

1 **The Accuracy of the Patient Health Questionnaire-9 (PHQ-9) for Screening to Detect Major**
2 **Depression: an Individual Participant Data Meta-analysis**

3
4 **Authors:**

5 Brooke Levis, Andrea Benedetti, Brett D. Thombs and the DEPRESSion Screening Data
6 (DEPRESSD) Collaboration.

7
8 **DEPRESSD Collaboration:**

9 Kira E. Riehm, Nazanin Saadat, Alexander W. Levis, Marleine Azar, Danielle B. Rice, Matthew J.
10 Chiovitti, Tatiana A. Sanchez, Jill Boruff, Pim Cuijpers, Simon Gilbody, John P.A. Ioannidis, Lorie
11 A. Kloda, Dean McMillan, Scott B. Patten, Ian Shrier, Roy C. Ziegelstein, Dickens H. Akena, Bruce
12 Arroll, Liat Ayalon, Hamid R. Baradaran, Murray Baron, Charles H. Bombardier, Peter
13 Butterworth, Gregory Carter, Marcos H. Chagas, Juliana C. N. Chan, Kerrie Clover, Yeates
14 Conwell, Janneke M. de Man-van Ginkel, Jaime Delgadillo, Jesse R. Fann, Felix H. Fischer, Daniel
15 Fung, Bizu Gelaye, Felicity Goodyear-Smith, Catherine G. Greeno, Brian J. Hall, John Hambridge,
16 Patricia A. Harrison, Martin Härter, Ulrich Hegerl, Leanne Hides, Stevan E. Hobfoll, Marie
17 Hudson, Masatoshi Inagaki, Khalida Ismail, Nathalie Jetté, Mohammad E. Khamseh, Kim M. Kiely,
18 Yunxin Kwan, Shen-Ing Liu, Manote Lotrakul, Sonia R. Loureiro, Bernd Löwe, Laura Marsh,
19 Anthony McGuire, Sherina Mohd Sidik, Tiago N. Munhoz, Kumiko Muramatsu, Flávia L. Osório,
20 Vikram Patel, Brian W. Pence, Philippe Persoons, Angelo Picardi, Katrin Reuter, Alasdair G.
21 Rooney, Iná S. Santos, Juwita Shaaban, Abbey Sidebottom, Adam Simning, Lesley Stafford, Sharon
22 C. Sung, Pei Lin Lynnette Tan, Alyna Turner, Christina M. van der Feltz-Cornelis, Henk C. van

23 Weert, Paul A. Vöhringer, Jennifer White, Mary A. Whooley, Kirsty Winkley, Mitsuhiro Yamada,
24 Yuying Zhang,

25

26 **Affiliations:**

27 Lady Davis Institute for Medical Research, Jewish General Hospital and McGill University,
28 Montréal, Québec, Canada

29 Brooke Levis (doctoral student)

30 Kira E. Riehm (research assistant)

31 Nazanin Saadat (research assistant)

32 Alexander W. Levis (masters student)

33 Marleine Azar (masters student)

34 Danielle B. Rice (doctoral student)

35 Matthew J. Chiovitti (research assistant)

36 Tatiana A. Sanchez (research assistant)

37 Ian Shrier (sport medicine physician)

38 Murray Baron (rheumatologist)

39 Marie Hudson (rheumatologist)

40 Brett D. Thombs (professor)

41

42 Department of Epidemiology, Biostatistics and Occupational Health, McGill University, Montréal,
43 Québec, Canada

44 Andrea Benedetti (associate professor)

45

46 Schulich Library of Physical Sciences, Life Sciences, and Engineering, McGill University,
47 Montréal, Québec, Canada
48 Jill Boruff (associate librarian)
49
50 Department of Clinical, Neuro and Developmental Psychology, Amsterdam Public Health Research
51 Institute, Vrije Universiteit, Amsterdam, the Netherlands
52 Pim Cuijpers (professor)
53
54 Hull York Medical School and the Department of Health Sciences, University of York, Heslington,
55 York, UK
56 Simon Gilbody (professor)
57 Dean McMillan (reader)
58 Christina M. van der Feltz-Cornelis (professor)
59
60 Department of Medicine, Department of Health Research and Policy, Department of Biomedical
61 Data Science, Department of Statistics, Stanford University, Stanford, California, USA
62 John P.A. Ioannidis (professor)
63
64 Library, Concordia University, Montréal, Québec, Canada
65 Lorie A. Kloda (senior librarian)
66
67 Department of Community Health Sciences, University of Calgary, Calgary, Alberta, Canada
68 Scott Patten (professor)

69

70 Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland, USA

71 Roy C. Ziegelstein (professor)

72

73 Department of Psychiatry, Makerere University College of Health Sciences, Kampala, Uganda

74 Dickens H. Akena (psychiatrist)

75

76 Department of General Practice and Primary Health Care, University of Auckland, New Zealand

77 Bruce Arroll (professor)

78 Felicity Goodyear-Smith (professor)

79

80 Louis and Gabi Weisfeld School of Social Work, Bar Ilan University, Ramat Gan, Israel

81 Liat Ayalon (professor)

82

83 Endocrine Research Center, Institute of Endocrinology and Metabolism, Iran University of Medical

84 Sciences, Tehran, Iran

85 Hamid R. Baradaran (professor)

86 Mohammad E. Khamseh (professor)

87

88 Department of Rehabilitation Medicine, University of Washington, Seattle, Washington, USA

89 Charles H. Bombardier (professor)

90

91 Centre for Research on Ageing, Health and Wellbeing, Research School of Population Health, The
92 Australian National University, Canberra, Australia
93 Kim M. Kiely (NHMRC Fellow)
94
95 Centre for Mental Health, Melbourne School of Population and Global Health, University of
96 Melbourne, Melbourne, Australia
97 Peter Butterworth (professor)
98
99 Centre for Brain and Mental Health Research, University of Newcastle, New South Wales, Australia
100 Gregory Carter (conjoint professor)
101 Kerrie Clover (clinical psychologist)
102
103 Department of Neurosciences and Behavior, Ribeirão Preto Medical School, University of São
104 Paulo, Ribeirão Preto, Brazil
105 Marcos H. Chagas (assistant professor)
106 Sonia R. Loureiro (professor)
107 Flávia L. Osório (teacher)
108
109 Department of Medicine and Therapeutics, Prince of Wales Hospital, The Chinese University of
110 Hong Kong, Hong Kong Special Administrative Region, China
111 Juliana C. N. Chan (professor)
112 Yuying Zhang (researcher)
113

114 Psycho-Oncology Service, Calvary Mater Newcastle, New South Wales, Australia
115 Kerrie Clover (clinical psychologist)
116 Adam Simning (assistant professor)
117
118 Department of Psychiatry, University of Rochester Medical Center, Rochester, New York, USA
119 Yeates Conwell (professor)
120
121 Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, University
122 Utrecht, Utrecht, the Netherlands
123 Janneke M. de Man-van Ginkel (assistant professor)
124
125 Clinical Psychology Unit, Department of Psychology, University of Sheffield, Sheffield, UK
126 Jaime Delgadillo (lecturer in clinical psychology)
127
128 Department of Psychiatry and Behavioral Sciences, University of Washington, Seattle, Washington,
129 USA
130 Jesse R. Fann (professor)
131
132 Department of Psychosomatic Medicine, Center for Internal Medicine and Dermatology, Charité -
133 Universitätsmedizin Berlin, corporate member of Freie Universität Berlin, Humboldt-Universität zu
134 Berlin, and Berlin Institute of Health, Berlin, Germany, Germany
135 Felix H. Fischer (research fellow)
136

137 Department of Child & Adolescent Psychiatry, Institute of Mental Health, Singapore
138 Daniel Fung (associate professor)
139
140 Department of Epidemiology, Harvard T. H. Chan School of Public Health, Boston, Massachusetts,
141 USA
142 Bizu Gelaye (assistant professor)
143
144 School of Social Work, University of Pittsburgh, Pittsburgh, Pennsylvania, USA
145 Catherine G. Greeno (associate professor)
146
147 Global and Community Mental Health Research Group, Department of Psychology, Faculty of
148 Social Sciences, University of Macau, Macau Special Administrative Region, China
149 Brian J. Hall (associate professor)
150
151 Liaison Psychiatry Department, John Hunter Hospital, Newcastle, Australia
152 John Hambridge (clinical psychologist)
153
154 City of Minneapolis Health Department, Minneapolis, Minnesota, USA
155 Patricia A. Harrison (director of research and evaluation)
156
157 Department of Medical Psychology, University Medical Center Hamburg-Eppendorf, Hamburg,
158 Germany
159 Martin Härter (professor)

- 160
- 161 Department of Psychiatry and Psychotherapy, German Depression Foundation, Leipzig, Germany
- 162 Ulrich Hegerl (professor)
- 163
- 164 School of Psychology, University of Queensland, Brisbane, Queensland, Australia
- 165 Leanne Hides (professor)
- 166
- 167 STAR-Stress, Anxiety & Resilience Consultants, Chicago, Illinois, USA
- 168 Stevan E. Hobfoll (managing member)
- 169
- 170 Department of Psychiatry, Faculty of Medicine, Shimane University, Shimane, Japan
- 171 Masatoshi Inagaki (professor)
- 172
- 173 Department of Psychological Medicine, Institute of Psychiatry, Psychology and
- 174 Neurosciences, King's College London Weston Education Centre, London, UK
- 175 Khalida Ismail (professor)
- 176
- 177 Department of Neurology, Ichan School of Medicine at Mount Sinai, New York, New York, USA
- 178 Nathalie Jetté (professor)
- 179
- 180 Department of Psychological Medicine, Tan Tock Seng Hospital, Singapore
- 181 Yunxin Kwan(psychiatrist)
- 182 Pei Lin Lynnette Tan, MMed (psychiatrist)

183
184 Department of Psychiatry, Mackay Memorial Hospital, Taipei, Taiwan
185 Shen-Ing Liu (professor)
186
187 Department of Psychiatry, Faculty of Medicine, Ramathibodi Hospital, Mahidol University,
188 Bangkok, Thailand
189 Manote Lotrakul (professor)
190
191 Department of Psychosomatic Medicine and Psychotherapy, University Medical Center Hamburg-
192 Eppendorf, Hamburg, Germany
193 Bernd Löwe (professor)
194
195 Baylor College of Medicine, Houston and Michael E. DeBakey Veterans Affairs Medical Center,
196 Houston, Texas, USA
197 Laura Marsh (professor)
198
199 Department of Nursing, St. Joseph's College, Standish, Maine, USA
200 Anthony McGuire (professor)
201
202 Cancer Resource & Education Centre, and Department of Psychiatry, Faculty of Medicine and
203 Health Sciences, Universiti Putra Malaysia, Serdang, Selangor, Malaysia
204 Sherina Mohd Sidik (professor)
205

