

Diagnostic/Biomarker Development Breakout

John J. Sninsky, PhD



September 19-20, 2019
Liver Forum
Washington D.C.

Breakout Session Objectives

1. Summarize the basic biomarker development principles including analytical performance, study design, biostatistics and levels of evidence
2. Recognize the frequent missteps in biomarker studies
3. Understand key elements for critique of biomarker manuscripts and peer-reviewed papers

Where Are We?

Clinical Chemistry 63:5
963-972 (2017)

Review

Waste, Leaks, and Failures in the Biomarker Pipeline

John P.A. Ioannidis^{1*} and Patrick M.M. Bossuyt²

BACKGROUND: The large, expanding literature on biomarkers is characterized by almost ubiquitous significant results, with claims about the potential importance, but few of these discovered biomarkers are used in routine clinical care.

CONTENTS: The pipeline of biomarker development includes several specific stages: discovery, validation, clinical translation, evaluation, implementation (and, in the case of nonutility, deimplementation). Each of these stages can be plagued by problems that cause failures of the overall pipeline. Some problems are nonspecific challenges for all biomedical investigation, while others are specific to the peculiarities of biomarker research. Discovery suffers from poor methods and incomplete and selective reporting. External independent validation is limited. Selection for clinical translation is often shaped by nonrational choices. Evaluation is sparse and the clinical utility of many biomarkers remains unknown. The regulatory environment for biomarkers remains weak and guidelines can teach biased or divergent recommendations. Removing inefficient or even harmful biomarkers that have been entrenched in clinical care can meet with major resistance.

SUMMARY: The current biomarker pipeline is too prone to failures. Consideration of clinical needs should become a starting point for the development of biomarkers. Improvements can include the use of more stringent methodology, better reporting, larger collaborative studies, careful external independent validation, preregistration, rigorous systematic reviews and umbrella reviews, pivotal randomized trials, and implementation and deimplementation studies. Incentives should be aligned toward delivering useful biomarkers.

© 2016 American Association for Clinical Chemistry

Progress in unraveling the molecular basis of diseases and advances in technology have fueled the search for novel

biomarkers in many diseases. There is hope that biomarkers will improve our ability to identify, manage, or prevent a wide range of conditions that jeopardize health.

Research in this field has expanded over the years to include measurements of increasing numbers of proteins (the more typical type of biomarker) and other types of molecules (metabolites, DNA genetic systems, different types of RNA molecules) that may serve as biomarkers. For proteins, mass spectrometry allows measurement of multiple analytes with possibly high sensitivity and selectivity, at impressive speed. Multiple peptides, proteins, and their isoforms can be analyzed simultaneously, seemingly allowing even more refined classifications of patients into different classes, or the monitoring of patients during the course of their disease or the management thereof (1). Similarly, recent technical advances in metabolome, genome, and transcriptome measurements have been impressive and raise new methodological and clinical challenges for harnessing this information (2, 3). Biomarkers are typically measured in biospecimens, but an expanded definition includes also other bioinformation, e.g., procured by imaging (4), sensors, or other measurement tools.

Despite enthusiasm and high prospects for biomarkers, this measurement revolution has not yet resulted into tangible health benefits. Several commentators have pointed to the fact that, despite massive investments of resources, the biomarker business has added very little to everyday clinical medicine so far (5–7). An evaluation of the indications and contraindications of all drugs considered by the European Medicines Agency (EMA) and published in 883 European public assessment reports and 2000 clinical trials found mentions of only 37 predictive biomarkers (8, 9).

Previous authors have already provided possible explanations for this failure to deliver. Some have offered potential accompanying solutions to remedy this situation (9–11). Most of the previous literature has focused on technical challenges pertaining to the analytical capacity and limitations of existing measurement methods or the difficulty in finding laboratory professionals with sufficient expertise in applying these laboratory methods. Calls are often made for the investment of even more resources into biomarker research, and several countries have launched new funding schemes, for biomarker research in general, or for programs for specific diseases and conditions. However, it is questionable whether overcoming the technical limitations or perpet-

SUMMARY: The current biomarker pipeline is too prone to failures. Consideration of clinical needs should become a starting point for the development of biomarkers. Improvements can include the use of more stringent methodology, better reporting, larger collaborative studies, careful external independent validation, preregistration, rigorous systematic reviews and umbrella reviews, pivotal randomized trials, and implementation and deimplementation studies. Incentives should be aligned toward delivering useful biomarkers.

¹Departments of Medicine, Health Research and Policy, and Statistics, and the Meta-Research Innovation Center at Stanford (METRICS), Stanford University, Stanford, CA; ²Department of Clinical Epidemiology, Biostatistics & Bioinformatics, Academic Medical Center, University of Amsterdam, Amsterdam, the Netherlands.

* Address correspondence to this author at: Stanford Prevention Research Center, Medical School Office Building, Room 306, 1265 Welch Rd., Stanford, CA 94305. E-mail: jioannid@stanford.edu

Received August 15, 2016; accepted November 30, 2016.
Previously published online at DOI: 10.1373/clinchem.2016.254649

© 2016 American Association for Clinical Chemistry

963

Common Missteps in Diagnostic Studies - 1

- Performance of test in Discovery set only (overfit test performance)
- Use 'normal' samples as comparator rather than differential diagnosis samples (exaggerated performance)
- Dissimilar Discovery, Validation and Clinical Use sets (inaccurate estimate of performance) or distribution of samples
- Mixture of Discovery and Validation sets (inaccurate estimate of performance, overfit; solely statistical cross-validation insufficient)
- Lack pre-specified clinical/statistical analysis plan (introduction of bias)
- Convenience or opportunistic samples (solely retrospective; not representative; inaccurate performance)
- Single center study rather than multi-center study (test robustness)
- Poorly validated analytical performance (inaccurate performance, robustness, transferability)

Common Missteps in Diagnostic Studies - 2

- Does not consider implications of pre-analytical variation of biomarker
- Samples tested with different versions of test (inaccurate performance)
- Small sample sets (likely bias and chance; lack generalizability)
- Provide clinical validity but not clinical utility (questionable reimbursement)
- Lacks attention to PPV or NPV for indication of test (actionability)
- Cost effectiveness not modeled (questionable reimbursement)
- Statistical analysis only includes ROC, or sensitivity and specificity (test performance but not patient performance)
- Lack actionable outcomes (what will clinician or patient do differently with information)
- Does not compare performance relative to single or combined routinely used tests or information (independence relative to presently used information)

Sea Change in Clinical Diagnostics

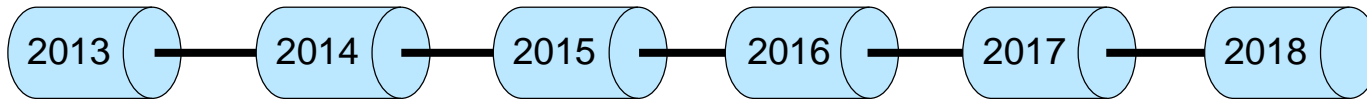
- Increased complexity of our understanding of disease
 - Multiple underlying etiologies
- Formal phased development of diagnostic tests similar to drug development has been adopted (AV, CV, CU, and Health Econ)
- High quality evidence needs to be provided by test service (LDT) or test kit (IVD) providers
- Clinical utility now required for reimbursement instead of only clinical validity as in past
- Evidence now understood to be a continuum and value-based
- Weave together CLIA (CLSI), FDA, NYSDOH, AMP, CAP and MoIDx NGS recommendations in some cases from related but distinct topic guidances to facilitate regulatory approvals
- Staged adoption of diagnostic tests considering indication benefit-risk ratio of managed patient



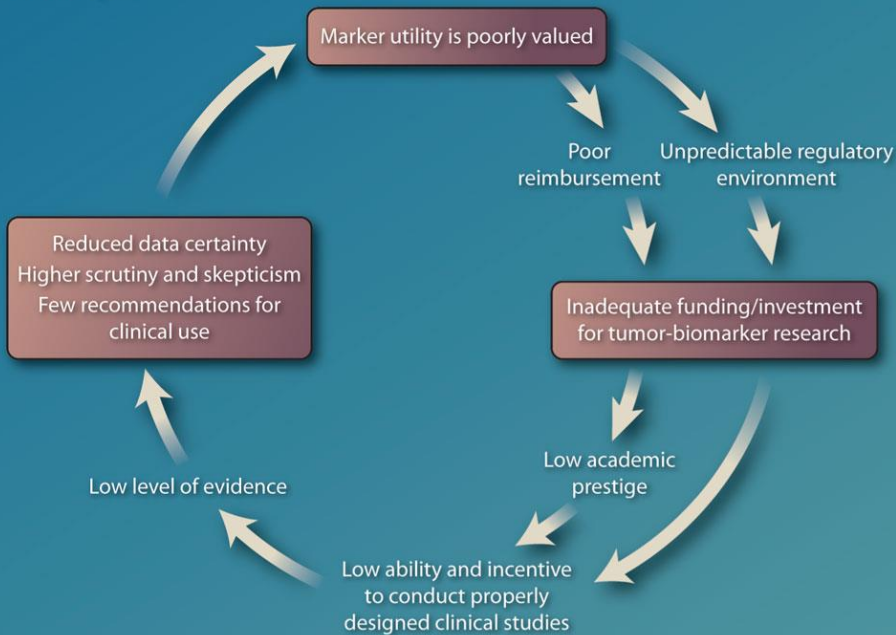
Why is Understanding Biomarker Regulatory Oversight and Reimbursement Essential?

- Even though diagnostics only makes up about 3% of healthcare expenditure, diagnostics informs how 65% of spend directed¹
- Concern that important medical insights are not being translated in a timely manner to patient care
- Translational, Clinical Development and Regulatory Sciences are evolving at a rapid pace
- Accelerated translation of discoveries into practice of medicine requires 'directed path' instead of 'exploratory walk'
- High quality, evidence-supported 'clinical-grade' biomarker assays require substantial investment
- If clinical-grade assays are not value priced, innovation from government and private industry will be stifled

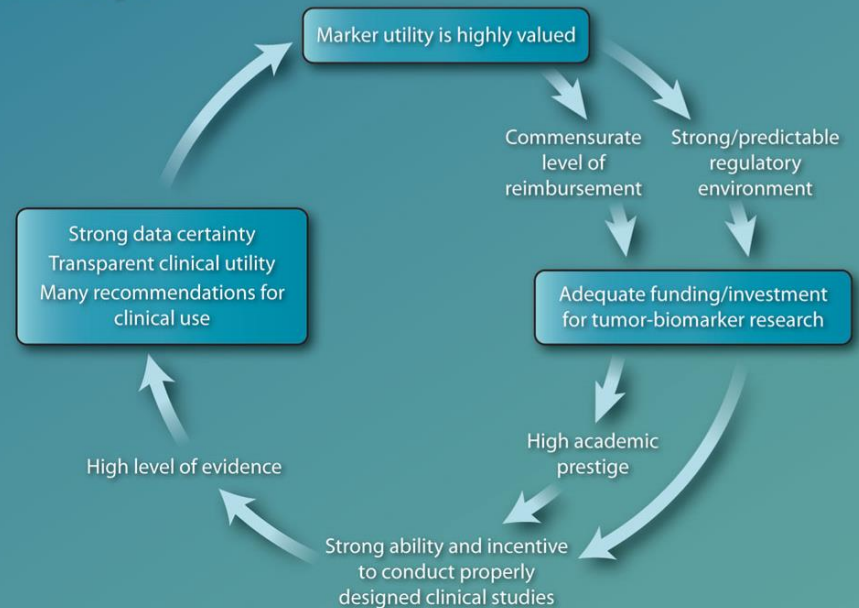
Appreciation of the Critical Importance of High Value Diagnostic Tests: Five years of change



A. Vicious cycle



B. Virtuous cycle



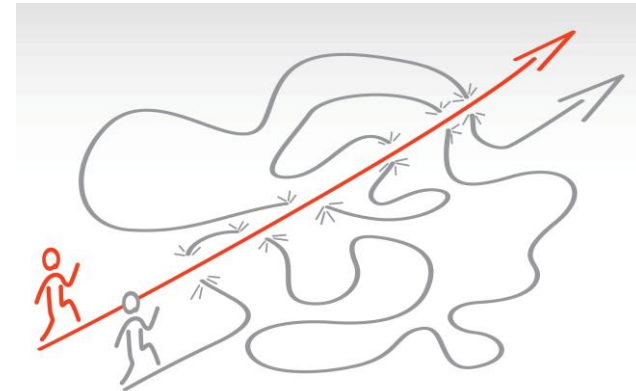
Biomarker 'Discovery' and 'Translation' Have Different and Discrete Objectives: equally valuable

