

From Measurement to Meaning: A Validity-Centered Approach to AI Evaluation

Olawale Salaudeen^{1*†} Anka Reuel^{2*} Ahmed Ahmed² Suhana Bedi²
Zachary Robertson² Sudharsan Sundar² Ben Domingue² Angelina Wang^{2,3‡}
Sanmi Koyejo^{2‡†}

¹ Massachusetts Institute of Technology ² Stanford University ³ Cornell Tech

Abstract

While the capabilities and utility of AI systems have advanced, rigorous norms for evaluating these systems have lagged. Grand claims, such as models achieving general reasoning capabilities, are evaluated with narrow benchmarks, like performance on graduate-level exam questions, which provide a limited and potentially misleading assessment. We provide a structured approach to reason about the types of claims that can be made about an evaluation given the available evidence. For instance, whether performance on a mathematical benchmark is an indication of merely math test scores or a broader kind of reasoning ability. Our framework is well-suited for the contemporary paradigm in machine learning, where various stakeholders provide measurements and evaluations that downstream users use to validate their claims and decisions. At the same time, our framework also informs the construction of evaluations that are valid for the claims that are intended to be made. We illustrate our framework through detailed case studies of vision and language model evaluations, highlighting how explicitly considering validity strengthens the connection between evaluation evidence and the claims being made.

1 Introduction

Suppose we *evaluate* an AI system’s ability to solve the International Mathematical Olympiad (IMO) by *measuring* its accuracy on such problems (Glazer et al., 2024). Then, we want to validate two distinct *claims* about the model’s capabilities with this evaluation:

Claim 1. The system can solve university-level math accurately.

Claim 2. The system has reached human-level reasoning.

Clearly, asserting Claim 2 requires a much greater inferential leap from the observed evidence (IMO accuracy) than Claim 1. If we claim that good performance on IMO problems demonstrates

*Equal contribution.

†Corresponding authors: olawale@mit.edu; sanmi@cs.stanford.edu.

‡Equal senior authorship.

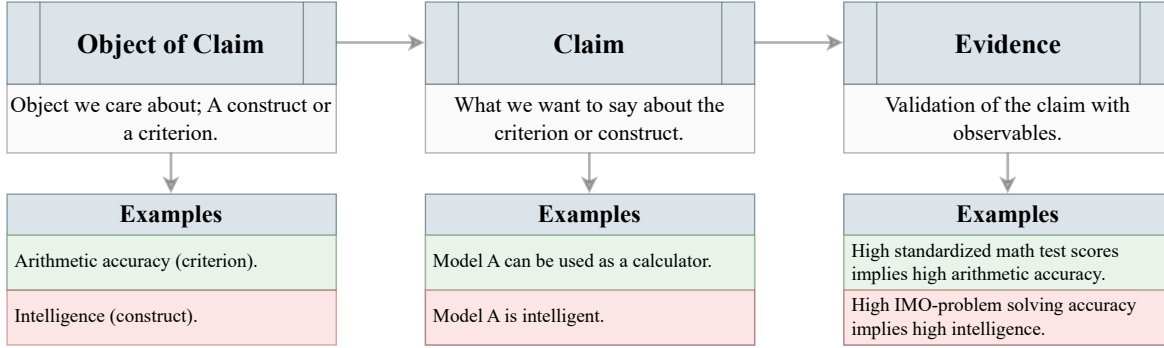


Figure 1: The three components are needed to begin the process of validation. First, we must decide what object our claim is about; is it a criterion, or is it a construct? Then, we must explicitly state the claim. Finally, we must identify our evidence and assess whether it supports the desired claim—i.e., do we have a valid claim based on the evidence? Here, a green background indicates that the claim-evidence pair is reasonably well supported. In contrast, a red background means the inferential leap between claim and evidence is larger and less well-supported.

competence in university-level math, the justification is reasonable: IMO problems often involve advanced undergraduate-level techniques, so proficiency in them provides reasonable evidence of university-level mathematical ability. However, if we claim that the system has reached human-level reasoning, the justification is much weaker. Solving IMO problems primarily requires mathematical problem-solving, but human reasoning encompasses a broader spectrum—including common sense, adaptability, and metacognition—which IMO performance alone does not assess. This difference highlights that we must scrutinize an evaluation and measurement in the context of the claim we wish to support. Thus, we consider validity. Validity refers to the degree to which evidence and theory support the interpretations of test scores for proposed uses of tests.

In measurement theory¹, ‘measurement validity’ refers to the extent to which a test accurately measures what it is intended to measure. Yet, validity is not inherently a property of the measurement itself—it also depends on the context of the evaluation it enables, the claims supported by that evaluation, and the potential real-world consequences of those claims (Messick, 1998, 1995; Shepard, 1993). Notably, this is counter to Borsboom et al. (2004)’s view that validity is solely a property of the test.

We first clarify some terms (we refer the reader to also carefully examine Table 1). In this work, measurement refers to assigning a quantitative or qualitative value to a specific property of a system (e.g., accuracy or usability), while evaluations are the broader process of interpreting these measurements to provide insights about the system. Claims, judgments, assertions, and decisions about the system can then be supported by measurements and evaluations. Measurement instruments are tools to quantify, assess, or categorize an object and include benchmarks, user studies, and expert assessments.

For instance, we can measure accuracy on question–answer problems. However, the context of applying this measurement to IMO problem-solving problems (our measurement instrument) makes it an evaluation of a system’s accuracy in answering IMO problem-type questions; we are not merely

¹Historically, measurement theory studies how latent psychological concepts are quantitatively assessed, ensuring that the chosen measures accurately capture the intended constructs while maintaining validity and reliability. We discuss validity and its treatment in other scientific disciplines in Appendix B.

Table 1: Table of terms and definitions.

Term	Definition	How does it relate to other terms?	Example
Measurement Instrument	A tool used to gather observations or assign values (e.g., a benchmark, user study, or survey).	Underlies the act of measurement. Evaluation and claims often hinge on data obtained via instruments.	A dataset of IMO math problems (the “instrument”) used to gather system accuracy scores.
Measurement	Assigning a quantitative or qualitative value to a property of a system (e.g., accuracy, usability).	Involves applying the instrument and recording results; informs subsequent evaluation.	“The system answers 15 of 20 IMO questions correctly” (accuracy = 75%).
Evaluation	The broader process of interpreting one or more measurements in context.	Translates raw measurement into insights (e.g., domain-specific analysis, comparisons to a baseline).	“Because the system can solve 75% of these IMO problems, it demonstrates proficiency in competition-level algebra.”
Claim	An assertion, judgment, or decision made about the system, potentially based on evaluation results.	Draws on evaluation evidence to generalize or conclude something about the system or its capabilities.	“The system exhibits human-level math reasoning skills.”
Criterion	A directly measurable or observable concept (e.g., ‘university-level math accuracy’).	Can be measured directly and often serves as a baseline or gold standard for evaluation.	“University-level exam accuracy” – the system’s performance on real university math exams.
Construct	An abstract concept not directly measurable (e.g., ‘mathematical reasoning’ or ‘trustworthiness’).	Requires an operational definition plus proxies or indicators to measure and evaluate indirectly.	“Mathematical reasoning” – a theoretical ability captured through various problem sets and expert assessments.

recording an accuracy score but interpreting it in a specific domain context (IMO problems) to gain some insight about the system’s capabilities. To belabor the point, while the measurement is an accuracy score, the interpretation of that measurement as an indicator of math problem-solving capability is an evaluation. Finally, one may then make claims—not necessarily correctly—about general reasoning capabilities to be supported by the evaluation.

For another example, we can measure the frequency of harmful outputs (e.g., misinformation or offensive responses) from a language model. However, the context of applying this measurement specifically to high-stakes medical advice scenarios (our measurement instrument) transforms it into an evaluation of the system’s safety in that domain; we are not merely counting harmful responses but interpreting their potential impact within a clinical setting. From there, one might make broad claims about the system’s trustworthiness or readiness for real-world deployment—claims that may be more or less justified depending on whether the measurement truly captures the range of potential harms and aligns with relevant medical standards. This full pipeline is relevant context for establishing validity.

We can measure without evaluating—for example, by collecting raw accuracy scores without concluding their implications. However, to evaluate, we must measure in some form (quantitatively or qualitatively) and then interpret those measurements in a domain-specific context. One might then ask, why measure if not to evaluate a system? We can measure as a means to develop new metrics for future evaluation; we can measure to observe or characterize phenomena before making judgments; we can measure for calibration. For instance, we may measure accuracy to identify patterns or outliers that guide future studies without evaluating whether the system’s performance is ‘good’ or meets real-world requirements.

A core limitation of current AI evaluation discourse is that validity—if considered—often focuses

Table 2: We provide an overview of the different forms of validity considered in this work, along with key questions to ask in their assessment. The standard of evidence for validity depends on the conceptual gap between the measurement and the object of the claim, with broader gaps demanding stronger justification. Certain forms of validity, such as criterion validity, encompass multiple facets that capture different aspects of the evaluation. We adopt a mix of [Cronbach and Meehl \(1955\)](#) and [Messick \(1995\)](#)’s views on validity. Like [Cronbach and Meehl \(1955\)](#), we do not categorize all validity uniquely as facets of construct validity; however, like [Messick \(1995\)](#), we consider external and consequential validity also critical.

Validity Type		Description	Example: IMO Problem Solving → Reasoning
Content Validity		Does your evaluation cover all relevant cases?	Does solving IMO problems sufficiently capture the content relevant to reasoning?
Criterion Validity		Does your evaluation correlate with a known validated standard?	Does IMO problems accuracy predict other external criteria of reasoning, e.g., common sense reasoning benchmarks?
	Predictive Validity	Can your evaluation predict downstream outcomes?	–
	Concurrent Validity	To what extent does your evaluation agree with another validated assessment under the exact same conditions?	–
Construct Validity		Does your evaluation truly measure the intended construct?	Does IMO problem solving capture all components of reasoning and only components of reasoning?
	Structural Validity	Does your evaluation capture the structure of the construct you are measuring?	–
	Convergent Validity	Does your evaluation correlate with other measures that assess the same construct?	–
	Discriminant Validity	Can your evaluation differentiate between constructs that should be distinct?	–
External Validity		Does your evaluation generalize across different environments or settings?	Does excelling at IMO problems translate to solving university-level math problems where the problems are provided in different formats?
Consequential Validity		Does your evaluation consider the real-world impact of test interpretation and use?	Does emphasizing IMO problem-solving in AI development narrow research focus in ways that overlook other essential reasoning skills?

on the measurement–evaluation relationship, i.e., designing measurements that sufficiently support a predefined evaluated object ([Gema et al., 2024](#); [Wallach et al., 2025](#))—an important aspect of validity. Consequently, one may then conclude that if a measurement does not fully meet the needs of the object it was designed to evaluate, then no valid claims can be made, and no insights can be gained; but, just because IMO accuracy does not sufficiently measure reasoning does not mean that it cannot support better-scoped claims. One may also conclude a need for strong ties between potentially distinct stakeholders who develop measurements, those who perform evaluations, and those who make decisions. Additionally, one may conclude an ordering where we must start with the claims we want to make, then decide what objects to evaluate to support such claims, and then design the appropriate measurements. While this process is perhaps ideal, it differs from real-world practices.

For instance, a policymaker might claim that an AI system is sufficiently safe for deployment based on apriori evaluations and measurements from corporations. Then, academic researchers may

develop new measurements and evaluations to test safety and support or refute the claim. Furthermore, even when measurements are limited, we can still derive some claims and insight—though they are often necessarily more narrow. Our framework allows claims and evaluations to come in any order, potentially from independent stakeholders, and allows for deriving as much valid information from measurements as possible.

To assess the validity of a claim derived from an evaluation (or vice-versa), Figure 1, we must:

1. Carefully consider the object of the claim. Is it a construct—an abstract object that cannot be measured directly, like ‘mathematical reasoning’? Or is it a criterion—a directly measurable object, such as ‘university-level math accuracy’?
2. Furthermore, does the claim refer to the same property measured, or does it extend beyond the specific evaluation to infer something about a different property? For example, one might measure IMO accuracy as part of an evaluation of mathematical ability, then use this evaluation to support a claim about university-level math ability, which is a different object.
3. Finally, is the claim supported directly by the measurement (e.g., IMO accuracy implies university-level math capability), or does it rely on an intermediate construct (e.g., IMO accuracy implies mathematical reasoning, which in turn implies university-level math accuracy)?

These distinctions determine the necessary standards of evidence required before an evaluation can meaningfully support a claim. The alignment between what is measured, how it is interpreted (evaluation), and the overarching claim is central to establishing validity.

Ensuring validity requires five forms of validity from psychometrics (Thorndike, 1949; Cronbach and Meehl, 1955; Messick, 1995; Borsboom et al., 2004), outlined in Table 2. Additionally, Appendix B Table 4 enumerates tools to investigate and establish each form. While other forms of validity exist², we identify this set as most relevant for current AI measurement validity gaps in Section 2.

The standard of evidence required to demonstrate validity depends on the conceptual gap between what is actually measured (and how it is evaluated) and the object of the desired claim. The greater the gap, the more arduous the task of establishing validity. Notably, different forms of validity are not independent and work together to demonstrate validity—we illustrate this in our case studies in Section 3 and Appendix D. Our proposed framework ties claims directly to the requisite standard of evidence provided by the evaluation. This alignment ensures the appropriate downstream use of evaluations while also guiding improvements that support more general claims.

This framework is pivotal because AI evaluations inform decisions with real-world consequences. For example, under Article 51 of the EU AI Act (5), benchmarks are explicitly referenced as indicators for classifying AI models according to their systemic risk. Developers of models deemed high-risk must comply with significantly more obligations. If we fail to consider validity in this context, such classification may become meaningless because we cannot be certain that what truly matters—namely, the risk posed by these models—is accurately captured by the chosen measurement instruments (e.g., benchmarks). This can lead to a false sense of security. Similarly, measurement instruments are often used within organizations (Hardy et al., 2024) to guide resource allocation and further training aimed at improving a model’s capabilities. Yet, if the chosen instrument does not accurately measure the capability developers care about, additional training may simply become

²Other forms of validity include, for example, face validity. Additional forms of validity are given in (Lim, 2024; Hughes, 2018).

an exercise in ‘teaching to the test’ (Jennings and Bearak, 2014) rather than leading to genuine improvements in the model.

Recognizing these challenges, we propose a structured framework to assess the validity of AI assessments, ensuring appropriate use and interpretation. Specifically, our contributions are:

1. We examine and identify limitations in how the relevance of different forms of validity has co-evolved with the progress of AI and the corresponding norms and practices of evaluation.
2. We propose a practical and structured claim-aware framework for identifying the necessary evidence to establish the validity of claims based on AI evaluations. We also enumerate adoptable practices to demonstrate validity.
3. We illustrate our framework through vision and language evaluation case studies, providing concrete, prescriptive examples of validating claims based on evaluations.

2 Validity Gaps in Current AI Evaluations and Related Work

AI evaluation has evolved alongside the complexity of AI tasks. Still, the gap between measurement and the object of claims about AI utility has widened (Appendix C), with deficiencies across all five forms of validity—content, criterion, construct, external, and consequential (Table 2). This widening gap results from changes in how AI systems are evaluated and the claims necessary to imply their utility. Early systems were tested on held-out samples from the same distribution (independent and identically distributed setting), ensuring content validity, i.e., relevant conditions were well-represented in the training data. Then, with pertaining, systems were first trained on large datasets like ImageNet (Deng et al., 2009; Russakovsky et al., 2014) and then finetuned (or transfer learning) for specific tasks, shifting some evaluation towards predicting the downstream performance learned pretraining representations and initializations, a form of criterion validity (Kornblith et al., 2018; Recht et al., 2019).

Subsequent studies on spurious correlations and out-of-distribution generalizability (Bai et al., 2025; Salaudeen and Koyejo, 2021, 2024; Lopez-Paz et al., 2016; Xiao et al., 2020; Arjovsky et al., 2019; Rosenfeld et al., 2020; Koh et al., 2020; Gulrajani and Lopez-Paz, 2020), causal representations (Schölkopf et al., 2021), biased representations (Gichoya et al., 2022), etc., emphasized external (Salaudeen and Hardt, 2024), consequential (Wang and Russakovsky, 2023), and some construct (Bell et al., 2024; Salaudeen et al., 2025) validity.

However, evaluations remained benchmark-driven, primarily supporting claims of technical progress (Hardt and Recht, 2021; Orr and Kang, 2024)—this use is surprisingly robust (Blum and Hardt, 2015; Salaudeen and Hardt, 2024). Now, while these evaluation norms align researchers and industry, accelerating advances (Donoho, 2023; Recht, 2024; Barocas et al., 2023), they fail to predict real-world reliability (Hardy et al., 2024).

The rise of foundation models, which can operate across diverse tasks without finetuning, further complicates this issue. Traditional evaluation methods increasingly fail to capture real-world AI behaviors that require investigating abstract capabilities like trustworthiness and reasoning (Wu et al., 2023; Wan et al., 2024; Mirzadeh et al., 2024).

Narrow datasets used for “general-purpose” evaluation raises content, construct, and external validity concerns, especially for complex tasks like reasoning (Bostrom et al., 2020; Alaa et al., 2025).

Furthermore, these evaluations, already saturating (Ott et al., 2022), lack criterion validity and fail to predict criteria of real-world applicability (Hardy et al., 2024). At the same time, the socio-technical gap between evaluation results and real-world needs undermines consequential validity (Liao and Xiao, 2023). Consequently, overgeneralized results erode evaluation credibility (Raji et al., 2021).

Prior work has demonstrated the need for validity frameworks (Jacobs and Wallach, 2021; Saxon et al., 2024; Subramonian et al., 2023; Xiao et al., 2023; Blodgett et al., 2021; Coston et al., 2023; Xiao et al., 2024; Reuel et al., 2024). Yet, much of it has focused on measurement validity and limitations of measurements, as well as conditions for perfect measurements of nicely defined concepts, which is far from practice. METRICEVAL (Xiao et al., 2023) raises validity concerns stemming from vague benchmark articulation and repurposed datasets in 16 natural language generation metrics. The important work of (Chouldechova et al., 2024; Wallach et al., 2025) applied (Adcock and Collier, 2001)’s measurement theory, critiquing ML evaluations for conflating systematizing a background concept with operationalizing a systemized concept³.

