



A Benchmark Model for Language Models towards Increased Transparency

Ayşe Kok Arslan

Abstract: One of the mostly advanced AI technologies in recent year has been language models (LM) which necessitate a comparison or benchmark among many LM to enhance transparency of these models. The purpose of this study is to provide a fuller characterization of LMs rather than to focus on a specific aspect in order to increase societal impact. After a brief overview of the constituents of a benchmark and features of transparency, this study explores main aspects of a model - scenario, adaptation, metric- required to provide a roadmap for how to evaluate language models. Given the lack of studies in the field it is a step towards the design of more sophisticated models and aims to raise awareness of the importance of developing benchmarks for AI models.

Introduction

The original promise of computing was to solve information overload in science. In his 1945 essay "As We May Think", Vannevar Bush (1945) proposed computers as a solution to manage the growing mountain of information. Licklider (1960) expanded on this with the vision of a symbiotic relationship between human-beings and machines so that computers would be "preparing the way for insights and decisions in scientific thinking" (Licklider, 1960).

Within this spirit, in the past couple of years, LMs have continued to push the limits of what is possible with deep neural networks. However, when it comes to topics such as understanding, reasoning, planning, and common sense, scientists are divided about how to assess LMs.

At its core, a LM is a box that takes in text and generates text (Figure 1). LMs are general purposes text interfaces that could be applied across a vast expanse of scenarios. For each scenario, there may be a broad set of desiderata such as accuracy, fairness, efficiency etc. among many others.

This rapid proliferation of LMs necessitates a comparison or benchmark among many language models. Benchmarks encode values and priorities (Ethayarajh and Jurafsky, 2020; Birhane et al., 2022) that specify directions for the AI community to be improved upon (Spärck Jones and Galliers, 1995; Spärck Jones, 2005; Kiela et al., 2021; Bowman and Dahl, 2021; Raji et al., 2021).

Overview of Benchmarks

Benchmarks are one of the thorniest problems of AI research. On the one hand, researchers need a way to evaluate and compare their models. On the other hand, some concepts are really hard to measure.

One of the major problems underpinning benchmarks is that we usually view them from a human intelligence perspective. As a simplified example, we consider chess as a complicated intelligence challenge because, on their way to mastering chess, human-beings must acquire a set of cognitive skills through hard work and talent. Yet, from a computational perspective, there can be a shortcut for finding good chess moves through a good algorithm, and the right inductive biases.

As this example demonstrates, even some of the most carefully crafted benchmarks can be prone to computational shortcuts. In other words, while benchmarks are a good tool to compare machine learning models against one another, they are not anthropomorphic measures of cognitive skills in machines.

When implemented and interpreted appropriately, benchmarks enable the broader community to better understand AI technology and influence its trajectory. In general, a benchmark involves three elements (Figure 1.):

- (1). **Broad coverage and recognition of incompleteness:** As it is not possible to consider all the scenarios and the desiderata that (could) pertain to LMs, a benchmark should provide a top-down taxonomy and make explicit all the major scenarios and metrics that are missing.
- (2). **Multi-metric measurement:** Societally beneficial systems reflect many values, not just accuracy. A benchmark should represent these plural desiderata, evaluating every desideratum for each scenario considered.
- (3). **Standardization:** As the object of evaluation is the LM, not a scenario-specific system, the strategy for adapting an LM to a scenario should be controlled for.



Metric		Metrics						
Scenarios		Accuracy	Calibration	Robustness	Fairness	Bias	Toxicity	Efficiency
	Natural Questions	✓ (Accuracy)	✓	✓	✓	✓	✓	✓
XSUM	✓ (Accuracy)	✓	✓	✓	✓	✓	✓	✓
AdversarialQA	✓ (Robustness)	✓	✓	✓	✓	✓	✓	✓
RealToxicity Prompts	✓ (Toxicity)	✓	✓	✓	✓	✓	✓	✓
BBQ	✓ (Bias)	✓	✓	✓	✓	✓	✓	✓

Fig. 1.0 Overview of an LM Benchmark Example

Overall, a benchmark builds transparency by assessing LMs in their totality. Rather than focusing on a specific aspect, the aim is to strive for a fuller characterization of LMs to improve scientific understanding and increase societal impact.

