# A Conceptual Framework for AI Capability Evaluations

**María Victoria Carro** [1 2]   **Denise Alejandra Mester** [2]   **Francisca Gauna Selasco** [2]   **Luca Nicolás Forziati Gangi** [2]
**Matheo Sandleris Musa** [3]   **Lola Ramos Pereyra** [2]   **Mario Leiva** [4 5]   **Juan Gustavo Corvalan** [3]
**Maria Vanina Martinez** [6]   **Gerardo Simari** [4 5 7]

## Abstract

As AI systems advance and integrate into society, well-designed and transparent evaluations are becoming essential tools in AI governance, informing decisions by providing evidence about system capabilities and risks. Yet there remains a lack of clarity on how to perform these assessments both comprehensively and reliably. To address this gap, we propose a conceptual framework for analyzing AI capability evaluations, offering a structured, descriptive approach that systematizes the analysis of widely used methods and terminology without imposing new taxonomies or rigid formats. This framework supports transparency, comparability, and interpretability across diverse evaluations. It also enables researchers to identify methodological weaknesses, assists practitioners in designing evaluations, and provides policymakers with an accessible tool to scrutinize, compare, and navigate complex evaluation landscapes.

## 1. Introduction

Evaluations are gaining significant attention in the field of AI, with substantial efforts dedicated to advancing this area. The rapid evolution of large language models (LLMs) and their growing integration into daily life have underscored the need for robust and rigorous processes to understand the state of frontier AI current capabilities, identify potential risks to improve safety, and comprehend its societal impact.

In particular, capability evaluations[1] are crucial for various stakeholders, including academia, industry, government, and end-users. They serve as essential tools for tracking and communicating progress within the AI community, and for demonstrating the improvements of newly proposed methods over prior baselines (Biderman et al., 2024). Additionally, they play a key role in defining the threshold for achieving artificial general intelligence (AGI) (Pfister & Jud, 2025; Bubeck et al., 2023).

Evaluations are key components of governance regimes (Reuel et al., 2025; Mökander et al., 2023). In regulatory contexts, they help to identify areas where intervention may be needed and inform decisions by providing evidence about system capabilities and risks (Reuel et al., 2025; Shevlane et al., 2023; Paskov et al., 2025). Notably, the EU AI Act incorporates benchmarks in several key provisions (Eriksson et al., 2025), becoming the world's first mover in government-mandated general-purpose AI (GPAI) evaluations (Paskov et al., 2024).

For evaluations to fulfill these roles effectively, they must be conducted and reported rigorously (Paskov et al., 2025), with sufficient context and transparency to allow for meaningful interpretation by a wide range of stakeholders, including policymakers (Staufer et al., 2025). However, there is a lack of clarity on how to perform these assessments both comprehensively and reliably (Reuel et al., 2024b), and existing practices are among the key factors shaping practitioners' evaluation choices (Zhou et al., 2022), which may not align with the informational needs of decision-

---

[1]Università degli Studi di Genova, GE, Italy [2]FAIR, IALAB, University of Buenos Aires, BA, Argentina [3]University of Buenos Aires, BA, Argentina [4]Dept. of Computer Science and Engineering, Universidad Nacional del Sur (UNS) [5]Inst. of Computer Science and Engineering (ICIC UNS-CONICET), Bahía Blanca, BA, Argentina [6]Artificial Intelligence Research Institute (IIIA-CSIC), Universidad Autónoma de Barcelona, Barcelona, España [7]School of Computing and Augmented Intelligence, Arizona State University, USA. Correspondence to: María Victoria Carro <6381013@studenti.unige.it>.

---

[1]Capability evaluations have been defined as those that comprehensively assess a system's overall capabilities, including planned, unplanned, emerging, or dangerous capabilities (Xia et al., 2024). Burden et al. (2025) note that capability evaluations are used by AI developers and regulators to determine whether a system is safe to deploy, and by AI adopters to assess whether a system can automate specific tasks within their organizations. In this paper, we adopt a broad definition of capability evaluations, aligned with the usage found in Reuel et al. (2024b). Accordingly, we exclude impact evaluations—also referred to as "real-world impact evaluations"(Burden et al., 2025)—which measure the effects of AI systems once deployed in real-world settings. These evaluations typically involve human subjects and treat the AI system's assistance as an "intervention," whose effect must be empirically quantified.

makers. Without well-designed evaluations, governments risk relying on incomplete, misleading, or selectively reported information, undermining efforts to ensure the safe and beneficial development and deployment of AI systems.

In this paper, we propose a conceptual framework for analysing capability evaluations, which includes key elements and sub-elements of evaluation processes, their interrelationships, and examples with relevant variables. Additionally, some of the key challenges associated with evaluations are identified and integrated into the framework in the extended version of Appendix 1. By creating a highly simplified representation of complex processes, essential features can be abstracted to permit systematic reasoning (Gupta et al., 2024). This, in turn, provides a clear overview and a structured approach that supports the standardization of concepts, methods and reports—recognized as a critical need in the field (Weidinger et al., 2023; Thurnherr, 2024)—, as well as the comparison and transparency of evaluation processes. Such transparency is especially important when evaluations are developed by independent third parties, as policymakers and other stakeholders may need to scrutinize the results, something that depends, in part, on how well the evaluation's development and execution are documented (Thurnherr, 2024).

We developed this framework through a careful analysis of the design of numerous evaluations, some of which are cited as examples throughout the paper. While not exhaustive, it is intended to be comprehensive—capturing the full range of aspects that an evaluation method may encompass, even though not all elements will be present in every individual evaluation (e.g., prompt techniques).

Our conceptual framework offers two main advantages. First, it adopts existing terminology and systematizes current methods, integrating the most widely used concepts in the field rather than introducing new terms. In doing so, it adapts to the established language of diverse stakeholders, rather than requiring them to conform to a novel taxonomy. The framework is thus descriptive rather than normative. However, in instances where recommendations or best practices are mentioned, these are drawn from the existing literature and included based on their broad acceptance within the field, rather than originating from the authors themselves. Second, it prioritizes simplicity and accessibility, avoiding unnecessary complexity or rigid formats. It aims to foster consensus around a broad, intuitive framework that remains easily understandable to all participants in the field of evaluations, including those without deep technical expertise.

This approach has the potential to support a wide range of stakeholders in navigating and interpreting AI capability evaluations:

- For **public bodies**, it offers a structured tool to scrutinize, audit and systematically compare evaluations, enabling more consistent assessments.

- For **researchers** and **third-party evaluators**, it provides a conceptual map to understand which dimensions matter for use cases, synthesize and critically analyze evaluation practices, identify methodological strengths, weaknesses, and gaps.

- For **industry practitioners** and **model developers**—who must make numerous decisions throughout the evaluation process (Gupta et al., 2024)—this highlights key elements to consider when designing or reporting evaluations, promoting greater transparency and comparability.

- For **newcomers to the field** and **policymakers**, it serves as an accessible entry point to understand evaluation structures and interpret results.

## 2. Related Work

Dow et al. (2024) present Dimensions of Generative AI Evaluation Design for mapping the state of evaluations in which many proposed dimensions—such as input source, task type, and metric type—overlap with elements of our own. However, their framework is limited to generative AI systems and safety evaluations, whereas ours may be applicable more broadly. While the underlying motivation is similar—highlighting the importance of structured evaluation—our framework goes further by decomposing the elements into subcomponents, systematically integrating evaluation challenges within each element, and providing a more granular and detailed perspective.

Additional related work includes Laskar et al. (2024), who survey the challenges of evaluation and structure their analysis by segmenting the evaluation process into distinct stages. Furthermore, Paskov et al. (2025) outline suggestions for enhancing the rigor of general-purpose AI (GPAI) evaluations, offering practical recommendations for each stage of the evaluation life cycle. While the latter primarily focuses on benchmarks and uplift studies, several of their insights—such as documentation practices and statistical safeguards—have broader applicability. However, both works rely on a sequential model of the evaluation process, prescribing a particular order of steps. In contrast, our framework adopts a descriptive approach: rather than prescribing what should happen first or next, it maps the key elements of evaluation without imposing a fixed sequence.

Finally, while Weidinger et al. (2023) propose and construct a repository of existing evaluations, their work is limited to risk evaluations, which are categorized according to a risk taxonomy introduced by the authors. Moreover, the

repository is a static resource, though they suggest the development of a living version as future work. In contrast, we propose a broader repository that includes, but is not limited to, risk evaluations, and organizes them according to multiple criteria. This repository is envisioned as an interactive resource—similar in spirit to the MIT Risk Repository (Slattery et al., 2024), but focused specifically on evaluations.

## 3. The Conceptual Framework

### 3.1. Evaluation Target

**Capability.** Capability evaluations focus on what AI systems can and cannot do. The specific capability under evaluation might range from mathematical reasoning (Didolkar et al., 2024) and causal inference (Kiciman et al., 2023) to legal reasoning (Guha et al., 2023), deception (Hagendorff, 2024), and many others.

**Objective.** The objective of the evaluation may be to quantify system capabilities, track progress, enable large-scale comparisons, understand behavioral patterns, identify and estimate potential risks (Weidinger et al., 2023), provide assurance of system safety (Burden et al., 2025; Weidinger et al., 2023), predict if future models can lead to catastrophic harm (Barnett & Thiergart, 2024) or evaluate evaluation methods, among others.

### 3.2. Task

The **task** refers to a particular specification of a problem (Raji et al., 2021), typically represented as a mapping between an input space and an output or action space (Schlangen, 2021; Eriksson et al., 2025).

**Task Mode.** The task mode is the format of the response required for the system, which can be closed-ended, when the subject is asked to identify an answer from a predetermined set of options, or open-ended, when the subject is required to generate a novel output (Burden et al., 2025). Some examples include question answering, text completion, text summarization, or error identification and correction.

**Steps.** Some evaluations require the system to perform multiple coordinated steps to complete a task, to simulate real-world activities (Kraprayoon et al., 2025). These intermediate steps can be autonomously defined by the system or dictated by a human subject. This is particularly relevant for AI agent evaluations, which may involve end-to-end processes that include planning, tool use, and iterative refinement to produce a coherent final output (Testini et al., 2025).

