**ChatGPT**

# Diagnostic Test Evaluation: From Basics to Advanced Concepts

## Resident-Friendly Overview

### Basic Definitions: Sensitivity, Specificity, and Predictive Values

**Sensitivity** is the ability of a test to correctly identify patients with the disease (true positive rate) [1]. In other words, it's the probability that the test is positive given the patient **does** have the disease (P[T+ | D+]) [2]. **Specificity** is the ability of a test to correctly identify those without the disease (true negative rate) [3] – the probability the test is negative given the patient **does not** have the disease (P[T– | D–]) [2]. High sensitivity means few false negatives, and high specificity means few false positives. However, no test is perfect: sensitivity and specificity often trade off with each other depending on where we set the threshold for a "positive" result (more on that below).

**Positive Predictive Value (PPV)** is the probability a patient actually has the disease given a positive test result [4]. **Negative Predictive Value (NPV)** is the probability a patient is disease-free given a negative test [5]. These predictive values are *not inherent test characteristics* – they depend on disease prevalence in the tested population [6]. For example, even a very sensitive and specific test can have a low PPV if used in a population where the disease is rare. One illustration showed that a new test with 95% sensitivity and 95% specificity had a PPV of only 15% when disease prevalence was 1% (versus 95% PPV at 50% prevalence) [6]. In low-prevalence settings, false positives will far outnumber true positives, meaning most positive results are actually false alarms [6]. This is why **context matters** – in practice, we must consider pre-test probability (how common the disease is or how likely based on clinical factors) when interpreting predictive values.

*Take-home for residents:* Sensitivity and specificity tell us about test accuracy in idealized conditions (given disease status), while PPV/NPV tell us what a result means for *this patient* in *this setting*. High sensitivity tests are useful to **rule out** disease (few false negatives – "SnNout"), and high specificity tests help to **rule in** disease (few false positives – "SpPin"). Remember that a "great" test on paper can still mislead you on the wards if you ignore how common the disease is in your patient population.

### Likelihood Ratios and Bayes' Theorem

**Likelihood ratios (LR)** bridge the gap between test accuracy and real-world probability. An LR tells us how much a test result changes the odds of disease. The **positive likelihood ratio (LR⁺)** is the factor by which the odds of disease increase when a test is positive. It is calculated as sensitivity / (1 – specificity) [7]. A high LR⁺ (e.g. 10 or more) means a positive result is much more likely in someone with disease than without, so it significantly raises the post-test probability of disease. The **negative likelihood ratio (LR⁻)** is (1 – sensitivity) / specificity [8]. A very low LR⁻ (e.g. 0.1 or below) means a negative result makes disease much less likely.

Clinically, we apply LRs using **Bayes' theorem**: combine the pre-test odds of disease (from clinical judgment or prevalence) with the LR to get post-test odds, which can be converted to post-test probability. This can be done mathematically or with a simple tool like **Fagan's nomogram**, which lets you draw a line through the pre-test probability and LR to read off the post-test probability [9] . For example, if you estimate a patient has a 30% pre-test probability of pulmonary embolism and the D-dimer test LR$^+$ is about 3, a positive D-dimer would raise the probability to roughly 60%, whereas a negative D-dimer (LR$^-$ maybe ~0.1) would reduce the probability to ~5% – low enough to confidently exclude PE. The key point is that likelihood ratios, unlike predictive values, do *not* depend on prevalence; they are an intrinsic summary of test performance that you can apply to **any** patient by plugging in that patient's pre-test likelihood [10] . This makes LRs very powerful for evidence-based decision-making – they are the "translation tool" between research accuracy and bedside use.

*Take-home for residents:* If you can remember just a few numbers: an LR$^+$ >10 or LR$^-$ <0.1 are usually considered strong evidence to rule in or rule out disease, respectively. Moderate LRs (~2-5 for LR$^+$, 0.2-0.5 for LR$^-$) give smaller shifts. Using LRs may seem abstract, but it's essentially what you already do intuitively (e.g. a highly abnormal test greatly increases your suspicion). Tools like nomograms or smartphone apps can help you apply LRs without crunching the math every time. Learning to use LRs will deepen your understanding of diagnostic reasoning beyond "positive or negative" thinking [11] [12] .

