
NeurIPS 2025 E2LM Competition : Early Training Evaluation of Language Models

Mouadh Yagoubi¹ Yasser Dahou¹ Billel Mokeddem¹ Younes Belkada¹
 Phuc H. Le-Khac¹ Basma El Amel Boussaha¹ Reda Alami¹ Jingwei Zuo¹
 Damiano Marsili² Mugariya Farooq¹ Mounia Lalmas³
 Georgia Gkioxari² Patrick Gallinari⁴ Philip Torr⁵ Hakim Hacid¹

Abstract

Existing benchmarks have proven effective for assessing the performance of fully trained large language models. However, we find striking differences in the early training stages of small models, where benchmarks often fail to provide meaningful or discriminative signals. To explore how these differences arise, this competition tackles the challenge of designing scientific knowledge evaluation tasks specifically tailored for measuring early training progress of language models. Participants are invited to develop novel evaluation methodologies or adapt existing benchmarks to better capture performance differences among language models. To support this effort, we provide three pre-trained small models (0.5B, 1B, and 3B parameters), along with intermediate checkpoints sampled during training up to 200B tokens. All experiments and development work can be run on widely available free cloud-based GPU platforms, making participation accessible to researchers with limited computational resources. Submissions will be evaluated based on three criteria: the quality of the performance signal they produce, the consistency of model rankings at 1 trillion tokens of training, and their relevance to the scientific knowledge domain. By promoting the design of tailored evaluation strategies for early training, this competition aims to attract a broad range of participants from various disciplines, including those who may not be machine learning experts or have access to dedicated GPU resources. Ultimately, this initiative seeks to make foundational LLM research more systematic and benchmark-informed from the earliest phases of model development.

Keywords

Early training analysis, Evaluation benchmarks, Scientific knowledge, Low-resource ML research, Language Models.

¹Technology Innovation Institute (TII), UAE

²California Institute of Technology, Pasadena, CA, USA

³Spotify, UK

⁴ISIR - Sorbonne University, France

⁵University of Oxford, UK

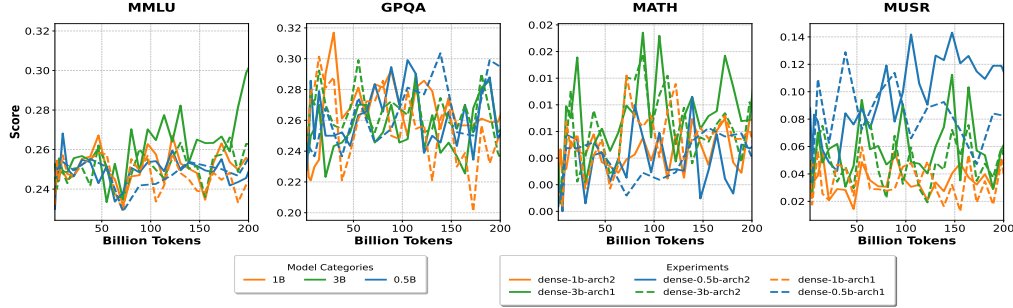


Figure 1: Noisy results obtained with state-of-the-art benchmarks with different models sizes.

1 Competition description

1.1 Background and impact

Large Language Models (LLMs) have advanced at an unprecedented pace in recent years [1, 2, 3, 4]. In response, the research community has increasingly focused on developing benchmarking frameworks to assess their capabilities in a variety of domains, including language, mathematics, and reasoning. Early evaluations primarily focused on measuring the language modeling skills [5, 6], a fundamental competence that appears to be increasingly well-acquired by newer LLM variants saturating such benchmarks. Recently, attention has shifted towards evaluating scientific knowledge, code generation and reasoning capabilities, using benchmarks such as MMLU [7] and its enhanced version MMLU-pro [8], GPQA [9], GSM8K [10], MATH [11], and LiveCodeBench [12]. A common theme across these benchmarks is the increasing complexity of tasks, aimed at testing the limits of frontier models.

Despite these advances, Small Language Models (SLMs) with lower than 7B parameters consistently underperform relative to their larger counterparts on these challenging benchmarks. This performance gap is particularly problematic when training SLMs from scratch, as these benchmarks fail to offer discriminative signals during early training phases (typically between 0-200 billion tokens). Consequently, they provide limited guidance for training dynamics regarding: hyperparameter search, data mixture and architectural sweeps. Developing more granular and targeted evaluation benchmarks in the scientific knowledge domain would enable more systematic investigation of fundamental research questions during the early training stages of SLMs.

Concretely, we provide three pre-trained SLMs (0.5B, 1B, and 3B) along with intermediate checkpoints up to 200 billion tokens. As shown in Figure 1, current benchmarks are often noisy and fail to provide a clear signal to measure progress, raising several fundamental questions: Are SLMs not learning any scientific knowledge early in the training? Or are the benchmark questions too difficult at this stage? Is the noise due to a lack of acquired knowledge, or could it be related to how the questions are formatted from a linguistic standpoint? Are there ways to refine current benchmarks to better extract meaningful signals?

The primary goal of this competition is to investigate these questions by inviting participants to design and develop novel evaluation tasks that enable a progressive assessment of scientific knowledge in SLMs, hence, enabling a progressive assessment of scientific knowledge in SLMs and producing more informative performance signals. Specifically, the expected impact of the challenge can be summarized as follows:

- **Enhancing SLMs research:** By enabling more rigorous and conclusive studies on smaller models, this competition will open opportunities to develop stronger SLMs in scientific knowledge, code, and mathematics—ultimately helping to bridge the gap with larger models through meaningful analysis across a large number of small-scale experiments.
- **Transfer to large models:** Lessons learned from studying SLMs can help inform better decisions when training larger models, particularly in terms of data mixtures, hyper-parameters, and architec-

tural choices. We refer the reader to [13, 14] for a comprehensive overview of the role of small models in advancing LLMs.