A benchmark of LM has two levels:

- (i) an abstract taxonomy of scenarios and metrics to define the design space for LM evaluation and
- (ii) a concrete set of implemented scenarios and metrics that were selected to prioritize coverage (e.g. different English varieties), value (e.g. user-facing applications), and feasibility (e.g. limited engineering resources).

When doing a benchmark some key considerations should be taken into account. To begin with, while standardizing a model evaluation, in particular by evaluating all models for the same scenarios, same metrics, models themselves may be more suitable for particular scenarios, particular metrics, and particular prompts/adaptation methods.

Moreover, while the evaluation itself may be standardized, the computational resources required to train these models may be very different (e.g. resource-intensive models generally fare better in our evaluation).

Furthermore, models may also differ significantly in their exposure to the particular data distribution or evaluation instances in use, with the potential for train-test contamination.

Even for the same scenario, the adaptation method that maximizes accuracy can differ across models which poses a fundamental challenge for what it means to standardize LM evaluation in a fair way across models.

Given the myriad scenarios where LMs could provide value, it would be appealing for many reasons if upstream perplexity on LM objectives reliably predicted downstream accuracy. Unfortunately, when making these comparisons across model families, even when using bits-per-byte (BPB)- which could provide more comparison than perplexity-, this type of prediction might not always work well.

Overview of LMs

LMs are a sub-category of NLP (neural language programming) within the field of AI. As in any other field of AI, the challenge of transparency in AI models and datasets continues to receive increasing attention from academia and industry.

When it comes to developing an AI model, *producers* are upstream creators of dataset and documentation, responsible for dataset collection, ownership, launch and maintenance.

Agents are stakeholders who read transparency reports, and possess the agency to use or determine how themselves or others might use the described datasets or AI systems.

Agents are distinct from *Users*, who are individuals and representatives who interact with products that rely on models trained on dataset. Users may consent to providing their data as a part of the product experience, and require a significantly different set of explanations and controls grounded within product experiences.

Dataset design also plays a crucial role for LM development. All data is processed in a common markdown format to blend knowledge between sources. For the interface, one can use task-specific tokens to support different types of knowledge. Uncurated data also means more tokens with limited transfer value for the target use-case; wasting compute budget.

Transparency refers to a *clear, easily understandable, and plain language explanation of what something is, what it does and why it does that*. The following table (Table 1.0) includes core aspects of transparency.



Transparency Characteristic	Description
Balance opposites	For example, disclosing information about AI systems without leaving creators vulnerable beyond reason, reporting fairness analyses without legitimizing inequitable or unfair systems, introducing standards for transparency that are wholly automated or become checklists.
Increase in expectations	Any information included in a transparency artifact can be expected to receive greater scrutiny.
Constant availability	Users want access to transparency information at multiple levels, even if they don't need to use it.
Require checks and balances	Transparency artifacts and their creation must be amenable to 3rd party evaluation, with the caveat that excessive transparency can open an AI system vulnerable to adversarial actors.
Subjective interpretations	Stakeholders have different definitions and unique ideas on what constitutes transparency.
Trust enabler	Accessible and relevant information about AI systems in-creases the the willingness of a data consumer or user to take a risk based on the expectation of benefits from the data, algorithms and the products they use.
Reduce knowledge asymmetries	Cross-disciplinary stakeholders are more effective when they possess a shared mental model and vocabulary to describe aspects of the AI system.
Reflects human values	It comes from both technical and non-technical disclosure about assumptions, facts and alternatives.

Table. 1.0 Core traits of transparency

Yet, attempts to introduce standardized and sustainable mechanisms for transparency is hindered by real world constraints of the diversity of goals, workflows, and backgrounds of individual stakeholders participating in the life cycles of datasets and AI systems.

In order to increase the transparency of NLP s, it might be useful to gain an understanding of the different tasks that they accomplish.

