



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

# Toward a Framework for Outcome-Based Analytical Performance Specifications: A Methodology Review of Indirect Methods for Evaluating the Impact of Measurement Uncertainty on Clinical Outcomes

### Citation for published version:

Smith, AF, Shinkins, B, Hall, PS, Hulme, CT & Messenger, MP 2019, 'Toward a Framework for Outcome-Based Analytical Performance Specifications: A Methodology Review of Indirect Methods for Evaluating the Impact of Measurement Uncertainty on Clinical Outcomes', *Clinical Chemistry*, vol. 65, no. 11, pp. 1363-1374. <https://doi.org/10.1373/clinchem.2018.300954>

### Digital Object Identifier (DOI):

[10.1373/clinchem.2018.300954](https://doi.org/10.1373/clinchem.2018.300954)

### Link:

[Link to publication record in Edinburgh Research Explorer](#)

### Document Version:

Peer reviewed version

### Published In:

Clinical Chemistry

### Publisher Rights Statement:

This is a pre-copyedited, author-produced version of an article accepted for publication in clinical chemistry following peer review. The version of record "Toward a Framework for Outcome-Based Analytical Performance Specifications: A Methodology Review of Indirect Methods for Evaluating the Impact of Measurement Uncertainty on Clinical Outcomes" is available online at: [doi:10.1373/clinchem.2018.300954](https://doi.org/10.1373/clinchem.2018.300954)

### General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

### Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



1 **Towards a framework for outcome-based analytical performance**  
2 **specifications: a methodology review of indirect methods for evaluating the**  
3 **impact of measurement uncertainty on clinical outcomes**

4

5 **Authors:** Alison F. Smith<sup>1,2</sup>, Bethany Shinkins<sup>1,2</sup>, Peter S. Hall<sup>4</sup>, Claire T. Hulme<sup>1,2,5</sup>, Mike  
6 P. Messenger<sup>2,3</sup>

7

8 **Affiliations:**

9 <sup>1</sup>Test Evaluation Group, Academic Unit of Health Economics, University of Leeds, Leeds,  
10 UK

11 <sup>2</sup>NIHR Leeds In Vitro Diagnostic (IVD) Co-operative, Leeds, UK

12 <sup>3</sup>Leeds Centre for Personalised Medicine & Health, University of Leeds, Leeds, UK

13 <sup>4</sup> Cancer Research UK Edinburgh Centre, MRC Institute of Genetics & Molecular Medicine,  
14 University of Edinburgh, Edinburgh, UK

15 <sup>5</sup>Health Economics Group, University of Exeter, Exeter, UK

16 **Corresponding Author Contact Details:**

17 Alison F. Smith

18 Research Fellow

19 Test Evaluation Group, Academic Unit of Health Economics, University of Leeds, Leeds,

20 UK

21

22 **Keywords**

23 Measurement performance specifications; measurement uncertainty; analytical error;

24 evidence-based laboratory medicine

25

26 **Journal Categories**

27 Evidence-Based Laboratory Medicine and Test Utilization (TUO)

28

29 **Manuscript details:**

30 Word count: 4,342

31 Number of tables: 3

32 Number of figures: 3

33 Supplemental material: Yes

34

35 **List of Abbreviations**

36 EFLM = European Federation of Clinical Chemistry and Laboratory Medicine

37 ROC = Receiver operator characteristic

38 AUC = Area under the curve

39 CV = coefficient of variation

40 SD = standard deviation

41 EQA = External Quality Assessment

42 QALY = quality adjusted life year

CONFIDENTIAL

43 **Abstract**

44 **Background:** For medical tests that have a central role in clinical decision-making, current  
45 guidelines advocate *outcome-based* analytical performance specifications. Given that  
46 empirical (clinical-trial style) analyses are often impractical or unfeasible in this context, the  
47 ability to set such specifications is expected to rely on indirect studies to calculate the impact  
48 of test measurement uncertainty on downstream clinical, operational and economic outcomes.  
49 Currently however, a lack of awareness and guidance concerning available alternative  
50 indirect methods is limiting the production of outcome-based specifications. Our aim  
51 therefore was to review available indirect methods and present an analytical framework to  
52 inform future outcome-based performance goals.

53 **Content:** A methodology review consisting of database searches and extensive citation  
54 tracking was conducted to identify studies using indirect methods to incorporate or evaluate  
55 the impact of test measurement uncertainty on downstream outcomes (including clinical  
56 accuracy, clinical utility and/or costs). Eighty-two studies were identified, most of which  
57 evaluated the impact of imprecision and/or bias on clinical accuracy. A common analytical  
58 framework underpinning the various methods was identified, consisting of three key steps:  
59 (1) calculation of “*true*” test values; (2) calculation of *measured* test values (incorporating  
60 uncertainty); and (3) calculation of the *impact* of discrepancies between (1) and (2) on  
61 specified outcomes. A summary of the methods adopted is provided, and key considerations  
62 discussed.

63 **Conclusions:** Various approaches are available for conducting indirect assessments to  
64 inform outcome-based performance specifications. This study provides an overview of  
65 methods and key considerations to inform future studies and research in this area.

66 **Introduction**

67 Although systematic and random variation around measured test values (henceforth,  
68 *measurement uncertainty*) is now routinely documented within the clinical laboratory, the  
69 potential impact of this uncertainty on downstream clinical, operational and economic  
70 outcomes is rarely quantified. Meanwhile, evaluation of the impact of measurement  
71 uncertainty on clinical outcomes has become a recurring recommendation in protocols for  
72 determining analytical performance specifications. In their recently updated guidance, for  
73 example, the European Federation of Clinical Chemistry and Laboratory Medicine (EFLM)  
74 stipulate that, for medical tests that “have a central role in the decision-making of a specific  
75 disease or clinical situation and where cut-off/decision limits are established”, specifications  
76 should be based on the effect of analytical performance *on the clinical outcome* [termed  
77 “Model 1”], as opposed to basing specifications on biological variation [“Model 2”] or state  
78 of the art measurements [“Model 3”] (1).

79 Two types of studies are suggested to inform specifications under Model 1: (i) *direct outcome*  
80 *studies* (i.e. analyses based solely on empirical data, such as randomised controlled trials  
81 evaluating the impact of varying analytical procedures on outcomes); or (ii) *indirect outcome*  
82 *studies* (i.e. analyses using non-empirical approaches, such as decision analytic modelling, to  
83 determine the impact of varying procedures on outcomes) (2). Since (i) is often unfeasible or  
84 impractical due to ethical, financial and time constraints associated with robust end-to-end  
85 test-outcome studies, the indirect methods of (ii) are expected to play the dominant role in  
86 this context (3).

87 Despite general agreement that outcome-based specifications provide the best mechanism to  
88 ensure tests best serve patients’ needs, studies in this area remain uncommon. A primary  
89 reason often cited for this concerns the inherent difficulties in conducting direct outcomes

90 studies (1, 3). It is likely, however, that a lack of awareness and specific guidance concerning  
91 alternative *indirect* methods that may be employed is also a key limiting factor. The aim of  
92 this study therefore was to review methodological approaches used in previous indirect  
93 assessments and outline an analytical framework to inform future outcome-based  
94 performance specifications.

## 95 **Methods**

96 A literature search was conducted in November 2017 across four databases (Ovid  
97 Medline(R), Embase, Web of Science (core collection) and Biosis Citation Index) and  
98 covering a 10 year publication period (2008 to November 2017). The search was  
99 subsequently updated in 2019 (covering the period 2008 to March 2019). The search strategy  
100 (provided in the **Supplemental Appendix**) combined key terms relating to (a) tests, (b)  
101 measurement uncertainty, and (c) simulation/ methodology. From those studies identified via  
102 the database searches, subsequent citation tracking (including extensive backwards and  
103 forwards tracking) was conducted to identify additional studies published on any date (i.e.  
104 including studies published before 2008).