- **Understanding the learning dynamics:** Current benchmarks often do not fully reflect the convergence or complexity of the model in a meaningful way. Benchmark difficulty should ideally evolve with a model’s level of training. Tackling this problem may offer insights into how models learn over time.
- **Expand the reach of LLM research:** We aim to encourage participation from researchers beyond the core machine learning community, inviting them to bring evaluation ideas and benchmarks rooted in their domain expertise. Such interdisciplinary collaboration has the potential to broaden the impact and relevance of LLM research.

To this end, participants can either adapt existing benchmarks or design entirely new tasks from scratch. Moreover, we provide a baseline approach that has been successfully tested on the provided checkpoints (see Section 1.4). This baseline draws inspiration from the work of [15], who proposed a variant of the well-known MMLU benchmark [7], referred to as MMLU-var. In this variant, prompts are reformulated in a completion/close-style format rather than as multiple-choice questions. The motivation for this adaptation stems from the observation that LLMs typically gain proficiency in multiple-choice questions only in later training stages (e.g., after 1 trillion tokens), which already represents a substantial computational budget.

Furthermore, the results on MMLU-var (Figure 2) demonstrate how simple prompt modifications (see Appendix B) can substantially improve the evaluation of SLMs. Grade-school-level benchmarks like ARC-easy [16] and SciQ [17] also provide useful signals (refer to Appendix A). That said, benchmark complexity is inherently multi-dimensional; while MMLU-var addresses linguistic complexity in prompts, participants are encouraged to explore additional dimensions such as semantic depth and domain-specific reasoning in areas like mathematics or code generation.

1.2 Novelty

Research on LLMs has been rapidly evolving across multiple fronts, including how to evaluate them effectively [18], how to improve their efficiency [19], their training scaling laws [20, 21], and how to combine or adapt them for diverse use cases [22]. In recent NeurIPS editions, several competitions have been introduced to address the high costs associated with training and deploying them: the LLM Merging Challenge explored how to combine existing models to create more powerful ones without additional training [23], while the Edge LLM Challenge [24] focused on developing efficient, optimized models capable of running on resource-constrained edge devices. A Single-GPU Fine-Tuning and Inference Challenge was also organized to promote the development of methods that reduce the computational costs of using LLMs in practice [25].

While these challenges target fully trained model’s applicability, our proposal shifts the focus towards SLMs and explores how best to evaluate their acquisition of scientific knowledge during the early phases of training. Drawing from foundational work in cognitive development and education [26, 27], we are inspired by the idea that evaluating learning requires sensitivity to both developmental stage and type of understanding. For example, [26] emphasized that learners acquire increasingly abstract forms of reasoning over time, suggesting that evaluation tasks should be developmentally aligned—not just more difficult, but also structured differently at each stage. Similarly, taxonomies in [27, 28] propose that comprehension, analysis, and knowledge application emerge progressively.

By analogy, we argue that evaluation strategies for SLMs should not focus solely on final performance, but rather take into account the learning dynamics and training trajectory. Just as students are not assessed on differential equations before mastering algebra, language models should be evaluated on tasks aligned with their stage of development. We argue that early-stage benchmarks should probe intermediate reasoning capabilities and the gradual acquisition of scientific knowledge in a way that reflects how understanding evolves with exposure and training.

This competition aims to provide the research community with new evaluation benchmarks that reflect the learning dynamics of Small Language Models. It will equip participants with tailored tools to assess SLMs, with a particular emphasis on reasoning and scientific knowledge domains.

By establishing dedicated reasoning- and scientific-knowledge-focused evaluation benchmarks for SLMs, we aim not only to improve the evaluation of SLMs but also to support broader LLM research. Fine-grained evaluation at a small scale enables more informative experimentation, potentially

improving decisions in large-scale trainings. Ultimately, we believe that addressing the challenge of evaluating SLMs in a competition setting will foster innovation in benchmark design and inform the development of metrics better suited to the evolving needs of the research and industry communities.

1.3 Data

To support participants and facilitate the validation of their benchmarks, we trained multiple models using two distinct data mixtures. Below, we provide a detailed description of each:

- **Web-only Dataset:** This consists of a random subset of FineWeb [29], a cleaned and deduplicated English web corpus derived from CommonCrawl⁶. Models trained on this mixture typically acquire broad linguistic fluency and general world knowledge.
- **Scientific Knowledge Data Mixture:** This mixture comprises selected subsets from the following publicly available datasets: FineWeb-edu [30] (50%), The Stack V1 [31] and The Stack V2 [32] (21.6%), InfIMM [33] (18.9%), and TxT360⁷ [34] (9.5%). Compared to the Web-only mixture, models trained on this subset are expected to demonstrate stronger performance in logical inference, mathematical problem-solving, and code generation.

Each model was trained on 1 trillion tokens using a custom tokenizer with a vocabulary size of 65,000 tokens. To maintain a clear scope for the competition and ensure a fair comparison across submissions, only English-language data was used. Additionally, we provide models at three scales—each instantiated with two distinct architectural variants for the same size: a deep variant (arch1) and a wide variant (arch2), assuming deeper models reason better [35], as shown in Appendix C. The choice is merely to differentiate models of the same size rather than to compare width vs depth, which is out of the scope of this work. Participants are provided with intermediate checkpoints up to 200 billion tokens, enabling evaluation at various stages of training.

1.4 Tasks and application scenarios

We begin our investigation using the baseline results obtained on the MMLU-var benchmark. Figure 2 presents a comparison of various model architectures and sizes using both the standard MMLU and its modified variant, MMLU-var [15].