206 Post-graduate Program in Epidemiology, Federal University of Pelotas, Pelotas, RS, Brazil
207 Tiago N. Munhoz (professor)
208 Iná S. Santos (professor)
209
210 Department of Clinical Psychology, Graduate School of Niigata Seiryō University, Niigata, Japan
211 Kumiko Muramatsu (psychiatrist)
212
213 Department of Global Health and Social Medicine, Harvard Medical School, Boston,
214 Massachusetts, USA
215 Vikram Patel (professor)
216
217 Department of Epidemiology, Gillings School of Global Public Health, The University of North
218 Carolina at Chapel Hill, Chapel Hill, North Carolina, USA
219 Brian W. Pence (associate professor)
220
221 Mind-Body Research, Department of Neurosciences, Katholieke Universiteit Leuven, Leuven,
222 Belgium
223 Philippe Persoons (assistant professor)
224
225 Centre for Behavioural Sciences and Mental Health, Italian National Institute of Health, Rome, Italy
226 Angelo Picardi (senior researcher)
227
228 Group Practice for Psychotherapy and Psycho-oncology, Freiburg, Germany

229 Katrin Reuter (psychologist)
230
231 Division of Psychiatry, Royal Edinburgh Hospital, University of Edinburgh, Edinburgh, Scotland,
232 UK
233 Alasdair G. Rooney (physician)
234
235 Department of Family Medicine, School of Medical Sciences, Universiti Sains Malaysia, Kelantan,
236 Malaysia
237 Juwita Shaaban, MMed (family medicine specialist)
238
239 Allina Health, Minneapolis, Minnesota, USA
240 Abbey Sidebottom (epidemiologist)
241
242 Melbourne School of Psychological Sciences, University of Melbourne, Australia
243 Lesley Stafford (associate professor)
244
245 Programme in Health Services & Systems Research, Duke-NUS Medical School, Singapore
246 Sharon C. Sung (assistant professor)
247
248 IMPACT Strategic Research Centre, School of Medicine, Deakin University, Geelong, Victoria,
249 Australia
250 Alyna Turner (senior lecturer)
251

252 Department of General Practice, Amsterdam Institute for General Practice and Public Health,
253 Amsterdam University Medical Centers, Location AMC, Amsterdam, the Netherlands

254 Henk C. van Weert (professor)

255

256 Millennium Institute for Depression and Personality Research (MIDAP), Ministry of Economy,
257 Macul, Santiago, Chile

258 Paul A. Vöhringer (adjunct researcher)

259

260 Monash University, Melbourne, Australia

261 Jennifer White (research fellow)

262

263 Department of Medicine, Veterans Affairs Medical Center, San Francisco, California, USA

264 Mary A. Whooley (professor)

265

266 Florence Nightingale Faculty of Nursing, Midwifery & Palliative Care, King's College London,
267 Waterloo Road, London, UK

268 Kirsty Winkley (reader)

269

270 Department of Neuropsychopharmacology, National Institute of Mental Health, National Center of
271 Neurology and Psychiatry, Ogawa-Higashi, Kodaira, Tokyo, Japan

272 Mitsuhiko Yamada (director)

273

274 **Corresponding author:**

275 Brett D. Thombs, PhD; Jewish General Hospital; 4333 Cote Ste Catherine Road; Montreal, Quebec
276 H3T 1E4; Tel: (514) 340-8222 ext. 25112; E-mail: brett.thombs@mcgill.ca; ORCID: 0000-0002-
277 5644-8432
278
279 **Word count:** 4,835

280 **ABSTRACT**

281 **Objective:** The Patient Health Questionnaire-9 (PHQ-9) has been recommended for screening to
282 identify patients with depression. Conventional meta-analyses have been limited by selective cutoff
283 reporting in primary studies and have not examined accuracy for different reference standards or
284 participant subgroups. This study aimed to determine PHQ-9 accuracy for detecting major
285 depression using individual participant data meta-analysis (IPDMA).

286 **Design:** IPDMA.

287 **Data sources:** Medline, Medline In-Process & Other Non-Indexed Citations, PsycINFO, and Web
288 of Science (January 2000-February 2015).

289 **Eligibility criteria for selecting studies:** Eligible studies compared PHQ-9 scores to major
290 depression diagnoses from validated diagnostic interviews. Primary study data and study-level data
291 extracted from primary reports were synthesized. For PHQ-9 cutoffs 5 to 15, bivariate random-
292 effects meta-analysis was used to estimate pooled sensitivity and specificity, separately, among
293 studies that used semi-structured diagnostic interviews, which are designed for clinician
294 administration; fully structured interviews, which are designed for lay administration; and the Mini
295 International Neuropsychiatric (MINI) diagnostic interviews, a brief fully structured interview.
296 Sensitivity and specificity were examined among participant subgroups and, separately, using meta-
297 regression, considering all subgroup variables in a single model.

298 **Results:** Data were obtained for 58 of 72 eligible studies (N participants = 17,357, N cases = 2,312).
299 Combined sensitivity and specificity was maximized at a cutoff of ≥ 10 among studies using a semi-
300 structured interview (N = 29 studies, 6,725 participants; sensitivity [95% CI] = 0.88 [0.83 to 0.92],
301 specificity [95% CI] = 0.85 [0.82 to 0.88]). Across cutoffs 5 to 15, sensitivity with semi-structured
302 interviews was 5 to 22% higher than for fully structured interviews (MINI excluded; N = 14 studies,

303 7,680 participants) and 2 to 15% higher than for the MINI (N = 15 studies, 2,952 participants).
304 Specificity was similar across diagnostic interviews. The PHQ-9 appears to be similarly sensitive
305 but may be less specific for younger patients than for older patients; a cutoff of ≥ 10 can be used
306 regardless of age. In studies conducted in primary care using a semi-structured interview (N = 9
307 studies, 3,163 participants), sensitivity [95% CI] = 0.94 [0.88 to 0.97] and specificity [95% CI] =
308 0.88 [0.79 to 0.93]).

309 **Conclusions:** PHQ-9 sensitivity compared to semi-structured diagnostic interviews was greater than
310 in previous conventional meta-analyses that combined reference standards. A cutoff of ≥ 10
311 maximized combined sensitivity and specificity overall and for subgroups.

312 **Funding:** Canadian Institutes of Health Research (KRS-134297; PCG-155468).

313 **Registration:** PROSPERO (CRD42014010673).

314 Depression screening refers to the use of a depression screening questionnaire to identify
315 patients who may have depression, but have not been identified. When screening programs are
316 recommended, clinicians are advised to administer a depression symptom questionnaire and to use a
317 pre-identified cutoff threshold to classify patients as having positive or negative screening results.
318 Those with positive screening results can then be evaluated to determine if they have depression
319 and, if appropriate, be offered treatment.^{1,2}

320 The Patient Health Questionnaire-9 (PHQ-9)³⁻⁵ is a nine-item questionnaire designed to screen
321 for depression in primary care and other medical settings.^{6,7} The standard cutoff for screening to
322 identify possible major depression is ≥ 10 ,³⁻⁷ which was established in the first study on the PHQ-9
323 (N total = 580, N major depression = 41).^{3,5}

324 A conventional PHQ-9 meta-analysis from 2015 (N studies = 36, N participants = 21,292),⁸
325 evaluated sensitivity and specificity for cutoffs 7 to 15 by combining accuracy results for each cutoff
326 that were published in included primary studies. Pooled sensitivity for the standard cutoff of 10 was
327 0.78 (95% confidence interval [CI] = 0.70 to 0.84), and pooled specificity was 0.87 (95% CI = 0.84
328 to 0.90). Incomplete reporting of results from cutoffs other than 10 in the primary studies that were
329 included, however, resulted in cutoff ranges where sensitivity implausibly increased as cutoff scores
330 increased.⁸ This suggested possible selective cutoff reporting in some primary studies to maximize
331 accuracy.^{8,9} Additional limitations included the inability to assess differences across patient
332 subgroups, since subgroup results were not reported in primary studies; the inability to exclude
333 participants already diagnosed or being treated for depression, who would not be screened in
334 practice, but were included in many primary studies;^{10,11} and the combining of accuracy estimates
335 without differentiating between reference standards.¹² Semi-structured diagnostic interviews (e.g.,
336 Structured Clinical Interview for DSM Disorders [SCID]¹³) are intended to be conducted by

337 experienced diagnosticians and require clinical judgment. Fully structured interviews (e.g.,
338 Composite International Diagnostic Interview [CIDI]¹⁴) are fully scripted and designed to be
339 administered by lay interviewers in order to reduce the cost of employing trained clinical
340 interviewers; they are intended to achieve a high level of standardization, but may sacrifice
341 accuracy.¹⁵⁻¹⁸ The Mini International Neuropsychiatric Interview (MINI) is fully structured, but was
342 designed for very rapid administration and described by its authors as being over-inclusive as a
343 result.^{19,20} In a recent analysis, controlling for depressive symptom scores, we found that the MINI
344 classified approximately twice as many participants with major depression as other fully structured
345 interviews. Compared to semi-structured interviews, fully structured interviews (MINI excluded)
346 classified more patients with low symptom scores but fewer patients with high symptom scores as
347 having major depression.¹²

348 Individual participant data meta-analysis (IPDMA) involves a standard systematic review,
349 then synthesis of participant-level data from primary studies rather than summary results from study
350 reports.²¹ Advantages include the ability to conduct subgroup analyses not reported in primary
351 studies, the ability to report results from all relevant cutoffs from all included studies, and the ability
352 to exclude already diagnosed or treated participants who would not be screened in practice.

353 The objectives of this study were to use IPDMA to evaluate the diagnostic accuracy of the
354 PHQ-9 screening tool (1) among studies using semi-structured, fully structured (MINI excluded),
355 and MINI diagnostic interviews as reference standards, separately, with priority given to semi-
356 structured interview results; (2) among participants not currently diagnosed or receiving treatment
357 for a mental health problem; and (3) among participant subgroups based on age, sex, country human
358 development index, and recruitment setting.

359 **METHOD**

360 This IPDMA was registered in PROSPERO (CRD42014010673), a protocol was published,²²
361 and results were reported following PRISMA-DTA²³ and PRISMA-IPD²⁴ reporting guidelines.