Our work complements this literature by explicitly identifying that validity depends not only on the measurement and evaluation but also on the claim intended to be made. Building on Wallach et al. (2025)—who underscore the importance of explicit systematization needed in AI, where concepts often emerge from practice rather than theory—our work clarifies this process in the context of nomological networks (Cronbach and Meehl, 1955). These networks represent not only the relationship between the background concept and systematized concept but also the broader relationships to other background and systematized concepts. In the sense of the Duhem-Quine thesis⁴, a nomological network serves as a map of empirical and theoretical relationships, helping manage the holistic nature of scientific testing by clarifying how constructs relate to each other as well as observable evidence. Consequently, we expand upon Wallach et al. (2025)’s view of systematization by arguing for the importance of broader nomological networks beyond just specifying which definition of a background concept will be used.

Liu et al. (2024) further challenged benchmarks’ ability to measure intended constructs, proposing the Evidence-Centered Benchmark Design (ECBD) framework to ensure rigorous metric selection. However, these works focus on the process of designing new measurement instruments for evaluations. While developing better measurement instruments for better evaluations is also important, and our framework also applies to this task, we find it of practical value to understand what claims can be made from existing evaluations and evidence, given the intractability of creating tailored evaluations for each claim, and the already unwieldy amount of existing benchmarks (Copilot, 2024).

Ultimately, our framework takes a practical approach, emphasizing that validity is not only a property of measurement and evaluation. As Cronbach and Meehl emphasize (Cronbach and Meehl, 1955): ‘In one sense, it is naive to inquire ‘Is this test valid?’ One does not validate a test, but only a principle for making inferences. If a test yields many different types of inferences, some can be valid and others invalid.’ We further enumerate risks, tools, and evidence exemplars to assess whether evaluations meet appropriate validity standards in Section 4 and Table 4.

³According to Adcock and Collier, a background concept is a “broad constellation of meanings and understandings associated with [the] concept,” and systematization describes the process of refining and explicitly defining a concept to create a structured and consistent foundation for measurement and analysis (the systematized concept) while operationalization the process of transforming a systematized concept into measurable indicators (Adcock and Collier, 2001).

⁴The Duhem-Quine thesis emphasizes that scientific claims are interconnected, meaning that rejecting or modifying one hypothesis affects others within the theoretical framework.

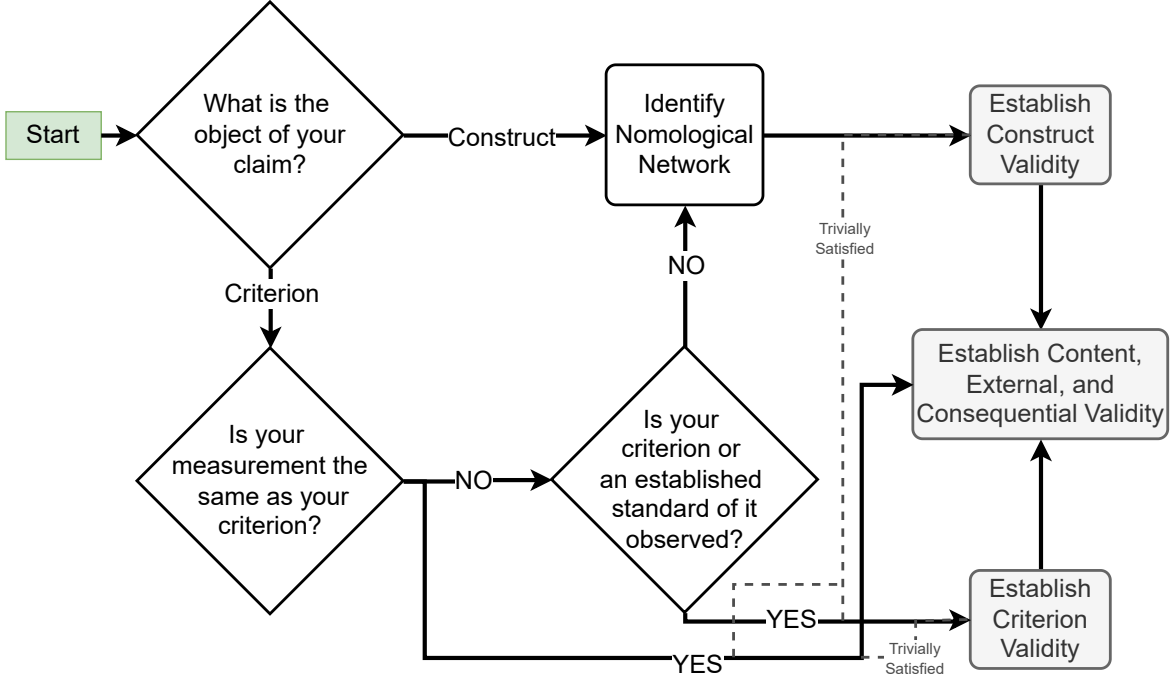


Figure 2: Decision tree for establishing validity. For the decision processes that do not directly go through establishing construct or criterion validity, our argument is not that those forms of validity are irrelevant, but rather that they may be trivially satisfied in the context of the measurement, evaluation, and claim.

3 A Framework for Claim-Centered Validity Assessment in AI Evaluation

In this section, we categorize when and how different forms of validity are most critical for supporting a claim with measurements and evaluation. While we maintain that all forms of validity are always necessary, some may be trivially satisfied depending on the measurement, evaluation, and claim context. Rather than applying uniform scrutiny to all forms of validity, we account for context-dependent nuances that make certain forms particularly significant in some cases; this is distinct from previous work (Adcock and Collier, 2001; Wallach et al., 2025).

Establishing validity is an iterative process (Cronbach and Meehl, 1955; Kuhn, 1997), and recognizing the limitations of measurements and evaluations—rather than outright rejecting them when certain validity criteria fall short—requires nuance. This approach is essential for practical utility, enabling us to extract meaningful claims even from evaluations that do not rigorously satisfy all conditions of validity.

A claim can (and perhaps should (Wang et al., 2024)) be supported by many evaluations and measurements. However, for simplicity and without loss of generality, we focus on a single measurement and evaluation.

Recall that the object of a claim can be a criterion—directly measurable—or a construct—abstract and not directly measurable. The primary considerations for investigating validity are determined

	Claude 3.5 Sonnet	Claude 3 Opus	GPT-4o	Gemini 1.5 Pro	Llama-400b (early snapshot)
Graduate level reasoning <i>GPQA, Diamond</i>	59.4%* 0-shot CoT	50.4% 0-shot CoT	53.6% 0-shot CoT	—	—
Undergraduate level knowledge <i>MMLU</i>	88.7%** 5-shot	86.8% 5-shot	—	85.9% 5-shot	86.1% 5-shot
	88.3% 0-shot CoT	85.7% 0-shot CoT	88.7% 0-shot CoT	—	—

Figure 3: A real-world example of a developer presenting the Graduate-Level Google-Proof Q&A Benchmark (GPQA) as a proxy for the graduate-level reasoning capabilities of their system. Additionally, SWE-bench is used as a proxy for agentic coding, and TAU-bench as a proxy for agentic tool use (Anthropic, 2024; Rein et al., 2023).




by the following (Figure 1):
















1. Is the object of the claim a criterion (e.g., math exam accuracy) or a construct (e.g., reasoning ability)?
2. Is the measurement the same as the object of the claim (e.g., evaluating IMO accuracy when IMO problems are also the object of the claim)?
3. Does the measurement directly imply the claim, or does it require a mediating construct (e.g., does IMO problem-solving imply university math exam problem-solving, OR does it imply mathematical reasoning, which implies math exam problem-solving)?

The five forms of validity we foreground are relevant in different ways. Ideally, we should validate a claim by directly measuring the object it is about. However, this may not be possible. When we perform a measurement but the object of the claim is a different criterion (e.g., evaluate IMO problem solving → make claims about university math problem solving), criterion validity is most important—ensuring the measurement reliably predicts the object of the claim (predictive validity) or an established external standard of the object of the claim (concurrent validity). When neither the object of claim nor an established external standard is available, we may validate the claim through an intermediate construct (evaluate IMO problem solving → infer mathematical reasoning → make claims about university math problem solving), requiring construct validity.

When the object of the claim is itself a construct, and we directly measure and evaluate its proxies (e.g., evaluate IMO accuracy → make claims about mathematical reasoning), construct validity is essential to determine whether the measurement genuinely measures the intended construct rather than an unrelated or superficial correlation. This is also necessary when we aim to validate a claim about a construct with measurements and evaluations of proxies of other constructs (e.g., evaluate IMO problem solving → infer logical reasoning → make claims about mathematical reasoning).

Importantly, a claim about a construct cannot be validated in isolation—instead, it gains meaning and validity through its relationships with other constructs and observable measures. Cronbach and Meehl’s nomological network (Cronbach and Meehl, 1955) provides a rigorous, albeit historically challenging to operationalize, way to reason about constructs within a broader abstract and

Table 3: A Graduate-Level Google-Proof Q&A Benchmark (GPQA) (Rein et al., 2023) Application. A subjective score for validity—the standard for “reasonable” is demonstrating that obvious risks to invalidity are addressed: : reasonable; : proceed with caution; : insufficient. Even for a score of “reasonable,” there will be weaknesses in the evidence. The score is given because the strengths outweigh the weaknesses in terms of determining the validity of the claim from that evidence. This is never a binary classification nor complete, and should rather be a cyclic process— for instance, as our forms of what constitutes graduate-level chemistry may evolve over time and from school to school.

Claims from Graduate-Level Google-Proof Q&A (GPQA) Benchmark Accuracy Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. AI systems can accurately answer <i>graduate-level specialized multiple-choice questions</i> in biology, physics, and chemistry.					
2. AI systems can accurately answer <i>graduate-level specialized questions</i> in specialized scientific domains.					
3. AI systems can exhibit <i>general reasoning abilities</i> that can transfer beyond current human specialization.					

empirical system, allowing for more scientifically robust validation. A nomological network is a conceptual framework that maps the relationships between constructs and criteria (Cronbach and Meehl, 1955). Figure 4 gives an example from human psychology. An explicit nomological network for AI constructs, while vital, remains a missing piece in efforts toward valid AI evaluations, limiting the establishment of validity when constructs are involved. Although a detailed treatment of nomological networks is beyond the scope of this work, we emphasize their importance in establishing validity and explicitly indicate where they are necessary in our framework. We refer the reader to Cronbach and Meehl’s seminal work for more detail (Cronbach and Meehl, 1955).

Next, we illustrate our framework for determining validity. We focus in the main text on GPQA (Graduate-Level Google-Proof Q&A) accuracy as our measurement and evaluation. We then investigate claims of varying generality commonly made from this evaluation (Buntz, 2025; Rein et al., 2023). Additional examples are in Appendix D. This clarifies that a given measurement may not support broad claims, yet it can still be highly useful for supporting more narrowly defined ones. This adds necessary nuance to the discourse on validity in AI assessment.

Table 3 summarizes the following example of applying our framework to assess the validity of claims from GPQA multiple-choice question-answering accuracy. We supplement this section with detailed case studies in the context of evaluating popular vision and/or language AI systems in Appendix D.

Sources of Validity Evidence. While we restrict ourselves to the evidence of validity provided in the GPQA paper for the previous analysis for brevity and simplicity, establishing validity can (and should) be done across multiple asynchronous studies and various stakeholders.

Furthermore, we can start from a claim and use our framework to determine the necessary types of measurements and evaluations to support it, but this can also be done in reverse. Adcock and Collier (2001)’s proposed framework for measurement validity starts with the construct, as does Wallach et al. (2025)’s framework for assessing generative AI. However, increasingly, measurements and evaluations are provided by some set of stakeholders, while other stakeholders are left to make

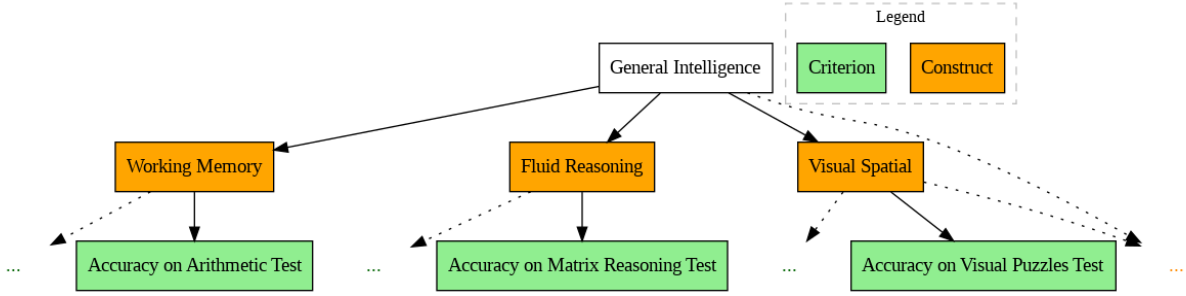


Figure 4: This figure illustrates a nomological network for human general intelligence according to the Wechsler Intelligence Scale, adapted from Canivez et al. (2017). *We note that this breakdown of general intelligence is for illustrative purposes only; it does not necessarily translate to artificial general intelligence.* The network consists of background concepts (blue w/ rounded corners) linked by hypothesized associations, reflecting abstract expectations. Observable indicators (green w/ sharp corners) represent criteria, measurable variables used to assess the constructs. Establishing a robust nomological network is critical to ensuring construct validity by demonstrating abstract coherence, convergent validity, and discriminant validity within empirical research.

sense of these relative to downstream uses of AI systems.

Additionally, when going through the process of validating claims from measurements and evaluations, the five forms of validity we foreground in Table 2 are not uniformly important or always necessary. In the following sections, we will describe situations where some types of validity are trivially satisfied.

3.1 Scenario 1: The object of the claim is a criterion

We first consider the setting where the object of the claim is a criterion, and the desired claim in this setting looks like this: ‘A higher measurement score predicts a higher/lower criterion score.’ ‘A higher GPQA accuracy predicts a higher PhD qualifying exam questions accuracy.’ In this case, the measurement and the criterion object of the claim can be (i) identical (GPQA accuracy vs. new GPQA question accuracy), (ii) proxies of the same underlying construct (GPQA accuracy vs. general scientific question-answering accuracy), or (iii) proxies of two different but related underlying constructs (GPQA accuracy (scientific reasoning) vs. surgical planning accuracy (medical reasoning)). In each case, we must justify why the measurement supports the claim.

In scenario (i), we are primarily concerned with content and external validity, i.e., does our measurement cover the relevant content of the criterion, and does it generalize to relevant contexts beyond that of the measurement? Construct validity and criterion are trivially satisfied as a consequence of directly measuring and evaluating the object of the claim. This could be because the hard work of systematization and operationalization of the construct we would have otherwise attempted to measure and evaluate has already been done (Adcock and Collier, 2001; Wallach et al., 2025).

Identical Measurement and Criterion. *Example.* Our object of claim is a criterion: multiple-choice questions accuracy in physics, chemistry, and biology⁵. We then aim to support claims about

⁵Thresholding is commonly used in real-world decision-making to transform continuous measurements (e.g., confidence scores) into binary or categorical outcomes (e.g., pass/fail, high risk/low risk). The choice of threshold can

an AI system’s accuracy on such questions by measuring and evaluating the system’s accuracy on the GPQA dataset; we must:

- Establish content validity: GPQA has expert-curated questions, which enhance content validity by ensuring relevance and rigor across biology, physics, and chemistry, with the performance gap between experts and non-experts indicating effective assessment of specialized knowledge. However, the construction criteria may inadvertently exclude certain relevant topics, potentially skewing subfield representation. Systematic content mapping and expert diversity analysis can strengthen validity by ensuring comprehensive coverage and mitigating selection biases. In modern AI, red-teaming can also help identify the overage gaps that hinder content validity (Perez et al., 2022).

If content validity holds and one does not expect that the context of measurement is different than the context of the claim, then the claim can be supported by the measurement. However, if the claim must hold in a different context than the measurement, one must also:

- Establish external validity: The GPQA measurement reflects real-world conditions, with human experts developing questions and a measurement format aligned with academic multiple-choice assessments, so the context of measurement is aligned with the context of the claim in this sense. However, the human assessment may not generalize beyond the measurement context, and without comparison to other multiple-choice science tests, its generalizability remains unverified. To strengthen external validity, validation against diverse question formats, question types, and other variations of context is necessary.

This setting is commensurate with traditional AI benchmarking practices. Many AI benchmarks have focused on these forms of generalization, including classical generalization (Vapnik and Chervonenkis, 1971) and out-of-distribution generalization (Shimodaira, 2000). By ensuring strong content and external validity, such benchmarks provide a solid foundation for validation claims for directly measurable criteria.

Different Measurement and Criterion. For (ii)-(iii), different criteria that are either proxies of the same or different but related mediating constructs. In this case, we ideally additionally directly establish criterion validity. That is, establish that the object that is measured is predictive of the desired criterion or an established standard. Then, the existence of these mediating constructs may inform how we establish criterion validity, but we do not need to reason about them directly to establish validity.

Example. Our object of claim is scientific question-answering, and we want to quantify general scientific question-answering using GPQA accuracy as evidence. We still need to demonstrate that the measurement covers relevant content and generalizes to all the contexts we want the claim to hold (content and external validity). However, we must additionally establish that the measurement of GPQA accuracy is predictive of the scientific question-answering criterion or a validated standard, i.e.:

profoundly affect both evaluation and claim validity: even a perfectly measured property may lead to an invalid claim if the threshold does not align with the intended context or the actual consequences of misclassification. This is known in psychometrics as standard testing (Cizek and Bunch, 2007).

- **Establish Criterion Validity.** Human expert accuracy provides a strong external criterion, supporting concurrent validity, while the AI-expert performance gap reinforces the benchmark’s credibility. However, there is no evidence of predictive validity, as accuracy has not been tested against future performance on specialized assessments, and concurrent validity remains incomplete without correlations to established external measures of expertise, such as standardized exams in other fields. To strengthen criterion validity, correlations should be established with real graduate program exams for concurrent validity, and predictive validity studies should track the system’s downstream performance across domains.

However, if criterion validity is implausible in this way, we may attempt to leverage our understanding of the underlying structure in constructs and their known mapping to observables, when it is available, to establish validity, i.e., a nomological network. Importantly, such a nomological network often does not exist in the current paradigm of AI assessment.

When a nomological network is unknown, establishing validity becomes significantly more difficult, as there is no agreed-upon basis for interpreting how abstract constructs like ‘reasoning’ map to measurable criteria (e.g., ‘GPQA accuracy’). In such cases, evaluations risk being narrow or misleading—a system might excel at scientific reasoning yet still lack medical reasoning, and evaluators could erroneously assume success in one facet implies overall ‘reasoning.’ Without explicit connections between sub-constructs and corresponding measurements, conflicting results may emerge, and different assessments might rely on unfounded inferences about a system’s capabilities. This lack of structure not only obscures whether a measurement provides meaningful evidence for a given claim but also undermines the reliability of validity assessments, leaving practitioners vulnerable to inflated claims and misguided deployment decisions.

