

Promise versus Reality: Optimism Bias in Package Inserts for Tuberculosis Diagnostics

Claudia M. Denkinger,^a Jasmine Grenier,^b Jessica Minion,^c and Madhukar Pai^{d,e}

Department of Medicine, Beth Israel Deaconess Medical Center, Boston, Massachusetts, USA^a; Faculty of Medicine, McGill University, Montreal, Quebec, Canada^b; Department of Medical Microbiology & Immunology, University of Alberta, Edmonton, Alberta, Canada^c; Department of Epidemiology, Biostatistics, and Occupational Health, McGill University, Montreal, Quebec, Canada^d; and Respiratory Epidemiology & Clinical Research Unit, Montreal Chest Institute, Montreal, Quebec, Canada^e

Laboratorians and clinicians often rely on package inserts of diagnostic tests to assess their accuracy. We compared test accuracy for tuberculosis diagnostics reported in 19 package inserts against estimates in published meta-analyses and found that package inserts generally report overoptimistic accuracy estimates. However, package inserts of most tests approved by the U.S. Food and Drug Administration (FDA) or endorsed by the World Health Organization provide more realistic estimates that agree with meta-analyses.

Most laboratory professionals anticipate that test performance when applied to patient care may be less impressive than what is reported in package inserts by test manufacturers. However, this gap between promise and reality has not been examined in a systematic way.

One approach for studying this gap is to compare sensitivity and specificity estimates in package inserts with pooled estimates from published systematic reviews and meta-analyses, which often include studies from diverse settings and may provide a more realistic assessment of test accuracy.

We conducted such a comparison using commercial tuberculosis (TB) diagnostics as a case study. The TB diagnostic pipeline has rapidly expanded, and a large number of meta-analyses have been published on various TB tests (24). Several tests are now endorsed by the World Health Organization (WHO). At the same time, there are examples of suboptimal and inaccurate TB diagnostic tests, and their use has been discouraged by the WHO (32). In fact, *in vitro* tests are often poorly regulated in many countries (13).

We searched PubMed, the Cochrane Library, and the Evidence-Based TB Diagnosis website (www.tbvidence.org) for systematic reviews on the accuracy of diagnostics for TB published through March 2012 (search terms are available from the authors upon request). We excluded meta-analyses that reported only performance characteristics other than sensitivity and specificity (e.g., reproducibility, likelihood ratios), as they were not comparable to information provided in package inserts.

We searched company websites for package inserts and contacted test manufacturers if package inserts were not available online. Diagnostic tests were not considered if they were not commercially available. Diagnostic tests for latent TB were not considered, because no good reference standard is available to determine accuracy.

A total of 19 TB tests were included in the final analysis, because they met our eligibility criteria and package inserts as well as meta-analyses where available. These included gamma interferon (IFN- γ) release assays (IGRAs) for active TB, antibody-based serological tests, antigen detection tests, nucleic acid amplification tests (NAAT), and culture-based tests.

As seen in Table 1, the quality of information provided in package inserts varies widely. Eighteen out of 19 package inserts in total

overestimated the test accuracy. In particular, technologies that were neither recommended by the WHO nor approved by the U.S. Food and Drug Administration (FDA) reported higher accuracy (i.e., serological tests for active TB, IGRAs for active TB, bacteriophage-based tests for active TB and drug-susceptibility testing, and urine lipoarabinomannan [LAM] antigen assays) (30, 32, 34). Claims made in package inserts ranged on average from 20 to 30% higher for sensitivity estimates than meta-analysis (comparing a range of findings in meta-analyses to estimates in package inserts). A direct comparison of absolute estimates of test accuracy was not feasible for most nonapproved tests, because the systematic reviews were unable to compute pooled sensitivity and specificity due to heterogeneity across studies.

In contrast, there is a better match between sensitivity and specificity estimates in package inserts and meta-analyses for tests that are approved by the FDA (i.e., Gen-Probe amplified MTD) or endorsed by the WHO (line probe assays, Xpert MTB/RIF, and MODS) (4, 31, 33, 34). For the WHO-endorsed tests that allowed a comparison of test accuracy between meta-analyses and package inserts (i.e., all except for the line probe assays), the package inserts overestimated sensitivity and specificity by at most 5%. This was also true for FDA-approved tests. IGRAs are FDA approved for latent TB but not active TB. The sensitivity for IGRAs for active TB was overestimated in package inserts by up to 20%. The comparison of the specificity of IGRAs was limited by the fact that meta-analyses of active TB often used TB suspects as controls, while package inserts reported specificity for latent TB infection among healthy, low-risk populations (6, 9, 19, 26).