362 **Search strategy**

363 A medical librarian searched Medline, Medline In-Process & Other Non-Indexed Citations via
364 Ovid, PsycINFO, and Web of Science (January 1, 2000 – February 7, 2015) on February 7, 2015,
365 using a peer-reviewed²⁵ search strategy (eMethods1). The search was limited to the year 2000
366 forward because the PHQ-9 was published in 2001.³ We also reviewed reference lists of relevant
367 reviews and queried contributing authors about non-published studies. Search results were uploaded
368 into RefWorks (RefWorks-COS, Bethesda, MD, USA). After de-duplication, unique citations were
369 uploaded into DistillerSR (Evidence Partners, Ottawa, Canada) for storing and tracking search
370 results.

371 **Identification of eligible studies**

372 Datasets from articles in any language were eligible for inclusion if they included diagnostic
373 classification for current Major Depressive Disorder (MDD) or Major Depressive Episode (MDE)
374 based on a validated semi-structured or fully structured interview conducted within two weeks of
375 PHQ-9 administration among participants ≥ 18 years who were not recruited from youth or
376 psychiatric settings or because they were identified as having symptoms of depression. We required
377 the diagnostic interviews and PHQ-9 to be administered within two weeks of each other because
378 Diagnostic and Statistical Manual of Mental Disorders (DSM) and International Classification of
379 Diseases (ICD) major depression diagnostic criteria specify that symptoms must have been present
380 in the last two weeks. We excluded patients from psychiatric settings or those already identified as
381 having symptoms of depression because screening is done to identify previously unrecognized
382 cases.

383 Datasets where not all participants were eligible were included if primary data allowed
384 selection of eligible participants. For defining major depression, we considered MDD or MDE
385 based on the DSM or ICD. If more than one was reported, we prioritized MDE over MDD, since
386 screening would attempt to detect depressive episodes and further interview would determine if the
387 episode is related to MDD or bipolar disorder, and DSM over ICD. Across all studies, there were 23
388 discordant diagnoses depending on classification prioritization (0.1% of participants).

389 Two investigators independently reviewed titles and abstracts for eligibility. If either deemed
390 a study potentially eligible, full-text review was done by two investigators, independently, with
391 disagreements resolved by consensus, consulting a third investigator when necessary. Translators
392 were consulted for languages other than those for which team members were fluent.

393 **Data extraction, contribution and synthesis**

394 Authors of eligible datasets were invited to contribute de-identified primary data. Country,
395 recruitment setting (non-medical, primary care, inpatient, outpatient specialty), and diagnostic
396 interview were extracted from published reports by two investigators independently, with
397 disagreements resolved by consensus. Countries were categorized as “very high”, “high”, or “low-
398 medium” development based on the United Nation’s human development index.²⁶ Participant-level
399 data included age, sex, major depression status, current mental health diagnosis or treatment, and
400 PHQ-9 scores. In two primary studies, multiple recruitment settings were included; thus recruitment
401 setting was coded at the participant-level. When datasets included statistical weights to reflect
402 sampling procedures, we used provided weights. For studies where sampling procedures merited
403 weighting, but the original study did not weight, we constructed weights using inverse selection
404 probabilities. Weighting occurred, for instance, when all participants with positive screens and a
405 random subset of participants with negative screens were administered a diagnostic interview.

406 Individual participant data were converted to a standard format and synthesized into a single
407 dataset with study-level data. We compared published participant characteristics and diagnostic
408 accuracy results with results from raw datasets and resolved any discrepancies in consultation with
409 the original investigators.

410 Two investigators assessed risk of bias of included studies independently, based on the
411 primary publications, using the Quality Assessment of Diagnostic Accuracy Studies-2 tool
412 (QUADAS-2; eMethods2).²⁷ Discrepancies were resolved by consensus.

413 **Statistical analyses**

414 We conducted three main sets of analyses. First, we estimated sensitivity and
415 specificity across PHQ-9 cutoffs 5 to 15 for studies with semi-structured (SCID¹³,
416 Schedules for Clinical Assessment in Neuropsychiatry²⁸, Depression Interview and
417 Structured Hamilton²⁹), fully structured (MINI excluded; CIDI¹⁴, Clinical Interview
418 Schedule-Revised³⁰, Diagnostic Interview Schedule³¹), and MINI^{19,20} reference standards,
419 separately. Second, for each reference standard category, we estimated sensitivity and
420 specificity across PHQ-9 cutoffs for all participants from primary studies, as has been done
421 in existing conventional meta-analyses and, separately, among only participants who could
422 be confirmed as not currently diagnosed or receiving treatment for a mental health problem
423 at the time of assessment. This was done because existing conventional meta-analyses have
424 all been based on primary studies that generally do not exclude patients already diagnosed
425 or receiving treatment. Since screening is done to identify previously unrecognized cases,
426 however, those patients would not be screened in practice, and their inclusion in diagnostic
427 accuracy studies could bias results.^{10,11} Third, for each reference standard category, we
428 estimated and compared sensitivity and specificity across PHQ-9 cutoffs among subgroups

429 based on age (<60 versus \geq 60 years), sex, country human development index (very high;
430 high; low-medium), and recruitment setting (non-medical; primary; inpatient specialty;
431 outpatient specialty). Among studies that used the MINI, we combined inpatient and
432 outpatient specialty care settings, as only one study included inpatient participants. In each
433 subgroup analysis, we excluded primary studies with no major depression cases, as this did
434 not allow application of the bivariate random effects model. This resulted in a maximum of
435 15 participants excluded from any subgroup analysis.

436 For each meta-analysis, for cutoffs 5 to 15 separately, bivariate random-effects models were
437 fitted via Gauss-Hermite adaptive quadrature.³² This 2-stage meta-analytic approach models
438 sensitivity and specificity simultaneously, accounting for the inherent correlation between them and
439 for precision of estimates within studies. For each analysis, this model provided estimates of pooled
440 sensitivity and specificity.

441 To compare results across reference standards and other subgroups, we constructed empirical
442 receiver operating characteristic (ROC) curves for each group based on the pooled sensitivity and
443 specificity estimates and calculated areas under the curve (AUC). We estimated differences in
444 sensitivity and specificity between subgroups at each cutoff by constructing confidence intervals for
445 differences via the cluster bootstrap approach,^{33,34} resampling at study and subject levels. For each
446 comparison, we ran 1000 iterations of the bootstrap. We removed iterations that did not produce
447 difference estimates for cutoffs 5 to 15 prior to determining confidence intervals and noted the
448 number of iterations removed.

449 In addition to categorical subgroup analyses, we compared sensitivity and specificity across
450 the different reference standards by conducting one-stage meta-regressions with interactions
451 between reference standard category (reference category = semi-structured interviews) and accuracy

452 coefficients (logit(sensitivity) and logit(specificity)), and we compared results to those seen in the
453 original two-stage bivariate random effects meta-analytic models. Additionally, within each
454 reference standard category, we conducted one-stage meta-regressions where we interacted all
455 subgrouping variables (age [measured continuously], sex [reference category = women], country
456 human development index [reference category = very high] and participant recruitment setting
457 [reference category = primary care]) with logit(sensitivity) and logit(specificity). Similar to our
458 main subgroup analyses, we once again determined which significant interactions replicated across
459 all three reference standard categories. For subgrouping variables that were significantly associated
460 with sensitivity or specificity coefficients for all three reference standard categories for all or most
461 cutoffs in the main one-stage meta-regression, we conducted additional one-stage meta-regressions
462 to produce accuracy estimates for the subgroups of interest, and we compared these results to those
463 seen in the original two-stage bivariate random effects meta-analytic models. Although age was
464 included as a continuous variable in the main meta-regression, we again dichotomized it (<60
465 versus ≥ 60 years) to estimate accuracy and compare to the bivariate model results.

466 To investigate heterogeneity, we generated forest plots of sensitivities and specificities for
467 cutoff 10 for each study, first for all studies in each reference standard category, and then separately
468 across participant subgroups within each reference standard category. We quantified cutoff 10
469 heterogeneity overall and across subgroups, by reporting estimated variances of the random effects
470 for sensitivity and specificity (τ^2) and estimating R, the ratio of the estimated standard deviation of
471 the pooled sensitivity (or specificity) from the random-effects model to that from the corresponding
472 fixed-effects model.³⁵ We used a complete case analysis since complete data for all subgrouping
473 variables were available for 17,357 participants (98% of eligible participants in the database).

474 To estimate positive and negative predictive values using cutoff 10 for different major
475 depression prevalence values, we generated nomograms for each reference standard category by
476 applying the cutoff 10 sensitivity and specificity estimates from the meta-analysis to hypothetical
477 major depression prevalence values of 5 to 25%.

478 In sensitivity analyses, for each reference standard category, we compared accuracy results
479 across subgroups based on QUADAS-2 items for all items with at least 100 major depression cases
480 among participants categorized as having “low” risk of bias and among participants with “high” or
481 “unclear” risk of bias.

482 We did not conduct sensitivity analyses that combined IPDMA accuracy results with
483 published results from studies that did not contribute IPD because among the 14 eligible studies that
484 did not contribute IPD, only two studies with a semi-structured reference standard (N total = 173, N
485 major depression = 29), one study with a fully structured reference standard (N total = 730, N MDD
486 = 32), and one study using the MINI (N total = 172, N MDD = 33) published accuracy results
487 eligible for the present IPDMA. The other studies had eligible datasets, but did not publish eligible
488 diagnostic accuracy results (eTable1b).

489 All analyses were run in R (R version R 3.4.1 and R Studio version 1.0.143) using the glmer
490 function within the lme4 package, which uses one quadrature point.

491 The only substantive deviations from our initial protocol were that we stratified accuracy
492 results by reference standard category and did not conduct sensitivity analyses that combined
493 IPDMA accuracy results with published results from studies that did not contribute IPD.

494 **Patient and public involvement**

495 Patients and members of the public were not involved in the study.

496 **RESULTS**

497 **Search results and inclusion of primary datasets**

498 Of 5,248 unique titles and abstracts identified from the database search, 5,039 were
499 excluded after title and abstract review and 113 after full-text review, leaving 96 eligible
500 articles with data from 69 unique participant samples, of which 55 (80%) contributed
501 datasets (eFigure1). Reasons for exclusion for the 113 articles excluded at full-text level
502 are provided in eTable1. In addition, authors of included studies contributed data from
503 three unpublished studies, for a total of 58 datasets (N participants = 17,357, N major
504 depression = 2,312 [13%]). Study characteristics of included studies and eligible studies
505 that did not provide datasets are shown in eTable2a and eTable2b. Excluding the three
506 unpublished studies, of 21,171 participants in 69 eligible published studies, 16,956
507 participants (80%) from 55 included published studies were included.