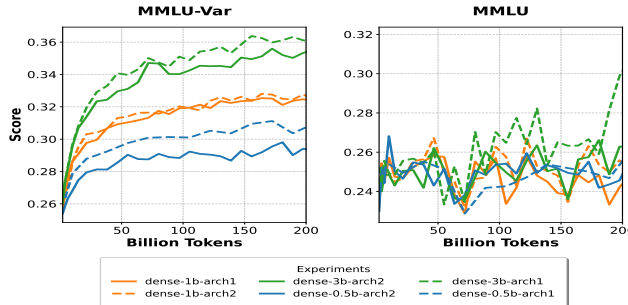


Figure 2: **MMLU-var (left) vs MMLU (right) comparison on different configurations and sizes of models.** Different colors for each size of the model (0.5B, 1B and 3B) and markers distinguish experiments within the same category. MMLU-var allows a clear comparison between variants while MMLU is giving noise (expect for 3B arch2).

We observe that results on MMLU do not consistently differentiate between models of varying sizes and architectures (i.e. 0.5B, 1B, and 3B), with no apparent scaling law or pattern, except at later stages of training (around 180 billion Tokens), where the *3B-Arch-2* model clearly outperforms the others.

In contrast, the MMLU-var (left plot) offers better informative results throughout the training process. On top of ranking models according to size, it also distinguishes between different architectural variants within the same size. These results suggest that early trained SLMs struggle to explicitly compare and reason over multiple-choices in the input space, but still gives a higher log-likelihood to the correct answer when prompted in a completion manner. Interestingly, we found that tasks requiring general language capabilities (i.e., HellaSwag [5], WinoGrande[6]) provide meaningful

⁶<https://commoncrawl.org>

⁷All subsets are included except for Common Crawl.

signals during early training (0 to 200 BT). Furthermore, a core objective of this competition is to identify and explore similar ideas that reliably evaluate scientific knowledge.

The effect of data mixtures: we trained a 1B model on a different datasets, one using only web data, and the other using a curated mix of knowledge-rich domains (as described in Section 1.3). As shown in Figure 3, both models perform similarly on HellaSwag. However, on MMLU-var, the model trained on the knowledge-enriched dataset significantly outperforms the web-only variant. This demonstrates that while HellaSwag is effective for evaluating general language understanding in SLMs, it does not adequately capture reasoning or knowledge capabilities. In contrast, MMLU-var clearly differentiates between the two, with the knowledge-trained model achieving higher scores. This eliminates the hypothesis that acquiring scientific knowledge is a fundamental limitation in SLMs’ early training and suggests that factors such as data, architecture, or hyperparameters could potentially enhance the model’s capacity at this stage.

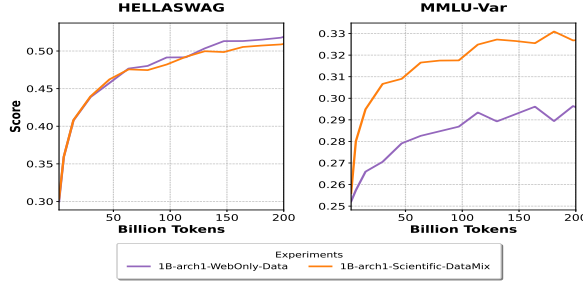


Figure 3: Comparison of results from a 0.5B parameter model trained on two datasets — Web-Only and Knowledge Data-Mixture. While both model configurations achieve similar performance on the HellaSwag benchmark, the model trained on the Knowledge Data-Mixture significantly outperforms the Web-Only model on MMLU-var.

1.5 Metrics

We propose to evaluate submitted solutions based on three main criteria, which will be combined into a global score used for the final ranking: Signal Quality (SQ), Ranking Consistency (RC), and Compliance with Scientific Knowledge Domains (CS). Additionally, two validation procedures will be systematically applied to all submissions: (i) verification of alignment with established scientific knowledge domains (see Appendix F.1), and (ii) detection of potential information leakage, specifically the presence of the answer within the question prompt (see Appendix F.2). The overall score is computed as a weighted sum:

$$\text{Score} = \alpha_1 \times \text{Score}_{\text{SQ}} + \alpha_2 \times \text{Score}_{\text{RC}} + \alpha_3 \times \text{Score}_{\text{CS}} \quad (1)$$

Here, α_{SQ} , α_{RC} and α_{CS} are weighting coefficients that reflect the relative importance of each criterion. We set the initial weights as $\alpha_1 = 0.5$, $\alpha_2 = 0.1$ and $\alpha_3 = 0.4$, thereby placing greater emphasis on signal quality and compliance to scientific knowledge, which we consider the most important metrics in evaluating submissions.

Participants will be able to compute the signal quality subscore (Score_{SQ}) locally using the provided model checkpoints (ranging from 0 to 200 BT) along with the accompanying scoring algorithm (provided in a notebook the starting kit). In contrast, the other two subscores cannot be computed independently, as the corresponding checkpoints—from 200 billion tokens to 1 trillion tokens, as well as the 0.5 billion parameter model trained exclusively on web data—will remain hidden throughout the competition. The global score will, however, be automatically computed upon submission via the Codabench platform, enabling participants to monitor their overall performance. This setup is designed to prevent overly customized solutions tailored specifically to the released checkpoints.