We also found that test accuracy estimates in package inserts are often derived from unpublished, in-house, case-control studies with small numbers of specimens. Confirmed TB cases and healthy controls are often used, which can introduce significant

Received 28 March 2012 Returned for modification 25 April 2012

Accepted 30 April 2012

Published ahead of print 9 May 2012

Address correspondence to Madhukar Pai, madhukar.pai@mcgill.ca.

Copyright © 2012, American Society for Microbiology. All Rights Reserved.

doi:10.1128/JCM.00842-12

TABLE 1 Comparison of test accuracy in package inserts versus meta-analyses^a

Indication	Specimen	Test, yr of package insert publication (if reported)	WHO endorsed	FDA approved	Package insert		Meta-analysis		Comment on meta-analysis
					Sensitivity (%)	Specificity (%)	Sensitivity (%)	Specificity (%)	
IPN- γ release assays	Blood	TB-SPOT.TB assay (Oxford Immunotec, Abingdon, United Kingdom), 2010	No	No	96	97	92	59	Sensitivity reported for culture-confirmed cases; specificity reported for TB suspects, which may in part explain the reduced estimate compared to the PI
			No	No	89	99	81	79	Sensitivity reported for culture-confirmed cases; specificity reported for TB suspects, which may in part explain the reduced estimate compared to the PI
Active PTB	Blood	QuantIFERON-TB Gold In-Tube (QFT-GIT, Cellistis, Chadstone, Australia), 2006	No	No	96	97	83 (68 HIV ⁺ , 88 HIV ⁻)	61 (52 HIV ⁺)	Stratified by HIV; sensitivity reported for culture-confirmed cases; specificity reported for TB suspects
			No	No	89	99	69 (65 HIV ⁺ , 84 HIV ⁻)	52 (50 HIV ⁺)	Stratified by HIV; sensitivity reported for culture-confirmed cases; specificity reported for TB suspects
Active PTB and EPTB	Blood	TB-SPOT	No	No	96	97	88–91	86–88	Specificity reported in healthy controls; specificity lower if assessed in TB suspects
			No	No	89	99	79–81	93–99	Specificity reported in healthy controls; specificity lower if assessed in TB suspects
Antigen-based tests	Urine	LAM-ELISA (Chemogen, Clearview, Now Alere TB-LAM ELISA, Waltham, MA), 2011	No	No	73–81 (for HIV ⁺ only)	70–88	47–51	94–96	Evaluated precommercial and commercial tests; concerns about methodology in majority of studies; specificity mostly tested on TB suspects but healthy controls included too; sensitivity approaches no. in package inserts in patients with advanced HIV
			No	No	Estimates not reported for HIV ⁻	Estimates not reported for HIV ⁻	14	97	Evaluated precommercial and commercial tests; concerns about methodology in majority of studies; specificity mostly tested on TB suspects but healthy controls included too; sensitivity approaches no. in package inserts in patients with advanced HIV

Assay	Sample	Manufacturer	WHO recommendation	Study range	Reported	PI details	PI range	Number of studies	Study range	Comments	
Serological antibody detection assays	Active S ⁺ PTB	Blood	anda-TB IgG (anda Biologicals, Strasbourg, France)	No	48-100	71-100	Not reported	76	92	11/1,570 (28)	All studies in the literature with serious methodological problems and concerns about study population not being representative
				WHO recommended against	48-100	71-100	Not reported	59	91	11/1,570 (28)	All studies in the literature with serious methodological problems and concerns about study population not being representative
				Commercially available	48-100	71-100	Not reported	81	85	11/1,570 (28)	All studies in the literature with serious methodological problems and concerns about study population not being representative
Active PTB	Blood	Commercially available	WHO recommended against	No	70-100 (all tests but anda)	90-100 (all tests but anda)	No, often not reported or <50	60-88	50-98	2 MAs: 54/3,696; 8/mean SS per study 250, total no. not available (10, 28)	All studies with methodological problems as described above
				WHO recommended against	70-100 (all tests but anda)	90-100 (all tests but anda)	No, often not reported or <50	60-88	50-98	2 MAs: 54/3,696; 8/mean SS per study 250, total no. not available (10, 28)	All studies with methodological problems as described above