508 Of 58 included studies, 29 used semi-structured reference standards, 14 used fully structured
509 reference standards, and 15 used the MINI (Table 1). The SCID was the most common semi-
510 structured interview (26 studies, 4,733 participants), and the CIDI was the most common fully
511 structured interview (11 studies, 6,272 participants). Among studies that used semi-structured, fully
512 structured, and MINI diagnostic interviews, mean sample sizes were 232, 549, and 197, and mean
513 number (%) with major depression were 32 (14%), 60 (11%), and 37 (19%; Table 2).

514 **PHQ-9 accuracy by reference standard**

515 Comparisons of sensitivity and specificity estimates by reference standard category are shown
516 in Table 3. Cutoff 10 maximized combined sensitivity and specificity among studies using semi-
517 structured interviews (sensitivity [95% CI] = 0.88 [0.83 to 0.92], specificity [95% CI] = 0.85 [0.82
518 to 0.88]). Based on cutoff 10, sensitivity and specificity [95% CI] were 0.70 [0.59 to 0.80] and 0.84
519 [0.77 to 0.89] for fully structured interviews, and 0.77 [0.68 to 0.83] and 0.87 [0.83 to 0.91]) for the

520 MINI. Across cutoffs, specificity estimates were similar across reference standards; however,
521 sensitivity estimates for semi-structured interviews were 5 to 22% higher than for fully structured
522 interviews (median difference = 18%, at cutoff 10) and 2 to 15% higher than for the MINI (median
523 difference = 11%, at cutoff 10). ROC curves and AUC values are shown in eFigure2.

524 Heterogeneity analyses suggested moderate heterogeneity across studies, which improved in
525 some instances when subgroups were considered. Cutoff 10 sensitivity and specificity forest plots
526 are shown in eFigure3, with τ^2 and R values shown in eTable3.

527 Nomograms of positive and negative predictive values for cutoff 10 for each reference
528 standard category are shown in Figure 1. For hypothetical major depression prevalence values of 5
529 to 25%, estimates of positive predictive values based on summary sensitivity and specificity values
530 ranged from 24 to 66% for semi-structured interviews, 19 to 59% for fully structured interviews,
531 and 24 to 66% for the MINI; estimates of negative predictive values ranged from 96 to 99% for
532 semi-structured interviews, 89 to 98% for fully structured interviews, and 92 to 99% for the MINI.

533 When examined with meta-regression analysis, consistent with our main results, we found that
534 PHQ-9 sensitivity estimates for semi-structured interviews were significantly higher than for fully
535 structured interviews or the MINI (eTable4). The significant interactions corresponded to
536 differences in sensitivity that across cutoffs were 4 to 22% higher for semi-structured interviews
537 than for fully structured interviews (median = 18%) and 1 to 16% higher for semi-structured
538 interviews than the MINI (median = 11%). Across all cutoffs, the magnitude of the differences
539 estimated based on meta-regression were within 1% of those estimated using the original two-stage
540 bivariate random effects meta-analytic models.

541 **PHQ-9 accuracy among participants not diagnosed or receiving treatment for a mental health**
542 **problem compared to all participants**

543 Sensitivity and specificity estimates were not statistically significantly different for any
544 reference standard category when restricted to participants not currently diagnosed or receiving
545 treatment for a mental health problem compared to all participants. See eTable5 for results and
546 eFigure4 for ROC curves and AUC values.

547 **PHQ-9 accuracy among subgroups**

548 For each reference standard category, comparisons of sensitivity and specificity estimates
549 based on bivariate models across PHQ-9 cutoffs 5 to 15 among subgroups based on age, sex,
550 country human development index and participant recruitment setting are shown in eTable5, with
551 forest plots shown in eFigure3, ROC curves and AUC values shown in eFigure4, and τ^2 and R
552 values shown in eTable3.

553 Of the total of 484 categorical subgroup analyses that were done (22 subgroups x 11 cutoff
554 thresholds for sensitivity and specificity) using the bivariate model, 4 comparisons excluded the null
555 value of zero difference for cutoffs 7 to 15. No comparisons that were significantly different in one
556 reference standard category were statistically significant in either of the other two reference
557 standard categories. Subgroup analyses are shown in eTable5.

558 In the meta-regression analyses, on the other hand, older age (measured continuously) was
559 associated with higher specificity for all reference standards (eTable4). The significant interaction
560 corresponded to specificity estimates that were 2 to 14% higher for participants aged ≥ 60 versus
561 < 60 among participants based on semi-structured interviews (median = 6%), 2 to 14% based on
562 fully structured interviews (median = 8%), and 1 to 8% based on the MINI (median = 5%; eTable4).
563 Across all cutoffs, the magnitudes of the differences estimated based on meta-regression with
564 dichotomous age were within 2% of those estimated using the original two-stage bivariate random
565 effects meta-analytic models.

566 **Risk of bias sensitivity analyses**

567 eTable6 shows QUADAS-2 ratings for each included primary study, while comparisons of
568 PHQ-9 accuracy across individual items for each reference standard category are shown in eTable5.
569 For the item on blinding of the reference standard to PHQ-9 results, specificity was significantly
570 greater for studies and participants with high or unclear vs. low risk of bias for semi-structured
571 interviews, but significantly greater for low vs. high or unclear risk of bias for fully structured
572 interviews and the MINI. For the item on recruiting a consecutive or random sample of participants,
573 specificity was significantly greater for low vs. high or unclear risk of bias for fully structured
574 interviews and the MINI. No other statistically significant differences were found, and no
575 significant differences replicated across all reference standards.

576 **DISCUSSION**

577 **Principal findings**

578 We compared the accuracy of scores on the PHQ-9 to detect major depression, separately, to
579 semi-structured diagnostic interviews, fully structured diagnostic interviews (MINI excluded), and
580 the MINI. Based on results from the semi-structured interviews, which most closely replicate
581 clinical interviews done by trained professionals, the PHQ-9 was more sensitive than has been
582 reported in previous meta-analyses that combined reference standards.^{8,36} Specificity was similar to
583 previous studies and across reference standards. Based on semi-structured interviews, the standard
584 cutoff of 10 maximized combined sensitivity and specificity. There was evidence from
585 multivariable meta-regression that the PHQ-9 may be more sensitive among older patients
586 compared to younger patients, but this would not require that a different cutoff be used. Results did
587 not differ depending on whether studies that did not explicitly exclude already diagnosed patients

588 were included or excluded. Among studies conducted in primary care settings, approximately half
589 of patients who screen positive on the PHQ-9 had major depression.

590 **Findings in context**

591 This is the first meta-analysis that has analyzed diagnostic accuracy for the PHQ-9 separately
592 for different diagnostic interviews. Diagnostic interviews that are used to classify major depression
593 case status are imperfect reference standards. Semi-structured interviews, such as the SCID,¹³ most
594 closely approximate an expert diagnosis. They are set up to replicate a guided diagnostic
595 conversation with standardized questions, but the option for interviewers to make additional queries
596 and use clinical judgment to decide whether symptoms are present.^{16,17} Semi-structured interviews
597 involve lengthy processes that must be conducted by skilled diagnosticians and, thus, are expensive.
598 Fully-structured interviews, such as the CIDI,¹⁴ are designed to replicate as closely as possible
599 expert-administered semi-structured interviews, but are not expected to have the same level of
600 validity and reliability. Fully structured interviews can be administered by lay interviewers and
601 involve fully scripted standardized interview protocols that are read verbatim without additional
602 probes or interpretation. Fully structured interviews are designed to increase reliability with
603 administration by lay interviewers who are not trained to independently carry out diagnostic
604 interviews at the possible cost of validity.^{16,17} The MINI is a specific fully structured interview that
605 was designed to be administered in a fraction of the time compared to other interviews and
606 described by its developers as intentionally over-inclusive.^{19,20} Test-retest reliability for diagnosis of
607 current major depression has been reported to be $\kappa = 0.74$ for the SCID (N = 51; mean = 9
608 days)³⁷ and $\kappa = 0.52$ for the CIDI (N = 60, mean = 2 days).³⁸

609 Consistent with the design features and rigour of each type of diagnostic interview, we
610 previously reported that compared to semi-structured interviews, fully structured interviews

611 (excluding the MINI) classify more people with low symptoms as having major depression but
612 fewer people with high symptoms.¹² We also found that the MINI identified approximately twice as
613 many cases as other fully structured interviews.¹² The finding in the present study that sensitivity
614 was greater among studies with semi-structured rather than fully structured reference standards is
615 consistent with both the design features and rigor of the different types of diagnostic interviews and
616 with our previous findings. It is possible that the lower sensitivity among fully structured interviews
617 may have been due to overdiagnosis of major depression among participants with low depressive
618 symptom levels when fully structured interviews were used. In the present meta-analysis, most
619 participants did not have major depression (87%), thus misclassification of major depression among
620 participants with sub-threshold depressive symptom levels based on fully structured interviews
621 might explain the lower sensitivity compared to semi-structured interviews if the PHQ-9 were less
622 likely to identify “false positive” classifications based on fully structured interviews. The same logic
623 would apply to the lower sensitivity for the MINI.

624 Among studies that used semi-structured reference standards, sensitivity was also greater than
625 reported in previous traditional meta-analyses, where studies with semi- and fully structured
626 reference standards and the MINI were combined without adjustment. Using IPD data from the 29
627 studies that used a semi-structured interview as the reference standard, we found that at cutoff 10,
628 sensitivity and specificity were 0.88 and 0.85 compared to 0.78 and 0.87 in a 2015 conventional
629 meta-analysis of 34 studies that combined reference standards.⁸ In primary care settings, we found
630 sensitivity and specificity of 0.94 and 0.88 (9 studies with a semi-structured interview) compared to
631 0.82 and 0.85 in a 2016 conventional meta-analysis of 20 studies that combined reference
632 standards.³⁶

633 For semi-structured interviews, major depression prevalence in our dataset was 14%. Using
634 our cutoff 10 accuracy estimates (sensitivity = 0.88, specificity = 0.85), positive predictive value
635 would only be 49%; thus 51% of all positive screens would be false positives. For primary care
636 settings, where accuracy was even higher, major depression prevalence was 12%. Using our
637 accuracy estimates for cutoff 10 (sensitivity = 0.94, specificity = 0.88, positive predictive value =
638 52%), 22% of patients in primary care would screen positive at this cutoff, but only approximately
639 half would be true positives.