105 Studies were included if they met the inclusion criteria shown in **Table 1**. Studies were  
106 required to include an assessment of downstream outcomes including: clinical *accuracy* (the  
107 ability of a test to distinguish between patients with and without a specified condition, or  
108 identify a change in condition), *clinical utility* (the ability of a test to impact on healthcare  
109 management decisions or patient health outcomes) and/or *cost-effectiveness* (the ability of a  
110 test to produce an efficient impact on health outcomes in relation to cost). Note that studies  
111 using indirect methods *at any stage of the analysis* were eligible for inclusion; this means, for  
112 example, that several method-comparison studies (an essentially empirical study design) were

113 nevertheless included in cases where an indirect method was subsequently used to assess the  
114 impact of identified measurement discrepancies on outcomes.

115 <<**Table 1**>>

116 All screening (including initial title/abstract screening, full text screening, and citation  
117 tracking) was conducted by the primary reviewer (AS). A data extraction form was developed  
118 (including items on key study, test, and method details) and piloted on the first 10% of  
119 included studies. Subsequent full data extraction of included studies was conducted by the  
120 primary reviewer and double checked by one of four secondary reviewers (BS, MM, CH and  
121 PH). Regular meetings with all authors were conducted to review the ongoing study findings  
122 and resolve (via group consensus) any inclusion and/or extraction uncertainties.

CONFIDENTIAL



123 **Results**

124 **Study characteristics**

125 A total of 82 studies were identified (see **Figure 1**). Regarding data extraction checking, 35  
126 papers (43%) were checked by BS; 16 (20%) by CH; 16 (20%) by MM; and 15 (18%) by PH.  
127 Agreement between reviewers across extraction items was >99%.

128 Study characteristics are summarized in **Table 2**, and details of measurement uncertainty  
129 components and test outcomes evaluated are provided in **Table 3**. Most studies focused on  
130 evaluating tests or devices used for the purposes of monitoring, diagnosis and/or screening  
131 across four key disease areas: diabetes or glycemic control, cardiovascular diseases, cancer  
132 and metabolic or endocrine disorders. Imprecision was most commonly addressed, followed  
133 by bias and total error, and studies primarily evaluated clinical accuracy outcomes.

134 <<**Figure 1**>>

135 <<**Table 2**>>

136 <<**Table 3**>>

137 **Aim of analyses**

138 Most studies were conducted with the objective of either: (i) determining/ informing  
139 analytical performance specifications (4-22); (ii) exploring the impact of uncertainty allowed  
140 by *current* performance specifications (23-34); or (iii) evaluating the potential impact of  
141 measurement uncertainty on outcomes (without explicitly defining specifications) (35-78). A  
142 final group of studies consisted of “incidental” analyses, in which the impact of measurement  
143 uncertainty on outcomes was incorporated within the analysis but was not part of the primary  
144 study aim (79-85).

145 **Methodology Framework**

146 Based on the included studies, a common analytical framework underpinning the various  
147 approaches to evaluating the impact of measurement uncertainty on outcomes was identified.  
148 This framework consists of three key steps: (1) calculation of “*true*” test values; (2)  
149 calculation of *measured* test values (i.e. incorporating measurement uncertainty); and (3)  
150 calculation of the *impact* of discrepancies between (1) and (2) on the outcome(s) under  
151 consideration. An outline of the various methods adopted within this framework is provided  
152 below and summarized in **Figure 2**. A summary table detailing the methods used in each  
153 individual study is provided in **Supplemental Table 1**.

154 **1. Step one: calculation of “true” test values**

155 Calculation of “true” test values was based either on *empirical* data values (5, 7, 9-11, 18, 21,  
156 26, 30-32, 34-37, 39-42, 45, 49-53, 56-58, 60, 61, 64, 66-69, 71, 74, 77, 78, 85) and/or  
157 *simulated* values (4-6, 8, 12-17, 19, 20, 22-25, 27-29, 33, 36, 38, 43, 44, 46-48, 54, 55, 59,  
158 62, 63, 65, 70, 72-76, 79-84).

159 Studies using empirical data here included: (i) method comparison and external quality  
160 assessment (EQA) studies, which utilized indirect methods to determine the impact of  
161 discrepancies between empirical reference (i.e. “true”) test measurements vs. index (i.e.  
162 uncertain) test measurements on specified outcomes (e.g. using the “error grid” approach  
163 outlined in Step 3) (35, 37, 41, 42, 51, 53, 56-58, 60, 64, 66-69, 71, 75, 78); and (ii) studies  
164 which derived uncertain measurements from “true” empirical data values using various (non-  
165 empirical) approaches outlined in Step 2 (5, 7, 9-11, 18, 21, 26, 30-32, 34, 36, 39, 40, 45, 48-  
166 50, 52, 61, 77, 85).

167 Studies using simulation methods here used a range of approaches – the simplest of which  
168 was to assume a *fixed set* of individual “true” values specified across the measurement range

169 and simulate uncertainty around these values (see Step 2) (12, 16, 27, 33, 36, 38, 79, 83, 84).  
170 Whilst this approach does not require any simulation for the “true” measurements per se, the  
171 values here are nevertheless generated rather than using real-world data directly. An  
172 extension of this approach is to assume a *uniform distribution* to describe the “true”  
173 frequency distribution(s): that is, assume a constant probability of occurrence for each test  
174 value along a specified measurement range, and draw from this distribution within the  
175 simulation (14, 17, 19, 44, 55). Alternatively, the expected likelihood of test values was often  
176 modelled using *Gaussian* (i.e. normal) or *log-Gaussian* frequency distributions, specified  
177 using published or empirical data on the expected mean and variance of test values (4-6, 8,  
178 13-15, 20, 46, 47, 59, 63, 65). Other infrequently adopted parameterizations included mixed  
179 Gaussian distributions (54, 62), multivariate Gaussian distributions (where correlations  
180 between tests are known (43)) and the exponential distribution (82). Non-parametric  
181 simulation approaches were also used, based on sampling with replacement from an  
182 empirical dataset (18, 30). Finally, several studies used simulation techniques (22, 23, 70, 74,  
183 75), or utilized findings from previously published simulation studies (24, 25, 73, 76), but did  
184 not clearly report details regarding the calculation of “true” baseline values.

185 An important issue with respect to the estimation of “true” test values concerns how well the  
186 underlying data may be considered a reliable proxy for the truth. A handful of studies  
187 attempted to directly address this issue, by “stripping” known measurement uncertainty from  
188 baseline “true” test values via statistical adjustment: imprecision, for example, can be  
189 removed from the variance term of a specified Gaussian/log-Gaussian distribution using a  
190 reverse form of the “sum of squares rule”; whilst bias can be removed from the mean term (7-  
191 10, 13, 15, 31). In general, however, the likelihood that the adopted “true” test values would  
192 in fact be representative of the truth was either implicitly assumed or not discussed.

193       **2. Step two: calculation of measured test values (incorporating measurement**  
194       **uncertainty)**

195       Approaches to the calculation of measured test values predominantly fell into four broad  
196       categories: (1) *empirical assessment* (35, 37, 41, 42, 51, 53, 56-58, 60, 64, 66-69, 71, 74, 78),  
197       (2) *graphical assessment* (5, 7, 9-11, 36), (3) *computer simulation* (4-6, 8, 12, 14-25, 27-31,  
198       34, 38, 39, 44, 46, 49, 50, 52, 54, 55, 59, 61-63, 65, 70, 72-77, 79-85), or *regression analysis*  
199       (26, 32, 43, 47).

200       Studies using empirical assessment here included method-comparison studies (35, 37, 41, 42,  
201       53, 56-58, 60, 64, 66-69, 71, 75, 78) and an EQA study (51) which based “true” test values  
202       on the specified reference test and measured values on the index test measurements.