- **Biomarker Discovery (Exploratory Walk)**

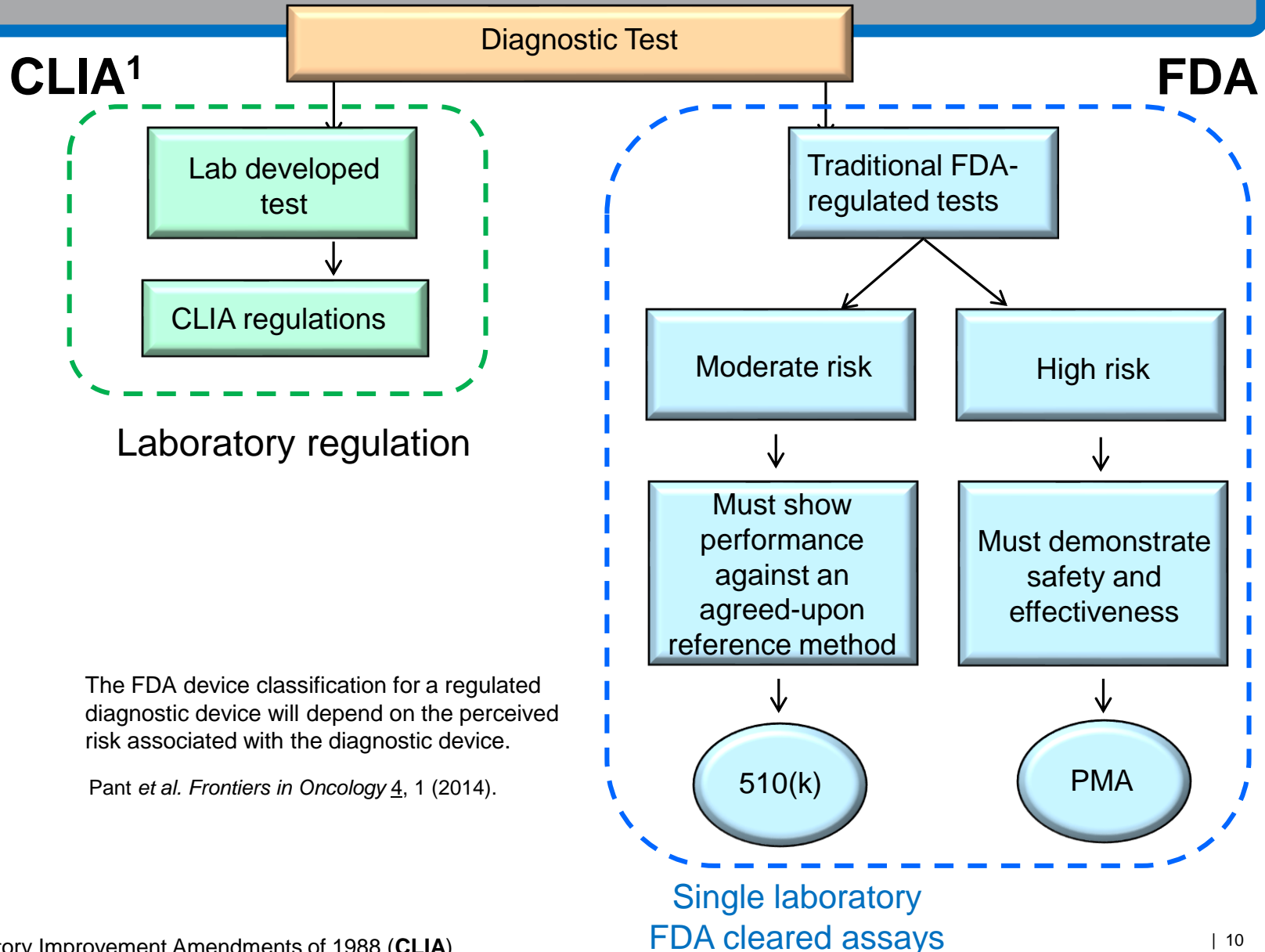
- Biomarker- or biology-centric
- Promises key insights into fundamental underlying pathophysiology
- Plethora of biologically plausible biomarkers
- Benefits from deep understanding of biology
- Correlations and group diagnostic metrics suffice

- **Biomarker Translation for Clinical Practice (Directed Path)**

- Clinical question-centric
- Promises improved patient management
- Few biomarkers that merit prioritization
- Benefits from translation and diagnostic development path knowledge
- Predictive values are most important for individual patients



Two Paths for Regulatory Oversight

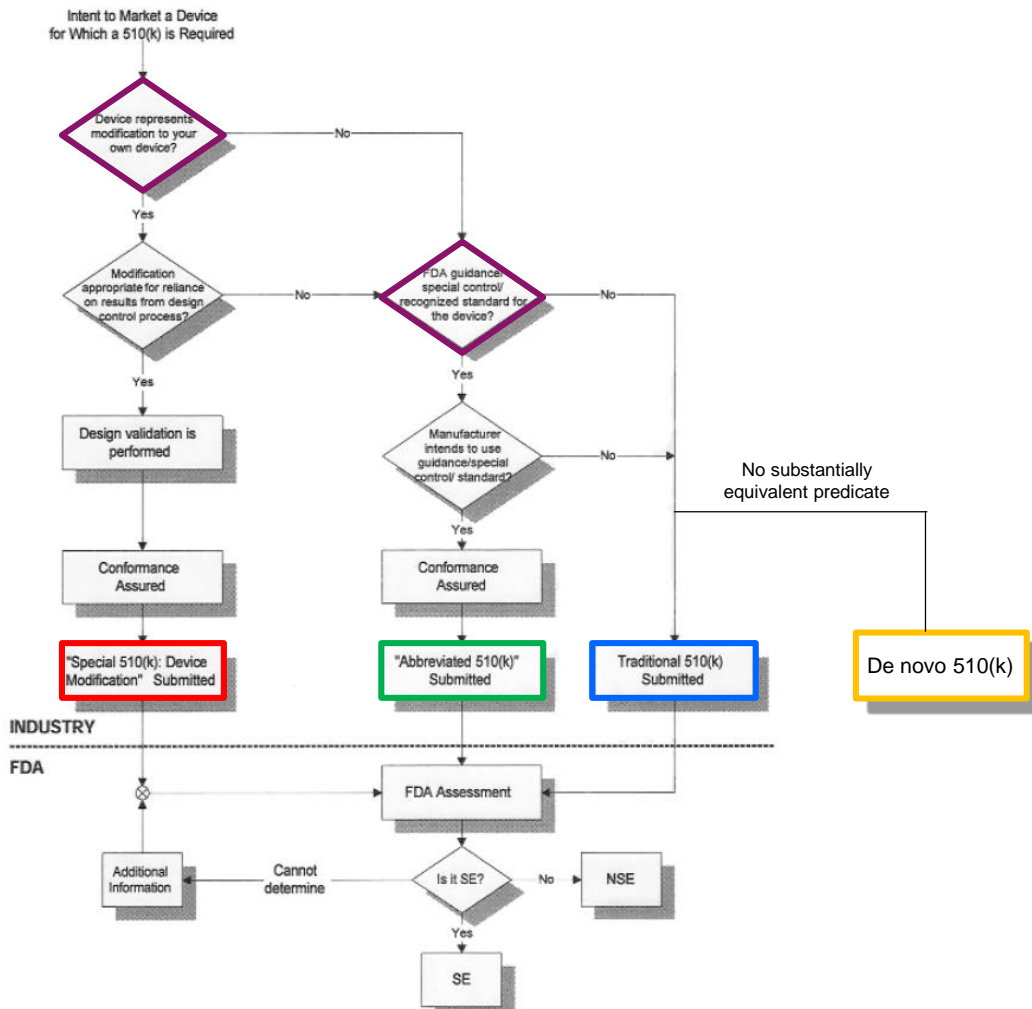


The FDA device classification for a regulated diagnostic device will depend on the perceived risk associated with the diagnostic device.

Pant *et al.* *Frontiers in Oncology* 4, 1 (2014).

¹ Clinical Laboratory Improvement Amendments of 1988 (CLIA)

The 510(k) Paradigm Continues to Evolve



- Traditional 510(k): substantial equivalent predicate prior 1976
- Special 510(k): Modification of vendors prior cleared product
- Abbreviated 510(k): 510(k) guidance/special controls or recognized standard available
- De novo 510(k): no substantial equivalent predicate; guidance/special controls not available; devices that are classified through the de novo process may be marketed and used as predicates for future 510(k) submissions

Biomarker Guidelines

Guideline Acronym	Guideline	Area	Reference
GRIPS	Genetic Risk Prediction Studies	genetic risk studies	Janssens <i>et al. Ann Inter Med</i> (2011).
STREGA	Strengthening the Reporting of Genetic Association Studies	genetic association studies	Little <i>et al. PLoS Med</i> (2009).
STROBE	Strengthening the Reporting of Observational Studies in Epidemiology	observational studies	Von Elm <i>et al. PLoS Med</i> (2007)
STARD	Standards for Reporting Diagnostic accuracy studies	diagnostic studies	Bossuyt <i>et al. Clin Chem</i> (2015).
REMARK	Reporting Recommendations for Tumor Marker Prognostic Studies	tumor marker prognostic studies	McShane <i>et al. Nat Clin Prac Urol</i> (2005).
EGAPP	Evaluation of Genomic Applications in Practice and Prevention; National Institutes of Health [NIH] (United States). Secretary's Advisory Committee on Genetic Testing [SACGT]; ACCE Framework (CDC: ACCE: a CDC-sponsored project (2000–2004)); http://www.cdc.gov/genomics/gtesting/ACCE/acce_proj.htm#T1 .	systematic process for assessing the available evidence regarding the validity and utility of rapidly emerging genetic tests for clinical practice	Teutsch <i>et al. Genetics in Medicine</i> (2009); Andrea Ferreira-Gonzalez <i>et al. Pers Med</i> (2010); Godard <i>et al. Genetics in Medicine</i> (2013)
		Pre-specified statistical analysis plans	Gamble <i>et al. JAMA</i> (2017); Ioannidis <i>JAMA</i> (2019); Yuan <i>et al. Ped Anesth</i> (2017).
		Catalog of reporting guidelines	Simera <i>et al. Eur J Clin Invest</i> (2010).
		Link to guidelines	http://www.equator-network.org/reporting-guidelines/stard/

Why Most Published Research Biomarker Studies Are Not Reproducible Nor Advance Field

- **Corollary 1:** The smaller the studies conducted in a scientific field, the less likely the research findings are to be true (reproducible).
- **Corollary 2:** The smaller the effect sizes in a scientific field, the less likely the research findings are to be true.
- **Corollary 3:** The greater the flexibility in designs, definitions, outcomes, and analytical modes in a scientific field, the less likely the research findings are to be true.
- **Corollary 4:** The assay does not address a clear unmet actionable diagnostic need.
- **Corollary 5:** The study does not accurately reflect the eventual intended use population.
- **Corollary 6:** The level of evidence is insufficient to be used in a clinical setting with confidence.

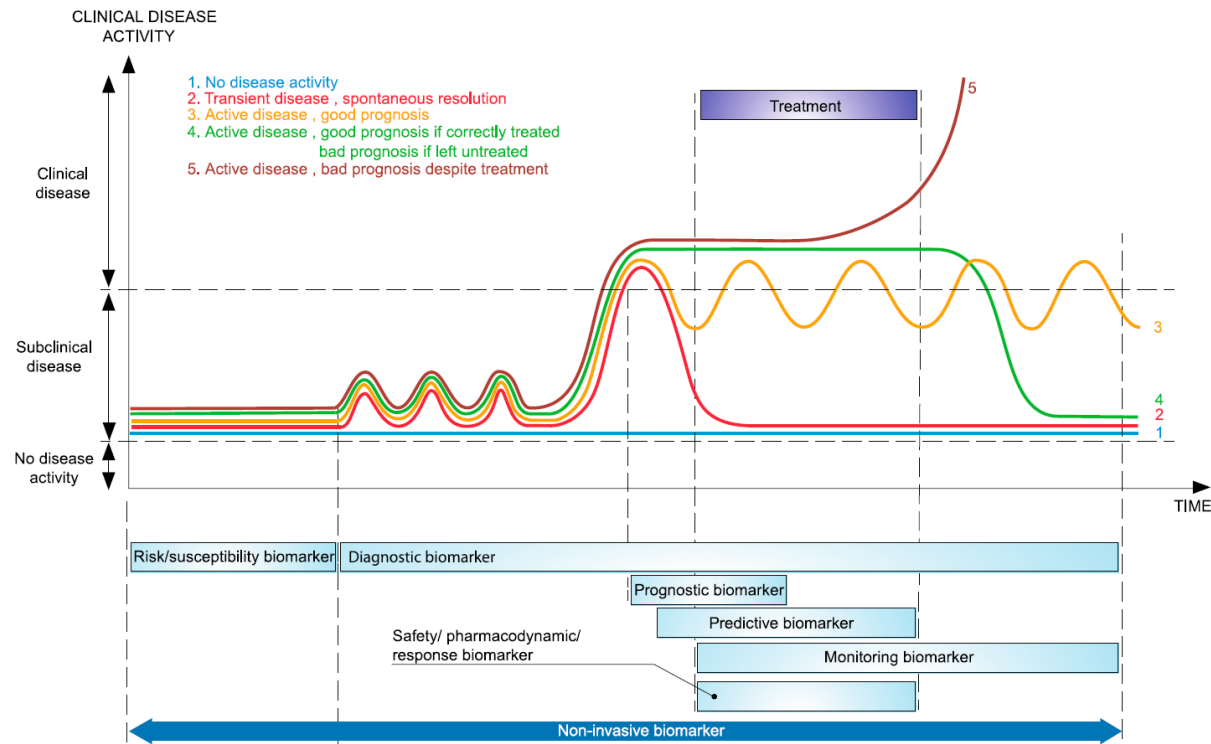
Types of Reproducibility

- Reproducibility of methods: the ability to understand or repeat as exactly as possible the experimental and computational procedures.
- Reproducibility of results: the ability to produce corroborating results in a new study, having followed the same experimental methods.
- Reproducibility of inferences: the making of knowledge claims of similar strength from a study replication.

CLSI and peer-reviewed assay precedent
inform assay development

Time Frames of Biomarkers

- Different biomarkers have value in distinct time frames
- Important to understand biological variation of a biomarker
- Biological variation may be due to temporary 'homeostatic disruption'
- Biomarkers for managing treatment are a compelling unmet need
- Statistical tools vary across types of biomarkers



Start in the Right Place

Identify the Right Question

- The need to answer a relevant clinical question. Make sure your solution will address a clinical question that will change what happens next for the patient. This may sound simple, but, looking backward, the diagnostics landscape is littered with companies that failed to take this point into account and instead started with a technology that never found a viable problem.

Understand the Needed Evidence

- Begin with the end result in mind. Impactful diagnostics efforts identify the critical sample sets upfront rather than address as an after-thought. You should determine your clinical utility study protocols as you develop your validation trials in order to maximize efficiency and increase your likelihood of receiving reimbursement earlier upon commercialization. You should decide on requisite evidence for reimbursement and how you will collect.

Commit to High Quality Studies

- Make an investment in high-quality studies that compare test performance against accepted reference and clinical truth (outcome) and publish in peer-reviewed journals. Cutting corners to save time or money when it comes to validating diagnostic tests simply won't work.