640 **Clinical implications**

641 Screening for depression in primary care is recommended in the United States,³⁹ but national
642 guidelines from Canada and the United Kingdom recommend against routine depression
643 screening.^{40,41} Those guidelines cite the lack of evidence of benefit from well-conducted randomised
644 controlled trials, as well as concerns about high false positive rates, overdiagnosis, and substantial
645 resource utilization and opportunity costs.⁴⁰⁻⁴¹ Well-conducted and adequately powered trials
646 designed specifically to assess the effects of depression screening are needed.^{1,2,40-43} If screening is
647 to be done clinically based on recommendations in the United States, the cutoff that maximizes
648 sensitivity and specificity is the standard cutoff of 10 or greater. It is not known, however, if using
649 this standard cutoff would maximize the likelihood that screening would successfully improve
650 mental health and minimize unnecessary resource use and adverse outcomes if tested in a trial.
651 Ideally, robust trials that are sufficiently powered to evaluate the effects of screening across a range
652 of cutoffs will be conducted. Clinical trials provide the best possible evidence to inform both the
653 decision on whether or not depression screening should be implemented as part of routine care and,
654 if so, on thresholds for intervening or what steps might be taken for patients with borderline
655 screening results.⁴⁴

656 **Strengths and limitations**

657 This was the first study to use IPDMA to assess diagnostic accuracy of the PHQ-9 or any
658 other depression screening tool. Strengths include the large sample size, the ability to include results
659 from all cutoffs from all studies (rather than just those published), the ability to examine participant
660 subgroups, and the ability to assess accuracy separately across reference standards, which had not
661 been done previously. There are also limitations to consider. First, we were unable to include
662 primary data from 14 of 69 published eligible datasets (20% of eligible datasets and participants),
663 and we restricted our analyses to those with complete data for all variables used in our various
664 analyses (98% of available data). Nonetheless, for all cutoffs other than 10, our sample was much
665 larger than previous traditional meta-analyses of the PHQ-9. Second, despite the large sample size,
666 there was substantial heterogeneity across studies, although it did improve in some instances when
667 subgroups were considered. We were not able to conduct subgroup analyses based on specific
668 medical comorbidities or cultural aspects such as country or language because comorbidity data
669 were not available for over half of participants, and many countries and languages were represented
670 in few primary studies. However, we were able to compare participant subgroups based on age, sex,
671 country human development index, and participant recruitment setting category, which has not been
672 done previously. Third, while we categorized studies based on the diagnostic interview
673 administered, interviews are sometimes adapted and thus not always used in the way that they were
674 originally designed. Although we coded for interviewer qualification for all semi-structured
675 interviews as part of our QUADAS-2 rating, two studies used interviewers who did not meet typical
676 standards, and approximately half of studies were rated as unclear on this item.

677 Although our original two-stage bivariate random effects meta-analytic models did not find
678 significant differences in accuracy estimates across participant subgroups, our meta-regressions

679 suggested that specificity might be somewhat higher among older participants whether measured
680 continuously or dichotomously. This difference in significance may be due to the differences
681 between the analytical approaches. Whereas statistical significance of the interactions between
682 covariates and accuracy estimates in the meta-regressions were based on parametric standard errors,
683 statistical significance of subgroup comparisons in the two-stage bivariate random effects models
684 was based on non-parametric bootstrap methods. Moreover, whereas the meta-regression models
685 provide a within-study interpretation, the two-stage bivariate random effects models did not link
686 study clusters across subgroups and thus focused more on between-study comparisons.

687 **Conclusions and policy implications**

688 In summary, we found that PHQ-9 sensitivity compared to semi-structured reference
689 standards was substantially greater than when compared to fully structured reference standards or
690 the MINI. It was also substantially higher than previously reported in conventional meta-analyses
691 which combined reference standards.^{8,36} The standard cutoff of 10 or greater maximized combined
692 sensitivity and specificity. However, in primary care, approximately half of patients with positive
693 screens would be false positives if used in practice, a concern that has been emphasized by the
694 Canadian Task Force on Preventive Health Care, UK National Screening Committee, and UK
695 National Institute for Health and Care Excellence, given the resources that would be required for
696 additional assessment and the possibility that some of these patients might be treated without
697 benefit.^{40,41,43} Future research on the PHQ-9 should ideally be based on semi-structured diagnostic
698 interviews, should consider estimating probabilities of depression across the full spectrum of PHQ-9
699 screening scores (rather than dichotomizing scores at a cutoff), and should combine screening
700 scores with individual characteristics to generate individualized probabilities of major depression.

701 **Contributions:**

702 BLevis, AB, JB, PC, SG, JPAI, LAK, DM, SBP, IS, RCZ and BDT were
703 responsible for the study conception and design. JB and LAK designed and conducted
704 database searches to identify eligible studies. DHA, BA, LA, HRB, MB, CHB, PB, GC,
705 MHC, JCNC, KC, YC, JMG, JD, JRF, FHF, DF, BG, FGS, CGG, BJH, JH, PAH,
706 MHärter, UH, LH, SEH, MHudson, MI, KI, NJ, MEK, KMK, YK, SL, ML, SRL, BLöwe,
707 LM, AM, SMS, TNM, KM, FLO, VP, BWP, PP, AP, KR, AGR, ISS, JS, ASidebottom,
708 ASimming, LS, SCS, PLLT, AT, CMvdFC, HCvW, PAV, JW, MAH, KW, MY, YZ, and
709 BDT contributed primary datasets that were included in this study. BLevis, KER, NS, MA,
710 DBR, MJC, TAS, and BDT contributed to data extraction and coding for the meta-analysis.
711 BLevis, AB, AWL, and BDT contributed to the data analysis and interpretation. BLevis,
712 AB, and BDT contributed to drafting the manuscript. All authors provided a critical review
713 and approved the final manuscript. AB and BDT are the guarantors; they had full access to
714 all the data in the study and take responsibility for the integrity of the data and the accuracy
715 of the data analyses. The corresponding author (BT) attests that all listed authors meet
716 authorship criteria and that no others meeting the criteria have been omitted.

717

718 **Copyright for authors:**

719 The Corresponding Author has the right to grant on behalf of all authors and does grant on
720 behalf of all authors, a worldwide licence to the Publishers and its licensees in perpetuity, in all
721 forms, formats and media (whether known now or created in the future), to i) publish, reproduce,
722 distribute, display and store the Contribution, ii) translate the Contribution into other languages,
723 create adaptations, reprints, include within collections and create summaries, extracts and/or,

724 abstracts of the Contribution, iii) create any other derivative work(s) based on the Contribution, iv)
725 to exploit all subsidiary rights in the Contribution, v) the inclusion of electronic links from the
726 Contribution to third party material where-ever it may be located; and, vi) licence any third party to
727 do any or all of the above.

728 The Corresponding Author has the right to grant on behalf of all authors and does grant on
729 behalf of all authors, an exclusive licence (or non exclusive for government employees) on a
730 worldwide basis to the BMJ Publishing Group Ltd to permit this article (if accepted) to be published
731 in BMJ editions and any other BMJ PGL products and sublicences such use and exploit all
732 subsidiary rights, as set out in our licence.

733

734 **Funding:**

735 This study was funded by the Canadian Institutes of Health Research (CIHR; KRS-134297).
736 Ms. Levis was supported by a CIHR Frederick Banting and Charles Best Canada Graduate
737 Scholarship doctoral award. Drs. Benedetti and Thombs were supported by Fonds de recherche du
738 Québec - Santé (FRQS) researcher salary awards. Ms. Riehm and Ms. Saadat were supported by
739 CIHR Frederick Banting and Charles Best Canada Graduate Scholarship master's awards. Mr. Levis
740 and Ms. Azar were supported by FRQS Masters Training Awards. Ms. Rice was supported by a
741 Vanier Canada Graduate Scholarship. Collection of data for the study by Arroll et al. was supported
742 by a project grant from the Health Research Council of New Zealand. Data collection for the study
743 by Ayalon et al. was supported from a grant from Lundbeck International. The primary study by
744 Khamseh et al. was supported by a grant (M-288) from Tehran University of Medical Sciences. The
745 primary study by Bombardier et al. was supported by the Department of Education, National
746 Institute on Disability and Rehabilitation Research, Spinal Cord Injury Model Systems: University

747 of Washington (grant no. H133N060033), Baylor College of Medicine (grant no. H133N060003),
748 and University of Michigan (grant no. H133N060032). Dr. Butterworth was supported by
749 Australian Research Council Future Fellowship FT130101444. Collection of data for the primary
750 study by Zhang et al. was supported by the European Foundation for Study of Diabetes, the Chinese
751 Diabetes Society, Lilly Foundation, Asia Diabetes Foundation and Liao Wun Yuk Diabetes
752 Memorial Fund. Dr. Conwell received support from NIMH (R24MH071604) and the Centers for
753 Disease Control and Prevention (R49 CE002093). Collection of data for the primary study by
754 Delgadillo et al. was supported by grant from St. Anne's Community Services, Leeds, United
755 Kingdom. Collection of data for the primary study by Fann et al. was supported by grant RO1
756 HD39415 from the US National Center for Medical Rehabilitation Research. The primary studies by
757 Amoozegar and by Fiest et al. were funded by the Alberta Health Services, the University of
758 Calgary Faculty of Medicine, and the Hotchkiss Brain Institute. The primary study by Fischer et al.
759 was funded by the German Federal Ministry of Education and Research (01GY1150). Data for the
760 primary study by Gelaye et al. was supported by grant from the NIH (T37 MD001449). Collection
761 of data for the primary study by Gjerdingen et al. was supported by grants from the NIMH (R34
762 MH072925, K02 MH65919, P30 DK50456). The primary study by Eack et al. was funded by the
763 NIMH (R24 MH56858). Collection of data for the primary study by Hobfoll et al. was made
764 possible in part from grants from NIMH (RO1 MH073687) and the Ohio Board of Regents. Dr. Hall
765 received support from a grant awarded by the Research and Development Administration Office,
766 University of Macau (MYRG2015-00109-FSS). Collection of data provided by Drs. Härter and
767 Reuter was supported by the Federal Ministry of Education and Research (grants No. 01 GD 9802/4
768 and 01 GD 0101) and by the Federation of German Pension Insurance Institute. The primary study
769 by Hides et al. was funded by the Perpetual Trustees, Flora and Frank Leith Charitable Trust, Jack