## ROC Curves and the AUC (C-Statistic)

Many diagnostic tests yield a **continuous result** (e.g. troponin level, blood glucose, risk score). The choice of what cutoff defines "positive" vs "negative" will affect sensitivity and specificity. **Receiver Operating Characteristic (ROC) curves** provide a global picture of test performance across all possible cut-offs [13] . To construct an ROC curve, you plot the sensitivity (true positive rate) on the Y-axis against 1 – specificity (false positive rate) on the X-axis for various threshold values [13] . Each point on the curve is a sensitivity/ specificity pair corresponding to a particular cutoff. An ROC curve that climbs quickly toward the top-left corner indicates the test can achieve high sensitivity with only a small trade-off in false positives – that's a good test. Conversely, a curve near the 45° diagonal line means the test is no better than flipping a coin.

The **Area Under the Curve (AUC)** of the ROC graph (also called the **c-statistic** for *concordance statistic*) is a single summary measure of overall test accuracy [14] [15] . AUC represents the probability that a randomly chosen patient with disease will have a higher test result (more "positive") than a randomly chosen patient without disease [14] . An AUC of 1.0 means a perfect test (100% sensitive and 100% specific at some threshold), whereas an AUC of 0.5 means the test has no discriminative ability (equivalent to random guessing) [16] . In practice, we interpret AUC values in ranges: **>0.9** is considered *excellent* discrimination (high accuracy), **0.7–0.9** is *moderate/good*, **0.5–0.7** is *low/poor* [15] . For example, a troponin assay with AUC 0.95 for diagnosing myocardial infarction is outstanding, whereas a clinical risk score with AUC 0.60 for predicting stroke is fairly weak (only slightly better than chance).

One advantage of the ROC/AUC approach is that it is *independent of prevalence* (it assesses intrinsic discrimination ability) [17] . This makes AUC especially popular for comparing prediction models or diagnostic biomarkers. For instance, if one biomarker's ROC curve lies above another's, it has better overall accuracy for distinguishing disease vs no-disease [18] . ROC analysis is also used to choose an **optimal cutoff**. Common methods include selecting the point **closest to the top-left corner** (maximizing sensitivity and specificity simultaneously), or maximizing the **Youden Index (J)**, defined as Sensitivity + Specificity – 1 [19] . The Youden J identifies the threshold that gives the best trade-off between true positives and false

positives – essentially the point on the curve farthest from the diagonal line of no-discrimination [20] . In many studies, you'll read something like "Using the Youden index, an optimal cutoff of X was found, with sensitivity Y% and specificity Z%." That just means they chose the threshold that maximized (sens + spec). There isn't always a clear "elbow" in the ROC curve, so this quantitative approach helps make the choice more objective.

**C-statistic vs. AUC:** For binary outcomes (disease vs no disease), the c-statistic is literally the ROC AUC. The term "c-statistic" is often used in the context of risk prediction models (like a logistic regression model or clinical scoring system) as a measure of how well the model discriminates outcomes [21] . In survival analysis (time-to-event), a generalized version of the c-statistic (Harrell's C-index) is used to account for time ordering. But conceptually it's the same idea – higher c means better model discrimination (it can tell apart high-risk vs low-risk patients). For a new test or model, doctors often look at the AUC/c-statistic to get a sense of "overall" performance. Just keep in mind, two tests can have the same AUC yet perform differently at specific decision thresholds relevant in practice. That's why we also consider other metrics and clinical context (see advanced section).

*Take-home for residents:* ROC curves and AUC give you the "big picture" of a test's accuracy. They're especially useful when dealing with continuous tests or comparing multiple tests. If you see an AUC of 0.5, don't bother – the test is essentially useless. An AUC of 0.8 means the test is pretty good overall, but still overlap between disease and non-disease results (there will be some false pos/negatives no matter where you cut it). And remember, the **optimal cutoff depends on what you value**: in a critical illness, you might favor a threshold that gives higher sensitivity (catch everyone who's sick, at expense of some false alarms); for a costly or risky treatment, you might choose a more specific threshold. ROC analysis helps to inform that choice systematically.

## Other Useful Measures for Test Evaluation

Beyond the big four (sens, spec, PPV, NPV) and ROC/AUC, there are a few other measures and concepts worth knowing:

- **Accuracy**: simply the overall proportion of correct results (true positives + true negatives) / total tests. It's easy to calculate but can be misleading if prevalence is extreme (e.g. 95% accuracy might just reflect that almost everyone is disease-free and the test calls everyone negative). So accuracy is less informative than sensitivity/specificity.