**Interactions.** Number of interactions between the evaluated subject and the user or evaluator, which can be classified as single-turn or multi-turn (Wang et al., 2024; Cheng et al., 2024; Sirdeshmukh et al., 2025).
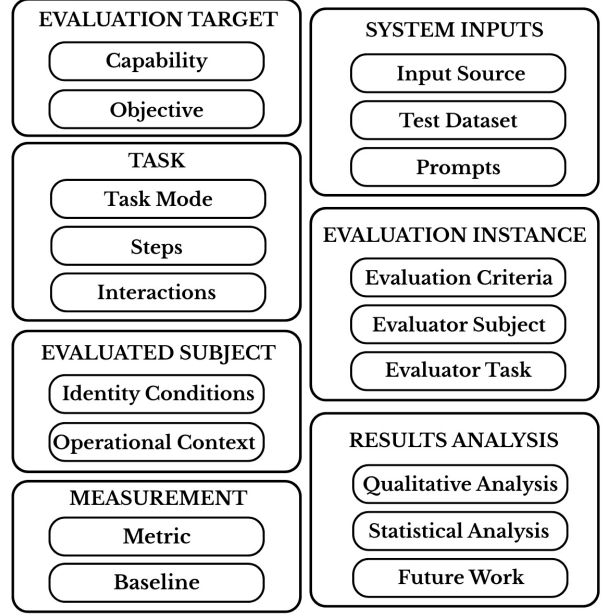


*Figure 1.* Overview of the proposed Conceptual Framework.

### 3.3. Evaluated Subject

The **evaluated subject** is the subject of the capability claim—in most cases, an AI system, and in the case of meta-evaluations, an evaluation method itself (e.g. Hong et al. (2024)).

**Identity Conditions.** The evaluated system should be individualized by defining its identity conditions (Harding & Sharadin, 2024; Biderman et al., 2024). This includes, where possible, fixed attributes such as details of the architecture and parameters (Paskov et al., 2025), the system's name and version, and whether it is a foundational model or a fine-tuned variant.

**Operational Context.** The operational context is the background conditions of the evaluation (Harding & Sharadin, 2024), including access methods (e.g., raw API or assistant interface), configuration settings (e.g. temperature), auxiliary tools (e.g. plugins), inference budget (e.g. token or query limits), and hardware or compute conditions where applicable.

### 3.4. System Inputs

The **input source** refers to the origin of the inputs used in the evaluation (Dow et al., 2024). These may come from a pre-existing test dataset, or user-generated inputs, including those produced dynamically interacting with the system. For example, in red teaming, attacks on the target model are often conducted through trial and error to identify effective

strategies that elicit abnormal behavior (Lin et al., 2025). The input can also take various modalities, such as image, audio, video, text, or multimodal combinations.

**Test Dataset.** A separate dataset, distinct from the training data, consisting of data instances can be used to evaluate the AI system on a given task. However, in contexts where generalization is not prioritized, the test data may overlap with the training data. In the case of benchmarks and reference-based evaluations, the dataset also includes the desired output and annotations for each instance (Eriksson et al., 2025; Liu et al., 2024). Key dimensions to consider include the dataset's size and construction method—for example, whether it was manually developed by researchers or domain experts, or synthetically generated, for instance, through an LLM.

**Prompt Techniques.** A prompting strategy serves to turn each data instance into an input which can be processed by an AI system. In this context, most prompting strategies can be viewed as 'templates', which generate the model input by embedding the query in a pre-written prompt format (Harding & Sharadin, 2024). Furthermore, prompts may include model-specific instructions designed to steer behavior and optimize performance (Schulhoff et al., 2024). Some of the prompt techniques include for example, few shots, one shot, chain-of-thought (CoT), self-consistency or emotion prompting (Sahoo et al., 2024; Schulhoff et al., 2024).

### 3.5. Evaluation Instance

**Evaluation Criteria.** This refers to specific aspects of the output that define what constitutes a correct, appropriate, or high-quality response. Some examples are correctness, fluency, generalization, reasoning, robustness (Hu & Zhou, 2024), or logical coherence.

**Evaluator Subject.** An evaluator subject—human or AI system—is needed when the evaluation criteria are not mathematically defined and no prior annotation of the dataset exists. It is key to consider how the evaluator is selected—e.g., human evaluators may be recruited via crowdsourcing, while, ideally, LLM evaluators should not be the same model under evaluation or from the same model family, to avoid self-preference bias (Panickssery et al., 2024; Bai et al., 2023).

**Evaluator Task.** The nature and instructions of the evaluator's task must be specified, such as pairwise comparison, as in Chatbot Arena (Chiang et al., 2024), rating individual outputs on a Likert scale or assigning a score. The order in which outputs are presented is also important to avoid position bias (Van Der Lee et al., 2019; Gao et al., 2025).

Note that in this framework, we distinguish between the concepts of annotators and evaluators. An *annotator*—whether a human or an AI system—is responsible for labeling a test dataset prior to the evaluation of the evaluated subject. This may involve selecting the correct answer in a multiple-choice item, assigning categories to text fragments, or specifying the desired answer to a question. In contrast, an *evaluator* assesses and annotates the outputs of the evaluated subject, for instance by rating the creativity of a poem generated by an LLM using a Likert scale. While the term 'annotator' is often used in the literature to refer to both roles (e.g. Weidinger et al. (2023)), we maintain this distinction for clarity.

### 3.6. Measurement

**Metric.** A metric is a mathematically defined function used to quantify a specific evaluation criterion. Metrics provide a standardized and objective way to measure aspects of system outputs, as they can summarize model performance on a set of tasks and datasets as a single number or score (Raji et al., 2021; Eriksson et al., 2025). They can be classified as reference-based or reference-free (Ito et al., 2025; Li et al., 2024), and as general-purpose (e.g. Accuracy, Recall, Precision) or task-specific metrics (e.g. Translation Edit Rate (TER) (Snover et al., 2006)).

Evaluation practitioners may confuse evaluation criteria with the metrics used to operationalize them, but this confusion risks losing the ability to assess the metric's appropriateness and distinguish between the metric and the original goal (Zhou et al., 2022). As Goodhart's Law states, 'When a measure becomes a target, it ceases to be a good measure'.

**Baselines.** Baselines are typically used as reference points or standards for comparison when the goal of an evaluation is to compare the performance of different systems, approaches, or evaluation methods.

### 3.7. Result Analysis

**Qualitative Analysis.** Ranking systems according to a single quality number is easy and actionable, but often, it is much more important to understand when and why models fail (Eriksson et al., 2025). Therefore, it is also beneficial to conduct a deeper analysis of the evaluation results (e.g., analyzing results in relation to dataset complexity levels, if a dataset composition analysis was conducted).

**Statistical Analysis.** When applicable, statistical analysis should be performed to distinguish genuine performance improvements from random variability. Best practices emphasize treating evaluation processes as structured experiments, applying significance testing and quantifying uncertainty through confidence intervals or resampling techniques (Miller, 2024; Sun et al., 2024; Zhan et al., 2024).

**Future Work.** Results can reveal valuable insights and guide future research directions (Hämäläinen & Alnajjar, 2021) as well as inform system development and highlight

specific avenues for improvement.

## 4. Future Work and Conclusion

In this paper, we present a conceptual framework designed to map and structure capability evaluations of AI systems. A potential direction for future work is to analyze a range of evaluations using this framework as a set of use cases and to iteratively refine it. To demonstrate its practical application, we conducted such an analysis on three papers, presented in Appendix 2. Given the exponential growth in the number of evaluations being conducted and published, it is becoming increasingly difficult for stakeholders to stay up to date. This framework could serve as the basis for an interactive and dynamic repository of evaluations, enabling users to track trends in the field across different elements and challenges. For example, a researcher could explore all proposed solutions to data contamination or consult a catalog of metrics to identify those best suited to a particular evaluation. While this approach could contribute to a degree of standardization that has been recognized as necessary in the field (Weidinger et al., 2023), it would also preserve flexibility by making relevant and up-to-date information more accessible. Ultimately, this would facilitate cumulative knowledge-building and allow the community to make more efficient progress on the field's open challenges.

## References

AISI. Early insights from developing question-answer evaluations for frontier AI, 2024. Accessed May 10, 2025.

Bai, Y., Ying, J., Cao, Y., Lv, X., He, Y., Wang, X., Yu, J., Zeng, K., Xiao, Y., Lyu, H., et al. Benchmarking foundation models with language-model-as-an-examiner. *Advances in Neural Information Processing Systems*, 36: 78142–78167, 2023.

Barnett, P. and Thiergart, L. Declare and justify: Explicit assumptions in ai evaluations are necessary for effective regulation. *arXiv preprint arXiv:2411.12820*, 2024.

Benedetto, L. A quantitative study of nlp approaches to question difficulty estimation. In *International Conference on Artificial Intelligence in Education*, pp. 428–434. Springer, 2023.

Bhambhoria, R., Dahan, S., Li, J., and Zhu, X. Evaluating ai for law: Bridging the gap with open-source solutions. *arXiv preprint arXiv:2404.12349*, 2024.

Biderman, S., Schoelkopf, H., Sutawika, L., Gao, L., Tow, J., Abbasi, B., Aji, A. F., Ammanamanchi, P. S., Black, S., Clive, J., et al. Lessons from the trenches on reproducible evaluation of language models. *CoRR*, 2024.

Bieger, J., Thórisson, K. R., Steunebrink, B. R., and Thorarensen, T. Evaluation of general-purpose artificial intelligence: Why, what & how. In *Proceedings of the IJCAI Workshop on Evaluating General-Purpose Artificial Intelligence (EGPAI 2016)*, 2016. URL https://alumni.media.mit.edu/~kris/ftp/EGPAI_2016_paper_9.pdf. Accessed May 10, 2025.

Bubeck, S., Chadrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S., et al. Sparks of artificial general intelligence: Early experiments with gpt-4, 2023.

Burden, J. Evaluating ai evaluation: Perils and prospects. *arXiv preprint arXiv:2407.09221*, 2024.

Burden, J., Tešić, M., Pacchiardi, L., and Hernández-Orallo, J. Paradigms of ai evaluation: Mapping goals, methodologies and culture. *arXiv preprint arXiv:2502.15620*, 2025.

Cao, J., Chan, Y.-K., Ling, Z., Wang, W., Li, S., Liu, M., Qiao, R., Han, Y., Wang, C., Yu, B., He, P., Wang, S., Zheng, Z., Lyu, M. R., and Cheung, S.-C. How should we build a benchmark? revisiting 274 code-related benchmarks for llms, 2025. URL https://arxiv.org/abs/2501.10711.

Casper, S., Ezell, C., Siegmann, C., Kolt, N., Curtis, T. L., Bucknall, B., Haupt, A., Wei, K., Scheurer, J., Hobbhahn, M., et al. Black-box access is insufficient for rigorous ai audits. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pp. 2254–2272, 2024.

Cheng, Y., Georgopoulos, M., Cevher, V., and Chrysos, G. G. Leveraging the context through multi-round interactions for jailbreaking attacks. *arXiv preprint arXiv:2402.09177*, 2024.

Chiang, W.-L., Zheng, L., Sheng, Y., Angelopoulos, A. N., Li, T., Li, D., Zhu, B., Zhang, H., Jordan, M., Gonzalez, J. E., et al. Chatbot arena: An open platform for evaluating llms by human preference. In *Forty-first International Conference on Machine Learning*, 2024.

