# ACCURACY IS NOT INTELLIGENCE

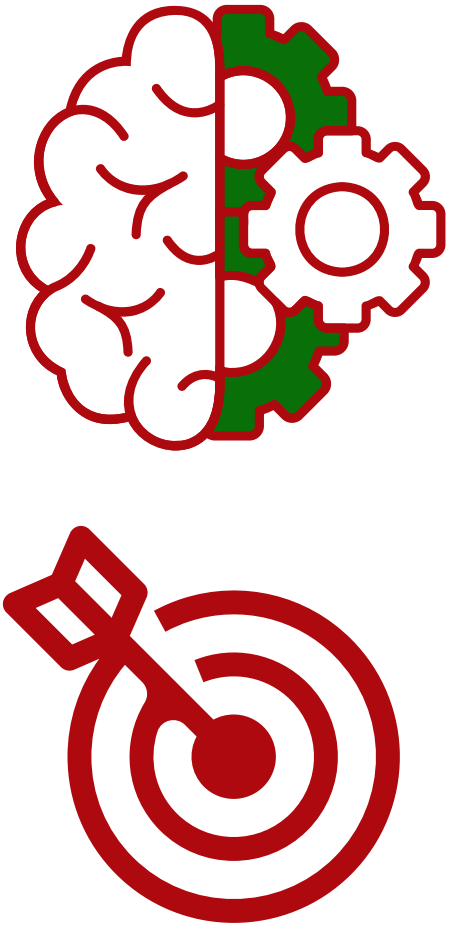Accuracy as an Upper Bound, Not an Indicator

**C O O R D I N A T O R S**

**FACULTY CO-ORDINATOR**
Dr. BHARGAVI PEDDIREDDY
(ASSOCIATE PROFESSOR)

**STUDENT CO-ORDINATORS**
DHADI SAI PRANEETH REDDY (1602-23-733-038)
SANDEEP GUNDU (1602-23-733-043)

# WHEN ACCURACY MISLEADS INTELLIGENCE: ACCURACY AS AN UPPER BOUND, NOT AN INDICATOR: A FAILURE-DRIVEN REINTERPRETATION OF INTELLIGENCE EVALUATION IN MODERN AI SYSTEMS



For much of modern machine learning, accuracy has served as the dominant indicator of progress. Models that achieved higher accuracy were considered more capable, more intelligent, and more suitable for deployment. This assumption was both pragmatic and productive: it enabled rapid benchmarking, standardized comparison, and large-scale optimization. However, as contemporary AI systems approach saturation on many established benchmarks, accuracy has begun to obscure more than it reveals.

This article advances a central claim: accuracy is not a measure of intelligence, but an upper bound on observable performance under narrowly defined conditions. As model capacity increases, accuracy increasingly conceals structural weaknesses related to robustness, reasoning, and failure behavior. Intelligence, we argue, must be evaluated through failure-aware metrics that expose how systems behave when assumptions break, not merely how often they succeed when assumptions hold.

# A FAILURE-AWARE EVALUATION STACK



To address this blindness, we propose a layered evaluation perspective that treats accuracy as necessary but insufficient.
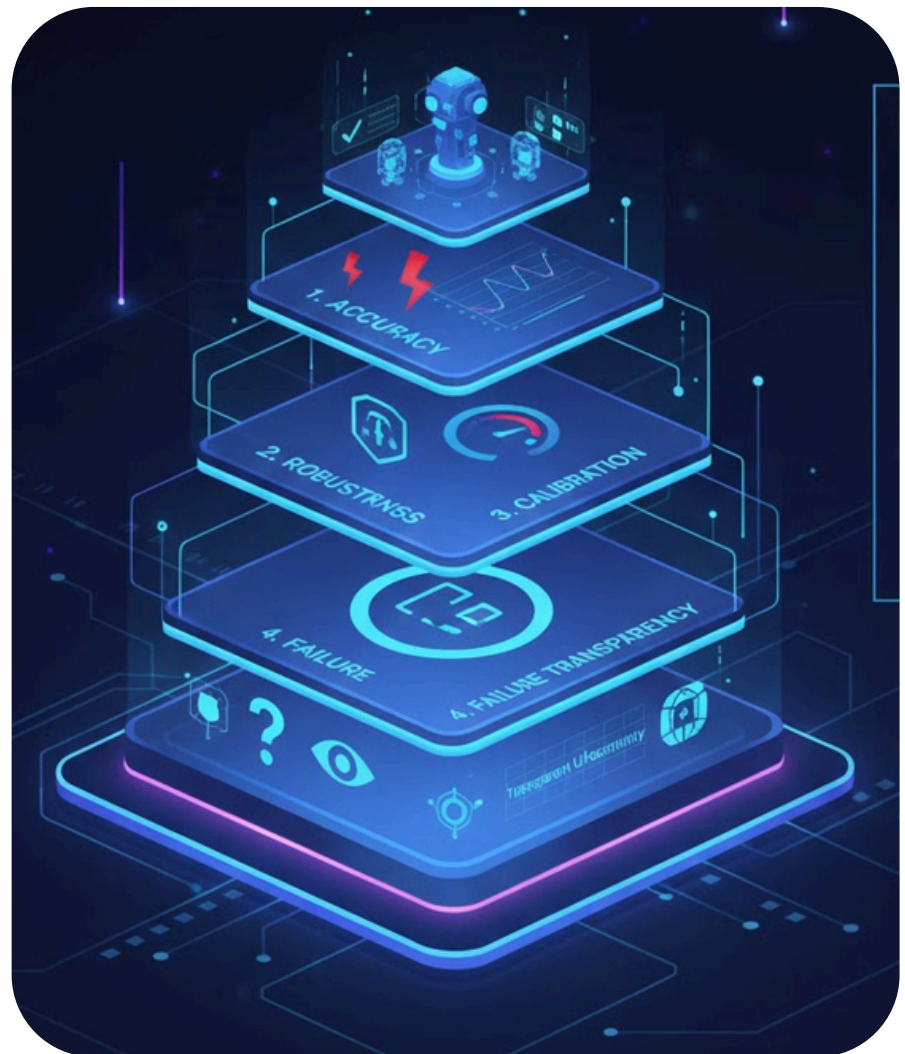
At the first layer, Accuracy measures correctness under standard conditions.