- **Diagnostic Odds Ratio (DOR)**: a single number summary = (sens/(1–sens)) ÷ ( (1–spec)/spec ) = $LR^+$ / $LR^-$. A DOR > 1 means the test is better than chance; higher values indicate better discriminatory test performance. For an excellent test, $LR^+$ is large and $LR^-$ is tiny, so DOR will be very high. While useful for meta-analysis, DOR is not very intuitive clinically compared to LRs or predictive values.

- **F1-score** (from machine learning): the harmonic mean of precision (PPV) and sensitivity. It's more relevant in contexts like algorithm performance, especially if you want to balance false negatives and false positives and the class distribution (prevalence) is skewed. In medicine, you'll see it used occasionally in studies of AI diagnostic tools.

- **Kappa coefficient**: if you're assessing a *screening or diagnostic agreement* beyond chance (e.g. two doctors reading the same X-ray), kappa measures agreement adjusted for chance. This is more for inter-rater reliability than test accuracy per se (though related if no clear gold standard).

For now, focus on mastering sensitivity, specificity, LRs, and AUC – these will cover most scenarios. And always think about what outcome matters: ruling out life-threatening diseases? Avoiding over-treatment of false positives? The "best" test metric to pay attention to can change with the clinical question.

### Resources for Further Learning (Resident-Friendly)

- **NCBI "StatPearls" Handbook – *Diagnostic Testing Accuracy***: A free online chapter that works through examples of calculating sensitivity, specificity, predictive values, and LRs in a step-by-step manner [22] [7] . It's a great refresher with clinical examples.
- **Bland & Altman's BMJ Statistics Notes**: Short, accessible articles by two famous statisticians. For example, *"Statistics Notes: Diagnostic tests 1 – sensitivity and specificity"* (BMJ 1994) and *"Diagnostic tests 2 – predictive values"* (BMJ 1994). These classic one-page notes succinctly explain these concepts in plain English (with examples) [6] .
- **YouTube – "Sensitivity and Specificity" (Osmosis)** or **"Diagnostic Testing: ROC Curves"** (Univ. of Colorado Biostatistics) videos. Video explanations can be very helpful. There are high-yield videos on the basics of test characteristics and on understanding ROC curves in many online lectures – these often use visuals and mnemonics that make the concepts stick.
- **Blogs by Clinician-Educators**: For instance, the blog **Statistically Funny** (by Hilda Bastian) occasionally covers evidence-based medicine concepts with cartoons, or **EBM cartoons by Ibrahim** – these can make learning less dry. Another resource is the **JAMA Guide to Statistics and Methods**, which often has user-friendly explanations (if you have access via your institution).
- **Fagan's Nomogram Tool**: Try an online Fagan's nomogram calculator or app. Input a pre-test probability and LR, and it outputs post-test probability. This will reinforce how LRs work. The concept is explained in sources like the Merck Manual (with an example for MI) [9] .

By mastering these basics, you'll be well prepared for more advanced topics like prognostic models, which we'll tackle in a future session. Keep these references handy – they are high-yield for boards and for practical diagnosis!

## Advanced Deep Dive for Professors (and the Curious Resident)

### Pitfalls of Traditional Metrics and Dichotomization

While sensitivity, specificity, and related metrics are the cornerstone of diagnostic test evaluation, they come with important caveats. A major criticism is that they **force a binary decision** (disease vs no disease, "positive" vs "negative") in situations that may be inherently continuous or probabilistic [11] [23] . Dichotomizing a continuous marker (like turning a lab value into "high vs normal") throws away information. As Frank Harrell emphasizes, the act of declaring a test "positive" or "negative" often ignores the nuanced gradients of risk [23] . For example, a fasting glucose of 126 mg/dL is "diabetes" but 124 mg/dL is "not diabetes," yet biologically there's no meaningful difference – patients just on either side of the cut-off are treated as if they're in completely different states [23] . This can mislead clinicians, especially if the cut-point is somewhat arbitrary or optimized in one sample but applied universally.

Moreover, **sensitivity and specificity are not fixed properties of a test** like physical constants; they can **vary with patient mix and disease definition** [24] . If a "disease" cohort in one study is at a more advanced stage, a test might appear more sensitive (because advanced cases are easier to detect). Patient factors (age, comorbidities) can also affect test performance. Thus, sens/spec estimated in one setting may not hold in another. They also have a **backwards time logic** – clinicians think in terms of "Given the test result, what is the chance of disease?" but sensitivity/specificity answer the reverse question [2] [25] . That's why predictive values or LRs often feel more natural clinically.