When such a network is available, to establish construct validity, we utilize its facets: structural, convergent, and discriminant validity:

- **Establish Construct Validity:** For brevity, please refer to the subsequent discussion in section 3.2 on when the object of the claim is a construct.

Importantly, when the measurement and criterion are proxies of different constructs, we must also validate relationships between constructs in addition to their relationships to observables, i.e., how is scientific reasoning related to medical reasoning. Doing this also requires knowledge of a nomological network. For example, in Figure 4, the accuracy on the arithmetic test and accuracy on the matrix reasoning test must go between working memory and fluid reasoning.

Example. Suppose we want to claim that AI systems can reason about the outcomes of different surgical plans. To evaluate this, we must first define what reasoning entails. Suppose a model evaluator interprets reasoning as scientific reasoning and uses the GPQA benchmark to measure it. However, the object of the claim is most related to medical reasoning. At this stage, defining reasoning in this way is neither inherently valid nor invalid.

Now, suppose we want to claim that strong GPQA performance translates into accurate surgical planning; this requires several inferential steps. GPQA (insufficiently) assesses scientific reasoning, while surgical planning likely relies on medical reasoning, potentially a different subspace of reasoning—according to one’s nomological network. Establishing structural validity requires examining whether GPQA captures the key components of general reasoning relevant to surgical decision-making. Without showing that GPQA performance reflects the same underlying capabilities as medical reasoning, claims about AI outperforming surgeons based on GPQA remain unverified.

Establishing convergent validity through latent variable modeling and item-response theory can help demonstrate if the measurement captures variance in the latent subspace shared between the two constructs that determine the outcome criterion.

3.2 Scenario 2. The object of the claim is a construct

In many cases, we want to validate a claim about a construct by evaluating its proxies. This looks like: ‘A higher measurement score implies a higher latent capability, e.g., GPQA accuracy to scientific reasoning.’ Then, construct validity is paramount.

Example. Suppose the object of the claim is general reasoning and we want to make a claim about a system’s general reasoning ability by measuring GPQA accuracy. Here, we must establish all five of our forms of validity, especially construct validity (recall it is composed of structural, convergent, and discriminant validity):

- **Establish Construct Validity:** Performance on GPQA aligns with success in structured question-answering tasks, suggesting some reasoning component. However, structural validity is unclear, as the test may not sufficiently capture the rank of reasoning. Convergent validity is unverified since GPQA accuracy has not been correlated with other explicit reasoning assessments. Discriminant validity is also uncertain, as it remains unclear whether GPQA measures genuine scientific reasoning or simply domain-specific knowledge and memorization. Comparing performance to humans with access to Google is an attempt to do this. To address these concerns, methods like factor analysis (Kim and Mueller, 1979) should be conducted to distinguish reasoning from memorization, and performance should be validated against other dedicated reasoning assessments while ensuring it diverges from pure knowledge recall.

Additionally, content and external validity must be established to confirm that the essential aspects of the construct are accurately measured and that findings generalize to unmeasured components. Moreover, criterion validity, when a construct-relevant criterion or established standard is available, can support construct validity since well-designed measurements should reliably predict external outcomes related to the same construct.

Consequential Validity. Consequential validity examines whether the real-world outcomes of decisions based on an assessment align with its intended purpose. In the case of GPQA, if the benchmark effectively measures scientific reasoning, AI models that perform well on it could support decision-making in scientific research or education. However, there is a risk of overgeneralization—high GPQA accuracy might lead to misinterpreting AI as possessing broad reasoning abilities when it may only excel at structured multiple-choice problems. In this case, there could be harmful consequences like replacing human workers with ill-suited technology.

For strong consequential validity, GPQA measurement must align with the reasoning skills they intend to measure, ensuring AI performance is interpreted within its actual capabilities. Clear performance guidelines should distinguish validated reasoning abilities from speculative claims, preventing misapplications of AI in scientific decision-making.

4 Evidence of Validity

Next, we examine common risks to validity claims and discuss existing tools and methodologies for assessing and strengthening validity in AI assessment. Appendix A Table 4 categorizes in detail key risks, investigation tools, and evidence exemplars across multiple forms of validity in assessment.

Risks to content validity include coverage deficiency, where important aspects of the construct are missing, and construct irrelevance, where extraneous factors influence scores (AERA, 2014; Messick, 1995). Imbalanced content can lead to assessments overemphasizing certain skills while neglecting others. These issues can be examined through expert review, adversarial scrutiny, and synthetic data generation, with supporting evidence from explicit content mapping and coverage analysis.

Risks to external validity include sample bias, where the test is validated on a narrow or unrepresentative population (Henrich et al., 2010), and unrealistic testing conditions, which may not reflect real-world scenarios (Donald T. Campbell, 1963). Temporal variability and interaction effects can also distort results if performance shifts over time or due to specific environmental factors (Anderson et al., 2023). These issues can be investigated through stress testing, A/B testing, transfer testing, and population-stratified assessments, with evidence from performance comparisons across different conditions and sensitivity analyses.

Risks to criterion validity include criterion contamination, where extraneous factors influence assessment, and criterion deficiency, where relevant aspects of performance are omitted (Brogden and Taylor, 1950; Austin and Villanova, 1992). Restricted range limits the ability to detect meaningful relationships if all scores are too similar. These issues can be addressed through real-world longitudinal studies, validated criterion studies, and behavioral testing, with evidence from correlations with gold-standard benchmarks and predictions of real-world utility.

Risks to construct validity can come from structural, convergent, and discriminant validity risks. Structural validity is compromised by poor factor structure, where test items fail to group in expected ways (Clark and Watson, 1995; Elhami Athar, 2023), and complex measurement range, where constructs are not well captured across different levels of ability (Messick, 1995). Convergent validity can suffer from high measurement error (Cheung et al., 2024), which reduces reliability, while discriminant validity can be compromised by construct overlap, where different abilities are not clearly distinguished (Shaffer et al., 2016). These risks can be investigated using hypothesis testing, factor modeling, and benchmark suites, with supporting evidence from item-test correlations and demonstrated non-significant overlap with unrelated constructs.

Risks to consequential validity include bias and fairness issues, where results systematically disadvantage certain groups (Messick, 1995; Randall, 2023). While bias and fairness can themselves be constructs of interest, they are also important to consider in any measurement. Further, unintended incentives can distort behavior if assessment criteria encourage gaming rather than genuine learning (Nichols and Berliner, 2007). Policy consequences may emerge if flawed assessments influence high-stakes decisions. These risks can be assessed through anticipatory ethics methods (Umbrello et al., 2023), societal impact audits, and ethical stress testing, with evidence from stakeholder feedback, improvements in fairness and reliability, and documented real-world impacts.

While this framework highlights key risks and mitigation strategies, additional risks may arise in different contexts, necessitating continuous assessment and refinement.

5 Conclusion

Historically, AI evaluation has been benchmark-driven, focusing on narrow technical progress without critically assessing the validity of broader claims. This was fine in a regime where generative AI was a primarily research endeavor with less far-reaching consequences. However, as general-purpose generative AI systems continue to emerge, these traditional evaluation norms fail to predict real-world utility and risk irresponsible deployment and decision-making for AI systems. To address this, we enumerate five key forms of validity—content, external, criterion, construct, and consequential validity—each playing a critical role in determining whether a measurement and evaluation truly substantiates a given claim.

A fundamental challenge in AI evaluation is the conceptual gap between measured performance and real-world capability. Our claim-centered validity framework systematically bridges this gap, ensuring that AI assessments are rigorous, contextually appropriate, and scientifically robust. By explicitly mapping the relationship between measurements, evaluations, and the claims they are used to support, this framework prevents the overgeneralization of evaluation results and promotes a more accurate understanding of AI capabilities.

This work serves as a call for greater scientific rigor in AI evaluation, emphasizing that AI assessments must be claim-aware, evidence-driven, and methodologically sound. Whether formally articulated or not, the relationship between measurements, evaluations, and the claims they aim to validate implies an underlying nomological network—a structured web of relationships between constructs and criteria. However, the lack of explicit articulation of these networks has hindered progress in systematically aligning measurements and evaluations with AI capabilities. Without a deeper understanding of these relationships, evaluations risk misrepresenting AI capabilities, leading to flawed conclusions and misguided applications.

By adopting a principled approach to AI evaluation validity, we can move beyond surface-level benchmarking and towards a more scientifically grounded, transparent, and reliable assessment of AI systems. This shift is essential for ensuring that AI technologies are developed, evaluated, and deployed responsibly, ultimately fostering more trustworthy AI systems that align with real-world needs and expectations.

This work offers a theoretical foundation for validity-centered AI evaluation, setting the stage for more practical applications and empirical investigations. By clarifying how measurements, evaluations, and claims interact—and by emphasizing nomological networks—we provide a framework that can be adapted to diverse AI domains and tasks. We have enumerated some existing tools for probing validity. Future research includes a focus on operationalizing this framework in high-stakes contexts and on building robust nomological networks that systematically map AI constructs to measurable variables. Such endeavors will help ensure that evaluations not only gauge narrow performance but also yield trustworthy insights into real-world utility and risk.

References

Robert Adcock and David Collier. Measurement validity: A shared standard for qualitative and quantitative research. *Am. Polit. Sci. Rev.*, 95(3):529–546, September 2001.

AERA. Standards for educational & psychological testing (2014 edition). <https://www.aera.net/publications/books/>

- [standards-for-educational-psychological-testing-2014-edition](#), 2014. Accessed: 2025-3-10.
- Shivani Agarwal and Dan Roth. Learning a sparse representation for object detection. In *Computer Vision — ECCV 2002*, Lecture notes in computer science, pages 113–127. Springer Berlin Heidelberg, Berlin, Heidelberg, 2002.
- Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C Lawrence Zitnick, Dhruv Batra, and Devi Parikh. VQA: Visual question answering. *arXiv preprint*, May 2015.
- Ahmed Alaa, Thomas Hartvigsen, Niloufar Golchini, Shiladitya Dutta, Frances Dean, Inioluwa Deborah Raji, and Travis Zack. Medical large language model benchmarks should prioritize construct validity. *arXiv preprint arXiv:2503.10694*, 2025.
- D I Andonov, B Ulm, M Graessner, A Podtschaske, M Blobner, B Jungwirth, and S M Kagerbauer. Impact of the covid-19 pandemic on the performance of machine learning algorithms for predicting perioperative mortality. *BMC Med. Inform. Decis. Mak.*, 23(1):67, April 2023.
- Anthropic. Claude 3.5 sonnet. <https://www.anthropic.com/news/claude-3-5-sonnet>, 2024. [Accessed 05-04-2025].
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint*, July 2019.
- James T Austin and Peter Villanova. The criterion problem: 1917–1992. *J. Appl. Psychol.*, 77(6): 836–874, December 1992.
- Xuechunzi Bai, Angelina Wang, Ilia Sucholutsky, and Thomas L Griffiths. Explicitly unbiased large language models still form biased associations. *Proceedings of the National Academy of Sciences*, 122(8):e2416228122, 2025.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and machine learning: Limitations and opportunities*. MIT press, 2023.
- Samuel J Bell, Diane Bouchacourt, and Levent Sagun. Reassessing the validity of spurious correlations benchmarks. *arXiv preprint arXiv:2409.04188*, 2024.
- Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. Stereotyping norwegian salmon: An inventory of pitfalls in fairness benchmark datasets. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1004–1015, Stroudsburg, PA, USA, 2021. Association for Computational Linguistics.
- Avrim Blum and Moritz Hardt. The ladder: A reliable leaderboard for machine learning competitions. *arXiv preprint*, February 2015.
- Denny Borsboom, Gideon J Mellenbergh, and Jaap van Heerden. The concept of validity. *Psychol. Rev.*, 111(4):1061–1071, October 2004.
- Nick Bostrom, Allan Dafoe, and Carrick Flynn. Public policy and superintelligent AI: A vector field approach. In *Ethics of Artificial Intelligence*, pages 293–326. Oxford University Press New York, September 2020.

- Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated corpus for learning natural language inference. *arXiv preprint*, August 2015.
- Hubert E Brogden and Erwin K Taylor. The theory and classification of criterion bias. *Educ. Psychol. Meas.*, 10(2):159–183, July 1950.
- Brian Buntz. Eureka 2.0: AI is beginning to ace grad-level science, but can you trust it? <https://www.rdworldonline.com/eureka-2-0-ai-is-beginning-to-ace-grad-level-science-but-can-you-trust-it/>, February 2025. Accessed: 2025-3-10.
- D T Campbell and D W Fiske. Convergent and discriminant validation by the multitrait-multimethod matrix. *Psychol. Bull.*, 56(2):81–105, March 1959.
- D T Campbell and J C Stanley. *Experimental and Quasi-Experimental Designs for Research*. Ravenio Books, 2015.
- Gary L Canivez, Marley W Watkins, and Stefan C Dombrowski. Structural validity of the wechsler intelligence scale for Children-Fifth edition: Confirmatory factor analyses with the 16 primary and secondary subtests. *Psychol. Assess.*, 29(4):458–472, April 2017.
- Gordon W Cheung, Helena D Cooper-Thomas, Rebecca S Lau, and Linda C Wang. Reporting reliability, convergent and discriminant validity with structural equation modeling: A review and best-practice recommendations. *Asia Pac. J. Manag.*, 41(2):745–783, June 2024.
- Alexandra Chouldechova, Chad Atalla, Solon Barocas, A Feder Cooper, Emily Corvi, P Alex Dow, Jean Garcia-Gathright, Nicholas Pangakis, Stefanie Reed, Emily Sheng, Dan Vann, Matthew Vogel, Hannah Washington, and Hanna Wallach. A shared standard for valid measurement of generative AI systems’ capabilities, risks, and impacts. *arXiv preprint*, December 2024.
- Gregory J Cizek and Michael B Bunch. *Standard setting: A guide to establishing and evaluating performance standards on tests*. SAGE Publications Ltd, 2007.
- Lee Anna Clark and David Watson. Constructing validity: Basic issues in objective scale development. *Psychol. Assess.*, 7(3):309–319, September 1995.
- Alexis Conneau and Douwe Kiela. SentEval: An evaluation toolkit for universal sentence representations. *arXiv preprint*, March 2018.
- Paper Copilot. NeurIPS 2024 statistics: Datasets & benchmarks track. <https://papercopilot.com/statistics/neurips-statistics/neurips-2024-statistics-datasets-benchmarks-track/>, April 2024. Accessed: 2025-3-12.
- Amanda Coston, Anna Kawakami, Haiyi Zhu, Ken Holstein, and Hoda Heidari. A validity perspective on evaluating the justified use of data-driven decision-making algorithms. In *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pages 690–704. IEEE, February 2023.
- L J Cronbach and P E Meehl. Construct validity in psychological tests. *Psychol. Bull.*, 52(4):281–302, July 1955.

- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. IEEE, June 2009.
- Julian Stanley Donald T. Campbell. Experimental and Quasi-Experimental designs for research. *Cengage Learning*, 1963.
- David Donoho. Data science at the singularity. *arXiv preprint*, October 2023.
- Mojtaba Elhami Athar. The pitfalls of untested assumptions and unwarranted/oversimplistic interpretation of cultural phenomenon: a commentary on sajjadi et al. (2023). *Front. Psychol.*, 14: 1248246, September 2023.
- Mark Everingham, Luc Van Gool, Christopher K I Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (VOC) challenge. *Int. J. Comput. Vis.*, 88(2):303–338, June 2010.
- Li Fei-Fei, R Fergus, and P Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *2004 Conference on Computer Vision and Pattern Recognition Workshop*, pages 178–178. IEEE, 2005.
- Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani, Claire Barale, Robert McHardy, Joshua Harris, Jean Kaddour, Emile van Krieken, and Pasquale Minervini. Are we done with MMLU? *arXiv preprint*, June 2024.
- Judy Wawira Gichoya, Imon Banerjee, Ananth Reddy Bhimireddy, John L Burns, Leo Anthony Celi, Li-Ching Chen, Ramon Correa, Natalie Dullerud, Marzyeh Ghassemi, Shih-Cheng Huang, et al. Ai recognition of patient race in medical imaging: a modelling study. *The Lancet Digital Health*, 4(6):e406–e414, 2022.
- Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego Chicharro, Evan Chen, Alex Gunning, Caroline Falkman Olsson, Jean-Stanislas Denain, Anson Ho, Emily de Oliveira Santos, Olli Järvinemi, Matthew Barnett, Robert Sandler, Matej Vrzala, Jaime Sevilla, Qiuyu Ren, Elizabeth Pratt, Lionel Levine, Grant Barkley, Natalie Stewart, Bogdan Grechuk, Tetiana Grechuk, Shreeparnav Varma Enugandla, and Mark Wildon. FrontierMath: A benchmark for evaluating advanced mathematical reasoning in AI. *arXiv preprint*, November 2024.
- Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. *arXiv preprint*, July 2020.
- Vipul Gupta, David Pantoja, Candace Ross, Adina Williams, and Megan Ung. Changing answer order can decrease MMLU accuracy. *arXiv preprint*, June 2024.
- Moritz Hardt and Benjamin Recht. Patterns, predictions, and actions: A story about machine learning. *arXiv preprint*, February 2021.
- Amelia Hardy, Anka Reuel, Kiana Jafari Meimandi, Lisa Soder, Allie Griffith, Dylan M Asmar, Sanmi Koyejo, Michael S Bernstein, and Mykel J Kochenderfer. More than marketing? on the information value of AI benchmarks for practitioners. *arXiv preprint*, December 2024.

- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*, September 2020.
- Joseph Henrich, Steven J Heine, and Ara Norenzayan. The weirdest people in the world? *Behav. Brain Sci.*, 33(2-3):61–83; discussion 83–135, June 2010.
- David J Hughes. *Psychometric validity: Establishing the accuracy and appropriateness of psychometric measures*. Wiley Online Library, 2018.
- Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A Mong, Safwan S Halabi, Jesse K Sandberg, Ricky Jones, David B Larson, Curtis P Langlotz, Bhavik N Patel, Matthew P Lungren, and Andrew Y Ng. CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison. January 2019.
- Abigail Z Jacobs and Hanna Wallach. Measurement and fairness. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, New York, NY, USA, March 2021. ACM.
- Jennifer L Jennings and Jonathan Marc Bearak. “teaching to the test” in the NCLB era: How test predictability affects our understanding of student performance. *Educ. Res.*, 43(8):381–389, November 2014.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421, 2021.
- Kushal Kaffe and Christopher Kanan. An analysis of visual question answering algorithms. *arXiv preprint*, March 2017.
- Jae-On Kim and Charles W Mueller. *Factor analysis: Statistical methods and practical issues*. Quantitative Applications in the Social Sciences. SAGE Publications, Thousand Oaks, CA, February 1979.
- Jennifer L Kobrin, Brian F Patterson, Emily J Shaw, Krista D Mattern, and Sandra M Barbuti. Validity of the SAT® for predicting First-Year college grade point average. research report no. 2008-5. *College Board*, 2008.
- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton A Earnshaw, Imran S Haque, Sara Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. WILDS: A benchmark of in-the-wild distribution shifts. *arXiv preprint*, December 2020.
- Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better ImageNet models transfer better? *arXiv preprint*, May 2018.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.(2009). *University of Toronto*, 18268744, September 2009.
- T S Kuhn. *The structure of scientific revolutions*, volume 962. Chicago: University of Chicago press, 1997.

