

1 The Performance of Digital Technologies for 2 Measuring Tuberculosis Medication 3 Adherence: A Systematic Review

4
5 Miranda Zary¹, Mona Salaheldin Mohamed¹, Cedric Kafie¹, Chimweta Ian Chilala², Shruti Bahukudumbi³,
6 Nicola Foster², Genevieve Gore⁴, Katherine Fielding², Ramnath Subbaraman^{3,5}, and Kevin Schwartzman¹

7
8 ¹McGill International Tuberculosis Centre; Research Institute of the McGill University Health Centre,
9 Montréal, Canada

10 ²TB Centre, London School of Hygiene and Tropical Medicine, London, UK

11 ³Department of Public Health and Community Medicine, Tufts University School of Medicine, Boston,
12 USA

13 ⁴Schulich Library of Physical Sciences, Life Sciences, and Engineering; McGill University, Montréal,
14 Canada

15 ⁵Division of Geographic Medicine and Infectious Diseases, Tufts Medical Center, Boston, USA

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

Word count (without abstract): 4,590

Correspondence: Dr. Kevin Schwartzman
Centre for Outcomes Research and Evaluation
Research Institute of the McGill University Health Centre
5252 boulevard de Maisonneuve Ouest, Room 3D.63
Montreal, Quebec H4A 3S5
CANADA
Email kevin.schwartzman@mcgill.ca

38 **ABSTRACT** (296 words)

39 **Introduction:** Digital adherence technologies (DATs), such as phone-based technologies, and digital
40 pillboxes, can provide more person-centric approaches to support tuberculosis (TB) medication
41 adherence. We synthesized evidence addressing the performance of DATs for measuring tuberculosis
42 medication adherence.

43 **Methods:** We conducted a systematic review (PROSPERO - CRD42022313526) which identified relevant
44 published literature from January 2000 through April 2023 in five databases, and pertinent preprints.
45 Studies reporting quantitative data on the performance of DATs for measuring adherence to
46 medications for TB disease or infection, against a reference standard, with at least 20 participants using
47 the DAT were included. Study characteristics and performance outcomes (e.g., sensitivity, specificity,
48 positive and negative predictive values) were extracted. Article quality was assessed using the QUADAS-
49 2 tool for diagnostic accuracy studies.

50 **Results:** Of 5692 studies initially identified by our systematic search, 13 met our inclusion criteria. These
51 studies addressed the performance of medication sleeves with phone calls [branded as “99DDOTS”;
52 N=4], digital pillboxes [N=5], ingestible sensors [N=2], artificial intelligence-based video observed
53 therapy [N=1], and multifunctional mobile applications [N=1]. All but one involved persons with TB
54 disease. For medication sleeves with phone calls, compared to urine analysis, reported sensitivity and
55 specificity was 70-94% and 0-61%, respectively. For digital pillboxes, compared to pill count, reported
56 sensitivity and specificity was 25-99% and 69-100%, respectively. For ingestible sensors, the sensitivity of
57 dose detection was $\geq 95\%$ in comparison to directly observed ingestion. Participant selection was the
58 most frequent potential source of bias across articles.

59 **Conclusion:** Limited available data suggest suboptimal and variable performance of DATs for dose
60 monitoring, with significant evidence gaps, notably in real-world programmatic settings. Future research

61 should aim to improve understanding of the relationships of specific technologies, settings, user
62 characteristics, and user engagement with DAT performance, and should measure and report
63 performance in a more standardized manner.

64 **Key words:** digital adherence technology, systematic review, accuracy, tuberculosis, performance,
65 treatment observation, medication adherence

66

67

68

69

70

71

72

73 **KEY MESSAGES**

74 **What is already known on this topic:** Several cohort studies have suggested that digital adherence
75 technologies (DATs) can both underestimate and overestimate medication ingestion among persons
76 treated for tuberculosis. No previous review has synthesized available evidence in this regard.

77 **What this study adds:** Reports of DAT (medication sleeves with phone calls, digital pillboxes)
78 implementation in real-world treatment settings consistently indicate suboptimal performance for
79 measuring medication adherence. However, available evidence is limited in scope and quality.

80 **How this study might affect research, practice, or policy:** Suboptimal dose reporting from DATs
81 potentially compromises their effectiveness, and program efficiency. Future clinical practice will be
82 strengthened by rigorous technology evaluations that reflect more consistent use of reference
83 standards, and clearer benchmarks for medication adherence.

84 **INTRODUCTION** (511 words)

85 Tuberculosis (TB) disease requires treatment with medication regimens involving varying pill
86 burdens lasting at least 4 months. Adherence to these treatment regimens—which involve daily
87 medication intake—is crucial, as non-adherence can lead to treatment failure, relapse, development of
88 drug resistance, ongoing TB transmission, and death.(1,2) Treatment for TB infection is between 1 and 9
89 months. Adherence to treatment for TB infection is essential to reduce the risk of developing TB disease
90 but is often difficult to achieve in routine care, particularly since treated persons are asymptomatic.(3,4)
91 Directly observed therapy (DOT) has been recommended for TB treatment support.(5) DOT involves
92 healthcare workers, or community workers, watching up to 100% of prescribed medication doses.(6)
93 However, DOT is logistically challenging, expensive, potentially intrusive, raises ethical concerns, and can
94 have limited or varying effectiveness for improving treatment outcomes.(6–8)

95 Digital adherence technologies (DATs) have been increasingly studied and used in routine care
96 as an alternative or adjunctive approach for supporting TB treatment. DATs include mobile
97 communication and other innovations that can remind people with TB to take their medication, digitally
98 observe doses taken, compile dosing histories, triage people who may be at higher risk for unfavorable
99 treatment outcomes, and enable differentiated (i.e., intensified or individualized) care.(2) DATs include
100 a range of technologies that may facilitate more person-centered approaches for monitoring adherence,
101 potentially improving treatment outcomes.(2)

102 For example, one of the most widely used DATs, branded as 99DOTS, involves wrapping a paper
103 sleeve over a medication blister pack. Dispensation of a medication dose then reveals a hidden phone
104 number. By calling this number, the person with TB can report dose ingestion, creating a digital dosing
105 history that allows early identification by healthcare providers of people who may be nonadherent.(9–
106 12) However, people with TB may call the phone number on the medication sleeve without taking their
107 doses (i.e., over-reporting adherence) or take medication doses without calling (i.e., under-reporting
108 adherence), which may hinder the ability to identify people experiencing nonadherence. Other DATs

109 that have been used in routine care by TB programs (13,14)—such as two-way short messaging service
110 (SMS), digital pillboxes, or video-supported therapy—similarly involve an SMS response, pillbox opening,
111 or remote visualization of dosing by video, respectively, serving as a proxy for a treatment dose taken.
112 All may face challenges related to over-reporting or under-reporting of adherence.

113 An initial systematic review examining DATs for TB treatment support was published in 2018.
114 However, since 2018—particularly with the COVID-19 pandemic—interest in and experience with these
115 technologies have expanded substantially.(15) While other reviews have been published subsequently
116 (2,16,17), none has examined the performance of DATs for measuring TB medication adherence, which
117 is crucial for healthcare providers to identify people with TB who may be experiencing nonadherence, so
118 they can be given intensified or individualized support. If they do not yield accurate assessments of
119 adherence, any resulting actionable information is of questionable value, resulting in limited public
120 health impact.(2) The present systematic review synthesizes evidence on the performance of digital
121 adherence technologies for measuring TB medication adherence, among persons treated for
122 tuberculosis disease and infection.

123 **METHODS** (1,178 words)

124 **Design.** Our systematic review protocol was registered in PROSPERO, the International
125 Prospective Register of Systematic Reviews (CRD42022313526).(18) This review follows the Preferred
126 Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) guidelines.

127 **Search and screening strategy.** The search for relevant literature was conducted on April 28th,
128 2023 (updated from April 14, 2022) in MEDLINE/Ovid, Embase, CENTRAL, CINAHL, and Web of Science
129 Core Collection, Europe PMC preprints (including MedRxiv) and clinicaltrials.gov from January 1, 2000, to
130 April 28, 2023. Key search concepts included TB (disease or infection), digital technologies (such as
131 mobile phone, smartphone, video observation, digital pillboxes, and text messaging), and accuracy (such

132 as sensitivity, specificity, and area under the curve). The complete search strategy can be found in Table
133 S1 of the supplemental appendix. The database searches were conducted by a health librarian (GG).
134 Separately, we hand-searched the Union World Conference on Lung Health for relevant abstracts on
135 DATs and performance from 2004 to 2022 inclusively. There were no language restrictions.

