials* 2013;**11**:336–43. doi:10.1177/1740774513514145
- 23 33 Biglan KM, Brocht A, Raca P, *et al.* Implementing Risk-Based Monitoring (RBM) in
24 STEADY-PD III, A Phase III Multi-site Clinical Drug Trial for Parkinson Disease. *Mov*
25 *Disord* 2016;**31**:E10–E10.<http://onlinelibrary.wiley.com/doi/10.1002/mds.26749/full>
- 26 34 Knott C, Valdes-Marquez E, Landray M, *et al.* Improving efficiency of on-site
27 monitoring in multicentre clinical trials by targeting visits. *Trials* 2015;**16**:1.
28 doi:10.1186/1745-6215-16-s2-o49
- 29 35 Lindblad AS, Manukyan Z, Purohit-Sheth T, *et al.* Central site monitoring: Results from
30 a test of accuracy in identifying trials and sites failing Food and Drug Administration
31 inspection. *Clin Trials* 2014;**11**:205–
32 17.<http://ctj.sagepub.com/content/11/2/205.abstract>
- 33 36 O’Kelly M. Using statistical techniques to detect fraud: a test case. *Pharm Stat*
34 2004;**3**:237–46. doi:10.1002/pst.137
- 35 37 van den Bor RM, Vaessen PW, Oosterman BJ, *et al.* A Computationally Simply Central
36 Monitoring Procedure was Proposed and Effectively Applied to Empirical Trial data
37 with Known Fraud. *J Clin Epidemiol* 2017;**0**. doi:10.1016/j.jclinepi.2017.03.018
- 38 38 Wu X, Carlsson M. Detecting data fabrication in clinical trials from cluster analysis
39 perspective. *Pharm Stat* 2011;**10**:257–64. doi:10.1002/pst.462
- 40 39 Desmet L, Venet D, Doffagne E, *et al.* Linear mixed-effects models for central
41 statistical monitoring of multicenter clinical trials. *Stat Med* 2014;**33**:5265–79.
42 doi:10.1002/sim.6294

- 1 40 Desmet L, Venet D, Doffagne E, *et al.* Use of the Beta-Binomial Model for Central
2 Statistical Monitoring of Multicenter Clinical Trials. *Stat Biopharm Res* 2017;**9**:1–11.
3 doi:10.1080/19466315.2016.1164751
- 4 41 Pogue JM, Devereaux P, Thorlund K, *et al.* Central statistical monitoring: Detecting
5 fraud in clinical trials. *Clin Trials* 2013;**10**:225–35. doi:10.1177/1740774512469312
- 6 42 Diani CA, Rock A, Moll P. An evaluation of the effectiveness of a risk-based monitoring
7 approach implemented with clinical trials involving implantable cardiac medical
8 devices. *Clin Trials* 2017;:1740774517723589. doi:10.1177/1740774517723589
- 9 43 Agrafiotis DK, Lobanov VS, Farnum MA, *et al.* Risk-based Monitoring of Clinical Trials:
10 An Integrative Approach. *Clin Ther* 2018;**40**:1204–12.
11 doi:10.1016/j.clinthera.2018.04.020
- 12 44 Brueton V, Tierney JF, Stenning S, *et al.* Identifying additional studies for a systematic
13 review of retention strategies in randomised controlled trials: making contact with
14 trials units and trial methodologists. *Syst Rev* 2017;**6**:167. doi:10.1186/s13643-017-
15 0549-9
- 16 45 Oba K. Statistical challenges for central monitoring in clinical trials: a review. *Int J Clin*
17 *Oncol* 2016;**21**:28–37.
- 18 46 Timmermans C, Venet D, Burzykowski T. Data-driven risk identification in phase III
19 clinical trials using central statistical monitoring. *Int J Clin Oncol* 2016;**21**:38–45.
20 doi:10.1007/s10147-015-0877-5
- 21 47 Venet D, Doffagne E, Burzykowski T, *et al.* A statistical approach to central monitoring
22 of data quality in clinical trials. *Clin Trials* 2012;**9**:705–13.
23 doi:10.1177/1740774512447898
- 24 48 Buyse M, Evans SJW, Buyse M, *et al.* Fraud in Clinical Trials. In: *Encyclopedia of*
25 *Biostatistics*. Chichester, UK: : John Wiley & Sons, Ltd 2005.
26 doi:10.1002/0470011815.b2a01026
- 27 49 Peto R, Collins R, Sackett D, *et al.* The trials of Dr. Bernard Fisher: A European
28 perspective on an American episode. *Control Clin Trials* 1997;**18**:1–13.
29 doi:10.1016/S0197-2456(96)00225-5
- 30 50 Seachrist L. Scientific misconduct. NIH tightens clinical trials monitoring. *Science (80-)*
31 1994;**264**:499–499. doi:10.1126/science.8160006
- 32 51 George SL, Buyse M. Data fraud in clinical trials. *Clin Invest (Lond)* 2015;**5**:161–73.
33 doi:10.4155/CLI.14.116
- 34 52 van den Bor RM, Vaessen PWJ, Oosterman BJ, *et al.* A computationally simple central
35 monitoring procedure, effectively applied to empirical trial data with known fraud. *J*
36 *Clin Epidemiol* 2017;**87**:59–69. doi:10.1016/j.jclinepi.2017.03.018
- 37 53 Hurley C, Sinnott C, Clarke M, *et al.* Perceived barriers and facilitators to Risk Based
38 Monitoring in academic-led clinical trials: a mixed methods study. *Trials* 2017;**18**:423.
39 doi:10.1186/s13063-017-2148-4
- 40 54 Williamson PR, Altman DG, Blazeby JM, *et al.* Developing core outcome sets for
41 clinical trials: Issues to consider. *Trials* 2012;**13**:132. doi:10.1186/1745-6215-13-132

