

---

# Position: Evaluations Should Acknowledge Model Multifacetedness in the Era of Large Language Models

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

1 The rapid evolution of Artificial Intelligence (AI), particularly Large Language  
2 Models (LLMs), marks a significant departure from earlier machine learning (ML)  
3 paradigms. This advancement has exposed critical misconceptions in our under-  
4 standing of the “model” itself, especially evident in evaluation methodologies  
5 that often rely on narrow observational windows to assess overall model quality.  
6 This paper argues that a fundamental reconceptualization of the “model” itself is  
7 necessary to address this evaluative crisis. We introduce a five-tiered hierarchical  
8 framework. Specifically, we divide models into: *Noumenal*, *Conceptual*, *Instan-*  
9 *tiated*, *Reachable*, and *Observable* ones. Using this framework, we examine the  
10 historical development of how models have been conceptualized and evaluated  
11 within the ML field, analyzing the roles of experiments, ablation studies, and  
12 datasets. The paper further argues that LLMs’ current development fundamen-  
13 tally challenges these long-standing evaluation patterns, as existing benchmarks  
14 and metrics increasingly fail to capture the true capabilities and limitations of  
15 these complex models. Our primary contribution is to consolidate and structure  
16 many of these historical insights and evolving challenges. By organizing these  
17 often fragmented pieces of understanding into the proposed five-tiered hierarchical  
18 framework, we aim to offer a more cohesive and systematic lens for approaching  
19 AI model evaluation. We believe that such a structured approach, which encourages  
20 assessment strategies to be explicitly contextualized by a model’s position within  
21 this hierarchy and informed by its preceding layer, can help cultivate a more robust  
22 and meaningful comprehension of these increasingly complex LLM systems.

23 

## 1 Introduction

24 Artificial intelligence has undergone several phases of rapid advancement. Yet, the recent emergence  
25 and widespread adoption of Large Language Models mark a fundamental shift [11, 17, 33, 36, 132].  
26 This development significantly challenges established approaches, not only in how AI is created  
27 but, crucially, in how it is evaluated [5, 21, 67, 71, 110, 127]. The sophisticated capabilities of  
28 contemporary AI models have surpassed the existing conceptual and methodological tools previously  
29 used to understand and evaluate ML systems [11]. This paper will explore the characteristics of this  
30 significant change, by proposing a hierarchical perspective for assessing AI models, to help navigate  
31 the present difficulties in LLMs’ evaluation. Historically, earlier ML systems were typically designed  
32 for clearly defined, narrow tasks, such as classifying images or detecting spam [61]. While their  
33 internal structures could be complex, they were often more transparent. Evaluation metrics could  
34 frequently provide a direct measure of the model’s usefulness for its intended function. For example,  
35 the accuracy of a classification model was a relatively clear indicator of its performance [10]. In  
36 contrast, LLMs are moving beyond restricted, task-specific roles towards more general abilities and  
37 often display emergent behaviors [121]. These new emergent skills—such as learning from examples

38 within the prompt (in-context learning [16]), step-by-step reasoning (chain-of-thought [122]), and  
39 even behaviors that appear to be creative or strategic [45, 128], were not the primary goals of their  
40 design, nor were they easily observable in older ML models. Often, the full range of their potential  
41 behaviors is not known in advance, even by their creators [11]. The appearance of unexpected  
42 abilities signifies that a system’s overall properties cannot be fully understood simply by examining  
43 its individual components or its initial design specifications [3].

44 Currently, the most common approach to assessing LLMs is through what might be called a “small  
45 observation window”. This method usually involves testing models on standardized benchmarks.  
46 However, by their very nature, these benchmarks can only examine a tiny fraction of a model’s  
47 potential range of behaviors [96]. Such limited observation can lead to an erroneous understanding  
48 regarding a model’s true quality, its ability to generalize to new situations, and its potential risks,  
49 including safety concerns [49]. For instance, high performance on a specific benchmark might result  
50 from issues like data contamination [72] (where test data was accidentally part of the training data) or  
51 it might simply show that the model has overfitted to the particular behavior of the benchmark tasks,  
52 rather than demonstrating a genuinely robust and widely applicable capability. **Therefore, evaluating**  
53 **LLMs as if they were merely more powerful versions of traditional, fully understandable ML**  
54 **models is a fundamental error in categorization** [101]. An over-reliance on these narrow evaluation  
55 windows can inadvertently create a superficial or misleading impression of understanding, akin to  
56 a “simulacrum” [6]. In this situation, reported benchmark scores can become disconnected from  
57 the model’s actual abilities. They may represent a performance specifically *manufactured* for that  
58 benchmark, rather than an intrinsic, generalizable quality of the model itself. This practice can result  
59 in a hyperreal [125] assessment environment within the research community, where the benchmark  
60 score is treated as more significant or *real*. Such a scenario risks skewing the research agenda towards  
61 optimizing performance on these limited benchmarks, rather than pursuing a more comprehensive  
62 understanding or development of AI capabilities.

63 While much of the existing literature has concentrated on the design and refinement of evaluation  
64 benchmarks [13, 21, 22, 25, 26, 52, 59, 64, 66, 76, 86, 93, 96, 98, 99, 108, 110, 116, 130, 133],  
65 this paper seeks to complement these efforts by focusing on the underlying conceptualization of  
66 the *model* itself. We observe that effectively addressing the current challenges in evaluating LLMs  
67 **benefits from a clearer and more structured understanding of what constitutes a “model” in**  
68 **this evolving landscape**. Our work aims to synthesize various perspectives by proposing that models  
69 can be understood across multiple levels of abstraction and concrete realization. To this end, we  
70 introduce a five-tier hierarchical framework, which forms the conceptual backbone of this paper:

- 71 • **Noumenal Model:** The ultimate, and perhaps inherently unknowable, generative principles  
72 or reality that the AI system is intended to approximate or capture.
- 73 • **Conceptual Model:** The intended design, underlying theories, and architectural blueprints.
- 74 • **Instantiated Model:** The actual implemented algorithmic artifact with an initialization state.
- 75 • **Reachable Model:** The optimized model, with the full spectrum of its potential behaviors.
- 76 • **Observable Model:** A subset of behaviors that are actually witnessed during specific evalua-  
77 tion procedures and interactions.

78 The subsequent sections will discuss the definition of these five tiers, drawing inspiration from  
79 established traditions of modeling and abstraction in both philosophy and science [37]. We will then  
80 trace the historical development and relevance of these conceptual layers within the field of machine  
81 learning, highlighting how the emergence of LLMs overturns the interrelation between these layers.  
82 Finally, we propose that robust assessment should involve a more deliberate and structured approach,  
83 where evaluations conducted at a specific model tier are explicitly defined and constrained by an  
84 understanding of the preceding, more fundamental tier, contributing to a more systematic framework  
85 for meaningful, comprehensive, and reliable evaluations of LLM systems.

## 86 2 Hierarchical Ontology of Models

87 To navigate the complexities of modern AI systems, particularly LLMs, it is proposed that the very  
88 concept of “model” be deconstructed and reassembled into a hierarchical ontology. Before the  
89 detailed definitions, we first briefly introduce our inspiration drawn from established principles in  
90 system theory, cognitive science, AI planning, and machine learning.

91 Hierarchical analysis of complex systems is a well-established paradigm. General system theory [12, 114] and hierarchy theory [106] explain that layered structures improve understanding of system 92 components, interactions, scales, and observer roles, helping manage complexity due to differing 93 interactions and emergent properties across levels. Such analysis is also applied in cognitive science 94 for AI frameworks. For example, David Marr proposes three levels for understanding information- 95 processing systems, the computational theory (goal and logic of computation), representation and 96 algorithm (how the theory is implemented), and hardware implementation (physical realization) [75]. 97 Similarly, Allen Newell distinguished between the knowledge level, which describes a system in terms 98 of its goals and the knowledge it rationally employs, and the symbol level, which details the specific 99 symbolic representations and processes that mechanize this knowledge [84, 85]. These frameworks 100 highlight the importance of differentiating between abstract purpose, procedural specification, and 101 concrete instantiation. Furthermore, fundamental concepts from ML also inform our hierarchical 102 view. For instance, the initial design of an AI system often implicitly defines a hypothesis space 103 for all the possible functions or solutions [10, 81], thus an algorithm could then search this space 104 within a specific instantiation. Crucially, these ML theories also emphasize the distinction between a 105 model’s true generalization capability (on unseen data) and its observed performance on finite test 106 sets. Building upon these diverse theoretical foundations, our proposed five-tier hierarchical ontology 107 aims to provide a specialized framework tailored to the nuances of modern AI systems, particularly 108 LLMs, and the challenges they pose for evaluation.

110 **Definition 2.1 ( $\mathcal{M}_N$ : *Noumenal*)** *The Noumenal Model represents the ultimate, perhaps intrinsi- 111 cally unknowable, generative principles or the “true” underlying structure of the reality that an AI 112 system aims to capture or approximate. The Noumenal Model is the ideal form of knowledge or the 113 perfect causal understanding of a domain.*

114 Philosophically, this concept draws inspiration from Immanuel Kant’s notion of the *noumenon* or 115 *thing-in-itself* [57], particularly his distinction between *phenomena* and *noumena* (*Critique of Pure 116 Reason*, A235/B294–A260/B315). We can conceptualize a theoretical machine learning model that 117 remains fundamentally unrecognizable to human beings, and which we can only ever imperfectly 118 apprehend through phenomena [73]. Such a model would not be a *black box* whose mechanisms are 119 too complex for us to trace, but rather one whose fundamental operational principles and cognitive 120 architecture have no common standard of human thoughts and empirical observation.