Deutsch, D., Dror, R., and Roth, D. On the limitations of reference-free evaluations of generated text. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 10960–10977, 2022.

Didolkar, A., Goyal, A., Ke, N. R., Guo, S., Valko, M., Lillicrap, T., Jimenez Rezende, D., Bengio, Y., Mozer, M. C., and Arora, S. Metacognitive capabilities of llms: An exploration in mathematical problem solving. *Advances in Neural Information Processing Systems*, 37: 19783–19812, 2024.

Dong, Y., Jiang, X., Liu, H., Jin, Z., Gu, B., Yang, M., and Li, G. Generalization or memorization: Data contamination and trustworthy evaluation for large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pp. 12039–12050, 2024.

Dow, A., Vaughan, J. W., Barocas, S., Atalla, C., Chouldechova, A., and Wallach, H. Dimensions of generative ai evaluation design. In *Proceedings of the NeurIPS 2024 Workshop on Evaluating Evaluations: Examining Best Practices for Measuring Broader Impacts of Generative AI*, 2024. URL https://neurips.cc/virtual/2024/104211. Accessed May 11, 2025.

Eriksson, M., Purificato, E., Noroozian, A., Vinagre, J., Chaslot, G., Gomez, E., and Fernandez-Llorca, D. Can we trust ai benchmarks? an interdisciplinary review of current issues in ai evaluation. *arXiv preprint arXiv:2502.06559*, 2025.

European Commission. Second draft of the general purpose AI code of practice, April 2024. Written by independent experts. Accessed May 10, 2025.

Frontier Model Forum. Issue brief: Early best practices for frontier AI safety evaluations, 2024. Accessed May 10, 2025.

Gao, M., Hu, X., Yin, X., Ruan, J., Pu, X., and Wan, X. Llm-based nlg evaluation: Current status and challenges. *Computational Linguistics*, pp. 1–28, 2025.

Gehrmann, S., Clark, E., and Sellam, T. A case for better evaluation standards in nlg. In *Workshop on Setting up ML Evaluation Standards to Accelerate Progress at ICLR 2022*, 2022. URL https://iclr.cc/virtual/2022/7328. Poster presentation.

Gehrmann, S., Clark, E., and Sellam, T. Repairing the cracked foundation: A survey of obstacles in evaluation practices for generated text. *Journal of Artificial Intelligence Research*, 77:103–166, 2023.

Guha, N., Nyarko, J., Ho, D., Ré, C., Chilton, A., Chohlas-Wood, A., Peters, A., Waldon, B., Rockmore, D., Zambrano, D., et al. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. *Advances in Neural Information Processing Systems*, 36:44123–44279, 2023.

Gupta, N. R., Hullman, J., and Subramonyam, H. A conceptual framework for ethical evaluation of machine learning systems. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, volume 7, pp. 534–546, 2024.

Hagendorff, T. Deception abilities emerged in large language models. *Proceedings of the National Academy of Sciences*, 121(24):e2317967121, 2024.

Hagendorff, T., Dasgupta, I., Binz, M., Chan, S. C. Y., Lampinen, A., Wang, J. X., Akata, Z., and Schulz, E. Machine psychology, 2024. URL https://arxiv.org/abs/2303.13988.

Hämäläinen, M. and Alnajjar, K. Human evaluation of creative NLG systems: An interdisciplinary survey on recent papers. In Bosselut, A., Durmus, E., Gangal, V. P., Gehrmann, S., Jernite, Y., Perez-Beltrachini, L., Shaikh, S., and Xu, W. (eds.), *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pp. 84–95, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.gem-1.9. URL https://aclanthology.org/2021.gem-1.9/.

Harding, J. and Sharadin, N. What is it for a machine learning model to have a capability? *The British Journal for the Philosophy of Science*, 2024. Advance online publication. Available at https://doi.org/10.1086/732153.

Hofstätter, F., Teoh, J., van der Weij, T., and Ward, F. R. The elicitation game: Stress-testing capability elicitation techniques. In *Workshop on Socially Responsible Language Modelling Research*, 2024. URL https://openreview.net/forum?id=zy6LB5t62f.

Hong, Z.-W., Shenfeld, I., Wang, T.-H., Chuang, Y.-S., Pareja, A., Glass, J. R., Srivastava, A., and Agrawal, P. Curiosity-driven red-teaming for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=4KqkizXgXU.

Hu, T. and Zhou, X.-H. Unveiling llm evaluation focused on metrics: Challenges and solutions. *arXiv preprint arXiv:2404.09135*, 2024.

Huang, H., Qu, Y., Zhou, H., Liu, J., Yang, M., Xu, B., and Zhao, T. On the limitations of fine-tuned judge models for llm evaluation. *arXiv preprint arXiv:2403.02839*, 2024.

Imani, S., Du, L., and Shrivastava, H. MathPrompter: Mathematical reasoning using large language models. In Sitaram, S., Beigman Klebanov, B., and Williams, J. D. (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pp. 37–42, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-industry.4. URL https://aclanthology.org/2023.acl-industry.4/.

Ito, T., van Deemter, K., and Suzuki, J. Reference-free evaluation metrics for text generation: A survey. *arXiv preprint arXiv:2501.12011*, 2025.

Ivanova, A. A. Running cognitive evaluations on large language models: The do's and the don'ts. *arXiv preprint arXiv:2312.01276*, 2023.

Jacovi, A., Caciularu, A., Goldman, O., and Goldberg, Y. Stop uploading test data in plain text: Practical strategies for mitigating data contamination by evaluation benchmarks. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 5075–5084, 2023.

Jin, Z., Chen, Y., Leeb, F., Gresele, L., Kamal, O., Lyu, Z., Blin, K., Gonzalez, F., Kleiman-Weiner, M., Sachan, M., and Schölkopf, B. Cladder: assessing causal reasoning in language models. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.

Jin, Z., Liu, J., LYU, Z., Poff, S., Sachan, M., Mihalcea, R., Diab, M. T., and Schölkopf, B. Can large language models infer causation from correlation? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=vqIH0ObdqL.

Khan, A., Hughes, J., Valentine, D., Ruis, L., Sachan, K., Radhakrishnan, A., Grefenstette, E., Bowman, S. R., Rocktäschel, T., and Perez, E. Debating with more persuasive llms leads to more truthful answers. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 23662–23733, 2024.

Kiciman, E., Ness, R., Sharma, A., and Tan, C. Causal reasoning and large language models: Opening a new frontier for causality. *Transactions on Machine Learning Research*, 2023.

Kraprayoon, J., Williams, Z., and Fayyaz, R. Ai agent governance: A field guide. *arXiv preprint arXiv:2505.21808*, 2025.

Kuhn, L., Gal, Y., and Farquhar, S. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=VD-AYtP0dve.

Lalor, J. P., Wu, H., and Yu, H. Building an evaluation scale using item response theory. In *Proceedings of the conference on empirical methods in natural language processing. Conference on empirical methods in natural language processing*, volume 2016, pp. 648, 2016.

Laskar, M. T. R., Alqahtani, S., Bari, M. S., Rahman, M., Khan, M. A. M., Khan, H., Jahan, I., Bhuiyan, A., Tan, C. W., Parvez, M. R., et al. A systematic survey and critical review on evaluating large language models: Challenges, limitations, and recommendations. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 13785–13816, 2024.

Lee, S., Kim, M., Cherif, L., Dobre, D., Lee, J., Hwang, S. J., Kawaguchi, K., Gidel, G., Bengio, Y., Malkin, N., and Jain, M. Learning diverse attacks on large language models for robust red-teaming and safety tuning. In *Red Teaming GenAI: What Can We Learn from Adversaries?*, 2025. URL https://openreview.net/forum?id=nnGITZZw9N.

Li, Z., Xu, X., Shen, T., Xu, C., Gu, J.-C., Lai, Y., Tao, C., and Ma, S. Leveraging large language models for nlg evaluation: Advances and challenges. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 16028–16045, 2024.

Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A., Newman, B., Yuan, B., Yan, B., Zhang, C., Cosgrove, C., Manning, C. D., Re, C., Acosta-Navas, D., Hudson, D. A., Zelikman, E., Durmus, E., Ladhak, F., Rong, F., Ren, H., Yao, H., WANG, J., Santhanam, K., Orr, L., Zheng, L., Yuksekgonul, M., Suzgun, M., Kim, N., Guha, N., Chatterji, N. S., Khattab, O., Henderson, P., Huang, Q., Chi, R. A., Xie, S. M., Santurkar, S., Ganguli, S., Hashimoto, T., Icard, T., Zhang, T., Chaudhary, V., Wang, W., Li, X., Mai, Y., Zhang, Y., and Koreeda, Y. Holistic evaluation of language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=iO4LZibEqW. Featured Certification, Expert Certification, Outstanding Certification.

Liao, Q. V. and Xiao, Z. Rethinking model evaluation as narrowing the socio-technical gap. *arXiv preprint arXiv:2306.03100*, 2023.

Liao, T., Taori, R., Raji, I. D., and Schmidt, L. Are we learning yet? a meta review of evaluation failures across machine learning. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. URL https://openreview.net/forum?id=mPducS1MsEK.

Lin, L., Mu, H., Zhai, Z., Wang, M., Wang, Y., Wang, R., Gao, J., Zhang, Y., Che, W., Baldwin, T., et al. Against the achilles' heel: A survey on red teaming for generative models. *Journal of Artificial Intelligence Research*, 82: 687–775, 2025.

Liu, Y., Cao, J., Liu, C., Ding, K., and Jin, L. Datasets for large language models: A comprehensive survey. *arXiv preprint arXiv:2402.18041*, 2024.

McIntosh, T. R., Susnjak, T., Arachchilage, N., Liu, T., Watters, P., and Halgamuge, M. N. Inadequacies of large language model benchmarks in the era of generative artificial intelligence. *arXiv preprint arXiv:2402.09880*, 2024.

Meyes, R., Lu, M., de Puiseau, C. W., and Meisen, T. Ablation studies in artificial neural networks. *arXiv preprint arXiv:1901.08644*, 2019.

Miller, E. Adding error bars to evals: A statistical approach to language model evaluations. *arXiv preprint arXiv:2411.00640*, 2024.

Mökander, J., Schuett, J., Kirk, H., and Floridi, L. Auditing large language models: a three-layered approach. *AI and Ethics*, 4(4), 2023.

Moniri, B., Hassani, H., and Dobriban, E. Evaluating the performance of large language models via debates. In Chiruzzo, L., Ritter, A., and Wang, L. (eds.), *Findings of the Association for Computational Linguistics: NAACL 2025*, pp. 2040–2075, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-195-7. URL https://aclanthology.org/2025.findings-naacl.109/.

