

Optimism bias, inflated accuracy estimates, and contradicted findings in TB diagnostic research



Madhukar Pai, MD, PhD [madhukar.pai@mcgill.ca]
Jessica Minion, MD

McGill University, Montreal



McGill



L'Institut de recherche du Centre universitaire de santé McGill
The Research Institute of the McGill University Health Centre
Les meilleurs soins pour la vie
The Best Care for Life

1

Context

There is some evidence that:

- Initially stronger effects and subsequent contradictions are not infrequent in highly cited research of clinical interventions and their outcomes.
- Claims from highly cited observational studies persist and continue to be supported in the medical literature despite strong contradictory evidence from randomized trials.
- Newly discovered true (non-null) associations often have inflated effects compared with the true effect sizes.
- Publication bias is a major concern and may be more widespread than we think; some have also challenged the conventional publishing model
- Even within published studies, selective reporting of positive outcomes in randomized trials as well as observational studies appears to be frequent
- Lack of replication of research findings and over-interpretation of findings are other concerns, especially in some fields

2

Context

- All of these likely result in "optimism bias" —unwarranted belief in the efficacy of new therapies, and overinterpretation of the applicability of findings
- Optimism bias is more likely in industry-supported research
- Optimism bias and conflicting study findings appear to be eroding the public's faith in research
- Even among some researchers, there is concern that most published research findings may be false!

3

Contradicted and Initially Stronger Effects in Highly Cited Clinical Research

John P. A. Ioannidis, MD

CLINICAL RESEARCH ON IMPORTANT questions about the efficacy of medical interventions is sometimes followed by subsequent studies that either reach opposite conclusions or suggest that the original claims were too strong. Such disagreements may upset clinical practice and acquire publicity in both scientific circles and in the lay press. Several empirical investigations have tried to address whether specific types of studies are more likely to be contradicted and to explain observed controversies. For example, evidence exists that small studies may sometimes be refuted by larger ones.^{1,2}

Similarly, there is some evidence on disagreements between epidemiological studies and randomized trials.^{3,4} Prior investigations have focused on a variety of studies without any particular attention to their relative importance and scientific impact. Yet, most research publications have little impact while a small minority receives

Context: Controversy and uncertainty arise when the results of clinical research on the effectiveness of interventions are subsequently contradicted. Controversies are most prominent when high-impact research is involved.

Objectives: To understand how frequently highly cited studies are contradicted or find effects that are stronger than in other similar studies and to discern whether specific characteristics are associated with such refutation over time.

Design: All original clinical research studies published in 3 major general clinical journals or high-impact factor specialty journals in 1990–2003 and cited more than 1000 times in the literature were examined.

Main Outcome Measure: The results of highly cited articles were compared against subsequent studies of comparable or larger sample size and similar or better controlled designs. The same analysis was also performed comparatively for matched studies that were not so highly cited.

Results: Of 49 highly cited original clinical research studies, 45 claimed that the intervention was effective. Of these, 7 (16%) were contradicted by subsequent studies. Authors (16%) had found effects that were stronger than those of subsequent studies, 20 (44%) were replicated, and 11 (24%) remained largely unchallenged. Five of 6 highly cited nonrandomized studies had been contradicted or had found stronger effects vs 9 of 39 randomized controlled trials ($P = .008$). Among randomized trials, studies with contradicted or stronger effects were smaller ($P < .001$) than replicated or unchallenged studies, although there was no statistically significant difference in their early or overall citation impact. Matched control studies did not have a significantly different share of refuted results than highly cited studies, but they included more studies with "negative" results.

Conclusions: Contradicted and initially stronger effects are not unusual in highly cited research of clinical interventions and their outcomes. The extent to which high citations may provoke contradictions and vice versa needs more study. Controversies are most common with highly cited nonrandomized studies, but even the most highly cited randomized trials may be challenged and refuted over time, especially small ones. *JAMA*. 2008;294:218–228.

www.jama.com

Why Most Discovered True Associations Are Inflated

John P. A. Ioannidis

Abstract: Newly discovered true (non-null) associations often have inflated effects compared with the true effect sizes. I discuss here the main reasons for this inflation. First, theoretical considerations prove that when true discovery is claimed based on crossing a threshold of statistical significance and the discovery study is underpowered, the observed effects are expected to be inflated. This has been demonstrated in various fields ranging from early stopped clinical trials to genome-wide associations. Second, flexible analyses coupled with selective reporting may inflate the published discovered effects. The inflation into the rate of the largest to smallest effect in the same association approached with different analytic choices can be very large. Third, effects may be inflated at the stage of interpretation due to diverse conflicts of interest. Discovered effects are not always inflated, and under some circumstances may be deflated—for example, in the setting of late discovery of associations in sequentially accumulated component evidence, in some types of misclassification from measurement error, and in conflicts among review teams. Finally, I discuss potential approaches to this problem. These include being cautious about newly discovered effect sizes, considering more rational down-weighting, using meta-analytic methods that correct for the anticipated inflation, spreading the magnitude of the effect of an association, conducting large studies in the discovery phase, using strict protocols for analysis, pursuing complex and transparent reporting of all results, placing emphasis on replication, and being wary with interpretation of results. *JAMA*. 2008;294:1640–1649.

Persistence of Contradicted Claims in the Literature

Mikina Tsalikis, MD

Nicholas C. Bressler, MD

John P. A. Ioannidis, MD

Context: Some research findings based on observational epidemiology are contradicted by randomized trials, but may nevertheless still be supported in some scientific circles.

Objectives: To evaluate the change over time in the context of citations for 2 highly cited epidemiological studies that proposed major cardiovascular benefits associated with vitamin E in 1993, and to understand how these benefits continued being defended in the literature, despite strong contradictory evidence from large randomized clinical trials (RCTs). To examine the generalizability of these findings, we also examined the extent of persistence of supporting citations for the highly cited and contradicted protective effects of beta-carotene on cancer and of estrogen on Alzheimer disease.

Data Sources: For vitamin E, we sampled articles published in 1997, 2001, and 2006. Before, early, and late after publication of refuting evidence that refuted the highly cited epidemiological studies and separately sampled articles published in 2005 and refuting the major contradictory RCT in 2005. We also sampled articles published in 2006 that refuted highly cited articles proposing benefits associated with beta-carotene for cancer (published in 1981) and contradicted long ago by RCTs in 1994–1996) and estrogen for Alzheimer disease (published in 1996) and contradicted recently by RCTs in 2004.

Data Extraction: The stance of the citing articles was rated as favorable, equivocal, and unfavorable to the intervention. We also measured the range of counterarguments raised to defend effectiveness against contradictory evidence.

Results: For the 2 vitamin E epidemiological studies, more than 2000, 50% of citing articles remained favorable. A favorable stance was independently less likely in more recent articles, specifically in articles that also cited the RCT trial (odds ratio for 2001, 0.06 [95% confidence interval, 0.02–0.19; $P < .001$) and the odds ratio for 2005, 0.04 [95% confidence interval, 0.02–0.24; $P < .001$], as compared with 1997, and in general refuted evidence in specialty journals. Among articles citing the RCT trial in 2005, 41.4% were unfavorable in 2006, 62.5% of articles refuting the highly cited article that had proposed beta-carotene and 81.7% of those refuting the highly cited article on estrogen effectiveness were still favorable. 100% and 96%, respectively, of the citations appeared to specify potential, and citations were significantly less favorable ($P < .001$ and $P < .009$, respectively) when the major contradicting trials were also mentioned. Counterarguments defending vitamin E or estrogen included diverse selection and reference biases and genuine differences across studies in participants, interventions, confounders, and outcomes. Favorable citations to beta-carotene, long after evidence contradicted its effectiveness, did not consider the contradictory evidence.

Conclusion: Claims from highly cited observational studies persist and continue to be supported in the medical literature despite strong contradictory evidence from randomized trials. *JAMA*. 2007;297:2017–2026.

www.jama.com

Such debate offers opportunities to

prospective studies, and so forth. I start here with the assumption that a research finding is indeed true (non-null), so, it reflects a genuine association that is not entirely due to chance or biases (confounding, misclassification, selection biases, selective reporting, or other). The question is, do the effect sizes for such associations, at the time they are first discovered and published in the scientific literature, accurately reflect the true effect sizes?