Intended Use Drives Evidentiary Studies

MSK-IMPACT Tumor Profiling Intended Use



Qualitative,
Targeted NGS



The MSK-IMPACT assay is a **qualitative** in vitro diagnostic test that uses **targeted** next generation sequencing of

Specimen type(s)



formalin-fixed paraffin-embedded tumor tissue matched with normal specimens from

Target population:
patients previously diagnosed



patients with solid malignant neoplasms to detect tumor gene alterations in a broad multi gene panel.

Variant types



The test is intended to provide information on **somatic mutations (point mutations and small insertions and deletions)** and **microsatellite instability**

Indication
(must include this statement)



for use by qualified health care professionals in accordance with professional guidelines, and is not conclusive or prescriptive for labeled use of any specific therapeutic product.

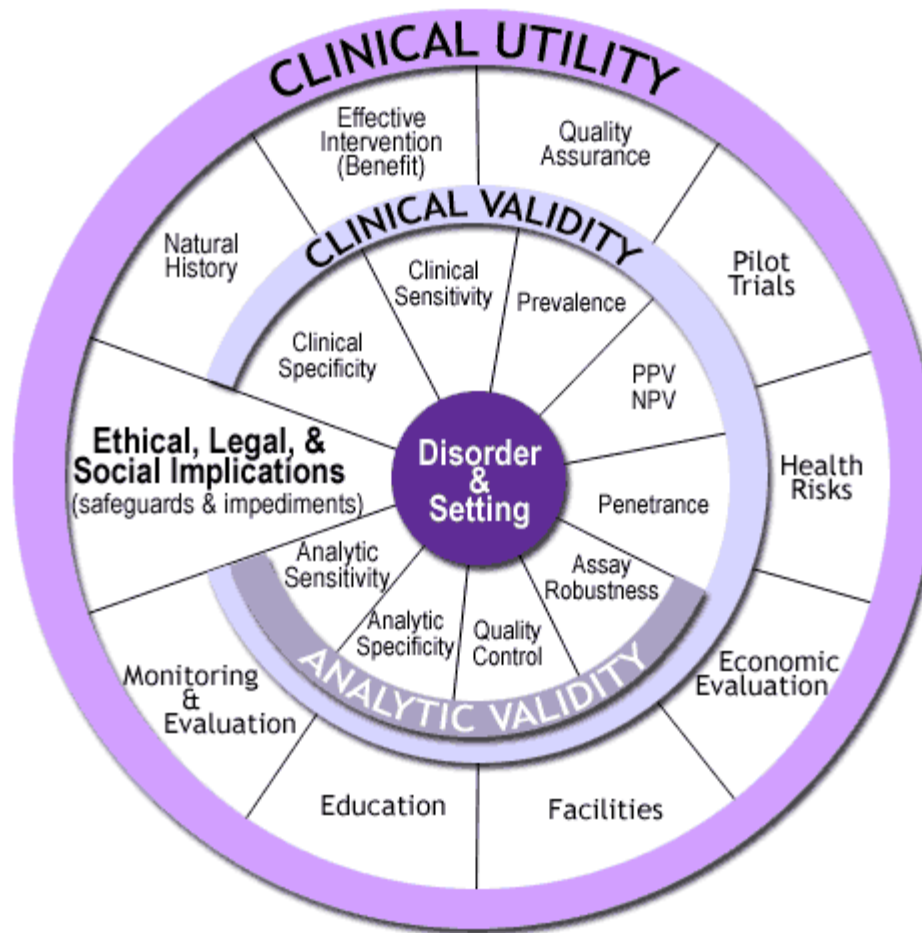
Single site and name



MSK-IMPACT is a **single-site assay performed at Memorial Sloan Kettering Cancer Center.**

www.fda.gov

ACCE Model



Steps in Diagnostic Test Development

- **Analytical Validity** refers to how well the test predicts the presence or absence of a biomarker. In other words, can the test accurately detect whether a specific biomarker is present or absent?
- **Clinical Validity** refers to how well the biomarker being analyzed is related to the presence, absence, or risk of a specific disease.
- **Clinical Utility** refers to whether the biomarker can provide clinically relevant information about diagnosis, treatment, management, or prevention of a disease that will be helpful to a patient, healthcare provider, or family member.
- **Cost Effectiveness** is the comparative analysis of two or more alternative interventions in terms of their health and economic consequences (Health Econ); factor in time horizon

Steps in Diagnostic Test Development

- **Analytical Validity** refers to how well the test predicts the presence or absence of a biomarker. In other words, can the test accurately detect whether a specific biomarker is present or absent?
Diagnostic biomarker assays are validated not biomarkers
Clinical-grade assays and software are critical, not research-grade versions
3 Rs of AV: repeatability, reproducibility and robustness
Follow CLSI documents
- **Clinical Validity** refers to how well the biomarker being analyzed is related to the presence, absence, or risk of a specific disease.
Specific intended uses are required rather than simple disease designation
Quality of evidence is critical
Training, Validation and Clinical use sets need to be independent with similar covariates
- **Clinical Utility** refers to whether the biomarker can provide clinically relevant information about diagnosis, treatment, management, or prevention of a disease that will be helpful to a patient, healthcare provider, or family member.
Clinical utility varies with stakeholder; payor critical due to reimbursement
Predictive values (NPV and/or PPV) are critical (prevalence determined), not Sen., Spec. and ROC
- **Cost Effectiveness** is the comparative analysis of two or more alternative interventions in terms of their health and economic consequences (Health Econ); factor in time horizon

Actionability: results that guide decision making

Actionability is an evolving concept and varies with patient, clinician, guideline committee, and payor

- Contextual for stage of disease (early vs advanced)
- Guidelines and FDA approved drug labels formally define accepted criteria
- Actionability is not binary but is best thought of as supported with a continuum of evidence
 - Fit-for-purpose (or matched) benefit – risk of managed patient group

A Question Driven Framework for Clinical Utility

- Who should be tested and under what circumstance?
- What does the test tell us that we did not know?
- Against what comparator is the test measured?
- Can we act on the information provided by the test?
- Will we act on the information provided by the test?
- What is the effectiveness of the action?
- Does the outcome of action change in a way in which we find value?

Assays: Clinical-grade vs Research-grade

- ‘Biomarkers’ are not validated, ‘biomarker assays’ are validated
- Clinical-grade assays are much more than just testing clinical samples
- Clinical-grade assays have to be of highest quality because they inform critical patient management decisions

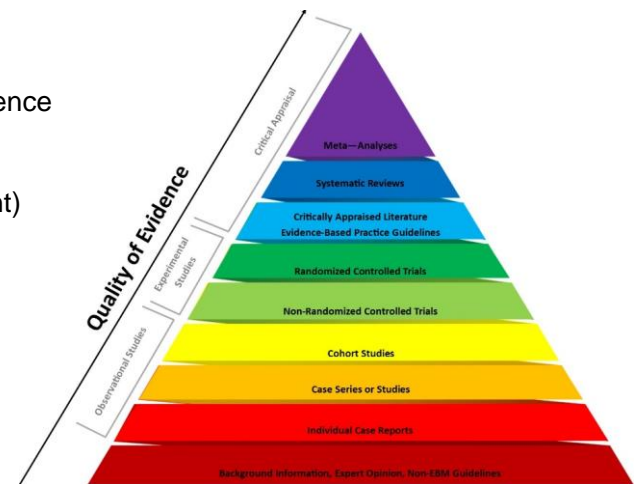
NGS assay	Research-grade	Clinical-grade	Comments
Reference materials	Internal specimens / External specimens	External standards; orthogonal technology validation	Ensures high test accuracy (Obtain reference standards through collaboration (e.g. Horizon Discovery))
Methods-based proficiency	Rarely used	Performed regularly	Ensures high test reproducibility (NIST-GIAB reference genome)
Information tracking systems	Sometimes used	Always use LIMS; some integration with EMRs	Ensures sample and reagent tracking; correct report for each patient sample
Bioinformatic analysis	Open source combined with subscription/license; frequently changing; & early adoption of new software/algorithms	Open source combined with subscription /license; use mature software and CDS Locked and change requires re-validation	Ensures test consistency and reproducibility (e.g. DNAnexus – platform also selected by FDA as part of precisionFDA initiative)
Validation of steps in process	Sometimes	Always	Follow applicable NGS recommendations/guidelines to ensure highest quality of the test
Documentation	No design control Little documentation	Yes Extensive	Formal methodology for test development (e.g. establish performance requirements, milestone progress reviews, documentation, etc.)

Software: Research-grade vs Clinical-grade

	Research-grade	Clinical-grade	Value
Software development Life Cycle (SDLC)	Not used	Follows SOP for SDLC	Development archive Includes phased design controls
Change control	Not used	Follows SOP; Documented	Archive of changes and verification
Design History File	No	Yes	Documented development
Documentation	Minimal and inappropriate for SW development	Extensive within-code documentation; Increased standardization and conventional for SW	Upgrade to commonly accepted practice for commercial use SW
Source code	R&D code gradually modified without complete cleaning	Production quality	Upgrade to commonly accepted architecture for commercial use SW
	Commonly have intertwined functions and logic	Modular design, clear logical flow	Simplifies maintenance, code inspection and targeted upgrading
	Commonly have opaque functions	Transparent functions	Simplifies maintenance, code inspection and targeted upgrading
	May have variables hard coded in (input or configuration data embedded directly into code)	Variables are soft coded that can be changed without going into the code (references data sources external to code)	Simplifies trouble shooting and facilitates upgrades
	Redundant codes are common	No redundancy	Simplifies trouble shooting, code inspection and facilitates upgrades
Computation	Non-parallelized computation is common	Discrete computation (parallelization through multi-threading or multiple instances)	Permits parallel processing to increase speed Allows inclusion of multiple types of tests for future updates
Processing	Mostly batch processing	Stream processing to have real time analysis	Rapid result turnaround time
Naming conventions	Cryptic and ambiguous	Standardized and conventional	Improves understanding code and facilitates code inspection
Public software tools (versions)	Dated versions frequently used	Up-to-date versions	Increased robustness with added features
	Usually includes obsolete code	No obsolete code	Makes code readable for review and inspection
Coding	Non-selective in language use, Research languages are commonly used (such as R)	Languages that are readable, computationally efficient and memory conscious (Python, C, C++)	Improves code versatility, speed and readability
Cloud integration	Difficult or requires rewrite	Flexibility to integrate into cloud (Platform as a Service (PaaS))	Accommodates scaling

Hierarchy of Evidence: Dated view of value

- **Meta-analysis of randomized control trials**
 - Highest level of evidence
- **Randomized control trial**
 - Prospective in design
 - High level of evidence (e.g. probability-based inference such as p-values and confidence intervals easily interpretable).
 - *Post hoc* analysis possible (e.g. pre-specified, avoid subgroups, use primary endpoint)
- **Observational cohort**
 - Prospective in design
 - Less likely to have masked bias
- **Case control study**
 - Retrospective in design
 - Susceptible to masked bias (e.g. survivorship, selection, ascertainment, drug treatment)
- **Anecdotal study**
 - Replication rarely reported



Where are 'adaptive' trials, observational registries, and EMR data (RWD) positioned in this hierarchy?

New Appreciation of Study Designs

- Randomized controlled trials can have compromised value
 - Include only narrowly defined, less ill patients (general validity in question)
 - Difficult to find time and funding for all trials desired
 - Not ‘real world’ studies
- Registries bring value to evidence collection
 - Permits collection of real world data to complement and extend RCT data
 - Facilitates collection of comprehensive and unbiased data on diagnostic tests to enhance the available body of evidence for informed patient management decisions
 - Provides insights into short and long-term outcomes
 - Allows health systems, clinicians, and patients to work together to create a setting for generating evidence in practice

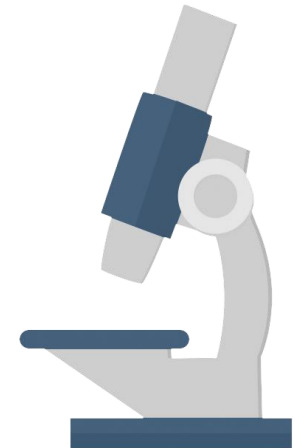
Reengineered Evidence Paradigm

Clinical trial design	Advantages	Disadvantages
Registry studies and observational studies	<ul style="list-style-type: none"> Ideal for description of standards Unselected patient populations (generalizable cohorts) Large number of events allows for the identification of rare events Inexpensive 	<ul style="list-style-type: none"> Data quality is variable and questionable Cannot be used for comparative outcomes research Confounding factors cannot be adjusted for, despite advanced statistical models
Randomized clinical trials	<ul style="list-style-type: none"> Well-designed studies with adequate power (gold-standard clinical design) Removes confounding factors 	<ul style="list-style-type: none"> Highly selected populations owing to specific inclusion and exclusion criteria Often performed at specialized study centres Often include surrogate end points Requires long time to plan and complete Expensive Often sponsored by industry (only studies with economic interest will be performed)
Registry-based randomized clinical trials	<ul style="list-style-type: none"> Randomization removes potential confounding Less-selected patient populations Large number of events allows for the identification of rare events Simple design Inexpensive 	<ul style="list-style-type: none"> Data quality might be variable and questionable Variables might not be well-defined Limited possibility for collection of detailed safety reporting, biospecimens, and pharmacokinetics or pharmacodynamic indices

James *et al.* *Nature* 12, 312 (2015).

Levels of Evidence: more nuanced perspective

- Similarity of inclusionary and exclusionary criteria (homogenous vs heterogeneous) across tested sample sets including intended use population
- Number of patients and events in each sample set
- Expected 'effect size' of tested diagnostic
- Expected number of events (prevalence)
- Single center versus multi-center collection
- Study Design used (retrospective (selection criteria), chronological, prospective, prospective-retrospective, single-arm with historical control, etc.)
- Study Objectives—Non-inferiority vs. Superiority vs. Equivalence
- Critical that pre-specified statistical analysis plans be used for validation^{1,2}



¹ Gamble *et al.* *JAMA* 318, 2337 (2017).