770 Brockhoff Foundation, Grosvenor Settlement, Sunshine Foundation and Danks Trust. The primary
771 study by Henkel et al. was funded by the German Ministry of Research and Education. Data for the
772 study by Razykov et al. was collected by the Canadian Scleroderma Research Group, which was
773 funded by the CIHR (FRN 83518), the Scleroderma Society of Canada, the Scleroderma Society of
774 Ontario, the Scleroderma Society of Saskatchewan, Sclérodermie Québec, the Cure Scleroderma
775 Foundation, Inova Diagnostics Inc., Euroimmun, FRQS, the Canadian Arthritis Network, and the
776 Lady Davis Institute of Medical Research of the Jewish General Hospital, Montreal, QC. Dr.
777 Hudson was supported by a FRQS Senior Investigator Award. Collection of data for the primary
778 study by Hyphantis et al. was supported by grant from the National Strategic Reference Framework,
779 European Union, and the Greek Ministry of Education, Lifelong Learning and Religious Affairs
780 (ARISTEIA-ABREVIATE, 1259). The primary study by Inagaki et al. was supported by the
781 Ministry of Health, Labour and Welfare, Japan. Dr. Jetté was supported by a Canada Research Chair
782 in Neurological Health Services Research. Collection of data for the primary study by Kiely et al.
783 was supported by National Health and Medical Research Council (grant number 1002160) and Safe
784 Work Australia. Dr. Kiely was supported by funding from a Australian National Health and Medical
785 Research Council fellowship (grant number 1088313). The primary study by Lamers et al. was
786 funded by the Netherlands Organisation for Health Research and development (grant number 945-
787 03-047). The primary study by Liu et al. was funded by a grant from the National Health Research
788 Institute, Republic of China (NHRI-EX97-9706PI). The primary study by Lotrakul et al. was
789 supported by the Faculty of Medicine, Ramathibodi Hospital, Mahidol University, Bangkok,
790 Thailand (grant number 49086). Dr. Bernd Löwe received research grants from Pfizer, Germany,
791 and from the medical faculty of the University of Heidelberg, Germany (project 121/2000) for the
792 study by Gräfe et al. The primary study by Mohd Sidik et al. was funded under the Research

793 University Grant Scheme from Universiti Putra Malaysia, Malaysia and the Postgraduate Research
794 Student Support Accounts of the University of Auckland, New Zealand. The primary study by
795 Santos et al. was funded by the National Program for Centers of Excellence
796 (PRONEX/FAPERGS/CNPq, Brazil). The primary study by Muramatsu et al. was supported by an
797 educational grant from Pfizer US Pharmaceutical Inc. Collection of primary data for the study by
798 Dr. Pence was provided by NIMH (R34MH084673). The primary studies by Osório et al. were
799 funded by Reitoria de Pesquisa da Universidade de São Paulo (grant number 09.1.01689.17.7) and
800 Banco Santander (grant number 10.1.01232.17.9). Dr. Osório was supported by Productivity Grants
801 (PQ-CNPq-2 -number 301321/2016-7). The primary study by Picardi et al. was supported by funds
802 for current research from the Italian Ministry of Health. Dr. Persoons was supported by a grant from
803 the Belgian Ministry of Public Health and Social Affairs and a restricted grant from Pfizer Belgium.
804 Dr. Shaaban was supported by funding from Universiti Sains Malaysia. The primary study by
805 Rooney et al. was funded by the United Kingdom National Health Service Lothian Neuro-Oncology
806 Endowment Fund. The primary study by Sidebottom et al. was funded by a grant from the United
807 States Department of Health and Human Services, Health Resources and Services Administration
808 (grant number R40MC07840). Simning et al.'s research was supported in part by grants from the
809 NIH (T32 GM07356), Agency for Healthcare Research and Quality (R36 HS018246), NIMH (R24
810 MH071604), and the National Center for Research Resources (TL1 RR024135). Dr. Stafford
811 received PhD scholarship funding from the University of Melbourne. Collection of data for the
812 studies by Turner et al were funded by a bequest from Jennie Thomas through the Hunter Medical
813 Research Institute. The study by van Steenberg-Weijnenburg et al. was funded by Innovatiefonds
814 Zorgverzekeraars. Dr. Vöhringer was supported by the Fund for Innovation and Competitiveness of
815 the Chilean Ministry of Economy, Development and Tourism, through the Millennium Scientific

816 Initiative (grant number IS130005). Collection of data for the primary study by Williams et al. was
817 supported by a NIMH grant to Dr. Marsh (RO1-MH069666). The primary study by Thombs et al.
818 was done with data from the Heart and Soul Study (PI Mary Whooley). The Heart and Soul Study
819 was funded by the Department of Veterans Epidemiology Merit Review Program, the Department
820 of Veterans Affairs Health Services Research and Development service, the National Heart Lung
821 and Blood Institute (R01 HL079235), the American Federation for Aging Research, the Robert
822 Wood Johnson Foundation, and the Ischemia Research and Education Foundation. The primary
823 study by Twist et al. was funded by the UK National Institute for Health Research under its
824 Programme Grants for Applied Research Programme (grant reference number RP-PG-0606-1142).
825 The study by Wittkamp et al. was funded by The Netherlands Organization for Health Research
826 and Development (ZonMw) Mental Health Program (nos. 100.003.005 and 100.002.021) and the
827 Academic Medical Center/University of Amsterdam. No other authors reported funding for primary
828 studies or for their work on the present study.

829

830 **Declaration of Competing Interests:**

831 All authors have completed the ICJME uniform disclosure form at
832 www.icmje.org/coi_disclosure.pdf and declare: no support from any organisation for the submitted
833 work; no financial relationships with any organisations that might have an interest in the submitted
834 work in the previous three years with the following exceptions: Drs. Jetté and Patten declare that
835 they received a grant, outside the submitted work, from the University of Calgary Hotchkiss Brain
836 Institute, which was jointly funded by the Institute and Pfizer. Pfizer was the original sponsor of the
837 development of the PHQ-9, which is now in the public domain. Dr. Chan is a steering committee
838 member or consultant of Astra Zeneca, Bayer, Lilly, MSD and Pfizer. She has received

839 sponsorships and honorarium for giving lectures and providing consultancy and her affiliated
840 institution has received research grants from these companies. Dr. Hegerl declares that within the
841 last three years, he was an advisory board member for Lundbeck and Servier; a consultant for Bayer
842 Pharma; a speaker for Roche Pharma and Servier; and received personal fees from Janssen, all
843 outside the submitted work. Dr. Inagaki declares that he has received a grant from Novartis Pharma,
844 and personal fees from Meiji, Mochida, Takeda, Novartis, Yoshitomi, Pfizer, Eisai, Otsuka, MSD,
845 Technomics, and Sumitomo Dainippon, all outside of the submitted work. All authors declare no
846 other relationships or activities that could appear to have influenced the submitted work. No funder
847 had any role in the design and conduct of the study; collection, management, analysis, and
848 interpretation of the data; preparation, review, or approval of the manuscript; and decision to submit
849 the manuscript for publication.

850

851 **Ethical Approval:** As this study involved secondary analysis of anonymized previously
852 collected data, the Research Ethics Committee of the Jewish General Hospital declared that
853 this project did not require research ethics approval. However, for each included dataset,
854 we confirmed that the original study received ethics approval and that all patients provided
855 informed consent.

856

857 **Transparency Declaration:** The manuscript's guarantor affirms that the manuscript is an honest,
858 accurate, and transparent account of the study being reported; that no important aspects of the study
859 have been omitted; and that any discrepancies from the study as planned (and, if relevant,
860 registered) have been explained.

861

862 **Data Sharing:** Requests to access data should be made to the corresponding author at

863 brett.thombs@mcgill.ca.

864

865 **What is already known on this topic:**

- 866 • The PHQ-9 is the most commonly used depression screening tool in primary care.
- 867 • Previous meta-analyses on the diagnostic test accuracy of the PHQ-9 have been limited by
- 868 selective cutoff reporting in primary studies; the inability to assess differences across patient
- 869 subgroups; the inability to exclude participants already diagnosed or being treated for
- 870 depression, who would not be screened in practice; and the combining of accuracy estimates
- 871 without differentiating between reference standards.

872 **What this study adds:**

- 873 • PHQ-9 diagnostic accuracy when compared to diagnoses made by semi-structured
- 874 diagnostic interviews is greater than when compared to diagnoses made by other reference
- 875 standards and greater than reported in previous meta-analyses, which did not distinguish
- 876 between different diagnostic standards.
- 877 • PHQ-9 diagnostic accuracy does not differ substantively across participant subgroups except
- 878 for age, where it may be more specific among older patients.
- 879 • The standard cutoff of 10 or greater maximizes combined sensitivity and specificity overall
- 880 and for subgroups.

881

882 **PRINT ABSTRACT**

883 **Study question:** What is the diagnostic accuracy of the Patient Health Questionnaire-9 (PHQ-9) for
884 screening to detect major depression?

885 **Methods:** Individual participant data meta-analysis was used to synthesize results from studies that
886 compared PHQ-9 scores to major depression diagnoses from validated diagnostic interviews. For
887 PHQ-9 cutoffs 5 to 15, bivariate random-effects meta-analysis was used to estimate pooled
888 sensitivity and specificity among studies that used semi-structured diagnostic interviews, fully
889 structured interviews, and the Mini International Neuropsychiatric (MINI), separately. Sensitivity
890 and specificity were examined among participant subgroups and, separately, using meta-regression,
891 considering all subgroup variables in a single model.

892 **Study answer and limitations:** Data were obtained for 58 of 72 eligible studies (N participants =
893 17,357, N cases = 2,312). Combined sensitivity and specificity was maximized at a cutoff of ≥ 10
894 among studies using a semi-structured interview; sensitivity [95% CI] was 0.88 [0.83 to 0.92],
895 specificity [95% CI] was 0.85 [0.82 to 0.88]). Across cutoffs, sensitivity with semi-structured
896 interviews was higher than for fully structured interviews (MINI excluded) and for the MINI.
897 Specificity was similar across diagnostic interviews. In studies conducted in primary care using a
898 semi-structured interview (major depression prevalence = 12%), sensitivity [95% CI] was 0.94 [0.88
899 to 0.97] and specificity [95% CI] was 0.88 [0.79 to 0.93]). Study limitations include the inability to
900 obtain data for 14 eligible studies, substantial heterogeneity across included studies, and the
901 inability to conduct subgroup analyses based on specific medical comorbidities or cultural aspects.