121 On one hand, the existence of  $\mathcal{M}_N$  is in the fundamental assumptions in the philosophy of science, 122 which posit an objective reality governed by (perhaps not fully) discoverable and comprehensible 123 natural laws (through systematic observation and experimentation). On the other hand, though 124 *wholly unknowable*, recognition of a  $\mathcal{M}_N$  carries practical weight, compelling critical examination 125 of AI’s fundamental goals. For instance, contemporary LLMs are primarily trained to predict the 126 next token in a sequence, implicitly adopting the data’s statistical distribution as their learning target. 127 However, if  $\mathcal{M}_N$  truly incorporates profound principles such as “core knowledge” [60, 109] or 128 “causal structures” [90], then merely mimicking surface-level statistical patterns in data may be 129 insufficient, resulting in the *brittleness* of LLMs. Consequently, holding the idea that any scientific 130 system can only provide an approximation of the  $\mathcal{M}_N$ , encourages a re-evaluation of AI’s ultimate 131 objectives and the methodologies used for designing the learning tasks.

132 **Definition 2.2 ( $\mathcal{M}_C$ : *Conceptual*)** *The Conceptual Model comprises the intended design and speci- 133 fied architecture, underlying theory and theoretical assumptions, chosen algorithms and blueprint 134 of the system, and finally, the high-level goals the system is meant to achieve, as envisioned by its 135 human creators.*

136 Following the Kantian inspiration, the human mind actively structures experience through a priori 137 categories of understanding (e.g., causality, unity) to make sense of the phenomenal world (*Critique of Pure 138 Reason*, B1-B2, A70/B95-A83/B109).  $\mathcal{M}_C$ , therefore, imposes a conceptual structure onto 139 a problem domain or desired functionality, from the observed phenomena. More specifically, it 140 contains i) the system’s high-level objectives (e.g., the form of loss functions), ii) the theoretical 141 assumptions guiding its operation (e.g., assumptions about the data, learning processes), iii) the 142 selected algorithms and data structures, iv) the overall formal description of the system which act as 143 Kantian schemata that mediate between pure concepts and observations.

144  $\mathcal{M}_C$  is a necessary abstraction (e.g., “attention”), with logic formalizing it in AI systems (e.g., “is all 145 you need” [112]). The logical framework enables structured human thought to engage with complex

146 realities, allowing designers to specify an AI’s intended knowledge, reasoning, and behaviors [14, 147 100]. Although the logical formalisms of the abstracted  $\mathcal{M}_C$  may not fully predict or constrain the 148 complex behaviors of these systems in operation (especially, LLMs’ actual behaviors can largely 149 diverge from an (expectative) logical rigor design, see Section 3). Nevertheless, acknowledging 150 its limitations does not diminish the importance of  $\mathcal{M}_C$ ; it constitutes the logical starting point, 151 becoming the vitally important reference benchmark for evaluating behavior deviation, diagnosing 152 system failures, and understanding unexpected problems.

153 **Definition 2.3 ( $\mathcal{M}_I$ : Instantiated)** *The Instantiated Model refers to the actual, concrete algorithmic 154 artifact that has been implemented in code and exists as a computational entity, encompassing the 155 specific implementation of algorithms, the precise network architecture, the initialized parameter 156 values, and the exact software and hardware environment in which the model operates.*

157 We *intentionally* define the concept of *initialization parameters* more vaguely and expansively, 158 encompassing a potential pre-training phase (at any specific checkpoint, but before task-specific 159 fine-tuning), not just a single random initialization of an established network. This is because 160 the initialization scheme itself also constitutes a concrete instantiation of the  $\mathcal{M}_C$ ’s abstracted 161 content. For instance, a neural network could be initialized (and further optimized) randomly [102], 162 orthogonally [54], or self-supervisedly with a large-scale dataset [16, 31]. These initialized parameter 163 values define the model’s specific state at a particular stage, directly influencing its subsequent 164 learning trajectory and potential capabilities (of the Reachable Model). For instance, a pre-trained 165  $\mathcal{M}_I$  can be highly structured, with parameters encoding significant general-purpose knowledge and 166 representations. Indeed, parameters taken from any specific checkpoint during or after a training 167 process also define a distinct  $\mathcal{M}_I$ , a snapshot of its learned state. However, it is crucial to distinguish 168  $\mathcal{M}_I$  from merely *a pre-trained model*, despite being a key example due to their structured initial 169 parameters.  $\mathcal{M}_I$  more broadly signifies the model’s tangible, concrete configuration at any defined 170 starting point that serves as the foundation prior to the specific optimization process designed to 171 evolve it towards its Reachable counterpart.

172 Furthermore, the specific characteristics of  $\mathcal{M}_I$  play a crucial role in constraining and shaping the 173 subsequent Reachable Model. The journey from the  $\mathcal{M}_C$  (e.g., the idea of attention mechanism) to the 174  $\mathcal{M}_I$  (e.g., the specific code with initial weights of a Transformer) involves numerous design choices 175 and initial conditions. Small variations in architecture or minor differences in initialization can send 176 the model down different optimization paths, leading to distinct Reachable Models ( $\mathcal{M}_R$ ) with varying 177 capabilities and biases. This is a critical juncture, as these early decisions and their non-obvious 178 influences on the model’s development represent the first steps in *a gradual departure from the original* 179 *concept*, significantly contributing to the well-known “black-box” problem [68]. Nevertheless, gaining 180 a better understanding of the  $\mathcal{M}_I$ ’s intrinsic properties (its architecture, representational style, and 181 initial state) is critical for anticipating the characteristics of the final, trained Reachable Model.

182 **Definition 2.4 ( $\mathcal{M}_R$ : Reachable)** *The Reachable Model is the Instantiated Model after its optimization 183 on a specific learning dataset (i.e., the set of finalized learned parameters). More broadly, it 184 encompasses the full spectrum of potential behaviors and internal stochastic processes (e.g., sampling 185 strategies) that the optimized model could exhibit across all possible valid inputs.*

186 In general,  $\mathcal{M}_R$  signifies more than just a *post-trained model*. While the “Reachable” materializes 187 after an optimization process acting upon an  $\mathcal{M}_I$ , its defining characteristic is the representation of 188 the model’s complete potential capabilities, a direct consequence of its specific learned parameters. 189 Thus, the focus is on this entire accessible behavioral repertoire, rather than merely the model’s 190 status as having completed a training phase. While  $\mathcal{M}_R$  represents the totality of what the model can 191 ultimately do, much of this capacity may not be immediately apparent from its static components 192 or the original design intentions. Meanwhile, the inscrutable nature of the training process further 193 intensifies the departure of  $\mathcal{M}_R$  from the initial concept. Consequently,  $\mathcal{M}_R$  becomes more akin to 194 what is typically understood as a “black-box model.” Furthermore, it is within  $\mathcal{M}_R$  that *emergent* 195 *abilities* manifest, which were not explicitly designed into  $\mathcal{M}_C$  nor readily predictable from the 196  $\mathcal{M}_I$  alone, but arise from the interplay of scale, data, and the optimization process. Such behavior 197 is indeed central to the essence that the term “black box” seeks to embody, while significant prior 198 research in this domain has already been dedicated to understanding  $\mathcal{M}_R$ . Examples include work 199 on adversarial testing [43], red-teaming [39], and frameworks for predicting emergent abilities [121]. 200 This underscores that critical aspects like AI safety [49] and alignment [4] are, at their core, attributes

201 of  $\mathcal{M}_R$ , necessitating evaluation strategies far more comprehensive than current standard practices  
202 and equipped to grapple with its inherent complexity and opacity.

203 **Definition 2.5 ( $\mathcal{M}_O$ : *Observable*)** *The Observable Model constitutes the subset of the Reachable*  
204 *Model’s behaviors that are actually witnessed, measured, and documented through available/existing*  
205 *evaluation protocols, datasets, and metrics. The Observable Model is the empirical manifestation of*  
206 *the AI systems’ performance under particular inspection.*

207 The observable manifestation is precisely what current AI benchmarks aim to capture, for instance, in  
208 the natural language process scenarios, we use MMLU [48, 120] for general knowledge, GLUE [116]  
209 and SuperGLUE [115] for natural language understanding, and more comprehensive frameworks  
210 like HELM [66]. However, the critical issue is that the choice of what to observe profoundly shapes  
211 our perception of an AI’s capabilities. This is because how convincing (*plausibility* [32, 79]) an  
212 explanation of an observed behavior is to a human user is often based on interactions with, and  
213 interpretations of, the Observable Model. For example, if an LLM is observed to perform well  
214 on simple problems presented in a benchmark but fails on more complex versions of the same  
215 underlying task, then we probably recognize this LLM as having only primary capabilities on this  
216 task, which can be a total misunderstanding about the potentiality resided within the Reachable  
217 Model. Unfortunately, essentially, even though current benchmarks have been working hard on  
218 providing a better observation window. For instance, HELM strives for "Broad coverage... Multi-  
219 metric measurement... Standardization" to improve how the Observable Model is captured. They  
220 still need to *explicitly* acknowledge the inherent incompleteness of any such observation. In this way,  
221 the Observable Model can become a skewed or unrepresentative sample of the Reachable Model’s  
222 true nature, and optimizing for it does not necessarily translate to the underlying Reachable Model  
223 having improved in a broadly generalizable manner, nor does it guarantee closer approximation to  
224 the Conceptual or Noumenal ideals.