² Ioannidis *JAMA* (2019).

Level of Evidence: Scorecard

		SIZE OF TREATMENT EFFECT												
		CLASS I <i>Benefit >>> Risk</i> Procedure/Treatment SHOULD be performed/administered	CLASS IIa <i>Benefit >> Risk</i> <i>Additional studies with focused objectives needed</i> IT IS REASONABLE to perform procedure/administer treatment	CLASS IIb <i>Benefit ≥ Risk</i> <i>Additional studies with broad objectives needed; additional registry data would be helpful</i> Procedure/Treatment MAY BE CONSIDERED	CLASS III <i>No Benefit or CLASS III Harm</i>									
					<table border="1"> <thead> <tr> <th></th> <th>Procedure/ Test</th> <th>Treatment</th> </tr> </thead> <tbody> <tr> <td>COR III: No Benefit</td> <td>Not Helpful</td> <td>No Proven Benefit</td> </tr> <tr> <td>COR III: Harm</td> <td>Excess Cost w/o Benefit or Harmful</td> <td>Harmful to Patients</td> </tr> </tbody> </table>		Procedure/ Test	Treatment	COR III: No Benefit	Not Helpful	No Proven Benefit	COR III: Harm	Excess Cost w/o Benefit or Harmful	Harmful to Patients
	Procedure/ Test	Treatment												
COR III: No Benefit	Not Helpful	No Proven Benefit												
COR III: Harm	Excess Cost w/o Benefit or Harmful	Harmful to Patients												
ESTIMATE OF CERTAINTY (PRECISION) OF TREATMENT EFFECT	LEVEL A Multiple populations evaluated* Data derived from multiple randomized clinical trials or meta-analyses	<ul style="list-style-type: none"> Recommendation that procedure or treatment is useful/effective Sufficient evidence from multiple randomized trials or meta-analyses 	<ul style="list-style-type: none"> Recommendation in favor of treatment or procedure being useful/effective Some conflicting evidence from multiple randomized trials or meta-analyses 	<ul style="list-style-type: none"> Recommendation's usefulness/efficacy less well established Greater conflicting evidence from multiple randomized trials or meta-analyses 	<ul style="list-style-type: none"> Recommendation that procedure or treatment is not useful/effective and may be harmful Sufficient evidence from multiple randomized trials or meta-analyses 									
	LEVEL B Limited populations evaluated* Data derived from a single randomized trial or nonrandomized studies	<ul style="list-style-type: none"> Recommendation that procedure or treatment is useful/effective Evidence from single randomized trial or nonrandomized studies 	<ul style="list-style-type: none"> Recommendation in favor of treatment or procedure being useful/effective Some conflicting evidence from single randomized trial or nonrandomized studies 	<ul style="list-style-type: none"> Recommendation's usefulness/efficacy less well established Greater conflicting evidence from single randomized trial or nonrandomized studies 	<ul style="list-style-type: none"> Recommendation that procedure or treatment is not useful/effective and may be harmful Evidence from single randomized trial or nonrandomized studies 									
	LEVEL C Very limited populations evaluated* Only consensus opinion of experts, case studies, or standard of care	<ul style="list-style-type: none"> Recommendation that procedure or treatment is useful/effective Only expert opinion, case studies, or standard of care 	<ul style="list-style-type: none"> Recommendation in favor of treatment or procedure being useful/effective Only diverging expert opinion, case studies, or standard of care 	<ul style="list-style-type: none"> Recommendation's usefulness/efficacy less well established Only diverging expert opinion, case studies, or standard of care 	<ul style="list-style-type: none"> Recommendation that procedure or treatment is not useful/effective and may be harmful Only expert opinion, case studies, or standard of care 									
Suggested phrases for writing recommendations		should is recommended is indicated is useful/effective/beneficial	is reasonable can be useful/effective/beneficial is probably recommended or indicated	may/might be considered may/might be reasonable usefulness/effectiveness is unknown/unclear/uncertain or not well established	COR III: No Benefit is not recommended is not indicated should not be performed/administered/other is not useful/beneficial/effective	COR III: Harm potentially harmful causes harm associated with excess morbidity/mortality should not be performed/administered/other								
Comparative effectiveness phrases [†]		treatment/strategy A is recommended/indicated in preference to treatment B treatment A should be chosen over treatment B	treatment/strategy A is probably recommended/indicated in preference to treatment B it is reasonable to choose treatment A over treatment B											

Statistical Metrics for Test Performance

- Prioritize individual classification over group averages
- No single statistical measure provides sufficient insight
- Predictive values (NPV and PPV) are more important than sensitivity and specificity (clinically relevant)
- ROC curves are informative but not directly clinically relevant
- Multivariate analysis with standard measures are critical
- Methods based on risk stratification have recently been proposed to compare models
 - reclassification calibration statistic
- Bayesian models for diagnostic test performance provide key insights (conditional probabilities; likelihood ratios)
- Explore integration of conventional factors and molecular biomarkers

Redefined Statistical Threshold

- Set statistical threshold at 0.005
- More focus on effect sizes and confidence intervals, treating the P value as a continuous measure
- Proposal should not be used to reject publications of novel findings with $0.005 < P < 0.05$ properly labelled as suggestive evidence
- Reminder that failing to reject the null hypothesis does not mean accepting the null hypothesis

Redefine statistical significance

We propose to change the default P -value threshold for statistical significance from 0.05 to 0.005 for claims of new discoveries.

Daniel J. Benjamin, James O. Berger, Magnus Johannesson, Brian A. Nosek, E.-J. Wagenmakers, Richard Berk, Kenneth A. Bollen, Björn Brembs, Lawrence Brown, Colin Camerer, David Cesarini, Christopher D. Chambers, Merlise Clyde, Thomas D. Cook, Paul De Boeck, Zoltan Dienes, Anna Dreber, Kenny Easwaran, Charles Efferson, Ernst Fehr, Fiona Fidler, Andy P. Field, Malcolm Forster, Edward I. George, Richard Gonzalez, Steven Goodman, Edwin Green, Donald P. Green, Anthony Greenwald, Jarrod D. Hadfield, Larry V. Hedges, Leonhard Held, Teck Hua Ho, Herbert Hoijtink, Daniel J. Hruschka, Kosuke Imai, Guido Imbens, John P. A. Ioannidis, Minjeong Jeon, James Holland Jones, Michael Kirchner, David Laibson, John List, Roderick Little, Arthur Lupia, Edouard Machery, Scott E. Maxwell, Michael McCarthy, Don Moore, Stephen L. Morgan, Marcus Munafo, Shinichi Nakagawa, Brendan Nyhan, Timothy H. Parker, Luis Pericchi, Marco Perugini, Jeff Rouder, Judith Rousseau, Victoria Savalei, Felix D. Schönbrodt, Thomas Sellke, Betsy Sinclair, Dustin Tingley, Trisha Van Zandt, Simine Vazire, Duncan J. Watts, Christopher Winship, Robert L. Wolpert, Yu Xie, Cristobal Young, Jonathan Zinman and Valen E. Johnson

The lack of reproducibility of scientific studies has caused growing concern over the credibility of claims of new discoveries based on 'statistically significant' findings. There has been much progress toward documenting and addressing several causes of this lack of reproducibility (for example, multiple testing, P -hacking, publication bias and under-powered studies). However, we believe that a leading cause of non-reproducibility has not yet been adequately addressed: statistical standards of evidence for claiming new discoveries in many fields of science are simply too low. Associating statistically significant findings with $P < 0.05$ results in a high rate of false positives even in the absence of other experimental, procedural and reporting problems.

For fields where the threshold for defining statistical significance for new discoveries is $P < 0.05$, we propose a change to $P < 0.005$. This simple step would immediately improve the reproducibility of scientific research in many fields. Results that would currently be called significant but do not meet the new threshold should instead be called suggestive. While statisticians have known the relative weakness of using $P \approx 0.05$ as a threshold for discovery and the proposal to lower it to 0.005 is not new^{1,2}, a critical mass of researchers now endorse this change.

We restrict our recommendation to claims of discovery of new effects. We do

not address the appropriate threshold for confirmatory or contradictory replications of existing claims. We also do not advocate changes to discovery thresholds in fields that have already adopted more stringent standards (for example, genomics and high-energy physics research; see the 'Potential objections' section below).

We also restrict our recommendation to studies that conduct null hypothesis significance tests. We have diverse views about how best to improve reproducibility, and many of us believe that other ways of summarizing the data, such as Bayes factors or other posterior summaries based on clearly articulated model assumptions, are preferable to P values. However, changing the P value threshold is simple, aligns with the training undertaken by many researchers, and might quickly achieve broad acceptance.

Strength of evidence from P values

In testing a point null hypothesis H_0 against an alternative hypothesis H_1 based on data x_{obs} , the P value is defined as the probability, calculated under the null hypothesis, that a test statistic is as extreme or more extreme than its observed value. The null hypothesis is typically rejected — and the finding is declared statistically significant — if the P value falls below the (current) type I error threshold $\alpha = 0.05$.

From a Bayesian perspective, a more direct measure of the strength of evidence for H_1 relative to H_0 is the ratio of their

probabilities. By Bayes' rule, this ratio may be written as:

$$\frac{\Pr(H_1 | x_{\text{obs}})}{\Pr(H_0 | x_{\text{obs}})} = \frac{f(x_{\text{obs}} | H_1)}{f(x_{\text{obs}} | H_0)} \times \frac{\Pr(H_1)}{\Pr(H_0)} \quad (1)$$

= BF \times (prior odds)

where BF is the Bayes factor that represents the evidence from the data, and the prior odds can be informed by researchers' beliefs, scientific consensus, and validated evidence from similar research questions in the same field. Multiple-hypothesis testing, P -hacking and publication bias all reduce the credibility of evidence. Some of these practices reduce the prior odds of H_1 relative to H_0 by changing the population of hypothesis tests that are reported. Prediction markets³ and analyses of replication results⁴ both suggest that for psychology experiments, the prior odds of H_1 relative to H_0 may be only about 1:10. A similar number has been suggested in cancer clinical trials, and the number is likely to be much lower in preclinical biomedical research⁵.

There is no unique mapping between the P value and the Bayes factor, since the Bayes factor depends on H_0 . However, the connection between the two quantities can be evaluated for particular test statistics under certain classes of plausible alternatives (Fig. 1).

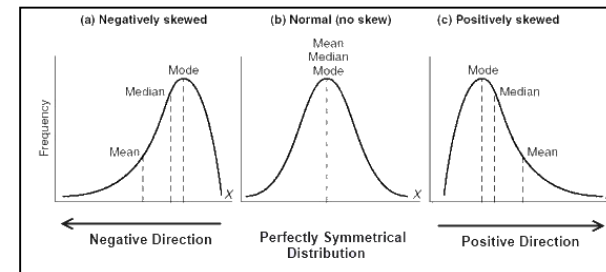
Parametric vs Non-parametric Analyses

- **Parametric**

- Uses the mean¹ of a sample set
- Normally distributed features or covariates
- Statistical tests (e.g. Two-sample t-test, Paired t-test, Analysis of variance (ANOVA, Pearson coefficient of correlation, etc.)
- Particularly worrisome for small sample sizes
- More power for same sample size
- If the data deviate strongly from the assumptions of a parametric procedure, using the parametric procedure could lead to incorrect conclusions.

- **Non-parametric**

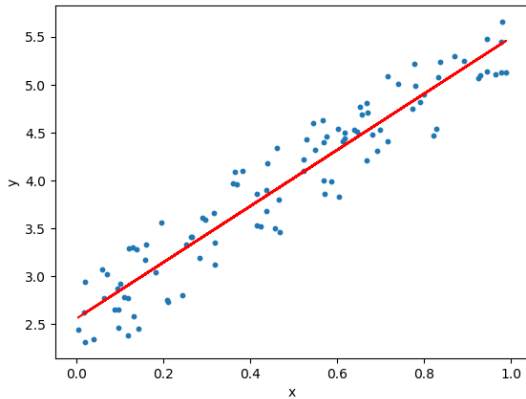
- Uses the median¹ of a sample set
- Unknown or not normally distributed features (covariates)
- Statistical tests (e.g. Wilcoxon rank-sum test, Wilcoxon signed-rank test, Kruskal-Wallis test, Spearman's rank correlation, etc.)