Narayanan, A. and Kapoor, S. Gpt-4 and professional benchmarks: The wrong answer to the right question. *AI Snake Oil* (Substack), March 2023. URL https://www.aisnakeoil.com/p/gpt-4-and-professional-benchmarks. Accessed May 10, 2025.

OECD. Oecd framework for the classification of ai systems. *OECD Digital Economy Papers, No. 323*, 2022.

Orr, W. and Kang, E. B. Ai as a sport: On the competitive epistemologies of benchmarking. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1875–1884, 2024.

Panickssery, A., Bowman, S., and Feng, S. Llm evaluators recognize and favor their own generations. *Advances in Neural Information Processing Systems*, 37:68772–68802, 2024.

Paskov, P., Berglund, L., Smith, E. T., and Soder, L. Gpai evaluations standards taskforce: towards effective ai governance. In *Workshop on Socially Responsible Language Modelling Research*, 2024.

Paskov, P., Byun, M., Wei, K., and Webster, T. Preliminary suggestions for rigorous gpai model evaluations. Technical Report PEA3971-1, RAND Corporation, April 2025. URL https://www.rand.org/pubs/perspectives/PEA3971-1.html. Accessed May 10, 2025.

Perez, E., Ringer, S., Lukosiute, K., Nguyen, K., Chen, E., Heiner, S., Pettit, C., Olsson, C., Kundu, S., Kadavath, S., et al. Discovering language model behaviors with model-written evaluations. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 13387–13434, 2023.

Pfister, R. and Jud, H. Understanding and benchmarking artificial intelligence: Openai's o3 is not agi. *arXiv preprint arXiv:2501.07458*, 2025.

Piktus, A., Akiki, C., Villegas, P., Laurençon, H., Dupont, G., Luccioni, S., Jernite, Y., and Rogers, A. The roots search tool: Data transparency for llms. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pp. 304–314. Association for Computational Linguistics, 2023.

Raji, I. D., Denton, E., Bender, E. M., Hanna, A., and Paullada, A. AI and the everything in the whole wide world benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. URL https://openreview.net/forum?id=j6NxpQbREA1.

Rasheed, Z., Waseem, M., Systä, K., and Abrahamsson, P. Large language model evaluation via multi ai agents: Preliminary results. In *International Conference on Learning Representations*, pp. 1–12, 2024.

Rauh, M., Marchal, N., Manzini, A., Hendricks, L. A., Comanescu, R., Akbulut, C., Stepleton, T., Mateos-Garcia, J., Bergman, S., Kay, J., et al. Gaps in the safety evaluation of generative ai. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, volume 7, pp. 1200–1217, 2024.

Reuel, A., Hardy, A., Smith, C., Lamparth, M., Hardy, M., and Kochenderfer, M. Betterbench: Assessing AI benchmarks, uncovering issues, and establishing best practices. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024a. URL https://openreview.net/forum?id=hcOq2buakM.

Reuel, A., Soder, L., Bucknall, B., and Undheim, T. A. Position: technical research and talent is needed for effective ai governance. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024b.

Reuel, A., Bucknall, B., Casper, S., Fist, T., Soder, L., Aarne, O., Hammond, L., Ibrahim, L., Chan, A., Wills, P., Anderljung, M., Garfinkel, B., Heim, L., Trask, A., Mukobi, G., Schaeffer, R., Baker, M., Hooker, S., Solaiman, I., Luccioni, A. S., Rajkumar, N., Moës, N.,

Ladish, J., Bau, D., Bricman, P., Guha, N., Newman, J., Bengio, Y., South, T., Pentland, A., Koyejo, S., Kochenderfer, M. J., and Trager, R. Open problems in technical ai governance, 2025. URL https://arxiv.org/abs/2407.14981.

Ruan, J., Pu, X., Gao, M., Wan, X., and Zhu, Y. Better than random: reliable nlg human evaluation with constrained active sampling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 18915–18923, 2024.

Sahoo, P., Singh, A. K., Saha, S., Jain, V., Mondal, S., and Chadha, A. A systematic survey of prompt engineering in large language models: Techniques and applications. *CoRR*, abs/2402.07927, 2024. URL https://doi.org/10.48550/arXiv.2402.07927.

Sainz, O., Campos, J., García-Ferrero, I., Etxaniz, J., de Lacalle, O. L., and Agirre, E. Nlp evaluation in trouble: On the need to measure llm data contamination for each benchmark. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 10776–10787, 2023.

Schlangen, D. Targeting the benchmark: On methodology in current natural language processing research. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pp. 670–674, 2021.

Schulhoff, S., Ilie, M., Balepur, N., Kahadze, K., Liu, A., Si, C., Li, Y., Gupta, A., Han, H., Schulhoff, S., Dulepet, P. S., Vidyadhara, S., Ki, D., Agrawal, S., Pham, C., Kroiz, G. C., Li, F., Tao, H., Srivastava, A., Costa, H. D., Gupta, S., Rogers, M. L., Goncearenco, I., Sarli, G., Galynker, I., Peskoff, D., Carpuat, M., White, J., Anadkat, S., Hoyle, A. M., and Resnik, P. The prompt report: A systematic survey of prompting techniques. *CoRR*, abs/2406.06608, 2024. URL https://doi.org/10.48550/arXiv.2406.06608.

Sclar, M., Choi, Y., Tsvetkov, Y., and Suhr, A. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=RIu5lyNXjT.

Shevlane, T., Farquhar, S., Garfinkel, B., Phuong, M., Whittlestone, J., Leung, J., Kokotajlo, D., Marchal, N., Anderljung, M., Kolt, N., et al. Model evaluation for extreme risks. *arXiv preprint arXiv:2305.15324*, 2023.

Shi, J., Li, J., Ma, Q., Yang, Z., Ma, H., and Li, L. CHOPS: CHat with customer profile systems for customer service

with LLMs. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=9Wmdk94oKF.

Sirdeshmukh, V., Deshpande, K., Mols, J., Jin, L., Cardona, E.-Y., Lee, D., Kritz, J., Primack, W., Yue, S., and Xing, C. Multichallenge: A realistic multi-turn conversation evaluation benchmark challenging to frontier llms. *arXiv preprint arXiv:2501.17399*, 2025.

Slattery, P., Saeri, A. K., Grundy, E. A., Graham, J., Noetel, M., Uuk, R., Dao, J., Pour, S., Casper, S., and Thompson, N. The ai risk repository: A comprehensive meta-review, database, and taxonomy of risks from artificial intelligence. *CoRR*, 2024.

Snover, M., Dorr, B., Schwartz, R., Micciulla, L., and Makhoul, J. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pp. 223–231, 2006.

Staufer, L., Yang, M., Reuel, A., and Casper, S. Audit cards: Contextualizing ai evaluations. *arXiv preprint arXiv:2504.13839*, 2025.

Sun, K., Wang, R., Liu, H., and Søgaard, A. Comprehensive reassessment of large-scale evaluation outcomes in llms: A multifaceted statistical approach. *CoRR*, abs/2403.15250, 2024. URL https://doi.org/10.48550/arXiv.2403.15250.

Testini, I., Hernández-Orallo, J., and Pacchiardi, L. Measuring data science automation: A survey of evaluation tools for ai assistants and agents. *arXiv preprint arXiv:2506.08800*, 2025.

Thurnherr, B. C. Who should develop which ai evaluations?, April 2024. Accessed May 10, 2025.

Van Der Lee, C., Gatt, A., Van Miltenburg, E., Wubben, S., and Krahmer, E. Best practices for the human evaluation of automatically generated text. In *Proceedings of the 12th International Conference on Natural Language Generation*, pp. 355–368, 2019.

Vania, C., Htut, P. M., Huang, W., Mungra, D., Yuanzhe Pang, R., Phang, J., Liu, H., Cho, K., and Bowman, S. R. Comparing test sets with item response theory. In *Annual Meeting of the Association for Computational Linguistics*, 2021.

Wang, X., Wang, Z., Liu, J., Chen, Y., Yuan, L., Peng, H., and Ji, H. Mint: Evaluating llms in multi-turn interaction with tools and language feedback. In *12th International Conference on Learning Representations, ICLR 2024*, 2024.

Weidinger, L., Rauh, M., Marchal, N., Manzini, A., Hendricks, L. A., Mateos-Garcia, J., Bergman, S., Kay, J., Griffin, C., Bariach, B., et al. Sociotechnical safety evaluation of generative ai systems. *arXiv preprint arXiv:2310.11986*, 2023.

Weidinger, L., Raji, I. D., Wallach, H., Mitchell, M., Wang, A., Salaudeen, O., Bommasani, R., Ganguli, D., Koyejo, S., and Isaac, W. Toward an evaluation science for generative ai systems. *arXiv preprint arXiv:2503.05336*, 2025.

Xia, B., Lu, Q., Zhu, L., and Xing, Z. An ai system evaluation framework for advancing ai safety: Terminology, taxonomy, lifecycle mapping. In *Proceedings of the 1st ACM International Conference on AI-Powered Software*, pp. 74–78, 2024.

Yang, L., Clivio, O., Shirvaikar, V., and Falck, F. A critical review of causal inference benchmarks for large language models. In *AAAI 2024 Workshop on "Are Large Language Models Simply Causal Parrots?"*, 2023. URL https://openreview.net/forum?id=mRwgczYZFJ.

Zhan, J., Wang, L., Gao, W., Li, H., Wang, C., Huang, Y., Li, Y., Yang, Z., Kang, G., Luo, C., Ye, H., Dai, S., and Zhang, Z. Evaluatology: The science and engineering of evaluation. *BenchCouncil Transactions on Benchmarks, Standards and Evaluations*, 4(1):100162, 2024. ISSN 2772-4859. doi: https://doi.org/10.1016/j.tbench.2024.100162. URL https://www.sciencedirect.com/science/article/pii/S2772485924000140.

Zhang, A. K., Klyman, K., Mai, Y., Levine, Y., Zhang, Y., Bommasani, R., and Liang, P. Language model developers should report train-test overlap. *arXiv preprint arXiv:2410.08385*, 2024a.

Zhang, A. Q., Shaw, R., Anthis, J. R., Milton, A., Tseng, E., Suh, J., Ahmad, L., Kumar, R. S. S., Posada, J., Shestakofsky, B., et al. The human factor in ai red teaming: Perspectives from social and collaborative computing. In *Companion Publication of the 2024 Conference on Computer-Supported Cooperative Work and Social Computing*, pp. 712–715, 2024b.

Zhang, H., Lin, Y., and Wan, X. Pacost: Paired confidence significance testing for benchmark contamination detection in large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 1794–1809, 2024c.

Zhang, Y. and Kanayet, F. Genai evaluation maturity framework (gemf) to assess and improve genai evaluations. In *Workshop on Evaluating Evaluations: Examining Best Practices for Measuring Broader Impacts of Generative AI at NeurIPS 2024*, 2024. URL https://neurips.cc/virtual/2024/104201. Accessed May 10, 2025.