The article has the following sections: a brief literature review on inflated early-effect sizes based on theoretical and empirical considerations; a description of the major reasons why early discovered effects are inflated and the major countervailing forces that may occasionally lead to deflated effects (underestimates); and suggestions on how to deal with these problems.

Evidence About Inflated Early-Effect Sizes

Table 1 cites articles suggesting that early studies give (on average) inflated estimates of effect.^{1–14} I list here only selected evaluations that cover either many different articles effects or a whole research domain or method. This list is nowhere close to exhaustive. For some topics, such as the inflation of regression coefficients for variables selected through stepwise statistical significance-based procedures, the literature is

4

Even within published studies, selective reporting of outcomes

Empirical Evidence for Selective Reporting of Outcomes in Randomized Trials Comparison of Protocols to Published Articles

An-Wen Chen, MD, PhD
Adrian Haidich, MD, PhD
Meta T. Haidich, RN
Peter C. Gøtzsche, MD, DSc
Douglas G. Altman, DSc

Context Selective reporting of outcomes within published studies based on the bias or direction of their results has been widely suspected, but direct evidence of such bias is currently limited to case reports.

Objective To study empirically the extent and nature of outcome reporting bias in a cohort of randomized trials.

Design Cohort study using protocols and published reports of randomized trials approved by the Scientific Ethics Committee for Copenhagen and Frederiksberg, Denmark, in 1994-1995. The number and characteristics of reported and unreported trial outcomes were recorded from protocols, journal articles, and a survey of trials. An outcome was considered incompletely reported if sufficient data were presented in the published articles for meta-analysis. Odds ratios relating the completeness of outcome reporting to statistical significance were calculated for each trial and then pooled to provide an overall estimate of bias. Protocols and published articles were also compared to identify discrepancies in primary outcomes.

Main Outcome Measures Completeness of reporting of efficacy and harm outcomes and of statistically significant or nonsignificant outcomes; consistency between primary outcomes defined in the most recent protocols and those defined in published articles.

Results One hundred two trials with 122 published journal articles and 3736 outcomes were identified. Overall, 50% of efficacy and 65% of harm outcomes per trial were incompletely reported. Statistically significant outcomes had a higher odds of being fully reported compared with nonsignificant outcomes for both efficacy (pooled odds ratio, 2.4; 95% confidence interval [CI], 1.4-4.0) and harm (pooled odds ratio, 4.7; 95% CI, 1.8-12.0) data. In comparing published articles with protocols, 42% of trials had at least 1 primary outcome that was changed, introduced, or omitted. Eighty-one percent of survey responders (42/49) denied the existence of unreported outcomes despite clear evidence to the contrary.

Conclusions The reporting of trial outcomes is not only frequently incomplete but also biased and inconsistent with protocols. Published articles, as well as reviews that incorporate them, may therefore be unreliable and compromise the benefits of an intervention. To ensure transparency, planned trials should be registered and protocols should be made publicly available prior to trial completion.

METHODS

Selection in Reported Epidemiological Risks: An Empirical Assessment

Fatmi K. Karvonen¹, George Liberopoulos¹, John P. A. Ioannidis^{1,2*}

¹ Clinical and Molecular Epidemiology Unit, Department of Epidemiology and Biostatistics, University of Toronto School of Medicine, Toronto, Ontario, Canada; ² Department of Medicine, Yale University School of Medicine, Boston, Massachusetts, United States of America

Background The authors assessed the extent of selective reporting of outcomes in a large sample of recent articles.

Methods The authors searched for specific findings for the study. The authors had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests The authors have declared that no competing interests exist.

Authors' Contributions Fatmi K. Karvonen, MD, PhD, and George Liberopoulos, MD, PhD, conceived of the study, participated in its design and coordination, drafted the manuscript, participated in the sequence alignment, and read and approved the final manuscript. John P. A. Ioannidis, MD, PhD, participated in the design and coordination, drafted the manuscript, participated in the sequence alignment, and read and approved the final manuscript.

Published: August 12, 2008
Copyright: © 2008 Karvonen et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Keywords: MEDLINE, ranges of outcomes, CI, confidence interval, OR, interquartile range

* To whom all correspondence should be addressed. E-mail: jioannidis@yale.edu

ABSTRACT

Background

Epidemiological studies may be subject to selective reporting, but empirical evidence thereof is limited. We empirically evaluated the extent of selection of significant results and large effect sizes in a large sample of recent articles.

Methods and Findings

We evaluated 369 articles of epidemiological studies that reported, in their respective abstracts, at least one relative risk for a continuous risk factor in contrast to median, quartile, or quartile categorizations. We examined the proportion and correlates of reporting statistically significant and nonsignificant results in the abstract and whether the magnitude of the relative risk presented seemed to be consistently >1.00 (odds depending on the type of contrast used for the risk factor). In 342 articles (87.0%), >1 statistically significant relative risk was reported in the abstract, while only 169 articles (45.8%) reported <1 statistically nonsignificant relative risk in the abstract. Reporting of statistically significant results was more common with structured abstracts, and was less common in US-based studies and in cancer outcomes. Among 32 randomly selected articles in which the full text was examined, a median of nine interquartile range 3-10 statistically significant and six interquartile range 3-10 statistically nonsignificant relative risks were presented ($p < 0.001$). Paradoxically, the smallest presented relative risks were based on the contrasts of extreme quartiles, extreme tertiles, and above-versus-below median values, respectively ($p < 0.001$).

Conclusions

Published epidemiological investigations almost universally highlight significant associations between risk factors and outcomes. For continuous risk factors, investigators selectively present contrasts between more extreme groups, when relative risks are inherently lower.

In RCTs

In observational studies

Essay

Why Most Published Research Findings Are False

John P. A. Ioannidis

Summary

There is increasing concern that most current published research findings are false. The probability that a research claim is true may depend on study power and bias, the number of other studies on the same question, and, importantly, the ratio of true to no relationships among the relationships probed in each scientific field. In this framework, a research finding is less likely to be true when the studies conducted in a field are smaller; when effect sizes are smaller; when there is a greater number and lesser preselection of tested relationships; where there is greater flexibility in designs, definitions, outcomes, and analytical modes; when there is greater financial and other interest and prejudices; and when more teams are involved in a scientific field in chase of statistical significance. Simulations show that for most study designs and settings, it is more likely for a research claim to be false than true. Moreover, for many current scientific fields, claimed research findings may often be simply accurate measures of the prevailing bias. In this essay, I discuss the implications of these problems for the conduct and interpretation of research.

factors that influence this problem and some corollaries thereof.

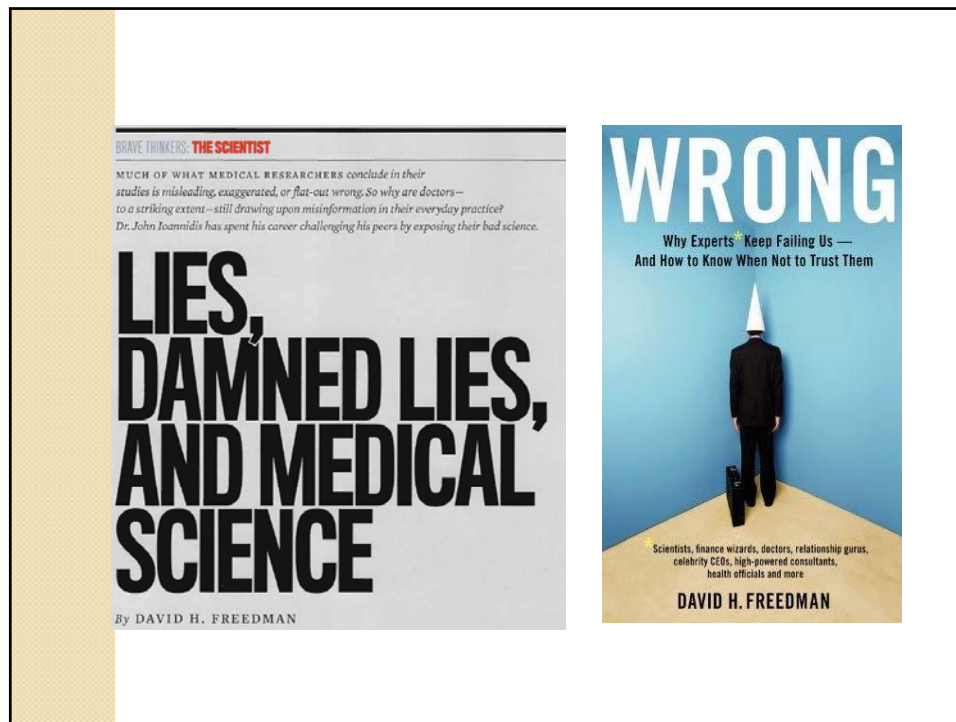
Modeling the Framework for False Positive Findings