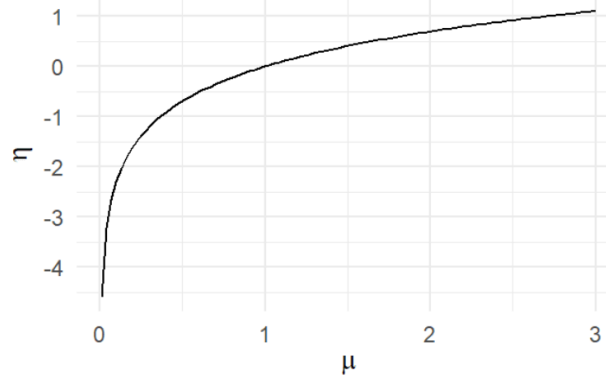
¹Mean and median are different for a sample set if distribution is skewed

Consider Nature of Interaction of Covariate and Outcome

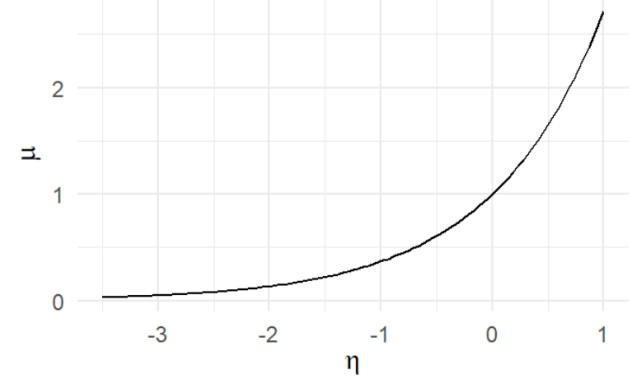
Linear



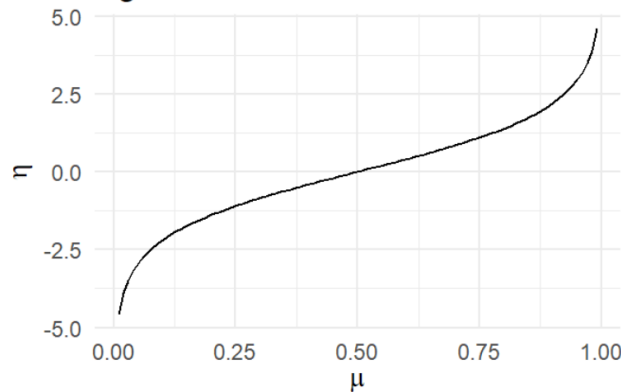
log link function



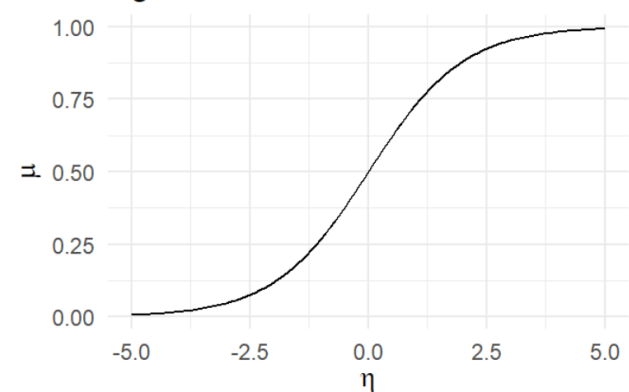
exponential mean function



logit link function

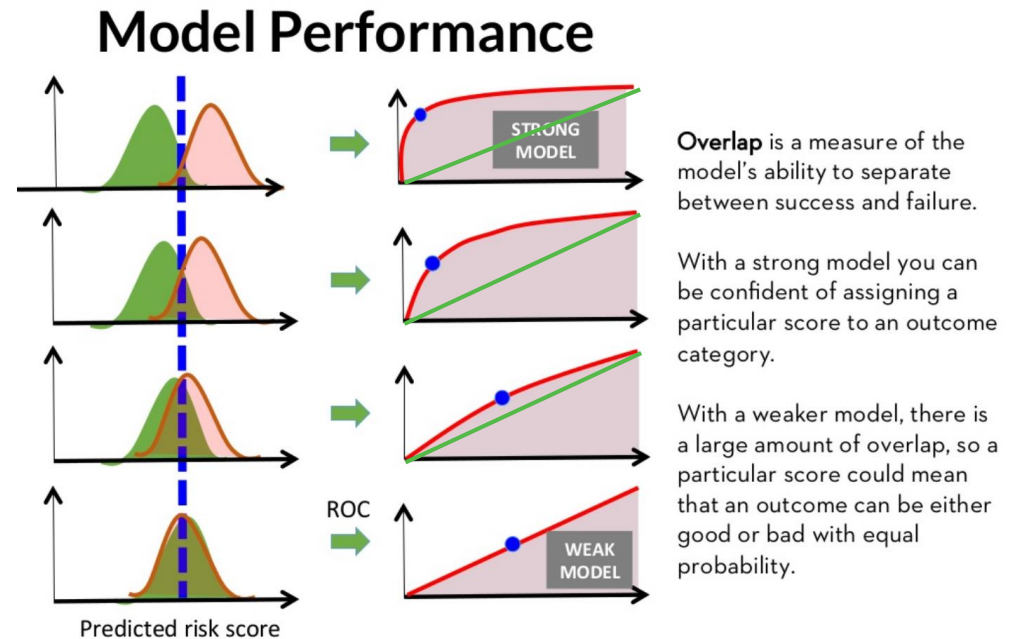
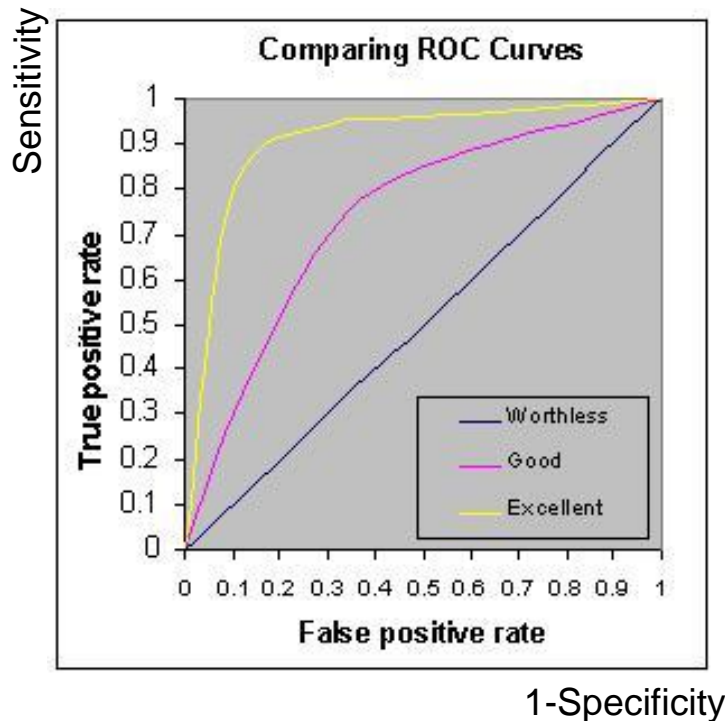


logistic mean function



Incorrect transform can lead to inaccurate results
Difficult to identify transform with small sample sets

Receiver Operator Curves (ROC)(AUC) or c-statistic



A rough guide for classifying the accuracy of a diagnostic test is:

- 0.90-1.00 = excellent
- 0.80-.90 = good
- 0.70-.80 = fair
- 0.60-.70 = poor
- 0.50-.60 = likely random

Threshold independent technique to visualize dichotomous diagnostic test performance

Permits selection of cutpoints for dichotomous categorization

Not a Means to an End

Same means, different variance

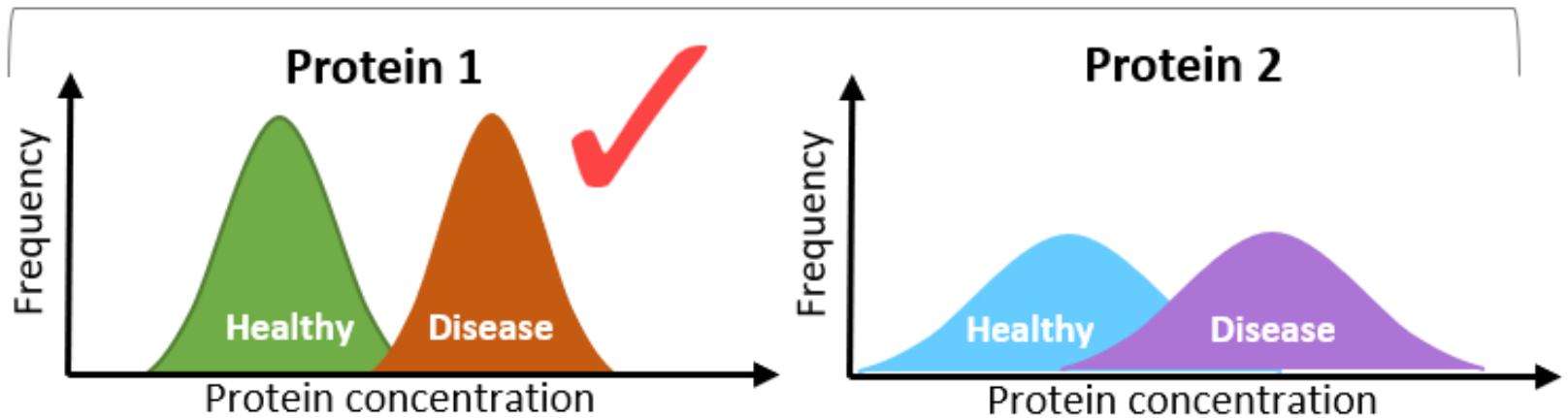
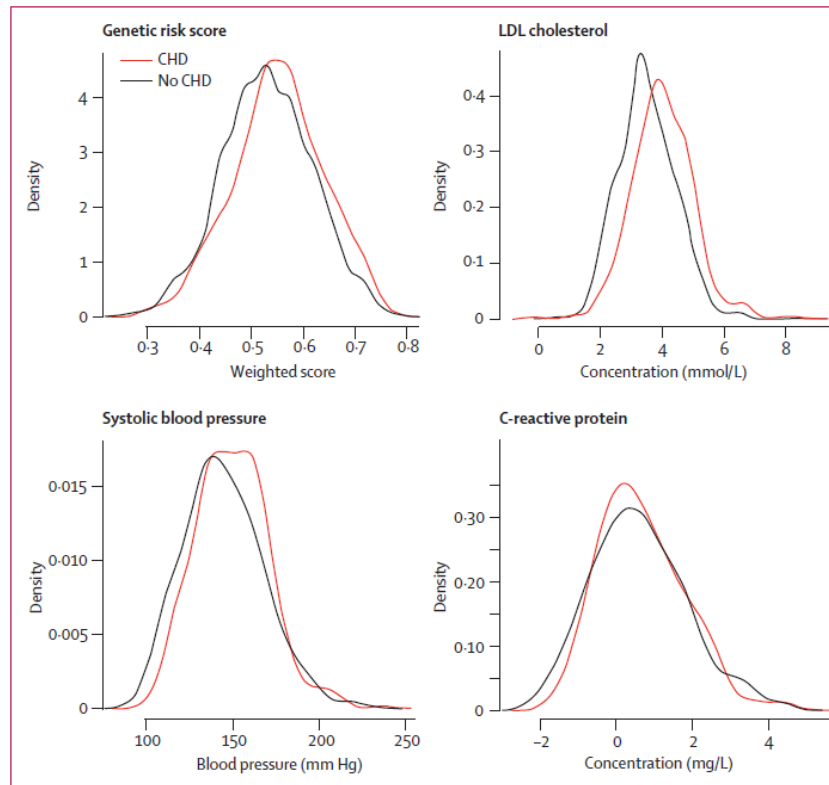


Figure 1. Overlapping histogram plots for concentrations of protein 1 in different populations.

Figure 2. Overlapping histogram plots for concentrations of protein 2 in different populations.

Use Perspective to Decide on ROC Performance

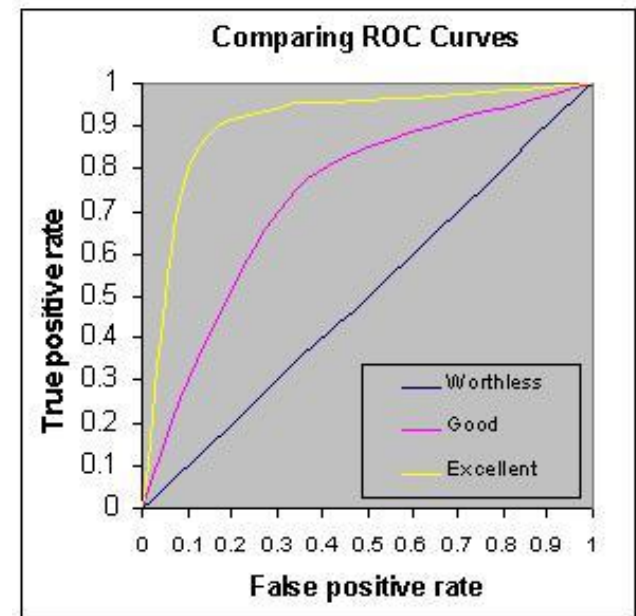


Distributions at baseline of genetic risk score, LDL cholesterol, systolic blood pressure, and log-transformed C-reactive protein by 10-year incident coronary heart disease event status in FINRISK 1992 and 1997 cohorts

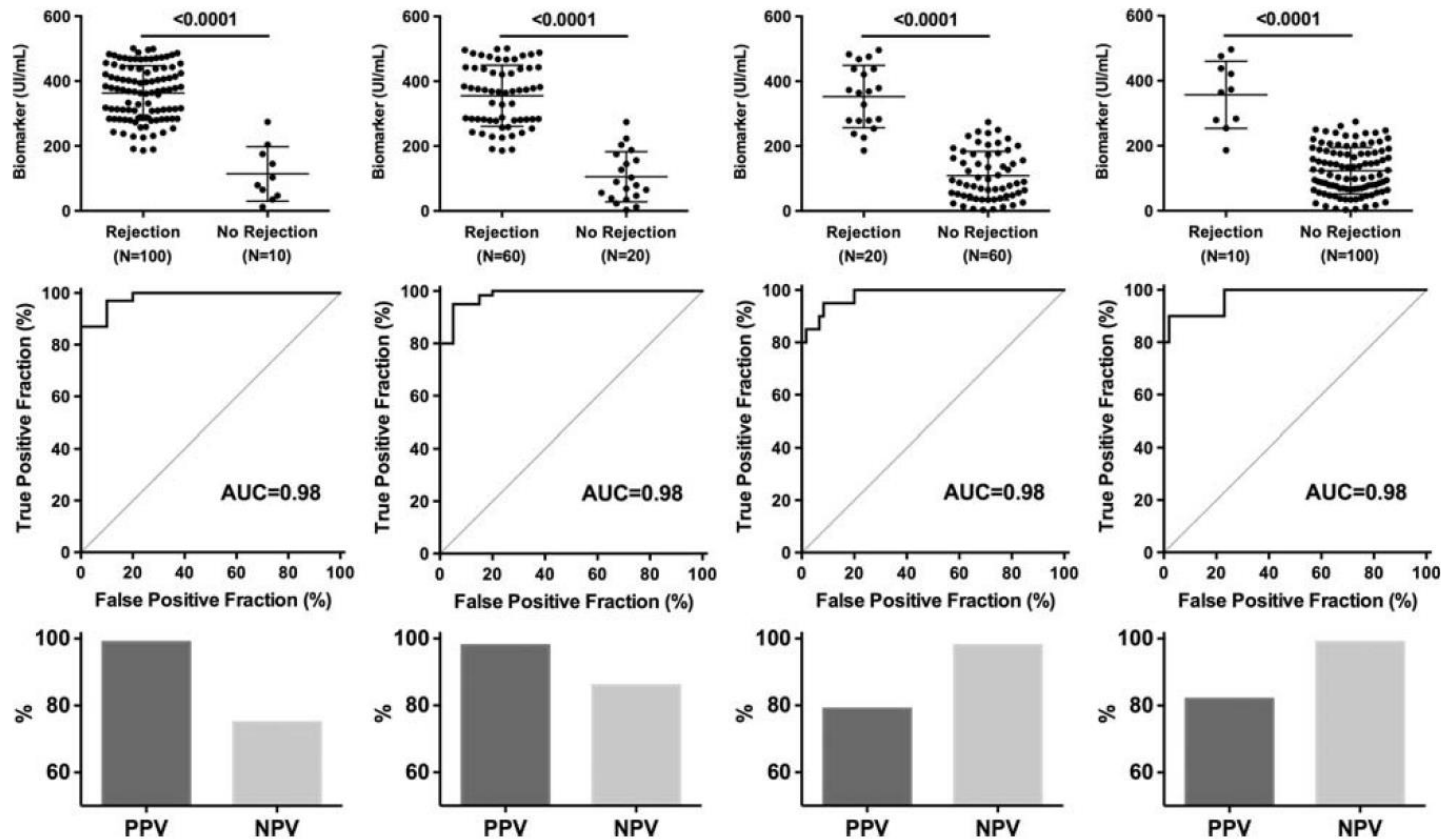
Ripatti *et al. Lancet* [376](#),1393 (2010).