203       An alternative method, first appearing in 1980, is based on applying hypothetical  
204       measurement uncertainty to “true” values via graphical manipulation (5, 7, 9-11, 36). This  
205       approach centers on plotting the cumulative percentage frequency of “true” values on the  
206       probit scale (x-axis) as a function of “true” values on the logarithmic scale (y-axis); assuming  
207       that the log-transformed data are Gaussian, then in the bimodal case (where healthy and  
208       diseased populations are modeled separately), cumulating the healthy (diseased) population  
209       from high (low) values results in two straight lines sloping in opposite directions for each  
210       population (i.e. forming an ‘X’ on the plot). The addition of negative (positive) bias is then  
211       explored by shifting the straight lines to the left (right) on the x-axis; whilst the addition of  
212       imprecision is explored by rotating each line around their mean value (i.e. broadening the  
213       95% confidence interval of the values on the probit scale). Given a specified cut-off  
214       threshold, the proportion of false positives and negatives at a particular level of bias and  
215       imprecision can be read off directly from this plot, by observing the point at which  
216       healthy/diseased populations cross the threshold line.

217 In response to modern computational capabilities, the graphical method has been superseded  
218 by computer simulation approaches which can accommodate more complex specifications of  
219 the measurand distribution and measurement uncertainty. The most flexible and widely  
220 adopted approach in the identified studies was based on iterative simulation, with uncertainty  
221 added on to “true” test values according to a specified *error model* – a function relating  
222 measured test values to baseline “true” values plus specified components of measurement  
223 uncertainty (14, 17-19, 28-30, 34, 54, 62, 79, 82-84). This method is largely attributed to the  
224 seminal 2001 paper by Boyd and Bruns (14) – the first study of this kind to clearly specify  
225 the error model as a mathematical function (as opposed to earlier (4-6) and later (21-25, 44,  
226 49, 52, 70, 72, 73, 76, 77, 80, 81, 85) studies limited to textual descriptions or indirect  
227 referencing). An example of a typical error model is as follows:

$$228 \quad \mathbf{Test}_{measured} = \mathbf{Test}_{true} + [ \mathbf{Test}_{true} * \mathbf{N}(0,1) * \mathbf{CV} ] + \mathbf{Bias} \quad (1)$$

229 where  $\mathbf{Test}_{true}$  is the “true” measurement value;  $\mathbf{Test}_{measured}$  is the observed test value  
230 measured with imprecision (coefficient of variation [CV%]) and absolute bias (Bias); and  
231  $\mathbf{N}(0,1)$  is a normal distribution (mean = 0, standard deviation [SD] = 1) applied with the  
232 CV% value in order to produce a spread of Gaussian-distributed results around  $\mathbf{Test}_{true}$ .

233 The error model iterative simulation approach works as follows: (i) a random draw is taken  
234 from the distribution of “true” values to generate a value for  $\mathbf{Test}_{true}$ ; (ii) components of  
235 measurement uncertainty are applied to  $\mathbf{Test}_{true}$  according to the error model formula to  
236 simulate a value for  $\mathbf{Test}_{measured}$  (this may require random number draws – for example in  
237 equation (1) a random draw from  $\mathbf{N}(0,1)$  is required for the application of imprecision); (iii)  
238 points (i) and (ii) are repeated (e.g. 10,000 times to simulate 10,000  $\mathbf{Test}_{true}$  and  $\mathbf{Test}_{measured}$   
239 values) for a given level of measurement uncertainty (e.g. CV% = 5% and Bias = 5%); and  
240 (iv) points (i) to (iii) are repeated for varying levels of measurement uncertainty (e.g. CV%

241 ranging from 0-20% and Bias ranging from +/-10% in 1% increments). This iterative process  
242 can be efficiently implemented using standard statistical software, such as Excel or R.

243 Rather than iteratively adding on uncertainty via error model simulation, an alternative  
244 approach is to incorporate uncertainty directly within a specified probability distribution (e.g.  
245 incorporating bias within the mean term, and imprecision within the variance term of a  
246 Gaussian or log-Gaussian distribution). This distribution can be applied iteratively around  
247 individual “true” values (12, 16, 18, 27, 30, 38, 46, 59, 61), or at a population level, by  
248 adjusting a specified “true” population distribution to include additional uncertainty (8, 15,  
249 31, 63, 65).

250 The remaining studies used regression analysis (26, 32, 43, 47), other one-off methods (12,  
251 13, 33, 40, 45, 48), or reported insufficient details regarding simulation techniques to  
252 determine the exact method employed (74, 75). Within the identified regression analyses,  
253 bias or total error was applied as a multiplicative factor to baseline measurements within a  
254 specified regression model, with the resulting impact on the regression output (e.g. likelihood  
255 ratio) explored. Details of studies using other one-off/ indeterminate methods can be found in  
256 **Supplemental Table 1.**

### 257 **3. Step three: calculation of the impact on test outcomes**

258 The final step is to assess the impact of deviations between “true” and measured values on the  
259 outcome(s) of interest.

260 Most studies focused on evaluating clinical accuracy (4-13, 15, 16, 20, 26-29, 31-33, 38, 39,  
261 43, 45-52, 55, 59, 61-63, 65, 79-85). In this case the calculation is generally straightforward:  
262 the rate of change in mis-categorizations (e.g. false positive/negative diagnoses) is  
263 determined according to the change in the proportion of measured values pushed above or

264 below the given test cut-off threshold(s) used to define disease status or inform treatment  
265 decisions, compared to the “true” value classifications. This was the typical approach taken in  
266 studies using the graphical and simulation approaches outlined in Step 2, for example.

267 Several studies evaluated the impact of measurement uncertainty on treatment management  
268 decisions (14, 18, 21, 30, 35, 37, 41, 42, 51, 53, 56-58, 60, 64, 66-69, 71, 74, 75, 78). Most of  
269 these were method-comparison studies which determined the impact of measurement  
270 deviations on treatment decisions using error grid analysis (35, 37, 41, 42, 53, 56-58, 60, 64,  
271 66-69, 71, 74, 78). Two studies similarly employed the error grid approach, but used  
272 simulated (rather than empirical) reference and index test measurements (74, 75). First  
273 developed in the 1980s, the original error grid aimed to evaluate the potential impact of  
274 measurement discrepancies between self-monitoring blood glucose devices and laboratory  
275 reference measurements in terms of insulin dosing errors (35). Using a scatter plot of  
276 reference vs. index test measurements, the plot was divided into five error grid “zones”  
277 according to assumed severity of associated dosing errors (from zone A = clinically accurate  
278 results; to zone E = erroneous results leading to dangerous failure to detect and treat). More  
279 recently studies have attempted to build on this approach, for example by expanding on the  
280 small sample of experts used to define the initial error grid (37, 74, 75), accounting for  
281 temporal aspects of measurement (41), or applying the same methodology to alternative  
282 clinical settings (64).

283 Others have attempted to incorporate the impact of measurement uncertainty on patient health  
284 outcomes (17, 19, 22, 23, 44, 54, 70, 72). All of these studies related to evaluations of  
285 monitoring devices for glycemic control, in which health outcomes such as hypoglycemia  
286 and hyperglycemia were determined using decision analytic models based around sequential  
287 glucose measurements (incorporating measurement uncertainty via the error model  
288 simulation approach, for example). Combined with data on insulin dose administrations

289 (resulting from measured values), and additional factors such as patient insulin sensitivity and  
290 gluconeogenesis, these models were used to track patients' response to administered doses  
291 and resulting health outcomes.