(Continued on following page)

TABLE 1 (Continued)

Indication	Specimen	Test, yr of package insert publication (if reported)	WHO endorsed	FDA approved	Package insert		No. of samples included	Meta-analysis		No. of studies/no. of participants included (reference)	Comment on meta-analysis
					WHO recommended against	Sensitivity (%)		Sensitivity (%)	Specificity (%)		
Active PTB	Blood	Commercially available (ELISA +ICT (MycDot, Mossman Blackstone, MA; <i>Mycobacterium tuberculosis</i> IgG, IBL, Hamburg, Germany), 2011; ActiveTBDetect (InBios International, Seattle, WA), 2008; SEVA (Mahatma Gandhi Institute of Medical Sciences, India; Pathozyme, Omega Diagnostics, Alva, Scotland), 2009; Hexagon (Human Gesellschaft Biochemica und Diagnostica, Wiesbaden, Germany), 2011; Serocheck-MTB (Zephyr, Biomedicals, Goa, India)	No	No	48–100	71–100	No, often not reported or <50	43	93	4/604 (10)	All studies of moderate quality
Bacteriophage-based tests Active PTB	Sput direct	FASTPlaque-TB (Biotec, Kentford, United Kingdom), 2004	No	No	73–82 (S ⁻ , 49–67; S ⁺ , 87)	98–99 (S ⁻ , 98–100; S ⁺ , 83–88)	>2,000	21–94 (S ⁻ , 13–78; S ⁺ , 75–87)	83–100 (S ⁻ , 89–99; S ⁺ , 60–88)	13/5,820 (14)	Data in meta-analysis not pooled due to heterogeneity
Active PTB; RIF	Sput S ⁺	FASTPlaque RIF, FASTPlaque Response (Biotec, Kentford, United Kingdom), 2005	No	No	96–100 (only S ⁺ cases)	98–100	374	96	95	31/3,085 (22)	Combines older and newer tests; 3–16% uninterpretable results
Nucleic acid-based test MTB detection Active PTB	Sput direct	Xpert MTB/Rif (Cepheid, Sunnyvale, CA), 2011	Yes	No	92 (S ⁺ , 98; S ⁻ , 73)	99	1,335	90 (S ⁺ , 99; S ⁻ , 75)	98 (S ⁺ , 98; S ⁻ , 98)	18/10,224 (5)	Study performed simple pooling of sensitivity and specificity
Active PTB	Sput direct	Amplified MTD (Gen-Probe, San Diego, CA), 2001	No	Yes	86 (S ⁺ , 97; S ⁻ , 72)	99 (S ⁺ , 100; S ⁻ , 99)	206	88 (S ⁺ , 97–100; S ⁻ , 70–76)	96 (S ⁺ , 96–98; S ⁻ , 95–97)	3 MAAs: 25/mean SS per study 362; 14/median SS 410; 40/mean SS 715 (10, 12, 16)	Results more consistent for specificity; BD Probe Tec data combined older and newer versions
	Sput direct	Probe Tec ET (BD, Franklin Lakes, NJ), 2010	No	Yes	91 (S ⁺ , 99; S ⁻ , 75)	97	986	86–88 (S ⁺ , 98; S ⁻ , 71)	98–99 (S ⁺ , 89; S ⁻ , 97)	3 MAAs: 3/213 mean SS per study; 12/median SS 410; 9/mean SS 715 (10, 12, 16)	Results more consistent for specificity; BD Probe Tec data combined older and newer versions