Signal Quality metric (Score_{SQ}) This score rewards evaluation tasks that produce smooth and informative learning curves throughout training. Let $X_{200\text{BT}} = \{(i_j, x_j) \mid j = 1, 2, \dots, n\}$ denote the set of training iterations i_j and corresponding benchmark scores x_j measured up to 200 billion tokens. Score_{SQ} is composed of two subcomponents:

- **Monotonicity Score:** This component uses Spearman’s rank correlation to measure the degree of monotonic improvement over time. Given rank differences d_j between iteration indices and their

associated scores, the score is computed as:

$$\text{Score}_{\text{Monotonicity}} = \max \left(0, 1 - \frac{6 \sum d_j^2}{n(n^2 - 1)} \right) \quad (2)$$

Negative trends yield a score of 0. This formulation ensures that learning curves showing consistent progress receive higher scores.

- **Autocorrelation Strength:** This captures temporal coherence, rewarding signals that are stable over time rather than noisy. For each lag $\ell \in \{1, 2, \dots, L\}$ with $L = \lfloor n/4 \rfloor$, we compute the Pearson correlation coefficient between the original score sequence and its lagged version (with $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$):

$$\rho_\ell = \frac{\sum_{i=1}^{n-\ell} (x_i - \bar{x})(x_{i+\ell} - \bar{x})}{\sqrt{\sum_{i=1}^{n-\ell} (x_i - \bar{x})^2} \times \sqrt{\sum_{i=1}^{n-\ell} (x_{i+\ell} - \bar{x})^2}} \quad (3)$$

The autocorrelation score is the average of the absolute correlations across all lags:

$$\text{Score}_{\text{AutoCorr}} = \frac{1}{L} \sum_{\ell=1}^L |\rho_\ell| \quad (4)$$

The final signal quality score (SQ) is a weighted combination of the two components:

$$\text{Score}_{\text{SQ}} = \beta_1 \times \text{Score}_{\text{Monotonicity}} + \beta_2 \times \text{Score}_{\text{AutoCorr}} \quad (5)$$

This dual approach ensures that evaluation tasks reward both directional learning progress and temporal stability, minimizing noisy or erratic signals. Accordingly, we assign equal weight to the two sub-metrics by setting β_1 to 0.5 and β_2 to 0.5.

Ranking Consistency metric (Score_{RC}) This metric evaluates how well an evaluation task preserves a consistent ranking of models as training progresses, specifically after processing a large number of tokens (1 trillion tokens). The more stable and reproducible the model rankings over time, the higher the score. Ranking consistency is assessed separately for each model configuration (0.5B, 1B, and 3B), and the final score (Score_{RC}) is computed as the average across these configurations.

To quantify consistency, we use Kendall’s Tau coefficient [36], a widely adopted metric for measuring ordinal correlation between ranked lists. It has been successfully applied in prior work [37] to evaluate the robustness of model comparisons across training steps. The process proceeds as follows:

1. **Baseline Ranking (at 200 BT):** For each model size $s \in \{0.5\text{B}, 1\text{B}, 3\text{B}\}$, we compute the difference between the average performance across the checkpoints between 100 and 200 billion tokens of the two architectures $A \in \{\text{Arch1}, \text{Arch2}\}$, see the details about the different architectures in Table 2:

$$\text{rank}_{200\text{BT}}(s) = \begin{cases} 1 & \text{if } r(s) > 0 \\ 0 & \text{otherwise} \end{cases} \quad r(s) = \frac{1}{|K|} \sum_{j \in K} x_j^{\text{Arch1}}(s) - x_j^{\text{Arch2}}(s) \quad (6)$$

where $K = \{k \in \mathbb{N} \mid 100 \leq k < 200\}$, and x_j^{Arch1} and x_j^{Arch2} are the scores of models *Arch1* and *Arch2* at checkpoint j respectively. If *Arch1* is in average better than *Arch2*, $\text{rank}_{200\text{BT}}(s) = 1$, Otherwise $\text{rank}_{200\text{BT}}(s) = 0$.

2. **Ranking Consistency Evaluation (between 200 BT and 1TT):** Let $P = \{p \in \mathbb{N} \mid 220 \leq p \leq 1000\}$ representing the evaluation points post-200BT. At each point $p \in P$, we compare the current model ranking to the baseline. For each model size s , we define

$$\text{Score}_{\text{RC}} = \tau(s) = \frac{1}{|P|} \sum_{p \in P} \tau_p(s) \quad \tau_p(s) = \begin{cases} 1 & \text{if } \text{rank}_{200\text{BT}}(s) = \text{rank}_p(s) \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

where $\text{rank}_p(s) \in \{0, 1\}$, and $\text{rank}_p(s) = 1$ if $x_p^{\text{Arch1}}(s) > x_p^{\text{Arch2}}(s)$, otherwise $\text{rank}_p(s) = 0$.

A higher Score_{RC} reflects more stable rankings over time, indicating that the evaluation task provides a reliable signal even as training progresses.

Compliance to Reasoning and Knowledge Domains (Score_{CS}) This metric assesses how well the proposed evaluation task aligns with reasoning and scientific knowledge, rather than general language or commonsense capabilities. As illustrated in Figure 7, tasks designed to probe scientific knowledge (such as MMLU-var) clearly differentiate between models trained on a curated, knowledge-rich dataset and those trained solely on web data. In contrast, benchmarks like HellaSwag, which emphasize commonsense reasoning, tends to yield similar results across both models types, making them less effective for this competition’s objectives.

To quantify domain compliance, we compare the average performance of two 1B models across training steps: a model trained on the scientific knowledge-focused datamixture ($x^{\text{SciKW-DS}}$) and a model trained on web-only data ($x^{\text{Web-DS}}$). The score is define as the normalized performance gap:

$$\text{Score}_{\text{CS}} = \max \left(0, \frac{1}{n} \sum_{i=1}^n (x_i^{\text{SciKW-DS}} - x_i^{\text{Web-DS}}) \right) \quad (8)$$

This formulation ensures a fair comparison across evaluation tasks of varying difficulty levels, rewarding models that are sensitive to knowledge-intensive learning signals.