### 225 **3 Evolution of Model Conceptualizations**

226 The conceptualization of the “model” in machine learning has not been static; rather, it has undergone  
227 a continuous process of evolution and enrichment. The hierarchical structure situated above the  
228 Noumenal Model, was not an instantaneous creation, nor did it arise *spontaneously* with current  
229 advanced systems like LLMs; rather, it reflects a gradual process of differentiation. Specifically, when  
230 an AI system has significantly expanded the scope of its capabilities and conceptual complexity, a  
231 more concrete model tier would be “crystallized” from the lower one. Below, we will demonstrate this  
232 change through a rough definition of  $>$  and  $\simeq$  between models of different tiers. Briefly, Model Tier  
233 A  $>$  Model Tier B (A is broader/encompasses B) signifies A is more fundamental, B is a constrained  
234 version or subset of A, and the A-to-B transition involves reduction or constraint. Model Tier A  $\simeq$   
235 Model Tier B (A is similar/equal to B) signifies no significant *practical gap* between them; they  
236 largely capture each other reciprocally, and transitioning between them doesn’t substantially alter  
237 information or their core nature.

238 **Differentiation from the Conceptual Model:** For models such as Naive Bayes and Decision Trees,  
239 which possess relatively simple structures and clear theoretical underpinnings, their hierarchical  
240 relationship can be expressed as:  $\mathcal{M}_N > \mathcal{M}_C \simeq \mathcal{M}_I \simeq \mathcal{M}_R \simeq \mathcal{M}_O$ . For instance, if the ultimate  
241 true principles of the target domain, *e.g.*, the true biological mechanisms for disease prediction,  
242 are represented in  $\mathcal{M}_N$ . But, the Naive Bayes classifier based on selected features for disease  
243 prediction [46] represents a simplified concept, capturing only a limited, abstracted view of the  
244 observation, often with strong independence assumptions [10]. Meanwhile, these simpler models’  
245  $\mathcal{M}_C$  can be generally translated into  $\mathcal{M}_I$ ’s implementation faithfully, since there are fewer degrees  
246 of freedom that would lead to significant deviations. For instance, the recursive partitioning logic  
247 and splitting criteria for decision tree algorithms like ID3 or C4.5 [81, 95] strictly follow the concept  
248 of a tree structure. The training process then fully determines  $\mathcal{M}_I$ ’s final form and behavior, such  
249 as calculating conditional probabilities for Naive Bayes from data, or selecting splits and growing  
250 branches for a decision tree. Since these models operate based on explicit, inspectable rules or clearly  
251 defined probabilistic inferences [82], the space of potential outputs for any given input is constrained  
252 and directly calculable from the  $\mathcal{M}_R$ . Finally, due to this deterministic and transparent nature,  $\mathcal{M}_R$ ’s  
253 full spectrum of potential behaviors can be comprehensively captured by standard evaluation metrics  
254 (*e.g.*, precision, recall, F1-score, ROC curves) on representative test sets. Therefore, we conclude that  
255  $\mathcal{M}_O$  derived from such evaluations is thus a reliable and sufficiently complete representation of  $\mathcal{M}_I$   
256 and  $\mathcal{M}_R$ ’s capabilities and limitations for the defined problem scope.

257 **Differentiation from the Instantiated Model:** For models like K-Nearest Neighbors (KNN), Support  
258 Vector Machines (SVM), and Linear Regression, their hierarchical relationship shows a subtle  
259 change:  $\mathcal{M}_N > \mathcal{M}_C > \mathcal{M}_I \simeq \mathcal{M}_R \simeq \mathcal{M}_O$ . While  $\mathcal{M}_N$  of ultimate reality transcends any  
260 human-designed  $\mathcal{M}_C$ , a key distinction is that the theoretical ideals within the  $\mathcal{M}_C$  are more  
261 abstract than their practical implementation. For instance, SVM's maximum-margin hyperplane  
262 and the kernel trick [10, 28], Linear Regression's best-fitting plane achieved by minimizing a loss  
263 function [46], or KNN's neighbor-based decision principle [29, 46]. This separation occurs because  
264 instantiation necessitates specific, constraining choices that are not contained in underlying concepts,  
265 such as particular SVM kernel functions (*e.g.*, RBF, polynomial) and regularization parameter [104],  
266 optimization algorithms and loss functions like SMO [92] for SVMs or gradient descent with L2  
267 regularized mean squared error for linear regression, or defined K-values and distance metrics (*e.g.*,  
268 Euclidean, Minkowski) for KNN. These choices make the implemented algorithmic artifact ( $\mathcal{M}_I$ )  
269 a particular, constrained realization of the broader conceptual theory. Despite this  $\mathcal{M}_C > \mathcal{M}_I$   
270 distinction, once these models are trained and their parameters are finalized (*e.g.*, support vectors  
271 identified, regression coefficients determined, or training samples stored for KNN), their behavior  
272 becomes fully determined by this learned state, since there are generally no further complex emergent  
273 abilities beyond what is directly implied by the chosen structure and learned parameters. Furthermore,  
274 these instantiations, even with specific choices, are still highly structured and predictable. Their  
275 mechanisms are transparent enough (*e.g.*, linear coefficients, support vector locations, distance  
276 calculations) to allow standard evaluation methods to comprehensively capture their performance on  
277 test data, making  $\mathcal{M}_O$  a faithful and reasonably complete representation of  $\mathcal{M}_R$ .

278 **Differentiation from the Reachable Model:** For models like Shallow Neural Networks (Shallow  
279 NN), Multilayer Perceptrons (MLP), and Restricted Boltzmann Machines (RBM), the distinctions  
280 between tiers intensify further, typically expressed as:  $\mathcal{M}_N > \mathcal{M}_C > \mathcal{M}_I > \mathcal{M}_R \simeq \mathcal{M}_O$ . The  
281 gap widens from  $\mathcal{M}_I$  to  $\mathcal{M}_R$ , as this tier critically encompasses not only the specific architectural  
282 implementation (*e.g.*, topology of a three-layer MLP and choice of activation functions) but also  
283 the initial parameter values (*e.g.*, random initializations [41, 47]), which are vital for the subsequent  
284 optimization trajectory as they set the starting point in a complex, non-convex loss landscape [69,  
285 77]. Consequently, the complex optimization process of training a neural network transforms the  
286 initial states ( $\mathcal{M}_I$ ) to ones with significantly different capabilities and behaviors ( $\mathcal{M}_R$ ). Different  
287 initialization seeds [91] or minor variations in the optimization process [20] can lead the network to  
288 converge to different local minima in the loss landscape, resulting in distinct  $\mathcal{M}_R$  even from nearly  
289  $\mathcal{M}_I$  Models [42]. Nevertheless, for these shallow networks, although their internal representations  
290 may begin to exhibit the opacity characteristic of deep learning (*i.e.*, having global non-local sub-  
291 representations) [87], their overall behavioral complexity is generally considered sufficiently bounded.  
292 Practically, it is often assumed that standard, diverse benchmarks and evaluation metrics can still  
293 capture their core capabilities and generalization performance reasonably well [46], making  $\mathcal{M}_O$  a  
294 fair, albeit perhaps not exhaustive, representation of  $\mathcal{M}_R$ 's overall performance.

295 **Differentiation from the Observable Model:** As network depth and complexity increase, Deep  
296 Neural Networks (DNNs) exhibit more intricate hierarchical relationships, summarized as:  $\mathcal{M}_N >$   
297  $\mathcal{M}_C > \mathcal{M}_I > \mathcal{M}_R > \mathcal{M}_O$ . While the distinctions established in shallow NNs persist, a critical  
298 new divergence distinguishing DNNs arises between  $\mathcal{M}_R$  and  $\mathcal{M}_O$ . Due to their vast parameter  
299 counts, deep architectures, and extensive training on large datasets, DNNs learn extremely complex  
300 functions, resulting in an  $\mathcal{M}_R$  with an enormous potential behavioral space not explicitly programmed  
301 nor easily predictable from  $\mathcal{M}_I$  alone. However, our current methods of observation, standard  
302 evaluation protocols and benchmarks such as ImageNet [30], GLUE [116], or even comprehensive  
303 frameworks like HELM [66], can only access a limited subset of this vast behavioral repertoire.  
304  $\mathcal{M}_O$  is frequently reported to fail in fully presenting the true scope of  $\mathcal{M}_R$ 's capabilities. For  
305 instance, models' brittleness is easily demonstrated when faced with out-of-distribution inputs or  
306 slight adversarial paraphrases, which exposes superficial "shortcut" learning rather than robust  
307 understanding [40, 49, 86, 96].