Several methodologists have pointed out [9-11] that the high rate of nonreplication (lack of confirmation) of research discoveries is a consequence of the convenient, yet ill-founded strategy of claiming conclusive research findings solely on the basis of a single study assessed by formal statistical significance, typically for a p -value less than 0.05. Research is not most appropriately represented and summarized by p -values, but, unfortunately, there is a widespread notion that medical research articles

It can be proven that most claimed research findings are false.

should be interpreted based only on p -values. Research findings are defined here as any relationship reaching formal statistical significance, e.g., effective interventions, informative predictors, risk factors, or associations. "Negative" research is also more useful

is characteristic of the field and can vary a lot depending on whether the field targets highly likely relationships or searches for only one or a few true relationships among thousands and millions of hypotheses that may be postulated. Let us also consider, for computational simplicity, circumscribed fields where either there is only one true relationship (among many that can be hypothesized) or the power is similar to find any of the several existing true relationships. The pre-study probability of a relationship being true is $R/(R+1)$. The probability of a study finding a true relationship reflects the power $1 - \beta$ (one minus the Type II error rate). The probability of claiming a relationship when none truly exists reflects the Type I error rate, α . Assuming that c relationships are being probed in the field, the expected values of the 2×2 table are given in Table 1. After a research finding has been claimed based on achieving formal statistical significance, the post-study probability that it is true is the positive predictive value, PPV. The PPV is also the complementary probability of what Wacholder et al. have called the false positive report probability [10]. According to the 2



As researchers in tuberculosis we asked the question:

“is there evidence for ‘optimism bias’ in TB diagnostic research?”

We present several case studies to answer this question

Several new diagnostics are in the pipeline
But do they work? Will optimism bias prove to be a big issue?

Case study I:

How much evidence is sufficient for commercialization?

11

Promising new Point of Care test: LAM antigen detection



Journal of Microbiological Methods 45 (2001) 41–52

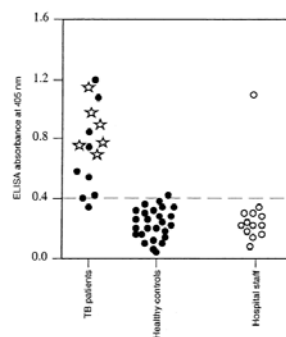
Journal
of Microbiological
Methods
www.elsevier.com/locate/jmicmeth

Rapid diagnosis of tuberculosis by detection of mycobacterial
lipoarabinomannan in urine

Beston Hamasur ^a, Judith Bruchfeld ^a, Melles Haile ^a, Andrzej Pawlowski ^a,
Bjarne Bjorvatn ^c, Gunilla Källentius ^{a,d}, Stefan B. Svenson ^{a,b,*}

Sensitivity 93%

Specificity 95%



12

Early data lead to rapid commercialization and marketing of a urine LAM assay

Transactions of the Royal Society of Tropical Medicine and Hygiene (2005) 99, 893–900



www.elsevierhealth.com/journals/trst

Detection of mycobacterial lipoarabinomannan with an antigen-capture ELISA in unprocessed urine of Tanzanian patients with suspected tuberculosis

C. Boehme^{a,*}, E. Molokova^b, F. Minja^c, S. Geis^a, T. Loscher^a, L. Maboko^d, V. Koulchin^b, M. Hoelscher^a

231 patients with suspected pulmonary TB and 103 healthy volunteers were screened with standard TB tests and with the new LAM-ELISA.

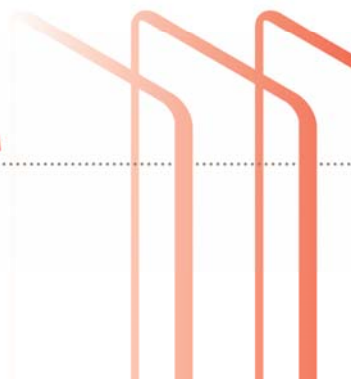
Of 132 patients with positive sputum culture, 106 were positive using the LAM-ELISA (sensitivity 80.3%)

To define the specificity of the assay, urine samples from 103 healthy volunteers were also screened using LAM-ELISA. All but one had an optical density below the cut-off (specificity 99%)

13

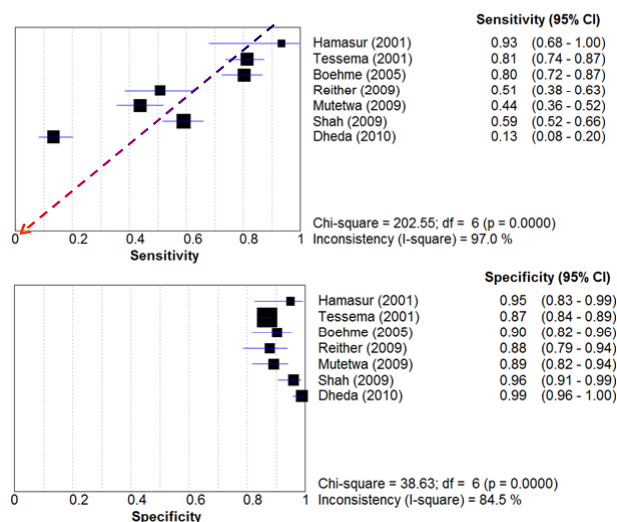
Marketed in 2007/08 by Inverness Medical Innovations

 **Clearview[®] TB ELISA**
LAM Specific Direct Urinary Antigen Test



14

Subsequent evidence from field studies in India, S Africa, Zimbabwe, Tanzania



Minion et al. Under review

15

Lessons

- Rapid commercialization on the basis of early data may be problematic (especially case-control studies that can exaggerate accuracy estimates)
- Thorough field evaluation in diverse settings (e.g. varying HIV prevalence) should have been done
- This case study raises an interesting question: at what point in time after a test is introduced should meta-analyses be done?

16

Case study 2:

How should we design and analyze diagnostic studies?

17

Serologic (antibody) tests for TB

- Attractive ... especially if point of care (POC) option
- >80 antigenic targets evaluated and several commercial assays developed
- All existing serologic tests have failed to demonstrate adequate accuracy

A systematic review of commercial serological antibody detection tests for the diagnosis of extrapulmonary tuberculosis

Karen R. Steingart,^{1*} Megan Henry,² Suman Laal,^{3,4} Philip C. Hopewell,² Andrew Ramsay,² Dick Menzies,⁵ Jane Cunningham,⁷ Karin Weldingh,¹⁰ Madhukar Pai^{6,8}

Thorax 2007

PLoS Medicine 2007

Commercial Serological Antibody Detection Tests for the Diagnosis of Pulmonary Tuberculosis: A Systematic Review

Karen R. Steingart^{1,2}, Megan Henry², Suman Laal^{3,4}, Philip C. Hopewell^{2,5}, Andrew Ramsay², Dick Menzies⁶, Jane Cunningham⁷, Karin Weldingh¹⁰, Madhukar Pai^{6,8}

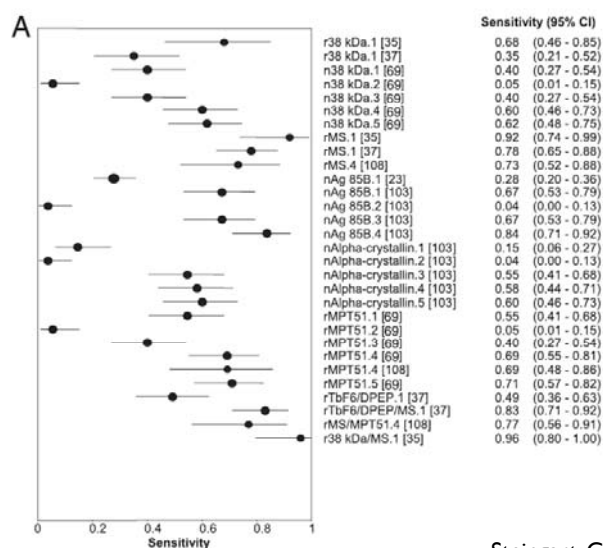
Performance of Purified Antigens for Serodiagnosis of Pulmonary Tuberculosis: a Meta-Analysis^{17,†}