Zhang, Y., Zhang, M., Yuan, H., Liu, S., Shi, Y., Gui, T., Zhang, Q., and Huang, X. Llmeval: A preliminary study on how to evaluate large language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17):19615–19622, Mar. 2024d. doi: 10.1609/aaai.v38i17.29934. URL https://ojs.aaai.org/index.php/AAAI/article/view/29934.

Zhou, K., Blodgett, S. L., Trischler, A., Daumé III, H., Suleman, K., and Olteanu, A. Deconstructing nlg evaluation: Evaluation practices, assumptions, and their implications. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 314–324, 2022.

Zhou, K., Zhu, Y., Chen, Z., Chen, W., Zhao, W. X., Chen, X., Lin, Y., Wen, J.-R., and Han, J. Don't make your llm an evaluation benchmark cheater. *arXiv preprint arXiv:2311.01964*, 2023.

Zhou, X., Weyssow, M., Widyasari, R., Zhang, T., He, J., Lyu, Y., Chang, J., Zhang, B., Huang, D., and Lo, D. Lessleak-bench: A first investigation of data leakage in llms across 83 software engineering benchmarks. *arXiv preprint arXiv:2502.06215*, 2025.

# A. Appendix 1: The Conceptual Framework for Capability Evaluations

In this expanded version of the framework, several well-known challenges and limitations of evaluations are integrated within the elements and sub-elements; additional examples of specific practices or methodologies and their applications are included; and aspects of the sub-elements that the relevant literature identifies as reporting best practices are highlighted.

## A.1. Evaluation Target

Capability evaluations focus on what AI systems can and cannot do. The specific capability under evaluation might range from mathematical reasoning (e.g.Didolkar et al. (2024); Imani et al. (2023)) and causal inference (e.g. Kiciman et al. (2023); Jin et al. (2024)) to legal reasoning (e.g. Guha et al. (2023); Bhambhoria et al. (2024)), deception (e.g. Hagendorff (2024)), and many others.

A philosophical account defines a machine learning (ML) system as having a capability to X just when it would reliably succeed at doing X if it 'tried' (Harding & Sharadin, 2024; Burden, 2024). A further useful distinction is between competence and performance: while the former refers to the possession of a capability, performance refers to its successful manifestation. Thus, a failure in performance does not necessarily imply the absence of the underlying competence (Harding & Sharadin, 2024; Hagendorff et al., 2024).

The objective of the evaluation may be to quantify system capabilities, track progress, enable large-scale comparisons -as in traditional benchmarks (Reuel et al., 2024a), understand behavioral patterns, identify and estimate potential risks, provide assurance of system safety (Burden et al., 2025), predict future behavior (Burden, 2024) e.g. if upcoming models can lead to catastrophic harm (Barnett & Thiergart, 2024) or evaluate evaluation methods, among others.

An evaluation may therefore aim to compare the performance of different LLMs on document classification, or to assess a system's capacity to generate toxic content in order to identify potential misuse risks.

It is also important to recognize that the interpretation and implications of an evaluation may depend on the specific capability or objective, as well as on the context in which the system is implemented or intended to be implemented. For instance, the same evaluation objective—such as assessing a system's classification capability—may carry different significance depending on whether the system is deployed in a critical domain (e.g., healthcare, criminal justice) or a non-critical setting (e.g., entertainment recommendations) (OECD, 2022; Thurnherr, 2024).

## A.2. Task

As in human evaluations, there are many scenarios and tasks for evaluating a given capability. The task refers to a particular specification of a problem (Raji et al., 2021), typically represented as a mapping between an input space and an output or action space (Schlangen, 2021; Eriksson et al., 2025).

Part of the task specification involves the format of the response required from the model, which Burden et al. (2025) called the task mode, explaining that in some cases the subject is asked to identify an answer from a predetermined set of options (e.g., selecting a multiple-choice answer or a class label) or to generate a novel output (e.g., a numerical value in a continuous range or free-form text). This distinction between closed-ended and open-ended generation tasks has a direct impact on the design of the measurement component and the evaluation instance, as will be examined in section A5. Further examples of task configurations include question answering, text completion, text summarization, debate or error identification and correction, among others.

Furthermore, particularly in the case of agent evaluations, the task may involve performing multiple coordinated steps to accomplish complex objectives, simulating real-world scenarios (Kraprayoon et al., 2025). These intermediate sub-tasks can be autonomously defined by the system or dictated by a human subject. For instance, when evaluating agents in data science settings, some studies assess performance on end-to-end tasks that involve planning, generating code and plots, and producing coherent outputs and insights (Testini et al., 2025).

The rationale for selecting the task should be described, explaining their relevance to the evaluation objective and the specific capabilities being assessed (Cao et al., 2025). This demonstrates that the tasks are not arbitrary, but rather carefully chosen to enhance transparency and enable public scrutiny (Reuel et al., 2024a; Staufer et al., 2025; Liang et al., 2023).

One dimension to consider is the number of interactions between the evaluated subject and the user or evaluator. Along this dimension, interactions can be classified as single-turn—as in most traditional benchmarks—or multi-turn, when

they involve multiple exchanges (Wang et al., 2024; Cheng et al., 2024; Sirdeshmukh et al., 2025). A typical example of multi-turn interactions is found in evaluation methods based on debate (e.g. Moniri et al. (2025)).

The Second Draft of the GPAI Code of Practice (European Commission, 2024) highlights two types of validity. External validity refers to the extent to which evaluation results can be used as a proxy for model behavior in contexts outside of the evaluation environment (Weidinger et al., 2023; Paskov et al., 2024; Biderman et al., 2024; Reuel et al., 2025; Burden, 2024). This aligns with what Burden (2024) describes as situational external validity, which he distinguishes from external validity across subjects—the question of whether pre-existing tests remain valid when applied to AI systems. Internal validity, on the other hand, refers to the extent to which the observed evaluation results represent the truth in the evaluation setting and are not due to methodological shortcomings (Weidinger et al., 2023; Paskov et al., 2024; Reuel et al., 2025).

These principles have a direct impact on task selection. For instance, multiple-choice questions offer a fast and scalable method for evaluating model capabilities, as their closed-ended nature makes it straightforward to compute correctness. However, their external validity may be limited, as they rarely mirror the interactive, open-ended, or multimodal nature of real-world AI applications (Yang et al., 2023; AISI, 2024), as well as their dynamic characteristics (McIntosh et al., 2024; Eriksson et al., 2025). Additionally, as Bieger et al. (2016) point out, while focusing solely on final outputs—such as a completed multiple-choice test—can simplify black-box comparisons, it also risks discarding rich information about the system's behavior and reasoning processes over time.

While these principles resonate in many critical works on benchmarks and evaluations (Biderman et al., 2024; Reuel et al., 2024a; Weidinger et al., 2025; Cao et al., 2025; Liao & Xiao, 2023), they may not be absolute and often depend on the specific goals of the evaluation. For instance, there are entire evaluation paradigms that do not prioritize real-world representativeness. As Burden et al. (2025) describe, the Construct-Oriented Paradigm carefully designs tasks, controlling for confounding factors and often adapting tasks from cognitive science literature. Conversely, Narayanan & Kapoor (2023) critique the use of professional benchmarks by OpenAI to evaluate GPT-4, noting that *'it's not like a lawyer's job is to answer bar exam questions all day'* (Eriksson et al., 2025). These design decisions are informed by well-established and accepted methodologies from other disciplines. While such methods may not be the most appropriately representative of real-world performance, they prioritize different foundational goals. Therefore, the field of AI evaluation should be open to considering these diverse approaches.

## A.3. Evaluated Subject

The evaluated subject is the subject of the capability claim—in most cases, an AI system. Depending on the objective of the evaluation, this may involve a single system or multiple systems, particularly when the goal is to compare the performance of different systems with respect to the capability under evaluation. However, in the case of meta-evaluations, the subject of evaluation is an evaluation method itself —such as its design, reliability, or alignment with its intended purpose (e.g. Lee et al. (2025); Hong et al. (2024); Zhang et al. (2024d)).

An important consideration is to individuate the system—that is, to define its identity conditions (Harding & Sharadin, 2024). This is essential for ensuring transparency and enabling fair comparisons—or "apples to apples"—across systems (Biderman et al., 2024). Although models are often characterized by their architecture and parameters (Harding & Sharadin, 2024; Paskov et al., 2025), this information is not always available in practice[2]. Therefore, where possible, fixed attributes that help individuate the system should be reported, including its name and version, whether it is open- or closed-source, and whether it is a foundational model or a fine-tuned variant (Harding & Sharadin, 2024).

Because evaluations can be highly sensitive to seemingly minor implementation details—such as interface design, prompting techniques, access method, and whether elicitation techniques were employed to probe latent capabilities (Hofstätter et al., 2024; Paskov et al., 2025)—that may significantly influence model behavior (Biderman et al., 2024; McIntosh et al., 2024; Liang et al., 2023), it is also important to report the surrounding scaffolding and operational context (Harding & Sharadin, 2024). This includes documenting background conditions that may vary across evaluations, such as how the system was accessed (e.g., via a raw API or an assistant interface), configuration settings (e.g., temperature, context window, random seeds), any auxiliary tools or capabilities (e.g., plugins, memory, or browsing), and the inference budget, including the compute or hardware conditions under which the system was executed, where relevant (Harding & Sharadin, 2024; Staufer

---

[2]The level of access granted to the evaluation practitioner shapes the range of information available. Access can take the form of black-box (limited to inputs and outputs), grey-box (partial internal insights), white-box (full access to weights, gradients, and activations), or, as introduced by Casper et al. (2024), outside-the-box (contextual details such as training data, source code, and documentation).

et al., 2025; Paskov et al., 2025).

Additionally, it is important to consider and report intra-model variance, the variation in performance that can arise from repeated evaluations of the same model under different conditions (e.g. sampling temperatures and random seeds, as mentioned above) (Reuel et al., 2024a). Reporting this, can help distinguish genuine performance differences between systems from noise introduced by the evaluation (Reuel et al., 2024a).

### A.4. System Inputs

The input source refers to the origin of the inputs used in the evaluation (Dow et al., 2024). These may come from a pre-existing test dataset, or user-generated inputs, including those produced dynamically interacting with the system. For example, in red teaming, attacks on the target model are often conducted through trial and error to identify effective strategies that elicit abnormal behavior (Lin et al., 2025). The input can also take various modalities, such as image, audio, video, text, or multimodal combinations.

The test dataset refers to a separate dataset, distinct from the training data, consisting of data instances used to evaluate the AI system on a given task. In other cases, when generalization is not a primary concern, such separation between test dataset and training data may not be required. In the case of benchmarks and reference-based evaluations, the dataset also includes the desired output and annotations for each instance (Eriksson et al., 2025; Liu et al., 2024).