136 **Inclusion/Exclusion Criteria.** Studies were included if they reported a quantitative outcome
137 addressing the performance of DATs for measuring TB medication adherence (i.e., any of sensitivity,
138 specificity, positive predictive value [PPV], negative predictive value [NPV], area under the receiver
139 operating characteristic curve [AUC], likelihood ratio, accuracy, or agreement). These are defined below
140 (Table 1). Articles were included if the number of participants using the DAT was at least 20, and the
141 study design included the comparison of adherence reports generated by a DAT with a reference
142 standard such as urine drug metabolite testing, pill count, direct observation of medication ingestion, or
143 other such information. DAT interventions included but were not limited to smartphone-based
144 technologies such as phone-based dosing records, SMS or video-supported treatment, digital pillboxes,
145 and ingestible sensors. We defined a DAT as an intervention with a digital component (which could be
146 part of a multi-component intervention) with the intention to measure and promote treatment
147 adherence and/or reduce missed visits and/or reduce losses to follow-up. Examples of DATs included
148 and excluded with this definition can be found in Table S2.

149 We included studies of persons treated for TB disease or infection, including subgroups such as
150 people with drug-resistant TB and people with human immunodeficiency virus (HIV). Eligible study
151 designs included all observational studies except case-control studies.⁽¹⁹⁾ We excluded reports if they
152 did not, in fact, involve a DAT (Table S2), or if they were reviews, abstracts (other than from the Union
153 World Conference on Lung Health), editorials, commentaries, news articles, or study protocols. Relevant
154 grey literature (such as ministry reports, technical papers, and preprints) was accepted if it met the

155 eligibility criteria. A complete list of the outcomes, population, intervention, and control groups of
156 interest, and all inclusion/exclusion criteria can be found in Table S3.

157 **Study Selection.** After de-duplication using EndNote (Version 20.2.1 – Clarivate, London UK),
158 five reviewers (M.Z., M.S., C.C., S.B., C.K., and N.F.) independently screened all titles and abstracts for
159 their relevance to DATs for TB treatment support, supported by Rayyan.ai (Rayyan, Cambridge
160 USA).(20,21) Potentially relevant studies underwent independent full-text review by the same
161 reviewers, for eligibility according to the inclusion criteria above. Each screening stage was conducted in
162 duplicate, by two reviewers blinded to each other’s assessment, with conflicts resolved by a senior
163 investigator (K.S., R.S., or K.F.). All references from and citations of each included publication were also
164 screened for inclusion using Google Scholar (Alphabet, Mountain View USA).(22)

165 **Data extraction.** For each included study, data were extracted into a pre-specified Excel
166 (Microsoft, Redmond USA) template by two independent reviewers (M.Z. and M.S., or C.K.) in parallel,
167 and subsequently compared for any discrepancies. Conflicts were resolved by consensus and discussion
168 with a third reviewer (K.S.) when necessary. Extracted data included study characteristics e.g., study
169 design, study setting (i.e., geographic location; inpatient or outpatient), participant characteristics, DAT
170 used, reference standard used, the approaches to classifying adherence for both the DAT and reference
171 standard (e.g., time frame and frequency of adherence assessment), and any important gaps noted by
172 the reviewers. For each study, we extracted all reported performance parameters e.g., true positives
173 (TP), true negatives (TN), false positives (FP), false negatives (FN), sensitivity, specificity, PPV, NPV,
174 accuracy, and AUC.

175 **Performance parameter definitions.** The performance of digital technologies for measuring TB
176 medication adherence can include an assessment of the technical performance or functioning of the
177 DAT, or the performance of the DAT for detecting human behaviour. Technical performance is assessed
178 in controlled conditions to determine whether the designated hardware or software of the DAT can

179 adequately detect doses known to be taken. The performance of a DAT for detecting human behaviour
 180 is assessed in real-world conditions to determine its ability to detect a person’s adherence to their
 181 medication during their treatment course. In either condition, DAT performance can be assessed per
 182 person or per medication dose. Table 1 contains the definitions of the performance parameters used in
 183 this review by the unit of assessment: dose vs. person.

184 **Table 1: Definitions of parameters used to describe the performance of DATs for measuring**
 185 **tuberculosis medication adherence.** Parameters are defined against a reference standard and by unit of
 186 assessment.

Performance		
Parameter	Definition – Dose	Definition - Person
Sensitivity	The percentage of doses that were classified as taken by the DAT, among those that were taken by the reference standard.	The percentage of participants who were classified as adherent by the DAT, among those who were adherent by the reference standard.
Specificity	The percentage of doses that were classified as not taken by the DAT, among those that were not taken by the reference standard.	The percentage of participants who were classified as nonadherent by the DAT, among those who were non-adherent by the reference standard.
Positive predictive value	Likelihood of a dose being taken by the reference standard among those classified as taken by the DAT.	Likelihood of a participant being adherent by the reference standard among those classified as being adherent by the DAT.
Negative predictive value	Likelihood of a dose not being taken by the reference standard among those classified as not being taken by the DAT.	Likelihood of a participant being non-adherent by the reference standard among those classified as being non-adherent by the DAT.
Accuracy	The percentage of doses that were classified as taken or not taken by both the DAT and the reference standard.	The percentage of participants who were classified as adherent or non-adherent by both the DAT and the reference standard.

187 **Data synthesis.** The extracted data were summarized in tabular form. Pre-specified subgroup
 188 analyses were performed where appropriate, addressing specific DATs, TB disease vs. infection, and
 189 groups at risk of unfavorable outcomes. Sensitivity and specificity estimates were displayed according to
 190 DAT type in forest plots created using RevMan (Version 5.4 – Cochrane, London UK).(23) We calculated
 191 pertinent performance parameters that were not directly reported when these could be derived from
 192 the underlying data. For parameters that were reported without binomial 95% confidence intervals (CI),

193 they were calculated, when possible, using the Clopper-Pearson Exact Method in R (Version 4.1.2 – GNU
194 Project).(24) This method was also used for articles that had repeated observations per individual, and
195 therefore does not account for within-individual clustering.

196 Publication bias was assessed qualitatively using Deek’s test for diagnostic accuracy studies.(25)
197 We created a funnel plot of the association between the diagnostic odds ratio (DOR) and the effective
198 sample size (ESS) of each study.(25) Quantitative assessment for publication bias using the associated
199 regression test of asymmetry could not be performed, given the small number of included studies.

200 **Quality assessment.** The included articles were assessed for risk of bias and applicability
201 concerns using the Quality Assessment of Diagnostic Accuracy Studies 2 (QUADAS-2) tool for primary
202 diagnostic accuracy studies.(26) This tool addresses participant selection, the performance of the index
203 test, the performance of the reference test, and the flow and timing of the tests. For each feature, three
204 scores could be used: low, unclear, or high. Two reviewers (M.Z. and M.S., or C.K.) independently
205 assessed each study for quality. Any conflicts were resolved by consensus and discussion with a third
206 reviewer (K.S.) when necessary. Quality assessment results are displayed in graphic and summary format
207 using RevMan (Version 5.4 – Cochrane, London UK).(23)

208 **GRADE assessment.** Finally, we rated the robustness of evidence using the Grading of
209 Recommendations, Assessment, Development, and Evaluation (GRADE) approach for diagnostic tests
210 and strategies.(27) The GRADEpro Guideline Development Tool (McMaster University, Canada &
211 Evidence Price, EU) was used for recommendations.(28) The risk of bias and indirectness was assessed
212 using the relevant elements of the QUADAS-2 checklist. Inconsistency and imprecision were assessed
213 based on outcome variability and confidence interval ranges across articles.

214 **Patient and public involvement:** Patients and the public were not specifically involved in the
215 design, conduct, reporting, or dissemination plans of our research.

216

217 **RESULTS** (1,775 words)

218 **Study selection.** Figure 1 illustrates the PRISMA 2020 flow chart, beginning with the records
219 identified by our search. After de-duplication, there were 5692 records identified, of which 5290 titles
220 and abstracts were not relevant to TB and DATs. Of the remaining 402 reports which underwent full-
221 text review, nine met our inclusion criteria. Four additional reports were identified by hand-searching,
222 for a total of 13 included reports.