- C H Lawshe. A QUANTITATIVE APPROACH TO CONTENT VALIDITY¹. *Pers. Psychol.*, 28 (4):563–575, December 1975.
- Y Lecun, L Bottou, Y Bengio, and P Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE Inst. Electr. Electron. Eng.*, 86(11):2278–2324, 1998.
- Q Vera Liao and Ziang Xiao. Rethinking model evaluation as narrowing the socio-technical gap. *arXiv preprint*, May 2023.
- Weng Marc Lim. A typology of validity: content, face, convergent, discriminant, nomological and predictive validity. *Journal of Trade Science*, 12(3):155–179, September 2024.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C Lawrence Zitnick, and Piotr Dollár. Microsoft COCO: Common objects in context. *arXiv preprint*, May 2014.
- Yu Lu Liu, Su Lin Blodgett, Jackie Chi Kit Cheung, Q Vera Liao, Alexandra Olteanu, and Ziang Xiao. ECBD: Evidence-Centered benchmark design for NLP. *arXiv preprint*, June 2024.
- David Lopez-Paz, Robert Nishihara, Soumith Chintala, Bernhard Schölkopf, and Léon Bottou. Discovering causal signals in images. *arXiv preprint*, May 2016.
- Samuel Messick. Validity of psychological assessment: Validation of inferences from persons’ responses and performances as scientific inquiry into score meaning. *Am. Psychol.*, 50(9):741–749, September 1995.
- Samuel Messick. Test validity: A matter of consequence. *Soc. Indic. Res.*, 45(1/3):35–44, November 1998.
- Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. GSM-symbolic: Understanding the limitations of mathematical reasoning in large language models. *arXiv preprint*, October 2024.
- C I Mosier. A critical examination of the concepts of face validity. *Educ. Psychol. Meas.*, 7(2): 191–205, July 1947.
- Sharon L Nichols and David C Berliner. Collateral damage: How High-Stakes testing corrupts america’s schools. In *Harvard Education Press*. Harvard Education Press. 8 Story Street First Floor, Cambridge, MA 02138. Tel: 888-437-1437; Tel: 617-495-3432; Fax: 978-348-1233; e-mail: hepg@harvard.edu; Web site: http://hepg.org/hepg-home/home, March 2007.
- Mohammad Sadegh Norouzzadeh, Anh Nguyen, Margaret Kosmala, Alexandra Swanson, Meredith S Palmer, Craig Packer, and Jeff Clune. Automatically identifying, counting, and describing wild animals in camera-trap images with deep learning. *Proc. Natl. Acad. Sci. U. S. A.*, 115(25): E5716–E5725, June 2018.
- Will Orr and Edward B Kang. AI as a sport: On the competitive epistemologies of benchmarking. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, New York, NY, USA, June 2024. ACM.
- Simon Ott, Adriano Barbosa-Silva, Kathrin Blagec, Jan Brauner, and Matthias Samwald. Mapping global dynamics of benchmark creation and saturation in artificial intelligence. *arXiv preprint*, March 2022.

- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models. *arXiv preprint arXiv:2202.03286*, 2022.
- Inioluwa Deborah Raji, Emily M Bender, Amandalynne Paullada, Emily Denton, and Alex Hanna. AI and the everything in the whole wide world benchmark. *arXiv preprint*, November 2021.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. *arXiv preprint*, June 2016.
- Jennifer Randall. It ain’t near ‘bout fair: Re-envisioning the bias and sensitivity review process from a justice-oriented antiracist perspective. *Educ. Assess.*, 28(2):68–82, April 2023.
- Benjamin Recht. The mechanics of frictionless reproducibility. *Harvard Data Science Review*, 6(1), 2024.
- Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? *PMLR*, pages 5389–5400, May 2019.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. GPQA: A Graduate-Level Google-Proof Q&A benchmark. *arXiv preprint*, November 2023.
- Richard Ren, Steven Basart, Adam Khoja, Alice Gatti, Long Phan, Xuwang Yin, Mantas Mazeika, Alexander Pan, Gabriel Mukobi, Ryan H Kim, Stephen Fitz, and Dan Hendrycks. Safetywashing: Do AI safety benchmarks actually measure safety progress? *arXiv preprint*, July 2024.
- Anka Reuel, Amelia Hardy, Chandler Smith, Max Lamparth, Malcolm Hardy, and Mykel J Kochenderfer. BetterBench: Assessing AI benchmarks, uncovering issues, and establishing best practices. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, November 2024.
- Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. The risks of invariant risk minimization. *Int Conf Learn Represent*, abs/2010.05761, October 2020.
- Yangjun Ruan, Chris J Maddison, and Tatsunori Hashimoto. Observational scaling laws and the predictability of language model performance. *arXiv preprint*, May 2024.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C Berg, and Li Fei-Fei. ImageNet large scale visual recognition challenge. *arXiv preprint*, September 2014.
- Olawale Salaudeen and Moritz Hardt. ImageNot: A contrast with ImageNet preserves model rankings. *arXiv [cs.LG]*, 2024.
- Olawale Salaudeen and Sanmi Koyejo. Causally inspired regularization enables domain general representations. In *International Conference on Artificial Intelligence and Statistics*, pages 3124–3132. PMLR, April 2024.
- Olawale Salaudeen, Nicole Chiou, Shiny Weng, and Sanmi Koyejo. Are domain generalization benchmarks with accuracy on the line misspecified? *arXiv preprint arXiv:2504.00186*, 2025.

- Olawale Elijah Salaudeen and Oluwasanmi O Koyejo. Exploiting causal chains for domain generalization. *NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and Applications*, 2021.
- Michael Saxon, Ari Holtzman, Peter West, William Yang Wang, and Naomi Saphra. Benchmarks as microscopes: A call for model metrology. *arXiv preprint*, July 2024.
- Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- Jonathan A Shaffer, David DeGeest, and Andrew Li. Tackling the problem of construct proliferation: A guide to assessing the discriminant validity of conceptually related constructs. *Organ. Res. Methods*, 19(1):80–110, January 2016.
- Lorrie A Shepard. Chapter 9: Evaluating test validity. *Rev. Res. Educ.*, 19(1):405–450, January 1993.
- Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *J. Stat. Plan. Inference*, 90(2):227–244, October 2000.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askeel, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubakaran, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D Manning, Christopher Potts, Cindy Ramirez, Clara E Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh

Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B Tenenbaum, Joshua S Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón, Luke Metz, Lütü Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Farooqi, Manaal Faruqi, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael A Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michal Swkedrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimeo Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun Peng, Nathan A Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millièvre, Rhythm Garg, Richard Barnes, Rif A Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R Bowman, Samuel S Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima, Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T Piantadosi, Stuart M Shieber, Summer Misherghe, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen,

- Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint*, June 2022.
- Arjun Subramonian, Xingdi Yuan, Hal Daumé, III, and Su Lin Blodgett. It takes two to tango: Navigating conceptualizations of NLP tasks and measurements of performance. *arXiv preprint*, May 2023.
- Robert L Thorndike. *Personnel selection; test and measurement techniques*. J. Wiley, New York, 1949.
- Steven Umbrello, Michael J Bernstein, Pieter E Vermaas, Anaïs Resseguier, Gustavo Gonzalez, Andrea Porcari, Alexei Grinbaum, and Laurynas Adomaitis. From speculation to reality: Enhancing anticipatory ethics for emerging technologies (ATE) in practice. *Technol. Soc.*, 74(102325): 102325, August 2023.
- V N Vapnik and A Ya Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory Probab. Appl.*, 16(2):264–280, January 1971.
- Hanna Wallach, Meera Desai, A Feder Cooper, Angelina Wang, Chad Atalla, Solon Barocas, Su Lin Blodgett, Alexandra Chouldechova, Emily Corvi, P Alex Dow, Jean Garcia-Gathright, Alexandra Olteanu, Nicholas Pangakis, Stefanie Reed, Emily Sheng, Dan Vann, Jennifer Wortman Vaughan, Matthew Vogel, Hannah Washington, and Abigail Z Jacobs. Position: Evaluating generative AI systems is a social science measurement challenge. *arXiv preprint*, February 2025.
- Yuxuan Wan, Wenxuan Wang, Yiliu Yang, Youliang Yuan, Jen-Tse Huang, Pinjia He, Wenxiang Jiao, and Michael R Lyu. LogicAsker: Evaluating and improving the logical reasoning ability of large language models. *arXiv preprint*, January 2024.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint*, April 2018.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. *arXiv preprint*, May 2019.
- Angelina Wang and Olga Russakovsky. Overwriting pretrained bias with finetuning data. *arXiv preprint*, March 2023.
- Angelina Wang, Aaron Hertzmann, and Olga Russakovsky. Benchmark suites instead of leaderboards for evaluating AI fairness. *Patterns (N. Y.)*, 5(11):101080, November 2024.
- Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim, Jacob Andreas, and Yoon Kim. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks. *arXiv preprint*, July 2023.
- Kai Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander Madry. Noise or signal: The role of image backgrounds in object recognition. *arXiv preprint*, June 2020.
- Ziang Xiao, Susu Zhang, Vivian Lai, and Q Vera Liao. Evaluating evaluation metrics: A framework for analyzing NLG evaluation metrics using measurement theory. *arXiv preprint*, May 2023.

Ziang Xiao, Wesley Hanwen Deng, Michelle S Lam, Motahhare Eslami, Juho Kim, Mina Lee, and Q Vera Liao. Human-centered evaluation and auditing of language models. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pages 1–6, New York, NY, USA, May 2024. ACM.

Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav Raja, Charlotte Zhuang, Dylan Slack, Qin Lyu, Sean Hendryx, Russell Kaplan, Michele Lunati, and Summer Yue. A careful examination of large language model performance on grade school arithmetic. *arXiv preprint*, May 2024.

Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ADE20K dataset. *arXiv preprint*, August 2016.

Wangchunshu Zhou, Yan Zeng, Shizhe Diao, and Xinsong Zhang. VLUe: A multi-task benchmark for evaluating vision-language models. *arXiv preprint*, May 2022.

Appendix Table of Contents

A Evidence of Validity	28
B Validity	32
C The (Co)Evolution of evaluations and claims	34
C.1 Vision	34
C.2 Language	36
D Case Studies	39
D.1 GPQA	41
D.2 MMLU	46
D.3 ImageNet	51
D.4 MedQA	56
D.5 SQuAD	61

A Evidence of Validity

Table 4: Common risks to validity, investigation tools, and evidence exemplars.

Validity	Common risks	Investigation Tools	Evidence Exemplar
Content Validity	<input type="checkbox"/> Coverage deficiency <input type="checkbox"/> Construct irrelevance <input type="checkbox"/> Imbalanced mixture of content	<input type="checkbox"/> Expert review <input type="checkbox"/> Red-teaming / adversarially designed evaluations <input type="checkbox"/> Synthetic data generation or edge cases	<input type="checkbox"/> Documentation of how test items comprehensively cover the construct <input type="checkbox"/> Explicit mapping of test content to abstract frameworks or industry standards <input type="checkbox"/> Coverage analysis
Criterion Validity	Predictive and Concurrent Validity <input type="checkbox"/> Criterion contamination <input type="checkbox"/> Criterion deficiency <input type="checkbox"/> Restricted range <input type="checkbox"/> Temporal/other shifts	<input type="checkbox"/> Real-world longitudinal studies <input type="checkbox"/> Real-world behavioral testing <input type="checkbox"/> Scaling-law predictive models <input type="checkbox"/> Validated criterion studies <input type="checkbox"/> Periodic post-deployment testing	<input type="checkbox"/> Correlation with an existing validated benchmark or gold standard <input type="checkbox"/> Evidence that higher scores in evaluation metrics predict real-world utility

Table continues on the next page

Table 4 continued from previous page

Validity	Common risks	Investigation Tools	Evidence Exemplar
Construct Validity	Structural:		
	<input type="checkbox"/> Rank deficiency	<input type="checkbox"/> Theory building and hypothesis testing	<input type="checkbox"/> Observed changes in test performance under controlled conditions
	<input type="checkbox"/> Poor factor structure	<input type="checkbox"/> Factor modeling	<input type="checkbox"/> Item-test correlations
	<input type="checkbox"/> Item interdependence	<input type="checkbox"/> Studies of process	<input type="checkbox"/> Emergent substructures in model behavior
	<input type="checkbox"/> Response format bias		
	<input type="checkbox"/> Complex measurement range		
	Convergent:		
	<input type="checkbox"/> Irrelevant or weakly related evaluations	<input type="checkbox"/> Benchmark suites for a construct (e.g., reasoning)	<input type="checkbox"/> High correlation with other measures that assess the same construct
	<input type="checkbox"/> High measurement error in scoring	<input type="checkbox"/> Representation probing (e.g., causal mediation analysis of embeddings)	<input type="checkbox"/> Empirical clustering of model behaviors that align with constructs
	<input type="checkbox"/> Restricted range (ceiling/floor effects)		
	<input type="checkbox"/> Confounding (e.g., memorization, format)		

Table continues on the next page

Table 4 continued from previous page

Validity	Common risks	Investigation Tools	Evidence Exemplar
	Discriminant:	<input type="checkbox"/> Orthogonal datasets	<input type="checkbox"/> Low or non-significant correlation with measures of distinct constructs
	<input type="checkbox"/> Construct overlap	<input type="checkbox"/> Decomposable metrics	<input type="checkbox"/> Evidence that evaluation does not overlap with unrelated dimensions
	<input type="checkbox"/> Format-induced correlations		
External Validity	<input type="checkbox"/> Sample bias	<input type="checkbox"/> Red-teaming	<input type="checkbox"/> Performance comparisons across different populations, environments, or settings
	<input type="checkbox"/> Unrealistic testing conditions	<input type="checkbox"/> Stress testing	
	<input type="checkbox"/> Temporal variability	<input type="checkbox"/> A/B testing	<input type="checkbox"/> Sensitivity analysis showing consistent performance under varying conditions
	<input type="checkbox"/> Interaction effects	<input type="checkbox"/> Transfer testing	
	<input type="checkbox"/> Experimenter effects	<input type="checkbox"/> Population-stratified evaluations	<input type="checkbox"/> Independent replication of results in different contexts or regions
	<input type="checkbox"/> Task-specific bias		

Table continues on the next page

Table 4 continued from previous page

Validity	Common risks	Investigation Tools	Evidence Exemplar
Consequential Validity	<input type="checkbox"/> Bias / Fairness <input type="checkbox"/> Adaptive overfitting <input type="checkbox"/> Misuse of results <input type="checkbox"/> Unintended incentives <input type="checkbox"/> Policy and systematic consequences <input type="checkbox"/> Temporal and other shift	<input type="checkbox"/> Stakeholder interviews and feedback loops <input type="checkbox"/> Societal impact audits <input type="checkbox"/> Ethical stress testing <input type="checkbox"/> Stakeholder feedback	<input type="checkbox"/> Documented instances of evaluation-driven improvements in safety, reliability, and fairness <input type="checkbox"/> Impact studies

B Validity

Validity refers to the extent to which a test accurately measures what it is intended to measure. Validity has a rich history, originally developed in the context of drawing valid conclusions from tests—much like how we now aim to draw valid conclusions from AI evaluations. One of the earliest forms of validity is face validity, which refers to the extent to which a test appears to measure what it claims to, based on intuitive judgment. For instance, one may ask if symbolic regression from BigBench (Srivastava et al., 2022) even appears to measure reasoning. However, relying on face validity alone can be misleading. As Charles Mosier (Mosier, 1947) famously observed:

“This form [face validity] is also gratifying to the ego of the unwary test constructor. It implies that his knowledge and skill in the area of test construction are so great that he can unerringly design a test with the desired degree of effectiveness in predicting job success or in evaluating defined personality characteristics, and that he can do this so accurately that any further empirical verification is unnecessary. So strong is this ego complex that if statistical verification is sought and found lacking, the data represent something to be explained away by appeal to sampling errors or other convenient rationalization, rather than by scientific evidence which must be admitted into full consideration.”

A more structured form of validity emerged with content validity, which ensures that a test comprehensively covers all relevant aspects of the construct it aims to measure. For instance, one may ask if mathematical problem-solving benchmarks cover all relevant aspects of reasoning. Content validity is also typically assessed through expert judgment rather than statistical validation. Charles Lawshe (Lawshe, 1975) later formalized this concept with the Content Validity Ratio (CVR), a method for quantifying expert agreement on test content.

Moving toward empirical rigor, predictive validity assesses a test’s ability to forecast an outcome of interest, typically a future outcome. This concept, introduced by Robert Thorndike in the mid-20th century during the rise of standardized testing, became central to fields like educational assessment, employment testing, and aptitude measurement (Thorndike, 1949). For example, the predictive validity of SAT scores for college GPA or cognitive ability tests for job performance has led to their widespread use for other outcomes (Kobrin et al., 2008). In the context of AI evaluation, one may ask “does accuracy on IMO benchmarks predict accuracy in grading university math exams?”

While predictive validity is useful for assessing direct correlations between tests and outcomes, its limitations became apparent when evaluating theoretical abstract constructs rather than simple outcome-based predictions. In their seminal on construct validity, (Cronbach and Meehl, 1955) highlighted these limitations. For example, while SAT scores may predict GPA, they may not reliably measure intelligence, as GPA is influenced by grading biases and other factors. Recognizing the risks of relying solely on criterion-based validity, Cronbach and Meehl introduced construct validity, which assesses the extent to which a test truly captures the theoretical construct it purports to measure.

Two key sources of evidence necessary for construct validity introduced by Campbell and Fiske (1959) are (Campbell and Fiske, 1959):

- Convergent validity—the degree to which a test correlates with other measures of the same construct.
- Discriminant validity—the degree to which a test does not correlate with measures of unrelated constructs.