AUC-ROC is not a Directly Clinically Relevant Diagnostic Metric

- As with any statistical metric, paucity of data compromises confidence of result
- ROC plots false positives (1-specificity) versus true positives (sensitivity) for every possible cutoff including regions not clinically relevant
- Requires highly accurate and related reference method to be informative
- A test with high sensitivity may have an identical or similar AUC to a test with high specificity
- Binary interpretation compromised (“Dichotomania”)
- Weights false positives and false negatives equally
- Does not address predictive values critical to ruling-in and ruling-out a diagnosis
- Insensitive to changes in absolute risk of tests compared

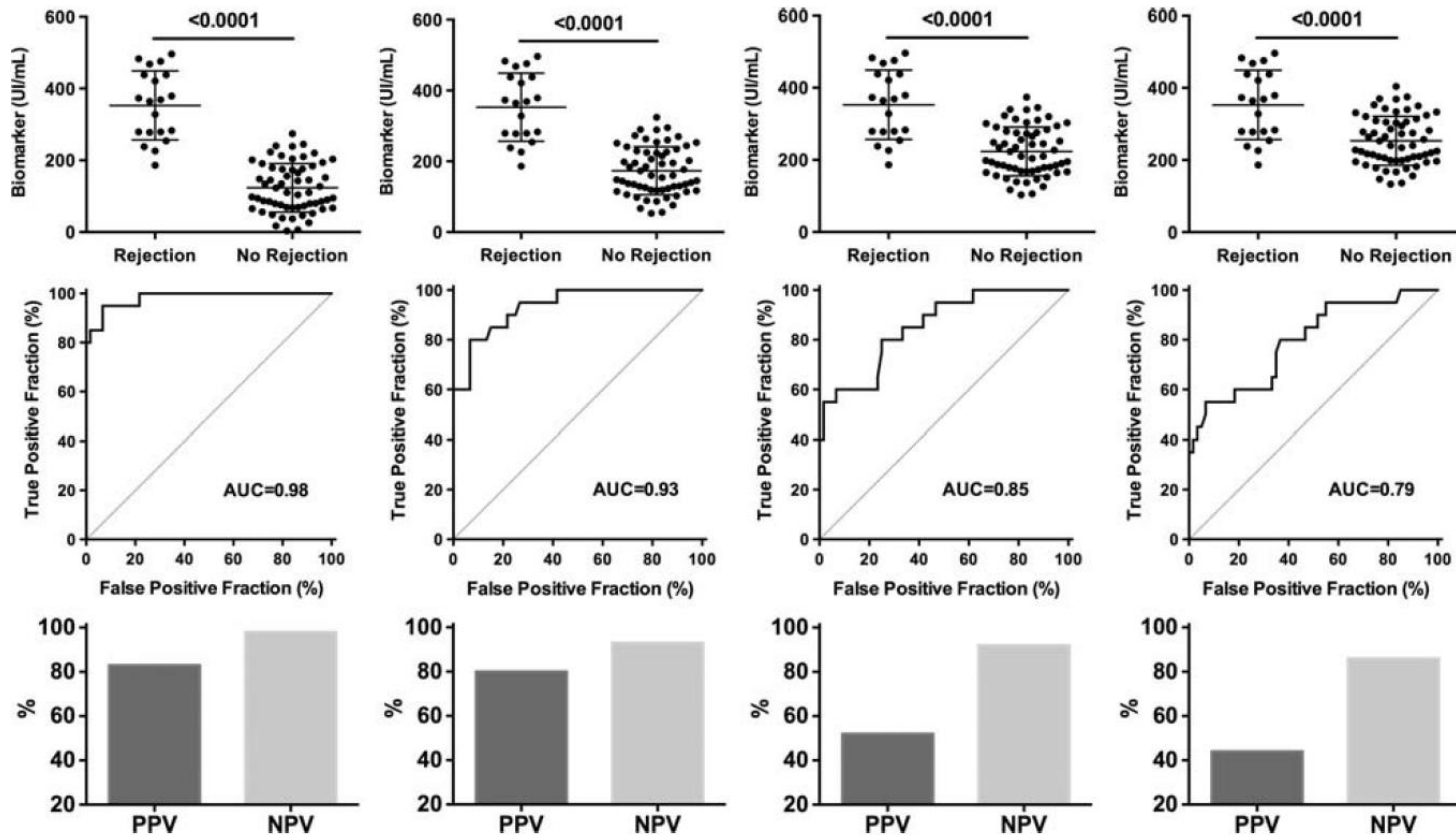


Predictive Values are Dependent on Prevalence of Disease



This figure illustrates how the prevalence of the disease can affect the predictive values of the biomarker, whereas the ROC appears similar in all conditions.

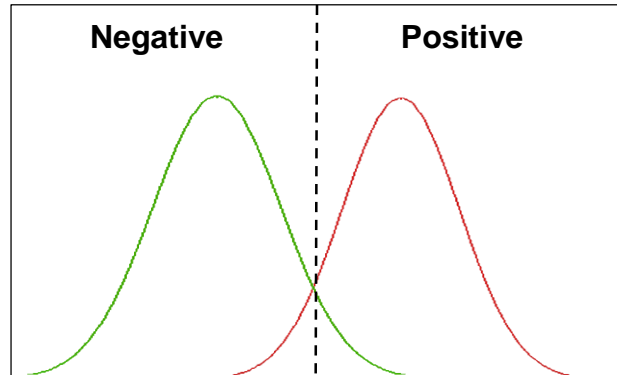
Statistical Group Differences are Not Diagnostic Accuracy



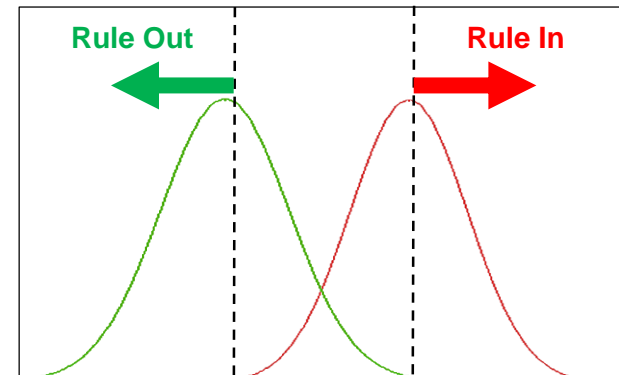
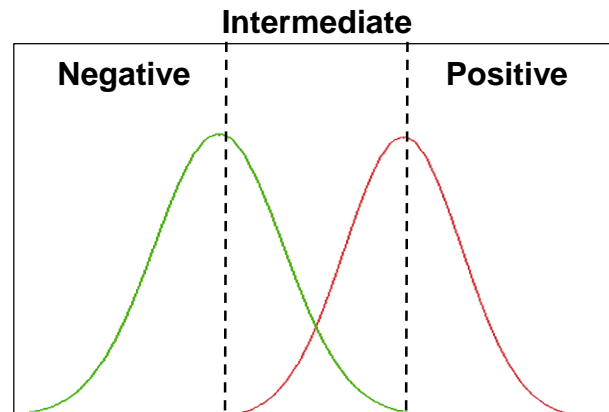
This figure illustrates 4 conditions in which a biomarker is highly significantly associated with the diagnosis of disease condition but has a highly variable diagnostic accuracy and predictive values.

Dichotomania

Single Threshold (Dichotomous)

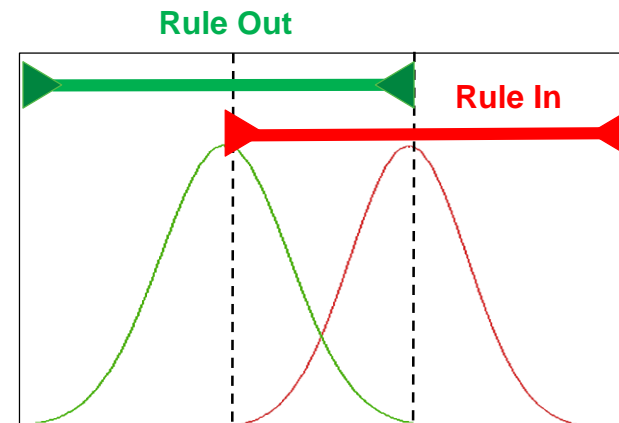


Dual Threshold



Disadvantages of dichotomous threshold

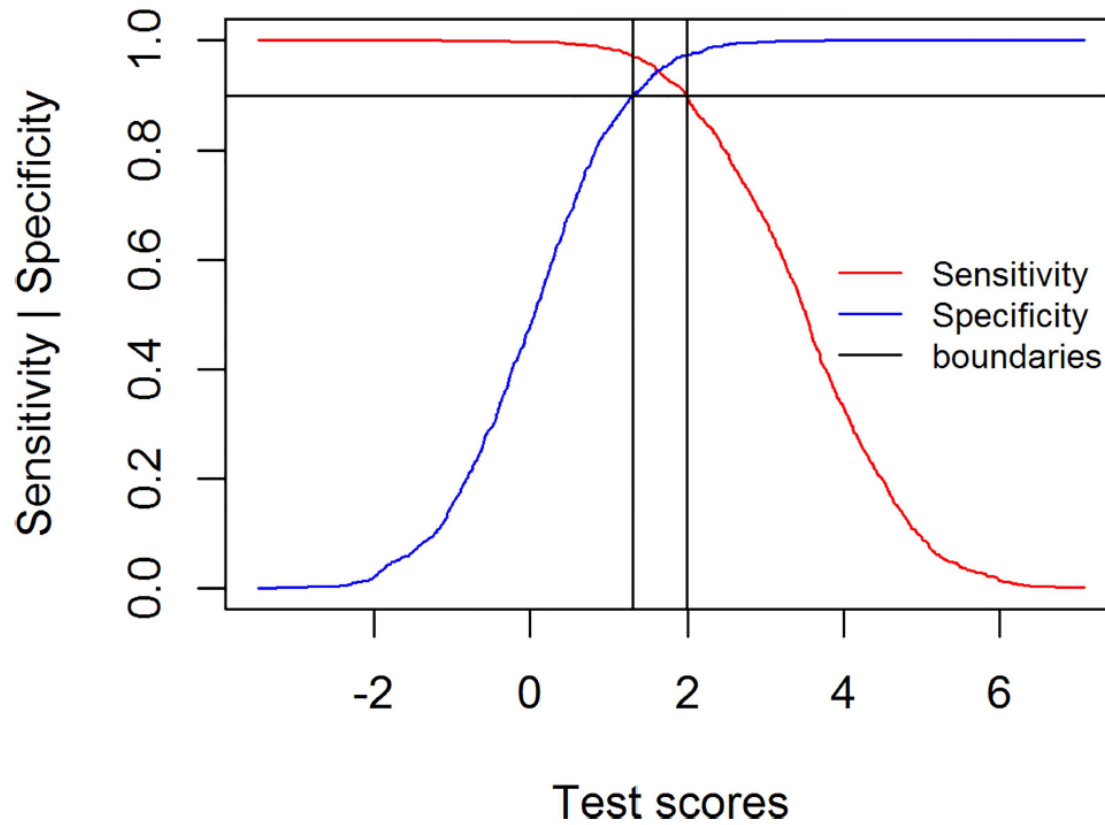
- Information loss
- Smaller difference between negative and positive groups
- Threshold significantly impacted by population distribution
- Intended use rarely represents a step function
- Less flexibility for intended use
- Practical use considers subjects at threshold differently anyway
- Critically dependent on ground truth accuracy of reference



Both single and dual threshold approaches have value but choice dependent on context of use

Two Graph (TG)-ROC to Set Thresholds

TG-ROC



Dichotomous Test Comparison

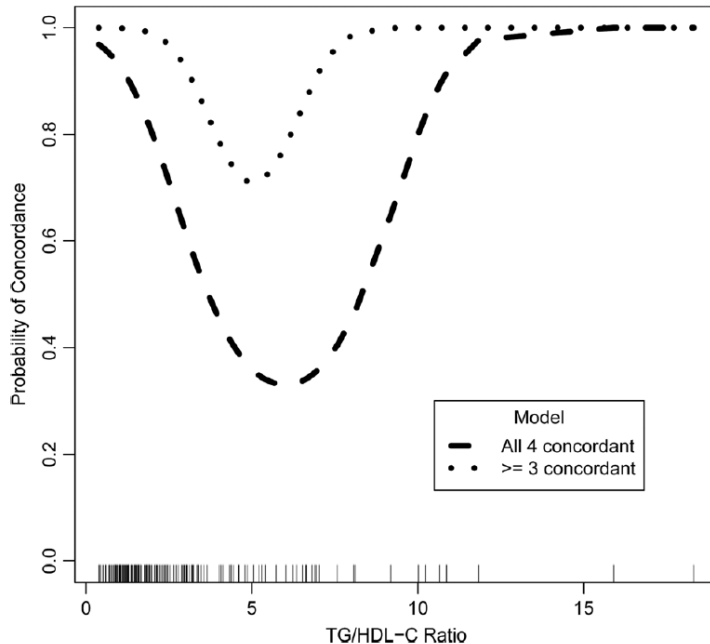


FIGURE 1. The estimated probability of concordance in the LDL phenotype = B versus not B among VAP, sGGE, IM, and NMR is plotted as a function of TG/HDL-C levels; dotted line represents the expected probability of concordance when concordance is defined as at least 3 methods in agreement based on the logistic regression equation where $\log \text{ odds of concordance} = 10.68 - 3.88 * (\text{TG}/\text{HDL-C}) + 0.39 * (\text{TG}/\text{HDL-C})^2$; dashed line represents the expected probability of concordance when concordance is defined as all 4 methods in agreement based on the logistic regression equation where $\log \text{ odds of concordance} = 4.01 - 1.57 * (\text{TG}/\text{HDL-C}) + 0.131 * (\text{TG}/\text{HDL-C})^2$; vertical dashes along x axis represent observed values of TG/HDL-C ratio.