- 1 55 Ćurković M, Košec A. Bubble effect: including internet search engines in systematic
2 reviews introduces selection bias and impedes scientific reproducibility. *BMC Med*
3 *Res Methodol* 2018;**18**:130. doi:10.1186/s12874-018-0599-2
- 4 56 Arksey H, O'Malley L. Scoping studies: towards a methodological framework. *Int J Soc*
5 *Res Methodol* 2005;**8**:19–32. doi:10.1080/1364557032000119616
- 6 57 Almkhatar T, Glassman A. Monitoring of adverse events reporting in multicenter
7 clinical trials using a mixed effect regression model. In: *36th annual meeting of the*
8 *Society for Clinical Trials*. 2015. P45.
- 9 58 Atanu B, Savanur S, Suman K. Risk based monitoring in clinical trial: an application
10 with neural networking. *Biostat Biometrics* 2017;**3**:555624.
- 11 59 Bailey L, Brudenell Straw FK, George SE. Implementing a risk based monitoring
12 approach in the early phase myeloma portfolio at Leeds CTRU. *Trials* 2017;**18**(Suppl
13 **1**:220.
- 14 60 Bengtsson SKS. *Risk based monitoring in clinical studies - improving data quality*.
15 2017.<https://lup.lub.lu.se/student-papers/search/publication/8928535>
- 16 61 Dress J, U. H, C. W, *et al*. High quality risk management for clinical trials: Use the data
17 at your hands to manage risk in your clinical trials. *Clin Trials* 2011;**8**:469.
18 doi:<http://dx.doi.org/10.1177/1740774511413037>
- 19 62 Kodama A, Cavalcanti A, Buehler A, *et al*. Application of simple statistical methods to
20 evaluated possible data fabrication: An example using the act trial. *Clin Trials*
21 2010;**7**:445. doi:10.1177/1740774510374795
- 22 63 Taylor RN, McEntegart DJ, Stillman EC, *et al*. Statistical techniques to detect fraud and
23 other data irregularities in clinical questionnaire data. *Drug Inf J* 2002;**36**:115–
24 25.[http://ovidsp.ovid.com/ovidweb.cgi?T=JS&PAGE=reference&D=emed8&NEWS=N](http://ovidsp.ovid.com/ovidweb.cgi?T=JS&PAGE=reference&D=emed8&NEWS=N&AN=34208207)
25 [&AN=34208207](http://ovidsp.ovid.com/ovidweb.cgi?T=JS&PAGE=reference&D=emed8&NEWS=N&AN=34208207)
- 26 64 Timmermans C, Doffagne E, Venet D, *et al*. Statistical monitoring of data quality and
27 consistency in the Stomach Cancer Adjuvant Multi-institutional Trial Group Trial.
28 *Gastric Cancer* 2016;**19**:24–30. doi:10.1007/s10120-015-0533-9
- 29 65 Valdez-Marquez E, H.C. J, L. M. A key risk indicator approach to central statistical
30 monitoring in multicentre clinical trials: Method development in the context of an
31 ongoing large-scale randomized trial. *Trials* 2011;**12**.
32 doi:<http://dx.doi.org/10.1186/1745-6215-12-S1-A135>
- 33 66 Valdes-Marquez E, J.C. H, J. A. Central statistical monitoring in multicentre clinical
34 trials: Developing statistical approaches for analysing key risk indicators. *Trials*
35 2013;**14**:267.
- 36 67 Whitham D, Turzanski J, Bradshaw L, *et al*. Development of a standardised set of
37 metrics for monitoring site performance in multicentre randomised trials: a Delphi
38 study. *Trials* 2018;**19**:557. doi:10.1186/s13063-018-2940-9
- 39 68 Zink RC, Dmitrienko A, Dmitrienko A. Rethinking the Clinically Based Thresholds of
40 TransCelerate BioPharma for Risk-Based Monitoring. *Ther Innov Regul Sci*
41 2018;**52**:560–71. doi:10.1177/2168479017738981