308 **A note on large language models:** For contemporary LLMs, the hierarchical gaps between conceptual  
309 tiers of models are widening dramatically, with the largest and most significant divide occurring  
310 between  $\mathcal{M}_R$  and  $\mathcal{M}_O$ , broadly expressed as:  $\mathcal{M}_N > \mathcal{M}_C > \mathcal{M}_I > \mathcal{M}_R >> \mathcal{M}_O$ . While  
311 significant gaps separate Noumenal model goals (*e.g.*, representing human language, knowledge,  
312 and reasoning) from Conceptual designs and Instantiations (*e.g.*, Transformers [112], Mamba [44]).  
313 The pre-existing distinctions are amplified in  $\mathcal{M}_R$  created by extensive post-training.  $\mathcal{M}_R$  of an

314 LLM exhibits an immensely vast potential behavioral space, featuring (sometimes) unpredictable  
315 emergent abilities like in-context learning, instruction following, and complex reasoning [119, 121].  
316 Concurrently, significant risks are reported, such as generating hallucinations [55], amplifying  
317 biases [7], or producing harmful content [123]. However, the combinatorial nature of language  
318 and the sheer scale of these models create a serious mismatch; what we learn from  $\mathcal{M}_O$  is a very  
319 incomplete picture of an LLM’s true overall abilities and the hidden dangers within its  $\mathcal{M}_R$ . This  
320 mismatch is a fundamental reason for the current problems in testing LLMs, the major challenges in  
321 making them behave safely and as intended, and the troublesome practice of “SOTA chasing”.  
322 This evaluation challenge appears to be significantly compounded by the field’s tendency to rely  
323 on an evaluation paradigm inherited from earlier ML. In those earlier and simpler systems, relative  
324 transparency and tighter coupling between the tiers characterized the confidence in standard metrics  
325 and experimental setups. These approaches became deeply ingrained and are now being somewhat  
326 uncritically applied to LLMs. With LLMs, the relationships between the tiers have become signifi-  
327 cantly more complex, opaque, and divergent. The historical success of these evaluation norms with  
328 simpler models established certain “patternized experiments” and expectations about what constitutes  
329 “good evaluation.” These established practices were then naturally carried over when LLMs emerged,  
330 despite them possessing vastly different characteristics, particularly in the complexity and opacity of  
331 their Instantiated and Reachable tiers. **This “historical muscle memory” from evaluating simpler**  
332 **models, when applied to the new context of LLMs, can be seen as a significant contributor to**  
333 **the current evaluation challenges.** In many ways, the field might have been attempting to navigate  
334 new, complex terrains using guides developed for older, more familiar landscapes.

## 335 4 Towards Hierarchically-Informed Evaluation Strategies

336 The advent of LLMs necessitates a significant re-evaluation of established methodologies for ex-  
337 perimental design, ablation studies, and dataset curation, which form the bedrock of traditional  
338 “patternized experiments.” This situation demands a shift in the overarching goal of LLM evalua-  
339 tion: moving away from rendering final, summative judgments based predominantly on leaderboard  
340 rankings ( $\mathcal{M}_O$ ) towards an ongoing, iterative process of *model cartography*. To achieve this, we  
341 first review how LLMs break the patterns, then we introduce the paradigm that hierarchically probes  
342 different aspects across model tiers, by assessing each in explicit relation to its antecedents. A detailed  
343 proposal for this framework is provided in the appendix.

### 344 4.1 Rethinking Experimental Design, Ablation, and Dataset Curation

345 **Experimental design** aimed to test hypotheses derived from an  $\mathcal{M}_C$  (e.g., the utility of a specific  
346 feature) or to generate an  $\mathcal{M}_O$  by measuring a trained  $\mathcal{M}_I$ ’s performance. While typically, experi-  
347 mental designs for early ML assume variables could be isolated and their effects clearly measured,  
348 this paradigm is inadequate for LLMs. Their profound sensitivity to subtle variations in prompt-  
349 ing [70] and the immense difficulty in controlling for confounding variables when assessing complex,  
350 generative behaviors make controlled experiments challenging. Static benchmarks, forming the  
351 traditional  $\mathcal{M}_O$ , often fail to capture the dynamic and vast capabilities of an LLM’s  $\mathcal{M}_R$ . Moreover,  
352 such benchmarks cannot explicitly probe inter-tier relationships (e.g., assessing the fidelity between  
353 a Conceptual design choice, like an architectural modification for long-context reasoning, and its  
354 actual manifestation in  $\mathcal{M}_I$ ’s representations), or systematically explore hypothesized regions of the  
355 Reachable behavior space to understand emergent capabilities [121] and their operational boundaries.

356 **Ablation studies** have historically served to understand component contributions within an  $\mathcal{M}_I$  to its  
357 Reachable performance, helping to validate or refine  $\mathcal{M}_C$  choices [42]. This reductionist approach,  
358 however, faces severe limitations with LLMs. In these highly complex, non-linear systems, represen-  
359 tations are often distributed, and components (like neurons or attention heads) can be polysemantic,  
360 contributing to multiple functions [34, 88]. Consequently, removing (or changing) a component from  
361 an LLM doesn’t simply isolate its original function; it can yield a fundamentally different system with  
362 altered internal dynamics and potentially different emergent properties. Interpreting such changes  
363 as the “contribution” of the ablated part becomes problematic, akin to challenges in intervening on  
364 complex causal systems [90]. Simple ablations may therefore offer superficial or even misleading  
365 insights into an LLM’s  $\mathcal{M}_I$ ’s functional architecture or the mechanisms generating its Reachable  
366 behaviors. Therefore, more cautious interpretations are essential, and complementary approaches  
367 like perturbation studies (assessing sensitivity to small changes) or influence studies (tracing impacts

368 of training data/features) might offer more reliable, albeit still partial, insights into these deeply  
369 interconnected systems.

370 **Dataset curation** focused on gathering data samples presumed to be representative of some aspect of  
371 the *Noumenal Model*. The goal was to train an  $\mathcal{M}_I$  into a generalizable  $\mathcal{M}_R$ , carefully navigating  
372 the bias-variance tradeoff [10]. However, in the LLM scenarios, the scale and often uncurated nature  
373 of web-derived training corpora make it exceptionally difficult to ensure true representativeness or,  
374 critically, to prevent contamination with data that might overlap with evaluation benchmarks. Such  
375 contamination can grossly inflate  $\mathcal{M}_O$  performance, providing a misleading picture of an LLM’s real  
376 generalization capabilities within its Reachable space [72]. Therefore, datasets should essentially be  
377 tools for characterizing  $\mathcal{M}_R$  itself, verifying alignment between Conceptual and Instantiated tiers.

## 378 4.2 Paradigm of Model Cartography

379 **Beyond Current Benchmarks for the Observable Model:**  $\mathcal{M}_O$ ’s performance and behaviors should  
380 be interpreted not in isolation, but in direct relation to the known (or reasonably estimated) properties  
381 of  $\mathcal{M}_R$ . The key question becomes: Does the observed performance on a benchmark or specific  
382 task reflect a robust and generalizable capability within the broader Reachable space, or is it an  
383 isolated success, perhaps an artifact of the evaluation setup, data contamination, or a highly specific  
384 and narrow competence? This involves probing the gap between observed performance and latent  
385 potential. Therefore, the primary goal is to determine if observed behaviors are indicative of broader,  
386 stable capabilities within the Reachable space or are mere artifacts. To achieve this, we should  
387 consider the following:

- 388 • Utilize metrics like the Model Utilization Index [19] to assess if high Observable performance  
389 stems from robust, general mechanisms (indicating broad Reachable capability) or overused,  
390 narrow circuits (implying fragile, benchmark-specific success).
- 391 • Employ adaptive testing, dynamic benchmarks, and interactive protocols to explore the  
392 Reachable space, especially around areas of success or failure identified in static tests [18].
- 393 • Focus on “construct validity” [2]: investigate if Observable benchmark performance correlates  
394 with diverse, real-world task performance when both theoretically use the same underlying  
395 (Reachable) abilities [78].
- 396 • Test if Observable performance holds under slight perturbations of inputs, changes in prompt-  
397 ing style, or minor variations in context. Robustness suggests the observed behavior taps into  
398 a stable region of the Reachable space [63, 107].

399 **Probing the Reachable Model:** The characteristics of  $\mathcal{M}_R$ , including its potential for beneficial  
400 emergent capabilities or undesirable harmful behaviors, should be assessed by investigating the  
401 properties of  $\mathcal{M}_I$ . Rather than passively waiting for behaviors to manifest in the Observable tier, the  
402 goal is to proactively use interpretability techniques, or formal methods on both  $\mathcal{M}_R$  and  $\mathcal{M}_I$  to  
403 predict and characterize its potential behavioral repertoire. To this end, consider:

- 404 • Develop and apply advanced mechanistic interpretability to probe the  $\mathcal{M}_I$ ’s internals, map  
405 its Reachable behaviors and circuits (beyond input-output analysis), acknowledging field  
406 limitations and progress [105].
- 407 • Employ intervention methods (e.g., from causal inference) to see how internal components  
408 affect Reachable behaviors, going beyond simple ablation to controlled perturbations and  
409 counterfactual analyses [129].
- 410 • Characterize the “behavioral envelope” or “capability manifold” by exploring model responses  
411 to diverse, structured inputs aimed at revealing a wide range of latent skills [8, 65].
- 412 • Analyze training dynamics and learning trajectories, check if the model learns representations  
413 aligned with the instantiated framework, or it finds “shortcuts,” exploiting spurious data  
414 correlations, and developing misaligned internal concepts [94, 97].