223 **Overview of studies.** Of the 13 reports included, 4 involved medication sleeves with phone calls,
224 5 involved digital pillboxes, 2 involved ingestible sensors, and 1 each involved a mobile application and
225 an artificial intelligence (AI)-based technology for viewing videos of medication ingestion (known as
226 video observed therapy [VOT]). The study characteristics of the included articles can be found in Table
227 2. Ingestible sensors, VOT (AI result), and the mobile application were assessed for their technical
228 performance under controlled conditions, while medication sleeves with phone calls and digital
229 pillboxes were assessed under real-world conditions for their ability to detect adherence during TB
230 treatment. Detailed descriptions of the DATs and reference standards used for performance assessment
231 for all articles can be found in Table S4. Briefly, included articles used either isoniazid (INH) urine
232 metabolite tests, rifampicin (RIF) urine colour tests, pill counts, DOT, or healthcare provider-based
233 adherence reports to assess the performance of the DAT (Table 2). 12 articles reported on persons
234 treated for TB disease, while one considered the performance of a digital pillbox in persons treated for
235 TB infection. Two studies compared pill counts and urine drug tests with the digital pillbox as their
236 reference standard; we back-calculated performance parameters for the pillbox compared with pill
237 counts and urine tests using the primary data they reported (Table 2). One study investigating digital
238 pillboxes (Huan et al. 2012) was translated from Simplified Chinese to English using Google Translate
239 (Alphabet, Mountain View USA).(29)

Table 2: Characteristics of included studies assessing the performance DATs for measuring TB medication adherence.

Study ID	Country	Participants				Intervention			Reference Test		
		Participants ^{††}	Age (years)	N ^{‡‡}	HIV%	DAT	Duration	Adherence classification	Test	Frequency	Adherence classification
Browne et al. 2019 (30)	United States****	DS TB disease	Mean: 43 SD: 17	77	NR	Ingestible sensors	2-3 weeks	1 detected dose	DOT	9x/person	1 dose ingested ^{††}
Belknap et al. 2013 (31)	United States****	TB disease [‡]	Median: 44 Range: 22-79	30	10%	Ingestible sensors	2-3 weeks (10 visits)	1 detected dose	DOT	10x/person	1 dose ingested ^{††}
Scott et al. 2023 (32)	China*** South Africa*** Spain**** USA****	TB infection	Median: 36 IQR: 27-49	665	1%	Digital pillboxes	3 months	≥ 11/12 doses recorded	Pill count	1x/person	≤1 dose remaining
Bionghi et al. 2018 [§] (33)	South Africa***	MDR TB disease	Median: 42 IQR: 33.5-54	21	100%	Digital pillboxes	3 weeks	1 dose recorded	Pill count	3x/person	1 dose missing ^{††}
Huan et al. 2012 (34)	China***	DS TB disease	Mean: 46 SD: 17	319	NR	Digital pillboxes	1-6 months	1 dose recorded in prior 24h	RIF urine colour test	1x/person	Red colour change
van den Boogaard et al. 2011 (35)	Tanzania**	PTB and EPTB disease	Mean: 41 SD: 14	37	44%	Digital pillboxes ^{‡‡}	6 months	100% or ≥95% of doses recorded	Pill count INH urine test RIF urine colour test	12x/person 4x/person 4x/person	0 or ≤5% doses remaining Purple/blue colour change [^] Orange urine colour
Ruslami et al. 2008 (36)	Indonesia**	PTB disease	Median: 32 Range: 16-84	30	NR	Digital pillboxes ^{‡‡}	4 weeks	100% of doses recorded ⁺	Pill count	2x/person	0 doses remaining
Subbaraman et al. 2021 ^{§§} (12)	India**	DS PTB and EPTB disease	Median: 35 IQR: 25-45	608	47%	Medication sleeves with phone calls	1-6 months	Adherence: 2 or 3 reported doses over the 2 days prior to the urine test and the day of the test Nonadherence: 0 or 1 reported doses over the 2 days prior to urine test and the day of the test	INH urine test	1x/person	Purple/blue or green colour change [^]
Thomas et al. 2020 (9)	India**	DS PTB and EPTB disease	Median: 35 Range: 18-83	597	48%	Medication sleeves with phone calls	1-6 months	Adherence: 1 dose reported in the prior 6h to 48h Nonadherence: No dose reported in the prior 72h	INH urine test	1x/person	Purple/blue or green colour change [^]
Efo et al. 2021 (10)	Tanzania**	DS TB disease	Range: 25-44	197	NR	Medication sleeves with phone calls	6 months	Adherence: 1 dose reported in the prior 48h Nonadherence: No dose reported in the prior 48h	INH urine test	1x/person	Purple/blue or green colour change [^]
Alacapa et al. 2020 [†] (11)	Philippines**	NR	NR	103	NR	Medication sleeves with phone calls	~3 months	Adherence: Dose reported in the prior 48h ⁺ Nonadherence: No dose reported in the prior 48h	INH urine test	1x/person	Purple/blue or green colour change [^]
Sekandi et al. 2023 (37)	Uganda*	DS TB disease	Mean: 31 Range: 19-50	51	28%	Other: VOT (AI result)	N/A [‡]	1 detected dose	VOT (Provider result)	10x/person	Research team identifies ingestion in VOT ^{††}
Goodwin et al. 2022 [†] (38)	Argentina***	DS TB disease	NR	NR	NR	Other: Mobile App (Software result)	N/A [‡]	1 detected dose	Mobile App (Provider result)	NR	Treatment supporter/ research staff detects colour change in photo of INH Urine Test. ^{††}

*low income country **lower middle-income country, ***upper middle-income country, ****high-income country

AI artificial intelligence, DOT directly observed therapy, DS drug-sensitive, DR drug-resistant, EPTB extra-pulmonary tuberculosis, INH isoniazid, IQR interquartile range, MDR multidrug-resistant, N number of participants,

N/A not applicable, NR not reported, PTB pulmonary tuberculosis, RIF rifampicin, SD standard deviation, VOT video observed therapy

§ Inpatient cohort

§§ Same cohort as Thomas et al. 2020 (9)

† Abstract from the World Conference on Lung Health of the International Union Against Tuberculosis and Lung Disease (The Union).

†† If not specified, resistance/susceptibility details or TB type were not indicated.

‡ Potentially INH-resistant TB in some clients.

249
250
251
252
253
254

¥N indicates the number of participants analyzed as part of DAT performance. Participant characteristics may be different among this specific group but were not reported in the text.

Not applicable due to batch submitted videos or photos being used to assess performance of the AI or software

Digital pillbox was used as the reference standard of the published study. Performance outcomes of the pillbox were calculated from originally published values on sensitivity, specificity, PPV, and NPV.

+ Assumed based on study information.

++ Study considered number of total medication doses across participants rather than adherence per participant's treatment cycle.

^Colour change indicates the time of last medication ingestion: Purple/blue (<24 hours ago), Green (24-48 hours ago), Yellow (>48 to 72 hours ago)

255 Assessment of adherence by the reference standard was conducted anywhere from once to
256 twelve times over participants' treatment and the duration of DAT use ranged from two weeks to six
257 months (Table 2). For two studies, authors were contacted with requests to add missing data, but they
258 did not reply. Six of the full articles reported some but not all applicable performance parameters
259 (sensitivity, specificity, PPV, NPV, accuracy; Table S5). While two studies estimated area under the
260 receiver operating characteristic (ROC) curve (Table S5), none provided ROC curves, likelihood ratios, or
261 estimated agreement. When possible, we calculated missing performance parameters from the data
262 provided (Table S5). A meta-analysis pooling performance estimates was not conducted given the small
263 number and marked heterogeneity of the studies retrieved. Only one (30) of five studies
264 (30,31,33,37,38) that analyzed repeated observations of individuals accounted for clustering in their
265 primary analysis. Summary ROC curves to demonstrate joint sensitivity and specificity were not created
266 due to the small number of studies per DAT.

267 **Performance of DATs under controlled conditions.**

268 *Ingestible sensors.* Two reports from the US addressed the performance of ingestible sensors for
269 measuring TB medication adherence in controlled conditions (Table 2).(30,31) This involved determining
270 if an on-body wearable sensor could detect the ingestible sensor placed on the medication, following its
271 ingestion which was directly observed. The percentage of doses that were correctly classified as taken
272 by the sensor's technology (sensitivity) was at least 95% across studies (Figure 2a, Table S6). Only
273 sensitivity could be estimated from these studies.