Implicitly, this framework also includes structural validity (Cronbach and Meehl, 1955; Messick, 1995), which examines whether a test’s internal structure aligns with the theoretical construct it is designed to measure. This is often assessed using factor analysis or other dimensionality evaluations.

Cronbach and Meehl categorize validity into three primary forms:

1. *Content validity*—ensuring a test comprehensively represents the concept it aims to measure.
2. *Criterion validity*—evaluating how well a test correlates with external measures, which include predictive and concurrent validity. Concurrent validity refers to a test’s agreement with a validated measure applied at the same time under the same conditions.
3. *Construct validity*—assessing the theoretical alignment between a test and its intended construct.

Beyond these core types, external validity refers to the extent to which a study’s findings can be generalized beyond its specific conditions. External validity examines whether results hold across different populations, settings, and time periods. Campbell and Stanley (Campbell and Stanley, 2015) were among the first to systematically define external validity, identifying factors like selection bias and situational specificity as risks to generalizability.

In response to Cronbach and Meehl’s framework, which emphasized the theoretical and statistical relationships between measures, (Messick, 1995, 1998) introduced consequential validity on the basis that validity is not just about measurement accuracy but also about the real-world impact of test interpretation and use.

We primarily follow Cronbach and Meehl’s classification of validity types. However, we adopt Messick’s view that validity extends beyond measurement properties. Thus, we incorporate consequential validity as an essential consideration in our framework. We, however, diverge from (Borsboom et al., 2004)’s view that validity is only a property of the test.

Borsboom (Borsboom et al., 2004) offers a different view: validity is a property of the test itself, and a test is valid if and only if it measures the construct it purports to measure. In this view, questions of use or consequence are orthogonal to validity; what matters is whether the test causally reflects variation in the construct. This perspective draws a clear boundary between measurement and interpretation, placing the burden of validity squarely on the psychometric relationship between construct and test score. While theoretically clean, this stance omits considerations critical to our context, namely, how test outputs are used to make decisions. We, therefore, depart from Borsboom’s definition, instead adopting a broader view in which validity also encompasses downstream consequences and use cases, particularly when evaluating AI systems deployed in high-stakes settings.

While these validity concepts were originally developed for psychological and educational testing, they provide a powerful lens for evaluating AI models. In the next section, we examine how these classical validity forms translate into the context of modern AI evaluation.

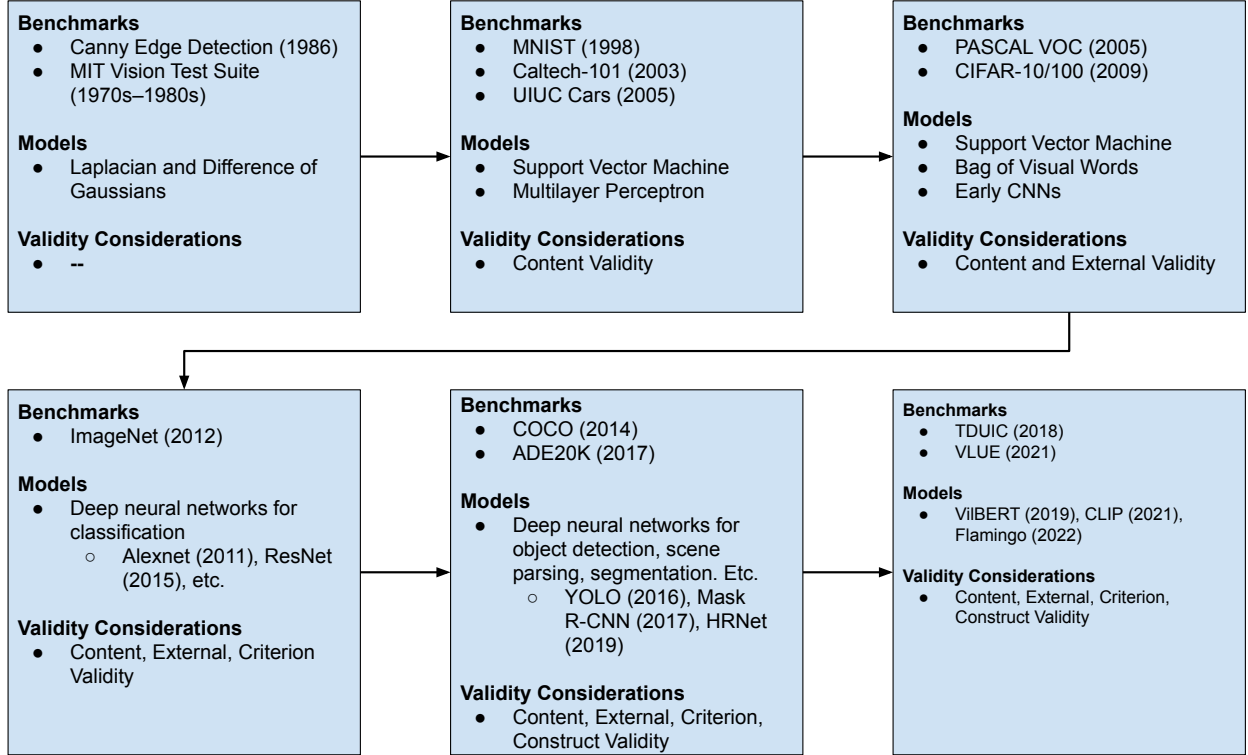


Figure 5: Coevolution of benchmarks, models, and the type of validity necessary for common conclusions for vision.

C The (Co)Evolution of evaluations and claims

C.1 Vision

The evolution of AI benchmarks has been closely tied to the kinds of conclusions researchers aimed to draw and the evidence available at the time—Figure 5. In the 1960s to 1980s, benchmarks were hyper-localized, focusing on narrowly defined technical tasks like edge detection and simple shape recognition. The goal was primarily technical exploration—improving algorithmic efficiency—so the scope of conclusions was very narrow and directly supported by the evaluations carried out.

In the 1990s, AI benchmarks became more structured and began incorporating more applied tasks. A notable example is MNIST (Lecun et al., 1998) for handwritten digit classification, which provided a standardized way to evaluate machine learning models. This trend continued into the early 2000s, with datasets such as UIUC Cars (Agarwal and Roth, 2002) for vehicle detection and Caltech-101 (2003) (Fei-Fei et al., 2005) for object recognition. While these benchmarks remained narrow in scope, they represented a step toward evaluating AI on more applied tasks, bridging the gap between theoretical research and practical applications. However, evaluations were still primarily designed for well-defined technical interests, with conclusions remaining local—focused on determining which techniques were most effective for the specific task being evaluated. During this period, researchers also became increasingly aware of content validity, recognizing that different datasets captured different aspects of classification tasks, which in turn influenced dataset design and evaluation methodologies.

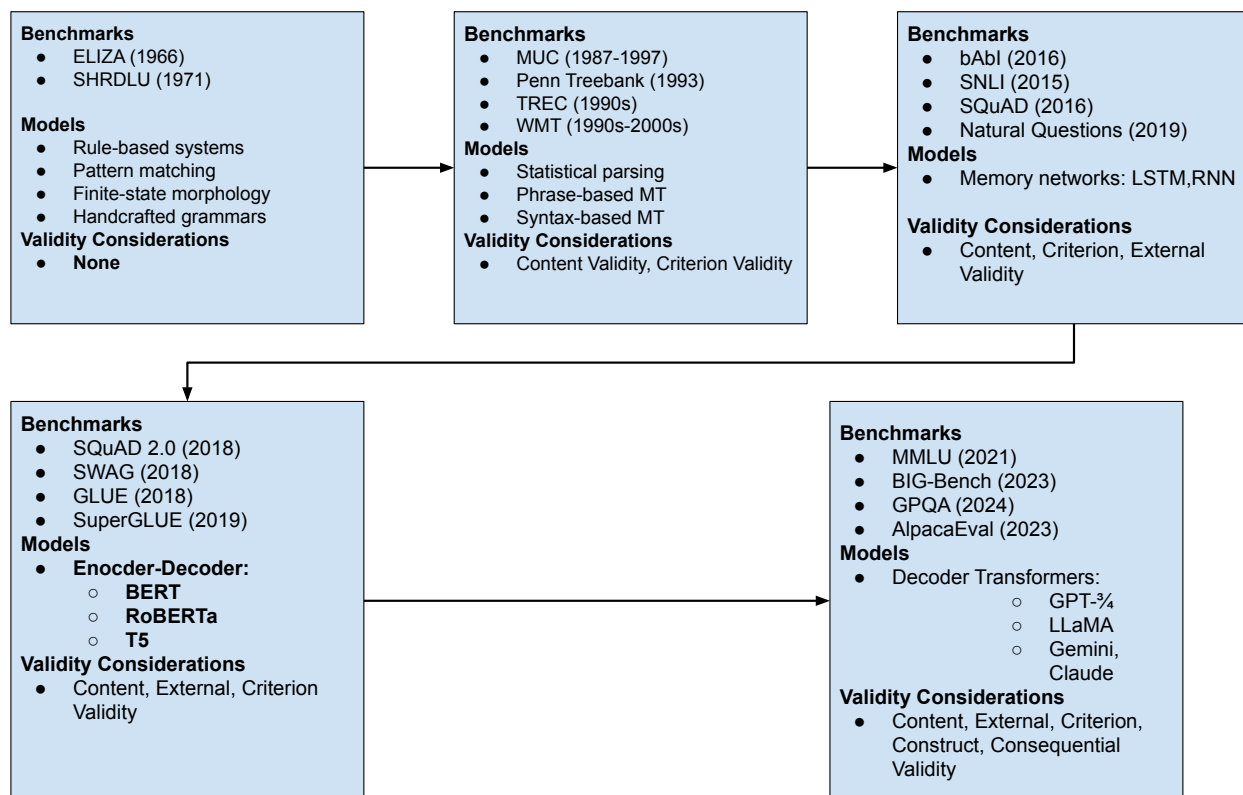


Figure 6: Coevolution of benchmarks, models, and the type of validity necessary for common conclusions for language.

By the mid-2000s, large-scale benchmarks such as PASCAL VOC (2007) (Everingham et al., 2010) introduced greater complexity, expanding evaluation beyond simple classification tasks. Later, in the late 2000s, CIFAR-10 and CIFAR-100 (Krizhevsky and Hinton, 2009) further pushed the field toward standardized comparisons in object recognition. The shift from classification to more structured tasks like object detection and segmentation moved benchmarks toward broader contexts, with a stronger emphasis on generalization. Content validity and external validity became increasingly relevant as researchers began evaluating models across multiple datasets and questioning whether benchmark performance was a meaningful proxy for real-world vision tasks.

During this period, criterion validity also gained prominence, as benchmark results were increasingly used to compare models in ways that suggested performance rankings carried external significance. However, construct validity remained largely unexplored—models were evaluated based on their outputs rather than on the reasoning processes behind their decisions. As a result, while evaluations became more sophisticated, they remained focused on performance metrics rather than deeper insights into model behavior. By this stage, the focus of AI evaluation began shifting from isolated dataset-specific improvements to broader claims about model robustness and transferability across different domains.

The 2010s marked a turning point with the ImageNet revolution. The introduction of ImageNet (Deng et al., 2009) and the ILSVRC (Russakovsky et al., 2014) competition (2010) provided large-scale, diverse, and complex benchmarks that dramatically reshaped AI research. During the early 2010s, the focus remained on improving accuracy in image classification and object detection.

However, by the mid-2010s, AI evaluation expanded beyond leaderboards to real-world applications, particularly in medical imaging and autonomous driving. Researchers increasingly recognized the importance of content validity and external validity, leading to the widespread practice of testing models across multiple datasets to assess robustness.

As benchmark results gained influence, criterion validity became central—accuracy on ImageNet was frequently treated as a proxy for general AI capabilities in vision. However, construct validity remained largely unaddressed in the early years. By the mid-2010s, early concerns emerged as researchers identified shortcut learning, adversarial vulnerabilities, and spurious correlations, leading to growing interest in understanding how models made decisions beyond raw accuracy. The rise of segmentation (COCO (Lin et al., 2014), ADE20K (Zhou et al., 2016)) and video analysis benchmarks (Kinetics, AVA) reflected an effort to capture more complex real-world tasks, but fundamental concerns about model robustness and bias persisted.

In the 2020s, the rise of multimodal and foundation models introduced even greater evaluation challenges. Benchmarks such as VQA (Agrawal et al., 2015), VLUE (Zhou et al., 2022), and TDIUC (Kafle and Kanan, 2017) attempted to assess multimodal reasoning, but defining what these benchmarks truly measured became increasingly difficult. Construct validity became a major concern as researchers debated whether these benchmarks genuinely assessed reasoning and understanding or merely exposed a model’s ability to exploit statistical correlations in large datasets. Unlike earlier benchmarks, which primarily focused on accuracy, modern benchmarks aim to evaluate the latent properties of AI systems. However, fundamental questions about the validity of these evaluations remain unresolved, particularly in assessing generalization, robustness, and true reasoning ability.

Across these decades, benchmarks evolved alongside the conclusions researchers sought to make. Early benchmarks required little discussion of validity because they were purely technical exercises. As AI models became more ambitious and claims about their capabilities expanded, benchmarks had to keep up—introducing concerns about content, external, and criterion validity. More recently, as AI systems move toward multimodal reasoning and foundation models, discussions of construct validity have become central. As models grow in complexity, the challenge is no longer just about designing better benchmarks—it’s about defining what those benchmarks are actually supposed to measure in the first place.

C.2 Language

Language model benchmarks have seen an evolution from focusing on primarily basic questions of criterion validity against human performance to more nuanced considerations of other validity in more recent years—Figure 6. In the Blocks World Era (1960s-1980s), NLP evaluation was primarily qualitative and demonstration-based, lacking standardized metrics entirely. Systems like ELIZA (1966) and SHRDLU (1971) were evaluated through anecdotal observations of how users interacted with them in highly constrained environments. ELIZA simulated a psychotherapist using simple pattern matching, while SHRDLU operated in a “blocks world” where users could issue commands to manipulate virtual objects. Validity considerations during this era were minimal and largely implicit. Content validity was severely limited by extremely narrow domains, criterion validity was nonexistent without standardized measurements and construct validity wasn’t addressed as researchers weren’t attempting to measure specific capabilities like “reasoning” or “understanding.” External validity was particularly weak as systems couldn’t generalize beyond their constrained environments. Success was measured simply by the system’s ability to maintain seemingly intelligent

conversations or follow instructions rather than through quantitative performance metrics or validity criteria. The North Star Era (1990s-2000s) marked a paradigm shift toward empirical evaluation with standardized benchmarks inspired by information retrieval traditions, where benchmarks with quantitative metrics and clearly defined train, validation, and test split gave the field a proverbial “North Star” to aim towards. Initiatives like the Message Understanding Conferences (MUC) and the Penn Treebank established common datasets, clearly defined tasks, and metrics such as precision, recall, and F-score for comparing systems. This era introduced the first rigorous validity considerations, though still narrow in scope. Benchmarks like TREC and WMT established improved criterion validity through standardized metrics that allowed consistent measurement across systems and time. Content validity improved but remained limited to specific linguistic tasks. Nascent construct validity concerns emerged as researchers began considering what abilities their tasks were actually measuring. However, external validity remained largely unaddressed as benchmarks weren’t designed to generalize beyond their specific contexts. Consequential validity still wasn’t a major consideration, as NLP applications weren’t yet widely deployed with significant societal impact.

In the early 2010s, many language benchmarks, such as SQuAD (Rajpurkar et al., 2016) and SNLI (Bowman et al., 2015), focused on individual tasks such as reading comprehension or natural language claims such as entailment or contradiction. The primary focus was on establishing baseline comparisons against human performance to create criterion validity for the benchmarks. However, such benchmarks had limitations to other aspects such as content validity due to limited focus on specific linguistic tasks and face validity due to narrow objectives and methods used to solve the task (both SQuAD and SNLI can be cast as relatively simple classification problems for which we can measure a gold standard of correctness). Other validity types were not heavily considered at this time.


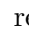

In the mid to late 2010s, the field began to focus more on multi-task evaluation, which was represented by benchmarks such as GLUE (Wang et al., 2018) and SentEval (Conneau and Kiela, 2018). During this time, emerging validity concerns became prominent. More sophisticated human baselines were required to maintain criterion validity, and broader task coverage led to great content validity. However, concerns about the underlying mechanisms that could explain performance began to emerge, which reflects early concerns about construct validity.












































































In the late 2010s there were key changes in language model evaluation. Benchmarks like SuperGLUE (Wang et al., 2019) aimed to resolve validity concerns with rigorous multi-annotator baselines, broader task selection, more attention to the demographics of annotators, and the first considerations of social impact and gaming. However, the lack of structural validity evidence and external validation remained as challenges. There were also few analyses of convergent/discriminant validity in studies.

The 2020s marked a shift toward comprehensive knowledge evaluation with benchmarks like MMLU (Hendrycks et al., 2020), reflecting a growing recognition that language models were advancing beyond narrow linguistic tasks to broader knowledge and reasoning capabilities. MMLU introduced several innovations in validity considerations: it established expert-level performance as the criterion validity benchmark rather than average human performance, expanded content validity through coverage of 57 subjects across multiple domains, and highlighted crucial external validity concerns through studies showing sensitivity to answer ordering and other conditions that should not have an effect on the downstream performance for an “intelligent” agent (as measured with respect to an expert). The evolution of MMLU reflects broader trends in the field’s approach to validity. Earlier benchmarks like SQuAD primarily focused on criterion validity through human performance




comparisons, while MMLU attempted to address multiple validity types simultaneously. However, new challenges emerged: convergent validity became more complex as models showed inconsistent performance across related tasks (e.g., philosophy versus morality questions), and discriminant validity concerns arose around distinguishing between memorization and reasoning capabilities. This progression has led to the current state of language model evaluation, characterized by greater sophistication in validity considerations but also a clearer recognition of inherent limitations. Recent work has highlighted the need for better convergent validity across benchmarks and more robust methods for assessing reasoning abilities. The field has moved from treating benchmarks as simple performance metrics to viewing them as complex instruments requiring multiple types of validation evidence ([Ruan et al., 2024](#)).
















D Case Studies

Table 5: Case Studies Summary. A subjective score for validity—the standard for “reasonable” is demonstrating that obvious risks to invalidity are addressed: : reasonable; : proceed with caution; : insufficient. Even for a score of “reasonable,” there will be weaknesses in the evidence. The score is given because the strengths outweigh the weaknesses in terms of determining the validity of the claim from that evidence. This is never a binary classification nor complete, and should rather be a cyclic process— for instance, as our forms of what constitutes graduate-level chemistry may evolve over time and from school to school.