The dataset may originate from widely used benchmarks or alternative sources and can undergo preparation and curation processes prior to evaluation (Cao et al., 2025). Key dimensions to consider include the dataset's size and construction method—for example, whether it was manually developed by researchers or domain experts, or synthetically generated by an LLM—and, where applicable, considerations related to informed consent, permissions (Piktus et al., 2023) and license compliance. Transparent reporting of these factors enables a better understanding of the dataset's context and limitations, and supports the consideration of relevant ethical implications (Reuel et al., 2024a; Staufer et al., 2025).

Zhang & Kanayet (2024) argue that a good evaluation should include a mix of item difficulties. Transparency is also enhanced by analyzing the composition of the dataset in terms of task difficulty and task type, among other aspects. Often, such analysis depends on the specific capability being evaluated or the domain to which it belongs. For example, Bai et al. (2023) use Bloom's Taxonomy to classify the questions in their dataset and calculate the distribution of question forms based on interrogative words (e.g., how, what, why). Similarly, in the field of causal reasoning, researchers have applied the Ladder of Causation to structure tasks according to levels of causal difficulty (Jin et al., 2023). Beyond domain-specific frameworks, there have also been efforts to estimate task difficulty a priori using textual features or model-based uncertainty (Benedetto, 2023; Kuhn et al., 2023). To support more principled test construction, Item Response Theory (IRT) has been applied in NLP to describe characteristics of individual items—their difficulty and discriminating power—and to account for these characteristics in the estimation of system performance (Lalor et al., 2016). This approach has been used to identify datasets that are best suited for distinguishing among state-of-the-art models, as well as those that are largely saturated and unlikely to detect further progress (Vania et al., 2021).

According to Harding & Sharadin (2024) the dataset can be thought of as a particular set of queries. A prompting strategy then serves to turn each query into an input which can be processed by an LLM. In this context, most prompting strategies can be viewed as 'templates', which generate the model input by embedding the query in a pre-written prompt format (Harding & Sharadin, 2024). In addition, prompts may also include model-specific instructions designed to steer the model's behavior toward the desired outcome (Schulhoff et al., 2024). Some of the prompt techniques include for example, few shots, one shot, chain-of-thought (CoT), self-consistency or emotion prompting (Sahoo et al., 2024; Schulhoff et al., 2024). Since even minor variations in prompts can significantly affect model performance (Perez et al., 2023; Sclar et al., 2024; Frontier Model Forum, 2024; McIntosh et al., 2024; Liang et al., 2023), it is important to report the full prompts used (Reuel et al., 2024a) in order to ensure that all data is available to support the reproducibility of evaluations (Biderman et al., 2024; Cao et al., 2025).

In general, the greater the number of inputs used in an evaluation, the more reliable and robust the results will be as evidence for claims about a model's capabilities (Harding & Sharadin, 2024). Additionally, because of the nondeterministic nature of LLMs, experiments should be repeated, and randomization strategies should be used to mitigate the effects of randomness and parameter configuration biases (Cao et al., 2025; Liang et al., 2023).

An important challenge in this dimension is data contamination -also known as data leaking (Zhou et al., 2025), which occurs when the test dataset overlaps with the training data (McIntosh et al., 2024), leading to an overestimation of the

system's performance (Sainz et al., 2023; Zhou et al., 2023; Laskar et al., 2024). Numerous strategies have been proposed to identify and mitigate this issue, making evaluations "google-proof" (Ivanova, 2023) (e.g. Sainz et al. (2023); Jacovi et al. (2023); Dong et al. (2024); Piktus et al. (2023)) and placing responsibility on the evaluation practitioner. However, it has also been argued that model developers share this responsibility (Reuel et al., 2024a), particularly in transparently reporting train-test overlap statistics along with evaluation results (Zhang et al., 2024a) or planning the evaluation benchmark prior to training (Piktus et al., 2023).

### A.5. Evaluation Instance

System outputs can be judged from various angles, such as correctness or fluency, each corresponding to a distinct evaluation criterion (Zhang et al., 2024d; Zhou et al., 2022), also referred to as a normative baseline (Weidinger et al., 2023). An evaluation criterion refers to a specific aspect of the system's output, necessary to define what constitutes a correct, appropriate, or high-quality response. The choice of criteria depends on both the nature of the task under evaluation (Van Der Lee et al., 2019) and the evaluation target element. Additional criteria may include generalization, reasoning, robustness (Hu & Zhou, 2024), logical coherence, harmlessness (Zhang et al., 2024d), and naturalness (Van Der Lee et al., 2019), among others.

If the evaluation criterion is not mathematically defined and no prior annotation of the dataset exists, it typically needs to be applied to the system's outputs by an evaluator subject, which may be either a human or another AI system. In all cases, it is recommended that evaluation criteria be specified as precisely as possible (Zhou et al., 2022; Van Der Lee et al., 2019) to minimize subjectivity and potential sources of bias.

When an evaluator subject is involved, several additional factors must be considered. These include how the evaluator is selected—for instance, if the evaluator is human, they may be recruited via crowdsourcing platforms; if the evaluator is another LLM, preferably, it should not be the same model under evaluation due to self-preference bias (Panickssery et al., 2024), also referred to as egocentric bias (Li et al., 2024), and it should not belong to the same model family, as this can introduce bias toward similar linguistic styles (Bai et al., 2023). Moreover, the evaluator system is typically expected to be more capable than the one being assessed (Li et al., 2024).

Additionally, the evaluator's level of expertise must be considered. For human evaluators, this last factor may involve distinguishing between laypeople and domain experts; for LLM evaluators, an analogous consideration is whether the model has been fine-tuned specifically for domain-specific evaluation tasks (e.g. Huang et al. (2024)).

The nature and instructions of the evaluator's task must also be specified, such as pairwise comparison -as in Chatbot Arena (Chiang et al., 2024)-, rating individual outputs on a Likert scale or assigning a score (Li et al., 2024). The order in which outputs are presented is also important, as both human evaluators (Van Der Lee et al., 2019), and AI systems are susceptible to position bias (Gao et al., 2025). Finally, in human evaluations, it is considered good practice to always report the number of participants along with relevant demographic data (e.g., gender, nationality, age, fluency in the target language, academic background) (Weidinger et al., 2023) to enhance replicability and allow readers to assess the significance of the results (Van Der Lee et al., 2019).

Note that in this framework, we distinguish between the concepts of annotators and evaluators. An *annotator*—whether a human or an AI system—is responsible for labeling a test dataset prior to the evaluation of the evaluated subject. This may involve selecting the correct answer in a multiple-choice item, assigning categories to text fragments, or specifying the desired answer to a question. In contrast, an evaluator assesses and annotates the outputs of the evaluated subject, for instance by rating the creativity of a poem generated by an LLM using a Likert scale. While the term 'annotator' is often used in the literature to refer to both roles (e.g. Weidinger et al. (2023)), we maintain this distinction for clarity.

### A.6. Measurement

A metric is a mathematically defined function used to quantify a specific evaluation criterion. Metrics provide a standardized and objective way to measure aspects of system outputs—such as accuracy, coherence, fluency, or similarity to a reference—as they can summarize model performance on a set of tasks and datasets as a single number or score (Raji et al., 2021; Eriksson et al., 2025).

Evaluation practitioners may confuse evaluation criteria with the metrics used to operationalize them, but this confusion risks losing the ability to assess the metric's appropriateness and distinguish between the metric and the original goal (Zhou et al., 2022). As Goodhart's Law states, 'When a measure becomes a target, it ceases to be a good measure.'

Broadly, evaluation metrics can be categorized as either reference-based or reference-free (Ito et al., 2025; Li et al., 2024). Reference-based metrics are used to compare the generated output to a predefined reference, which may be produced manually by human annotators or by another AI system, typically focusing on accuracy, relevance, coherence and similarity to the reference (Li et al., 2024). However, because generating such references can be costly—and in some cases, a ground truth may not be available—reference-free metrics estimate the quality of the output without relying on a reference, often concentrating on its intrinsic qualities, such as fluency or context relevance[3] (Li et al., 2024). Nevertheless, due to certain limitations—such as biases against higher-quality outputs—this kind of metrics are often better suited as diagnostic tools, rather than as definitive measures of task performance (Deutsch et al., 2022).

Validity concerns also apply to evaluation metrics. Internal validity may be compromised when inconsistencies in implementation or parameter choices introduce variation that undermines the reliability of metric scores (Liao et al., 2021). External validity issues, in turn, can arise from a mismatch in the evaluation metric of interest (Liao et al., 2021; Gehrmann et al., 2023). To address these limitations, automatic metrics are often used to complement human evaluation, which has long been considered the gold standard (Liao & Xiao, 2023; Ruan et al., 2024). This combination allows practitioners to assess the degree of correlation between automatic and human judgments, and to prioritize which models should undergo more resource-intensive human evaluation processes such as crowdsourcing (Zhou et al., 2022). For this reason it is also important to justify the choice of automatic metrics (Reuel et al., 2024a; Zhou et al., 2022).

Another distinction is between general-purpose metrics and task-specific metrics. General-purpose metrics are widely used across a variety of tasks to assess model performance and can be further categorized following Hu & Zhou (2024). Multiple-classification metrics evaluate how effectively an AI system classifies items into multiple groups, examples include Accuracy, Recall, Precision, and F1 score. Token-similarity metrics measure the similarity between texts generated by AI systems and corresponding reference texts; examples include Perplexity, BLEU, ROUGE, BERTScore, and METEOR. Question-answering metrics, such as Strict Accuracy, Lenient Accuracy, and Mean Reciprocal Rank, are used to evaluate performance on QA tasks.

In contrast, task-specific metrics are designed to meet the unique requirements of particular tasks. For example, Translation Edit Rate (TER) can be used to evaluate machine translation output (Snover et al., 2006). Other examples include metrics developed for debate experiments, such as Aggregate Rating and Win Rate (Khan et al., 2024), as well as the toxicity metric used in red teaming evaluations (Hong et al., 2024) or Attack Success Rate(Lin et al., 2025).

It is important to report the exact formulas or processes used to calculate these metrics, along with any parameters (Reuel et al., 2024a) and hyperparameters to allow fair comparisons (Gehrmann et al., 2022).

If the evaluator's task involves selecting the preferred response from a set, a ranking system is required to aggregate individual judgments and determine which system is preferred overall. In the case of pairwise comparisons, common approaches include the Elo rating system or Points Scoring System (Zhang et al., 2024d).