274 *Other DATs.* Two reports assessed other DATs for their ability to measure TB medication
275 adherence under controlled conditions. In one, a deep convolutional neural network (DCNN) was used
276 to determine whether a participant had ingested medication in a video submitted as part of a VOT
277 intervention.(37) The performance of the artificial intelligence algorithm was compared to the research
278 team's interpretation of the same submitted video and assessed using 5-fold cross validation, for five

279 types of convolutional neural networks to extract features of the videos. The best performing
280 convolutional neural network had a mean sensitivity of 95% (standard deviation, SD 2.6) but a specificity
281 of 55% (SD 6.5) (Table S6).

282 In the other report., a mobile application involved the participant's uploading a photo of their
283 INH urine test, and associated software then analyzed the images to determine whether or not they
284 indicated the presence of INH metabolites.(38) The output from this reader software was compared to
285 treatment supporter and research staff interpretation of the same urine test images. The sensitivity was
286 81% and 86% when compared to the treatment supporter and research staff interpretations,
287 respectively, while the percentage of doses correctly classified by the technology as not taken
288 (specificity) was 95% and 91% (Table S6).

289 **Performance of DATs under real-world conditions.**

290 *Medication sleeves with phone calls ["99DOTS"]*. Three studies investigated the performance of
291 medication sleeves with phone calls for measuring TB medication adherence, compared with
292 unannounced INH urine tests either in the clinic or at a home visit (Table 2). They generally classified
293 adherence as any dose reported within 48 hours before the urine test and nonadherence as no reported
294 dose within the previous 48 or 72 hours before the urine test. (9–11) All studies investigating this DAT
295 were conducted in lower middle-income countries. The sensitivity of participants' dose reports (patient-
296 reported doses) for detecting adherence ranged from 70% to 94%, while the specificity of participants'
297 dose reports for detecting nonadherence ranged from 0% to 61% (Figure 2b).

298 Two articles on medication sleeves with phone calls, using data from the same cohort, showed
299 that sensitivity and specificity vary depending on the approach used to classify DAT adherence. For
300 example, when "nonadherence" was defined as a person with TB not reporting any medication dose
301 during the 72 hours prior to INH urine testing (9), the specificity was lower than when "nonadherence"

302 was defined as a person reporting less than either two or three medication doses in the 72 hours prior
303 to the home visit (Table 3).(12)

Table 3: Performance of medication sleeves with phone calls and digital pillboxes under real-world conditions against a reference standard. Digital pillboxes were assessed for persons with tuberculosis disease and tuberculosis infection. If not reported, values and binomial confidence intervals were calculated from other reported performance data (sensitivity, specificity, PPV, NPV, or number of TP, FP, FN, and TN).

Study ID	DAT Adherence classification	Reference Standard	N	TP	FP	FN	TN	Sensitivity (95% CI)	Specificity (95% CI)	PPV (95% CI)	NPV (95% CI)	Accuracy (95% CI)
Medication sleeves with phone calls ["99DOTS"]												
Subbaraman et al. 2021 [§] (12)	2 or 3 reported doses in 72h prior	INH urine test	608	330	20	211	47	61% (57-65%)	70% (58-81%)	94% (91-96%)	18% (14-23%)	62% (58-66%)
			608 [†]	482	45	59	22	89% (86-91%)	33% (22-45%)	91% (89-94%)	27% (18-38%)	83% (80-86%)
Efo et al. 2021 (10)	1 reported dose in prior 48h	INH urine test	197	138	7	41	11	77% (70-83%)	61% (36-83%)	95% (90-98%)	21% (12-34%)	76% (69-81%)
Thomas et al. 2020 (9)	1 reported dose in prior 6-48h	INH urine test	597	368	28	157	44	70% (66-74%)	61% (48-72%)	93% (91-95%)	21% (18-25%)	69% (65-72%)
			597 [†]	446	44	79	28	85% (81-88%)	39% (28-52%)	91% (90-93%)	26% (20-33%)	79% (76-82%)
			310 [‡]	155	17	83	32	65% (59-71%)	66% (51-79%)	90% (86-93%)	28% (23-34%)	65% (59-71%)
Alacapa et al. 2020 ^{§§} (11)	1 reported dose in prior 48h	INH urine test	103	95	2	6	0	94% (88-98%)	0% (0-84%)	98% (93-100%)	0% (0-46%)	92% (85-97%)
Digital pillboxes												
TB Disease												
Bionghi et al. 2018 ^{¥¥} (33)	1 dose recorded	Pill count	21 [†] (489)	480	0	6	3	99% (97-100%)	100% (29-100%)	100% (99-100%)	33% (7-70%)	99% (97-100%)
Huan et al. 2012 ^{§§} (34)	1 dose recorded in prior 24h	RIF urine colour test	319	208	5	1	105	100% (97-100%)	95% (90-99%)	98% (95-99%)	99% (95-100%)	98% (96-99%)
van den Boogaard et al. 2011 [†] (35)	≥ 95% doses recorded	Pill count	37	19	11	5	2	79% (58-93%)	15% (2-45%)	63% (44-80%)	29% (4-71%)	57% (39-73%)
		INH urine test	37	18	12	4	3	82% (60-95%)	20% (4-48%)	60% (41-77%)	43% (10-82%)	57% (39-73%)
	RIF urine colour test	37	8	22	1	6	89% (52-100%)	21% (8-41%)	27% (12-46%)	86% (42-100%)	38% (22-55%)	
	100% doses recorded	Pill count	37	6	4	18	9	25% (10-47%)	69% (39-91%)	60% (26-88%)	33% (17-54%)	41% (25-58%)
		INH urine test	37	4	6	18	9	18% (5-40%)	60% (32-84%)	40% (12-74%)	33% (17-54%)	35% (20-53%)
RIF urine colour test	37	2	8	7	20	22% (3-60%)	71% (51-87%)	20% (3-56%)	74% (54-89%)	59% (42-75%)		
Ruslami et al. 2008 [#] (36)	100% doses recorded ^{##}	Pill count	30	15	2	5	8	75% (51-91%)	80% (44-97%)	88% (64-99%)	62% (32-86%)	77% (58-90%)
TB Infection												
Scott et al. 2023 (32)	≥ 11/12 doses recorded	Pill count	665 (Total)	497	0	40	128	93% (90-95%)	100% (97-100%)	100% (99-100%)	76% (69-82%)	94% (92-96%)
			513 (USA)	395	0	17	101	96% (93-98%)	100% (96-100%)	100% (99-100%)	86% (78-91%)	97% (95-98%)
			57 (SA)	25	0	21	11	54% (39-69%)	100% (72-100%)	100% (86-100%)	34% (19-53%)	63% (49-76%)
			65 (Spain)	50	0	2	13	96% (87-100%)	100% (75-100%)	100% (93-100%)	87% (60-98%)	97% (89-100%)
			30 (China)	27	0	0	3	100% (87-100%)	100% (29-100%)	100% (87-100%)	100% (29-100%)	100% (88-100%)

CI confidence interval, DAT digital adherence technology, FN false negatives, FP false positives, INH isoniazid, N number of participants, NPV negative predictive value, PPV positive predictive value, RIF rifampicin, SA South Africa, TN true negatives, TP true positives, USA United States of America

§ Same cohort as Thomas et al. 2020

§§ Abstract from the World Conference on Lung Health of the International Union Against Tuberculosis and Lung Disease (The Union)

† Patient and provider reported doses (Patient-reported doses relies only on calls made by the participant. For patient and provider reported doses, if no call is made after 24h, healthcare workers are supposed to contact the participant and report whether these doses were taken based on verbal report)

†† Subpopulation of people with HIV within this cohort study.

¥ Subpopulation of people without HIV in within this cohort study.

316 ¥¥ Inpatient cohort
317 # The digital pillbox was used as the original reference standard. Number of true positives, false positives, false negatives, and true negatives were back-calculated from the article's primary data on
318 pill count or urinalysis sensitivity and specificity, and number of events.
319 ## Assumed based on study information
320 + 489 total doses, shown in parentheses, were assessed for performance across 21 participants.

321 The same authors also examined dose verification added by providers, who telephoned persons
322 who had not phoned in their doses on a given day. When provider-reported doses were added, the
323 sensitivity of combined patient-and provider dose reports for detecting adherence increased from 70%
324 [95% CI 66-74%] to 85% [95% CI 81-88%], but specificity for detecting nonadherence decreased from
325 61% [95% CI 48-72%] to 39% [95% CI 28-52%] (Table 3).(9)

326 Estimated overall accuracy ranged from 69-92% (Table 3). Nonetheless, the negative predictive
327 value, or likelihood of a participant being non-adherent by the urine test among those classified as being
328 non-adherent by the DAT, was consistently low (range 0-26%). In other words, non-reporting of a dose
329 by a person with TB often did not mean that they had in fact missed it—and instead often signaled
330 limited ongoing engagement with the technology.