Another issue: By themselves, sensitivity and specificity **do not tell us the clinical value**. A very sensitive test might detect disease but could be impractical if it has extremely low specificity (yielding too many false alarms). A notorious example is PSA screening for prostate cancer or mammography in certain age groups – highly sensitive for cancer but many false positives, leading to anxiety and further invasive tests. For this reason, some advocate moving away from the rigid "positive/negative" paradigm altogether: report the **estimated risk or probability** of disease and let the clinician/patient decide what probability is high enough to act upon [26] . This approach retains the uncertainty information. As one article put it, we should "do away with 'positive' and 'negative'" and provide continuous risk estimates [26] . For instance, instead of saying "troponin positive vs negative" at a fixed cutoff, report the troponin level and the corresponding risk of myocardial infarction – a 0.04 ng/mL might be 5% risk, 0.4 ng/mL might be 90% risk. This integrates pre-test risk and other factors better than a yes/no threshold [12] .

**Workup (Verification) Bias:** An advanced nuance is that sensitivity and specificity can be inflated if the "gold standard" is applied differentially. For example, suppose a new noninvasive test is positive and then you send those patients for the definitive invasive test (gold standard), but you don't send test-negative patients for confirmation. This **verification bias** means you're mainly confirming disease in those the test already flagged, potentially missing false negatives (making sensitivity look higher than it really is). Good study design (testing a representative sample with gold standard regardless of the initial test) or statistical correction is needed to mitigate this. Always be cautious when reading *diagnostic accuracy studies* for such biases.

In summary, traditional metrics are useful but **must be applied thoughtfully**. Relying blindly on a single "accuracy" number or a dichotomized result can lead to misdiagnosis or overdiagnosis. As a professor, you can challenge residents to think: what happens to predictive values if prevalence changes? How would choosing a different threshold change patient outcomes (false negatives vs false positives trade-off)? The goal is to encourage a nuanced understanding that diagnosis is a probability game, not a black-and-white verdict [27] [28] .

## Discrimination vs Calibration: AUC Isn't Everything

When we say a model or test has good discrimination, we mean it does well at **separating** those with and without the outcome – this is what the ROC AUC or c-statistic measures. However, discrimination is only part of the story. **Calibration** is equally important in prediction models: it is the agreement between predicted probabilities and actual outcomes. A test or model is **well-calibrated** if, among all patients it gives a 30% predicted risk, about 30% actually have the disease (and so on for other risk levels). You can have a model with an impressive AUC that is poorly calibrated – for example, consistently overestimating risk for everyone. Conversely, you might calibrate a model well (no systematic bias) but it could still have mediocre discrimination (it assigns similar probabilities to sick and healthy).

Traditional diagnostic test reporting (sens, spec, AUC) focuses on discrimination. **Why calibration matters:** Suppose two COVID mortality risk scores both have AUC ~0.85 (good discrimination). If Score A tends to over-predict death (everyone's risk is exaggerated by 2x) and Score B is perfectly calibrated, Score B is more useful clinically – you can trust the percentage it gives you. Calibration can be assessed with calibration plots or statistics (like Hosmer-Lemeshow test in logistic regression, though that test has its own issues with sample size). One combined metric is the **Brier Score**, which directly measures the mean squared error of predicted probabilities vs outcomes (0 = no event, 1 = event) [29] . The Brier score thus incorporates both discrimination and calibration – smaller Brier means the predictions are, on average, closer to reality [29] . A perfectly accurate probabilistic predictor would have Brier score 0; a trivial predictor that says "risk = prevalence for everyone" will have a Brier equal to p*(1-p); a coin-flip 50% predictor has Brier 0.25 for a binary outcome with balanced classes.

However, it's important to note that *lower Brier isn't always "clinically" better* – it's an overall performance measure. In fact, **Brier score is sensitive to prevalence**: when disease becomes very rare, a strategy of predicting "nobody has it" yields a great Brier score (because it's right most of the time) but obviously that strategy has zero clinical utility [30] . Researchers have pointed out scenarios where using Brier to choose between models can be misleading [30] [31] . For example, a more *clinical* metric like **net benefit** might favor a model that identifies a few treatable cancers (accepting some false positives) over a model that simply says "no one has cancer" (no false positives but misses all cancers). Brier might favor the latter in a low-prevalence situation, because it gets the majority of cases (the negatives) correct [30] [31] .