902 **What this study adds:** PHQ-9 sensitivity compared to semi-structured diagnostic interviews was
903 greater than in previous meta-analyses that combined reference standards. A cutoff of ≥ 10
904 maximized combined sensitivity and specificity overall and for subgroups.

905 **Funding:** The Canadian Institutes of Health Research (KRS-134297; PCG-155468).

906 **Registration:** PROSPERO (CRD42014010673).

907

908 **REFERENCES**

- 909 1. Thombs BD, Ziegelstein RC. Does depression screening improve depression outcomes in
910 primary care? *BMJ*. 2014;**348**:g1253.
- 911 2. Thombs BD, Coyne JC, Cuijpers P, et al. Rethinking recommendations for screening for
912 depression in primary care. *CMAJ* .2012;184:413-8.
- 913 3. Kroenke K, Spitzer RL, Williams JB. The PHQ-9: validity of a brief depression severity
914 measure. *J Gen Intern Med*. 2001;**16**:606-13.
- 915 4. Kroenke K, Spitzer RL. The PHQ-9: a new depression diagnostic and severity measure.
916 *Psychiatr Ann*. 2002;**32**:1-7.
- 917 5. Spitzer RL, Kroenke K, Williams JB. Validation and utility of a self-report version of PRIME-
918 MD: the PHQ primary care study. Primary Care Evaluation of Mental Disorders. Patient Health
919 Questionnaire. *JAMA*. 1999;**282**:1737-44.
- 920 6. Wittkamp KA, Naeije L, Schene AH, et al. Diagnostic accuracy of the mood module of the
921 Patient Health Questionnaire: a systematic review. *Gen Hosp Psychiatry*. 2007;**29**:388-95.
- 922 7. Gilbody S, Richards D, Brealey S, et al. Screening for depression in medical settings with the
923 Patient Health Questionnaire (PHQ): a diagnostic meta-analysis. *J Gen Intern Med*.
924 2007;**22**:1596-1602.
- 925 8. Moriarty AS, Gilbody S, McMillan D, Manea L. Screening and case finding for major
926 depressive disorder using the Patient Health Questionnaire (PHQ-9): a meta-analysis. *Gen Hosp*
927 *Psychiatry*. 2015;**37**:567-76.
- 928 9. Levis B, Benedetti A, Levis AW, et al. Selective cutoff reporting in studies of diagnostic test
929 accuracy: a comparison of conventional and individual-patient-data meta-analysis of the Patient
930 Health Questionnaire-9 depression screening tool. *Am J Epidemiol*. 2017;**185**:954-64.

- 931 10. Thombs BD, Arthurs E, El-Baalbaki G, et al. Risk of bias from inclusion of already diagnosed
932 or treated patients in diagnostic accuracy studies of depression screening tools: A systematic
933 review. *BMJ*. 2011;**343**:d4825.
- 934 11. Rice DB, Thombs BD. Risk of bias from inclusion of currently diagnosed or treated patients in
935 studies of depression screening tool accuracy: A cross-sectional analysis of recently published
936 primary studies and meta-analyses. *PLOS ONE*. 2016;11:e0150067.
- 937 12. Levis B, Benedetti A, Riehm KE, et al. Probability of major depression diagnostic classification
938 using semi-structured vs. fully structured diagnostic interviews. *Br J Psychiatry*. 2018;**212**:377-
939 85.
- 940 13. First MB. Structured clinical interview for the DSM (SCID). John Wiley & Sons, Inc. 1995.
- 941 14. Robins LN, Wing J, Wittchen HU, et al. The Composite International Diagnostic Interview: an
942 epidemiologic instrument suitable for use in conjunction with different diagnostic systems and
943 in different cultures. *Arch Gen Psychiatry*. 1988;**45**:1069-77.
- 944 15. Brugha TS, Jenkins R, Taub N, Meltzer H, Bebbington PE. A general population comparison of
945 the Composite International Diagnostic Interview (CIDI) and the Schedules for Clinical
946 Assessment in Neuropsychiatry (SCAN). *Psychol Med*. 2001;**31**:1001-13.
- 947 16. Brugha TS, Bebbington PE, Jenkins R. A difference that matters: comparisons of structured and
948 semi-structured psychiatric diagnostic interviews in the general population. *Psychol Med*.
949 1999;**29**:1013-20.
- 950 17. Nosen E, Woody SR. Chapter 8: Diagnostic Assessment in Research. In, McKay D. Handbook
951 of research methods in abnormal and clinical psychology. Sage; 2008.
- 952 18. Kurdyak PA, Gnam WH. Small signal, big noise: performance of the CIDI depression module.
953 *Can J Psychiatry*. 2005;**50**:851-6.

- 954 19. Lecrubier Y, Sheehan DV, Weiller E et al. The Mini International Neuropsychiatric Interview
955 (MINI). A short diagnostic structured interview: reliability and validity according to the CIDI.
956 *Eur Psychiatry*. 1997;**12**:224-31.
- 957 20. Sheehan DV, Lecrubier Y, Sheehan KH et al. The validity of the Mini International
958 Neuropsychiatric Interview (MINI) according to the SCID-P and its reliability. *Eur Psychiatry*.
959 1997;**12**:232-41.
- 960 21. Riley RD, Lambert PC, Abo-Zaid G. Meta-analysis of individual participant data: rationale,
961 conduct, and reporting. *BMJ*. 2010;**340**:c221.
- 962 22. Thombs BD, Benedetti A, Kloda LA, et al. The diagnostic accuracy of the Patient Health
963 Questionnaire-2 (PHQ-2), Patient Health Questionnaire-8 (PHQ-8), and Patient Health
964 Questionnaire-9 (PHQ-9) for detecting major depression: protocol for a systematic review and
965 individual patient data meta-analyses. *Syst Rev*. 2014;**27**;3:124.
- 966 23. McInnes MDF, Moher D, Thombs BD, et al. Preferred Reporting Items for a Systematic Review
967 and Meta-analysis of Diagnostic Test Accuracy Studies: The PRISMA-DTA Statement. *JAMA*.
968 2018;**319**(4):388-96.
- 969 24. Stewart LA, Clarke M, Rovers M, et al. Preferred Reporting Items for Systematic Review and
970 Meta-Analyses of individual participant data: the PRISMA-IPD Statement. *JAMA*.
971 2015;**313**(16):1657-65.
- 972 25. PRESS – Peer Review of Electronic Search Strategies: 2015 Guideline Explanation and
973 Elaboration (PRESS E&E). Ottawa: CADTH; 2016.
- 974 26. United Nations. International Human Development Indicators. <http://hdr.undp.org/en/countries>.
975 Accessed February 1, 2019.

- 976 27. Whiting PF, Rutjes AW, Westwood ME, et al. QUADAS-2: a revised tool for the quality
977 assessment of diagnostic accuracy studies. *Ann Intern Med.* 2011;**155**:529-36.
- 978 28. World Health Organization. Schedules for clinical assessment in neuropsychiatry: manual. Amer
979 Psychiatric Pub Inc. 1994.
- 980 29. Freedland KE, Skala JA, Carney RM, et al. The Depression Interview and Structured Hamilton
981 (DISH): rationale, development, characteristics, and clinical validity. *Psychosom Med.*
982 2002;**64**:897-905.
- 983 30. Lewis G, Pelosi AJ, Araya R, Dunn G. Measuring psychiatric disorder in the community: a
984 standardized assessment for use by lay interviewers. *Psychol Med.* 1992;**22**:465-86.
- 985 31. Robins LN, Helzer JE, Croughan J, Ratcliff KS. National Institute of Mental Health Diagnostic
986 Interview Schedule: Its history, characteristics, and validity. *Arch Gen Psychiatry.* 1981;**38**:381-
987 9.
- 988 32. Riley RD, Dodd SR, Craig JV, et al. Meta-analysis of diagnostic test studies using individual
989 patient data and aggregate data. *Stat Med.* 2008;**27**:6111-36.
- 990 33. van der Leeden R, Busing FMTA, Meijer E. Bootstrap methods for two-level models. Technical
991 Report PRM 97-04, Leiden University, Department of Psychology, Leiden, The Netherlands,
992 1997.
- 993 34. van der Leeden R, Meijer E, Busing FMTA. Chapter 11: Resampling multilevel models. In:
994 Leeuw J, Meijer E, eds. *Handbook of multilevel analysis* New York, NY: Springer; 2008:401-
995 33.
- 996 35. Higgins JP, Thompson SG. Quantifying heterogeneity in a meta-analysis. *Stat Med.*
997 2002;**21**:1539-58.

- 998 36. Mitchell AJ, Yadegarfar M, Gill J, Stubbs B. Case finding and screening clinical utility of the
999 Patient Health Questionnaire (PHQ-9 and PHQ-2) for depression in primary care: a diagnostic
1000 meta-analysis of 40 studies. *BJPsych Open*. 2016;**2**:127-38.
- 1001 37. Shankman SA, Funkhouser CJ, Klein DN, et al. Reliability and validity of severity dimensions
1002 of psychopathology assessed using the Structured Clinical Interview for DSM-5 (SCID). *Int J*
1003 *Methods Psychiatr Res*. 2018;**27**:e1590.
- 1004 38. Semler G, Wittchen HU, Joschke K, et al. Test-retest reliability of a standardized psychiatric
1005 interview (DIS/CIDI). *Eur Arch Psychiatry Neurol Sci*. 1987;**236**:214-22
- 1006 39. Siu AL, and the US Preventive Services Task Force (USPSTF). Screening for Depression in
1007 Adults: US Preventive Services Task Force Recommendation Statement. *JAMA*. 2016;**315**:380-
1008 7.
- 1009 40. Allaby M. Screening for depression: A report for the UK National Screening Committee
1010 (Revised report). London, United Kingdom: UK National Screening Committee; 2010.
- 1011 41. Joffres M, Jaramillo A, Dickinson J, et al. Recommendations on screening for depression in
1012 adults. *CMAJ*. 2013;**185**:775-82.
- 1013 42. Thombs BD, Ziegelstein RC, Roseman M, Kloda LA, Ioannidis JP. There are no randomized
1014 controlled trials that support the United States Preventive Services Task Force guideline on
1015 screening for depression in primary care: A systematic review. *BMC Med*. 2014;**12**:13.
- 1016 43. National Institute for Health and Care Excellence. Depression in Adults: treatment and
1017 management. Consultation draft (May 2018). [https://www.nice.org.uk/guidance/gid-](https://www.nice.org.uk/guidance/gid-cgwave0725/documents/full-guideline-updated)
1018 [cgwave0725/documents/full-guideline-updated](https://www.nice.org.uk/guidance/gid-cgwave0725/documents/full-guideline-updated). Accessed February 1, 2019.
- 1019 44. Ferrante di Ruffano L, Hyde CJ, McCaffery KJ, Bossuyt PM, Deeks JJ. Assessing the value of
1020 diagnostic tests: a framework for designing and evaluating trials. *BMJ*. 2012;**344**:e686.