- Extremes of dichotomous tests agree with each other a large fraction of time
- Dichotomous test comparisons are more discordant at thresholds
- Raises question of ground truth

V-plot Methodology

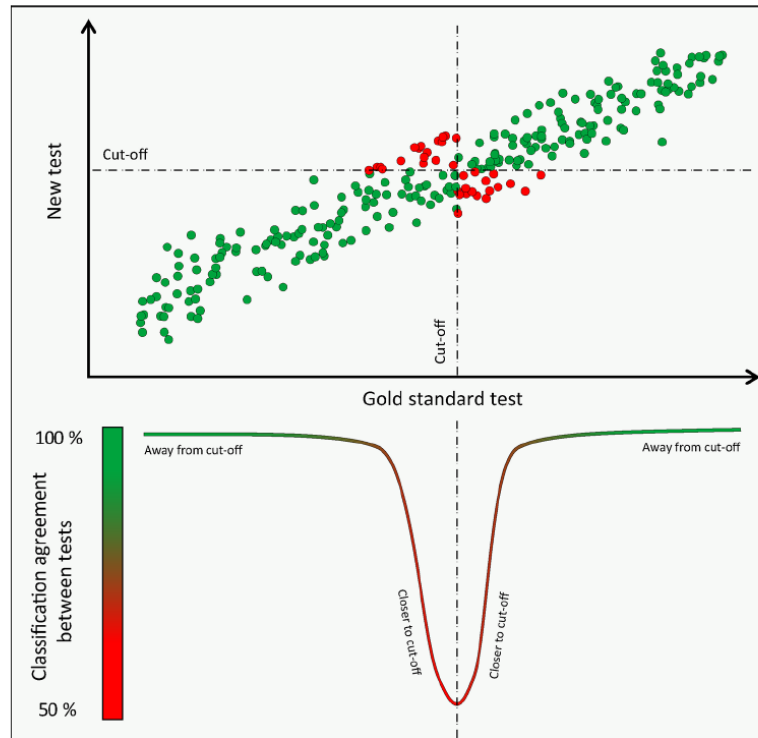


Figure 1 Disease severity and classification agreement between methods: schematic representation of the principle that classification agreement between two methods of measurement (or diagnostic accuracy if one is seen as a reference gold standard) varies across the range of disease severity. At the extremes of disease and health agreement is 100%. Close to the classification cut-off, around the intermediate range of disease severity, agreement falls, reaching a nadir close to 50%.

No single value of diagnostic accuracy can be determined in a dichotomous test comparison if the underlying sample distribution varies

Importance of Biomarker Distribution for Validation

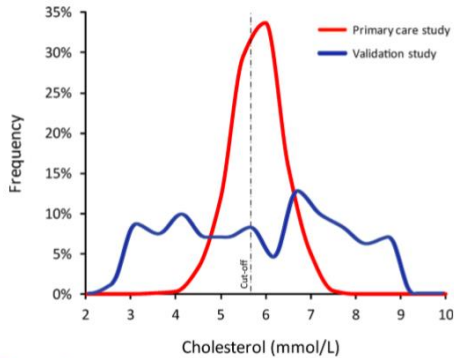
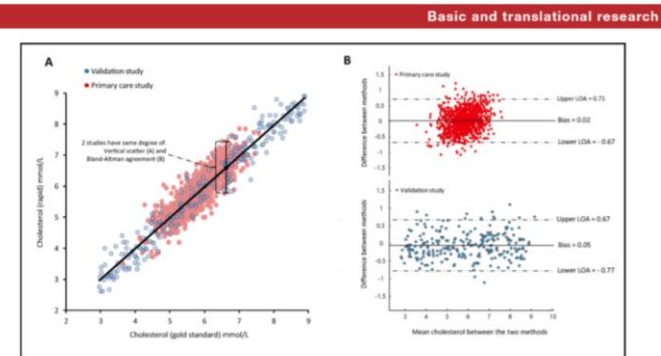
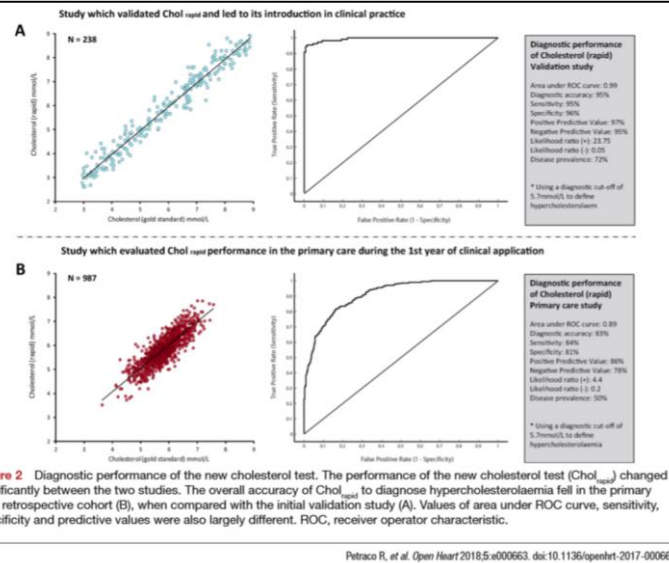


Figure 4 Histograms of cholesterol values from both studies. While the validation study included patients with a wide range of cholesterol values, the primary care cohort was formed predominantly of patients with intermediate values of cholesterol. This difference was responsible for the significant drop in $\text{Chol}_{\text{rapid}}$ accuracy reported in the primary care study.

- Generally, more discordance between two comparative tests occurs at the selected cutoff(s)
- Sample sets with distributions that differ from intended use population, therefore, will not serve as relevant validation test sets
- Rather than report overall accuracy, best to determine agreement across portions (quantiles) of the biomarker continuum



V-plot Methodology

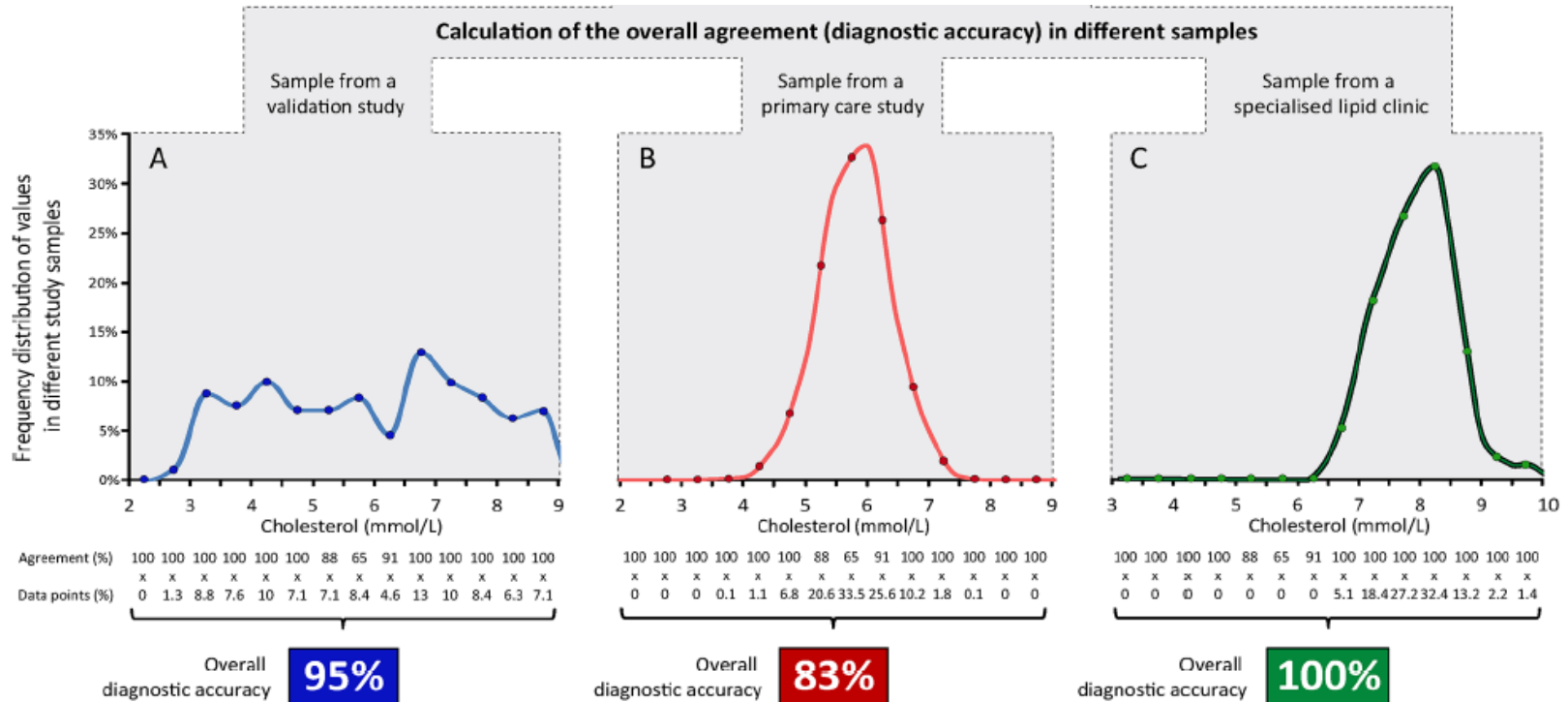


Figure 7 Calculating the overall accuracy in different samples using the V-plot. The V-plot agreement between $\text{Chol}_{\text{rapid}}$ and $\text{Chol}_{\text{gold}}$ can be derived from any study that compared the two methods (top panel). It can be used as a fingerprint of classification agreement to calculate the overall agreement between $\text{Chol}_{\text{rapid}}$ and $\text{Chol}_{\text{gold}}$ in any sample in which the distribution of cholesterol values is known (samples A, B and C).

Prevalence and Predictive Value

		Condition (as determined by "gold standard")		
		Condition positive	Condition negative	
Test outcome	Test outcome positive	True positive	False positive (Type I error)	Positive predictive value = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Test outcome positive}}$
	Test outcome negative	False negative (Type II error)	True negative	Negative predictive value = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Test outcome negative}}$
		Sensitivity = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Condition positive}}$	Specificity = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$	

The mathematical relationship between the predictive value of a biomarker, sensitivity, specificity and prevalence is defined by Bayes Theorem, which mathematically can be reduced to the following equations:

$$\text{PPV} = \frac{(\text{sensitivity})(\text{prevalence})}{(\text{sensitivity})(\text{prevalence}) + (1 - \text{specificity})(1 - \text{prevalence})}$$

$$\text{NPV} = \frac{(\text{specificity})(1 - \text{prevalence})}{(\text{specificity})(1 - \text{prevalence}) + (1 - \text{sensitivity})(\text{prevalence})}$$

Cautionary note that prevalence of intended use testing may vary from sample set tested

Critical Role of Prevalence

If the sample sizes in the positive (disease present) and the negative (disease absent) groups do not reflect the real prevalence of the disease, then the Positive and Negative Predicted Values, and Accuracy cannot be estimated and you should ignore those values.