1 Tables

2 Table 1: general characteristics of included studies

Characteristic	Number (total=30)	Percentage (%)
Publication year		
1996-2000	0	0%
2000-2005	2	7%
2006-2010	2	7%
2010-2015	13	43%
2016-2018	13	43%
Type of source		
Peer-reviewed paper	21	70%
Conference abstract or poster	8	27%
Thesis	1	3%
Disease setting of trial involved		
Cardiovascular disease	4	13%
Emergency medicine	1	3%
Haematology	1	3%
Infectious diseases	1	3%
Mental health	3	10%
Neurology	1	3%
Oncology	3	10%
Ophthalmology	1	3%
Renal disease	1	3%
Respiratory disease	1	3%
Not known or no specific trial involved	13	43%
Geographical setting of trials involved		
Brazil	1	3%
International	7	23%
Japan	1	3%
North America	4	13%
UK	2	7%
Not known or no specific trial involved	15	50%
Use of Investigational Medicinal Product (IMP) in involved trials		
Involves IMP	14	47%
No IMP	1	3%
Not known or no specific trial involved	15	50%
Phase of trials involved		
Phase I	0	0%
Phase II	1	3%
Phase II and III	1	3%
Phase III	9	30%
Not known or no specific trial involved	19	63%
Status of investigational medicinal product used^a		
Unlicensed	0	0%
Licensed, used outside of its licensed indication	5	17%
Licensed, used within its licensed indication	4	13%
Not known or no specific trial involved	22	73%

Focus of work^a		
Central statistical monitoring, focus on fraud or misconduct	7	23%
Central statistical monitoring, general	13	43%
Triggered monitoring	9	30%
Other method(s) for highlighting sites at risk	2	7%
Scope of work		
Description or development of method	9	30%
Some assessment of methods' effectiveness	21	70%

1

2 ^a Categories not mutually exclusive

Table 2: full listing of all included reports

First author, publication year	Type of source	Focus of work	Scope of work
Agrafiotis, 2018 ⁴³	Peer-reviewed paper	Triggered monitoring	Some assessment of methods' effectiveness
Almukhtar, 2015 ⁵⁷	Conference abstract/poster	Central statistical monitoring, general	Description or development of method
Atanu, 2017 ⁵⁸	Peer-reviewed paper	Central statistical monitoring, general	Description or development of method
Bailey, 2017 ⁵⁹	Conference abstract/poster	Triggered monitoring	Description or development of method
Bengtsson, 2017 ⁶⁰	Thesis	Central statistical monitoring, general	Some assessment of methods' effectiveness
Biglan, 2016 ³³	Conference abstract/poster	Triggered monitoring	Some assessment of methods' effectiveness
Desmet, 2014 ³⁹	Peer-reviewed paper	Central statistical monitoring, general	Some assessment of methods' effectiveness
Desmet, 2017 ⁴⁰	Peer-reviewed paper	Central statistical monitoring, general	Some assessment of methods' effectiveness
Diani, 2017 ⁴²	Peer-reviewed paper	Triggered monitoring	Some assessment of methods' effectiveness
Djali, 2010 ³¹	Peer-reviewed paper	Other method(s) for highlighting sites at risk (combines site metric scores directly to flag sites of concern)	Some assessment of methods' effectiveness
Dress, 2011 ⁶¹	Conference abstract/poster	Triggered monitoring	Description or development of method
Edwards, 2014 ³²	Peer-reviewed paper	Central statistical monitoring with triggered monitoring	Some assessment of methods' effectiveness
Kirkwood, 2013 ¹⁷	Peer-reviewed paper	Central statistical monitoring, general	Some assessment of methods' effectiveness
Knepper, 2016 ¹⁹	Peer-reviewed paper	Central statistical monitoring, focus on fraud or misconduct	Some assessment of methods' effectiveness
Knott, 2015 ³⁴	Conference abstract/poster	Central statistical monitoring, general	Some assessment of methods' effectiveness
Kodama, 2010 ⁶²	Conference abstract/poster	Central statistical monitoring, focus on fraud or misconduct	Some assessment of methods' effectiveness