So, while AUC and Brier are "proper" statistical scores (with Brier being a proper scoring rule that is minimized by the true probabilities [32] ), they **do not directly inform clinical consequences**. This is where newer paradigms come in:

- **Decision-analytic metrics (Net Benefit):** Net benefit (NB) incorporates a chosen threshold probability and weighs false positives vs false negatives in clinical terms. It's like putting a "value" on identifying a true case vs the harm of a false alarm. Andrew Vickers and colleagues popularized decision curve analysis, which shows the net benefit of using a model at various probability cut-offs [33] . In their analysis, they found net benefit reliably identified the better test/model in scenarios where Brier score rankings were inconsistent [33] . In short, if you care about patient outcomes, you may want to use decision curves to decide whether a model is worth using in practice (Does it help us treat more of the right patients without too many unnecessary interventions?).

- **Clinical Utility and Threshold Metrics:** This includes things like **sensitivity at a given specificity** (or vice versa) to reflect a clinically relevant point. For instance, a screening test might be evaluated by "specificity when sensitivity is set at 95%" – if the specificity is still decent, it's a good sign. **Partial AUC** in a high-sensitivity region is another method used in some research.

- **Precision-Recall Curve (PR AUC):** When we deal with very imbalanced outcomes (e.g. a rare disease), ROC curves can paint an overly optimistic picture because the false positive rate on the X-axis doesn't penalize having false positives among a huge pool of negatives. In such cases, a precision (PPV) vs recall (sensitivity) curve can be more informative. The area under the PR curve focuses on how well the test does for the positive class specifically. A classic example is screening for cancer in a low-prevalence population: a modest AUC might correspond to a very low PPV. PR curve analysis would highlight that issue.

**Bottom line:** As a teacher or advanced learner, emphasize that **no single metric tells the whole story** [21] . AUC (c-statistic) is great for discrimination ranking, but one should also ask: is the model well-calibrated? Does it improve decision-making enough to be worth it? Two models with the same AUC could have different clinical implications if one systematically overestimates risk. And an improvement in AUC might appear small yet still be worthwhile if it correctly reclassifies some patients at a critical threshold (though beware of overhyped reclassification metrics – see below). In modern practice, we aim to use *both* discrimination and calibration measures, and increasingly decision-analytic measures, to get a comprehensive picture of a test's performance [34] [35] .

## Newer Metrics and Methods: NRI, IDI, and Reclassification

To address the perceived "insensitivity" of AUC to improvements, statisticians proposed metrics like **Net Reclassification Improvement (NRI)** and **Integrated Discrimination Improvement (IDI)** in 2008 [34] . The motivation was that adding a new biomarker to a risk model often yields tiny increases in AUC, even if the biomarker has a strong association with outcome. NRI attempts to quantify how many individuals are reclassified to more correct risk categories with the new model. For example, if adding a biomarker moves some patients from "intermediate risk" to "high risk" and they indeed had events, that's a useful reclassification. IDI looks at the change in the model's discrimination slope (difference in mean predicted risk between events and non-events).

While intuitively appealing, these metrics have faced **significant criticism** [36] . NRI in particular can be **misleadingly large** even when a new model is not truly better, especially if risk categories are arbitrarily chosen or if the model is not well-calibrated [35] . In fact, NRI *encourages* making extreme predictions – a poorly calibrated model that confidently pushes people into 0% or 100% predicted risk could score a high NRI but be clinically nonsensical. Scholars (like Kerr et al. 2014) have published critiques debunking NRI's purported advantages [37] . They showed that NRI can inflate Type I error (finding "improvement" where none exists) [35] . IDI is somewhat more straightforward (it's basically the difference in *Yates' slope* or discrimination slope), but it also lacks a direct clinical interpretation.

The consensus in the statistical community is that **improvements in AUC, though often small, are still meaningful when considered with other evidence** [34] . Instead of relying on NRI/IDI alone, one should check calibration, confidence intervals, and conduct decision analysis. If a new marker raises AUC from 0.75 to 0.78, NRI might paint it as a big jump in reclassification, but what does that mean for patients? Perhaps only a small subset changed management. So, use NRI/IDI judiciously; they can be reported for completeness, but they are not "magic" metrics.