Nucleic-acid based test resistance detection	Active PTB; R	Active PTB; R	RIF	Sput direct	Xpert MTB/Rif	Yes	No	97	98	567	94	97	18/10,224 (5)	Study performed simple pooling of sensitivity and specificity
Active PTB; R	Yes	No	No	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	98-99	99	2 MAs: 5/767; 4/931 (2, 17)	One meta-analysis analyzed studies that tested only directly on smear-positive sputum; the other evaluated both direct and indirect testing
Active PTB; R	Yes	No	No	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	PI without any information about test accuracy	89-96	99-100	2 MAs: 5/767; 5/981 (2, 17)	One meta-analysis analyzed studies that tested only directly on smear-positive sputum; the other evaluated both direct and indirect testing
Active PTB; R	Yes	No	No	99	100	289	Retrospective, unpublished; only done indirectly on culture	92-100	15/1,738 (23)	Discrepancy in sensitivity estimates compared to PI likely in parts due to the inclusion of studies that tested directly on sputum in meta-analysis				
Microscopic observation drug susceptibility testing	Active PTB	Yes	No	98 for MTB detection (R 100% for RIF, 97% for INH)	100 for MTB detection	No data	PI refers to one major published study without further elaborating on study findings	92	14/3,731 for sensitivity; 12/7,226 for specificity (15)			96		
Active PTB; R	Yes	No	No	98 for MTB detection (R 100% for RIF, 97% for INH)	100 for MTB detection	No data	PI refers to one major published study without further elaborating on study findings	96-98	2 MAs: 6/1,187 and 9/1,474 (2, 20)	Sensitivity, specificity for INH varies slightly at different cutoff (0.1 or 0.4 µg/ml)				
Active PTB; R	Yes	No	No	98 for MTB detection (R 100% for RIF, 97% for INH)	100 for MTB detection	No data	PI refers to one major published study without further elaborating on study findings	92	2 MAs: 6/1,187 and 9/1,474 (2, 20)	Sensitivity, specificity for INH varies slightly at different cutoff (0.1 or 0.4 µg/ml)				

^a PI, package insert; R, resistance; RIF, rifampin; INH, isoniazid; PTB, pulmonary TB; EPTB, extrapulmonary TB; S⁺, smear positive; S⁻, smear negative; SS, sample size; MAs, meta-analyses; MTB, *Mycobacterium tuberculosis*.

selection (spectrum) bias (25). The data in the package inserts often are not stratified based on important predictors of performance, including prevalence of TB or HIV and adults versus children, which may contribute to the overestimation of accuracy (3, 8, 19). In contrast, meta-analyses were often based on a fairly large number of studies that used cross-sectional or prospective designs and often were conducted in clinical settings with TB suspects that had a confirmed alternative final diagnosis serving as controls. Results were often stratified based on clinically relevant subgroups.

In general, involvement of industry and test developers in diagnostic evaluations has been associated with an overestimation of test accuracy (1). With TB tests, this has been documented with bacteriophage-based tests and urine lipoarabinomannan assays (11, 14, 21, 22). Users in real-world clinical settings may lack the same degree of expertise and skill as test developers. Also, quality control and assurance in routine clinical and laboratory settings may not match that of the industry. While data included in package inserts are almost always funded by industry, a proportion of studies included in the meta-analyses also are industry supported or conducted by test developers. This may then spuriously narrow the gap between package insert and meta-analyses estimates.

Our study has limitations. We were unable to compute numeric differences in the estimates of meta-analyses versus package inserts because pooling of data was often not possible due to heterogeneity between studies and the presence of several meta-analyses that included partially overlapping studies. We acknowledge that real-world performance of tests, especially when tests are scaled up in public health programs, may be worse than those reported in research studies, including meta-analyses (27). Thus, the real gap between package insert estimates and real-world performance may be even wider than what we document here. Pragmatic trials and implementation research are needed to overcome this problem (18). We also acknowledge that tests that measure the immune response to TB (i.e., serology, IGRAs) rather than products of *Mycobacterium tuberculosis* (i.e., Xpert MTB/RIF) might be more prone to variability in the results; however, this underlines the fact that accuracy data should always be stratified based on clinically relevant subgroups (i.e., HIV positive).

In summary, this case study of TB diagnostics suggests that package inserts often report overoptimistic estimates of test accuracy, especially if the products are not FDA approved (provided that approval was solicited) or WHO endorsed. These data provide some reassurance that independent review by credible agencies such as the FDA and WHO may serve as a yardstick for judging new TB technologies. However, not all TB tests are reviewed by the FDA or WHO, and most developing countries have weak regulatory systems for diagnostics. It is important that these countries create systems for in-country validation of all TB tests, guided by their national TB programs. Also, an expansion of the WHO prequalification of diagnostic programs to TB diagnostics will help countries procure quality-assured TB tests.