Metric Scores Analysis Table 1 presents the evaluation scores for the top-performing benchmarks. MMLU-var achieves the highest overall score, followed by ARC-Easy and SciQ, both of which perform well on all metrics. Our analysis suggests that ARC-Easy is particularly effective for assessing scientific knowledge in the early training stages and can serve as a second baseline alongside MMLU-var in this competition. However, in 82% of SciQ questions, the correct answer appears verbatim in the prompt, likely making the task easier for models. To mitigate this, we will apply a leakage check to all submissions and exclude such questions from benchmark evaluations, as detailed in Appendix F.2. Benchmarks targeting scientific knowledge consistently outperform HellaSwag, despite the latter scoring higher on the Consistency and Signal Quality metrics. HellaSwag receives a zero on the Compliance metric, which evaluates whether the benchmark is classified as scientific knowledge (central to this competition). To validate our scoring metric, we ranked all tested benchmarks (see Appendix E and G for full results). Although commonsense reasoning tasks such as HellaSwag, Winnogrande and other related benchmarks are included in the full score table, similar submitted benchmarks will not appear on the official leaderboard, as they fail the scientific compliance pre-check. Their scores are reported solely to assess the robustness of our evaluation metric.

Benchmark	Score _{SQ}	Score _{RC}	Score _{CS}	Total Score (%)
MMLU-var	0.959	0.837	0.384	71.7
ARC-Easy	0.832	0.822	0.394	65.5
SciQ	0.846	0.772	0.316	62.7
HellaSwag	0.992	0.936	0.000	59.0
GSM8K	0.655	0.915	0.369	56.6

Table 1: Benchmark scores across different metrics, with Total Score calculated using Equation (1).

1.6 Baselines, code, and material provided

The provided starting kit ⁸ includes a set of Jupyter notebooks designed to help participants get started with the competition materials. These notebooks support the development of evaluation tasks by providing access to model checkpoints and allowing prompting and interaction with the provided LLMs. A complete list of available notebooks is included in Appendix D. Since the largest model used in the competition has 3 billion parameters, approximately 6GB of GPU memory is sufficient for inference. As a result, the experiments and development work can be run on a free-tier Google Colab GPU (NVIDIA T4 16GB), ensuring broad accessibility for participants. In addition, to support participants without a strong background in machine learning, a tutorial session will be organized to explain how the lm-evaluation-harness framework works.

1.7 Website, tutorial and documentation

The competition website ⁹ will serve as the central hub for all key information, including: (i) competition overview, (ii) rules and terms, (iii) competition timeline and announcements, (iv) links to

⁸Competition starting kit

⁹<https://et-slm-evaluation.github.io/>

the starting kit netbooks, (v) tutorial section that will host the competition webinars, and (vi) contact information for the organizers along with a link to the discussion channel. Additionally, we will also host practical webinar sessions focused on best practices for evaluating Large Language Models.

2 Organizational aspects

The competition will be hosted on the Codabench platform [38], with the link available on the competition website. Participants must: (1) create an account, (2) download the starting kit to prepare their submission, and (3) upload their code to Codabench following the provided interface specifications. Submissions will be used to evaluate LLM models on the provided checkpoints to compute the scores. Finally, the global scores will be regularly displayed on the Codabench leaderboard, along with a detailed metrics page for each submission.

2.1 Rules and Engagement

- **Submissions requirement:** each submission must include a set of question–answer pairs and a clearly defined evaluation metric, implemented in Python code, to assess model performance.
- **Size of submitted solutions:** submissions must include 100–15,000 samples. The upper limit matches MMLU’s 14,042-sample test split—one of the largest benchmarks—supporting compatibility and discouraging unnecessarily large, combined benchmarks. The 100-sample minimum ensures accessibility for resource-constrained teams exploring focused ideas.
- **Competition phases:** this challenge is open to anyone and runs in 3 phases: during phases 1 & 2, participants can submit their code and view their results on the regularly updated leaderboard; specifically during the warmup phase, the organizers may adjust the global score formula.
- **Submission limits:** each participant or team may submit up to 10 entries per day.
- **Final ranking:** the final ranking will be based on the global score, calculated by the organizers and shared with all participants. The scoring algorithm will be made available so participants can evaluate their own solutions locally.
- **Team accounts:** teams must use a single group account. The use of multiple accounts is not permitted.
- **Prize eligibility:** To be eligible for a prize, a team must open-source its code at least two weeks before the NeurIPS competition workshop.

2.2 Schedule and readiness

The competition will consist of three phases to allow for a smooth participation/organization:

In the *Warm-up phase*, participants get familiarized with the provided materials and the proposed format to integrate their evaluation benchmarks. They may submit preliminary solutions and provide feedback to the organizers, which will help refine the competition setup for the next phase.

In the *Development phase*, participants will test and refine their evaluation tasks using the provided models and checkpoints. They may use their own resources to compute the signal quality score supported by the starting kit. To obtain a global score, participants must submit their benchmarks via the Codabench platform, where evaluations will be conducted using the competition’s resources.

In the *Final phase*, the organizers will validate the final rankings from the development phase. This includes verifying the submitted benchmarks and assessing the robustness of the top-performing solutions through multiple test runs. At the end of this phase, the organizers will also reveal the hidden checkpoints used during the competition to compute the final scores.

The proposed schedule is the following:

- **The competition will run from June 16 to October 17, 2025, comprising a 5-week warm-up phase, a 10-week development phase, and a 3-week final phase.**
- Announcement of results: 20 October, 2025.
- Fact sheets and code release by winners due: 22 November, 2025.
- Presentation of results at NeurIPS competition workshop: 6/7 December 2025.