Finally, when the goal of an evaluation is to compare the performance of different systems, approaches, or evaluation methods, one or more baselines are typically used as reference points or standards for comparison (see, for example Hong et al. (2024); Shi et al. (2024)). Baselines help contextualize the performance of the models under evaluation. By comparing a model's results against a baseline, researchers can better understand whether improvements are meaningful and how significant they are relative to established expectations. Additionally, human performance in the evaluation may be included as a floor, allowing readers to put the system's performance into perspective (Reuel et al., 2024a).

### A.7. Results Analysis

Ranking models according to a single quality number is easy and actionable, but, in many circumstances, it is much more important to understand when and why models fail (Eriksson et al., 2025). Therefore, it is also beneficial to conduct a deeper analysis of the evaluation results, as this can reveal valuable insights and guide future research directions (Hämäläinen & Alnajjar, 2021). This typically involves providing a concise summary of key findings—both quantitative and qualitative (e.g., analyzing results in relation to dataset complexity levels, if a dataset composition analysis was conducted)—, and outlining potential areas for future work—explaining how the results inform model development and suggest avenues for improvement. Interpreting results in light of contextual thresholds for what constitutes 'acceptable' model performance—that may vary across use cases, domains, or applications—is a component of this stage (Weidinger et al., 2023).

---

[3]The examples of accuracy, relevance, coherence, fluency and context relevance reflect typical focuses in natural language generation evaluation; other domains may emphasize different intrinsic or comparative qualities.

Additionally, when relevant, statistical analysis should be performed to distinguish genuine performance improvements from random variability. Best practices emphasize treating evaluation processes as structured experiments, applying significance testing and quantifying uncertainty through confidence intervals or resampling techniques (Miller, 2024; Sun et al., 2024). Methods such as ANOVA, Tukey's HSD, Student's t-test, Mann–Whitney U test, and McNemar's test can be adopted to compare models under various settings, enabling more robust conclusions across different architectures and training regimes (Sun et al., 2024). Moreover, careful statistical assessment is useful for detecting biases or artifacts in benchmark performance, ensuring that reported gains reflect true advances rather than contamination or evaluation artifacts (see for example, Zhang et al. (2024c)). Incorporating these practices helps build more reliable and reproducible evaluations, particularly in high-stakes or competitive settings where small metric differences may lead to significant decisions.

In human evaluations, it is considered good practice to calculate inter-annotator agreement (Reuel et al., 2024a) as a means of assessing the consistency of the judgments.

Finally, it is recommended to publish system outputs and evaluation code, enabling others to reproduce the findings (Eriksson et al., 2025), re-score outputs under different criteria, or conduct their own statistical analyses (Biderman et al., 2024). Facilitating reproducibility is particularly important, as it remains a major concern in this field (Reuel et al., 2024a).

## A.8. Transversal Elements of the Conceptual Framework

**Evaluation Practitioners.** The human subject responsible for designing and conducting the evaluation may come from industry, government, academia, or act as an independent third-party evaluator. The degree of human involvement can vary depending on the extent to which components of the evaluation are automated (Eriksson et al., 2025). While this may seem self-evident, highlighting this role is important. For example, in human evaluations, the person designing the evaluation should not be the same individual who assesses the system's outputs, as preexisting hypotheses or expectations may bias their judgment. To mitigate this risk, output assessments are often delegated to external raters, such as those recruited through crowdsourcing platforms.

Additional considerations include the importance of multidisciplinary and diversity. With regard to the former, the natural limitations of the evaluation designers' knowledge on a potentially infinitely large number of domains and tasks, may lead to generalist approaches that often fail to address the subtle requirements of critical sectors or approaches which could hinder innovation (McIntosh et al., 2024; Eriksson et al., 2025). As for the latter, evaluation is never neutral (Weidinger et al., 2023; Rauh et al., 2024): many decisions implicitly reflect specific epistemological perspectives regarding how the world is ordered (Orr & Kang, 2024; Eriksson et al., 2025). Moreover, the social and cultural environments in which evaluations are conducted shape them through shared and often arbitrary assumptions, commitments, and dependencies (Eriksson et al., 2025).

**Ethical Considerations.** Gupta et al. (2024) present a conceptual framework for balancing information gain and ethical harm in the evaluation of AI systems. They argue that ethical issues arising during evaluation are often overlooked or inadequately predicted, which can diminish the acceptability of the results and the overall value and utility of the information obtained. For example, adversarial testing practices and red teaming may introduce ethically harmful impacts on practitioners or data annotators during the model output labeling process (Zhang et al., 2024b; Gupta et al., 2024; Weidinger et al., 2023).

**Pilot Tests.** Pilot tests refer to preliminary trials conducted prior to the full-scale evaluation to refine the evaluation setup and detect potential problems or vulnerabilities early in the process. This may include verifying the clarity of task instructions and identifying ambiguities or biases in the data or scoring procedures. Pilot testing helps ensure the robustness of the evaluation design and is especially valuable when human evaluators are involved (Hämäläinen & Alnajjar, 2021) or when novel methods are proposed (see for example, Bai et al. (2023)).

**Ablation Studies.** Ablation studies refer to a research methodology used to investigate the contributions of individual components within a system or process. This involves systematically removing or altering specific elements of a system to observe the effect of their absence or modification on the overall performance or behavior. By isolating the impact of each component, researchers can identify which factors are essential for achieving desired outcomes and which are less influential (for examples in evaluations, see Shi et al. (2024); Hong et al. (2024)). This approach is particularly useful for refining theories, improving system design, and studying the robustness and efficiency of complex models or processes (Meyes et al., 2019).

# B. Appendix 2: Use Cases

This appendix presents three application cases of our evaluation framework, demonstrating its applicability across a diverse range of existing evaluation settings. The first case involves LLMs acting as evaluators (Bai et al., 2023), the second focuses on a meta-evaluation of red-teaming methods (Lee et al., 2025), and the third examines the evaluation with AI agents (Rasheed et al., 2024). While the selection is not exhaustive and some details may be simplified or open to interpretation, the aim is to illustrate how the proposed framework can be applied to analyze real-world evaluation practices.

| Elements of the Conceptual Framework | Paper 1: *"Benchmarking Foundation Models with Language-Model-as-an-Examiner"* | Paper 2: *"Learning Diverse Attacks on Large Language Models for Robust Red Teaming and Safety Tuning"* | Paper 3: *"Large Language Model Evaluation Via Multi AI Agents: Preliminary Results"* |
|---|---|---|---|
| **Evaluation Target** | | | |
| **Capability** | Examiners: the ability of LLMs to serve as examiners and their potential biases in a decentralized setting. Evaluated LLMs: question answering. | Capacity to elicit undesirable responses from a target LLM, through the generation of diverse and effective attacks. | Code generation. |
| **Objective** | To evaluate if the proposed decentralized LLM-based evaluation method mitigate the lack of diversity and scope of the generated questions and bias during evaluation compared to a centralized evaluation. | To evaluate whether the proposed method for improving the diversity and effectiveness of attacks performs better, worse, or comparably to existing approaches. | To evaluate and compare the code generation capabilities of various LLMs using a novel multi-agent AI framework with automatic verification. |
| **Task** | | | |
| **Task Mode** | Open-ended question answering. | Attacker Model Task: Prompt generation. Target Model Task: Open-ended question answering. | Open-ended code generation. |
| **Steps** | Single-step | Single-Step | Multi-step |
| **Interactions** | Both single-turn and multi-turn. | Single-turn. | Single-turn interaction (however, the authors reported several rounds of interaction for the OpenAI series). |
| **Evaluated Subject** | | | |

| Elements of the Conceptual Framework | Paper 1: *"Benchmarking Foundation Models with Language-Model-as-an-Examiner"* | Paper 2: *"Learning Diverse Attacks on Large Language Models for Robust Red Teaming and Safety Tuning"* | Paper 3: *"Large Language Model Evaluation Via Multi AI Agents: Preliminary Results"* |
|---|---|---|---|
| **Evaluated Subject** | For the centralized evaluation 8 LLMs were evaluated: BLOOMZ, Flan-T5, Flan-UL2, GLM- 130B, LLaMA, Vicuna-13B, and ChatGPT. For the descentralized evaluation 4 LLMs were evaluated: ChatGPT, Claude, Vicuna-13B, Bard. | Since this is a meta-evaluation, the primary subject of evaluation is the red-teaming method itself, specifically the attacker LLM. The paper mentions the evaluation of two attacker models: GPT-2 small and Llama-3.2-1B. To validate the method's effectiveness, five target LLMs were red-teamed: GPT-2, Gemma-2b-it, Dolly-v2-7b, Llama-2-7b-chat, and Llama-3.1-8B-Instruct. | 7 LLMs were evaluated: ChatGPT-4 Turbo, GPT-4, GPT-3.5 Turbo, GPT-3.5, Google Bard, Llama and Hugging Face (CodeBERT). |

| Elements of the Conceptual Framework | Paper 1: *"Benchmarking Foundation Models with Language-Model-as-an-Examiner"* | Paper 2: *"Learning Diverse Attacks on Large Language Models for Robust Red Teaming and Safety Tuning"* | Paper 3: *"Large Language Model Evaluation Via Multi AI Agents: Preliminary Results"* |
|---|---|---|---|
| **Identity Conditions** | Number of parameters and training procedure were reported for the centralized evaluation: BLOOMZ (the 176B model), Flan-T5 (the XXL model, 11B), Flan-UL2 (20B), GLM-130B, LLaMA (the 13B model and the 65B model), Vicuna-13B, and ChatGPT. These models are categorized based on their training procedure: whether they have undergone Supervised Fine-Tuning (SFT) or not. The first 6 models are trained without SFT, whereas the last 2 models are fine-tuned. | The proposed method use GFlowNet fine-tuning, followed by a secondary smoothing phase, to train the attacker model to generate diverse and effective attack prompts. The attacker model is defined by its probabilistic formulation—specifying its objective, sampling policy, training constraints, and fine-tuning procedure —as detailed in Section 3.2 of the paper. The attackers models include GPT-2 Small (124M parameters) and Llama-3.2-1B (1.2B parameters)// The target models evaluated include GPT-2 (124M parameters), Gemma-2b-it (2B parameters), Dolly-v2-7b (7B parameters), Llama-2-7b-chat (7B parameters), and Llama-3.1-8B-Instruct (8B parameters). | Number of parameters: ChatGPT-4 Turbo - 1.96 trillion, GPT-4 - 1.76 trillion, GPT-3.5 Turbo - 154 billion, GPT-3.5 - 125 billion, Google Bard - 1.56 trillion, Llama - 70 billion and Hugging Face (CodeBERT) - 355M. |
| **Operational Context** | Centralized evaluation: the temperature was set to 0 for both the examiner and the subject models to ensure reproducibility. | All target models are accessed in a black-box setting, meaning only prompt-response behavior is observed. | LLMs were accessed directly via raw APIs. 7 AI Agents were responsible for interacting with a different evaluated LLM to retrieve programming code. |
| **System Inputs** | | | |