292 Nine final studies included an assessment of costs or cost-effectiveness (7, 8, 11, 24, 25, 40,  
293 73, 76, 77). Four were based on a simple assignment of expected costs of misdiagnoses to  
294 rates of false positive/negative results (7, 8, 11), or expected costs of adverse events applied  
295 to simulated health outcomes data (77). One study included a more comprehensive costing  
296 analysis, in which the potential financial implications of calibration bias in serum calcium  
297 testing was explored (40). The remaining four studies all utilized the previous work of Breton  
298 and Kovatchev (2010), in which the impact of reduced glucose meter imprecision on  
299 glycemic events was simulated using a published simulation platform (23). Two studies  
300 constructed simple cost-consequence decision models, combining the Breton and Kovatchev  
301 (2010) findings with data on patient population numbers, glucose meter costs, and the rate of  
302 myocardial infarctions resulting from glycemic outcomes, to estimate annual cost savings  
303 associated with improved meter precision (73, 76). Two more recent studies conducted full  
304 cost-effectiveness analyses, using cohort Markov (i.e. state-transition) models to link the data  
305 on improved glycemic control and reduced glycemic event rates, with data on diabetes  
306 complication rates, patient health-related quality of life and health service costs (24, 25).  
307 Using these models the authors were able to estimate the incremental cost per additional  
308 quality adjusted life year (QALY) associated with reduced device error.

309 <<Figure 2>>

310



311 **Discussion**

312 **Review findings**

313 Based on our methodology review findings, a three-step analytical framework underpinning  
314 the various approaches to determining the impact of measurement uncertainty on outcomes  
315 was identified (see **Figure 2**). Key points for consideration within this framework are  
316 discussed below.

317 With regards to Step 1 (calculation of “true” test values), the primary advantage of using  
318 either empirical data or informed parametric distributions is that, by accounting for the  
319 expected frequency of values, population-level conclusions (such as analytical performance  
320 specifications) may be derived. In contrast, the primary drawback of the fixed-values  
321 approach, and by extension the uniform distribution approach (assuming this is not a realistic  
322 parameterization), is that population-level conclusions cannot be derived. Nevertheless, such  
323 approaches may be useful for exploring the impact of measurement uncertainty in specific  
324 scenarios – for example, to explore the impact of uncertainty on test values close to the test  
325 cut-off threshold.

326 A question that must be considered when using either empirical or parametric distributions, is  
327 how well the underlying data may be considered to represent the truth. If values used to  
328 inform the “true” distributions are themselves subject to measurement uncertainty (even if  
329 this uncertainty is expected to be small), then all subsequent analyses may be affected by this  
330 confounding factor and care should be taken when asserting absolute maximum bounds for  
331 imprecision and bias. A handful of studies did attempt to address this issue using statistical  
332 adjustment methods however this approach depends on having reliable information on the  
333 expected measurement uncertainty contained in the baseline “true” measurement values and  
334 can only be used when modelling test values as parametric distributions (7-10, 13, 15, 31).

335 A second consideration in the adoption of parametric distributions concerns the  
336 appropriateness of the assumed parametric form. Whilst a minority of studies provided some  
337 form of justification for the parametric choice (e.g. using the Kolmogorov–Smirnov test for  
338 normality), a common implicit assumption was that data would be likely to be Gaussian or  
339 log-Gaussian distributed. The validity of this assumption is not always clear, however.

340 Within Step 2 (calculation of *measured* test values) computer simulation methods offer the  
341 most flexible approach for exploring alternative specifications and levels of measurement  
342 uncertainty. In the context of setting performance goals, studies based on method-comparison  
343 analyses are of limited use given the fact that alternative levels of measurement uncertainty  
344 cannot be efficiently explored, and analyses using the graphical method suffer from the issue  
345 that non-Gaussian parameterisations or non-constant/ non-linear specifications of bias or  
346 imprecision cannot be accommodated. The error model approach is particularly useful in this  
347 respect. While the example formula provided in Equation (1) specifies one CV% element  
348 representing total imprecision, additional elements of imprecision (e.g. pre-analytical,  
349 analytical and biological) may be separately specified. Alternative characterisations of  
350 imprecision may also be defined: for example, using (i) a fixed SD, (ii) different SD/CV  
351 values for different sections of the measurement range, or (iii) imprecision defined as a  
352 linear/ non-linear function of  $\text{Test}_{\text{true}}$ . Similarly bias may also be characterised in alternative  
353 ways.

354 With regards to Step 3 (calculation of the impact on *outcomes*), a further advantage of the  
355 simulation approach is that, by sampling over a range of bias and imprecision values, the  
356 joint impact of these components on outcomes can be clearly explored. In particular, several  
357 studies used *contour plots* to present their findings (14-19, 21, 30, 34, 62): an example,  
358 provided in **Figure 3**, represents a hypothetical case in which bias and imprecision have been  
359 applied (according to equation (1)) to normally distributed healthy [N(30,5)] and diseased

360 [N(60,10)] populations. The plotted lines indicate at which values of imprecision and bias a  
361 given value of clinical sensitivity/specificity is maintained. For example in this case, at  
362 imprecision=0, increasing positive bias decreases clinical specificity and increases clinical  
363 sensitivity, whilst negative bias has the opposite effect. Based on this plot, we expand on the  
364 typical contour plot to show how maximum allowable bounds for imprecision and bias can be  
365 identified according to specified minimum requirements for clinical accuracy. Suppose, for  
366 example, that we require sensitivity to remain above 90% and specificity to remain above  
367 80% in order to maintain expected health utility gains. The region of acceptable analytical  
368 bias and imprecision values for this specification of clinical accuracy is illustrated by the  
369 shaded region of the contour plot – from this we can see that, if bias is zero we can tolerate  
370 up to 20% imprecision, whilst if imprecision is zero we can tolerate -8 to +6 units of absolute  
371 bias. Plots such as this one offer an effective means of highlighting acceptable bounds for  
372 measurement uncertainty.

373 <<Figure 3>>

374 Whilst most studies focused on the intermediate outcome of clinical accuracy, ideally  
375 technologies should be evaluated in terms of their influence on “end-point” outcomes i.e.  
376 health outcomes (clinical utility), operational and/or cost-effectiveness outcomes. Several of  
377 the identified studies utilized analytic decision modeling techniques to determine the impact  
378 of measurement uncertainty on health outcomes: while these all related to the context of  
379 glycemic control devices, decision models can feasibly be used to explore any clinical  
380 pathway of interest, subject to data availability. Within the field of health technology  
381 assessment, for example, decision models are routinely employed to evaluate the expected  
382 clinical utility and cost-effectiveness of novel tests, by linking data on disease prevalence and  
383 test clinical accuracy (e.g. the proportion of correct and incorrect diagnoses), with  
384 downstream data on the expected change in patient management, patient compliance to

385 treatment and treatment effectiveness (often referred to as the “linked-evidence approach”)  
386 (86-88). Although this approach is more resource- and data-intensive, and care must be taken  
387 to ensure that the model structure appropriately reflects key aspects of the clinical pathway, it  
388 nevertheless has the advantage of explicitly capturing the impact of additional parameters  
389 (e.g. treatment effectiveness) on end-point outcomes (which may not always produce  
390 expected or intuitive results) and uncertainty around the exact values of these parameters can  
391 be quantitatively characterised in the model framework (89). We identified two recent studies  
392 which utilized health-economic models to estimate the cost-effectiveness of improved  
393 analytical performance (24, 25). These studies explored a limited set of fixed imprecision  
394 levels relating to pre-existing performance specifications: future studies could extend this  
395 methodology to explore a broader range of measurement uncertainty values (e.g. by linking  
396 error-model simulations with the downstream health-economic modelling) and derive de  
397 novo performance specification based on maintaining or optimizing cost-utility and cost-  
398 effectiveness outcomes.

399 **Strengths and limitations:**

400 In light of the sustained international focus on outcome-based analytical performance  
401 specifications, it is expected that the indirect approaches outlined in this study will become  
402 increasingly important. The analytical framework presented in this study provides a useful  
403 starting point to inform future studies in this area, by clearly outlining available methods in  
404 sufficient detail to enable practical implementation, and highlighting possible advantages and  
405 limitations to consider under each approach. Whereas previous studies have provided  
406 commentaries and general reviews of various approaches to setting analytical performance  
407 specifications (3, 90, 91), this is the first methodology review to focus specifically on indirect  
408 methods for setting outcome-based performance specifications.