To overcome the problem of optimism bias, studies evaluating diagnostics under routine clinical and programmatic conditions, independent of industry sponsorship or test developers, are needed, as they provide more useful and realistic evidence to guide laboratorians, clinicians, and decision-makers. Furthermore, studies must go beyond accuracy and assess clinical impact of tests on decision-making and patient outcomes and collect operational and cost-effectiveness data in programmatic settings (7, 18).

ACKNOWLEDGMENTS

This study was supported by grants from the Canadian Institutes of Health Research (CIHR MOP-88918) and European and Developing Countries Clinical Trials Partnership (EDCTP) (TB-NEAT). Madhukar Pai is supported by a CIHR New Investigator Award and a career award from the Fonds de Recherche du Québec—Santé. None of these agencies were involved in the conduct or review of this study or the decision to publish its results.

We thank Daphne Ling for her contributions to this project.

The authors have no other relevant affiliations or financial involvement with any organization or entity with a financial interest in or financial conflict with the subject matter or materials discussed.

REFERENCES

1. Bekelman JE, Li Y, Gross CP. 2003. Scope and impact of financial conflicts of interest in biomedical research: a systematic review. *JAMA* 289:454–465.
2. Bwanga F, Hoffner S, Haile M, Joloba ML. 2009. Direct susceptibility testing for multi drug resistant tuberculosis: a meta-analysis. *BMC Infect. Dis.* 9:67.
3. Cattamanchi A, et al. 2011. Interferon-gamma release assays for the diagnosis of latent tuberculosis infection in HIV-infected individuals: a systematic review and meta-analysis. *J. Acquir. Immune Defic. Syndr.* 56:230–238.
4. Centers for Disease Control and Prevention (CDC). 2009. Updated guidelines for the use of nucleic acid amplification tests in the diagnosis of tuberculosis. *MMWR Morb. Mortal. Wkly. Rep.* 58:7–10.
5. Chang K, et al. 27 February 2012. Rapid and effective diagnosis of tuberculosis and rifampicin resistance with Xpert MTB/RIF assay: a meta-analysis. *J. Infect.* 64:580–588.
6. Chang KC, Leung CC. 2010. Systematic review of interferon-gamma release assays in tuberculosis: focus on likelihood ratios. *Thorax* 65:271–276.
7. Cobelens F, et al. 2012. On behalf of the Evidence for Scale-up Group. Which new diagnostics for tuberculosis, and when? *J. Infect. Dis.* 205(Suppl 2): S191–S198.
8. Dheda K, van Zyl Smit R, Badri M, Pai M. 2009. T-cell interferon-gamma release assays for the rapid immunodiagnosis of tuberculosis: clinical utility in high-burden vs. low-burden settings. *Curr. Opin. Pulm. Med.* 15:188–200.
9. Diel R, Loddenkemper R, Nienhaus A. 2010. Evidence-based comparison of commercial interferon-gamma release assays for detecting active TB: a meta-analysis. *Chest* 137:952–968.
10. Dinnes J, et al. 2007. A systematic review of rapid diagnostic tests for the detection of tuberculosis infection. *Health Technol. Assess.* 11:1–196.
11. Flores LL, et al. 2011. Systematic review and meta-analysis of antigen detection tests for the diagnosis of tuberculosis. *Clin. Vaccine Immunol.* 18:1616–1627.
12. Greco S, Girardi E, Navarra A, Saltini C. 2006. Current evidence on diagnostic accuracy of commercially based nucleic acid amplification tests for the diagnosis of pulmonary tuberculosis. *Thorax* 61:783–790.
13. Grenier J, et al. 2012. Widespread use of serological tests for tuberculosis: data from 22 high-burden countries. *Eur. Respir. J.* 39:502–505.
14. Kalantri S, Pai M, Pascopella L, Riley L, Reingold A. 2005. Bacteriophage-based tests for the detection of *Mycobacterium tuberculosis* in clinical specimens: a systematic review and meta-analysis. *BMC Infect. Dis.* 5:59.
15. Leung E, Minion J, Benedetti A, Pai M, Menzies D. 2011. Microcolony culture techniques for tuberculosis diagnosis: a systematic review. *Int. J. Tuberc. Lung Dis.* 16:16–23.
16. Ling DI, Flores LL, Riley LW, Pai M. 2008. Commercial nucleic-acid amplification tests for diagnosis of pulmonary tuberculosis in respiratory specimens: meta-analysis and meta-regression. *PLoS One* 3:e1536. doi: 10.1371/journal.pone.0001536.
17. Ling DI, Zwerling AA, Pai M. 2008. GenoType MTBDR assays for the diagnosis of multidrug-resistant tuberculosis: a meta-analysis. *Eur. Respir. J.* 32:1165–1174.
18. Mann G, et al. 2010. Beyond accuracy: creating a comprehensive evidence base for TB diagnostic tools. *Int. J. Tuberc. Lung Dis.* 14:1518–1524.
19. Metcalfe JZ, et al. 2011. Interferon-gamma release assays for active pulmonary tuberculosis diagnosis in adults in low- and middle-income