## Putting It All Together – Teaching Points

In advanced discussions, I highlight a few **key teaching points**:

- **Contextualize Test Performance:** Always relate test characteristics back to clinical context. A "high accuracy" test (say 90% sensitivity, 90% specificity) can still lead to tons of false positives if used for mass screening in a low prevalence scenario [6] . Conversely, even a modest AUC model might be very useful if it identifies a small subgroup at very high risk who benefit from intervention.

- **Individualized Decisions vs Population Metrics:** Metrics like AUC are population-level – they don't directly tell a doctor what to do for an individual patient. David J. Hand famously pointed out that

when we integrate over the whole ROC curve, we implicitly weight false positives/negatives in a way that may not reflect any specific clinical situation [28] . The relative importance of false positive vs false negative must come from clinical judgment or patient preference, not from the data alone [28] . Thus, the "optimal" threshold from ROC (Youden or closest-to-(0,1)) might not align with what's optimal for a given patient or healthcare system. Encourage learners to think in terms of **expected value** or **utility** – sometimes a lower Youden index cutoff is chosen because missing a case is so costly that we tolerate more false positives (e.g. threshold for starting antibiotics in suspected meningitis).

- **Prospective Validation:** It's one thing to derive a model/test performance in a study, another to see it perform in practice. Emphasize the need for external validation and assessment of these metrics in different settings. Many initial AUCs shrink upon validation. Calibration often deteriorates over time if underlying prevalence or treatment patterns change. This is why dynamic monitoring of model performance (and possibly updating models) is a hot area now.

- **Advanced methods:** For completeness, you can mention emerging techniques like machine learning algorithms that use cross-validation to maximize AUC or other metrics, or techniques for *calibration* like Platt scaling and isotonic regression used in ML to improve probability estimates. But those might be beyond the scope of this session.

In conclusion, modern diagnostic test evaluation is moving toward a more holistic approach: combining **discrimination, calibration, and clinical utility**. As a professor, you can use real examples (some below) to illustrate these concepts – e.g., show how a risk score with c-stat 0.65 can still be the standard of care (CHA$_2$DS$_2$-VASc for atrial fibrillation stroke risk) because it's simple and "good enough," whereas a fancy model with c-stat 0.80 might not be adopted if it's too cumbersome or doesn't reclassify patients in a useful way. The goal for learners is to appreciate the strengths and limitations of each metric and to always ask **"How will this test or model help me improve patient outcomes?"** rather than just "Is its AUC high?".

## Five Notable Studies Illustrating C-Statistics in Action

To ground these concepts, here are five high-impact studies (across different areas of internal medicine) where **c-statistics (AUC)** were used to evaluate and compare diagnostic or prognostic tools:

1. **Adding Coronary Calcium to Risk Scores for Heart Disease – JAMA 2023:** Khan *et al.* studied whether incorporating a coronary artery calcium (CAC) score and a polygenic risk score could improve prediction of coronary events beyond the standard Atherosclerotic Cardiovascular Disease (ASCVD) risk factors [38] [39] . They found that adding CAC markedly improved discrimination (increase in C-statistic by ~0.09, p<0.001), whereas adding the genetic risk score did not significantly change the C-statistic [39] . For example, the base risk model had an AUC around the mid-0.7s, and including CAC pushed it into the low-0.8s. They also reported **Net Reclassification Improvement (NRI)** and **Integrated Discrimination Improvement (IDI)** – the continuous NRI was significant for CAC addition, underscoring that many patients were correctly reclassified into higher or lower risk categories [40] [39] . This study shows how c-statistics are used to quantify incremental benefit of a new test: CAC (a CT scan measure) provided a meaningful boost in predictive accuracy, whereas the genetic score (PRS) did not, in this cohort. It reinforces that *effect size matters* – a risk factor must have a fairly large impact to noticeably lift the AUC [41] .