2.3 Competition promotion and incentives

The challenge will be promoted through various channels, including announcements at leading conferences (NeurIPS, ICLR, ICML), relevant mailing lists, and social media platforms. Monetary prizes will be awarded to the top-performing teams: \$6,000 for 1st place, \$4,000 for 2nd place, and \$2,000 for 3rd place. A special award of \$2,000 will also be granted to each of the two best student solutions. A public leaderboard will be kept displaying the top evaluation tasks across a broad range of large language models (LLMs).

References

- [1] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. [arXiv preprint arXiv:2407.21783](#), 2024.
- [2] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. [arXiv preprint arXiv:2412.15115](#), 2024.
- [3] Quentin Malartic, Nilabhra Roy Chowdhury, Ruxandra Cojocaru, Mugariya Farooq, Giulia Campesan, Yasser Abdelaziz Dahou Djilali, Sanath Narayan, Ankit Singh, Maksim Velikanov, Basma El Amel Boussaha, et al. Falcon2-11b technical report. [arXiv preprint arXiv:2407.14885](#), 2024.
- [4] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. [arXiv preprint arXiv:2303.18223](#), 1(2), 2023.
- [5] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? [arXiv preprint arXiv:1905.07830](#), 2019.
- [6] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. [Communications of the ACM](#), 64(9):99–106, 2021.
- [7] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. [Proceedings of the International Conference on Learning Representations \(ICLR\)](#), 2021.
- [8] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhramil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In [The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track](#), 2024.
- [9] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In [First Conference on Language Modeling](#), 2024.
- [10] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. [arXiv preprint arXiv:2110.14168](#), 2021.
- [11] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. [arXiv preprint arXiv:2103.03874](#), 2021.
- [12] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. [arXiv preprint arXiv:2403.07974](#), 2024.
- [13] Lihu Chen and Gaël Varoquaux. What is the role of small models in the llm era: A survey. [arXiv preprint arXiv:2409.06857](#), 2024.

- [14] Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer. arXiv preprint arXiv:2203.03466, 2022.
- [15] Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, et al. Olmoe: Open mixture-of-experts language models. arXiv preprint arXiv:2409.02060, 2024.
- [16] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. arXiv:1803.05457v1, 2018.
- [17] Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. In Leon Derczynski, Wei Xu, Alan Ritter, and Tim Baldwin, editors, Proceedings of the 3rd Workshop on Noisy User-generated Text, pages 94–106, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.
- [18] Laura Weidinger, Deb Raji, Hanna Wallach, Margaret Mitchell, Angelina Wang, Olawale Salaudeen, Rishi Bommasani, Sayash Kapoor, Deep Ganguli, Sanmi Koyejo, et al. Toward an evaluation science for generative ai systems. arXiv preprint arXiv:2503.05336, 2025.
- [19] Zhongwei Wan, Xin Wang, Che Liu, Samiul Alam, Yu Zheng, Jiachen Liu, Zhongnan Qu, Shen Yan, Yi Zhu, Quanlu Zhang, et al. Efficient large language models: A survey. arXiv preprint arXiv:2312.03863, 2023.
- [20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. arXiv preprint arXiv:2001.08361, 2020.
- [21] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. arXiv preprint arXiv:2203.15556, 2022.
- [22] Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. A survey on large language models: Applications, challenges, limitations, and practical usage. Authorea Preprints, 3, 2023.
- [23] Derek Tam, Margaret Li, Prateek Yadav, Rickard Br  el Gabrielsson, Jiacheng Zhu, Kristjan Greenewald, Mikhail Yurochkin, Mohit Bansal, Colin Raffel, and Leshem Choshen. Llm merging: Building llms efficiently through merging. In NeurIPS 2024 Competition Track, 2024.
- [24] Shiwei Liu, Kai Han, Adriana Fernandez-Lopez, AJAY KUMAR JAISWAL, Zahra Atashgahi, Boqian Wu, Edoardo Ponti, Callie Hao, Rebekka Burkholz, Olga Saukh, et al. Edge-llms: Edge-device large language model competition. In NeurIPS 2024 Competition Track, 2024.
- [25] Neurips large language model efficiency challenge: 1 llm + 1gpu + 1day. <https://llm-efficiency-challenge.github.io/>, 2024.
- [26] Kurt W Fischer and Louise Silvern. Stages and individual differences in cognitive development. Annual review of psychology, 36(1):613–648, 1985.
- [27] Robert J Marzano and John S Kendall. The new taxonomy of educational objectives. Corwin Press, 2006.
- [28] David R Krathwohl, B Bloom, and B Masia. Taxonomy of educational objectives the classification of educational goals. Handbook II: Affective Domain. (New York: McKay, 1964.), 1964.
- [29] Guilherme Penedo, Hynek Kydl  cek, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale. In The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track, 2024.