415 **Assessing the Instantiated Model:** The properties of the  $\mathcal{M}_I$  should be evaluated for their fidelity to  
416 the specifications, goals, and theoretical underpinnings of  $\mathcal{M}_C$ . This involves asking: How faithful  
417 do interpretations of observed behavior map onto the model’s actual internal computations and  
418 mechanisms? Are there significant or unintended deviations that arose during implementation or  
419 training? The objective is to assess the fidelity of the actual implementation against its original design  
420 specifications and theoretical goals. Achieving this requires considering:

- 421 • Conduct rigorous audits of the implemented architecture and core algorithmic components  
422 (e.g., attention mechanisms, layer structures, activation functions) against the detailed specific-  
423 cations and theoretical assumptions documented in  $\mathcal{M}_C$  [80].
- 424 • Analyze chosen hyperparameters and embedded architectural constraints for consistency with  
425  $\mathcal{M}_C$ ’s design rationale and implicit theoretical foundations [50, 53, 58, 126].
- 426 • Developing techniques to trace the influence of pre-training data [23, 56], initialization  
427 strategies, and architectural components [113] on the model’s ultimate potential (its Reachable  
428 state) is crucial.
- 429 • For pre-trained foundation models, assess if their initial representations and zero-shot ca-  
430 pabilities on relevant basic tasks align with the objectives and knowledge domains of their  
431 conceptual pre-training design [118].

432 **Validating Alignment with the Conceptual Model:**  $\mathcal{M}_C$  itself can be subjected to evaluation by  
433 investigating its coherence and soundness, involving philosophical and theoretical critique: How well  
434 do the theories of language, reasoning, or intelligence embedded in  $\mathcal{M}_C$  align with deeper principles  
435 of true linguistic competence or general intelligence? Does  $\mathcal{M}_C$  adequately represent the problem it  
436 aims to solve or the reality it aims to model? We should therefore consider:

- 437 • Engage in critical analysis of the core concepts (e.g., *understanding, reasoning, creativity, or*  
438 *safety*), ensuring their definitions are adequate and well-grounded approximations of the true  
439 (Noumenal) nature of these complex phenomena [1].
- 440 • Apply principles from the philosophy of information and epistemology, guaranteeing  $\mathcal{M}_C$   
441 explicitly acknowledges its own *level of abstraction*, inherent simplifications, and limitations  
442 concerning the complexity of the Noumenal ideal it seeks to address [38, 124].
- 443 • Subject  $\mathcal{M}_C$  to scrutiny from experts in relevant fields beyond AI, such as cognitive science,  
444 linguistics, philosophy, and ethics, to assess the validity and potential blind spots of its  
445 foundational assumptions [9, 79, 103].

446 **Acknowledging the Noumenal Model in Evaluation Design:** While  $\mathcal{M}_N$  may be unknowable in  
447 its entirety, its consideration should inform evaluation design. This means designing evaluations that  
448 probe for “core knowledge” or fundamental understanding of underlying principles (e.g., intuitive  
449 physics, causality, basic logic), rather than solely testing task-specific pattern matching. The tasks  
450 included in  $\mathcal{M}_O$  should be critically assessed to determine whether they serve as good proxies for the  
451 more fundamental principles believed to constitute  $\mathcal{M}_N$  for a given *domain of intelligence*, regardless  
452 of how they are conceptualized (e.g., Turing Test [111], Winograd Schema Challenge [62] and their  
453 variations [15, 24, 51, 74]).

## 454 5 Conclusion

455 The rapid ascent of LLMs has outpaced traditional methods of understanding and evaluating artificial  
456 intelligence. We argue that a core issue lies in a persistent mis-cognition of the “model” concept  
457 itself, often leading to an over-reliance on narrow, observable behaviors as proxies for overall model  
458 quality and capability. To help address this, we organize the historical insights into a proposed  
459 five-tiered hierarchical framework, distinguishing between the Noumenal (the ultimate generative  
460 principles), Conceptual (the intended design), Instantiated (the algorithmic artifact), Reachable (the  
461 full potential behavior space), and Observable (the witnessed behaviors) models, aiming to offer a  
462 more cohesive and systematic lens for approaching AI model evaluation. We have explored how  
463 LLMs’ vast behavioral potential, particularly when viewed through the lens of historically diverging  
464 conceptualizations of these tiers in machine learning, challenges established experimental patterns  
465 and underscores the value of evolving our evaluation systems. We suggest that developing evaluation  
466 strategies that explicitly consider the relationships between these tiers can lead to more insightful and  
467 robust assessments. Adopting such hierarchically-informed perspectives is not intended to propose  
468 an entirely new paradigm in isolation, but rather to encourage a more nuanced and contextualized  
469 approach by building upon the collective and structured understanding of the field. This way of  
470 thinking endeavors to cultivate a more meaningful comprehension of these complex systems, fostering  
471 responsible innovation and contributing to the development of AI that is beneficial, robust, and aligned  
472 with human values.

473 **References**

474 [1] Dokpesi Timothy ADIDI. Aristotle's concept of telos and artificial intelligence: Exploring  
475 the relevance of classical philosophy to contemporary ai development. *AMAMIHE Journal of*  
476 *Applied Philosophy*, 22(3), 2024.

477 [2] Ahmed Alaa, Thomas Hartvigsen, Niloufar Golchini, Shiladitya Dutta, Frances Dean, In-  
478 ioluwa Deborah Raji, and Travis Zack. Medical large language model benchmarks should  
479 prioritize construct validity. *arXiv preprint arXiv:2503.10694*, 2025.

480 [3] Philip W Anderson. More is different: Broken symmetry and the nature of the hierarchical  
481 structure of science. *Science*, 177(4047):393–396, 1972.

482 [4] Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy  
483 Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, et al. A general language assistant as a  
484 laboratory for alignment. *arXiv preprint arXiv:2112.00861*, 2021.

485 [5] Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy  
486 Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. A multitask, multilingual, multi-  
487 modal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint*  
488 *arXiv:2302.04023*, 2023.

489 [6] Jean Baudrillard and Sheila Faria Glaser. *Simulacra and simulation*, volume 312. University  
490 of Michigan press Ann Arbor, 1994.

491 [7] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On  
492 the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021*  
493 *ACM conference on fairness, accountability, and transparency*, pages 610–623, 2021.

494 [8] Jan Betley, Xuchan Bao, Martín Soto, Anna Sztyber-Betley, James Chua, and Owain  
495 Evans. Tell me about yourself: Llms are aware of their learned behaviors. *arXiv preprint*  
496 *arXiv:2501.11120*, 2025.

497 [9] Mark H Bickhard and Loren Terveen. *Foundational issues in artificial intelligence and*  
498 *cognitive science: Impasse and solution*, volume 109. Elsevier, 1995.

499 [10] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*.  
500 Springer, 2006.

501 [11] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von  
502 Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the  
503 opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

504 [12] Kenneth E Boulding. General systems theory—the skeleton of science. *Management science*,  
505 2(3):197–208, 1956.

506 [13] Samuel R Bowman and George E Dahl. What will it take to fix benchmarking in natural  
507 language understanding? *arXiv preprint arXiv:2104.02145*, 2021.

508 [14] Ronald Brachman and Hector Levesque. *Knowledge representation and reasoning*. Elsevier,  
509 2004.

510 [15] Selmer Bringsjord, Paul Bello, and David Ferrucci. Creativity, the turing test, and the (better)  
511 lovelace test. *The Turing test: the elusive standard of artificial intelligence*, pages 215–239,  
512 2003.

513 [16] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-  
514 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language  
515 models are few-shot learners. *Advances in neural information processing systems*, 33:1877–  
516 1901, 2020.

517 [17] Sébastien Bubeck, Varun Chadrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece  
518 Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general  
519 intelligence: Early experiments with gpt-4, 2023.

520 [18] John Burden, Marko Tešić, Lorenzo Pacchiardi, and José Hernández-Orallo. Paradigms of  
 521 ai evaluation: Mapping goals, methodologies and culture. *arXiv preprint arXiv:2502.15620*,  
 522 2025.

523 [19] Yixin Cao, Jiahao Ying, Yaoning Wang, Xipeng Qiu, Xuanjing Huang, and Yugang Jiang.  
 524 Revisiting llm evaluation through mechanism interpretability: a new metric and model utility  
 525 law. *arXiv preprint arXiv:2504.07440*, 2025.

526 [20] Ernie Chang, Hui-Syuan Yeh, and Vera Demberg. Does the order of training samples mat-  
 527 ter? improving neural data-to-text generation with curriculum learning. *arXiv preprint*  
 528 *arXiv:2102.03554*, 2021.

529 [21] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen,  
 530 Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language  
 531 models. *ACM transactions on intelligent systems and technology*, 15(3):1–45, 2024.

532 [22] Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li,  
 533 Dacheng Li, Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E Gonzalez, et al. Chatbot  
 534 arena: An open platform for evaluating llms by human preference. In *Forty-first International*  
 535 *Conference on Machine Learning*, 2024.