- countries: systematic review and meta-analysis. *J. Infect. Dis.* 204(Suppl 4):S1120–S1129.
20. Minion J, Leung E, Menzies D, Pai M. 2010. Microscopic-observation drug susceptibility and thin layer agar assays for the detection of drug resistant tuberculosis: a systematic review and meta-analysis. *Lancet Infect. Dis.* 10:688–698.
 21. Minion J, et al. 2011. Diagnosing tuberculosis with urine lipoarabinomannan: systematic review and meta-analysis. *Eur. Respir. J.* 38:1398–1405.
 22. Minion J, Pai M. 2010. Bacteriophage assays for rifampicin resistance detection in *Mycobacterium tuberculosis*: updated meta-analysis. *Int. J. Tuberc. Lung Dis.* 14:941–951.
 23. Morgan M, Kalantri S, Flores L, Pai M. 2005. A commercial line probe assay for the rapid detection of rifampicin resistance in *Mycobacterium tuberculosis*: a systematic review and meta-analysis. *BMC Infect. Dis.* 5:62.
 24. Pai M, Minion J, Steingart K, Ramsay A. 2010. New and improved tuberculosis diagnostics: evidence, policy, practice, and impact. *Curr. Opin. Pulm. Med.* 16:271–284.
 25. Rutjes AW, Reitsma JB, Vandenbroucke JP, Glas AS, Bossuyt PM. 2005. Case-control and two-gate designs in diagnostic accuracy studies. *Clin. Chem.* 51:1335–1341.
 26. Sester M, et al. 2011. Interferon-gamma release assays for the diagnosis of active tuberculosis: a systematic review and meta-analysis. *Eur. Respir. J.* 37:100–111.
 27. Stall N, et al. 2011. Does solid culture for tuberculosis influence clinical decision making in India? *Int. J. Tuberc. Lung Dis.* 15:641–646.
 28. Steingart KR, et al. 2011. Commercial serological tests for the diagnosis of active pulmonary and extrapulmonary tuberculosis: an updated systematic review and meta-analysis. *PLoS Med.* 8:e1001062. doi:10.1371/journal.pmed.1001062.
 29. Reference deleted.
 30. World Health Organization. 2009. Accessible quality-assured diagnostics—2009 annual report. WHO, Geneva, Switzerland.
 31. World Health Organization. 2011. Automated real-time nucleic acid amplification technology for rapid and simultaneous detection of tuberculosis and rifampicin resistance: Xpert MTB/RIF system. WHO, Geneva, Switzerland.
 32. World Health Organization. 2011. Commercial serodiagnostic tests for diagnosis of tuberculosis. WHO, Geneva, Switzerland.
 33. World Health Organization. 2008. Molecular line probe assays for rapid screening of patients at risk for multidrug-resistant tuberculosis (MDR-TB). WHO, Geneva, Switzerland.
 34. World Health Organization. 2011. Noncommercial culture and drug-susceptibility testing methods for screening patients at risk for multidrug-resistant tuberculosis. WHO, Geneva, Switzerland.