- [30] Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. Fineweb-edu: the finest collection of educational content, 2024.
- [31] Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro von Werra, and Harm de Vries. The stack: 3 tb of permissively licensed source code. Preprint, 2022.
- [32] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2: The next generation, 2024.
- [33] Xiaotian Han, Yiren Jian, Xuefeng Hu, Haogeng Liu, Yiqi Wang, Qihang Fan, Yuang Ai, Huaibo Huang, Ran He, Zhenheng Yang, and Quanzeng You. Infimm-webmath-40b: Advancing multimodal pre-training for enhanced mathematical reasoning, 2024.
- [34] Liping Tang, Nikhil Ranjan, Omkar Pangarkar, Xuezhi Liang, Zhen Wang, Li An, Bhaskar Rao, Linghao Jin, Huijuan Wang, Zhoujun Cheng, Suqi Sun, Cun Mu, Victor Miller, Xuezhe Ma, Yue Peng, Zhengzhong Liu, and Eric P. Xing. Txt360: A top-quality llm pre-training dataset requires the perfect blend, 2024.
- [35] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and extraction. [arXiv preprint arXiv:2309.14316](#), 2023.
- [36] M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1-2):81–93, 1938.
- [37] Hynek Kydlíček, Guilherme Penedo, Clémentine Fourier, Nathan Habib, and Thomas Wolf. Finetasks: Finding signal in a haystack of 200+ multilingual tasks.
- [38] Zhen Xu, Sergio Escalera, Adrien Pavão, Magali Richard, Wei-Wei Tu, Quanming Yao, Huan Zhao, and Isabelle Guyon. Codabench: Flexible, easy-to-use, and reproducible meta-benchmark platform. *Patterns*, 3(7), 2022.
- [39] Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernandez. The LAMBADA dataset: Word prediction requiring a broad discourse context. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1525–1534, Berlin, Germany, August 2016. Association for Computational Linguistics.
- [40] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- [41] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023.
- [42] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. RACE: Large-scale ReAding comprehension dataset from examinations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.
- [43] Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on Freebase from question-answer pairs. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.

- [44] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, , and Jason Wei. Challenging big-bench tasks and whether chain-of-thought can solve them. [arXiv preprint arXiv:2210.09261](#), 2022.
- [45] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In [Proc. of NAACL](#), 2019.
- [46] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In [NAACL](#), 2019.
- [47] Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial nli: A new benchmark for natural language understanding. In [Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics](#). Association for Computational Linguistics, 2020.
- [48] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. [arXiv preprint arXiv:2109.07958](#), 2021.
- [49] Zayne Rea Sprague, Xi Ye, Kaj Bostrom, Swarat Chaudhuri, and Greg Durrett. MuSR: Testing the limits of chain-of-thought with multistep soft reasoning. In [The Twelfth International Conference on Learning Representations](#), 2024.

3 Resources

3.1 Organizing team

The team proposing this competition is composed of researchers and engineers in the domains of: Deep learning, NLP, data curation and management and software engineering. The short biography of each member is provided below. The organizers extend their gratitude to **Isabelle Guyon** for her valuable contributions to the discussions and insights related to the proposed challenge.

Mouadh Yagoubi (Lead organizer, Baseline method provider, Evaluator) is Lead Researcher at TII since November 2024. Prior to that, he was a Senior Researcher at IRT SystemX in France, where he contributed to several collaborative research projects in machine learning and optimization. He received his Ph.D. in Applied Mathematics from INRIA Saclay in 2012. His research interests include machine learning and evolutionary computation, with a particular focus on their application to industrial problems.

Yasser Dahou (Baseline method provider, Data provider, Evaluator) is a senior researcher at TII, working on multimodal vision-language research. Before joining TII, he earned a PhD in Computer Vision from Dublin City University in 2024, and a Master’s degree in Telecommunication Engineering in Algeria.

Billel Mokeddem (Data provider, Evaluator) is an AI Engineer at TII. Since he joined in February 2024, He worked on distributed text/audio data processing and model-based filtering, and LLMs pre-training. Prior to joining TII, he worked at Qatar Computing Research Institute as a Research Associate from 2021 to 2023, where he mainly focused on Multi-modality and the use of Variational Inference methods in Continuous Reinforcement Learning. He graduated from University of Science and technology Houari Boumediene (USTHB) in Algeria with a master’s degree in Artificial Intelligence in 2021.

Younes Belkada (Platform Administrator, Starting kit maintainer, Evaluator) is a Senior AI Engineer at TII, where he has been working since July 2024. He contributes to various projects involving pre-training, evaluation, and general tooling for AI and LLM-related initiatives. Prior to joining TII, he worked at Hugging Face from 2021 to 2024 as part of the open-source team, helping to develop a range of open-source deep learning libraries for the community. He graduated from ENS Paris-Saclay in 2021 with a master’s degree in Machine Learning and Computer Vision.

Phuc H. Le-Khac (Platform Administrator, Starting Kit maintainer, Evaluator) Phuc is a Deep Learning researcher at Technology Innovation Institute in Abu Dhabi and previously obtained his PhD in Representation Learning at Dublin City University. He is interested in both the technical details as well as the broader theory of learning system. He is currently applying his interest on developing next-generation Multimodal Model.

Basma El Amel Boussaha (Data provider, Evaluator) is a Lead Researcher at TII, specializing in building Large Language Models (LLMs) for Arabic. She earned her Ph.D. in Natural Language Processing from Nantes University in 2019. Prior to joining TII, she worked at Della AI, where she played a key role in developing AI-driven contract reviewing tools to support legal professionals. Her research continues to bridge cutting-edge AI with real-world impact across languages and domains with a focus on generative AI, question answering, distributed training, and advancing Arabic NLP.

Reda Alami (Evaluator) is currently a senior researcher at the Technology Innovation Institute in Abu Dhabi. His research interests include finetuning of large language models for reasoning tasks, test time compute scaling and reinforcement learning. He received his Msc degree with first class honours from IMT Atlantique, France. Then, he received his Phd degree from Paris Saclay University, France. He was with Orange Labs, Lannion, France, from 2016 to 2019 and then with TotalEnergies, Saclay, France from 2020 to 2021

Jingwei zuo (Evaluator) Jingwei Zuo is currently a Lead Researcher at the AI Center of the Technology Innovation Institute (TII). He and his team are working on the Falcon Foundation Models, with a focus on exploring next-generation Large Language Models (LLMs) that incorporate emerging architectural designs beyond or in combination with Transformers, as well as optimal training and model scaling strategies. Jingwei received his PhD in Computer Science from Université Paris-Saclay and was awarded the Plateau de Saclay Doctoral Prize for the best scientific production in Information and Communication Science and Technology (ICST) in 2022.