Karen R. Steingart,^{1*} Nandini Dendukuri,² Megan Henry,^{3,†} Ian Schiller,² Payam Nahid,⁴ Philip C. Hopewell,^{1,4} Andrew Ramsay,² Madhukar Pai,² and Suman Laal^{6,7,8}

Clin Vaccine Immunol 2009

18

Sensitivity varies from 0 – 100%

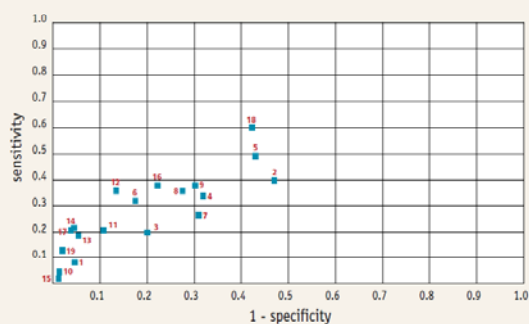


Steingart, CVI 2009

19

WHO evaluation of 19 commercial serological tests for TB

Figure 4. ROC curve of commercial rapid tests for the diagnosis of pulmonary tuberculosis (all patients, n=355)



1. ABP Diagnostics 2. Advanced Diagnostics 3. Products 6. Chembio Diagnostic Systems 7. CTK Biotech
American Bionostica 4. Ameritek USA 5. Bio-Medical 8. Hema Diagnostic Systems 9. Laboratorios Silanes



TDR/WHO Report 2008

20

Why do these tests fail in field studies?

TABLE 3. Characteristics of study quality

Characteristic	No. (%) of studies
Study design	
Cross-sectional	39 (15)
Case-control.....	208 (82)
Nested within observational study.....	7 (3)
Recruitment of participants	
Consecutive or random.....	20 (8)
Convenience or not reported.....	234 (92)
Selection criteria clearly described.....	141 (56)
Complete verification by use of the reference standard	107 (42)
Execution of test described in sufficient detail.....	253 (100)^a
Index test results blinded to reference standard?	
Yes.....	65 (26)
No	1 (0)
Not reported.....	188 (74)

^a The description of the test execution was deemed insufficient in one study.

A large % were case-control studies

Confirmed TB cases
Vs.
Healthy controls (often from low-incidence countries)

21

Spectrum bias (a form of selection bias)

- Population used for evaluating the test:
 - Extreme contrast
 - Case-control design
 - Normal contrast (Indicated population)
 - Consecutively recruited patients in whom the disease is suspected
- Extreme contrast (spectrum bias) can result in overestimation of test accuracy

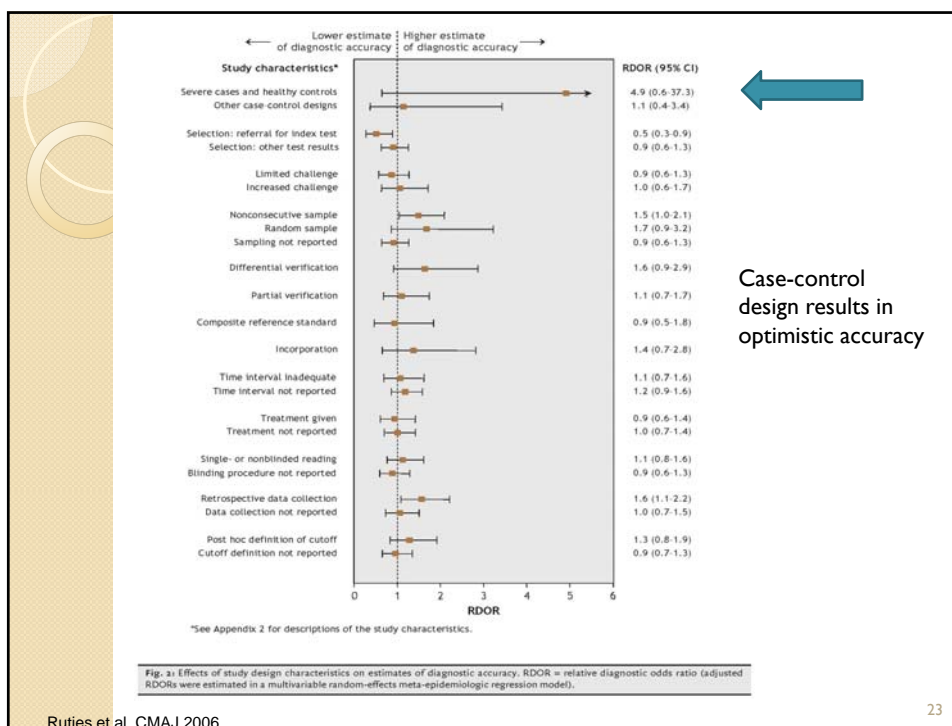
Clinical Chemistry 51:8
1335-1341 (2005)

Minireview

Case-Control and Two-Gate Designs in Diagnostic Accuracy Studies

ANNE W.S. RUTJES,^{1*} JOHANNES B. REITSMA,¹ JAN P. VANDENBROUCKE,² AFINA S. GLAS,³ and PATRICK M.M. BOSSUYT¹

22



Case-control design results in optimistic accuracy

We find this in TB as well: Example: PCR tests for TB meningitis

Diagnostic accuracy of nucleic acid amplification tests for tuberculous meningitis: a systematic review and meta-analysis

Madhukar Pai, Laura L Flores, Nikku Pai, Alan Hubbard, Lee W Riley, and John M Colford Jr

Case-control studies had a two-fold higher diagnostic odds ratios than cross-sectional studies

Table 4. Stratified analyses for the evaluation of heterogeneity among studies with in-house tests

Subgroup	Number of studies	Summary diagnostic odds ratio* (95% CI)	Test for heterogeneity† p value
Study design			
Case-control	19	86.5 (39.3, 190.2)	0.03
Cross-sectional	16	43.3 (22.5, 83.3)	0.94
Blinded interpretation of test and/or reference standard results			
Yes	21	46.9 (24.9, 88.6)	0.16
No	14	82.3 (39.8, 170.2)	0.70
Consecutive or random sampling of participants			
Yes	18	63.3 (32.8, 122.4)	0.20
No	17	46.8 (23.6, 92.8)	0.42
Prospective data collection			
Yes	18	59.9 (28.1, 127.6)	0.12
No	17	55.2 (29.9, 101.6)	0.59

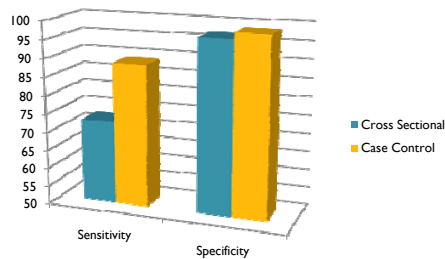
*Random effects model. †I² test for heterogeneity. CI=confidence interval.

LED microscopy for sputum examination

- Cross Sectional Studies
- Case Control Studies

Sensitivity 72.6%
(69.2, 75.8)
Specificity 96.9
(92.1, 98.8)

Sensitivity 88.7%
(81.4, 93.4)
Specificity: 98.6%
(97.3, 99.3)



Minion et al. Unpublished

25

Analysis of diagnostic studies

- It is not uncommon to see researchers:
 - Excluding patients or controls with no definitive diagnoses ("diagnostic myopia bias")
 - Excluding indeterminate or inconclusive results
 - Perform post-hoc "discrepant" analysis to move numbers within 2 x 2 tables
- Such analyses often result in spuriously inflated accuracy estimates

26

Example: exclusion of indeterminates can inflate accuracy estimates

OPEN ACCESS Freely available online



The Impact of HIV Infection and CD4 Cell Count on the Performance of an Interferon Gamma Release Assay in Patients with Pulmonary Tuberculosis

Martine G. Aabye^{1*}, Pernille Ravn², George PrayGod³, Kidola Jeremiah³, Apolinary Mugomela⁴, Maria Jepsen⁵, Daniel Faurholt⁵, Nyagosya Range⁶, Henrik Friis⁵, John Chagalucha⁷, Aase B. Andersen¹

¹ Department of Infectious Diseases, University of Copenhagen, Rigshospitalet, Copenhagen, Denmark, ² Unit for Infectious Diseases Q, University of Copenhagen, Herlev Hospital, Herlev, Denmark, ³ National Institute for Medical Research, Mwanza Medical Research Centre, Mwanza, Tanzania, ⁴ Zonal Tuberculosis Reference Laboratory, Bugando Medical Centre, Mwanza, Tanzania, ⁵ Department of Human Nutrition, Faculty of Life Sciences, University of Copenhagen, Frederiksberg, Denmark, ⁶ National Institute for Medical Research, Muhimbili Medical Research Centre, Dar Es Salaam, Tanzania

Abstract

Background: The performance of the tuberculosis specific Interferon Gamma Release Assays (IGRAs) has not been sufficiently documented in tuberculosis- and HIV-endemic settings. This study evaluated the sensitivity of the QuantiFERON TB-Gold In-Tube (QFT-IT) in patients with culture confirmed pulmonary tuberculosis (PTB) in a TB- and HIV-endemic population and the effect of HIV-infection and CD4 cell count on test performance.