331 *Digital pillboxes.* The five studies reporting on digital pillboxes used varying methods and were
332 conducted in countries with differing income levels (Table 2). Two used the digital pillbox as their
333 reference standard against multiple index tests, so we back-calculated performance estimates from
334 their reported values (Table S5).(35,36) One assessed performance in an inpatient cohort (33), while one
335 focused on TB infection (32). Four reports included pill counts conducted at clinic visits as at least one of
336 their reference standards (32,33,35,36), while one used unannounced urine color testing for rifampin at
337 home visits (34) (Table 2).

338 When compared against pill count for those with TB disease, three studies reported digital
339 pillbox sensitivities for overall perfect adherence ranging from 25% to 99%, while specificity for
340 detecting nonadherence was higher, ranging from 69% to 100% (Figure 2c).(33,35,36) The report which
341 only used urine color testing for rifampin estimated the sensitivity and specificity of digital pillboxes for
342 dose detection within 24h as 100% (95% CI 97-100%) and 95% (90-99%), respectively (Table 3).(34) In
343 one article, when compared against pill count among persons treated for TB infection, digital pillboxes

344 were very sensitive for adherence (93% [95% CI 90-95%]) and specific for non-adherence (100% [97-
345 100%]) and generally remained so in subgroup analyses involving the various study sites (Table 3).(32)

346 For digital pillboxes, overall accuracy estimates ranged from 35%-100% (Table 3). As with the
347 medication sleeves, negative predictive values were often poor, meaning that doses not reported by
348 pillbox opening were not necessarily missed.

349 **Subgroups and special populations.** In an inpatient cohort of 21 individuals with both HIV and
350 drug-resistant TB, digital pillboxes had high sensitivity and specificity (>99%) (Table 3).(33) In another
351 study investigating medication sleeves with phone calls, a subgroup analysis of people with HIV treated
352 for TB disease found that specificity was higher (66% [95% CI 51-79%]) among people with HIV (i.e.,
353 more nonadherence was correctly detected) than among people without HIV (48% [26-70%]) (Table
354 3).(9) Several other studies included persons living with HIV among their participants, but did not report
355 results separately for this group (Table 2).(31,35) No report addressed persons younger than 18.

356 **Quality Assessment.** The risk of bias was high or unclear in at least two categories for all but one
357 of the reports assessed. (Figure S1a). Two reports could not be assessed for quality as the limited
358 description of their methods precluded formal assessment of bias.(11,38) Participant selection was a
359 frequent source of potential bias since participants were rarely randomly sampled; in some cases, they
360 had also explicitly consented to take part in DAT implementation projects (10), which could also
361 introduce bias. Timing of DAT performance assessment relative to reference standard measures may
362 also have been an issue. There were fewer concerns about applicability with respect to participant
363 profiles, DATs used, or reference standards (Figure S1b). Details of the quality assessment criteria are
364 provided in Table S7.

365 **Publication bias.** There was visual asymmetry in the funnel plot displaying diagnostic odds
366 ratios and effective sample sizes, indicating possible publication bias (Figure S2). There were too few
367 studies to permit formal quantitative evaluation of publication bias. Additionally, the effective sample

368 sizes do not account for within-individual clustering and are therefore overestimated for these articles
369 (33,37). Several could not be included because of insufficient data to estimate the diagnostic odds
370 ratios.

371 **GRADE Results.** The findings suggested moderate certainty of evidence for the sensitivity of
372 ingestible sensors, because of consistent and precise results, a low risk of bias, but an inability to assess
373 publication bias. The certainty of evidence was very low for medication sleeves with phone calls because
374 of substantial risks of bias, and substantial variation in sensitivity and specificity across articles. Similarly,
375 the certainty of evidence was very low for digital pillboxes. More details on the quality of evidence for
376 each DAT can be found in Table S8-10 and Figure S3. A GRADE assessment was not conducted for the
377 DATs using AI and reader software since these involved only one report each.

378 **DISCUSSION** (1,126 words, including conclusion)

379 Available evidence addressing the performance of DATs for measuring TB medication adherence
380 is limited in scope and quality. Among 13 reports which considered five types of DATs, there was
381 substantial variation in reference standards, definitions of adherence, and adherence classifications,
382 which made it difficult to compare and summarize results. Ingestible sensors showed relatively high
383 sensitivity under experimental conditions in high-income countries, but this technology may be too
384 expensive to scale up in the near future and will face substantial barriers to real world implementation
385 in LMICs with high TB incidence. In general, however, studies of DAT implementation in real world
386 settings (medication sleeves with phone calls, digital pillboxes) found suboptimal performance of these
387 technologies for measuring medication adherence. These findings are concordant with a previous
388 review that assessed the accuracy of electronic monitoring in devices in antiretroviral medication
389 adherence.(39)

390 To understand the suboptimal performance of DATs during real world implementation, it is
391 important to consider that this reflects the ways in which people undergoing TB treatment do or do not

392 engage with the technology. Engagement may be influenced by technical attributes of the DAT and the
393 motivations, beliefs, social context, and structural barriers faced by people with TB. Suboptimal
394 engagement with DATs—which may result in under-reporting of true medication adherence (i.e.,
395 reduced sensitivity, reduced NPV)—can be influenced by limited access to mobile networks and
396 technology, shared cellphone use among multiple household members, and inadequate education on
397 the purpose and appropriate use of the DAT.(40–43) With substantial under-reporting of medication
398 adherence, healthcare providers may have difficulty identifying people who are truly experiencing
399 nonadherence, be unable to routinely reach out to all people with reported nonadherence (e.g., via
400 phone calls or home visits) (44), and may begin to ignore digital adherence data due to its limited
401 performance (42,44).

402 Conversely, over-reporting of adherence (i.e., reduced specificity, reduced PPVs)—in which
403 people with TB or healthcare providers may report adherence via phone call or pillbox opening despite
404 medication doses not actually ingested—may result from the desire by people with TB to conceal
405 nonadherence or by the desire of healthcare providers to report optimal TB outcomes.(45) Over-
406 reporting of adherence may be more concerning, as healthcare providers and health systems may miss
407 early identification of people who need intensified or personalized support to improve adherence and
408 ensure optimal treatment outcomes.

409 Several challenges limited our interpretation of DAT performance. The certainty around
410 evidence of performance for DATs under programmatic conditions was very limited. Additionally, the
411 use of different pill taking thresholds to define adherence, and of different windows for dose recording
412 vs. urine testing, adds substantial complexity and makes it even more difficult to summarize the
413 available evidence. One study evaluated digital pillboxes in an inpatient setting, which differs
414 substantially from the settings where DATs would most likely be used and limits the relevance of the
415 resulting data. Another used a method for documenting the presence or absence of rifampin in urine,

416 which was not used in other studies, highlighting the importance of consistent reference standards.
417 More generally, the use of urine tests as reference standards can be problematic, given known gaps in
418 their performance. This includes suboptimal sensitivity (based on witnessed pill ingestion) for doses
419 ingested over 24 hours before testing (46), and suboptimal specificity for doses ingested over 72 hours
420 before testing.(12,47) Similarly, while pill count is a standard method for determining adherence, pills
421 missing during routine pill counts may not have been ingested.

422 **Strengths and Limitations.** To our knowledge, this is the first systematic review focusing
423 specifically on the performance of DATs for measuring TB medication adherence. We used a wide-
424 ranging, pre-specified search strategy across multiple databases, without language restriction and with
425 the potential inclusion of suitable reports from grey literature, including preprints. Our search was
426 developed and conducted by an experienced health sciences librarian. Selection of reports and data
427 extraction were performed rigorously, by two independent reviewers at every step. Whenever possible,
428 we used primary data from each report to estimate any relevant performance parameters that the
429 authors had not reported. We also assessed the quality and robustness of the available evidence.

430 For logistical reasons, we could not systematically search for and retrieve abstracts other than
431 those presented at the conferences of the International Union Against Tuberculosis and Lung Disease,
432 which is the major venue for public health-oriented TB research. Some reports did not include sufficient
433 data to permit estimation of all performance parameters of interest. Additionally, there may be an
434 underestimation of the variance in instances where clustering was not considered (i.e., for repeated
435 measurements per individual). We did not identify any reports evaluating the performance of doses
436 observed by video against urine tests or pill counts, although video-supported treatment has been used
437 increasingly in high-income countries.(48–50) Finally, given the heterogeneity of study methods,
438 technologies, and results, we did not formally pool or meta-analyze their data. We have instead

439 summarized the key results from each study in tabular and/or graphic format, and provided an overview
440 of the results for each key technology.

