## 1 Study Design

▶ **Overview**

Study design is the foundational framework that determines how validation research is conducted. The choice between retrospective and prospective designs, single-center versus multi-center approaches, and internal versus external validation significantly impacts the generalizability and clinical applicability of AI models in medical imaging.

**Key Design Considerations:**

- **Retrospective advantages:** Rapid completion, cost-effective, large sample sizes readily available
- **Retrospective limitations:** Selection bias, missing data, variable

## 2 Ground Truth Establishment

▶ **Definition and Importance**

Ground truth represents the reference standard against which AI model predictions are compared. It is the "correct answer" that defines what the model should predict. The quality and reliability of ground truth directly determine the validity of all validation metrics.

| Method | Advantages | Limitations | Best Use Cases |
|--------|-----------|-------------|----------------|
| **Histopathology** | • Definitive diagnosis | • Invasive procedure | Cancer detection |

## 3 Reader Studies

▶ **Purpose and Design**

Reader studies compare the diagnostic performance of radiologists with and without AI assistance, providing evidence of the AI system's clinical utility. These studies simulate real-world clinical scenarios and assess whether AI improves diagnostic accuracy, efficiency, and reader confidence.

**Critical Design Elements:**

- **Sample size:** Use power analysis to determine adequate number of cases (typically 100-500 cases minimum)
- **Reader selection:** Include 3-6 radiologists with varying experience levels (junior, senior, subspecialty)

## 4 Statistical Analysis

▶ **Performance Metrics Overview**

Statistical analysis quantifies AI model performance using standardized metrics that enable comparison across studies and clinical contexts. Selecting appropriate metrics depends on the clinical task, dataset characteristics, and intended use case.

**Classification Metrics:**

- **Sensitivity (Recall):** TP / (TP + FN) - Proportion of actual positives correctly identified
- **Specificity:** TN / (TN + FP) - Proportion of actual negatives correctly identified
- **PPV (Precision):** TP / (TP + FP) - Proportion of positive predictions that are correct
- **NPV:** TN / (TN + FN) - Proportion of negative predictions that are correct
- **ROC-AUC:** Area under receiver operating characteristic curve (0.5 = chance, 1.0 = perfect)

**Segmentation Metrics:**

- **Dice Score (F1):** $2 \times (|A \cap B|)/(|A|+|B|)$ - Measures overlap between predicted and ground truth regions (0-1, higher better)

ly Protocol

aders (3-5+)
rience levels

ordering
AI results

pare:
s AI-assisted

Guidelines

stic accuracy
ction models
ansparency

us, biopsy

intervals and

TIONS

imaging protocols

- **Prospective advantages:** Standardized protocols, complete data collection, reduced bias
- **Prospective limitations:** Time-consuming, expensive, limited sample size
- **Multi-center validation:** Essential for demonstrating generalizability across different clinical settings, patient populations, and imaging equipment

**Clinical Example:**

A lung nodule detection AI trained at a single academic center achieved 95% sensitivity on internal validation. However, when tested at three community hospitals with different CT scanners and imaging protocols, sensitivity dropped to 78%, revealing the model's limited generalizability. Multi-center validation would have identified this

| Method | Advantages | Limitations | Best Use Cases |
|---|---|---|---|
| | • Objective standard • High accuracy | • Not always available • Sampling error possible | Tissue characterization Tumor classification |
| **Clinical Outcomes** | • Clinically relevant • Objective endpoints • Real-world evidence | • Long follow-up required • Loss to follow-up • Confounding factors | Prognosis prediction Risk stratification Treatment response |
| **Expert Consensus** | • Widely applicable • Feasible for large datasets • Non-invasive | • Inter-reader variability • Subjective interpretation • Potential for systematic bias | Image interpretation Lesion detection Classification tasks |

**Expert Consensus Best Practices:**

- **Multiple readers:** Use at least 2-3 independent expert radiologists to reduce individual bias
- **Blinding:** Readers should be blinded to clinical information and other readers' interpretations
- **Adjudication:** Establish clear protocols for resolving disagreements (third reader, discussion, majority vote)
- **Experience level:** Include readers with ≥5 years of subspecialty experience
- **Inter-rater reliability:** Report Cohen's kappa or intraclass correlation coefficient (ICC)

**Common Pitfall: Circular Reasoning**

Avoid using AI-assisted readings as ground truth when validating AI models. This creates circular reasoning and inflates performance metrics. Ground truth must be established independently of the AI system being validated.