To begin with, Question answering (QA) is a fundamental task in NLP that underpins many real-world applications including web search, chatbots, and personal assistants. QA is very broad in terms of the questions that can be asked and the skills that are required to arrive at the answer, covering general language understanding, integration of knowledge, and reasoning (Gardner et al., 2019; Rogers et al., 2021).

Information retrieval (IR) refers to the class of tasks concerned with searching large unstructured collections (often text collections), is central to numerous user-facing applications. IR has a long tradition of study (Salton and Lesk, 1965; Salton, 1971; Spärck Jones, 1972; Salton and McGill, 1983; Manning et al., 2008; Lin et al., 2021a) and is one of the most widely deployed language technologies.

Text summarization is an established research direction in NLP (Luhn, 1958; Mani, 1999; Spärck Jones, 1999; Nenkova and McKeown, 2012), with growing practical importance given the ever-increasing volume of text that would benefit from summarization.

One can formulate text summarization as an unstructured sequence-to-sequence problem, where a document (e.g. a CNN news article) is the input and the LM is tasked with generating a summary that resembles the reference summary (e.g. the bullet point summary provided by CNN with their article).

To evaluate model performance, the model-generated summary is compared against a human-authored reference summary using automated metrics for overall quality (Lin, 2004; Zhang et al., 2020b), faithfulness (Laban et al., 2022; Fabbri et al., 2022), and extractiveness (Grusky et al., 2018). Extractiveness refers to the extent to which model summaries involve copying from the input document.

Consequently, it is important to measure and improve the faithfulness of these systems since unfaithful systems may be harmful by potentially spreading misinformation, including dangerous, yet hard to detect errors, when deployed in real-world settings.



Sentiment analysis has blossomed into its own subarea in the field with many works broadening and deepening the study of sentiment from its initial binary text-classification framing (Wiebe et al., 2005; McAuley et al., 2012; Socher et al., 2013; Nakov et al., 2016; Potts et al., 2021).

Text classification has a long history in NLP (see Yang and Pedersen, 1997; Yang, 1999; Joachims, 1998; Aggarwal and Zhai, 2012) with tasks such as language identification, sentiment analysis, topic classification, and toxicity detection being some of the most prominent tasks within this family.

Focusing on fairness of models is essential to ensuring technology plays a positive role in social change (Friedman and Nissenbaum, 1996; Abebe et al., 2020; Bommasani et al., 2021). Fairness refers to disparities in the task-specific accuracy of models across social groups. One way to operationalize fairness is by means of counterfactual fairness (Dwork et al., 2012; Kusner et al., 2017) which refers to model behavior on counterfactual data that is generated by perturbing existing test examples (cf. Ma et al., 2021; Qian et al., 2022).

In contrast, bias refers to properties of model generations, i.e. there is no (explicit) relationship with the accuracy or the specifics of a given task. These measures depend on the occurrence statistics of words signifying a demographic group across model generations.

Toxicity detection (and the related tasks of hate speech and abusive language detection) is the task of identifying when input data contains toxic content, which originated due to the need for content moderation on the Internet (Schmidt and Wiegand, 2017; Rauh et al., 2022).

Critiques of the task have noted that (i) the study of toxicity is overly reductive and divorced from use cases (Diaz et al., 2022), (ii) standard datasets often lack sufficient context to make reliable judgments (Pavlopoulos et al., 2020; Hovy and Yang, 2021), and (iii) the construct of toxicity depends on the annotator (Sap et al., 2019a; Gordon et al., 2022).

Another crucial concept for ML models is toxicity used as an umbrella term for related concepts like hate speech, violent speech, and abusive language (see Talat et al., 2017).⁴⁸ To operationalize toxicity measurement, one can use the Perspective API (Lees et al., 2022)⁴⁹ to detect toxic content in model generations.

Given these features of LM the next section explores a conceptual framework for designing a LM benchmark.