441 **Changes to the protocol.** Changes were made to our pre-registered protocol in PROSPERO as
442 follows. We additionally searched Europe PMC for pre-print articles relevant to our review. We did not
443 anticipate the use of DOT, or provider-based DAT decisions as reference standards for quantifying DAT
444 performance, and therefore added them as possible comparators. We also did not anticipate articles
445 that used the DAT as the reference standard. For those reports, we back-calculated the published results
446 to provide us with the performance of the DAT as the index test.

447 **CONCLUSIONS**

448 Accurate dose reporting is fundamental to the use of digital technologies to aid treatment
449 adherence. Their use as alternatives to traditional direct observation is predicated on the concept that
450 providers and health programs can provide additional support to people who appear to be facing
451 challenges in this regard. Hence suboptimal performance of dose reports from DATs can potentially
452 compromise their effectiveness, as well as program efficiency. If the DAT fails to capture significant non-
453 adherence, this leads to missed opportunities for intervention, and potentially even poorer treatment
454 outcomes. On the other hand, if support and supervision are intensified for people wrongly labeled as
455 poorly adherent by the DAT, this is a waste of limited program resources. The future evidence base will
456 be strengthened by more consistent definitions and cutoffs for adherence, and by more consistent use
457 of one or more reference standards—e.g. validated methods for pill counts, standardized timing and
458 technique for urine tests. Future research should also address the interplay of specific technology,
459 setting, user characteristics, and user engagement with DAT performance. Additional studies examining
460 the performance of asynchronous video observation will also be important, as will further investigation
461 of digital technologies to support newer treatment regimens for TB disease and infection.

462 **FIGURE LEGENDS**

463 **Figure 1: PRISMA 2020 Flow Diagram of studies identified and included.** Studies were identified via
464 databases and registries, and via other methods related to the performance of DATs for measuring TB
465 medication adherence. Other methods included searching the references and citations of included
466 studies identified from databases and registries, searching the references and citations of the studies
467 included our linked systematic reviews, and searching the abstracts of the Union World Conference of
468 Lung Health and Tuberculosis from 2004 to 2022. Database and registry searches were conducted on
469 April 28, 2023 from January 1, 2000 to April 28, 2023.

470 **Figure 2: Forest plot of the sensitivity and specificity of a) ingestible sensors, b) Medication sleeves**
471 **with phone calls [“99DOTS”], and c) Digital pillboxes. a) DOT** was used as the reference standard.
472 Sensitivity was calculated using the reported positive detection accuracy of each study (# of ingestible
473 sensors detected/# of ingestible sensors ingested). Specificity was not estimable due to the lack of false
474 positives and true negatives. Browne et al. (30) confidence intervals account for within-individual
475 clustering. **b) INH urine test** was used as the reference standard. DAT adherence was recorded using
476 “patient-reported doses”, wherein dose reporting relies only on calls made by the person receiving
477 treatment. **c) Pill count** was used as the reference standard. Results are depicted for persons with
478 tuberculosis disease or infection. Bionghi et al. (33) was an inpatient cohort. Scott et al. (32) used a
479 cohort of persons with TB infection (TBI). In Ruslami et al. (36) and van den Boogaard et al. (35), the
480 digital pillbox was used as the reference standard. Number of true positives, false positives, false
481 negatives, and true negatives were back-calculated from the article’s primary data on pill count
482 sensitivity and specificity, and number of events.

483

484

485 **DECLARATIONS**

486 **Ethics approval and consent to participate:** Not applicable

487 **Consent for publication:** Not applicable

488 **Availability of data and materials:** The data and materials supporting the conclusions of this article that
489 are not already shared in the manuscript and appendix are available at

490 <https://borealisdata.ca/dataverse/mcgill>

491 **Competing interests:** Dr. Kevin Schwartzman reports research funding from the Canadian Institutes of
492 Health Research. He has also served as chair of the Data Safety and Monitoring Board for a COVID-19
493 therapeutic investigated by Laurent Pharmaceutical.

494 **Funding:** This review was supported by a grant from the Bill and Melinda Gates Foundation (grant INV-
495 038215). The funder had no role in the execution or reporting of this study, in manuscript preparation,
496 or in the decision to publish.

497 **Authors' contributions:** All authors contributed to the design of the research and approved the final
498 manuscript. MZ contributed to screening, data extraction and analysis, and drafted the manuscript. MS
499 contributed to the database search, screening, data extraction, and provided systematic review
500 expertise for the manuscript. CK contributed to screening and data extraction. GG contributed to the
501 database search. CC, SB, and NF contributed to screening. KF, RS, and KS provided support in all steps of
502 executing the review and revised the manuscript.

503 **Acknowledgements:** Part of this work was presented in abstract and poster format at the 2023 North
504 America Region Conference of the International Union of Tuberculosis and Lung Disease on Tuberculosis
505 (51).

506

507

508

509

510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535

REFERENCES

1. Global Tuberculosis Report [Internet]. World Health Organization. 2021 [cited 2022 Jul 27]. Available from: <https://www.who.int/publications/i/item/9789240037021>
2. Subbaraman R, de Mondesert L, Musiimenta A, Pai M, Mayer KH, Thomas BE, et al. Digital adherence technologies for the management of tuberculosis therapy: mapping the landscape and research priorities. *BMJ Glob Heal*. 2018;3(5):e001018.
3. Li J, Munsiff SS, Tarantino T, Dorsinville M. Adherence to treatment of latent tuberculosis infection in a clinical population in New York City. *Int J Infect Dis [Internet]*. 2010;14(4):e292–7. Available from: <https://www.sciencedirect.com/science/article/pii/S1201971209002124>
4. Chung SJ, Lee H, Koo GW, Min JH, Yeo Y, Park DW, et al. Adherence to nine-month isoniazid for latent tuberculosis infection in healthcare workers: a prospective study in a tertiary hospital. *Sci Rep*. 2020 Apr;10(1):6462.
5. World Health Organization. WHO consolidated guidelines on tuberculosis: Module 4: Treatment, Tuberculosis care and support [Internet]. WHO Press. 2022. 98 p. Available from: <https://tbksp.org/en/node/1898>
6. Karumbi J, Garner P. Directly observed therapy for treating tuberculosis. *Cochrane Database Syst Rev*. 2015;2015(5).
7. Sagbakken M, Frich JC, Bjune GA, Porter JDH. Ethical aspects of directly observed treatment for tuberculosis: a cross-cultural comparison. *BMC Med Ethics [Internet]*. 2013;14(1):25. Available from: <https://doi.org/10.1186/1472-6939-14-25>
8. Yellappa V, Lefèvre P, Battaglioli T, Narayanan D, Van der Stuyft P. Coping with tuberculosis and directly observed treatment: a qualitative study among patients from South India. *BMC Health Serv Res [Internet]*. 2016;16(1):283. Available from: <https://doi.org/10.1186/s12913-016-1545-9>