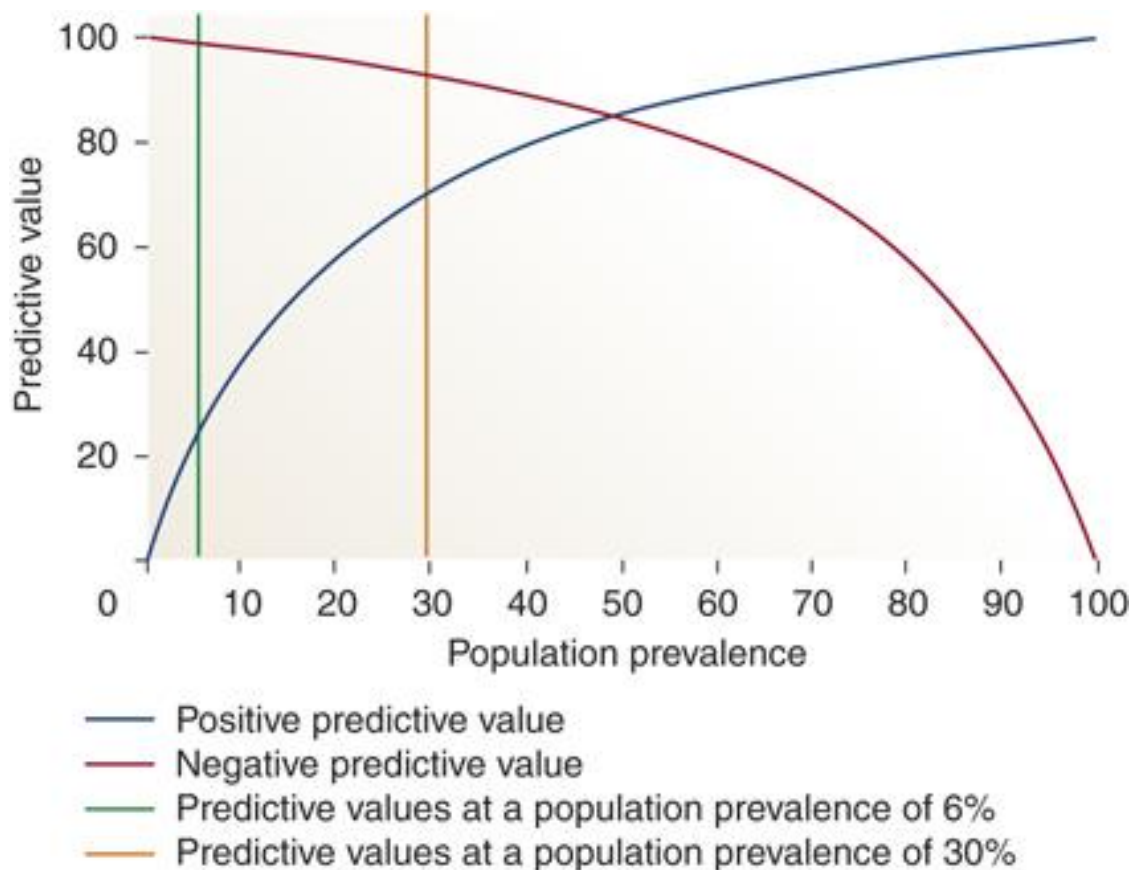
Alternatively, when the disease prevalence is known then the Positive and Negative Predictive Values can be calculated using the following formulas based on Bayes' theorem:

$$PPV = \text{Prev} \times \text{Sen} / (\text{Prev} \times \text{Sen} + (1 - \text{Prev}) \times (1 - \text{Spec}))$$

and

$$NPV = (1 - \text{Prev}) \times \text{Spec} / ((1 - \text{Prev}) \times \text{Spec} + \text{Prev} \times (1 - \text{Sen}))$$

Predictive Value: Impact of Prevalence



- Predictive value (probability that the patient actually has the disease) is typically more important to a doctor & patient than sensitivity and specificity *per se*.
 - Dependent on prevalence (“prior” probability) of disease in a population from which patient arises

Impact of Prevalence on Predictive Value

- Predictive Value is not intrinsic to the test - it depends on the prevalence of disease
- The results of a study may not apply to all situations if there are different prevalence rates between the discovery and validation studies or development and clinical practice populations
- If prevalence is very low even if sensitivity and specificity are high, test results will have high false positive rate
- Context of use determines whether PPV or NPV is critical

Disease Prevalence in the Intended Test Population	Probability of having the Disease if you have a Positive Result
0.1%	1.9%
1%	16%
10%	68%
20%	83%
50%	95%

Assumes a 95% sensitive and 95% specific test

Different Kinds of Diagnostic Tests (Context of Use)

Diagnostic

A biomarker that confirms or determines the presence of disease

Prognostic

A biomarker that predicts a clinical outcome regardless of treatment and includes element of time

Predictive

A biomarker that changes in response to treatment, and predicts a clinically relevant event or process, and could be used to identify subsets of patients who are most likely to respond to treatment

Clinical end point

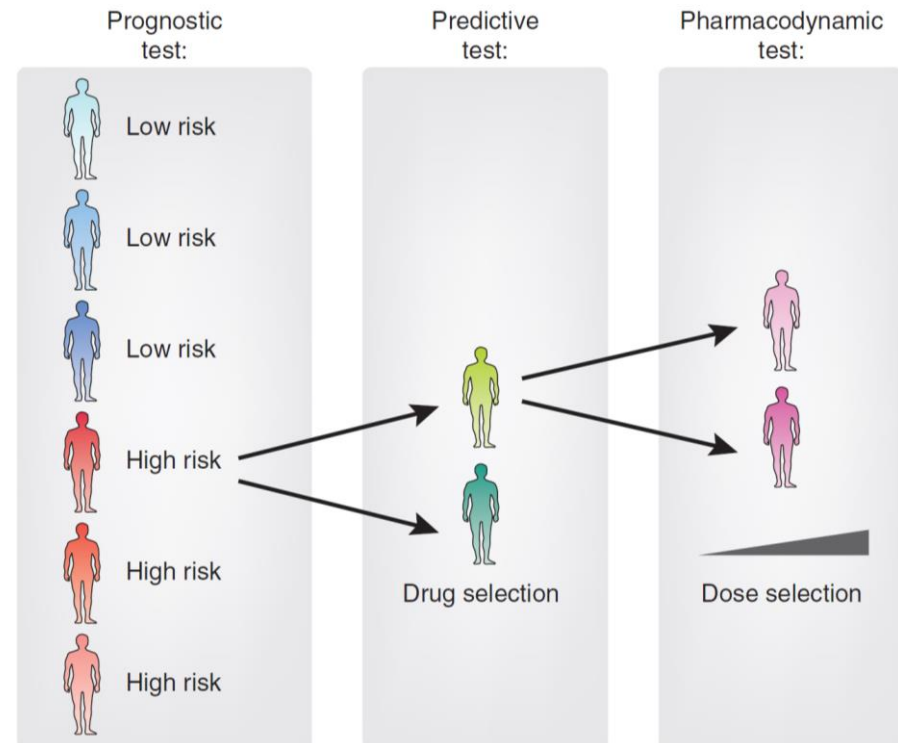
A characteristic or variable that reflects how a patient feels, functions, or survives

Surrogate end point (more likely 'proxy')

A biomarker that can substitute for a clinical end point based on biological rationale; accurately predicts a clinical end point and the effect of a given treatment on the clinical end point

Pharmacodynamic

A biomarker that provides information on drug performance

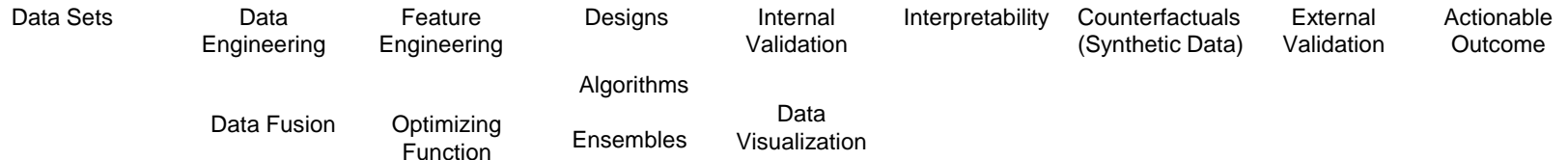


Context of Use drives Intended Use

Categories of Biomarkers for Drug Development

- **Pharmacodynamic** – Provides information on drug metabolism
- **Proof of Mechanism (PoM)** - Show that the candidate drug engages at a reliable and quantifiable level in humans, indicating a functional effect.
- **Proof of Principle (PoP)** - Show that the candidate drug results in a biological and/or clinical change associated with the disease and the mechanism of action.
- **Proof of Concept (PoC)** - Show that the candidate drug results in a clinical change on an accepted endpoint or surrogate, in patients with the disease, plus evidence of a high degree of confidence of success in phase III.
- **Predictive Biomarkers** (sometimes known as patient stratification, selection or enrichment biomarkers) – Biomarkers that can be used to pre-select patients most likely to respond to the agent or followed to determine ongoing efficacy
- **Safety Biomarkers** – Detect toxicity before symptoms appear

Steps in Machine Learning (ML) Pipeline



Start with high quality rather than opportunistic data sets

Careful to exclude bias and chance

Obtain multiple data sets for training and validation and do not mix

Consider feature merge issues and missing data across sets

Data preparation is the most time-consuming task

Data correction critical

Decide on categorical vs continuous data for each feature

Data set feature standardization is critical

Training can have several novel elements

Model selection is key

Comparative consideration of different models

Critical to use internal validation but understand that external validation required

Initial pass with simple algorithm before advancing to combined or ensembles

Graphical display of data often provides keen insights

Feature transforms to outcomes critical

Critical to use internal validation but understand that external validation required

Selected thresholds aligned with benefit - risk of actionability

Avoid dichotomous thresholds

Edit empiric data to explore synthetic data space

External validation is critical

External validation set must be similar to training data set and eventual intended use indication

- Important to start with context of use/intended use objective
- If not used properly, ML can replicate bad practices rather than improve them
- A novel combination of obvious elements may be patentable
- For healthcare and patents, ML analysis must be transparent and interpretable

Key Design Issues in Definitive Validation

- Size (and events) of Training and Validation sets
- Training and Validation sets need to be similar (e.g. prevalence, covariates, outcomes, co-morbidities, etc.)
- Study population needs to be same as intended clinical application
 - Sufficiently general; multiple institutions
- Marker well-defined in advance
 - Validation separate from Discovery
 - Locked assay (assays, analytes, model, and thresholds)
 - Same assay used to demonstrate Clinical Validity
- Pre-specified minimally acceptable performance criteria to be met
 - Describe justification
- Individual classification is critical, not group differences
- Anticipated/desirable performance drives sample size calculations

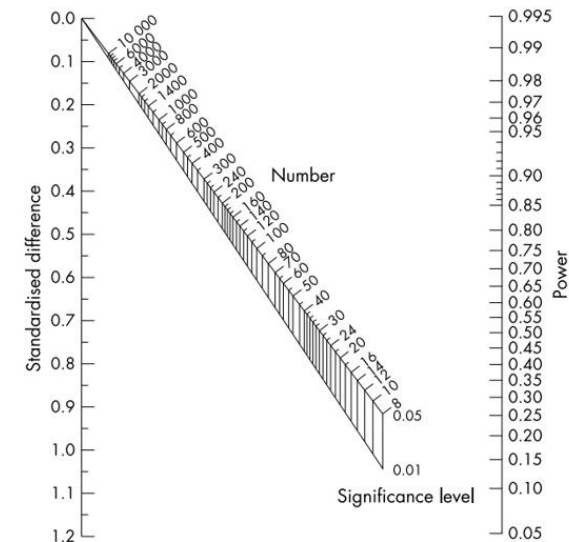


Figure 3 Nomogram for the calculation of sample size.

Critique of Biomarker Papers

- Are individual clinical validation training and test sets independent and matched with each other as well as with the intended use population?
- Is there a chance that bias or chance was introduced into sample sets being compared?
- Was the assay specifically locked (e.g. analyte(s), weighting, transform and thresholds) before validation testing?
- Was rigorous analytical validation of assay performed and published in peer-reviewed journal?
- Was a pre-specified statistical analysis plan put in place?
- What was the level of evidence collected (e.g. convenience, retrospective, prospective, single-center, multi-center, etc.)?
- Was a commonly accepted reference test used for comparison?
- Was potential of inaccuracy in reference test considered in analysis?
- Was test performance compared to and combined with conventional covariates for standard-of-care?

Not an exhaustive list but representative

Research Practices that Will Accelerate Research Findings into Clinical Practice

- Identify unmet clinical needs as primary objective
- Adoption of replication culture
- Start with high quality samples instead of samples of convenience
- Reward reproducibility studies
- More appropriate statistical methods
- Standardization of definitions and analyses
- More stringent thresholds for claiming discoveries or “successes”
- Improvement of study design standards
- Better training of scientific workforce in methods and statistical literacy

Common Missteps in Diagnostic Studies - 1

- Performance of test in Discovery set only (overfit test performance)
- Use 'normal' samples as comparator rather than differential diagnosis samples (exaggerated performance)
- Dissimilar Discovery, Validation and Clinical Use sets (inaccurate estimate of performance) or distribution of samples
- Mixture of Discovery and Validation sets (inaccurate estimate of performance, overfit; solely statistical cross-validation insufficient)
- Lack pre-specified clinical/statistical analysis plan (introduction of bias)
- Convenience or opportunistic samples (solely retrospective; not representative; inaccurate performance)
- Single center study rather than multi-center study (test robustness)
- Poorly validated analytical performance (inaccurate performance, robustness, transferability)

Common Missteps in Diagnostic Studies - 2

- Does not consider implications of pre-analytical variation of biomarker
- Samples tested with different versions of test (inaccurate performance)
- Small sample sets (likely bias and chance; lack generalizability)
- Provide clinical validity but not clinical utility (questionable reimbursement)
- Lacks attention to PPV or NPV for indication of test (actionability)
- Cost effectiveness not modeled (questionable reimbursement)
- Statistical analysis only includes ROC, or sensitivity and specificity (test performance but not patient performance)
- Lack actionable outcomes (what will clinician or patient do differently with information)
- Does not compare performance relative to single or combined routinely used tests or information (independence relative to presently used information)

Informative References

Barsanti-Innes *et al. The Oncologist* (2016). Missteps

Begley and Ioannidis *Circ Res* (2015). Reproducibility

Bossuyt *et al. BMJ* (2015). STARD

Button *et al. Nature Rev Neuroscience* (2013). Small sample sets

Gamble *et al. JAMA* (2017). Prespecified statistical analysis plans

Harrell (2015). Bad Biomarkers

Ioannidis *PLoS Med* (2005). Reproducibility

Ioannidis *JAMA* (2019). Prespecified statistical analysis plans

Leptak *et al. Science Trans Med* (2017). Biomarker qualification

Pant *et al. Frontiers in Oncol* (2014). FDA and CLIA

Parkinson *et al. Clin Cancer Res* (2014). AV, CV, and CU

Petraco *et al. Open Heart* (2017). V Plot

Ransohoff *Nature Rev Cancer* (2004). Bias and chance

Ransohoff *Nature Reviews Cancer* (2005). Bias and chance

Ransohoff *J Clin Epidem* (2007). Bias and chance

Scheerens *et al. Clin Trans Sci* (2017). Companion and Complementary tests