Claims from MMLU Benchmark Accuracy Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. Language models can demonstrate broad knowledge across diverse academic and professional subjects.					
2. Language models can perform expert-level reasoning across specialized domains.					
3. MMLU performance predicts a model’s general language understanding capabilities.					
Claims from Graduate-Level Google-Proof Q&A (GPQA) Benchmark Accuracy Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. AI systems can accurately answer <i>graduate-level specialized multiple-choice questions</i> in biology, physics, and chemistry.					
2. AI systems can accurately answer <i>graduate-level specialized questions</i> in specialized scientific domains.					
3. AI systems can exhibit <i>general reasoning abilities</i> that can transfer beyond current human specialization.					
Claims from ImageNet Validity Assessment Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. ImageNet tests how well models learn complex associations between images and labels.					
2. ImageNet gauges the ability to learn semantically general visual features for object classification.					
3. ImageNet measures overall visual understanding of a model.					
Claims from MedQA Benchmark Accuracy Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. AI models can accurately answer USMLE-style multiple-choice questions in core medical fields (e.g., internal medicine, pediatrics).					
2. AI models can accurately answer advanced specialized medical questions across diverse clinical subfields (e.g., oncology, psychiatry).					
3. AI models exhibit general (human-like) medical reasoning abilities.					
Claims from SQuAD Benchmark Accuracy Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. The evaluated model can accurately identify the most relevant snippet of a high quality encyclopedia passage for answering a question about the passage.					
2. The evaluated model can accurately identify the most relevant snippet of an online text passage for answering a question about the passage.					
3. The evaluated model exhibits human-level reading comprehension.					

D.1 GPQA

Table 6: A Graduate-Level Google-Proof Q&A Benchmark (GPQA) (Rein et al., 2023) Application. A subjective score for validity—the standard for “reasonable” is demonstrating that obvious risks to invalidity are addressed: : reasonable; : proceed with caution; : insufficient. Even for a score of “reasonable,” there will be weaknesses in the evidence. The score is given because the strengths outweigh the weaknesses in terms of determining the validity of the claim from that evidence. This is never a binary classification nor complete, and should rather be a cyclic process— for instance, as our forms of what constitutes graduate-level chemistry may evolve over time and from school to school.

Claims from Graduate-Level Google-Proof Q&A (GPQA) Benchmark Accuracy Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. AI systems can accurately answer <i>graduate-level specialized multiple-choice questions</i> in biology, physics, and chemistry.					
2. AI systems can accurately answer <i>graduate-level specialized questions</i> in specialized scientific domains.					
3. AI systems can exhibit <i>general reasoning abilities</i> that can transfer beyond current human specialization.					

Description of dataset. The GPQA (Graduate-Level Google-Proof Question Answering) benchmark is a challenging dataset comprising 448 multiple-choice questions crafted by domain experts in biology, physics, and chemistry (Rein et al., 2023). These questions are designed to be exceptionally difficult, with experts holding or pursuing PhDs in the respective fields achieving an accuracy of 65% (74% when excluding clear mistakes identified retrospectively). Notably, highly skilled non-expert validators, even with unrestricted web access and spending over 30 minutes per question, attained only 34% accuracy, underscoring the “Google-proof” nature of the dataset. State-of-the-art AI systems also find this benchmark challenging; for instance, a GPT-4 based model achieved 39% accuracy. The GPQA dataset serves as a valuable resource for developing scalable oversight methods, aiming to enable human experts to effectively supervise and extract truthful information from AI systems that may surpass human capabilities.

Object of Claim: Multiple-choice questions in biology, physics, and chemistry accuracy.

Claim 1: AI models can accurately answer graduate-level specialized multiple-choice questions in biology, physics, and chemistry — criterion is accuracy on such questions.

Evidence: Accuracy on [N] multiple-choice questions in biology, physics, and chemistry.

Validity of Claim from Evidence:

1. Content Validity

- *Strength:* Expert-curated questions ensure high-quality, relevant content across key topics in biology, physics, and chemistry. The performance gap between experts and non-experts confirms the questions assess specialized knowledge.
- *Weakness:* The dataset’s construction criteria may exclude some relevant questions, potentially leading to over- or underrepresentation of certain subfields.

- *Suggestions:* Conduct systematic content mapping across subfields to ensure balanced representation. Include expert diversity analysis to mitigate potential biases in question selection.

2. Criterion Validity

- *Strength:* Human expert accuracy provides a meaningful external criterion, reinforcing concurrent validity.
- *Weakness:* Criterion validity could be stronger with comparisons to other specialized science Q/A benchmarks. Predictive validity is untested—no evidence that GPQA accuracy predicts future performance on exams or coursework, for example.
- *Suggestions:* Compare performance with established science Q&A benchmarks. Conduct longitudinal studies tracking how benchmark performance predicts success on real graduate exams.

3. Construct Validity

- Since the claim is strictly about accuracy on a defined criterion, construct validity is not necessary to evaluate this specific claim.

4. External Validity

- *Strength:* The test mirrors a real-world setting—human experts develop the questions, and the evaluation format aligns with academic multiple-choice assessments. GPQA includes diverse topics within its disciplines.
- *Weakness:* Similar to the criterion validity gap, GPQA accuracy is not compared to other multiple-choice science tests, leaving external generalization unverified.
- *Suggestions:* Validate against different question formats and compare performance across multiple science benchmarks.

5. Consequential Validity

- *Strength:* The AI-expert performance gap prevents premature claims of AI superiority, mitigating risks of overestimating AI scientific knowledge. However, models have quickly improved in this benchmark⁶. GPQA-trained models could support science education as study tools.
- *Weakness:* If AI models reach high accuracy, stakeholders may overgeneralize their competence, assuming they have true expertise in physics, biology, and chemistry, despite lacking deeper scientific reasoning skills.
- *Suggestions:* Develop clear guidance for stakeholders on interpreting results. Create documentation explicitly distinguishing multiple-choice performance from broader scientific expertise.

Object of Claim: Q/A in biology, physics, and chemistry accuracy.

Claim 2: AI models can accurately answer graduate-level questions in specialized scientific domains — criterion is accuracy on such questions.

Evidence: Accuracy on [N] multiple-choice questions in biology, physics, and chemistry.

Validity of Claim from Evidence:

⁶<https://www.youtube.com/watch?v=ZANbujPTvOY>.

1. Content Validity ⚠️

- *Strength:* Expert-curated, high-quality questions covering key topics in biology, physics, and chemistry. Non-expert performance gap supports specialization.
- *Weakness:* Limited to three disciplines, excluding other specialized scientific domains (e.g., medicine, engineering). Only Q/A questions, excluding fill-in-the-blank or open-ended questions.
- *Suggestions:* Expand questions to include other scientific subdomains. Conduct systematic content mapping across subfields to ensure balanced representation. Include expert diversity analysis to mitigate potential biases in question selection.

2. Criterion Validity ⚠️

- *Strength:* Human expert accuracy serves as a strong external criterion (concurrent validity). AI-expert performance gap reinforces benchmark credibility.
- *Weakness:* No predictive validity—GPQA accuracy is not tested against future performance on other specialized assessments.
- *Suggestions:* Establish correlations with performance on real graduate program assessments. Develop predictive validity studies tracking model performance across time and domains.

3. Construct Validity ⚠️

- *Strength:* Expert-curated questions in biology, physics, and chemistry are designed to capture fundamental aspects of specialized scientific knowledge. This suggests that the construct measured—domain-specific scientific competence—has meaningful representation, and high accuracy should correlate with understanding key scientific principles.
- *Weakness:* GPQA’s focus on biology, physics, and chemistry limits its ability to capture the overall construct of “specialized scientific knowledge,” as other fields like medicine and engineering require different reasoning and knowledge structures. Moreover, the paper does not provide evidence linking GPQA performance to external measures of scientific competence (such as standardized test scores), leaving its alignment with related constructs unclear. Finally, the multiple-choice format may favor recognition or memorization over deeper analytical reasoning, potentially failing to capture key facets like synthesis and in-depth understanding.
- *Suggestions:* To improve construct validity, expand GPQA to include additional domains (e.g., medicine, engineering) and correlate its scores with independent standardized assessments to establish convergent and discriminant validity. Additionally, incorporating alternative formats like open-ended questions and problem-solving tasks will better capture deep analytical reasoning and synthesis skills.

4. External Validity ⚠️

- *Strength:* Real-world, expert-created multiple-choice questions ensure relevance. Coverage across multiple subfields increases generalization within biology, physics, and chemistry.
- *Weakness:* No evidence of generalization to other science assessments (e.g., (non-)multiple choice PhD qualifying exams).

- *Suggestions:* Test generalization to other assessment formats including written exams, oral defenses, and research proposal evaluations.

5. Consequential Validity ⚠️

- *Strength:* AI-expert performance gap prevents overstating AI’s scientific capabilities; models could support science education.
- *Weakness:* Risk of overgeneralization—high scores may be misinterpreted as broad scientific expertise beyond tested domains.
- *Suggestions:* Create clear limitations documentation highlighting specific domains where evidence supports or doesn’t support performance claims.

Object of Claim: Reasoning.

Claim 3: AI models exhibit general reasoning abilities.

Evidence: Accuracy on [N] multiple-choice questions in biology, physics, and chemistry.

Validity of Claim from Evidence:

1. Content Validity ⚠️

- *Strength:* Covers multiple scientific disciplines, requiring some level of reasoning beyond factual recall.
- *Weakness:* Multiple-choice format limits assessment of forms of reasoning like logical deduction, or abstract problem-solving.
- *Suggestions:* Develop specific reasoning-focused questions that isolate logical deduction from domain knowledge. Include diverse reasoning types (inductive, deductive, abductive).

2. Criterion Validity ❌

- *Strength:* Human expert accuracy serves as a real-world external criterion, and the AI-expert performance gap indicates a meaningful benchmark for reasoning capabilities.
- *Weakness:* GPQA tests factual and applied knowledge rather than abstract reasoning skills. No predictive validity—performance on GPQA is not tested against other established reasoning benchmarks (e.g., LSAT-style logical reasoning or problem-solving tests).
- *Suggestions:* Compare performance against established reasoning benchmarks like LSAT, GRE analytical, and domain-independent logical reasoning tests.

3. Construct Validity ❌

- *Strength:* AI performance on GPQA correlates with success in structured question-answering tasks, suggesting some reasoning component. Additionally, the dataset can distinguish between human experts and non-experts.
- *Weakness:* Does not separate reasoning from memorization—AI models may exploit dataset patterns rather than apply logical deduction. While non-experts with access to Google perform worse than experts, non-experts are given a limited time per question, which may not sufficiently show that models have not been trained on such questions.

No convergent validity—GPQA accuracy is not correlated with performance on explicit reasoning assessments. No discriminant validity—It is unclear whether GPQA measures reasoning ability or just domain-specific knowledge.

- *Suggestions:* Conduct factor analysis to distinguish reasoning from memorization. Demonstrate convergent validity with dedicated reasoning assessments and discriminant validity from pure knowledge recall.




4. External Validity ❌
















- *Strength:* GPQA questions require problem-solving across multiple disciplines, increasing the likelihood that some reasoning ability is being tested.
- *Weakness:* Reasoning should generalize across domains, but GPQA only includes three scientific fields. No evidence that AI models with high GPQA accuracy perform well on general reasoning tasks outside science (e.g., logical puzzles, mathematical proofs, legal or philosophical reasoning).
- *Suggestions:* Test performance on reasoning tasks across non-scientific domains including logic puzzles, mathematical proofs, and philosophical arguments.

5. Consequential Validity ⚠️

- *Strength:* If GPQA successfully measures reasoning, AI models excelling on it could serve as decision-support tools in scientific research or education.
- *Weakness:* Overgeneralization risk—high GPQA accuracy may lead to misinterpreting AI as possessing broad, human-like reasoning abilities when it may only excel at structured multiple-choice problems.
- *Suggestions:* Develop clear performance interpretation guidelines specifying which reasoning capabilities are supported by evidence versus which remain speculative.

D.2 MMLU

Table 7: A Massive Multitask Language Understanding (MMLU) [Hendrycks et al. \(2020\)](#) Application. A subjective score for validity—the standard for “reasonable” is demonstrating that obvious risks to invalidity are addressed: : reasonable; : proceed with caution; : insufficient. Even for a score of “reasonable,” there will be weaknesses in the evidence. The score is given because the strengths outweigh the weaknesses in terms of determining the validity of the claim from that evidence. This is never a binary classification nor complete, and should rather be a cyclic process.

Claims from MMLU Benchmark Accuracy Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. Language models can demonstrate broad knowledge across diverse academic and professional subjects.					
2. Language models can perform expert-level reasoning across specialized domains.					
3. MMLU performance predicts a model’s general language understanding capabilities.					

Description of dataset. Massive Multitask Language Understanding (MMLU) is a benchmark designed to test natural language understanding across 57 subjects spanning STEM, humanities, social sciences, and professional fields. It consists of multiple-choice questions (four options) drawn from standardized tests like the GRE and medical licensing exams, LSAT exams, and various exams oriented towards domain specific knowledge in the fields listed above

Object of Claim: Broad knowledge across diverse subjects.

Claim 1: Language models can demonstrate broad knowledge across diverse academic and professional subjects.

Evidence: Accuracy on [N] multiple-choice questions spanning 57 subjects across STEM, humanities, social sciences, and professional fields, drawing from practice questions for standardized tests such as the Graduate Record Examination and the United States Medical Licensing Examination.

Validity of Claim from Evidence:

1. Content Validity

- *Strength:* MMLU covers an extensive range of domains (57 subjects) spanning STEM, humanities, social sciences, and professional fields.
- *Weakness:* The multiple-choice format with only four options limits the depth of understanding that can be assessed, and some subjects may have inadequate representation.
- *Suggestions:* Conduct detailed content mapping to ensure proportional representation across domains and expand beyond multiple-choice to include open-ended responses.

2. Criterion Validity

- *Strength:* MMLU has been shown to correlate with downstream performance on other capability oriented tasks, demonstrating predictive validity. Related work on benchmarking measured correlation of MMLU scores with the aggregate of scores on MMLU and other capability benchmarks, and found that MMLU to have a very high correlation only behind MedQA and Arc Challenge ([Ren et al., 2024](#)).

- *Weakness:* There are inconsistencies in how well MMLU correlates with other measures of related capabilities (e.g., models performing well on philosophy but poorly on morality despite their relatedness).
- *Suggestions:* Conduct more systematic studies correlating MMLU performance with other established benchmarks of knowledge across domains.

3. Construct Validity ⚠️

- *Strength:* The benchmark draws from standardized tests designed to measure knowledge in respective fields.
- *Weakness:* MMLU doesn't effectively distinguish between recall and reasoning lacking discriminant validity; high performance could indicate mere memorization from training data scraped from the internet rather than deep understanding.
- *Suggestions:* Add questions that explicitly test reasoning or precision versus recall, and incorporate analysis of model explanations, not just final answers.

4. External Validity ⚠️

- *Strength:* Using questions from standardized tests provides some real-world grounding.
- *Weakness:* Significant issues undermine generalizability: labeling errors (57% of Virology questions contain errors), answer ordering effects, and the constrained multiple-choice format. This suggests an independent reproduction of MMLU might present different results.
- *Suggestions:* Implement rigorous quality control (as in MMLU-Pro), test with varied answer orderings, and expand beyond multiple-choice formats.

5. Consequential Validity ⚠️

- *Strength:* MMLU has successfully become a standard benchmark driving industry progress in language model development.
- *Weakness:* There is a risk of overoptimization as models are increasingly designed specifically to perform well on MMLU multiple choice, and might overfit to doing well on easily testable questions rather than broad subject knowledge (Goodhart's Law).
- *Suggestions:* Regularly update the benchmark with new questions and maintain clear documentation about what MMLU does and doesn't measure.

Object of Claim: Expert-level reasoning.

Claim 2: Language models can perform expert-level reasoning across specialized domains.

Evidence: MMLU compares model performance against estimated expert-level accuracy (89.8%) and measures performance across specialized domains from medicine to formal logic.

Validity of Claim from Evidence:

1. Content Validity ⚠️

- *Strength:* MMLU includes questions from specialized professional domains that require some domain expertise.

- *Weakness:* Multiple-choice questions as they are written within MMLU primarily tests factual knowledge rather than complex reasoning processes experts employ.
- *Suggestions:* Include multi-step reasoning problems and questions requiring application of principles to novel scenarios.

2. Criterion Validity

- *Strength:* Performance is benchmarked against estimated expert-level accuracy (89.8%) so MMLU has a good claim to concurrent validity
- *Weakness:* The benchmark cannot distinguish between memorized answers and expert reasoning. Error analysis shows 39% of incorrect answers on MMLU-Pro stem from reasoning errors despite correct knowledge, meaning the correlation with correct answers might be spurious.
- *Suggestions:* Incorporate expert validation of both answers and reasoning paths, perhaps through analysis of model explanations.

3. Construct Validity

- *Strength:* Some questions require application of domain knowledge rather than simple facts.
- *Weakness:* The benchmark doesn't capture expert reasoning processes, only the final answers lacking structural validity.
- *Suggestions:* Develop metrics to evaluate reasoning quality, not just answer correctness, and include questions that cannot be solved through memorization alone. Elicit experts per domain for their reasoning process, as well as suggestions for relevant question formats and protocols.

4. External Validity

- *Strength:* Using standardized test questions provides some grounding in real assessment practices. As mentioned earlier, there is some evidence MMLU performance is correlated with performance on other capability benchmarks.
- *Weakness:* Multiple-choice tests do not capture the open-ended, iterative nature of expert reasoning in real-world contexts. Changing answer ordering can also affect scores which an expert should be invariant to.
- *Suggestions:* Develop supplementary benchmarks with more authentic professional tasks and varied formats. Perhaps where the model provides reasoning chains and is evaluated with a reward model calibrated to expert preference.

5. Consequential Validity

- *Strength:* The benchmark has helped identify strengths and weaknesses in model capabilities across different domains.
- *Weakness:* High MMLU scores might create an illusion that models can replace domain expert judgement, leading to inappropriate applications.
- *Suggestions:* Provide clear guidance on the limitations of what MMLU scores indicate about true expert-level reasoning.

Object of Claim: Predictive power for general capabilities.

Claim 3: MMLU performance predicts a model’s general language understanding capabilities.

Evidence: MMLU has been highly correlated with downstream quality and capability, as noted by industry teams building large language models and supported by research on observational scaling laws.