409 As a methodology review, the aim of this study was not to systematically identify all  
410 evidence, but rather to ensure that key examples of relevant methods were identified. While  
411 we attempted to make the database search as sensitive as possible, due to the vast volume of  
412 literature in this area we necessarily had to focus the search strategy by: (i) concentrating on  
413 terms related to in-vitro biomarkers, (ii) including a filter for simulation and methodology  
414 terms, and (iii) restricting the initial database search period to 10 years. Extensive citation  
415 tracking was additionally conducted, extending into preceding years, in order to ensure that  
416 seminal papers informing modern practices would be identified in addition to current state-of-  
417 the-art methodology. Although we believe that this two-stage strategy will have captured key  
418 methodologies, not all relevant material relating to each method will have been identified and  
419 we cannot therefore draw definitive conclusions regarding the frequency that each method  
420 has been used. Nevertheless, we believe our findings provide a valuable overview of indirect  
421 study methods and an informative starting point for future studies in this area.

422

423 **Acknowledgements**

424 The authors would like to thank the following individuals for their feedback on the project  
425 plan and/or manuscript: Christopher Hyde (Exeter, UK), Christopher Bojke (Leeds, UK),  
426 Rebecca Kift (Leeds, UK), Joy Allen (Newcastle, UK), Jon Deeks (Birmingham, UK), James  
427 Turvill (York, UK), Natalie King (Leeds, UK) and the anonymous reviewers.

428

429 **Funding**

430 Alison Smith is supported by the NIHR Doctoral Research Fellowship programme (DRF-  
431 2016-09-084). Dr Bethany Shinkins and Dr Mike Messenger are also supported by the NIHR  
432 Leeds In Vitro Diagnostics Co-operative. The views expressed are those of the author(s) and  
433 not necessarily those of the NHS, the NIHR or the Department of Health.

434 **References**

- 435 1. Ceriotti F, Fernandez-Calle P, Klee GG, Nordin G, Sandberg S, Streichert T, et al. Criteria  
436 for assigning laboratory measurands to models for analytical performance  
437 specifications defined in the 1st EFLM Strategic Conference. *Clin Chem Lab Med*  
438 2017;55:189-94.
- 439 2. Sandberg S, Fraser CG, Horvath AR, Jansen R, Jones G, Oosterhuis W, et al. Defining  
440 analytical performance specifications: consensus statement from the 1st Strategic  
441 Conference of the European Federation of Clinical Chemistry and Laboratory  
442 Medicine. *Clin Chem Lab Med* 2015;53:833-5.
- 443 3. Horvath AR, Bossuyt PM, Sandberg S, St John A, Monaghan PJ, Verhagen-Kamerbeek  
444 WD, et al. Setting analytical performance specifications based on outcome studies—is  
445 it possible? *Clin Chem Lab Med* 2015;53:841-8.
- 446 4. Groth T, Hakman M, Hällgren R, Roxin L-E, Venge P. 4.5. Diagnosis, size estimation and  
447 prediction of acute myocardial infarction from S-myoglobin observations. A system  
448 analysis to assess the influence of various sources of variability. *Scand J Clin Lab*  
449 *Invest* 1980;40:Suppl:S111-24.
- 450 5. Hørder M, Petersen PH, Groth T, Gerhardt W. 4.3. Influence of analytical quality on the  
451 diagnostic power of a single S-CK B test in patients with suspected acute myocardial  
452 infarction. *Scand J Clin Lab Invest* 1980;40:Suppl:S95-100.
- 453 6. Jacobson G, Groth T, Verdier C-HD. 4.1. Pancreatic iso-amylase in serum as a diagnostic  
454 test in different clinical situations. A simulation study. *Scand J Clin Lab Invest*  
455 1980;40:Suppl:S77-84.
- 456 7. Petersen P, Rosleff F, Rasmussen J, Hobolth N. 4.2. Studies on the required analytical  
457 quality of TSH measurements in screening for congenital hypothyroidism. *Scand J*  
458 *Clin Lab Invest* 1980;40:Suppl:S85-93.
- 459 8. Groth T, Ljunghall S, De Verdier C-H. Optimal screening for patients with  
460 hyperparathyroidism with use of serum calcium observations. A decision-theoretical  
461 analysis. *Scand J Clin Lab Invest* 1983;43:699-707.
- 462 9. Nørregaard-Hansen K, Petersen PH, Hangaard J, Simonsen E, Rasmussen O, Horder M.  
463 Early observations of S-myoglobin in the diagnosis of acute myocardial infarction.  
464 The influence of discrimination limit, analytical quality, patient's sex and prevalence  
465 of disease. *Scand J Clin Lab Invest* 1986;46:561-9.
- 466 10. Wiggers P, Dalhøj J, Petersen PH, Blaabjerg O, Hørder M. Screening for  
467 haemochromatosis: Influence of analytical imprecision, diagnostic limit and  
468 prevalence on test validity. *Scand J Clin Lab Invest* 1991;51:143-8.
- 469 11. Arends J, Petersen PH, Nørgaard-Pedersen B. 6.1. 2.3 Prenatal screening for neural tube  
470 defects, quality specification for maternal serum alphafetoprotein analysis. *Ups J Med*  
471 *Sci* 1993;98:339-47.
- 472 12. Kjeldsen J, Lassen JF, Petersen PH, Brandslund I. Biological variation of International  
473 Normalized Ratio for prothrombin times, and consequences in monitoring oral  
474 anticoagulant therapy: computer simulation of serial measurements with goal-setting  
475 for analytical quality. *Clin Chem* 1997;43:2175-82.
- 476 13. von Eyben FE, Petersen PH, Blaabjerg O, Madsen EL. Analytical quality specifications  
477 for serum lactate dehydrogenase isoenzyme 1 based on clinical goals. *Clin Chem Lab*  
478 *Med* 1999;37:553-61.
- 479 14. Boyd JC, Bruns DE. Quality specifications for glucose meters: assessment by simulation  
480 modeling of errors in insulin dose. *Clin Chem* 2001;47:209-14.