Methodology/Principal Findings: 161 patients with sputum culture confirmed PTB were subjected to HIV- and QFT-IT testing and measurement of CD4 cell count. The QFT-IT was positive in 74% (119/161; 95% CI: 67–81%). Sensitivity was higher in HIV-negative (75/93) than in HIV-positive (44/68) patients (81% vs. 65%; $p=0.02$) and increased with CD4 cell count in HIV-positive patients (test for trend $p=0.03$). 23 patients (14%) had an indeterminate result and this proportion decreased with increasing CD4 cell count in HIV-positive patients (test for trend $p=0.03$). Low CD4 cell count (<300 cells/ μ l) did not account for all QFT-IT indeterminate nor all negative results. Sensitivity when excluding indeterminate results was 86% (95% CI: 81–92%) and did not differ between HIV-negative and HIV-positive patients (88 vs. 83%; $p=0.39$).

Conclusions/Significance: Sensitivity of the QFT-IT for diagnosing active PTB infection was reasonable when excluding indeterminate results and in HIV-negative patients. However, since the test missed more than 10% of patients, its potential as a rule-out test for active TB disease is limited. Furthermore, test performance is impaired by low CD4 cell count in HIV-positive patients and possibly by other factors as well in both HIV-positive and HIV-negative patients. This might limit the potential of the test in populations where HIV-infection is prevalent.

- If indeterminates are included:
 - Sens = 74%
- If indeterminates are excluded:
 - Sens = 86%

27

J Clin Epidemiol Vol. 52, No. 12, pp. 1231–1237, 1999
Published by Elsevier Science Inc.



0895-4356/99/5—see front matter
PII S0895-4356(99)00101-8

Discrepant Analysis: A Biased and an Unscientific Method for Estimating Test Sensitivity and Specificity

Alula Hadgne*

CENTERS FOR DISEASE CONTROL AND PREVENTION, DIVISION OF STD PREVENTION, ATLANTA, GEORGIA

ABSTRACT. Discrepant analysis is a widely used technique for estimating test performance indices (sensitivity, specificity, etc.) of DNA-amplification tests for detecting infectious diseases. It has recently been claimed that the discrepant analysis-based estimates of specificity are typically less biased than those based on culture and that the discrepant analysis-based specificity shows little appreciable bias. In this article, I show that those conclusions are incorrect. Using a typical example from the published literature, I show that the discrepant analysis-based estimates of sensitivity and specificity can generate a significant and clinically important overestimation of the true sensitivity and specificity values. Moreover, I demonstrate that the concept of discrepant analysis is profoundly flawed and unscientific. It violates a fundamental principle of diagnostic testing—the principle that the new test should not be used to determine the true disease status. Thus, the major problem with discrepant analysis is not only that it is biased but that it is unscientific. Therefore, discrepant analysis should not be adopted for the evaluation of any diagnostic or screening test. J CLIN EPIDEMIOL 52;12:1231–1237, 1999. Published by Elsevier Science Inc.

KEY WORDS. Discrepant analysis, sensitivity, specificity, DNA-amplification tests, *Chlamydia trachomatis*

28

Lessons

- Early case-control studies are often used to promote and market tests
- But a large proportion of tests fail, once they are used in real world settings (e.g. large number of failed commercial serological tests)
- Case-control studies exaggerate accuracy estimates, especially if the two-gate approach is used
- Certain data analytic approaches can also inflate accuracy estimates
- Diagnostic studies can begin as case-control studies, but need to move beyond that to prospective studies in clinically indicated populations
- Even accuracy data may be insufficient to decide on clinical impact
- Regulatory agencies should demand prospective data and not just rely on case-control accuracy studies

29

Case study 3:

Where should TB tests be evaluated and which populations are appropriate?

30

It is not uncommon to see TB test evaluations where:

- Cases come from a high-incidence country and controls from a low-incidence country
- Tests work well in a low-incidence country and fall apart in a high-incidence country
- Tests that work well in immunocompetent persons fail in populations with high HIV prevalence

31

Example: cases from Zambia and controls from England

Table 1. Response rates in the ex-vivo enzyme-linked immunospot assay to ESAT-6- and CFP-10-derived peptides, recombinant ESAT-6 antigen and purified protein derivative in Zambian pulmonary tuberculosis patients and healthy Zambian and British adults, and tuberculin skin test results in healthy Zambian adults.

	Tuberculosis patients		Healthy Zambian adults		Healthy British adults (n = 40)
	HIV- (n = 11)	HIV+ (n = 39)	HIV- (n = 54)	HIV+ (n = 21)	
Response rates (%)					
ESAT-6 peptides	11 (100)	34 (87)	28 (52)	6 (29)	0 (0)
CFP-10 peptides	8 (73)	25 ^a (66)	34 (63)	7 (33)	0 (0)
Combined ESAT-6/CFP-10 peptides	11 (100)	35 (90)	37* (69)*	9* (43)*	0 (0)
ESAT-6 antigen	9 (82)	18 (46)	23 (43)	4 (19)	0 (0)
PPD	11 (100)	28 (72)	45 (83)	6 (29)	33 (83)
TST	–	–	28/35** (80)	5/14** (36)	–

PPD, Purified protein derivative; TST, tuberculin skin test.

^an = 38.

*P value for difference 0.064.

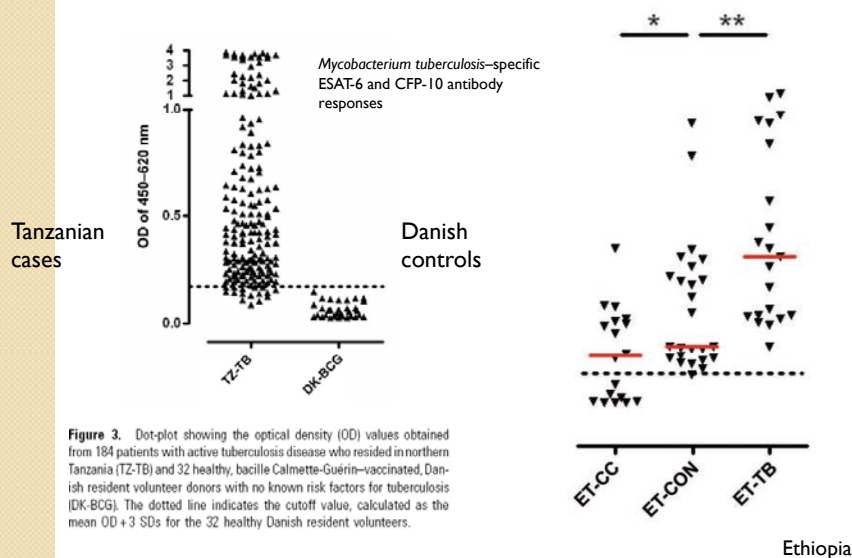
**P value for difference 0.0057.