Validity of Claim from Evidence:

1. Content Validity ⚠️

- *Strength:* MMLU covers a wide range of domains, providing breadth in assessment that is a non-trivial subset of understanding of “general” topics, if such topics are the enumeration of all academic topics.
- *Weakness:* It doesn’t cover all aspects of language understanding, particularly creative, open-ended, or interactive capabilities. It also doesn’t cover areas of knowledge that aren’t readily measured in academic settings.
- *Suggestions:* Supplement with other benchmarks measuring different facets of language understanding and areas that don’t easily map to academic fields of study such as humor.

2. Criterion Validity ⚠️

- *Strength:* Research on observational scaling laws notes that when running a PCA on evaluation performance of prominent benchmarks against downstream performance, variation in MMLU explains a large fraction of variation [Ruan et al. \(2024\)](#). As mentioned earlier in claim 1 and claim 2, research shows MMLU scores correlate well with performance on other tasks, supporting its use as a general predictor ([Ren et al., 2024](#)). Combined with the earlier observation that performance is benchmarked against estimated expert-level accuracy (89.8%), this gives MMLU a good claim to concurrent validity.
- *Weakness:* Correlation patterns are inconsistent across different types of tasks and domains ([Wang et al., 2024](#)).
- *Suggestions:* Develop a more nuanced framework showing which aspects of MMLU best predict which types of downstream capabilities or rely on the observational scaling laws framework.

3. Construct Validity ❌

- *Strength:* The benchmark captures some aspects of knowledge acquisition and application.
- *Weakness:* “Natural language understanding” as a construct encompasses much more than multiple-choice question answering, including discourse comprehension, pragmatics, and nuanced interpretation none of which are covered here.
- *Suggestions:* Clarify the specific sub-constructs of language understanding that MMLU actually measures.

4. External Validity ❌




- *Strength:* The breadth of subjects provides some basis for generalization, assuming we are focused on breadth and a more shallow definition of generality rather than depth.
















- *Weakness:* MMLU’s format and limitations (answer ordering effects, label errors) raise questions about how well scores generalize to real-world language understanding tasks ([Wang et al., 2024](#); [Gupta et al., 2024](#)).
- *Suggestions:*

5. Consequential Validity ⚠️

- *Strength:* MMLU has influenced productive research directions in language model development, such as BigBench, GPQA, GAIA and other benchmarks that test language models on a broad set of tasks.
- *Weakness:* Over-reliance on MMLU as a general capability metric could lead to narrowly optimized models for the benchmark rather than genuinely more capable ones. This can lead to overstating progress and capabilities of the latest models and systems, i.e. models such as Phi-1 and Mistral which overfits to GSM8k and saw large drops in performance when tested on a new private split ([Zhang et al., 2024](#)).
- *Suggestions:* Develop complementary metrics that capture aspects of language understanding not measured by MMLU, and emphasize a balanced assessment approach.

D.3 ImageNet

Table 8: An ImageNet (Deng et al., 2009; Russakovsky et al., 2014) Application. A subjective score for validity—the standard for “reasonable” is demonstrating that obvious risks to invalidity are addressed: : reasonable; : proceed with caution; : insufficient. Even for a score of “reasonable,” there will be weaknesses in the evidence. The score is given because the strengths outweigh the weaknesses in determining the validity of the claim from that evidence. This evaluation is an iterative process, acknowledging that both the benchmark and its interpretations may evolve over time.

Claims from ImageNet Validity Assessment Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. ImageNet tests how well models learn complex associations between images and labels.					
2. ImageNet gauges the ability to learn semantically general visual features for object classification.					
3. ImageNet measures overall visual understanding of a model.					

Description of dataset. ImageNet (Deng et al., 2009; Russakovsky et al., 2014) (specifically ILSVRC 2012) is a benchmark for predicting an image’s label from a fixed set of 1000 diverse categories. The dataset—curated primarily from Flickr with human annotation—is evaluated using accuracy/error rate and precision/recall metrics.

Object of Claim: Accuracy in labeling 1000 diverse image classes.

Claim 1: Model architectures can accurately predict labels for images.

Evidence: Performance on accuracy/error rate and precision/recall metrics.

Validity of Claim 1 from Evidence:

1. Content Validity

- *Strength:* The dataset covers 1000 diverse categories with extensive natural variability—including differences in poses, lighting, backgrounds, and fine-grained distinctions (e.g., different dog breeds)—making it well-suited to assess image-label associations.
- *Weakness:* It is confined to static, natural RGB images and does not include other modalities (e.g., grayscale medical images or hyperspectral data) or dynamic contextual information (e.g., actions or inter-object relationships).
- *Suggestions:* Clearly specify that ImageNet targets static natural images, and consider integrating supplementary datasets to represent additional image types or contextual settings.

2. Criterion Validity

- *Strength:* There is robust evidence that performance on ImageNet is both predictive of downstream task success (models excelling on ImageNet often perform well on benchmarks such as CIFAR or Caltech, and in real-world applications like wildlife classification (Norouzzadeh et al., 2018)) and concurrent with human-annotated labels under similar conditions.

- *Weakness:* Most supporting evidence arises from community practice rather than controlled, formal comparative studies.
- *Suggestions:* Undertake targeted comparative evaluations with other established benchmarks to better substantiate both the predictive and concurrent facets of criterion validity.

3. External Validity

- *Strength:* The dataset is representative of real-world natural images, and its utility has been demonstrated under varying conditions (differences in image quality, size, and even in applications to non-traditional domains such as medical imaging (Irvin et al., 2019) or adversarially constructed settings (Salaudeen and Hardt, 2024)).
- *Weakness:* Its effectiveness in generalizing to highly specialized or non-natural imaging contexts remains less certain.
- *Suggestions:* Expand external validation by testing models on a broader range of datasets that capture diverse imaging modalities and conditions.

4. Construct Validity

- Since the claim is strictly about accuracy on a defined criterion, construct validity is not necessary to evaluate this specific claim.

5. Consequential Validity

- *Strength:* The clear quantification of labeling accuracy offers a concrete performance metric, facilitating transparent and reproducible comparisons.
- *Weakness:* There is a risk that high ImageNet accuracy may be misinterpreted as reflecting comprehensive visual understanding, potentially leading to overconfident real-world deployments.
- *Suggestions:* Advise stakeholders that ImageNet performance should be interpreted strictly as a measure of static image classification and that complementary evaluations are necessary to assess broader aspects of visual intelligence.

Object of Claim: Learning of semantically general visual features.

Claim 2: ImageNet evaluates the ability of models to learn transferable visual features that are useful for object classification.

Evidence: Performance gains in fine-tuning tasks when using models pretrained on ImageNet, compared to those trained from scratch.

Validity of Claim 2 from Evidence:

1. Content Validity

- *Strength:* The wide coverage of natural image phenomena—including fine-grained details and numerous object classes—supports the learning of varied and versatile visual features.
- *Weakness:* It may not comprehensively represent features present in non-natural or synthetic environments, nor fully capture abstract contextual cues.

- *Suggestions:* Consider integrating supplementary datasets that include synthetic, non-natural, or contextually complex images to achieve a more comprehensive assessment.

2. Criterion Validity

- *Strength:* Empirical studies (e.g., Kornblith et al. (2018)) show that ImageNet pretraining is strongly predictive of improved fine-tuning and transfer learning outcomes and that performance is concurrent with established classification tasks, addressing both the predictive and concurrent dimensions.
- *Weakness:* Although the predictive correlation is robust, direct and extensive concurrent comparisons with alternative feature assessment methods are less common.
- *Suggestions:* Enhance validation by conducting side-by-side evaluations comparing learned features across different pretraining methods and downstream tasks.

3. Construct Validity

- *Strength:* The improvement in fine-tuning performance suggests that the learned features are semantically rich and transferable. This provides evidence of structural validity (as features capture fundamental visual components), convergent validity (via correlation with downstream task performance), and discriminant validity (in differentiating meaningful features from noise).
- *Weakness:* It is challenging to definitively establish that these benefits are due to genuine generalization of visual features rather than overfitting to ImageNet-specific patterns, leaving the discriminant aspect less clear.
- *Suggestions:* Perform in-depth analyses—such as saliency mapping or kernel visualization—to further elucidate the nature of the learned features and clarify the extent of structural, convergent, and discriminant validity.

4. External Validity

- *Strength:* The benefits of ImageNet pretraining have been observed across multiple downstream benchmarks, suggesting that the learned features generalize beyond the confines of natural images.
- *Weakness:* The degree of generalizability across diverse domains (e.g., synthetic or non-natural images) remains to be fully validated.
- *Suggestions:* Broaden external validation by pretraining on a more diverse set of data and assessing performance on cross-domain tasks.

5. Consequential Validity

- *Strength:* The transformative impact of ImageNet pretraining in advancing computer vision is well-documented, highlighting its practical benefits.
- *Weakness:* An overreliance on fine-tuning improvements may obscure limitations in the intrinsic quality of the learned features, risking overgeneralization regarding model capability.
- *Suggestions:* Clearly communicate that fine-tuning gains indicate enhanced performance in specific settings rather than a comprehensive measure of visual feature quality; encourage complementary evaluations focused specifically on feature robustness.

Object of Claim: Overall visual understanding.

Claim 3: ImageNet provides an indication of a model’s overall visual understanding beyond simple label prediction or isolated feature representation.

Evidence: Performance on the standard classification task under controlled evaluation conditions, independent of training context.

Validity of Claim 3 from Evidence:

1. **Content Validity** ✖

- *Strength:* The task of image classification is well-defined and widely used as a proxy for certain aspects of visual understanding.
- *Weakness:* Relying solely on classification does not capture the full range of visual understanding, which includes spatial reasoning, object detection, contextual awareness, and causal interpretation.
- *Suggestions:* Complement the classification task with additional evaluations—such as object detection, visual question answering, or spatial reasoning challenges—to more fully capture the construct.

2. **Criterion Validity** ✖

- *Strength:* Classification accuracy is a clear and quantifiable metric that enables direct comparison across models, addressing both predictive and concurrent aspects to some degree.
- *Weakness:* There is limited evidence that high performance on this narrow task reliably predicts the broader and deeper aspects of overall visual understanding.
- *Suggestions:* Compare ImageNet classification results with those from benchmarks explicitly designed to evaluate advanced visual reasoning and interpretative skills.

3. **Construct Validity** ✖

- *Strength:* Operationalizing visual understanding as performance on image labeling provides a measurable framework that reflects a basic structural organization of visual recognition. However, it offers only limited convergent evidence with tasks requiring integrated reasoning and does not fully differentiate (discriminant validity) between mere pattern recognition and comprehensive understanding.
- *Weakness:* This narrow operational approach may oversimplify the construct, favoring models that exploit dataset biases rather than achieving holistic visual comprehension.
- *Suggestions:* Introduce complementary evaluation tasks (e.g., visual question answering or spatial reasoning challenges) to capture additional dimensions of visual understanding and enhance assessments of structural, convergent, and discriminant validity.

4. **External Validity** ✖



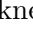
- *Strength:* ImageNet’s evaluation framework is reproducible, and similar performance trends have been observed across related image-based tasks.
- *Weakness:* Its ability to generalize to tasks requiring integrated reasoning, spatial awareness, and contextual interpretation remains unconfirmed.
















- *Suggestions:* Validate the broader aspects of visual understanding by employing a wider array of benchmarks that emphasize multidimensional reasoning and contextual evaluation.

5. Consequential Validity ✖

- *Strength:* The benchmark has stimulated important discussions on the limitations of measuring visual intelligence solely via classification, underscoring the need for more comprehensive evaluation methods.
- *Weakness:* High classification accuracy might be erroneously interpreted as evidence of complete visual understanding, potentially misleading real-world applications.
- *Suggestions:* Provide clear guidelines on the interpretative scope of ImageNet results and promote complementary measures to capture the full spectrum of visual intelligence.

D.4 MedQA

Table 9: MedQA Jin et al. (2021) Application. A subjective score for validity—the standard for “reasonable” is demonstrating that obvious risks to invalidity are addressed: : reasonable; : proceed with caution; : insufficient. Even for a score of “reasonable,” there will be weaknesses in the evidence. The score is given because the strengths outweigh the weaknesses in terms of determining the validity of the claim from that evidence. This is never a binary classification nor complete and should rather be a cyclic process.

Claims from MedQA Benchmark Accuracy Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. AI models can accurately answer USMLE-style multiple-choice questions in core medical fields (e.g., internal medicine, pediatrics).					
2. AI models can accurately answer advanced specialized medical questions across diverse clinical subfields (e.g., oncology, psychiatry).					
3. AI models exhibit general (human-like) medical reasoning abilities.					

Description of dataset. The MedQA benchmark is a large-scale, multilingual dataset crafted for open-domain question answering in the medical domain. It consists of multiple-choice questions drawn from professional medical board exams in English (12’723 questions), simplified Chinese (34’251 questions), and traditional Chinese (14’123 questions) testing complex clinical reasoning. Unlike prior QA datasets, MedQA emphasizes real-world diagnostic decision-making, requiring systems to retrieve and interpret evidence from extensive medical textbook corpora.

Object of Claim: Multiple-choice questions in USMLE core fields.

Claim 1: AI models can accurately answer USMLE-style multiple-choice questions in core medical fields (e.g., internal medicine, pediatrics).

Evidence: Accuracy on MedQA, a curated dataset containing 12,723 English USMLE-style multiple-choice questions, part of a larger multilingual collection that includes 34,251 simplified Chinese and 14,123 traditional Chinese questions. Original baseline models achieved only 36.7% accuracy on the English questions, while recent LLMs have reached 90% accuracy on this benchmark.

Validity of Claim from Evidence:

1. Content Validity

- *Strength:* The question set covers standard USMLE core areas (internal medicine, pediatrics, OB/GYN, surgery), curated by medical professionals. The gap between expert vs. non-expert performance helps confirm that the items do measure specialized knowledge.
- *Weakness:* Even “core” USMLE topics might be incomplete (e.g., narrower coverage of pediatrics vs. adult medicine).
- *Suggestions:* Do a content-mapping across subdomains to ensure each core field is represented proportionally. Include item analyses by domain experts to identify underrepresented subtopics.

2. Criterion Validity

- *Strength:* If the claim is specifically “accuracy on USMLE-style questions,” then MedQA directly measures that criterion. High performance against human experts or official pass thresholds bolsters concurrent validity.
- *Weakness:* There is limited predictive validity—we do not know if high scores on MedQA predict performance on subsequent medical assessments, other board certifications, or related medical knowledge evaluations.
- *Suggestions:* Compare MedQA performance to known USMLE pass rates or step scores. Conduct longitudinal or prospective studies to see if a model that excels on MedQA also performs robustly in real USMLE test trials.

3. Construct Validity ⚠️

- *Strength:* TBD
- *Weakness:* While USMLE questions are designed to test integrated medical knowledge and clinical decision-making, LLMs might leverage statistical patterns in their training data rather than demonstrating the intended construct. Many multiple-choice questions can be solved via pattern matching or memorization without genuine conceptual understanding.
- *Suggestions:* Further analyze how the model arrives at answers. Include open-ended or explanation-based items to confirm it is using medical reasoning (rather than memorized patterns).

4. External Validity ⚠️

- *Strength:* Because USMLE is a well-established exam format, it is somewhat representative of real licensing test questions.
- *Weakness:* The model’s performance is not tested in truly “real-world” situations (e.g., diagnosing patients with partial information). Variation in language, test format, or question style might degrade performance.
- *Suggestions:* Assess generalizability by testing with alternative question sources (e.g., NBME question banks, other medical boards), including different item formats (e.g., free-response, extended matching).

5. Consequential Validity ⚠️

- *Strength:* If model performance is below human experts, it prevents overestimation of AI’s clinical capabilities; the benchmark helps calibrate expectations.
- *Weakness:* If the model achieves high scores, there is a risk that stakeholders assume it can practice medicine or make reliable diagnoses—something USMLE-style Q&A alone does not prove.
- *Suggestions:* Provide guidance that warns against using MedQA results as a proxy for “clinical readiness.” Create disclaimers, ethics reviews, or guidelines so that high MedQA accuracy is not over-interpreted as real-world medical competency.

Object of Claim: Advanced specialized medical Q/A accuracy.

Claim 2: AI models can accurately answer advanced specialized medical questions across diverse clinical subfields (e.g., oncology, psychiatry, cardiology).

Evidence: The same MedQA multiple-choice items, which may include some specialized subtopics but are typically broad “licensing exam” style.

Validity of Claim from Evidence:

1. Content Validity ⚠️

- *Strength:* USMLE exams do include a range of subfields. If MedQA is properly sampled, it will have at least basic coverage in oncology, psychiatry, etc.
- *Weakness:* “Advanced specialized” questions in niche fields (e.g., transplant immunology, pediatric oncology) are usually not heavily represented in general licensing exams, so coverage may be thin.
- *Suggestions:* Evaluate how many questions truly belong to each advanced specialty. Expand the dataset or collect specialized question sets from relevant board exams (e.g., ABIM Oncology boards).

2. Criterion Validity ⚠️

- *Strength:* If specialists or specialized board pass rates are used as a reference, some measure of concurrent validity might be feasible.
- *Weakness:* We lack direct evidence that performance on these general medical exams transfers to in-depth specialty boards or practice.
- *Suggestions:* Correlate MedQA scores with actual performance on specialized board-style question sets. Conduct predictive analyses to see whether high performance in general med licensing implies success in more advanced specialties.

3. Construct Validity ⚠️

- *Strength:* High MedQA accuracy suggests some knowledge of specialized subfields. The benchmark captures certain aspects of clinical reasoning, including diagnostic pattern recognition, treatment, management and application of domain-specific knowledge in areas like oncology, cardiology, and psychiatry, though in a constrained multiple-choice format.
- *Weakness:* “Advanced specialized competence” is a broader construct than a general licensing exam can measure. True advanced knowledge typically requires deeper reasoning and domain-specific problem-solving, not just broad-spectrum test items. MCQ format doesn’t capture critical elements of specialized practice such as open-ended diagnostic reasoning, iterative decision-making based on evolving clinical information, managing uncertainty, and generating (rather than selecting) management plans.
- *Suggestions:* Use specialized test banks or real advanced clinical vignettes. Check convergent validity: does a model that excels at MedQA also excel at an oncology-focused question set, for example?

4. External Validity ⚠️

- *Strength:* If the subfield questions in MedQA truly reflect real exam conditions, there is some external relevance, as these exams are designed to assess knowledge that specialized experts need to demonstrate for certification, suggesting the model shares at least some abilities with trained specialists.