First author, publication year	Type of source	Focus of work	Scope of work
Lindblad, 2014 ³⁵	Peer-reviewed paper	Central statistical monitoring, general	Some assessment of methods' effectiveness
O'Kelly, 2004 ³⁶	Peer-reviewed paper	Central statistical monitoring, focus on fraud or misconduct	Some assessment of methods' effectiveness
Pogue, 2013 ⁴¹	Peer-reviewed paper	Central statistical monitoring, focus on fraud or misconduct	Some assessment of methods' effectiveness
Smith, 2012 ³⁰	Peer-reviewed paper	Other method(s) for highlighting sites at risk (use of "statistical process control methodology" to combine per-site risk indicator scores)	Description or development of method
Stenning, 2018 ¹⁴	Peer-reviewed paper	Triggered monitoring	Some assessment of methods' effectiveness
Taylor, 2002 ⁶³	Peer-reviewed paper	Central statistical monitoring, focus on fraud or misconduct	Some assessment of methods' effectiveness
Timmermans, 2016 ⁶⁴	Peer-reviewed paper	Central statistical monitoring, general	Some assessment of methods' effectiveness
Tudur Smith, 2014 ¹⁸	Peer-reviewed paper	Triggered monitoring	Description or development of method
Valdez-Marquez, 2011 ⁶⁵	Conference abstract/poster	Central statistical monitoring, general	Description or development of method
Valdez-Marquez, 2013 ⁶⁶	Conference abstract/poster	Central statistical monitoring, general	Description or development of method
van den Bor, 2017 ⁵²	Peer-reviewed paper	Central statistical monitoring, focus on fraud or misconduct	Some assessment of methods' effectiveness
Whitham, 2018 ⁶⁷	Peer-reviewed paper	Triggered monitoring	Description or development of method
Wu, 2011 ³⁸	Peer-reviewed paper	Central statistical monitoring, focus on fraud or misconduct	Some assessment of methods' effectiveness
Zink, 2018 ⁶⁸	Peer-reviewed paper	Central statistical monitoring, general	Some assessment of methods' effectiveness

Table 3: Types of assessments and evidence presented by reports that included some assessments of their methods' effectiveness

First author	Case studies	Illustration of method(s) on data with no known issues	Assessment of methods' ability to identify simulated problem sites	Assessment of methods' ability to identify known problems in real trial data	Methods used in ongoing trial, results of on-site monitoring reported	Methods used in ongoing trial, effects reported on trial in general (e.g. in terms of cost or data quality)	Prospectively designed, controlled study to assess methods' ability to target on-site monitoring visits to most problematic sites
Agrafiotis ⁴³					X	X	
Bengtsson ⁶⁰		X					
Biglan ³³					X	X	
Desmet (2014) ³⁹	X	X	X				
Desmet (2017) ⁴⁰		X	X				
Diani ⁴²						X	
Djali ³¹	X						
Edwards ³²	X						
Kirkwood ¹⁷		X	X				
Knepper ¹⁹			X				
Knott ³⁴					X		
Kodama ⁶²		X					
Lindblad ³⁵				X			
O'Kelly ³⁶			X				
Pogue ⁴¹		X		X			
Stenning ¹⁴							X
Taylor ⁶³		X					
Timmermans ⁶⁴		X					
van den Bor ⁵²				X			
Wu ³⁸			X	X			
Zink ⁶⁸		X					
Total	3	9	6	4	3	3	1