Rapid detection of active and latent tuberculosis infection in HIV-positive individuals by enumeration of *Mycobacterium tuberculosis*-specific T cells

Ann L.N. Chapman^a, Mwansa Munkanta^b, Katalin A. Wilkinson^c,
Ansar A. Pathan^a, Katie Ewer^a, Helen Ayles^{b,c}, William H. Reece^a,
Alynn Mwanga^b, Peter Godfrey-Faussett^{b,c} and Ajit Lalvani^a

32

Lack of discrimination in TB endemic settings: example

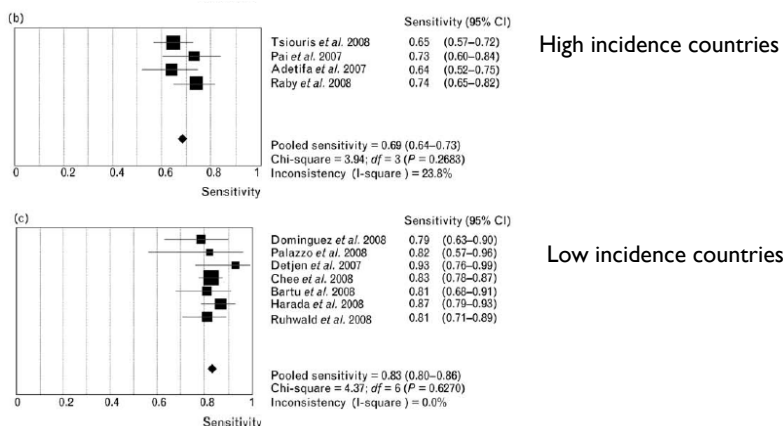


Hoff et al. Clinical Infectious Diseases 2007; 45:575–82

33

Variation in performance in high vs low endemic countries: example

T-cell interferon- γ release assays for the rapid immunodiagnosis of tuberculosis: clinical utility in high-burden vs. low-burden settings

Keertan Dheda^{a,b,c}, Richard van Zyl Smit^a, Motasim Badri^a and Madhukar Pai^d

34

HIV can prove to be the acid test for any test! Example of MycoDot

MOSSMAN ASSOCIATES
YOUR PARTNER IN BIOTECHNOLOGY

1 Village Circle
Sunnyvale, CA 94085
Phone: 916 438 4300
Email: info@mossmanassociates.com

Products: **Diagnostics** | **Chemistries** | **Equipment** | **Contract Services** | **Executive Placement**

Home | Products | About Us | Contact Us | Calendar | News | Site Map | Legal Page

Diagnostics

Immuno-Diagnostics

The MycoDot™ test employs liparabinomannan (LAM) antigen bound to plastic beads. When the beads are incubated in diluted serum, specific anti-LAM antibodies from the sample, if present, bind to the antigen. The beads are then washed to remove non-specific antibodies, and resuspended in a suspension of colored particles which bind to the bound anti-LAM antibodies. If a group of the specific antibodies are present in the serum sample, a colored spot will form when the antigen is exposed to the plastic beads. The sensitivity of the test is calibrated so that only cases of active mycobacterial disease such as tuberculosis will provide a positive reaction by MycoDot™.

Chemistries

MycoCheck™ is indicated for the semi-quantitative detection of mycobacterium tuberculosis complex (MTC) antigen in the verification of mycobacterium tuberculosis complex (MTC) antigen. The detection of mycobacterium tuberculosis complex (MTC) antigen can be estimated based on the color intensity of the reaction found on the MycoCheck™ test strips. Since it has been established that tuberculosis is one of the most significant causes of death and disease and that mycobacterium tuberculosis is a major cause of death and disease, MycoCheck™ provides a reliable indicator to the physician relative to the patient's potential risk level associated with these diseases and conditions.

See the MycoCheck™ insert and use as a screening or confirmatory product.
Obtain information on MycoCheck™
Download MycoCheck™ Package Insert
Download MycoCheck™ Package Insert
Download MycoCheck™ Package Insert
Download MycoCheck™ Package Insert

MycoDot was hailed to be a breakthrough because it was a simple dipstick test

Commercialized and marketed by Mossman Associates (with support of PATH)

Package insert: sensitivity of 70% and specificity of 95%

35

But when the test was evaluated in countries with high HIV prevalence, the performance was disastrous

Evaluation of the MycoDot™ test in patients with suspected tuberculosis in a field setting in Tanzania

G. R. Smit,¹ R. J. O'Brien,² G. S. Mfinanga,³ Y. A. Ispahani¹

¹National Institute for Medical Research, Dar Es Salaam, Tanzania; ²WHO Global Tuberculosis Programme, Geneva, Switzerland; ³National Tuberculosis and Leprosy Programme, Dar Es Salaam, Tanzania

SUMMARY

SETTING: Rapid, simple and inexpensive methods are needed to improve the diagnosis of tuberculosis in low-income countries. The MycoDot™ test has these characteristics.

OBJECTIVE: To assess the utility of the MycoDot™ test in screening patients with suspected tuberculosis.

DESIGN: Ambulatory patients presenting with symptoms of pulmonary tuberculosis were evaluated by physical examination and sputum acid-fast bacilli (AFB) microscopy. Separately, the MycoDot™ test was performed on whole blood. Patients with AFB-negative sputum were treated with a 10-day course of rifampicin. Those remaining symptomatic had a chest radiograph. All sputum specimens were cultured for mycobacteria. Patients with culture-negative tuberculosis and those without a tuberculosis diagnosis were re-evaluated at 2 months.

RESULTS: Among the 241 patients who were evaluated, the MycoDot™ test was positive in 26% of patients with AFB-positive/culture-positive tuberculosis, 7% with AFB-negative/culture-positive tuberculosis, 7% with culture-negative tuberculosis, 19% treated for tuberculosis who did not meet study case definition, and 16% without tuberculosis. Twenty-four patients did not complete the assessment. Test sensitivity was 16%, specificity 94%, and positive predictive value 45%. Sensitivity was highest (41%) in AFB-positive/HIV-negative patients and lowest (5%) in AFB-negative/HIV-positive patients.

CONCLUSION: The MycoDot™ test is not useful for the diagnosis of tuberculosis in sub-Saharan African countries, especially where HIV infection is prevalent.

KEY WORDS: tuberculosis, diagnosis, HIV, serology

Evaluation of a commercial immunodiagnostic kit incorporating liparabinomannan in the serodiagnosis of pulmonary tuberculosis in Ghana

S. D. Lartey,¹ E. H. Frimpong² and E. Nyarko³

¹Department of Medicine, School of Medical Sciences, University of Science and Technology, Kumasi, Ghana; ²Department of Microbiology, School of Medical Sciences, University of Science and Technology, Kumasi, Ghana; ³National Tuberculosis Control Programme, Ministry of Health, Accra, Ghana

Summary

We evaluated 'MycoDot', a commercially marketed immunodiagnostic test for tuberculosis which detects antibodies to liparabinomannan antigen. Serum was tested from 11 patients with newly diagnosed smear-positive pulmonary tuberculosis, of whom 10 were HIV-positive and 11 HIV-negative. Control sera were taken from 40 patients of whom 10 had active non-tuberculous lymphoproliferative disease and 30 patients had no respiratory disease. The test was found to have a very high specificity of 97.5% (95% CI 91.7–100%). However, the sensitivity in HIV-negative patients was 50% (95% CI 37.7–62.3%), and was substantially lower at 11% (95% CI 0–44%) in HIV-positive patients. In conclusion, 'MycoDot' was found to be a highly specific and easily performed assay. However, the poor sensitivity, especially in HIV-infected patients, renders it unlikely to be useful either as a primary or adjunctive diagnostic test for tuberculosis, particularly in countries with a high prevalence of HIV. A larger trial of this assay in Ghana was not deemed necessary.

Sens in HIV+ = 26%

Sens in HIV+ = 25%

Despite these results, the test is still available on the market!

36

Lessons

- TB evaluation studies must be done in high TB incidence countries, especially in high HIV prevalent settings
- Performance outcomes from low incidence countries may be deceptive and not reflect the performance in high incidence settings where the challenges include:
 - HIV
 - Severe TB
 - High background prevalence of TB infection
 - Widespread BCG vaccination
 - Malnutrition
 - Other diseases that can affect performance (e.g. worm infestations)
- If tests perform well in TB/HIV endemic countries, then they are likely to hold up well!

37

Case study 4:

Who should conduct TB diagnostic studies?