Conceptual Model

The study suggests to implement the following aspects for designing a LM benchmark (Figure 2.):

- (1). **Taxonomy:** One can taxonomize the vast design space of language model evaluation into scenarios and metrics. By stating this taxonomy, one can select systematically from this space, which makes explicit both priorities in benchmark design and the limitations in the benchmark at present.
- (2). **Broad coverage:** Given the taxonomy, one select and implement core scenarios, for which one can comprehensively measure major metrics (accuracy, calibration, robustness, fairness, bias, toxicity, efficiency).
- (3). **Evaluation of existing models:** One can evaluate existing Lms under the standardized conditions of the benchmark, ensuring models can now be directly compared across many scenarios and metrics. These models might vary in terms of their public accessibility: while some of them are open, others are limited-access, and a few might even be closed.
- (4). **Empirical findings:** The extensive evaluation will offer guidance for future language model development and ample opportunities for further analysis.

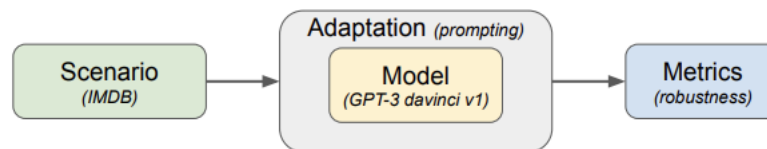


Fig. 2.0 Suggested LM Benchmark Process

As seen in Figure 2, the following aspects (scenario, adaptation, metric) are required to evaluate a LM to provide a roadmap for how to evaluate language models:

- **Scenarios:** A scenario instantiates a desired use case for a LM. Scenarios are what we want models to do. Each instance consists of (i) an input (a string) and (ii) a list of references. Each reference is a string annotated with properties relevant for evaluation (e.g. is it correct or acceptable?).



- **Adaptation:** Adaptation is the procedure that transforms a LM, along with training instances, into a system that can make predictions on new instances. Examples of adaptation procedures include prompting, lightweight-finetuning, and finetuning.

We define a language model to be a black box that takes as input a prompt (string), along with decoding parameters (e.g. temperature). The model outputs a completion (string), along with log probabilities of the prompt and completion. Viewing language models as text-to-text abstractions is important for two reasons:

1. First, while the prototypical LM is usually a dense Transformer trained on raw text, LMs could also use an external document store (Lewis et al., 2020c), issue search queries on the web (Nakano et al., 2021), or be trained on human preferences (Ouyang et al., 2022; Bai et al., 2022). An ideal model should be agnostic with regard to these implementation details.
 2. Second, the text-to-text abstraction is a convenient general interface that can capture all the (text-only) tasks of interest, an idea that was pioneered by McCann et al. (2018) and Raffel et al. (2019).
- **Metrics:** To determine how well the model performs, one can compute metrics over these completions and probabilities. Metrics concretely operationalize the abstract desiderata required for useful systems.

To evaluate a LM, a series of runs must be implemented, where each run is defined by a scenario, adaptation method and metric. Each of these scenarios, adaptation, and metrics define a complicated and structured space, which one implicitly navigates to make decisions in evaluating a LM.

Scenarios

One can taxonomize scenarios based on the following:

- (i) a task (e.g. question answering, summarization), which characterizes what we want a system to do;
- (ii) a domain (e.g. a Wikipedia 2018 dump), which characterizes the type of data we want the system to do well on; and
- (iii) the language or language variety (e.g. English).

Tasks, domains, and languages are not atomic or unambiguous constructs: they can be made coarser and finer. Given this structure, one can deliberately select scenarios based on main overarching principles:

- (i) coverage of the space,
- (ii) minimality of the set of selected scenarios, and
- (iii) prioritizing scenarios that correspond to user-facing tasks.

Given the ubiquity of natural language, the field of natural language processing (NLP) considers myriad tasks that correspond to language's many functions (Jurafsky and Martin, 2000). To generate this set, one can take the tracks at a major NLP conference (ACL 2022), and for each track, one can map the associated subarea of NLP to canonical tasks for that track.

Moreover, domains are a familiar construct in NLP, yet their imprecision complicates systematic coverage of domains. One can further decompose domains according to 3 W's:

- (1). **What (genre):** The type of text, which captures subject and register differences. Examples: Wikipedia, social media, news, scientific papers, fiction.
- (2). **When (time period):** When the text was created.
Examples: 1980s, pre-Internet, present day (e.g. does it cover very recent data?)
- (3). **Who (demographic group):** Who generated the data or who the data is about. Examples: Black/White, men/women, children/elderly.