- 536 9. Thomas BE, Kumar JV, Chiranjeevi M, Shah D, Khandewale A, Thiruvengadam K, et al. Evaluation
537 of the Accuracy of 99DOTS, a Novel Cellphone-based Strategy for Monitoring Adherence to
538 Tuberculosis Medications: Comparison of Digital Adherence Data with Urine Isoniazid Testing. *Clin*
539 *Infect Dis*. 2020;71(9):E513–6.
- 540 10. Efo E, Onjare B, Shilugu L, Levy J. Acceptability, feasibility and accuracy of 99DOTS adherence
541 technology in mining region of Tanzania. *2021 IST-Africa Conf IST-Africa 2021*. 2021;(5):1–12.
- 542 11. Alacapa J, Morales M, Levy J, Powers R, Villaneuva A. Evaluation of the accuracy of 99DOTS digital
543 adherence technology for tuberculosis in Metro Manila, the Philippines. *World Conf Lung Heal Int*
544 *Union Against Tuberc Lung Dis (The Union)*. 2020;24.
- 545 12. Subbaraman R, Thomas BE, Kumar JV, Lubeck-Schricker M, Khandewale A, Thies W, et al.
546 Measuring tuberculosis medication adherence: a comparison of multiple approaches in relation
547 to urine isoniazid metabolite testing within a cohort study in India. *Open Forum Infect Dis*.
548 2021;1–8.
- 549 13. Wang N, Shewade HD, Thekkur P, Zhang H, Yuan Y, Wang X, et al. Do electronic medication
550 monitors improve tuberculosis treatment outcomes? Programmatic experience from China. *PLoS*
551 *One* [Internet]. 2020;15(11 November):1–11. Available from:
552 <http://dx.doi.org/10.1371/journal.pone.0242112>
- 553 14. Chen AZ, Kumar R, Baria RK, Shridhar PK, Subbaraman R, Thies W. Impact of the 99DOTS digital
554 adherence technology on tuberculosis treatment outcomes in North India: a pre-post study. *BMC*
555 *Infect Dis* [Internet]. 2023;23(1):1–10. Available from: [https://doi.org/10.1186/s12879-023-](https://doi.org/10.1186/s12879-023-08418-2)
556 [08418-2](https://doi.org/10.1186/s12879-023-08418-2)
- 557 15. Ngwatu BK, Nsengiyumva NP, Oxlade O, Mappin-Kasirer B, Nguyen NL, Jaramillo E, et al. The
558 impact of digital health technologies on tuberculosis treatment: a systematic review. *Abubakar I,*
559 *Alipanah N, Bastos M, Boccia D, Chin D, Cohen T, et al., editors. Eur Respir J* [Internet]. 2018;51.

- 560 Available from: <https://erj.ersjournals.com/content/51/1/1701596>
- 561 16. Nglazi MD, Bekker LG, Wood R, Hussey GD, Wiysonge CS. Mobile phone text messaging for
562 promoting adherence to anti-tuberculosis treatment: a systematic review. *BMC Infect Dis.* 2013
563 Dec;13:566.
- 564 17. Ridho A, Alfian SD, van Boven JFM, Levita J, Yalcin EA, Le L, et al. Digital Health Technologies to
565 Improve Medication Adherence and Treatment Outcomes in Patients With Tuberculosis:
566 Systematic Review of Randomized Controlled Trials. *J Med Internet Res.* 2022 Feb;24(2):e33062.
- 567 18. Schwartzman K, Mappin-Kasirer B, Mohamed MS, Zary M, Kafie C, Subbaraman R, et al.
568 PROSPERO 2022 CRD42022313526: The Accuracy of Dose Reports generated by Digital
569 Adherence Technologies for Persons treated for Tuberculosis Disease or Infection: A Systematic
570 Review and Meta-Analysis [Internet]. 2022. Available from:
571 https://www.crd.york.ac.uk/prospero/display_record.php?ID=CRD42022313526
- 572 19. Rutjes AWS, Reitsma JB, Vandenbroucke JP, Glas AS, Bossuyt PMM. Case-control and two-gate
573 designs in diagnostic accuracy studies. *Clin Chem.* 2005;51(8):1335–41.
- 574 20. EndNote v20.2.1 [Internet]. Clarivate; 2023. Available from: <https://endnote.com/>
- 575 21. Rayyan - AI Powered Tool for Systematic Literature Reviews [Internet]. Rayyan; 2022. Available
576 from: <https://www.rayyan.ai/>
- 577 22. Google Scholar [Internet]. 2023 [cited 2022 Aug 2]. Available from: <https://scholar.google.com/>
- 578 23. Review Manager (RevMan). The Cochrane Collaboration; 2020.
- 579 24. RStudio Team. RStudio: Integrated Development for R. Boston, MA: RStudio, PBC; 2020.
- 580 25. Van Enst WA, Ochodo E, Scholten RJ, Hooft L, Leeflang MM. Investigation of publication bias in
581 meta-analyses of diagnostic test accuracy: A meta-epidemiological study. *BMC Med Res*
582 *Methodol.* 2014;14(1):1–11.
- 583 26. Whiting PF, Rutjes AWS, Westwood ME, Mallett S, Deeks JJ, Reitsma JB, et al. QUADAS-2: A

- 584 Revised Tool for the Quality Assessment of Diagnostic Accuracy Studies. *Ann Intern Med*
585 [Internet]. 2011 Oct 18;155(8):529–36. Available from:
586 <https://www.acpjournals.org/doi/abs/10.7326/0003-4819-155-8-201110180-00009>
- 587 27. Schünemann H, Brożek J, Guyatt G, Oxman A. The GRADE Working Group. 2013 [cited 2023 Apr
588 14]. GRADE Handbook for grading quality of evidence and strength of recommendations.
589 Available from: <https://gdt.gradeapro.org/app/handbook/handbook.html#h.f7lc8w9c3nh8>
- 590 28. McMaster University and Evidence Prime. GRADEpro GDT: GRADEpro Guideline Development
591 Tool [Internet]. 2022. Available from: gradeapro.org
- 592 29. Google Translate [Internet]. 2023 [cited 2023 Aug 2]. Available from: <https://translate.google.ca/>
- 593 30. Browne SH, Umlauf A, Tucker AJ, Low J, Moser K, Garcia JG, et al. Wirelessly observed therapy
594 compared to directly observed therapy to confirm and support tuberculosis treatment
595 adherence: A randomized controlled trial. *PLoS Med*. 2019;16(10):1–19.
- 596 31. Belknap R, Weis S, Brookens A, Au-Yeung KY, Moon G, DiCarlo L, et al. Feasibility of an Ingestible
597 Sensor-Based System for Monitoring Adherence to Tuberculosis Therapy. *PLoS One*. 2013;8(1):8–
598 12.
- 599 32. Scott NA, Sadowski C, Vernon A, Arevalo B, Beer K, Borisov A, et al. Using a medication event
600 monitoring system to evaluate self-report and pill count for determining treatment completion
601 with self-administered, once-weekly isoniazid and rifapentine. *Contemp Clin Trials*. 2023
602 Jun;129:107173.
- 603 33. Bionghi N, Daftary A, Maharaj B, Msibi Z, Amico KR, Friedland G, et al. Pilot evaluation of a
604 second-generation electronic pill box for adherence to Bedaquiline and antiretroviral therapy in
605 drug-resistant TB/HIV co-infected patients in KwaZulu-Natal, South Africa. *BMC Infect Dis*.
606 2018;18(1):1–9.
- 607 34. HUAN Shi-tong LIU Xiao-qiu, OU Xi-chao, JIANG Shi-wen, ZHAO Yan-lin, ZHANG Zhi-ying, ZHAN Si-

- 608 yan CR. Operational feasibility of medication monitors in monitoring treatment adherence
609 among TB patients [Internet]. Vol. 34, Chinese Journal of Antituberculosis. p. 419–24. Available
610 from: <http://www.zgflzz.cn>
- 611 35. van den Boogaard J, Lyimo RA, Boeree MJ, Kibiki GS, Aarnoutse RE. Electronic monitoring of
612 treatment adherence and validation of alternative adherence measures in tuberculosis patients:
613 A pilot study. *Bull World Health Organ.* 2011;89(9):632–9.
- 614 36. Ruslami R, Van Crevel R, Van De Berge E, Alisjahbana B, Aarnoutse RE. A step-wise approach to
615 find a valid and feasible method to detect non-adherence to tuberculosis drugs. *Southeast Asian
616 J Trop Med Public Health.* 2008;39(6):1083–7.
- 617 37. Sekandi JN, Shi W, Zhu R, Kaggwa P, Mwebaze E, Li S. Application of Artificial Intelligence to the
618 Monitoring of Medication Adherence for Tuberculosis Treatment in Africa: Algorithm
619 Development and Validation. *JMIR AI [Internet].* 2023;2:e40167. Available from:
620 <https://ai.jmir.org/2023/1/e40167>
- 621 38. Goodwin K, Liao Z, Iribarren S. Assessing and refining image analysis software to automate
622 reading of objective home-based TB adherence tests. *World Conf Lung Heal Int Union Against
623 Tuberc Lung Dis (The Union).* 2022;26.
- 624 39. Smith R, Villanueva G, Probyn K, Sguassero Y, Ford N, Orrell C, et al. Accuracy of measures for
625 antiretroviral adherence in people living with HIV. *Cochrane Database Syst Rev.* 2022;2022(7).
- 626 40. Thomas BE, Kumar JV, Periyasamy M, Khandewale AS, Hephzibah Mercy J, Raj EM, et al.
627 Acceptability of the Medication Event Reminder Monitor for Promoting Adherence to Multidrug-
628 Resistant Tuberculosis Therapy in Two Indian Cities: Qualitative Study of Patients and Health Care
629 Providers. *J Med Internet Res [Internet].* 2021;23(6):e23294. Available from:
630 <https://www.jmir.org/2021/6/e23294>
- 631 41. Thomas BE, Kumar JV, Onongaya C, Bhatt SN, Galivanche A, Periyasamy M, et al. Explaining