- 481 15. Petersen PH, Brandslund I, Jørgensen L, Stahl M, Olivarius NDF, Borch-Johnsen K.  
482 Evaluation of systematic and random factors in measurements of fasting plasma  
483 glucose as the basis for analytical quality specifications in the diagnosis of diabetes. 3.  
484 Impact of the new WHO and ADA recommendations on diagnosis of diabetes  
485 mellitus. *Scand J Clin Lab Invest* 2001;61:191-204.
- 486 16. Petersen PH, Jørgensen LG, Brandslund I, De Fine Olivarius N, Stahl M. Consequences  
487 of bias and imprecision in measurements of glucose and HbA1c for the diagnosis and  
488 prognosis of diabetes mellitus. *Scand J Clin Lab Invest* 2005;65:Suppl:S51-60.
- 489 17. Boyd JC, Bruns DE. Monte carlo simulation in establishing analytical quality  
490 requirements for clinical laboratory tests meeting clinical needs. *Methods Enzymol*  
491 2009;467:411-33.
- 492 18. Karon BS, Boyd JC, Klee GG. Glucose meter performance criteria for tight glycemic  
493 control estimated by simulation modeling. *Clin Chem* 2010;56:1091-7.
- 494 19. Boyd JC, Bruns DE. Effects of measurement frequency on analytical quality required for  
495 glucose measurements in intensive care units: assessments by simulation models. *Clin*  
496 *Chem* 2014;60:644-50.
- 497 20. Petersen PH, Klee GG. Influence of analytical bias and imprecision on the number of  
498 false positive results using guideline-driven medical decision limits. *Clin Chim Acta*  
499 2014;430:1-8.
- 500 21. Van Herpe T, De Moor B, Van den Berghe G, Mesotten D. Modeling of effect of glucose  
501 sensor errors on insulin dosage and glucose bolus computed by LOGIC-Insulin. *Clin*  
502 *Chem* 2014;60:1510-8.
- 503 22. Wilinska ME, Hovorka R. Glucose control in the intensive care unit by use of continuous  
504 glucose monitoring: what level of measurement error is acceptable? *Clin Chem*  
505 2014;60:1500-9.
- 506 23. Breton MD, Kovatchev BP. Impact of blood glucose self-monitoring errors on glucose  
507 variability, risk for hypoglycemia, and average glucose control in type 1 diabetes: an  
508 in silico study. *J Diabetes Sci Technol* 2010;4:562-70.
- 509 24. McQueen RB, Breton MD, Craig J, Holmes H, Whittington MD, Ott MA, Campbell JD.  
510 Economic value of improved accuracy for self-monitoring of blood glucose devices  
511 for type 1 and type 2 diabetes in England. *J Diabetes Sci Technol* 2018;12:992-1001.
- 512 25. McQueen RB, Breton MD, Ott M, Koa H, Beamer B, Campbell JD. Economic value of  
513 improved accuracy for self-monitoring of blood glucose devices for type 1 diabetes in  
514 Canada. *J Diabetes Sci Technol* 2016;10:366-77.
- 515 26. Turner MJ, Baker AB, Kam PC. Effects of systematic errors in blood pressure  
516 measurements on the diagnosis of hypertension. *Blood Press Monit* 2004;9:249-53.
- 517 27. Jorgensen LG, Petersen PH, Brandslund I. The impact of variability in the risk of disease  
518 exemplified by diagnosing diabetes mellitus based on ADA and WHO criteria as gold  
519 standard. *International Journal of Risk Assessment and Management* 2005;5:358-73.
- 520 28. Turner MJ, Irwig L, Bune AJ, Kam PC, Baker AB. Lack of sphygmomanometer  
521 calibration causes over- and under-detection of hypertension: a computer simulation  
522 study. *J Hypertens* 2006;24:1931-8.
- 523 29. Turner MJ, van Schalkwyk JM, Irwig L. Lax sphygmomanometer standard causes  
524 over-detection and under-detection of hypertension: a computer simulation study.  
525 *Blood Press Monit* 2008;13:91-9.
- 526 30. Karon BS, Boyd JC, Klee GG. Empiric validation of simulation models for estimating  
527 glucose meter performance criteria for moderate levels of glycemic control. *Diabetes*  
528 *Technol Ther* 2013;15:996-1003.



- 529 31. Kuster N, Cristol JP, Cavalier E, Bargnoux AS, Halimi JM, Froissart M, et al. Enzymatic  
530 creatinine assays allow estimation of glomerular filtration rate in stages 1 and 2  
531 chronic kidney disease using CKD-EPI equation. *Clin Chim Acta* 2014;428:89-95.
- 532 32. Åsberg A, Odsæter IH, Carlsen SM, Mikkelsen G. Using the likelihood ratio to evaluate  
533 allowable total error—an example with glycated hemoglobin (HbA1c). *Clin Chem Lab*  
534 *Med* 2015;53:1459-64.
- 535 33. Kroll MH, Garber CC, Bi C, Suffin SC. Assessing the impact of analytical error on  
536 perceived disease severity. *Arch Pathol Lab Med* 2015;139:1295-301.
- 537 34. Lyon ME, Sinha R, Lyon OA, Lyon AW. Application of a simulation model to estimate  
538 treatment error and clinical risk derived from point-of-care International Normalized  
539 Ratio device analytic performance. *J Appl Lab Med* 2017;2:25-32.
- 540 35. Clarke WL, Cox D, Gonder-Frederick LA, Carter W, Pohl SL. Evaluating clinical  
541 accuracy of systems for self-monitoring of blood glucose. *Diabetes care* 1987;10:622-  
542 8.
- 543 36. Petersen PH, de Verdier C-H, Groth T, Fraser CG, Blaabjerg O, Hørder M. The influence  
544 of analytical bias on diagnostic misclassifications. *Clin Chim Acta* 1997;260:189-206.
- 545 37. Parkes JL, Slatin SL, Pardo S, Ginsberg BH. A new consensus error grid to evaluate the  
546 clinical significance of inaccuracies in the measurement of blood glucose. *Diabetes*  
547 *care* 2000;23:1143-8.
- 548 38. Sölétormos G, Hyltoft Petersen P, Dombernowsky P. Progression criteria for cancer  
549 antigen 15.3 and carcinoembryonic antigen in metastatic breast cancer compared by  
550 computer simulation of marker data. *Clin Chem* 2000;46:939-49.
- 551 39. Rouse A, Marshall T. The extent and implications of sphygmomanometer calibration  
552 error in primary care. *J Hum Hypertens* 2001;15:587.
- 553 40. Gallaher MP, Mobley LR, Klee GG, Schryver P. The impact of calibration error in  
554 medical decision making. Washington: National Institute of Standards and  
555 Technology 2004.
- 556 41. Kovatchev BP, Gonder-Frederick LA, Cox DJ, Clarke WL. Evaluating the accuracy of  
557 continuous glucose-monitoring sensors: continuous glucose-error grid analysis  
558 illustrated by TheraSense Freestyle Navigator data. *Diabetes Care* 2004;27:1922-8.
- 559 42. Baum JM, Monhaut NM, Parker DR, Price CP. Improving the quality of self-monitoring  
560 blood glucose measurement: a study in reducing calibration errors. *Diabetes Technol*  
561 *Ther* 2006;8:347-57.
- 562 43. Nix B, Wright D, Baker A. The impact of bias in MoM values on patient risk and  
563 screening performance for Down syndrome. *Prenat Diagn* 2007;27:840-5.
- 564 44. Raine III C, Pardo S, Parkes J. Predicted blood glucose from insulin administration based  
565 on values from miscoded glucose meters. *J Diabetes Sci Technol* 2008;2:557-62.
- 566 45. Elloumi F, Hu Z, Li Y, Parker JS, Gulley ML, Amos KD, Troester MA. Systematic bias  
567 in genomic classification due to contaminating non-neoplastic tissue in breast tumor  
568 samples. *BMC Med Genomics* 2011;4:54.
- 569 46. Schlauch RS, Carney E. Are false-positive rates leading to an overestimation of noise-  
570 induced hearing loss? *J Speech Lang Hear Res* 2011;54:679-92.
- 571 47. Wright D, Abele H, Baker A, Kagan KO. Impact of bias in serum free beta-human  
572 chorionic gonadotropin and pregnancy-associated plasma protein-A multiples of the  
573 median levels on first-trimester screening for trisomy 21. *Ultrasound Obstet Gynecol*  
574 2011;38:309-13.
- 575 48. Drion I, Cobbaert C, Groenier KH, Weykamp C, Bilo HJ, Wetzels JF, Kleefstra N.  
576 Clinical evaluation of analytical variations in serum creatinine measurements: why  
577 laboratories should abandon Jaffe techniques. *BMC nephrology* 2012;13:133.