Models

When deployed in practice, models are confronted with the complexities of the open world (e.g. typos) that cause most current systems to significantly degrade (Szegedy et al., 2014; Goodfellow et al., 2015; Jia and Liang, 2017; Belinkov and Bisk, 2018; Madry et al., 2018; Ribeiro et al., 2020; Santurkar et al., 2020; Tsipras, 2021; Dhole et al., 2021; Koh et al., 2021; Yang et al., 2022).

One suggestion is to measure the robustness of different models by evaluating them on transformations of an instance. That is, given a set of transformations for a given instance, one can measure the worst-case performance of a model across these transformations.



On the one hand, measuring robustness to distribution or subpopulation shift (Oren et al., 2019; Santurkar et al., 2020; Goel et al., 2020; Koh et al., 2021) requires scenarios with special structure (i.e., explicit domain/subpopulation annotations) as well as information about the training data of the models.

On the other hand, measuring adversarial robustness (Biggio et al., 2013; Szegedy et al., 2014) requires many adaptive queries to the model in order to approximate worst-case perturbations, which might not always be feasible (Wallace et al., 2019a; Morris et al., 2020).

Moreover, the transformation/perturbation-based paradigm has been widely explored to study model robustness (e.g. Ribeiro et al., 2020; Goel et al., 2021; Wang et al., 2021a), in order to understand whether corruptions that arise in real use-cases (e.g. typos) affect the performance of the model significantly. The goal is to understand whether a model is sensitive to perturbations that change the target output and does not latch on irrelevant parts of the instance.

Metrics

To taxonomize the space of desiderata, one can begin by enumerating criteria that are necessary for developing useful systems. Yet, what does it mean for a system to be useful?

Too often in AI, this has come to mean the system should be accurate in an average sense. While (average) accuracy is an important, and often necessary, property for a system (Raji et al., 2022), accuracy is often not sufficient for a system to be useful/desirable.

Unfortunately, while many of the desiderata are well-studied by the NLP community, some are not codified in specific tracks/areas (e.g. uncertainty and calibration). Therefore, it is suggested to expand the scope to all AI conferences, drawing from a list of AI conference deadlines.

Recommendations

While reasoning is usually assumed to involve transitions in thought (Harman, 2013), possibly in some non-linguistic format, one typical way to assess reasoning abilities (e.g., in adult humans) is by means of explicit symbolic or linguistic tasks.

In order to distinguish reasoning from language and knowledge as much as possible, one can focus on relatively abstract capacities necessary for sophisticated text-based or symbolic reasoning.

To measure ampliative reasoning, one can use explicit rule induction and implicit function regression, which corresponds to making and applying claims about the likely causal structure for observations.

For rule induction, one can design and implement `rule_induct` inspired by the LIME induction tasks, where we provide two examples generated from the same rule string, and task the model with inferring the underlying rule.

For function regression, one can design and implement `numeracy_prediction`, which requires the model to perform symbolic regression given a few examples and apply the number relationship (e.g. linear) to a new input.

One can also evaluate language models on more complex and realistic reasoning tasks that require multiple primitive reasoning skills to bridge the gap between understanding reasoning in very controlled and synthetic conditions and the type of reasoning required in practical contexts.

Conclusions

Focusing on evaluation of models is essential to ensure that technology plays a positive role in social change. Within this regard, this study explored the main features of benchmarks - scenario, adaptation, metric-required to provide a roadmap for how to evaluate LMs. It also made recommendations of how to use model and metrics for fairness and transparency when it comes to developing LM.

This study conceptualized a benchmark model for evaluating NLP models. Given the lack of studies in the field it is a step towards the design of more sophisticated models and thus, right now, far from perfect. Nevertheless, it aims to raise awareness of the importance of developing benchmarks for AI models.