2. **qSOFA vs SIRS for Predicting Sepsis Outcomes – JAMA 2017:** In the validation of the Sepsis-3 criteria, a multicenter study compared the quick Sequential Organ Failure Assessment (qSOFA) score with older criteria like SIRS for prognosticating in-hospital mortality among infected patients. The **AUROC** for qSOFA was 0.80, significantly higher than the AUROC ~0.65 for SIRS criteria [42] . This indicated qSOFA had better discrimination for death (80% chance that a randomly chosen non-survivor had a higher qSOFA than a survivor) [42] . Based on these and similar findings, qSOFA (which uses only 3 clinical measures) was proposed as a quick bedside tool to identify high-risk infections. Interestingly, the full SOFA score had an AUROC ~0.77, slightly below qSOFA in that ED cohort [43] . This example shows how we use AUC to compare tools: even though SIRS was very sensitive (caught >90% of sepsis deaths) [44] , its discrimination was low (it also flagged a lot of survivors as positive, hence the low specificity and AUC). qSOFA achieved a better balance. It's a real-life demonstration that **"more sensitivity" doesn't always mean a better test** – we seek an optimal mix, quantified by AUC in such studies. (Of course, one must also consider that higher sensitivity might be desired in some settings – again, metrics vs context!)

3. **Kidney Failure Risk Equation (KFRE) Validation – JAMA 2016:** The KFRE is a prediction equation (using GFR, albuminuria, age, etc.) to estimate risk of end-stage renal disease. A meta-analysis of cohorts worldwide showed the 4-variable KFRE had a **c-statistic ~0.90 at 2 years** for predicting kidney failure [45] – which is exceptionally high discrimination for a clinical prediction tool. Even at 5 years, c-statistics ~0.88 were reported [45] . This means the model is very adept at ranking patients by risk (90% chance a randomly selected patient who developed ESRD had a higher predicted risk than one who didn't) [45] . It also validated well across age groups and regions. Notably, they also assessed calibration: the KFRE was well-calibrated in many settings after minor recalibration of baseline hazards [46]  [47] . The KFRE's high C-index and ease of use have made it a widely adopted tool in nephrology (e.g. to decide who to refer for transplant evaluation). This study is a great example for residents of a tool with both high discrimination and good calibration, highlighting what "excellent AUC" looks like in practice – and how it translates to useful risk stratification (e.g. identifying patients who have >40% 5-year risk of ESRD, where interventions like transplant or intensive management might be warranted).

4. **AI for Diabetic Retinopathy Screening – JAMA 2016:** Gulshan *et al.* developed a deep learning algorithm to detect diabetic retinopathy from retinal photos. The performance was evaluated against expert ophthalmologists. The results showed **AUCs in the 0.93–0.99 range** for detecting referable diabetic retinopathy across multiple validation sets [48] . At operating points chosen for high sensitivity, the algorithm achieved ~90% sensitivity and ~91% specificity in the primary validation (AUC 0.936) [48] ; for more severe disease, sensitivity was 100% (with high specificity) [49] . These extremely high AUCs illustrate how AI can approach expert-level diagnosis – an AUC of 0.94 means that almost every time a diseased eye and a healthy eye are shown, the model ranks the diseased eye more likely (94% of the time). Importantly, the study also highlights using **sensitivity/ specificity at certain cut-offs** (e.g. if the algorithm is set to 97% sensitivity, what is the corresponding specificity?) to meet clinical requirements [49] . For screening, they might favor near 100% sensitivity to not miss cases, and accept slightly lower specificity. This paper is a prototypical use of AUC for an AI diagnostic – demonstrating high discrimination – along with reporting PPV/NPV at given prevalence and the chosen operating threshold. It's also a reminder that when disease prevalence is low (they note only 3% had referable retinopathy in one setting [50] ), even a 90% sensitive test might have an appreciable false positive rate; hence, the need to tune and possibly

incorporate sequential human review. Nonetheless, a well-validated AUC ~0.95 gave confidence that the tool was ready for further real-world testing.

5. **Stroke Risk Scores in Atrial Fibrillation – multiple studies:** In cardiology, the $CHA_2DS_2$-VASc score is used to predict stroke risk in patients with atrial fibrillation and guide anticoagulation decisions. Despite its ubiquitous use, the score's discriminative ability is only **moderate, with c-statistics around 0.6–0.7** in various validations [51] . For example, one analysis noted $CHA_2DS_2$-VASc had an AUC ~0.64 for stroke, compared to ~0.58 for the older $CHADS_2$ score [52] . These values indicate *limited* predictive value – the score classifies patients only slightly better than chance. In fact, risk stratification in AF has a known ceiling with clinical factors alone. Yet, $CHA_2DS_2$-VASc is still very useful because it is simple and identifies clear risk categories for action (e.g. score 0 = low risk, $\geq 2$ = high risk). This scenario underscores a crucial point: **a model doesn't need a sky-high AUC to be clinically valuable**. If the decision threshold is low (essentially treat almost everyone above minimal risk because treatment (anticoagulation) is quite effective), then a rough stratification is enough. The score's **calibration** also matters – it does roughly calibrate to annual stroke rates (e.g. score 3 correlates to ~3% per year stroke risk in untreated patients). Moreover, the **clinical outcome and context** define utility: we tolerate a low AUC because the consequence of treating some moderate-risk patients unnecessarily (giving anticoagulation to someone who may not have stroked) is acceptable compared to not treating a truly high-risk patient. This example is a favorite teaching case to show that numerical metrics need interpretation: 0.6 AUC is usually "poor," but in context, $CHA_2DS_2$-VASc was an improvement over guessing and provided a basis for guidelines [53] [52] . It also challenges us to find better markers (e.g. adding biomarkers or imaging to improve prediction – an area of active research, hopefully to raise that AUC closer to 0.8 in the future).