- 578 49. Jin Y, Bies R, Gastonguay MR, Stockbridge N, Gobburu J, Madabushi R.  
579 Misclassification and discordance of measured blood pressure from patient's true  
580 blood pressure in current clinical practice: a clinical trial simulation case study. *J*  
581 *Pharmacokinet Pharmacodyn* 2012;39:283-94.
- 582 50. Sarno MJ, Davis CS. Robustness of ProsVue linear slope for prognostic identification of  
583 patients at reduced risk for prostate cancer recurrence: simulation studies on effects of  
584 analytical imprecision and sampling time variation. *Clin Biochem* 2012;45:1479-84.
- 585 51. Langlois MR, Descamps OS, van der Laarse A, Weykamp C, Baum H, Pulkki K, et al.  
586 Clinical impact of direct HDLc and LDLc method bias in hypertriglyceridemia. A  
587 simulation study of the EAS-EFLM Collaborative Project Group. *Atherosclerosis*  
588 2014;233:83-90.
- 589 52. Thomas F, Signal M, Harris DL, Weston PJ, Harding JE, Shaw GM, et al. Continuous  
590 glucose monitoring in newborn infants: how do errors in calibration measurements  
591 affect detected hypoglycemia? *J Diabetes Sci Technol* 2014;8:543-50.
- 592 53. De Block CE, Gios J, Verheyen N, Manuel-y-Keenoy B, Rogiers P, Jorens PG, et al.  
593 Randomized evaluation of glycemic control in the medical intensive care unit using  
594 real-time continuous glucose monitoring (REGIMEN Trial). *Diabetes Technol Ther*  
595 2015;17:889-98.
- 596 54. Krinsley JS, Bruns DE, Boyd JC. The impact of measurement frequency on the domains  
597 of glycemic control in the critically ill-a monte carlo simulation. *J Diabetes Sci*  
598 *Technol* 2015;9:237-45.
- 599 55. Bietenbeck A. Combining medical measurements from diverse sources: experiences from  
600 clinical chemistry. *Stud Health Technol Inform* 2016;228:58-62.
- 601 56. Shinotsuka CR, Brasseur A, Fagnoul D, So T, Vincent J-L, Preiser J-C. Manual versus  
602 Automated moNitoring Accuracy of Glucose II (MANAGE II). *Crit Care*  
603 2016;20:380.
- 604 57. Sutherland HL, Reynolds T. Technical and clinical accuracy of three blood glucose meters:  
605 clinical impact assessment using error grid analysis and insulin sliding scales. *J Clin*  
606 *Pathol* 2016;69:899-905.
- 607 58. Baumstark A, Jendrike N, Pleus S, Haug C, Freckmann G. Evaluation of accuracy of six  
608 blood glucose monitoring systems and modeling of possibly related insulin dosing  
609 errors. *Diabetes Technol Ther* 2017;19:580-8.
- 610 59. Bhatt IS, Guthrie On. Analysis of audiometric notch as a noise-induced hearing loss  
611 phenotype in US youth: data from the National Health And Nutrition Examination  
612 Survey, 2005–2010. *Int J Audiol* 2017;56:392-9.
- 613 60. Bochicchio GV, Nasraway S, Moore L, Furnary A, Nohra E, Bochicchio K. Results of a  
614 multicenter prospective pivotal trial of the first inline continuous glucose monitor in  
615 critically ill patients. *J Trauma Acute Care Surg* 2017;82:1049-54.
- 616 61. Chai JH, Ma S, Heng D, Yoong J, Lim WY, Toh SA, Loh TP. Impact of analytical and  
617 biological variations on classification of diabetes using fasting plasma glucose, oral  
618 glucose tolerance test and HbA1c. *Sci Rep* 2017;7:7.
- 619 62. Lyon AW, Kavsak PA, Lyon OA, Worster A, Lyon ME. Simulation models of  
620 misclassification error for single thresholds of high-sensitivity cardiac troponin I due  
621 to assay bias and imprecision. *Clin Chem* 2017;63:585-92.
- 622 63. Chung RK, Wood AM, Sweeting MJ. Biases incurred from nonrandom repeat testing of  
623 haemoglobin levels in blood donors: selective testing and its implications. *Biom J*  
624 2019;61:454-66.
- 625 64. Saugel B, Grothe O, Nicklas JY. Error grid analysis for arterial pressure method  
626 comparison studies. *Anesth Analg* 2018;126:1177-85.

- 627 65. Rodrigues Filho BA, Farias RF, dos Anjos W. Evaluating the impact of measurement  
628 uncertainty in blood pressure measurement on hypertension diagnosis. *Blood Press*  
629 *Monit* 2018;23:141-7.
- 630 66. Piona C, Dovc K, Mutlu GY, Grad K, Gregorc P, Battelino T, Bratina N. Non-adjunctive  
631 flash glucose monitoring system use during summer-camp in children with type 1  
632 diabetes: the free-summer study. *Pediatr Diabetes* 2018;19:1285-93.
- 633 67. Hansen EA, Klee P, Dirlewanger M, Bouthors T, Elowe-Gruau E, Stoppa-Vaucher S, et  
634 al. Accuracy, satisfaction and usability of a flash glucose monitoring system among  
635 children and adolescents with type 1 diabetes attending a summer camp. *Pediatr*  
636 *Diabetes* 2018;19:1276-84.
- 637 68. Freckmann G, Link M, Pleus S, Westhoff A, Kamecke U, Haug C. Measurement  
638 performance of two continuous tissue glucose monitoring systems intended for  
639 replacement of blood glucose monitoring. *Diabetes Technol Ther* 2018;20:541-9.
- 640 69. Hughes J, Welsh JB, Bhavaraju NC, Vanslyke SJ, Balo AK. Stability, accuracy, and risk  
641 assessment of a novel subcutaneous glucose sensor. *Diabetes Technol Ther*  
642 2017;19:S21-4.
- 643 70. Breton MD, Hinzmann R, Campos-Nanez E, Riddle S, Schoemaker M, Schmelzeisen-  
644 Redeker G. Analysis of the accuracy and performance of a continuous glucose  
645 monitoring sensor prototype: an in-silico study using the UVA/PADOVA type 1  
646 diabetes simulator. *J Diabetes Sci Technol* 2017;11:545-52.
- 647 71. Aberer F, Hajnsek M, Rumpler M, Zenz S, Baumann PM, Elsayed H, et al. Evaluation of  
648 subcutaneous glucose monitoring systems under routine environmental conditions in  
649 patients with type 1 diabetes. *Diabetes, Obesity and Metabolism* 2017;19:1051-5.
- 650 72. Kovatchev BP, Patek SD, Ortiz EA, Breton MD. Assessing sensor accuracy for non-  
651 adjunct use of continuous glucose monitoring. *Diabetes Technol Ther* 2015;17:177-  
652 86.
- 653 73. Schnell O, Erbach M. Impact of a reduced error range of SMBG in insulin-treated  
654 patients in Germany. *J Diabetes Sci Technol* 2014;8:479-82.
- 655 74. Kovatchev BP, Wakeman CA, Breton MD, Kost GJ, Louie RF, Tran NK, Klonoff DC.  
656 Computing the surveillance error grid analysis: procedure and examples. *J Diabetes*  
657 *Sci Technol* 2014;8:673-84.
- 658 75. Klonoff DC, Lias C, Vigersky R, Clarke W, Parkes JL, Sacks DB, et al. The surveillance  
659 error grid. *J Diabetes Sci Technol* 2014;8:658-72.
- 660 76. Schnell O, Erbach M, Wintergerst E. Higher accuracy of self-monitoring of blood glucose  
661 in insulin-treated patients in Germany: clinical and economical aspects. *J Diabetes Sci*  
662 *Technol* 2013;7:904-12.
- 663 77. Budiman ES, Samant N, Resch A. Clinical implications and economic impact of accuracy  
664 differences among commercially available blood glucose monitoring systems. *J*  
665 *Diabetes Sci Technol* 2013;7:365-80.
- 666 78. McGarraugh GV, Clarke WL, Kovatchev BP. Comparison of the clinical information  
667 provided by the FreeStyle Navigator continuous interstitial glucose monitor versus  
668 traditional blood glucose readings. *Diabetes Technol Ther* 2010;12:365-71.
- 669 79. Petersen PH, Soletormos G, Pedersen MF, Lund F. Interpretation of increments in serial  
670 tumour biomarker concentrations depends on the distance of the baseline  
671 concentration from the cut-off. *Clin Chem Lab Med* 2011;49:303-10.
- 672 80. Hu Y, Ahmed HU, Carter T, Arumainayagam N, Lecornet E, Barzell W, et al. A biopsy  
673 simulation study to assess the accuracy of several transrectal ultrasonography  
674 (TRUS)-biopsy strategies compared with template prostate mapping biopsies in  
675 patients who have undergone radical prostatectomy. *BJU Int* 2012;110:812-20.