- *Weakness:* Real clinicians in advanced specialties face more complex tasks than multiple-choice. We do not know if “exam success” generalizes to real specialist scenarios (e.g., reading labs, imaging).
- *Suggestions:* Compare the model performance to other specialized exams or real-world performance data, like mock boards or practical OSCE (Objective Structured Clinical Examination) tasks.

5. Consequential Validity ⚠️

- *Strength:* Since MedQA is drawn from a real medical practitioner exam (USMLE), it helps us assess whether models share one key aspect of medical expert competence, guarding against employing models that don’t pass this necessary (but not necessarily sufficient) bar for being able to address specialized medical questions.
- *Weakness:* There is a risk that stakeholders might over-interpret performance. Specifically, since MedQA cannot adequately cover the depth and diversity of all medical specialties, high overall performance might mistakenly be used as proof that the model is capable in particular specialties that weren’t well-represented in the benchmark. This could lead to inappropriate deployment in specialized domains where the model lacks adequate capabilities.
- *Suggestions:* Provide disclaimers, track real-world usage carefully, and ensure domain experts remain in the loop before trusting the system with advanced clinical decision-making.

Object of Claim: Reasoning.

Claim 3: AI models exhibit general (human-like) medical reasoning abilities.

Evidence: Same accuracy results on MedQA multiple-choice questions.

Validity of Claim from Evidence:

1. Content Validity ⚠️

- *Strength:* Medical licensing questions often require some reasoning (diagnostic logic, integrative thinking), so it is not purely rote.
- *Weakness:* While MedQA includes reasoning-oriented questions, it cannot cover the full breadth of medical reasoning scenarios. Important instances like reasoning about novel specialized cases, emergency decision-making with incomplete information, or longitudinal patient management are likely underrepresented content-wise.
- *Suggestions:* Include question types that explicitly test reasoning steps, causal inferences, or open-ended rationales, rather than single-select answers.

2. Criterion Validity ❌

- *Strength:* TBD
- *Weakness:* The claim ‘general medical reasoning’ is an abstract construct for which no single, universally accepted criterion measure exists. MedQA performance alone doesn’t provide evidence of correlation with any established reasoning assessments that might serve as imperfect but useful criteria.

- *Suggestions:* Compare performance on specialized “reasoning tests” (e.g., medical logic puzzles, case simulations). Show that high MedQA scorers also do well on validated reasoning exams for medical students or residents.

3. Construct Validity ❌

- *Strength:* TBD
- *Weakness:* The model’s question-answer patterns might rely on memorized knowledge or superficial cues—meaning it does not necessarily demonstrate the deeper mental processes clinicians use.
- *Suggestions:* Separate knowledge recall from genuine inference (factor analysis, requiring step-by-step justifications). Show strong correlations with tasks specifically designed to test “reasoning,” not just recall.




4. External Validity ❌
















- *Strength:* TBD
- *Weakness:* “General reasoning” implies an ability that transfers to any problem domain, but we only have evidence from a medical exam standpoint. No demonstration it generalizes to other fields or even beyond multiple-choice contexts.
- *Suggestions:* Evaluate the same model on other reasoning-heavy tasks (e.g., logic puzzles, legal reasoning sets, real-time patient simulations) to see if it truly exhibits domain-general reasoning.

5. Consequential Validity ⚠️

- *Strength:* If the model’s limitations are made explicit, at least we avoid the pitfall of claiming broad, “human-level” reasoning from a single test.
- *Weakness:* Over interpretation of high MedQA accuracy might lead to the illusion that the model “thinks like a doctor.” This could encourage unsupervised use in clinical settings.
- *Suggestions:* Provide disclaimers clarifying that test performance \neq human clinical reasoning. Develop ethical guidelines so that strong test scores do not lead to unqualified acceptance of AI’s medical judgments.

D.5 SQuAD

Table 10: SQuAD [Rajpurkar et al. \(2016\)](#) Application. A subjective score for validity—the standard for “reasonable” is demonstrating that obvious risks to invalidity are addressed: : reasonable; : proceed with caution; : insufficient. Even for a score of “reasonable,” there will be weaknesses in the evidence. The score is given because the strengths outweigh the weaknesses in terms of determining the validity of the claim from that evidence. This is never a binary classification nor complete, and should rather be a cyclic process.

Claims from SQuAD Benchmark Accuracy Report Card					
Claims	Content	Criterion	Construct	External	Consequential
1. The evaluated model can accurately identify the most relevant snippet of a high quality encyclopedia passage for answering a question about the passage.					
2. The evaluated model can accurately identify the most relevant snippet of an online text passage for answering a question about the passage.					
3. The evaluated model exhibits human-level reading comprehension.					

Description of benchmark. The SQuAD benchmark (v1.0) was initially released in fall 2016, before large pre-trained models like BERT or GPT-2 were developed (cite). It consists of over 100k questions drawn from over 500 selected English Wikipedia articles. Each problem in the benchmark consists of: 1) a passage (a paragraph from an article); 2) a question about the passage; 3) a span of the passage which contains the answer to the question. Answers to each problem are spans of the passage.

Object of Claim: Identifying the most relevant snippet of a high quality encyclopedia passage for answering a question about the passage.

Claim 1: The evaluated model can accurately identify the most relevant snippet of a encyclopedia passage for answering a question about the passage.

Evidence: The Wiki articles used are drawn from highly ranked pages according to Project Nayuki, and are therefore considered high quality. Around 500 articles on various subjects are used as passages for questions. 100,000 questions are used. The questions are original, human-constructed by vetted crowdsource workers, and are designed to be difficult. Both a strict metric (exact match accuracy) and a fuzzy metric (F1 overlap score) are used to determine model and human performance. A human baseline from vetted crowdworkers is included, which also serves as a measure of human annotator agreement on the correct answer to passages. The authors analyze answer types and find that they are varied, such as dates, locations, and quantities. A ‘dumb’ model (logistic regressor), capable of only surface-level pattern matching, is shown to perform much worse than the human baseline.

Validity of Claim from Evidence:

1. Content Validity

- *Strength:* Uses a large number and variety of questions and answer types. Exclusively uses highly ranked Wiki pages. Uses a relatively large number of Wiki pages to draw passages from.

- *Weakness:* TBD
- *Suggestions:* Using a larger number of Wiki pages would add to the content coverage of questions/answers.

2. Criterion Validity

- *Strength:* Annotator agreement on the answer snippets to passages is high and, therefore, demonstrates the answers to the questions are indeed correct/reliable. Crowdsourced workers are also vetted before being accepted as annotators, strengthening this point. As expected, human beings perform significantly better than a logistic regressor on SQuAD 1.0. This demonstrates that simple, surface-level pattern matching (the only thing a simple logistic regressor is capable of) is insufficient for solving questions on SQuAD 1.0.
- *Weakness:* No evidence provided for predictive validity.
- *Suggestions:* Some justification of basic predictive validity would be helpful here, such as a comparison of model performance on the benchmark to a downstream task, e.g. rate of question-answer completion of crowdsourced worker with and without access to the model as a tool. Still, I feel that there is generally reasonable evidence/justification provided.

3. Construct Validity

- Not applicable, since we are measuring a criterion, not a construct.

4. External Validity

- *Strength:* TBD
- *Weakness:* Although there is no clear evidence provided in the SQuAD paper of external validity (aside from the reassurance of its content validity via the diversity of topic and writing style of Wiki articles), the fact the claim is so narrow reduces the importance of this lack of evidence and warrants caution, rather than outright insufficiency.
- *Suggestions:* That being said, some evidence that would improve this borderline external validity is evaluation of a greater variety of encyclopedia formats, such as translated encyclopedia entries, which would build confidence in the likelihood that model performance on SQuAD would generalize across encyclopedic text.

5. Consequential Validity

- *Strength:* The authors mention that the size of the question and answer bank allows it to double as a large source of high-quality training data for question-answering systems, thereby offering a unique advantage in advancing model performance at this task.
- *Weakness:* The authors don't consider potential downstream harms of the subject matter included and excluded in the 500 Wiki articles used for generating the benchmark questions, such as the articles being biased toward a particular culture, country, discipline, etc. This could lead to model development disproportionately focusing on the overrepresented subject matter at the cost of excluded subjects.
- *Suggestions:* A basic consideration or analysis of the distribution of subject matter of the questions used in the benchmark would be useful for identifying if the concern mentioned in 'Weakness' is warranted or not.

Object of Claim: Identifying the most relevant snippet of an online text passage for answering a question about the passage.

Claim 2: The evaluated model can accurately identify the most relevant snippet of an online text passage for answering a question about the passage.

Evidence: The Wiki articles used are drawn from highly ranked pages according to Project Nayuki, and are therefore considered high quality. Around 500 articles on various subjects are used as passages for questions. 100,000 questions are used. The questions are original, human-constructed by vetted crowdsource workers, and are designed to be difficult. Both a strict metric (exact match accuracy) and a fuzzy metric (F1 overlap score) are used to determine model and human performance. A human baseline from vetted crowdworkers is included, which also serves as a measure of human annotator agreement on the correct answer to passages. The authors analyze answer types and find that they are varied, such as dates, locations, and quantities. A ‘dumb’ model (logistic regressor), capable of only surface-level pattern matching, is shown to perform much worse than the human baseline.

Validity of Claim from Evidence:

1. Content Validity

- *Strength:* A relatively large number of high-quality Wiki articles. A large number and variety of questions and answers.
- *Weakness:* Only high-quality encyclopedic text is used for passages, excluding many other types of online text (e.g. fiction, poetry, news articles, product catalogs, etc.).
- *Suggestions:* In order to cover content validity for online text (i.e. essentially any text), the benchmark would need to be augmented with many other types of texts, since high-quality Wiki articles are, when considered among ‘any text’, a quite niche subject/format/style. Hence, SQuAD 1.0 falls far short of solid coverage in this case.

2. Criterion Validity

- *Strength:* Annotator agreement on the answer snippets to passages is high and, therefore, demonstrates the answers to the questions are indeed correct/reliable. Crowdsource workers are also vetted before being accepted as annotators, strengthening this point. As expected, human beings perform significantly better than a logistic regressor on SQuAD 1.0. This demonstrates that simple, surface-level pattern matching (the only thing a simple logistic regressor is capable of) is insufficient for solving questions on SQuAD 1.0.
- *Weakness:* No evidence provided for predictive validity.
- *Suggestions:* A comparison of model performance on the benchmark to a downstream task, e.g., rate of question-answer completion of crowdsource worker with and without access to the model as a tool.

3. Construct Validity

- *Strength:* The structure of the task exactly matches that of the object of the claim: identifying the most relevant snippet of a passage to answer a given question. The fact that a simple sliding window method and logistic regressor significantly underperform a human at the task indicates the benchmark is able to tell apart rudimentary, surface-level pattern matching (i.e. the sliding window and logistic regressor’s capabilities) from

human-level comprehension and answering on the task. This serves as basic discriminant validity.

- *Weakness:* The paper doesn't provide comparisons with other snippet-based question-answering datasets, or other web text-based question-answering datasets in general, which might measure a similar construct.
- *Suggestions:* The construct validity of the benchmark could be reinforced with more and more detailed comparisons of model performance on SQuAD v1.0 with other span-based question-answering benchmarks, or other web text-based question-answering benchmarks in general.

4. External Validity ✖

- *Strength:* N/A
- *Weakness:* No evidence provided.
- *Suggestions:* When considering the much broader setting of 'any text', evidence of generalization and consistency of performance across varied settings is much more important. Since the SQuAD 1.0 paper doesn't provide direct evidence of this, we must count this validity insufficient. One way to address this validity type would be to perform small-scale experiments with question-answering (via identifying snippets of a given passage) on particularly rare or odd text, and see how well it matches up with the ranking and performance of models on the main benchmark. In fact, something to this effect was done in a later paper with adversarial text (cite), and it revealed that model performance was significantly lower on the adversarial text compared to the main benchmark.

5. Consequential Validity ✖

- *Strength:* The authors mention that the size of the question and answer bank allows it to double as a large source of high-quality training data for question-answering systems, thereby offering a unique advantage in advancing model performance at this task.
- *Weakness:* The authors don't consider potential downstream harms of only assessing performance on encyclopedic text when claiming performance on online text in general. For instance, models that perform well on non-fiction text may do poorly on fiction text. Unaware of this, a teacher may provide the 'high performing' model as a study aid for a fantasy novel, thereby harming educational outcomes for her students.
- *Suggestions:* In this case, including consideration of the potential ramifications of the poor content and external validity of the benchmark would be important for adequately addressing the consequences of the use of the benchmark. Likewise, the remedy to these concerns would be addressing and improving the content and external validity of the benchmark, since the downstream adverse impacts originate from these flaws in the benchmark design.

Object of Claim: Human-level reading comprehension.

Claim 3: The model exhibits human-level reading comprehension.

Evidence: The Wiki articles used are drawn from highly ranked pages according to Project Nayuki, and are therefore considered high quality. Around 500 articles on various subjects are used as passages for questions. 100,000 questions are used. The questions are original, human-constructed by vetted crowdsource workers, and are designed to be difficult. Both a strict metric

(exact match accuracy) and a fuzzy metric (F1 overlap score) are used to determine model and human performance. A human baseline from vetted crowdworkers is included, which also serves as a measure of human annotator agreement on the correct answer to passages. The authors analyze answer types and find that they are varied, such as dates, locations, and quantities. A ‘dumb’ model (logistic regressor), capable of only surface-level pattern matching, is shown to perform much worse than the human baseline.

Validity of Claim from Evidence:

1. Content Validity

- *Strength:* A relatively large number of high-quality Wiki articles. A large number and variety of questions and answers in the passage+question+snippet answer format.
- *Weakness:* The same content pitfalls of claim #2 apply here too: high-quality encyclopedic text is quite a niche category, and it doesn’t include many other major types of text, such as fiction, poetry, cooking recipes, etc. An important type of answer is also missing from the benchmark’s content (more to say on this in structural validity): the answer of ‘there isn’t enough information’. Knowing when you don’t know is a critical part of human-level reading comprehension, but this type of question is never asked in the SQuAD v1.0 benchmark.
- *Suggestions:* Similar to the previous claim, the inclusion of other types of text would be an important way to address content shortcomings. In addition, adding a greater diversity of question content (particularly, questions where there isn’t enough information provided to answer the question), as was done in SQuAD v2.0, would be critical to addressing content shortcomings. More plausibly content-related issues are addressed under construct validity, particularly structural validity, as those concerns more appropriately fall under structure.

2. Criterion Validity

- *Strength:* Annotator agreement on the answer snippets to passages is high and, therefore, demonstrates the answers to the questions are indeed correct/reliable. Crowdsource workers are also vetted before being accepted as annotators, strengthening this point. As expected, human beings perform significantly better than a logistic regressor on SQuAD v1.0 under matched conditions. This demonstrates that simple, surface-level pattern matching (the only thing a simple logistic regressor is capable of) is insufficient for solving questions on SQuAD 1.0.
- *Weakness:* No evidence provided for predictive validity.
- *Suggestions:* Although there are many shortcomings with SQuAD v1.0 as a benchmark for the given claim, when considering criterion validity isolated from other major issues (e.g. content, construct, and external validity), we see preliminary evidence that the evaluation results do coincide with a validated standard (mainly confirmation of expected results from human test-takers and ‘dumb’ models). Still, given the other shortcomings of the benchmark, the lack of predictive validity weighs more on the inadequacy of the overall criterion validity, leading it to be rated lower than for the other claims. One way to address this weakness is to, similar to previous suggestions, collect evidence on the usefulness of a high performing model as a tool to assist student reading comprehension by, say, measuring test-taker reading comprehension scores on an assessment with and without access to the model.

3. Construct Validity ✖

- *Strength:* One part of human-level reading comprehension involves being able to select the appropriate snippet of a passage to answer a given question, and SQuAD v1.0 covers the structure of this scenario well. The fact that a simple sliding window method and logistic regressor significantly underperform a human indicates the benchmark is able to tell apart rudimentary, surface-level pattern matching (i.e. the sliding window and logistic regressor’s capabilities) from human-level comprehension and answering on the task. This serves as basic discriminant validity.
- *Weakness:* Many aspects of the structure of human-level reading comprehension are unaccounted for. For example, open-ended short-response is not a type of capability tested, nor multiple-choice selection, nor decision-making based on external information, despite these activities being key aspects and expressions of human-level reading comprehension. Related to concerns of structural validity, there are important related but different constructs to human-level reading comprehension that SQuAD v1.0 is unable to tell apart. For example, human-level reading comprehension requires being able to synthesize an original response that isn’t contained in the passage or question provided. But, SQuAD v1.0 would be unable to tell apart a model capable of ‘original synthesis’ compared to a model that is merely capable of a ‘lesser’ reading comprehension and can only identify the correct answer if it sees it in the passage (i.e. is directly contained in a snippet of the passage). The paper also doesn’t provide comparisons with other reading comprehension benchmarks.
- *Suggestions:* Many of the risks to construct validity of this benchmark stem from the structural invalidities, particularly the lack of a variety of ‘answering paradigms’ assessed. Adding question formats that cover a broader range of reading comprehension tasks would help remedy this. In addition, providing examples/analysis of convergent validity i.e. comparisons of model results on SQuAD v1.0 to other reading comprehension benchmarks.

4. External Validity ✖

- *Strength:* N/A
- *Weakness:* No evidence provided.
- *Suggestions:* Given reading comprehension is such a general ability/area, external validity plays a key role in the overall validation of a human-level reading comprehension benchmark since it lends confidence to the fact that results from an evaluation will generalize to various unseen cases. For example, evidence of high performing models also being able to correctly answer questions on ‘trick’/adversarial text or oddly formatted text would be important in building confidence in the generalizability of reading comprehension performance based on SQuAD performance. Since there isn’t clear, direct evidence of this kind provided in the paper, there is little support for the claim’s external validity.

5. Consequential Validity ✖

- *Strength:* The authors mention that the size of the question and answer bank allows it to double as a large source of high-quality training data for question-answering systems, thereby offering a unique advantage in advancing model performance at this task.

- *Weakness:* Similar qualms as for claim #2: the weakness of content, construct, and external validity in particular make relying on this benchmark to assess human-level reading comprehension potentially harmful. For example, trusting a ‘high-performing’ model to serve as a tutor for struggling students in an English reading class could lead to miseducation of those students due to, say, the model being incapable of sufficiently good comprehension of non-encyclopedic information (like a fantasy novel).
- *Suggestions:* In this case, including consideration of the potential ramifications of the poor content, construct, and external validity of the benchmark would be important for adequately addressing the consequences of the use of the benchmark. Likewise, the remedy to these concerns would be addressing and improving the content, construct, and external validity of the benchmark, since the downstream adverse impacts originate from these initial flaws in the benchmark design.