38

Industry involvement in TB diagnostic studies: example from IGRA literature

Annals of Internal Medicine

REVIEW

Systematic Review: T-Cell–based Assays for the Diagnosis of Latent Tuberculosis Infection: An Update

Madhukar Pai, MD, PhD; Alice Zwerling, MSc; and Dick Menzies, MD, MSc

Of the 38 studies in the meta-analysis, 21 (55%) had some sort of industry involvement or support, such as sponsorship, donation of test kits, participation in advisory boards, involvement of test developers, or ownership of patents.

41

Industry involvement in TB, HIV, Malaria studies and likely impact: McGill-TDR/WHO study

Table 10 Multivariate logistic regression results using authors' conclusion (dependent variable) and industry involvement, disease of interest and quality assessment variables (independent variables)
[n = 153]

Variable	OR	95% CI
<i>Industry involvement</i>		
• No	1.0	Reference group
• Yes	4.28	1.83 - 10.02
• NR	5.11	1.77 - 14.74

Pai, Fontela, Dendukuri, Ramsay, et al.

42

Industry involvement in TB studies and likely impact: commercial IGRAs

- We searched for cost-effectiveness studies on commercial IFN-gamma release assays
- We found a total of 10 studies
- Of these 6 studies had industry involvement of some sort
 - 2 of 6 had CEO of a test making company as author!
- Of the 6 studies with industry involvement: ALL concluded in favor of the commercial test and claimed superior cost-effectiveness
- Of the 4 independent studies, two were in favor of the test, and two were cautious and recommended a more selective use of the test

43

Industry involvement in TB studies and likely impact: commercial IGRAs

Studies with industry involvement

Direct costs of three models for the screening of latent tuberculosis infection

P. Wrighton-Smith* and J-P. Zellweger*

Cost-optimisation of screening for latent tuberculosis in close contacts

R. Diehl*, A. Nienhaus*, C. Lange* and T. Schaberg*

Cost-effectiveness of interferon- γ release assay testing for the treatment of latent tuberculosis

R. Diehl*, P. Wrighton-Smith* and J-P. Zellweger*

Cost-effectiveness of Interferon- γ Release Assay Screening for Latent Tuberculosis Infection Treatment in Germany*

Richard Diehl, MD, MPH, Albert Nienhaus, MD, MPH, and Robert Lodenkemper, MD, FCCP

Enhanced cost-benefit analysis of strategies for LTBI screening and INH chemoprevention in Germany

R. Diehl*, T. Schaberg*, R. Lodenkemper*, T. Welte*, A. Nienhaus*

Targeted screening and treatment for latent tuberculosis infection using QuantiFERON®-TB Gold is cost-effective in Mexico

J. L. Burgos,* J. G. Kahn,* S. A. Strathdee,* A. Valencia-Mendoza,* S. Bautista-Arredondo,* R. Llanado-Laborin,* R. Castañeda,* R. Delis,* R. S. Garfalo*

Independent studies

Interferon-gamma release assays and TB screening in high-income countries: a cost-effectiveness analysis

O. Oxlade, K. Schwartzman, D. Menzies
Respiratory Epidemiology and Clinical Research Unit, Montreal Chest Institute, McGill University, Montreal, Canada

Cost-effectiveness of Interferon Gamma Release Assays vs Tuberculin Skin Tests in Health Care Workers

Marie A. de Perio, MD, Joel Tsvet, MD, MPH, Gary A. Roselle, MD, Stephen M. Kralovic, MD, MPH, Mark H. Eckman, MD, MS

Cost-effectiveness of a new interferon-based blood assay, QuantiFERON®-TB Gold, in screening tuberculosis contacts

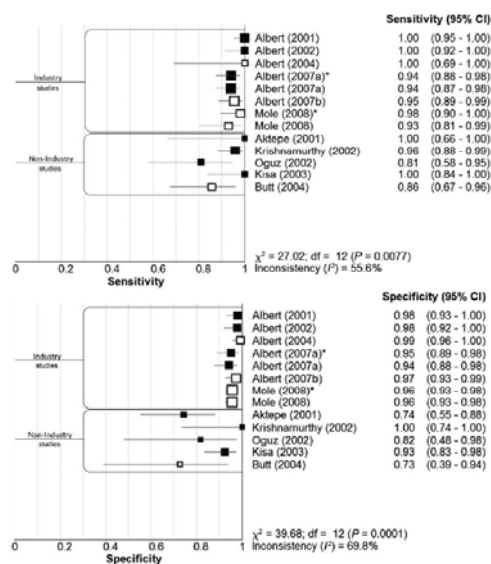
F. Marra,** C. A. Marra,** M. Sadatsafavi,* O. Morán-Mendoza,** V. Cook,** R. K. Elwood,** M. Marshall,** R. C. Brunham,** J. M. Fitzgerald**

Cost Effectiveness of Interferon- γ Release Assay for Tuberculosis Contact Screening in Japan

Akiko Kuroda,* Chisato Takahashi,* Takuro Shindo,* Sachiko Ohda,* Yasharu Takada,* and Tsuyoshi Fuku*

44

FASTPlaque tests for drug-resistant TB



Minion J et al. ITJLD 2010

45

Lessons

- When test developers do the studies, test performance is always good; performance is less optimal when others try to replicate the results
 - may be suppression of unfavourable data
 - may just be a learning curve issue (test developers, by definition, understand the test better and know how to make it work!)
- While industry is critical for test development and commercialization, test evaluations should, ideally, be done independent of industry support
- At the very least, industry involvement should be clearly disclosed in all publications and presentations
- Industry and test developers should definitely not be involved in guideline and policy development
 - At least 17 countries have guidelines and statements on IGRAs
 - Vast majority of these guidelines had no disclosures on conflicts of interest

46

Case study 5:

Can we trust the package insert?

47

Commercial package inserts always provide data on test accuracy: can we trust them?



**PACKAGE
INSERT**

For In Vitro Diagnostic Use



2009 package insert

TABLE 8. QuantiFERON®-TB Gold IT: Summary of results from clinical studies of subjects with culture-confirmed *M. tuberculosis* infection.

STUDY	QuantiFERON®-TB Gold IT			QuantiFERON®-TB Gold (liquid antigen)			TST (5mm)*	
	Pos	Neg	Ind	Pos	Neg	Ind	Pos	Neg
Australian	24	3	0	20	7	0	—	—
USA	47	11	3	34	10	6	60	19
Japanese	86	6	8	78	14	8	—	—
Overall Sensitivity	89% (157/177)			81% (132/163)			76% (60/79)	

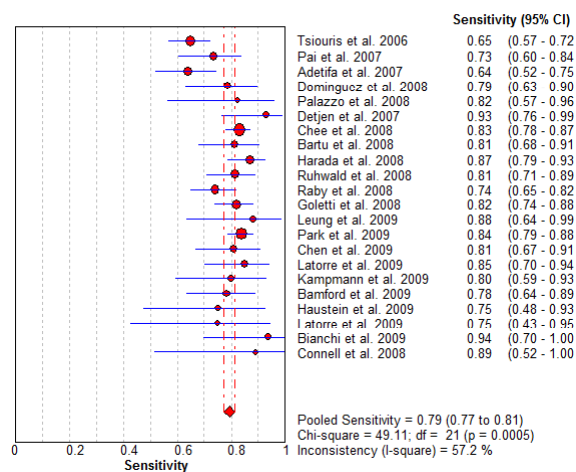
Pos – Positive; Neg – Negative; Ind – Indeterminate

* In the U.S. study of 86 *M. tuberculosis* patients, TST results were missing for 4 and invalid for 3.

According to the company, this test has 89% sensitivity in active TB

48

Updated meta-analyses on sensitivity of QuantiFERON-TB Gold In Tube



Pooled estimate was about 79%

49

More examples...

Test	Package insert sens	Package insert spec	Meta-analysis sens	Meta-analysis spec
FASTPlaque-Response	96 – 100%	99 – 100%	95%	95%
Anda-TB IgG	85 - 90%	85 - 100 %	60 - 75%	~90%
MycDot	70%	95%	26% - 76%	84% - 97%
Clearview TB ELISA	81% (HIV+)	93 – 98%	56% (HIV+)	95%
GenoType MDTBDrplus	99%	99%	98%	99%
Gen-Probe MTD	97% (S+) 72 (S-)	100% (S+) 99% (S-)	97% (S+) 76% (S-)	96% (S+) 95% (S-)

50

Lessons

- Company package inserts often present optimistic estimates based on small in-house evaluations that are usually sponsored by the companies
- Lab professionals and clinicians must be critical of advertised estimates of accuracy and performance
- Even when contradictory data are published, companies may not revise their package inserts or advertisements
- There is very little post-marketing surveillance of diagnostics and devices
- Regulatory agencies may not require companies to revise their package inserts
- Poorly performing tests may, in fact, never get pulled off the market

51

Case study 6:

Should we expect tests to be transferable and replicable?