- 676 81. Lecornet E, Ahmed HU, Hu Y, Moore CM, Nevoux P, Barratt D, et al. The accuracy of  
677 different biopsy strategies for the detection of clinically important prostate cancer: a  
678 computer simulation. *J Urol* 2012;188:974-80.
- 679 82. McCloskey LJ, Bordash FR, Ubben KJ, Landmark JD, Stickle DF. Decreasing the cutoff  
680 for Elevated Blood Lead (EBL) can decrease the screening sensitivity for EBL. *Am J*  
681 *Clin Pathol* 2013;139:360-7.
- 682 83. Lund F, Petersen PH, Pedersen MF, Abu Hassan SO, Soletormos G. Criteria to interpret  
683 cancer biomarker increments crossing the recommended cut-off compared in a  
684 simulation model focusing on false positive signals and tumour detection time. *Clin*  
685 *Chim Acta* 2014;431:192-7.
- 686 84. Abu Hassan SO, Petersen PH, Lund F, Nielsen DL, Tuxen MK, Sölétormos G.  
687 Monitoring performance of progression assessment criteria for cancer antigen 125  
688 among patients with ovarian cancer compared by computer simulation. *Biomark Med*  
689 2015;9:911-22.
- 690 85. Lin J, Fernandez H, Shashaty MG, Negoianu D, Testani JM, Berns JS, et al. False-  
691 positive rate of AKI using consensus creatinine-based criteria. *Clin J Am Soc Nephrol*  
692 2015;10:1723-31.
- 693 86. Merlin T, Lehman S, Hiller JE, Ryan P. The “linked evidence approach” to assess  
694 medical tests: a critical analysis. *Int J Technol Assess Health Care* 2013;29:343-50.
- 695 87. Schaafsma JD, van der Graaf Y, Rinkel GJ, Buskens E. Decision analysis to complete  
696 diagnostic research by closing the gap between test characteristics and cost-  
697 effectiveness. *J Clin Epidemiol* 2009;62:1248-52.
- 698 88. Trikalinos TA, Siebert U, Lau J. Decision-analytic modeling to evaluate benefits and  
699 harms of medical tests: uses and limitations. *Med Decis Making* 2009;29:E22-E9.
- 700 89. Bilcke J, Beutels P, Brisson M, Jit M. Accounting for methodological, structural, and  
701 parameter uncertainty in decision-analytic models: a practical guide. *Med Decis*  
702 *Making* 2011;31:675-92.
- 703 90. Klee GG. Establishment of outcome-related analytic performance goals. *Clin Chem*  
704 2010;56:714-22.
- 705 91. Panteghini M, Ceriotti F, Jones G, Oosterhuis W, Plebani M, Sandberg S. Strategies to  
706 define performance specifications in laboratory medicine: 3 years on from the Milan  
707 Strategic Conference. *Clin Chem Lab Med* 2017;55:1849-56.

708

709 **Tables**710 **Table 1. Review inclusion criteria**

<b>Population</b>	Any human population with any indication
<b>Intervention</b>	In-vitro test (excluding imaging) or any kind of medical device used for the purpose of screening, diagnosis, prognosis, monitoring or predicting treatment response
<b>Comparator</b>	Any
<b>Outcomes</b>	<p>(a) Clinical accuracy e.g.</p> <ul style="list-style-type: none"> <li>- Diagnostic sensitivity and/or specificity</li> <li>- Positive/negative predictive values</li> <li>- ROC curve/ AUC analysis</li> <li>- Relative risks</li> <li>- Likelihood ratios</li> </ul> <p>(b) Clinical utility</p> <ul style="list-style-type: none"> <li>- Impact on treatment management decisions</li> <li>- Impact on patient health outcomes</li> </ul> <p>(c) Costs</p> <p>(d) Cost-effectiveness</p>
<b>Method</b>	<p>Analysis includes indirect methods (i.e. excluding purely empirical analyses) to incorporate or assess the impact of one or more components of measurement uncertainty (below) on one or more outcomes (above):</p> <ul style="list-style-type: none"> <li>- Bias (e.g. calibration or method bias)</li> <li>- Imprecision (e.g. repeatability, within-laboratory or between-laboratory imprecision)</li> <li>- Pre-analytical or analytical effects</li> <li>- Summary metrics (e.g. total error [TE] or uncertainty of measurement [<math>U_M</math>])</li> </ul>
<b>Study type</b>	Full paper relating to an original study
<b>Language</b>	Full text in English
<b>Year of publication</b>	Database search: January 2008 – March 2019 Citation tracking: any data
ROC = Receiver operator characteristic; AUC = Area under the curve	

711

712 **Table 2. Study characteristics**

	N	%
<b>Year of publication</b>		
Pre-2008 (identified via citation tracking alone)	25	30%
2008 – 2009	3	4%
2010 – 2011	7	9%
2012 – 2013	9	11%
2014 – 2015	18	22%
2016 – 2017	13	16%
2018-2019	7	9%
<b>Clinical area<sup>a</sup></b>		
Diabetes & glycemic control	43	52%
Cardiovascular diseases	17	21%
Cancer	10	12%
Metabolic & endocrine disorders	8	10%
Kidney disorders	3	4%
Prenatal screening	3	4%
Noise induced hearing loss	2	2%
<b>Role of test<sup>a</sup></b>		
Monitoring	44	54%
Diagnosis	24	29%
Screening	11	13%
Prognosis	7	9%
<sup>a</sup> Several studies included a test or tests used in multiple clinical areas or roles (hence total percentages under these categories sum to >100%).		

713

714 **Table 3. Components of measurement uncertainty included and test outcomes assessed**

	N	%
<b>Component(s) of measurement uncertainty included<sup>a</sup></b>		
<b>Imprecision:</b>		
Analytical	31	38%
Pre-analytical / combined pre-analytical and analytical	8	10%
Non-specific	11	13%
<b>Total</b>	<b>50</b>	<b>61%</b>
<b>Bias:</b>		
Analytical	18	22%
Calibration bias	9	11%
Non-specific	9	11%
Pre-analytical / combined pre-analytical and analytical	2	2%
Between-method bias	1	1%
<b>Total</b>	<b>39</b>	<b>48%</b>
<b>Total error:</b>		
Method-comparison study	18	22%
EQA study	2	2%
Other	6	7%
<b>Total</b>	<b>26</b>	<b>32%</b>
<b>Biological variation included?</b>		
Yes - included as a separate element	13	16%
Yes - combined with imprecision	5	6%
<b>Total</b>	<b>18</b>	<b>22%</b>
<b>Primary test outcome assessed<sup>a</sup></b>		
<b>Clinical accuracy</b>	45	55%
<b>Clinical utility:</b>		
Impact on treatment management	23	28%
Impact on health outcomes	13	16%
<b>Costs</b>	7	9%
<b>Cost-effectiveness</b>	2	2%
<sup>a</sup> Several studies included multiple components of measurement uncertainty or assessed multiple test outcomes (hence total percentages under these categories sum to >100%).		

715

716 **Figure captions**

717 Figure 1. PRISMA flow diagram of included studies

718 Figure 2. Summary box outlining the three-step analytical framework, primary methods  
719 identified for each step in the framework, and key questions for consideration in future  
720 analyses

721 Figure 3. Example contour plot based on simulations using the error model approach (adding  
722 increasing magnitudes of bias and imprecision onto assumed “true” measurand values). The  
723 contour lines indicate what level of clinical accuracy is achieved across the range of bias and  
724 imprecision inputs explored: varying sensitivity levels as a function of bias and imprecision  
725 are represented by the solid contour lines, whilst varying specificity levels are represented by  
726 the dashed contour lines. The grey region represents an “acceptability region” for bias and  
727 imprecision, which maintains sensitivity  $\geq 90\%$  and specificity  $\geq 80\%$ .