1021 **FIGURES**

1022

1023 **Figure 1. Nomograms of positive and negative predictive value for cutoff 10 of the PHQ-9 for**
1024 **each reference standard category**

1025

1026 Nomograms of a) positive predictive value and b) negative predictive value for cutoff 10 of the
1027 PHQ-9, for major depression prevalence values of 5 to 25%, for semi-structured diagnostic
1028 interviews, fully structured diagnostic interviews, and the MINI.

TABLES

Table 1. Participant data by diagnostic interview

Diagnostic Interview	N Studies	N Participants	N (%) Major Depression
Semi-structured			
SCID	26	4,733	785 (17)
SCAN	2	1,892	130 (7)
DISH	1	100	9 (9)
Fully structured			
CIDI	11	6,272	554 (9)
DIS	1	1,006	221 (22)
CIS-R	2	402	64 (16)
MINI	15	2,952	549 (19)
Total	58	17,357	2,312 (13)

Abbreviations: CIDI: Composite International Diagnostic Interview; CIS-R: Clinical Interview Schedule-Revised; DIS: Diagnostic Interview Schedule; DISH: Depression Interview and Structured Hamilton; MINI: Mini International Neuropsychiatric Interview; SCAN: Schedules for Clinical Assessment in Neuropsychiatry; SCID: Structured Clinical Interview for DSM Disorders

Table 2. Participant data by subgroup

Participant Subgroup	Semi-Structured Diagnostic Interviews			Fully Structured Diagnostic Interviews			MINI		
	N Studies	N Participants	N (%) Major Depression	N Studies	N Participants	N (%) Major Depression	N Studies	N Participants	N (%) Major Depression
All participants	29	6,725	924 (14)	14	7,680	839 (11)	15	2,952	549 (19)
Participants not currently diagnosed or receiving treatment for a mental health problem	20	2,942	421 (14)	6	4,161	306 (7)	6	927	168 (18)
Age <60	26	4,132	629 (15)	14	5,504	645 (12)	14	1,958	310 (16)
Age ≥60	24	2,577	295 (11)	10	2,175	194 (9)	13	979	239 (24)
Women	28	3,906	573 (15)	14	4,285	463 (11)	15	1,666	337 (20)
Men	25	2,812	351 (12)	13	3,395	376 (11)	15	1,286	212 (16)
Very high country human development index	25	6,195	739 (12)	9	5,740	592 (10)	10	1,924	430 (22)
High country human development index	4	530	185 (35)	2	326	61 (19)	3	542	61 (11)
Low-medium country human development index	--	--	--	3	1,614	186 (12)	2	486	58 (12)
Non-medical care	2	567	105 (19)	2	963	74 (8)	2	299	72 (24)
Primary care	9	3,163	377 (12)	5	3,578	273 (8)	5	1,290	168 (13)
Inpatient specialty care	8	867	121 (14)	2	372	34 (9)	1	137	25 (18)
Outpatient specialty care	12	2,128	321 (15)	5	2,767	458 (17)	7	1,226	284 (23)

^aSome variables were coded at the study level, while others were coded at the participant level. Thus, number of studies does not always add up to total number in the reference category

Table 3a. Comparison of sensitivity and specificity estimates among semi-structured vs. fully structured reference standards

Cutoff	Semi-Structured Reference Standard ^a		Fully Structured Reference Standard ^b		Difference across reference standards (Semi-structured - Fully structured) ^c	
	Sensitivity (95% CI)	Specificity (95% CI)	Sensitivity (95% CI)	Specificity (95% CI)	Sensitivity (95% CI)	Specificity (95% CI)
5	0.98 (0.96 to 0.99)	0.55 (0.49 to 0.60)	0.93 (0.87 to 0.97)	0.54 (0.43 to 0.64)	0.05 (-0.01 to 0.13)	0.01 (-0.13 to 0.16)
6	0.98 (0.95 to 0.99)	0.63 (0.58 to 0.67)	0.91 (0.83 to 0.95)	0.61 (0.51 to 0.71)	0.07 (-0.01 to 0.18)	0.02 (-0.12 to 0.17)
7	0.98 (0.94 to 0.99)	0.69 (0.65 to 0.74)	0.86 (0.75 to 0.92)	0.69 (0.59 to 0.77)	0.12 (0.00 to 0.26)	0.00 (-0.10 to 0.15)
8	0.95 (0.91 to 0.97)	0.75 (0.71 to 0.79)	0.82 (0.71 to 0.89)	0.75 (0.66 to 0.82)	0.13 (0.00 to 0.28)	0.00 (-0.10 to 0.13)
9	0.91 (0.87 to 0.94)	0.80 (0.77 to 0.83)	0.74 (0.63 to 0.83)	0.79 (0.72 to 0.86)	0.17 (0.05 to 0.34)	0.01 (-0.08 to 0.12)
10	0.88 (0.83 to 0.92)	0.85 (0.82 to 0.88)	0.70 (0.59 to 0.80)	0.84 (0.77 to 0.89)	0.18 (0.04 to 0.36)	0.01 (-0.05 to 0.12)
11	0.84 (0.78 to 0.89)	0.89 (0.86 to 0.91)	0.62 (0.51 to 0.72)	0.87 (0.81 to 0.91)	0.22 (0.07 to 0.40)	0.02 (-0.04 to 0.10)
12	0.79 (0.73 to 0.83)	0.91 (0.89 to 0.93)	0.57 (0.45 to 0.68)	0.89 (0.85 to 0.93)	0.22 (0.05 to 0.40)	0.02 (-0.03 to 0.09)
13	0.70 (0.65 to 0.75)	0.93 (0.91 to 0.95)	0.49 (0.38 to 0.61)	0.92 (0.89 to 0.95)	0.21 (0.04 to 0.40)	0.01 (-0.03 to 0.07)
14	0.64 (0.58 to 0.70)	0.95 (0.93 to 0.96)	0.44 (0.32 to 0.56)	0.94 (0.91 to 0.96)	0.20 (0.03 to 0.40)	0.01 (-0.02 to 0.05)
15	0.56 (0.50 to 0.62)	0.96 (0.95 to 0.97)	0.35 (0.25 to 0.46)	0.96 (0.93 to 0.97)	0.21 (0.05 to 0.39)	0.00 (-0.02 to 0.04)

^a N Studies = 29; N Participants = 6,725; N major depression = 924

^b N Studies = 14; N Participants = 7,680; N major depression = 839

^c 1 bootstrap iteration (0.01%) did not produce a difference estimate for cutoff 5. This iteration was removed prior to determining the bootstrapped CI.

Abbreviations: CI: confidence interval

Table 3b. Comparison of sensitivity and specificity estimates among semi-structured vs. MINI reference standards

Cutoff	Semi-Structured Reference Standard ^a		MINI Reference Standard ^b		Difference across reference standards (Semi-structured - MINI)	
	Sensitivity (95% CI)	Specificity (95% CI)	Sensitivity (95% CI)	Specificity (95% CI)	Sensitivity (95% CI)	Specificity (95% CI)
5	0.98 (0.96 to 0.99)	0.55 (0.49 to 0.60)	0.96 (0.93 to 0.98)	0.57 (0.50 to 0.64)	0.02 (-0.02 to 0.07)	-0.02 (-0.14 to 0.11)
6	0.98 (0.95 to 0.99)	0.63 (0.58 to 0.67)	0.93 (0.87 to 0.97)	0.66 (0.59 to 0.72)	0.05 (-0.01 to 0.12)	-0.03 (-0.13 to 0.09)
7	0.98 (0.94 to 0.99)	0.69 (0.65 to 0.74)	0.90 (0.82 to 0.94)	0.72 (0.66 to 0.78)	0.08 (-0.00 to 0.16)	-0.03 (-0.12 to 0.08)
8	0.95 (0.91 to 0.97)	0.75 (0.71 to 0.79)	0.86 (0.78 to 0.91)	0.78 (0.73 to 0.83)	0.09 (-0.01 to 0.19)	-0.03 (-0.11 to 0.06)
9	0.91 (0.87 to 0.94)	0.80 (0.77 to 0.83)	0.82 (0.72 to 0.88)	0.84 (0.79 to 0.87)	0.09 (-0.02 to 0.22)	-0.04 (-0.09 to 0.05)
10	0.88 (0.83 to 0.92)	0.85 (0.82 to 0.88)	0.77 (0.68 to 0.83)	0.87 (0.83 to 0.90)	0.11 (-0.01 to 0.25)	-0.02 (-0.07 to 0.06)
11	0.84 (0.78 to 0.89)	0.89 (0.86 to 0.91)	0.70 (0.62 to 0.77)	0.90 (0.86 to 0.92)	0.14 (0.01 to 0.30)	-0.01 (-0.06 to 0.05)
12	0.79 (0.73 to 0.83)	0.91 (0.89 to 0.93)	0.65 (0.56 to 0.72)	0.92 (0.89 to 0.94)	0.14 (-0.01 to 0.28)	-0.01 (-0.05 to 0.05)
13	0.70 (0.65 to 0.75)	0.93 (0.91 to 0.95)	0.57 (0.49 to 0.65)	0.94 (0.91 to 0.96)	0.13 (-0.03 to 0.26)	-0.01 (-0.04 to 0.04)
14 ^c	0.64 (0.58 to 0.70)	0.95 (0.93 to 0.96)	0.49 (0.42 to 0.56)	0.96 (0.93 to 0.97)	0.15 (0.01 to 0.28)	-0.01 (-0.04 to 0.03)
15 ^c	0.56 (0.50 to 0.62)	0.96 (0.95 to 0.97)	0.42 (0.35 to 0.49)	0.97 (0.95 to 0.98)	0.14 (-0.01 to 0.27)	-0.01 (-0.03 to 0.02)

^a N Studies = 29; N Participants = 6,725; N major depression = 924

^b N Studies = 15; N Participants = 2,952; N major depression = 549

^c For these cutoffs, among studies that used the MINI as the reference standard, the default optimizer in glmer failed, thus bobyqa was used instead.

Abbreviations: CI: confidence interval; MINI: Mini International Neuropsychiatric Interview