- **Randomization:** Randomize case order between phases to prevent recall bias
- **Washout period:** Implement 4-8 week interval between reading sessions to minimize memory effects
- **Blinding:** Readers should be blinded to ground truth and previous interpretations
- **Data collection:** Record diagnosis, confidence level (1-5 scale), and reading time

| Design Type | Description | Applications |
|---|---|---|
| **Standalone AI** | AI operates independently without radiologist oversight | Screening programs, Triage systems, Worklist prioritization |
| **AI-Assisted (Concurrent)** | AI provides real-time suggestions during radiologist reading | Diagnostic reading, Lesion detection, Quality assurance |
| **AI as Second Reader** | Radiologist reads first, then reviews AI output | Double reading, Discrepancy detection, Teaching/training |

**Example: Mammography CAD Reader Study**

A reader study with 6 radiologists (2 breast fellowship-trained, 2 senior general, 2 junior general) evaluated 240 mammograms (120 cancer, 120 normal). Unassisted sensitivity ranged from 76-84%. With AI assistance, mean sensitivity improved to 88% (p=0.002), with greater improvement in junior radiologists (+15%) versus fellowship-trained (+6%). Reading time decreased by 23% with AI assistance.

**Automation Bias Warning**

Radiologists may over-rely on AI suggestions, potentially missing errors or accepting incorrect AI outputs without critical evaluation. Studies must assess both improvements AND potential negative

- **IoU (Jaccard Index):** |A∩B|/|A∪B| - Ratio of intersection to union (0-1, higher better)
- **Hausdorff Distance:** Maximum distance between boundary points - Measures worst-case boundary error (lower better)
- **Average Surface Distance:** Mean distance between surfaces - Overall boundary accuracy (lower better)

| Statistical Test | Use Case | Example |
|---|---|---|
| **McNemar's Test** | Compare paired binary outcomes (e.g., AI vs radiologist on same cases) | Test if AI and radiologist have different sensitivity on same 200 cases |
| **DeLong Test** | Compare two ROC curves (AUCs) | Compare AUC of two AI models: Model A (0.89) vs Model B (0.85) |
| **Bootstrap Method** | Calculate confidence intervals for any metric | 95% CI for Dice score: 0.82 (0.79-0.85) based on 1000 bootstrap samples |
| **GEE** | Reader studies with multiple readers and cases | Account for correlation when 5 readers evaluate 100 cases twice |

**Statistical Reporting Example:**

"The AI model achieved an AUC of 0.92 (95% CI: 0.88-0.95) compared to radiologist AUC of 0.84 (95% CI: 0.79-

issue before clinical
deployment.

### Best Practice Recommendation:

Aim for validation
across at least 3-5
independent centers
with diverse patient
demographics,
equipment vendors,
and imaging
protocols. Include
both academic and
community practice
settings to ensure
real-world
applicability.

effects of AI assistance, including false positive rate
changes and automation bias indicators.

0.88), p=0.003 by DeLong test.
Sensitivity improved from 78% (95% CI:
72-84%) to 88% (95% CI: 83-92%) with
AI assistance, p=0.008 by McNemar's
test. Inter-reader agreement was
substantial (ICC=0.78, 95% CI: 0.71-
0.84)."

### Multiple Comparisons Problem

When performing multiple statistical
tests, apply correction methods (e.g.,
Bonferroni, false discovery rate) to
control for Type I error inflation. If
testing 10 hypotheses at $\alpha=0.05$, expect
0.5 false positives by chance alone.