Table 4: best reported information on methods' classification ability, where available or deducible

First author	Available information on methods' classification abilities	Definition of 'positive' centres	'True' test status: real or simulated?	Test for 'true' centre status	Sensitivity ^a	Specificity ^b	Positive Predictive Value ^c	Negative Predictive Value ^d
Biglan ³³	Partial ('true' status known for only one centre; total number of centres not known)	Not clearly defined ^e	Real	On-site monitoring	Unavailable due to limited data; report states that one 'low-risk' centre was visited and considered to be misclassified (i.e. should have been 'medium risk' or 'high risk'). However, total number of sites classified and visited (overall and within each risk category) is not known.			
Desmet (2014) ³⁹	Explored through simulation	Presence of atypical data	Simulated	Known because simulated	Dependent on simulation scenario; no specific figure given			
	Detailed information (vital signs data used as illustrative example)	Presence of atypical data	Real	Unclear ('closer inspection')	Reported: 83% (10/(10+2))	Reported: 99% (204/(204+2))	Calculated: 83% (10/(10+2))	Calculated: 99% (204/(204+2))
Desmet (2017) ⁴⁰	Explored through simulation	Presence of atypical data	Simulated	Known because simulated	Reported: dependent on simulation scenario; no specific figure given	Reported: median specificity varied from 98%-100% depending on scenario	Not reported and not possible to calculate (results of many simulations presented)	

First author	Available information on methods' classification abilities	Definition of 'positive' centres	'True' test status: real or simulated?	Test for 'true' centre status	Sensitivity ^a	Specificity ^b	Positive Predictive Value ^c	Negative Predictive Value ^d
Knepper ¹⁹	Detailed information	Presence of fabricated data	Simulated with physician input	Known because simulated	Reported: best result from 4 scenarios (study 1): 86% (6/(6+1))	Reported: best result from 4 scenarios (study 1a): 87% (148/(148+23)) ^f	Reported: best result from 4 scenarios (study 2a): 27% (3/3+8)	Reported: best result from 4 scenarios (study 1): 99% (132/132+1)
Knott ³⁴	Partial (total number of sites not reported but likely more than number whose results reported; 'true' status of any unreported centres not known)	Presence of any findings	Real	On-site monitoring	Calculated: 85% (11/(11+2))	Calculated: 88% (7/(7+1))	Calculated: 92% (11/(11+1))	Calculated: 78% (7/(7+2))
		Presence of findings 'indicative of sloppy practice' (clearer definition not reported)	Real	On-site monitoring	Calculated: 83% (10/(10+2))	Calculated: 78% (7/(7+2))	Calculated: 83% (10/(10+2))	Calculated: 78% (7/(7+2))
		Presence of serious findings (clearer definition not reported)	Real	On-site monitoring	Calculated: 100% (1/1+0)	Calculated: 45% (9/(9+11))	Calculated: 8% (1/(1+11))	Calculated: 100% (9/(9+0))
Lindblad ³⁵	Partial ('true' status known only at 21/413 centres)	Presence of serious problems	Real	Regulatory inspection	Reported: 83% (5/(5+1))	Cannot be calculated without making assumptions about the 392/413 sites with unknown 'true' status		
		Presence of minor problems	Real	Regulatory inspection	Reported: 89% (8/(8+1))			
		Presence of any problems	Real	Regulatory inspection	Reported: 87% (13/(13+2))			