536 [23] Sang Keun Choe, Hwijeon Ahn, Juhan Bae, Kewen Zhao, Minsoo Kang, Youngseog Chung,  
 537 Adithya Pratapa, Willie Neiswanger, Emma Strubell, Teruko Mitamura, et al. What is  
 538 your data worth to gpt? llm-scale data valuation with influence functions. *arXiv preprint*  
 539 *arXiv:2405.13954*, 2024.

540 [24] François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.

541 [25] Pierre Colombo, Nathan Noiry, Ekhine Irurozki, and Stéphan Cléménçon. What are the best  
 542 systems? new perspectives on nlp benchmarking. *Advances in neural information processing*  
 543 *systems*, 35:26915–26932, 2022.

544 [26] Pierre Jean A Colombo, Chloé Clavel, and Pablo Piantanida. Infolm: A new metric to evaluate  
 545 summarization & data2text generation. In *Proceedings of the AAAI conference on artificial*  
 546 *intelligence*, pages 10554–10562, 2022.

547 [27] Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià  
 548 Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. *Ad-*  
 549 *vances in Neural Information Processing Systems*, 36:16318–16352, 2023.

550 [28] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20:273–297,  
 551 1995.

552 [29] Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE transactions on*  
 553 *information theory*, 13(1):21–27, 1967.

554 [30] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-  
 555 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*  
 556 *recognition*, pages 248–255. Ieee, 2009.

557 [31] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training  
 558 of deep bidirectional transformers for language understanding. In *Proceedings of the 2019*  
 559 *conference of the North American chapter of the association for computational linguistics:*  
 560 *human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.

561 [32] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine  
 562 learning. *arXiv preprint arXiv:1702.08608*, 2017.

563 [33] Yogesh K Dwivedi, Nir Kshetri, Laurie Hughes, Emma Louise Slade, Anand Jeyaraj, Arpan Ku-  
 564 mar Kar, Abdullah M Baabdullah, Alex Koohang, Vishnupriya Raghavan, Manju Ahuja, et al.  
 565 Opinion paper:“so what if chatgpt wrote it?” multidisciplinary perspectives on opportunities,  
 566 challenges and implications of generative conversational ai for research, practice and policy.  
 567 *International journal of information management*, 71:102642, 2023.

568 [34] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna  
 569 Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of  
 570 superposition. *arXiv preprint arXiv:2209.10652*, 2022.

571 [35] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,  
 572 Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for  
 573 transformer circuits. *Transformer Circuits Thread*, page 12, 2021.

574 [36] Tyna Eloundou, Sam Manning, Pamela Mishkin, and Daniel Rock. Gpts are gpts: Labor  
 575 market impact potential of llms. *Science*, 384(6702):1306–1308, 2024.

576 [37] Roman Frigg and Stephan Hartmann. Models in Science. In Edward N. Zalta and Uri  
 577 Nodelman, editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab,  
 578 Stanford University, Summer 2025 edition, 2025.

579 [38] Jean-Gabriel Ganascia. Abstraction of levels of abstraction. *Journal of Experimental &*  
 580 *Theoretical Artificial Intelligence*, 27(1):23–35, 2015.

581 [39] Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath,  
 582 Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. Red teaming language  
 583 models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint*  
 584 *arXiv:2209.07858*, 2022.

585 [40] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,  
 586 Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature*  
 587 *Machine Intelligence*, 2(11):665–673, 2020.

588 [41] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward  
 589 neural networks. In *Proceedings of the thirteenth international conference on artificial*  
 590 *intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings,  
 591 2010.

592 [42] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*. MIT  
 593 press Cambridge, 2016.

594 [43] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adver-  
 595 *sarial examples*. *arXiv preprint arXiv:1412.6572*, 2014.

596 [44] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces.  
 597 *arXiv preprint arXiv:2312.00752*, 2023.

598 [45] Daya Guo, Dejian Yang, Huawei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
 599 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in  
 600 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

601 [46] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements*  
 602 *of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.

603 [47] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
 604 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
 605 pages 770–778, 2016.

606 [48] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
 607 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*  
 608 *arXiv:2009.03300*, 2020.

609 [49] Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. Unsolved problems  
 610 in ml safety. *arXiv preprint arXiv:2109.13916*, 2021.

611 [50] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Hee-  
 612 woo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive  
 613 generative modeling. *arXiv preprint arXiv:2010.14701*, 2020.

614 [51] José Hernández-Orallo. *The measure of all minds: evaluating natural and artificial intelligence*.  
 615 Cambridge University Press, 2017.

616 [52] Anas Himmi, Ekhine Irurozki, Nathan Noiry, Stephan Clemenccon, and Pierre Colombo.  
 617 Towards more robust nlp system evaluation: Handling missing scores in benchmarks. *arXiv*  
 618 *preprint arXiv:2305.10284*, 2023.

619 [53] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
 620 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.  
 621 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.

622 [54] Wei Hu, Lechao Xiao, and Jeffrey Pennington. Provable benefit of orthogonal initialization in  
 623 optimizing deep linear networks. *arXiv preprint arXiv:2001.05992*, 2020.

624 [55] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,  
 625 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation.  
 626 *ACM computing surveys*, 55(12):1–38, 2023.

627 [56] Minhao Jiang, Ken Ziyu Liu, Ming Zhong, Rylan Schaeffer, Siru Ouyang, Jiawei Han, and  
 628 Sanmi Koyejo. Investigating data contamination for pre-training language models. *arXiv*  
 629 *preprint arXiv:2401.06059*, 2024.

630 [57] Immanuel Kant. *Critique of pure reason*, volume 6. Minerva Heritage Press, 2024.

631 [58] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon  
 632 Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural  
 633 language models. *arXiv preprint arXiv:2001.08361*, 2020.

634 [59] Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu,  
 635 Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, et al. Dynabench: Rethink-  
 636 ing benchmarking in nlp. *arXiv preprint arXiv:2104.14337*, 2021.

637 [60] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building  
 638 machines that learn and think like people. *Behavioral and brain sciences*, 40:e253, 2017.

639 [61] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444,  
 640 2015.

641 [62] Hector J Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge.  
 642 *KR*, 2012:13th, 2012.

643 [63] Thomas Lew and Marco Pavone. Sampling-based reachability analysis: A random set theory  
 644 approach with adversarial sampling. In *Conference on robot learning*, pages 2055–2070.  
 645 PMLR, 2021.

646 [64] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E  
 647 Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard  
 648 and benchbuilder pipeline. *arXiv preprint arXiv:2406.11939*, 2024.

649 [65] Xin Li and Anand Sarwate. Unraveling the localized latents: Learning stratified manifold struc-  
 650 tures in llm embedding space with sparse mixture-of-experts. *arXiv preprint arXiv:2502.13577*,  
 651 2025.

652 [66] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga,  
 653 Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of  
 654 language models. *arXiv preprint arXiv:2211.09110*, 2022.

655 [67] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic  
 656 human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.

657 [68] Zachary C Lipton. The mythos of model interpretability: In machine learning, the concept of  
 658 interpretability is both important and slippery. *Queue*, 16(3):31–57, 2018.

659 [69] Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-  
 660 parameterized non-linear systems and neural networks. *Applied and Computational Harmonic  
 661 Analysis*, 59:85–116, 2022.

662 [70] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig.  
 663 Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language  
 664 processing. *ACM computing surveys*, 55(9):1–35, 2023.

665 [71] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval:  
 666 NLg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*,  
 667 2023.

668 [72] Inbal Magar and Roy Schwartz. Data contamination: From memorization to exploitation.  
 669 *arXiv preprint arXiv:2203.08242*, 2022.

670 [73] Riya Manna and Rajakishore Nath. Kantian moral agency and the ethics of artificial intelli-  
 671 gence. *Problemos*, 100:139–151, 2021.

672 [74] Gary Marcus. What comes after the turing test. *The New Yorker*, 9, 2014.

673 [75] David Marr. *Vision: A computational investigation into the human representation and process-  
 674 ing of visual information*. MIT press, 2010.

675 [76] Mark Mazumder, Colby Banbury, Xiaozhe Yao, Bojan Karlaš, William Gaviria Rojas, Sudnya  
 676 Diamos, Greg Diamos, Lynn He, Alicia Parrish, Hannah Rose Kirk, et al. Dataperf: Bench-  
 677 marks for data-centric ai development. *Advances in Neural Information Processing Systems*,  
 678 36:5320–5347, 2023.

679 [77] Song Mei, Yu Bai, and Andrea Montanari. The landscape of empirical risk for nonconvex  
 680 losses. *The Annals of Statistics*, pages 2747–2774, 2018.

681 [78] Justin K Miller and Wenjia Tang. Evaluating llm metrics through real-world capabilities. *arXiv  
 682 preprint arXiv:2505.08253*, 2025.

683 [79] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial  
 684 intelligence*, 267:1–38, 2019.

685 [80] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben  
 686 Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model  
 687 reporting. In *Proceedings of the conference on fairness, accountability, and transparency*,  
 688 pages 220–229, 2019.

689 [81] Tom M Mitchell and Tom M Mitchell. *Machine learning*. McGraw-hill New York, 1997.

690 [82] Christoph Molnar. *Interpretable machine learning*. Lulu. com, 2020.