Transferability: technologies that work well in the hands of developers will not necessarily work well everywhere

Eg. MODS, phage assays

52

THE NEW ENGLAND JOURNAL of MEDICINE

ORIGINAL ARTICLE

Microscopic-Observation Drug-Susceptibility Assay for the Diagnosis of TB

David A.J. Moore, M.D., Carlton A.W. Evans, M.D., Ph.D., Robert H. Gilman, M.D., Luz Caviedes, B.Sc., Jorge Coronel, B.Sc., Aldo Vivar, M.D., Eduardo Sanchez, M.D., Yvette Pilledo, M.D., Juan Carlos Saravia, M.D., Cayo Salazar, M.D., Richard Oberhelman, M.D., Maria-Graciela Holm-Delgado, M.Sc., Doris LaChira, M.D., A. Roderick Escombe, M.D., Ph.D., and Jon S. Friedland, M.D., Ph.D.

MODS: developed in Peru – performs excellent

Sensitivity better than LJ (98 vs. 84%)

Fast turnaround time (1 week vs. 6 weeks+)

Implemented in India – performs poorly

Sensitivity 80%

Issues with contamination

Issues with reliability

Diagnostic accuracy of the microscopic observation drug susceptibility assay: a pilot study from India

J. S. Michael,* P. Daley,* S. Kalaiselvan,* A. Latha,* J. Vijayakumar,* D. Mathai,* K. R. John,* M. Pai†

INT J TUBERC LUNG DIS 14(4):482-488
© 2010 The Union

53

Simple, phage-based (FASTPlaque) technology to determine rifampicin resistance of *Mycobacterium tuberculosis* directly from sputum

H. Albert,* A. Trollip,* T. Seaman,* R. J. Mole*

* Biotec Laboratories Ltd, c/o National Health Laboratory Service, Cape Town, Western Cape, South Africa;
† Biotec Laboratories Ltd, Ipswich, Suffolk, United Kingdom

FASTPlaque phage assay – performed well when done by industry

100% sens
100% spec

Implemented in Kenya – performs poorly

Despite upgrading the lab:

Low accuracy (31% sens; 95% spec)

Issues with contamination (nearly have were not interpretable)

Evaluation of FASTPlaqueTB™ to diagnose smear-negative tuberculosis in a peripheral clinic in Kenya

M. Bonnet,* L. Gagnidze,* F. Varaine,* A. Ramsay,*† W. Githui,* P. J. Guerin*

* Epicentre, Paris; *Medecine Sans Frontières, Paris, France; †Liverpool School of Tropical Medicine, Liverpool, UK;
*United Nations Children's Fund/United Nations Development Programme/World Bank/World Health Organization Special Programme for Research and Training for Tropical Diseases (TDR), Geneva, Switzerland; †Centre for Respiratory Diseases Research, Kenya Medical Research Institute, Nairobi, Kenya

SUMMARY

OBJECTIVE: To evaluate the performance and feasibility of FASTPlaqueTB™ to smear-negative tuberculosis (TB) suspects in a peripheral clinic after laboratory upgrading.

DESIGN: Patients with cough >2 weeks, two sputum smear-negative results, no response to 1 week of amoxicillin and abnormal chest X-ray were defined as smear-negative suspects. One sputum sample was collected, decontaminated and divided into two: half was tested with FASTPlaqueTB in the clinic laboratory and the other half was cultured on Löwenstein-Jensen medium in the Kenyan Medical Research Institute. Test sensitivity and specificity were evaluated in all patients and in human immunodeficiency virus (HIV) infected patients. Feasibility was assessed by the contamination rate and the resources required to upgrade the laboratory.

RESULTS: Of 208 patients included in the study, 16.2% were HIV-infected. Of 203 FASTPlaqueTB tests, 98 (48.3%) were contaminated, which interfered with result interpretation and led to the interruption of the study. Sensitivity and specificity were respectively 31.2% (95%CI 12.1-58.0) and 94.8% (95%CI 86.8-98.4) in all patients and 33.3% (95%CI 9.9-65.1) and 93.9% (95%CI 63.1-98.7) in HIV-infected patients. Upgrading the laboratory cost £20,000.

CONCLUSION: FASTPlaqueTB did not perform satisfactorily in this setting. If contamination can be reduced, in addition to laboratory upgrading, its introduction in peripheral clinics would require further assessment in smear-negative and HIV co-infected patients and test adaptation for biohazard use.

KEY WORDS: tuberculosis; phage-based test; smear microscopy; diagnosis; developing countries

54

Replication

- There are many examples of novel tests for TB that show great promise, but do not get replicated
- Or subsequent results are disappointing and commercialization is abandoned
- Or test may be quite good, but impossible to develop and manufacture in a cost-effective way
- Results in a graveyard of inexplicably abandoned diagnostics

55

Example: MPB64 skin patch test (Sequella Inc.)



Early data in 1998:

Sensitivity: 98%

Specificity: 100%

In 2010, still not commercially available – plans have been abandoned

INT J TUBERC LUNG DIS 2007;10(5):543-546
© 1998 RIATLID

MPB64 mycobacterial antigen: a new skin-test reagent through patch method for rapid diagnosis of active tuberculosis

R. M. Nakamura,* M. A. Velmonte,[†] K. Kawajiri,* C. F. Ang,[†] R. A. Frias,[†] M. T. Mendoza,[†] J. C. Montoya,[†] I. Honda,* S. Haga,[†] I. Toida*

*Japan BCG Laboratory, Kyosai-shi, Tokyo, Japan, [†]Infectious Disease Section, Philippine General Hospital, Manila, Philippines, [‡]National Institute of Infectious Diseases, Toyama, Shizuoka-ku, Tokyo, Japan

SUMMARY

SETTING: A collaborative study between the Japan BCG Laboratory, Tokyo, Japan, and the Infectious Disease Section, Philippine General Hospital, Manila, the Philippines. Tuberculous patients from four clinics in the vicinity of Manila, Our Lady of Grace Parish, Sto. Nino de Tondo Parish, the Canossa Health and Social Center, and the Health Care Development Center, were examined.

OBJECTIVE: To develop a new, simple and rapid diagnostic method for active tuberculosis. Subjects were tested for skin reaction to a special antigen, MPB64, by the patch test method instead of intradermal injection of purified protein derivative (PPD).

DESIGN: Fifty-three active tuberculosis patients and 41 healthy PPD-positive controls were tested to determine whether or not the reaction to MPB64 was positive only in active tuberculosis patients.

RESULTS: Fifty-two of the 53 active tuberculosis patients showed a positive reaction to MPB64, while none of the 41 PPD-positive controls did. The specificity of MPB64 to active tuberculosis was 100%, and the sensitivity was 98.1%. The efficacy of the test was 99.0%.

CONCLUSION: The patch test with MPB64 is a promising method for the diagnosis of active tuberculosis, distinguishing tuberculous patients from those who are infected but have not developed the disease, and also from BCG-vaccinated individuals. This new skin test is a subject for further evaluation and it is important to compare the results with PPD Mantoux.

KEY WORDS: MPB64; patch skin test; rapid diagnosis; active TB.

56

Lessons

- Many novel tests and tools are bound to fail
- We need to appreciate the “failure rate” of new tests and interventions
- Replication, in diverse settings, is required, before proceeding with commercialization and clinical use
- Transferability of technologies must receive attention; tests need to be robust if they have to work well in all settings
- Tests that work well in the hands of developers may not work well in field settings, especially in resource-limited countries
- Single studies are never sufficient for policy and guideline development; we need more extensive evidence
- Even accuracy data are not sufficient for evidence-based policies

57



**Roadmap for rolling out Xpert MTB/RIF
for rapid diagnosis of TB and MDR-TB**

6 December 2010



58

Will Xpert MTB/RIF survive “optimism bias”?

- Validation data and early demonstration data look very good
- Not much “real world” experience in resource-limited and routine programmatic settings
- Impact of “point-of-treatment” use is not demonstrated