First author	Available information on methods' classification abilities	Definition of 'positive' centres	'True' test status: real or simulated?	Test for 'true' centre status	Sensitivity ^a	Specificity ^b	Positive Predictive Value ^c	Negative Predictive Value ^d
O'Kelly ³⁶	Detailed information, but sample of data from trial ^g	Presence of fabricated data	Simulated with physician input	Known because simulated	Calculated: 33% (1/(1+2))	Calculated: 95% (18/(18+1))	Calculated: 50% (1/(1+1))	Calculated: 90% (18/(18+2))
Pogue ⁴¹	POISE trial data: detailed information from all sites with >= 20 randomisations	Presence of fabricated data	Real	On-site monitoring	Reported: different models and different thresholds give different pros and cons in terms of classification. Models 1, 3 and 5 all have at least some scenarios where both specificity and sensitivity > 80%. (Model 1 and 5, risk score >=7; Model 3, risk score >=5)			
	HOPE trial data: summary information from all sites with >= 20 randomisations	Presence of fabricated data	Real	On-site monitoring	N/a (no true positives)	Reported: model 1: 99% (178/(178+2))	Calculated: all models: 0% (no true positives, so any positives are false)	Calculated: all models: 100% (no true positives, so all negatives are true negatives)
Stenning ¹⁴	Partial (only sample of negative-testing sites visited, although the study design aimed to control for this)	Presence of ≥ 1 serious (Major or Critical) finding	Real	On-site monitoring	Calculated: primary analysis: 52% (37/(37+34))	Calculated: primary analysis: 62% (8/(8+5))	Reported: primary analysis: 88% (37/(37+5))	Calculated: primary analysis: 19% (8/(8+34))
					Calculated: secondary analysis excluding re-consent findings: 59% (36/(36+25))	Calculated: secondary analysis excluding re-consent findings: 74% (17/(17+6))	Reported: secondary analysis excluding re-consent findings: 86% (36/(36+6))	Calculated: secondary analysis excluding re-consent findings: 40% (17/(17+25))

First author	Available information on methods' classification abilities	Definition of 'positive' centres	'True' test status: real or simulated?	Test for 'true' centre status	Sensitivity ^a	Specificity ^b	Positive Predictive Value ^c	Negative Predictive Value ^d
					Calculated: secondary analysis excluding all consent findings: 60% (29/(29+19))	Calculated: secondary analysis excluding all consent findings: 64% (23/(23+13))	Reported: secondary analysis excluding all consent findings: 69% (29/(29+13))	Calculated: secondary analysis excluding all consent findings: 55% (23/(23+19))
van den Bor ³⁷	Partial in paper, but authors confirmed that trial implemented source data verification for all sites (personal communication)	Presence of fabricated data	Real	On-site monitoring	<p>Various situations presented, with different implications for classification ability.</p> <p>Median false positives below 10% for all scenarios, lower with higher m-constant; in various situations (combinations of specific m-constants with specific scenarios), the fraudulent centre is flagged ≥ 3 times (authors' proposed threshold) 100% of the time.</p> <p>Some scenarios have 100% highlighting of fraudulent centre and very low false positive rate - e.g. scenario 1, m=20, scenario 2, m=2-, scenario 3, m=20 (all with false positive rate of 2%)</p>			
Wu ³⁸	Partial (15/17 sites have unknown 'true' status)	Presence of fabricated data	Real	Auditing	<p>Results presented narratively via a number of scenarios.</p> <p>For 'angular clustering', fourth scenario (correlation 0.7, 3 outliers) results in sensitivity, specificity, positive and negative predictive values all $\geq 98\%$. For 'neighborhood clustering', specificity in all scenarios is $\geq 94\%$ and second scenario (variances 0.45, 3 outliers, cluster size 27) results in sensitivity, specificity, positive and negative predictive values all $\geq 50\%$.</p>			

^a Number of correctly flagged problem sites / (number of correctly flagged problem sites + sites incorrectly *not* flagged as concerning); thick border used to highlight results more than or equal to 90%.

^b Number of sites correctly flagged as not concerning / (number of sites correctly flagged as not concerning + sites incorrectly flagged as concerning); thick border used to highlight results more than or equal to 90%.

^c Number of correctly flagged problem sites / (number of correctly flagged problem sites + sites incorrectly flagged as concerning); thick border used to highlight results more than or equal to 90%.

^d Number of sites correctly flagged as not concerning / (number of sites correctly flagged as not concerning + sites incorrectly *not* flagged as concerning); thick border used to highlight results more than or equal to 90%.