| Elements of the Conceptual Framework | Paper 1: *"Benchmarking Foundation Models with Language-Model-as-an-Examiner"* | Paper 2: *"Learning Diverse Attacks on Large Language Models for Robust Red Teaming and Safety Tuning"* | Paper 3: *"Large Language Model Evaluation Via Multi AI Agents: Preliminary Results"* |
|---|---|---|---|
| | *Continued from previous page* | | |
| **Amount** | Centralized evaluation: Single-Round: 10,000 questions. Multi-round: randomly selected 1000 questions. // Descentralized evaluation: each examiner model posed 100 questions. | Not reported. | 10 input prompts used across the 7 evaluated models. |
| **Modality** | Text. | Text. | Text. |
| **Input Source** | The test dataset, namely LMExamQA, was generated by an LLM. To ensure wide coverage of knowledge, Google Trends Categories were chosen as the domain taxonomy, and n domains from it were randomly selected. For each domain, the LLM was prompted to generate m distinct questions. Their designed prompt is formulated to ensure that the generated questions possess three essential characteristics: diversified question forms, varied cognitive levels, and most importantly, assurance that the LM has a comprehensive understanding of the knowledge surrounding the question it poses. The resulting dataset was analyzed according to Bloom's taxonomy. | Adversarial prompts generated by a separate attacker model. | Each LLM was provided with identical, high-level natural language descriptions of programming tasks. The test dataset was created ad-hoc, aiming at incorporating diverse coding tasks to comprehensively evaluate the model's coding capabilities. |

| Elements of the Conceptual Framework | Paper 1: *"Benchmarking Foundation Models with Language-Model-as-an-Examiner"* | Paper 2: *"Learning Diverse Attacks on Large Language Models for Robust Red Teaming and Safety Tuning"* | Paper 3: *"Large Language Model Evaluation Via Multi AI Agents: Preliminary Results"* |
|---|---|---|---|
| *Continued from previous page* | | | |
| **Prompt Techniques** | Centralized Evaluation: For models without SFT, their assess their 0-shot and 5-shot performance. | Not Reported. | The prompting technique used in the paper can be characterized as standardized zero-shot instruction prompting. The natural language descriptions of programming tasks were provided without examples or contextual demonstrations. |
| **Evaluation Instance** | | | |
| **Evaluation Criteria** | (1) Accuracy. This assesses the extent to which the provided response accurately answers the question. (2) Coherence. This evaluates the logical structure and organization of the response and the degree to which it can be comprehended by non-specialists. (3) Factuality. This examines whether the response contains factual inaccuracies. (4) Comprehensiveness. This gauges whether the response encompasses multiple facets of the question, thus providing a thorough answer. | Diversity and Toxicity. | Syntactic correctness, adherence to the prompt, computational efficiency in terms of execution and resource usage, code accuracy (functional correctness of the generated code as per the given description) and innovation demonstrated in the solutions. |

| Elements of the Conceptual Framework | Paper 1: *"Benchmarking Foundation Models with Language-Model-as-an-Examiner"* | Paper 2: *"Learning Diverse Attacks on Large Language Models for Robust Red Teaming and Safety Tuning"* | Paper 3: *"Large Language Model Evaluation Via Multi AI Agents: Preliminary Results"* |
|---|---|---|---|
| **Evaluator Subject** | The novel decentralized method incorporates multiple models to serve as examiners, namely Peer-examination. Centralized experiment: GPT-4. Descentralized experiment: GPT-4, Claude (Claude-instant), ChatGPT, Bard, and Vicuna-13B. To justify the reliability of the LM examiner, it was tasked with generating ground-truth answers, and a random sample of 100 questions was evaluated by human experts. | **RoBERTa** hate speech classifier (Vidgen et al., 2021) for GPT-2 and dolly-v2-7b, and **Llama-Guard** (LLM) for Llama-2-7b-chat, Llama-3.1- 8B-Instruct, and Gemma-2b-it. | The primary evaluator subject in this study is an automated verification agent, specifically designed to assess code generated by multiple LLMs. This agent employs the HumanEval benchmark and computes the pass@k metric (with k=1), allowing for consistent, objective, and reproducible evaluation of functional correctness across models. |
| **Evaluator Task** | Likert Scale Scoring and Ranking. | The evaluator's task is to assess the toxicity of the output generated by the target language models in response to adversarial prompts. Given a prompt–response pair (x,y), the evaluator—either the RoBERTa hate speech classifier or Llama-Guard—assigns a toxicity score or label, which determines whether the response is considered harmful. The output is binary variable denoting toxicity. A prompt is considered toxic if the toxicity classifier assigns a score greater than 0.5. | Rating of individual outputs via the HumanEval Benchmark. Additionally, a quality rating, depicted with stars in section 4.2 of the paper provides a subjective assessment of the code based on criteria such as readability, efficiency, and adherence to best practices. |
| **Measurement** | | | |

| Elements of the Conceptual Framework | Paper 1: *"Benchmarking Foundation Models with Language-Model-as-an-Examiner"* | Paper 2: *"Learning Diverse Attacks on Large Language Models for Robust Red Teaming and Safety Tuning"* | Paper 3: *"Large Language Model Evaluation Via Multi AI Agents: Preliminary Results"* |
|---|---|---|---|
| **Metric** | Likert scale scoring and a variant of pairwise comparison, namely ranking. Liket Scale: Each criteria is scored on a scale of 1 to 3, ranging from worst to best. The evaluator was also asked to provide an overall score ranging from 1 to 5, based on the scores assigned to the previous 4 criteria. This score serves as an indicator of the overall quality of the answer. In the pairwise comparison evaluators are given two responses and are tasked with determining which is superior, taking into account the four evaluation criteria. An evaluation of the metrics was conducted to know whether they correlate with human judgments. | **Toxicity: Toxicity rate** which is the percentage of generated prompts that are toxic. **Diversity: Average pairwise cosine distance** to measure diverity. A model named MiniLMv2 (sentence-transformer model) is used to embed the generated prompts so the distance can be calculated. | The primary metric used in this study is pass@k, from the HumanEval benchmark that quantifies functional correctness by measuring the proportion of generated code outputs that successfully pass predefined unit tests on the first attempt. In the experiments, a pass@1 approach was utilized, meaning that we evaluated the models based on their first attempt at generating a solution. This approach was chosen to reflect a more realistic scenario where a developer would use the model's first output. |
| **Baseline** | The centralized evaluation method was used as a baseline to compare against the proposed approach. | The method was compaired against some relevant red-teaming baselines: **Supervised Fine-tuning (SFT), In-Context Learning (ICL), REINFORCE, PPO + Novelty, GFlowNet, GFlowNet + MLE**. | Not used. |
| **Result Analysis** | | | |

| Elements of the Conceptual Framework | Paper 1: *"Benchmarking Foundation Models with Language-Model-as-an-Examiner"* | Paper 2: *"Learning Diverse Attacks on Large Language Models for Robust Red Teaming and Safety Tuning"* | Paper 3: *"Large Language Model Evaluation Via Multi AI Agents: Preliminary Results"* |
|---|---|---|---|
| **Qualitative Analysis** | They used a win-rate heatmap to visualize the results. Conclusions: 1. The results adhere to the scaling law of LMs; 2. Few-shot leads to more substantial improvement on higher cognitive-level questions; 3. SFT primarily plays a crucial role in aligning LM's responses for task adaptation, rather than enriching the model's knowledge — especially in the context of higher-level questions that demand more sophisticated answers; 4. LLMs can provide factually correct and coherent responses, but struggle for more comprehensive accurate answers; 5. They observe that excluding ChatGPT and Vicuna-13B, all examinee models exhibit a notable decrease in performance in the second round. This suggests that while these models initially demonstrated a robust understanding and knowledge base, their performance deteriorated when faced with more complicated questions. | The proposed method, GFlowNet + MLE, outperforms all baseline approaches in generating diverse and effective adversarial prompts across five target language models. It achieves a better balance between toxicity rate and prompt diversity, compared to prior methods like PPO + Novelty and REINFORCE, which either collapse to a few toxic prompts or fail to generate effective ones. Prompts generated by GFlowNet + MLE also show strong transferability, successfully attacking unseen models such as GPT-4o and larger LLaMA variants. Moreover, safety-tuning models using these prompts leads to more robust defenses against other red-teaming attacks without harming general performance. The second stage (MLE smoothing) is also computationally efficient, adding robustness with minimal additional training cost. | Conclusions: 1. Model Performance Is Not Strictly Correlated with Parameter Size. Despite having significantly fewer parameters than models like GPT-4 or Google Bard, GPT-3.5 Turbo outperformed all other models in terms of functional correctness; 2. Code Quality Varies Widely Across Models. The subjective star-based ratings revealed notable differences in code quality beyond correctness; 3. Instruction Adherence Is Uneven Across Systems. While all models received identical natural language prompts, their ability to faithfully implement the described functionality varied; 4. Evaluation Framework Enhances Comparative Interpretability. The multi-agent evaluation setup, combined with the verification agent and HumanEval benchmark, proved effective in revealing relative strengths and weaknesses across LLMs. |

| Elements of the Conceptual Framework | Paper 1: *"Benchmarking Foundation Models with Language-Model-as-an-Examiner"* | Paper 2: *"Learning Diverse Attacks on Large Language Models for Robust Red Teaming and Safety Tuning"* | Paper 3: *"Large Language Model Evaluation Via Multi AI Agents: Preliminary Results"* |
|---|---|---|---|
| **Statistical Analysis** | Not reported. | There is descriptive statistical analysis (means and standard deviations across multiple runs), but there's no inferential statistical testing like p-values or confidence intervals. | No Statistical Analysis was mentioned. Results were reported in terms of raw accuracy counts (i.e., number of correct code generations out of 10 per model) and qualitative star-based ratings. |
| **Future Work** | Expanding the framework to incorporate more domain-specific language models, or even vision language models, could potentially offer a more holistic evaluation. | The approach is not limited to text tokens and future work can explore the applicability to red-team multimodal models (e.g., text-to-image models ). Further, an interesting area of future work is extending the approach to the jailbreaking setting, where an attacker language model generates a suffix for an adversarial query prompt. Finally, in addition to red-teaming, it would be interesting to apply our method to generate prompts which can improve model performance on different tasks. | The future work outlined in the paper involves three key directions: (1) expanding the evaluation dataset from 10 to 50 input descriptions to enable broader and more robust analysis; (2) integrating the MBPP (Massively Multitask Benchmark for Python) to complement HumanEval and provide more diverse and challenging test cases; and (3) conducting real-world validation by sharing the model with 20 practitioners from diverse backgrounds, with the aim of collecting qualitative feedback to enhance the model's practical relevance and usability. |