691 [83] Yaswanth Narsupalli, Abhranil Chandra, Sreevatsa Muppirlala, Manish Gupta, and Pawan  
 692 Goyal. Refer: Improving evaluation and reasoning through hierarchy of models. *arXiv preprint  
 693 arXiv:2407.12877*, 2024.

694 [84] Allen Newell. The knowledge level. *Artificial intelligence*, pages 87–127, 1982.

695 [85] Allen Newell, Herbert Alexander Simon, et al. *Human problem solving*. Prentice-hall  
 696 Englewood Cliffs, NJ, 1972.

697 [86] Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela.  
 698 Adversarial nli: A new benchmark for natural language understanding. *arXiv preprint  
 699 arXiv:1910.14599*, 2019.

700 [87] Michael A Nielsen. *Neural networks and deep learning*, volume 25. Determination press San  
 701 Francisco, CA, USA, 2015.

702 [88] Chris Olah, Arvind Satyanarayanan, Ian Johnson, Shan Carter, Ludwig Schubert, Katherine Ye,  
 703 and Alexander Mordvintsev. The building blocks of interpretability. *Distill*, 3(3):e10, 2018.

704 [89] Anselm Paulus, Arman Zharmagambetov, Chuan Guo, Brandon Amos, and Yuandong Tian.  
 705 Advprompter: Fast adaptive adversarial prompting for llms. *arXiv preprint arXiv:2404.16873*,  
 706 2024.

707 [90] Judea Pearl. *Causality*. Cambridge university press, 2009.

708 [91] David Picard. Torch. manual\_seed (3407) is all you need: On the influence of random seeds in  
709 deep learning architectures for computer vision. *arXiv preprint arXiv:2109.08203*, 2021.

710 [92] John Platt. Sequential minimal optimization: A fast algorithm for training support vector  
711 machines. Technical report, Microsoft, 1998.

712 [93] José Pombal, Nuno M Guerreiro, Ricardo Rei, and André FT Martins. Zero-shot benchmarking:  
713 A framework for flexible and scalable automatic evaluation of language models. *arXiv preprint*  
714 *arXiv:2504.01001*, 2025.

715 [94] Xiangyu Qi, Ashwinee Panda, Kaifeng Lyu, Xiao Ma, Subhrajit Roy, Ahmad Beirami, Prateek  
716 Mittal, and Peter Henderson. Safety alignment should be made more than just a few tokens  
717 deep. *arXiv preprint arXiv:2406.05946*, 2024.

718 [95] J. Ross Quinlan. Induction of decision trees. *Machine learning*, 1:81–106, 1986.

719 [96] Inioluwa Deborah Raji, Emily M Bender, Amandalynne Paullada, Emily Denton, and  
720 Alex Hanna. Ai and the everything in the whole wide world benchmark. *arXiv preprint*  
721 *arXiv:2111.15366*, 2021.

722 [97] Yi Ren and Danica J Sutherland. Learning dynamics of llm finetuning. *arXiv preprint*  
723 *arXiv:2407.10490*, 2024.

724 [98] Pedro Rodriguez, Joe Barrow, Alexander Miserlis Hoyle, John P Lalor, Robin Jia, and Jordan  
725 Boyd-Graber. Evaluation examples are not equally informative: How should that change nlp  
726 leaderboards? In *Proceedings of the 59th Annual Meeting of the Association for Computational*  
727 *Linguistics and the 11th International Joint Conference on Natural Language Processing*  
728 (*Volume 1: Long Papers*), pages 4486–4503, 2021.

729 [99] Mark Rofin, Vladislav Mikhailov, Mikhail Florinskiy, Andrey Kravchenko, Elena Tutubalina,  
730 Tatiana Shavrina, Daniel Karabekyan, and Ekaterina Artemova. Vote’n’rank: Revision of  
731 benchmarking with social choice theory. *arXiv preprint arXiv:2210.05769*, 2022.

732 [100] Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. pearson, 2016.

733 [101] Gilbert Ryle and Julia Tanney. *The concept of mind*. Routledge, 2009.

734 [102] Wouter F Schmidt, Martin A Kraaijveld, Robert PW Duin, et al. Feed forward neural networks  
735 with random weights. In *International conference on pattern recognition*, pages 1–1. IEEE  
736 Computer Society Press, 1992.

737 [103] Jordan Richard Schoenherr. *Ethical artificial intelligence from popular to cognitive science:  
738 trust in the age of entanglement*. Routledge, 2022.

739 [104] Bernhard Schölkopf and Alexander J Smola. *Learning with kernels: support vector machines,  
740 regularization, optimization, and beyond*. MIT press, 2002.

741 [105] Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lindsey, Jeff Wu, Lucius Bushnaq, Nicholas  
742 Goldowsky-Dill, Stefan Heimersheim, Alejandro Ortega, Joseph Bloom, et al. Open problems  
743 in mechanistic interpretability. *arXiv preprint arXiv:2501.16496*, 2025.

744 [106] Herbert A Simon. The architecture of complexity. In *The Roots of Logistics*, pages 335–361.  
745 Springer, 2012.

746 [107] Charlotte Siska, Katerina Marazopoulou, Melissa Ailem, and James Bono. Examining the  
747 robustness of llm evaluation to the distributional assumptions of benchmarks. In *Proceedings*  
748 *of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
749 *Papers*), pages 10406–10421, 2024.

750 [108] Irene Solaiman, Zeerak Talat, William Agnew, Lama Ahmad, Dylan Baker, Su Lin Blodgett,  
751 Canyu Chen, Hal Daumé III, Jesse Dodge, Isabella Duan, et al. Evaluating the social impact  
752 of generative ai systems in systems and society. *arXiv preprint arXiv:2306.05949*, 2023.

753 [109] Elizabeth S Spelke and Katherine D Kinzler. Core knowledge. *Developmental science*, pages  
 754 89–96, 2007.

755 [110] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid,  
 756 Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al.  
 757 Beyond the imitation game: Quantifying and extrapolating the capabilities of language models.  
 758 *arXiv preprint arXiv:2206.04615*, 2022.

759 [111] Alan M Turing. *Computing machinery and intelligence*. Springer, 2009.

760 [112] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
 761 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information  
 762 processing systems*, 30, 2017.

763 [113] Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language  
 764 model. *arXiv preprint arXiv:1906.04284*, 2019.

765 [114] Ludwig Von Bertalanffy. General system theory. *New York*, 41973(1968):40, 1968.

766 [115] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix  
 767 Hill, Omer Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose  
 768 language understanding systems. *Advances in neural information processing systems*, 32,  
 769 2019.

770 [116] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.  
 771 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv  
 772 preprint arXiv:1804.07461*, 2018.

773 [117] An Wang, Xingwu Sun, Ruobing Xie, Shuaipeng Li, Jiaqi Zhu, Zhen Yang, Pinxue Zhao,  
 774 JN Han, Zhanhui Kang, Di Wang, et al. Hmoe: Heterogeneous mixture of experts for language  
 775 modeling. *arXiv preprint arXiv:2408.10681*, 2024.

776 [118] Thomas Wang, Adam Roberts, Daniel Hesslow, Teven Le Scao, Hyung Won Chung, Iz Beltagy,  
 777 Julien Launay, and Colin Raffel. What language model architecture and pretraining objective  
 778 works best for zero-shot generalization? In *International Conference on Machine Learning*,  
 779 pages 22964–22984. PMLR, 2022.

780 [119] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha  
 781 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in  
 782 language models. *arXiv preprint arXiv:2203.11171*, 2022.

783 [120] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo,  
 784 Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. Mmlu-pro: A more robust and  
 785 challenging multi-task language understanding benchmark. In *The Thirty-eight Conference on  
 786 Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.

787 [121] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani  
 788 Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large  
 789 language models. *arXiv preprint arXiv:2206.07682*, 2022.

790 [122] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,  
 791 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.  
 792 *Advances in neural information processing systems*, 35:24824–24837, 2022.

793 [123] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang,  
 794 Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of  
 795 harm from language models. *arXiv preprint arXiv:2112.04359*, 2021.

796 [124] Marty J Wolf, Frances S Grodzinsky, and Keith W Miller. *Artificial agents, cloud computing,  
 797 and quantum computing: Applying Floridi’s method of levels of abstraction*. Springer, 2012.

798 [125] Ryszard W Wolny. Hyperreality and simulacrum: Jean baudrillard and european postmod-  
 799 ernism. *European Journal of Interdisciplinary Studies*, 3(3):75–79, 2017.

800 [126] Xingyu Xie, Kuangyu Ding, Shuicheng Yan, Kim-Chuan Toh, and Tianwen Wei. Optimization  
801 hyper-parameter laws for large language models. *arXiv preprint arXiv:2409.04777*, 2024.

802 [127] Cheng Xu, Shuhao Guan, Derek Greene, M Kechadi, et al. Benchmark data contamination of  
803 large language models: A survey. *arXiv preprint arXiv:2406.04244*, 2024.

804 [128] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,  
805 Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,  
806 Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao,  
807 Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao  
808 Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan,  
809 Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv  
810 preprint arXiv:2412.15115*, 2024.