The second layer, Robustness, evaluates degradation under distribution shifts, corruptions, and constraint perturbations.

The third layer, Calibration, measures whether a model's confidence reflects its actual reliability—a prerequisite for safe decision-making.
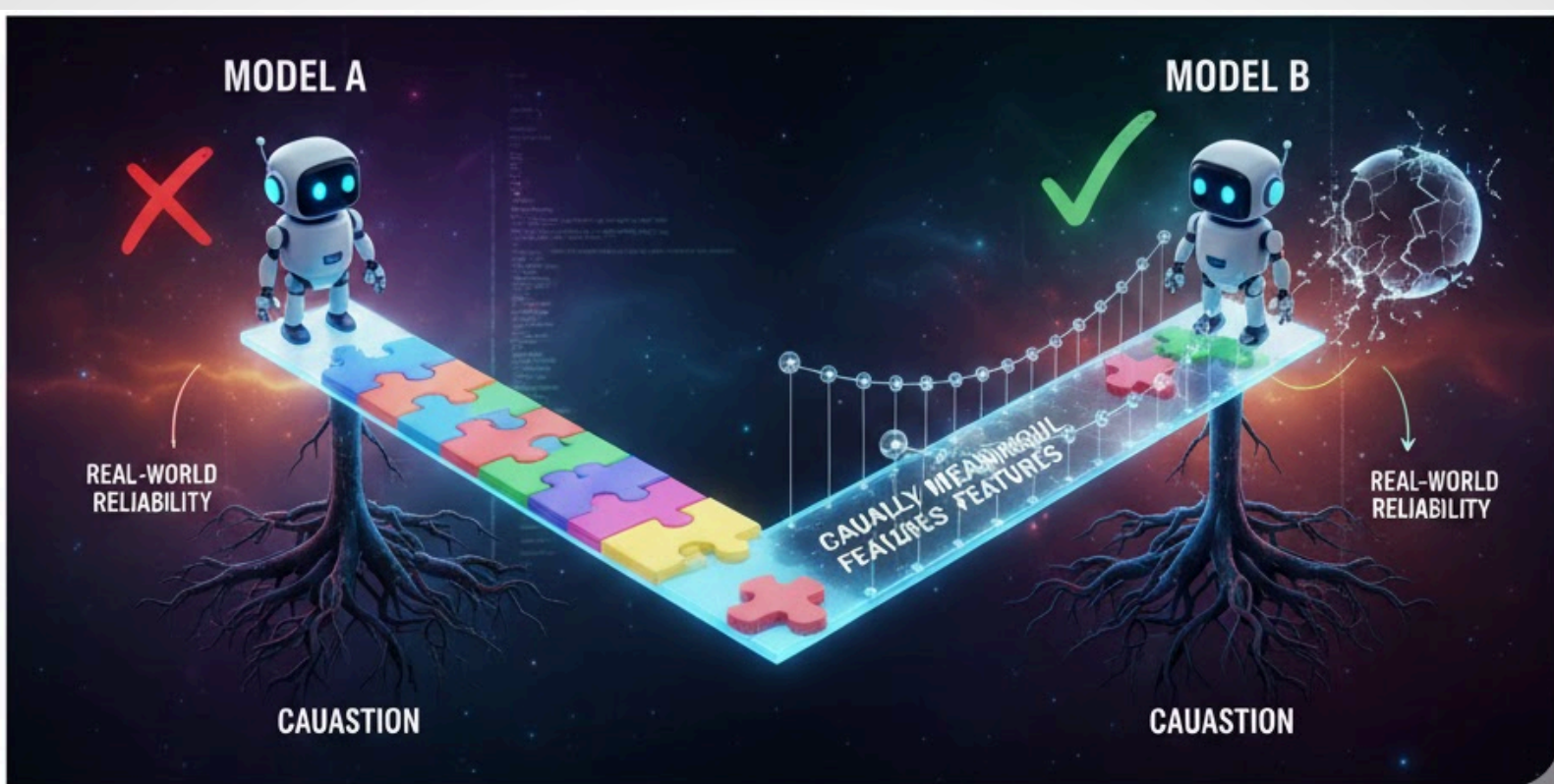
The fourth layer, Failure Transparency, assesses whether the system exposes uncertainty or fails silently with confident but incorrect outputs.

Under this framework, intelligence is not a scalar but a profile. Two systems with identical accuracy but different failure profiles are not equally intelligent.

# ACCURACY MEASURES OUTCOMES, NOT CAUSAL COMPETENCE



Accuracy evaluates whether a model's output matches a labeled target. It does not measure whether the model relies on causally meaningful features, whether it understands task constraints, or whether its internal representations generalize beyond the test distribution. As a result, two models with identical accuracy can differ radically in their internal mechanisms and real-world reliability.

This limitation is now well documented. Research on shortcut learning demonstrates that models often exploit superficial correlations that are predictive within a dataset but irrelevant to the underlying task [1]. Because accuracy does not penalize reliance on such shortcuts, it systematically rewards models that succeed for the wrong reasons.

The consequence is not merely academic. Systems optimized for correlation rather than causation may perform well during evaluation yet fail unpredictably when deployed in environments that violate hidden assumptions.