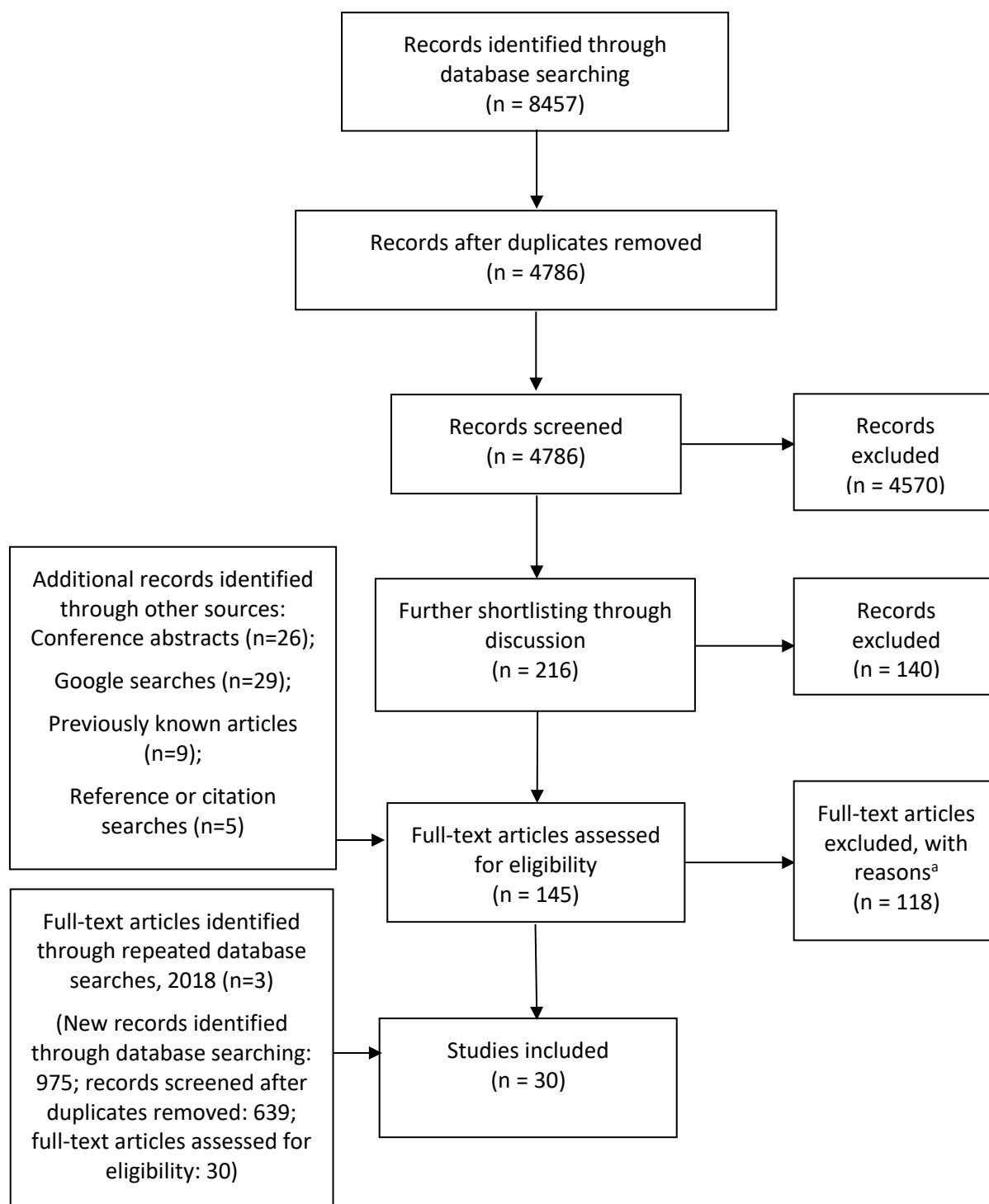
^e One 'positive' centre is described as "reveal[ing that] RBM was not assessing risk sufficiently to drive monitoring decisions"

^f Publication incorrectly rounds this to 86%.

^g Approximately one third of sites included from a trial; also some uncertainty about total number of sites (sometimes reported as 21, sometimes 22; used 22 for calculations given here as this is figure in Results section)

Figures

Figure 1 – PRISMA flow diagram



^a Reasons: no relevant methods presented (n=28); no novel methods presented (e.g. review article; n=28); method to measure variation between trial sites but no ‘flagging’ of sites of concern (n=25); abstract only and not enough detail to confirm relevance (n=10); duplicate or abstract where full paper also available (n=8); grey literature not considered to present reproducible methods (n=5); not about ‘monitoring’ according to ICH Good Clinical Practice definition (n=5); trial-level assessment only, not site-level (n=4); focus on consistency of outcome assessment only (n=4); method from observational study only, not clinical trial (n=1).

Figure 2 – Publications by year and type