References

- [1]. Allison Koenecke, Andrew Nam, Emily Lake, Joe Nudell, Minnie Quartey, Zion Mengesha, Connor Toups, John R Rickford, Dan Jurafsky, and Sharad Goel. 2020.
- [2]. Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres (2019) Quantifying the carbon emissions of machine learning. arXiv preprint arXiv:1910.09700.
- [3]. Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, and Richard Socher (2016) Ask me anything: Dynamic memory networks for natural language processing. In International conference on machine learning, pages 1378–1387. PMLR.



- [4]. Atoosa Kasirzadeh and Iason Gabriel (2022). In conversation with artificial intelligence: aligning language models with human values.
- [5]. Bernard Koch, Emily Denton, Alex Hanna, and Jacob Gates Foster (2021) Reduced, reused and recycled: The life of a dataset in machine learning research. In Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2).
- [6]. Christopher Potts, and Matei Zaharia (2021) A moderate proposal for radically better AI-powered Web search. Stanford HAI Blog.
- [7]. Daniel Khashabi, Yeganeh Kordi, and Hannaneh Hajishirzi (2022) Unifiedqa-v2: Stronger generalization via broader cross-format training. ArXiv, abs/2202.12359. Omar Khattab,
- [8]. Davis, D., Seaton, D., Hauff, C., &Houben, G. J. (2018, June). Toward large-scale learning design: Categorizing course designs in service of supporting learning outcomes. *Proceedings of the Fifth Annual ACM Conference on Learning at Scale* (pp. 1-10). <https://doi.org/10.1145/3231644.3231663>
- [9]. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). *Bert: Pre-training of deep bidirectional transformers for language understanding*. arXiv preprint arXiv:1810.04805. <https://arxiv.org/abs/1810.04805>
- [10]. Divyansh Kaushik, Eduard Hovy, and Zachary Lipton (2019) Learning the difference that makes a difference with counterfactually-augmented data. In International Conference on Learning Representations (ICLR).
- [11]. Fereshte Khani and Percy Liang. (2020) Feature noise induces loss discrepancy across groups. In International Conference on Machine Learning (ICML).
- [12]. González-Carvajal, S., & Garrido-Merchán, E. C. (2020). *Comparing BERT against traditional machine learning text classification*. arXiv preprint arXiv:2005.13012. <https://arxiv.org/abs/2005.13012>
- [13]. Grandini, M., Bagli, E., &Visani, G. (2020). *Metrics for multi-class classification: An overview*. arXiv preprint arXiv:2008.05756. <https://arxiv.org/abs/2008.05756>
- [14]. Hannah Rose Kirk, Abeba Birhane, Bertie Vidgen, and Leon Derczynski (2022) Handling and presenting harmful text in nlp research. Andy Kirkpatrick. 2020. The Routledge handbook of world Englishes. Routledge.
- [15]. Jared Kaplan, Sam McCandlish, T. J. Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeff Wu, and Dario Amodei. 2020. Scaling laws for neural language models. ArXiv, abs/2001.08361.
- [16]. Kingma, D. P., & Ba, J. (2014). *Adam: A method for stochastic optimization*. arXiv preprint arXiv:1412.6980. <https://arxiv.org/abs/1412.6980>
- [17]. Matt J Kusner, Joshua R Loftus, Chris Russell, and Ricardo Silva (2017) Counterfactual fairness. In Advances in Neural Information Processing Systems (NeurIPS), pages 4069–4079.
- [18]. Philippe Laban, Tobias Schnabel, Paul N. Bennett, and Marti A. Hearst (2022) SummaC: Re-Visiting NLI-based Models for Inconsistency Detection in Summarization. Transactions of the Association for Computational Linguistics, 10:163–177.
- [19]. Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. In Association for Computational Linguistics (ACL).
- [20]. Tomš Kočický, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gabor Melis, and Edward Grefenstette. 2017. The Narrative QA reading comprehension challenge. arXiv preprint arXiv:1712.07040.
- [21]. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, i. (2017). Attention is all you need. *31st Conference on Neural Information Processing Systems (NiPS 2017; pp. 5998-6008)*. Long Beach, USA. <https://arxiv.org/abs/1706.03762>