- 632 Differences in the Acceptability of 99DOTS, a Cell Phone–Based Strategy for Monitoring
633 Adherence to Tuberculosis Medications: Qualitative Study of Patients and Health Care Providers.
634 JMIR Mhealth Uhealth [Internet]. 2020;8(7):e16634. Available from:
635 <http://mhealth.jmir.org/2020/7/e16634/>
- 636 42. Bahukudumbi S, Chilala C, Mohamed M, Zary M, Kafie C, Foster N, et al. Contextual factors
637 impacting the implementation of TB digital adherence technologies: A scoping review. In: World
638 Conference on Lung Health 2023 of the International Union Against Tuberculosis and Lung
639 Disease (The Union) [Internet]. Paris, France; 2023. p. S124-125. Available from:
640 https://conf2023.theunion.org/wp-content/uploads/2023/12/UNION2023_Abstracts.pdf
- 641 43. Chilala C, Bahukudumbi S, Foster N, Mohamed M, Zary M, Kafie C, et al. Implementation
642 outcomes of TB digital adherence technologies: A scoping review using the RE-AIM Framework.
643 In: World Conference on Lung Health 2023 of the International Union Against Tuberculosis and
644 Lung Disease (The Union) [Internet]. Paris, France; 2023. p. S97–8. Available from:
645 https://conf2023.theunion.org/wp-content/uploads/2023/12/UNION2023_Abstracts.pdf
- 646 44. Thekkur P, Kumar AM V, Chinnakali P, Selvaraju S, Bairy R, Singh AR, et al. Outcomes and
647 implementation challenges of using daily treatment regimens with an innovative adherence
648 support tool among HIV-infected tuberculosis patients in Karnataka, India: a mixed-methods
649 study. Glob Health Action [Internet]. 2019 Jan 1;12(1):1568826. Available from:
650 <https://doi.org/10.1080/16549716.2019.1568826>
- 651 45. Thomas B, Kumar V, Chiranjeevi M, Ramachandran G, Murugesan P, Khandewale A, et al.
652 Understanding challenges TB patients face in using digital adherence technologies. In: 50th
653 World Conference on Lung Health of the International Union Against Tuberculosis and Lung
654 Disease (The Union). Hyderabad, India; 2019. p. S236.
- 655 46. Soobratty MR, Whitfield R, Subramaniam K, Grove G, Carver A, O’Donovan G V, et al.

- 656 Point-of-care urine test for assessing adherence to isoniazid treatment for tuberculosis. Eur
657 Respir J [Internet]. 2014 May 1;43(5):1519 LP – 1522. Available from:
658 <http://erj.ersjournals.com/content/43/5/1519.abstract>
- 659 47. Hanifa Y, Mngadi K, Lewis J, Fielding K, Churchyard G, Grant AD. Evaluation of the Arkansas
660 method of urine testing for isoniazid in South Africa. Int J Tuberc lung Dis Off J Int Union against
661 Tuberc Lung Dis. 2007 Nov;11(11):1232–6.
- 662 48. Chuck C, Robinson E, Macaraig M, Alexander M, Burzynski J. Enhancing management of
663 tuberculosis treatment with video directly observed therapy in New York City. Int J Tuberc lung
664 Dis Off J Int Union against Tuberc Lung Dis. 2016 May;20(5):588–93.
- 665 49. Lam CK, McGinnis Pilote K, Haque A, Burzynski J, Chuck C, Macaraig M. Using Video Technology
666 to Increase Treatment Completion for Patients With Latent Tuberculosis Infection on 3-Month
667 Isoniazid and Rifapentine: An Implementation Study. J Med Internet Res. 2018 Nov;20(11):e287.
- 668 50. Beeler Asay GR, Lam CK, Stewart B, Mangan JM, Romo L, Marks SM, et al. Cost of Tuberculosis
669 Therapy Directly Observed on Video for Health Departments and Patients in New York City; San
670 Francisco, California; and Rhode Island (2017-2018). Am J Public Health. 2020 Nov;110(11):1696–
671 703.
- 672 51. Zary M, Mohamed MS, Kafie C, Chilala CI, Bahukudumbi S, Fielding K, et al. The Accuracy of Dose
673 Reports Generated by Digital Adherence Technologies for Persons Treated for Tuberculosis
674 Disease or Infection: A Systematic Review. In: 27th Annual Conference of the Union-North
675 America Region Conference on Tuberculosis [Internet]. 2023. p. 18. Available from:
676 <https://bclung.ca/sites/default/files/2023 NAR Conference Abstracts.pdf>
677

Identification of studies via databases and registers

Identification of studies via other methods

Identification

Records identified from:
CENTRAL (n= 1190)
CINAHL (n = 408)
EMBASE (n= 2370)
MEDLINE (n= 1437)
MedRxiv via Europe PMC (n= 380)
WoS (n= 1577)
Clinicaltrials.gov (n= 384)

Records removed *before screening*:
Duplicate records removed
(n= 2054)

Records identified from:
Citation searching (n = 1)
Reference searching (n = 1)
Scoping review searching (n = 0)
Union Conference Abstracts (n = 2)

Screening

Records screened
(n = 5692)

Records excluded (n= 5290)

Reports sought for retrieval
(n = 402)

Reports not retrieved (n= 0)

Reports sought for retrieval
(n = 4)

Reports assessed for eligibility
(n = 402)

Reports excluded: (n= 393)
Did not meet DAT definition (n= 137)
Wrong publication type
(n= 157)
Systematic review (n= 15)
Protocol (n= 41)
Clinical Trial (n= 50)
Review (n= 24)
Abstract (n= 23)
Letter to the editor (n= 2)
Editorials (n= 2)
Duplicate Publication (n= 3)
Wrong outcome (n= 96)

Reports assessed for eligibility
(n = 4)

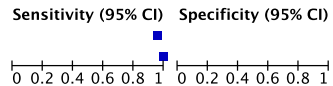
Included

Studies included in review (n=13)
Database Search (n= 9)
Supplementary Search (n = 4)

Studies included in review (n=13)
Database Search (n= 9)
Supplementary Search (n = 4)

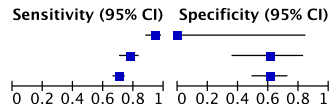
a) Ingestible sensors

Study	TP	FP	FN	TN	Reference Standard	DAT Adherence Classification	Sensitivity (95% CI)	Specificity (95% CI)
Belknap et al. 2013	1026	0	54	0	DOT	1 detected dose	0.95 [0.94, 0.96]	Not estimable
Browne et al. 2019	680	0	5	0	DOT	1 detected dose	0.99 [0.98, 1.00]	Not estimable



b) Medication sleeves with phone calls ["99DOTS"]

Study	TP	FP	FN	TN	Reference Standard	DAT Adherence Classification	Sensitivity (95% CI)	Specificity (95% CI)
Alacapa et al. 2020 (Abstract)	95	2	6	0	INH Urine Test	1 reported dose in prior 48h	0.94 [0.88, 0.98]	0.00 [0.00, 0.84]
Efo et al. 2021	138	7	41	11	INH Urine Test	1 reported dose in prior 48h	0.77 [0.70, 0.83]	0.61 [0.36, 0.83]
Thomas et al. 2020	368	28	157	44	INH Urine Test	1 reported dose in prior 6-48h	0.70 [0.66, 0.74]	0.61 [0.49, 0.72]



c) Digital pillboxes

Study	TP	FP	FN	TN	Reference Standard	DAT Adherence Classification	Sensitivity (95% CI)	Specificity (95% CI)
(TBI) Scott et al. 2023	497	0	40	128	Pill Count	≥ 11/12 doses recorded	0.93 [0.90, 0.95]	1.00 [0.97, 1.00]
Bionghi et al. 2018	480	0	6	3	Pill Count	1 dose recorded	0.99 [0.97, 1.00]	1.00 [0.29, 1.00]
Ruslami et al. 2008	15	2	5	8	Pill Count	100% doses recorded	0.75 [0.51, 0.91]	0.80 [0.44, 0.97]
van den Boogaard et al. 2011	6	4	18	9	Pill Count	100% doses recorded	0.25 [0.10, 0.47]	0.69 [0.39, 0.91]