Each of these studies highlights different facets of diagnostic/prognostic test evaluation – from assessing incremental improvement (Study 1), to comparing alternative criteria (Study 2), showing high discrimination success (Study 3 & 4), and recognizing limitations in current tools (Study 5). Reviewing such papers with residents can solidify understanding: they see how sensitivity, specificity, and AUC are reported and used in decision-making. Always encourage a critical eye: look at confidence intervals of AUC, look at calibration (if reported), and consider what these statistics mean for patient care. Armed with both conceptual knowledge and real-world examples, one can expertly navigate the modern landscape of diagnostic test evaluation [34] [27] .

[1] [3] [4] [5] [6] Understanding and using sensitivity, specificity and predictive values - PMC
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2636062/

[2] [11] [12] [23] [24] [25] [26] [27] [28] [19] Diagnosis – Biostatistics for Biomedical Research
https://hbiostat.org/bbr/dx

[7] [8] [10] [22] Diagnostic Testing Accuracy: Sensitivity, Specificity, Predictive Values and Likelihood Ratios - StatPearls - NCBI Bookshelf
https://www.ncbi.nlm.nih.gov/books/NBK557491/

[9] Image:Fagan Nomogram Used to Determine Need to Test-MSD Manual Professional Edition
https://www.msdmanuals.com/professional/multimedia/image/fagan-nomogram-used-to-determine-need-to-test

[13] [14] [15] [16] [17] [18] [19] [20] Understanding diagnostic tests 3: receiver operating characteristic curves
https://www.est.ufmg.br/~enricoc/pdf/medicina/artigos/roc_class.pdf

21  29  30  31  32  33  34  35  36  37  The Brier score does not evaluate the clinical utility of diagnostic tests or prediction models | Diagnostic and Prognostic Research | Full Text

https://diagnprognres.biomedcentral.com/articles/10.1186/s41512-017-0020-3

38  39  40  41  Coronary Artery Calcium Score and Polygenic Risk Score for the Prediction of Coronary Heart Disease Events | Cardiology | JAMA | JAMA Network

https://jamanetwork.com/journals/jama/fullarticle/2805138

42  43  44  Prognostic Accuracy of Sepsis-3 Criteria for In-Hospital Mortality Among Patients With Suspected Infection Presenting to the Emergency Department | Infectious Diseases | JAMA | JAMA Network

https://jamanetwork.com/journals/jama/fullarticle/2598268

45  46  47  Multinational Assessment of Accuracy of Equations for Predicting Risk of Kidney Failure: A Meta-analysis | Nephrology | JAMA | JAMA Network

https://jamanetwork.com/journals/jama/fullarticle/2481005

48  49  50  Development and Validation of a Deep Learning System for Diabetic Retinopathy and Related Eye Diseases Using Retinal Images From Multiethnic Populations With Diabetes | Diabetic Retinopathy | JAMA | JAMA Network

https://jamanetwork.com/journals/jama/fullarticle/2665775

51  Stroke risk in atrial fibrillation. Is there a missing piece?

https://www.journalofcmr.com/article/S1097-6647(24)01754-X/fulltext

52  Antithrombotic Therapy in Patients With Atrial Fibrillation …

https://www.ahajournals.org/doi/pdf/10.1161/CIRCINTERVENTIONS.111.965186

53  Evaluation of Risk Stratification Schemes for Ischaemic Stroke and …

https://www.acc.org/latest-in-cardiology/journal-scans/2012/01/25/16/35/evaluation-of-risk-stratification-schemes-for-ischaemic-stroke