811 [129] Longxuan Yu, Delin Chen, Siheng Xiong, Qingyang Wu, Dawei Li, Zhikai Chen, Xiaoze Liu,  
812 and Liangming Pan. Causaleval: Towards better causal reasoning in language models. In  
813 *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association  
814 for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*,  
815 pages 12512–12540, 2025.

816 [130] Guanhua Zhang and Moritz Hardt. Inherent trade-offs between diversity and stability in  
817 multi-task benchmarks. *arXiv preprint arXiv:2405.01719*, 2024.

818 [131] Zhehao Zhang, Jiaao Chen, and Difyi Yang. Darg: Dynamic evaluation of large language  
819 models via adaptive reasoning graph. *arXiv preprint arXiv:2406.17271*, 2024.

820 [132] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian  
821 Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models.  
822 *arXiv preprint arXiv:2303.18223*, 1(2), 2023.

823 [133] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao  
824 Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with  
825 mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–  
826 46623, 2023.

827 **A A Feasible Proposal for Explicitly Hierarchical Evaluation Frameworks**

828 Building on our summary of LLM assessment literature in Section 4.2, which highlighted key  
829 concerns and practical methods relevant to our suggested layer-wise framework, this section aims  
830 to synthesize these elements into a cohesive, operationalizable schema. Hierarchical evaluation  
831 frameworks are particularly promising, as exemplified by the ReFeR system [83]. ReFeR employs  
832 a “peer review” model where smaller AI “evaluators” provide initial assessments, and a more  
833 capable “area chair” model synthesizes these into final scores and reasoning. This hierarchical,  
834 multi-perspective meta-methodology strongly aligns with our goals. Consequently, we adopt a similar  
835 hierarchical system to embody our proposed *model cartography* as shown below.

836 •  $(\mathcal{M}_O)$  The assessment group that focuses on benchmark performance and output quality.

837 – One important approach involves employing adaptive testing and dynamic benchmarks,  
838 as exemplified by Zhang et al. (2024) [131]. Their method first extracts reasoning  
839 graphs from existing benchmark data points and then perturbs these graphs to generate  
840 novel test cases. Subsequently, a code-augmented LLM verifies the correctness of  
841 labels for the newly generated data. Using this framework, they observed that LLM  
842 performance declines with increasing task complexity, often revealing greater biases  
843 and excessive sensitivity to specific content. This research highlights the value of  
844 evaluating LLMs beyond conventional static benchmarks, offering a more dynamic  
845 perspective on assessing their capabilities and limitations.

846 – Another promising practice is measuring an LLM’s “usage of capacity,” as proposed  
847 by Cao et al. (2025) [19]. The central concept is that a comprehensive assessment of  
848 an LLM’s ability should consider the effort it expends to achieve an outcome. To this  
849 end, they introduce a measure called MUI to quantify how extensively a model utilizes  
850 its capabilities to complete tasks. This approach yields model rankings consistent  
851 with expert judgment and demonstrates robustness to variance. Their work offers  
852 a significant step towards addressing challenges in assessing model capacity and  
853 potentially mitigating the impact of data contamination on evaluations. Furthermore,  
854 combining this method with adaptive testing and dynamic benchmarks could allow  
855 for a more comprehensive “observation dynamics” of LLM performance.

856 •  $(\mathcal{M}_R)$  The assessment group that employs interpretability techniques and evaluates robustness  
857 to perturbations.

858 – Advanced mechanistic interpretability is vital for understanding and confirming a  
859 model’s reachable capabilities. For instance, circuit discovery is a key method for  
860 linking these observable behaviors to the model’s internally instantiated concepts.  
861 Seminal work by Elhage et al. (2021) [35] has shown how specific learned modules,  
862 like attention heads, are integral to behaviors such as in-context learning. Although  
863 applying such detailed techniques to larger, more complex models is challenging,  
864 emerging automated methods like Automatic Circuit DisCovery (ACDC) [27] offer  
865 a promising path. These tools aim to pinpoint underlying mechanisms and could  
866 provide more systematic ways, potentially even quantifiable measures, to assess how  
867 specific, predefined behaviors (*i.e.*, a standard dataset) are realized within a model’s  
868 architecture.

869 – Complementing interpretability techniques, exploring model responses to strategically  
870 designed inputs provides valuable insights. Adversarial prompting is a prominent ex-  
871 ample, involving inputs crafted to exploit model vulnerabilities or specific processing  
872 mechanisms. The generation of such prompts has become increasingly accessible; for  
873 instance, tools like AdvPrompter [89] can rapidly produce human-readable adversarial  
874 prompts, sometimes even while preserving the original prompt’s semantic meaning.  
875 A model’s susceptibility and characteristic responses to a suite of such prompts can  
876 serve as a crucial basis for metrics to evaluate its propensity for undesirable reachable  
877 behaviors, including generating misinformation, revealing sensitive information, or  
878 engaging in incorrect reasoning.

879 •  $(\mathcal{M}_I)$  The assessment group that investigates the faithfulness of explanations or analyzes  
880 internal activation patterns.

881 – For different instantiations of a concept, it’s crucial to assess whether their chosen  
 882 hyperparameters and architectural constraints align with the underlying design ratio-  
 883 nade and theoretical foundations. This is particularly true for models with pre-trained  
 884 initializations, where the pre-training data is fundamental to how well the model  
 885 instantiates the concept by leveraging its learned general features and representations.  
 886 These theoretical foundations are often informed by scaling laws [50, 53, 58, 126].  
 887 Scaling laws describe how a model’s performance (often measured by loss) predictably  
 888 relates to key factors such as model size, dataset size, and training compute, with this  
 889 relationship typically characterizable by mathematical functions. Consequently, by  
 890 comparing the actual performance of a series of model instantiations against predic-  
 891 tions derived from these scaling laws, we can confirm the quality of their instantiation.  
 892 This process also allows us to forecast the potential performance of significantly larger  
 893 models, thereby guiding strategic research directions.

894 – Another effective approach to assess how well concepts are instantiated within a model  
 895 is to employ visualizations of its internal modules. For instance, attention head visual-  
 896 izations [113] illustrate the patterns of attention, detailing how much consideration  
 897 each token gives to other tokens in the input (or across encoder-decoder interactions)  
 898 within specific attention heads and layers. This aids in understanding information  
 899 flow and identifying which tokens influence the representations of others. A further  
 900 example is the visualization of expert routing in Mixture-of-Experts models [117],  
 901 which demonstrates how different experts are activated in response to various input  
 902 tokens. Such examples underscore that visualization is crucial for gaining an intu-  
 903 itive understanding of how a model instantiates concepts and for debugging potential  
 904 conceptual issues.

905 •  $(\mathcal{M}_C)$  The assessment group that oversees the multi-tier assessments, comparing them against  
 906 the documented goals and assumptions of the concepts.

907 – Although directly measuring the inherent “quality” of a concept is challenging, trans-  
 908 lating qualitative concepts and theoretical underpinnings, such as *understanding*,  
 909 *reasoning*, *creativity*, or *safety*, into measurable outcomes is essential for genuinely  
 910 assessing conceptual results. To achieve this, there is a pressing need to develop and  
 911 implement precise metrics. Such metrics are indispensable for objectively evaluating  
 912 how effectively observed *phenomena* reflect these carefully defined concepts and  
 913 whether the system under scrutiny adheres to its acknowledged operational bound-  
 914 ries. This, in turn, enables a more robust and empirically grounded validation of that  
 915 system’s conceptual achievements.

## 916 B Border Impacts and Limitations

917 **Border Impacts:** Adopting the proposed five-tier hierarchical view of AI models carries significant  
 918 implications for AI development, research methodology, and AI epistemology. For development, it  
 919 encourages designing systems with evaluability in mind at each level, from clear concept articulation  
 920 and transparent implementation to better management of the optimized model’s scope. In research,  
 921 this framework calls for a shift from “SOTA-chasing” on narrow benchmarks to new experimental  
 922 designs and metrics that provide deeper insights into different model tiers, their interrelations,  
 923 behavior, generalization capabilities, and alignment. Epistemologically, the framework redefines  
 924 what it means to “understand” an AI model, challenging the adequacy of single-score evaluations for  
 925 complex systems and aligning with broader philosophical discussions about observing and inferring  
 926 truths about partially unobservable entities, such as those addressed by the Noumenal and Reachable  
 927 tiers, thereby contributing to the philosophy of AI.

928 **Limitations:** Operationalizing the proposed multi-tier evaluation framework presents considerable  
 929 challenges. A primary difficulty is defining clear, measurable, and meaningful criteria for the more  
 930 abstract Noumenal and Conceptual levels; for instance, assessing the “plausibility” of a Conceptual  
 931 Model becomes empirically perplexing when tied to an inherently unknowable Noumenon. Another  
 932 significant hurdle involves characterizing the full, combinatorially explosive spectrum of an LLM’s  
 933 potential behaviors, especially for models with billions of parameters and global representations.  
 934 Furthermore, existing substantial challenges in LLM evaluation, such as frequent benchmark updates,  
 935 the high cost of comprehensive assessment, and mitigating data contamination, are likely to be

936 amplified within such a demanding multi-tier regime. Finally, there's the inherent risk that any new,  
937 complex evaluation framework could itself become a target for "SOTA-chasing," diverting efforts to  
938 metric optimization without necessarily achieving genuine progress in underlying model quality or  
939 understanding.