Damiano Marsili (Evaluator) Damiano is a Ph.D student at Caltech advised by Georgia Gkioxari and Pietro Perona. Before his Ph.D, he was an Applied Science Intern at Amazon Robotics working on 3D spatial reasoning. Prior to that, he double-majored at Johns Hopkins in Computer Science and Mathematics.

Mugariya Farooq (Evaluator) is an Electronics Engineer with a Masters in Machine Learning and is currently an AI Engineer in Technology Innovation Institute. She has led evaluations for various models in Falcon family of models with prime focus on optimization and development of model agnostic evaluation pipeline and benchmarks.

Mounia Lalmas (Scientific Advisor) Mounia Lalmas is a Senior Director of Research at Spotify, and the Head of Tech Research in Personalisation, where she leads an interdisciplinary team of research scientists, working on personalization. Mounia also holds an honorary professorship at University College London. She also holds an additional appointment as a Distinguished Research Fellow at the University of Amsterdam. Before that, she was a Director of Research at Yahoo, where she led a team of researchers working on advertising quality. She also worked with various teams at Yahoo on topics related to user engagement in the context of news, search, and user-generated content. Prior to this, she held a Microsoft Research/RAEng Research Chair at the School of Computing Science, University of Glasgow. Before that, she was Professor of Information Retrieval at the Department of Computer Science at Queen Mary, University of London. She is regularly a senior programme committee member at conferences such as WSDM, KDD, WWW and SIGIR. She was programme co-chair for SIGIR 2015, WWW 2018 and WSDM 2020, and CIKM 2023.

Georgia Gkioxari (Scientific Advisor) Georgia is an assistant professor at the Computing + Mathematical Sciences at Caltech. She obtained her PhD in Electrical Engineering and Computer Science from UC Berkeley, where she was advised by Jitendra Malik. Prior to Berkeley, she earned her diploma from the National Technical University of Athens in Greece. After earning her PhD, she was a research scientist at Meta’s FAIR team. In 2021, she received the PAMI Young Researcher Award, which recognizes a young researcher for their distinguished research contribution to computer vision. She is the recipient of the PAMI Mark Everingham Award for the open-source software suite Detectron (2021), the Google Faculty Award (2024) and the Okawa Research Award (2024). In 2017, Georgia and her co-authors received the Marr Prize for “Mask R-CNN” published and presented at ICCV. She was named one of 30 influential women advancing AI in 2019 by ReWork and was nominated for the Women in AI Awards in 2020 by VentureBeat.

Patrick Gallinari (Scientific Advisor) is a professor in Computer Science at Sorbonne University in Paris. His research focuses on statistical learning with applications in different fields such as semantic data analysis and complex data modeling. He discovered the ML domain in the mid 80es when he started to work on Neural Networks, an emerging field at that time. He has been one of the pioneers of this research domain in France/ Europe and worked on NN and on other ML models since that. He investigated different application domains like Information Retrieval, Social Data analysis, User Modeling. Today his main focus is on Physics Aware Deep Learning and on some aspects of Natural Language Processing. He has been leading the Machine Learning team MLIA for some years. He has been director of the computer science lab. at Sorbonne University (LIP6) for 9 years (2005 to 2013) and vice director for 6 years before, He also acted as vice director of the scientific committee of the faculty of engineering at UPMC (2010 to 2021).

Philip Torr (Scientific Advisor) Professor Philip Torr did his PhD (DPhil) at the Robotics Research Group of the University of Oxford under Professor David Murray of the Active Vision Group. He left Oxford to work for six years as a research scientist for Microsoft Research, first in Redmond, USA, in the Vision Technology Group, then in Cambridge founding the vision side of the Machine Learning and Perception Group. He then became a Professor in Computer Vision and Machine Learning at Oxford Brookes University. In 2013, Philip returned to Oxford as full professor where he has established the Torr Vision group. He won several awards including the Marr prize (the highest honour in vision) in 1998. He is a Royal Society Wolfson Research Merit Award Holder. He was elected Fellow of the Royal Academy of Engineering (FREng) in 2019, and Fellow of the Royal Society (FRS) in 2021 for contributions to computer vision. In 2021 he was made Turing AI world leading researcher fellow.

Hakim Hacid (Scientific Advisor) Hakim Hacid is the Chief Researcher of the Artificial Intelligence and Digital Science Research Center in Technology Innovation Institute (TII), leading the diverse efforts around LLM's and Machine Learning. Prior to joining TII, he was an Associate Professor at Zayed University, a customer analytics head at Zain telecom, and a research department head at Bell Labs Research. He is a published author of many research articles in top journals conferences and holds several industrial patents to his credit. His research specialization includes machine learning, databases, natural language processing, security. He obtained his PhD in Data Mining/Databases and also a double master's in Computer Sc (Master by Research & Professional Master) from University of Lyon, France.

3.2 Resources provided by organizers

In order to evaluate submissions and compute the global score, a dozen H100 GPUs will be made available through GCP infrastructure to power the compute workers connected to Codabench. In addition to a webinar that will be organized to help participants get started, a training session on lm-eval-harness framework will also be provided and made available on the competition website to support participants who are not experts in machine learning.

3.3 Support requested

We would greatly appreciate the support of the NeurIPS 2025 Competition Track organizers, particularly in promoting our challenge through their official channels to help us reach participants worldwide.

A Results of state-of-the-art benchmarks with different model sizes

we provide in this section the results of state-of-the-art benchmarks

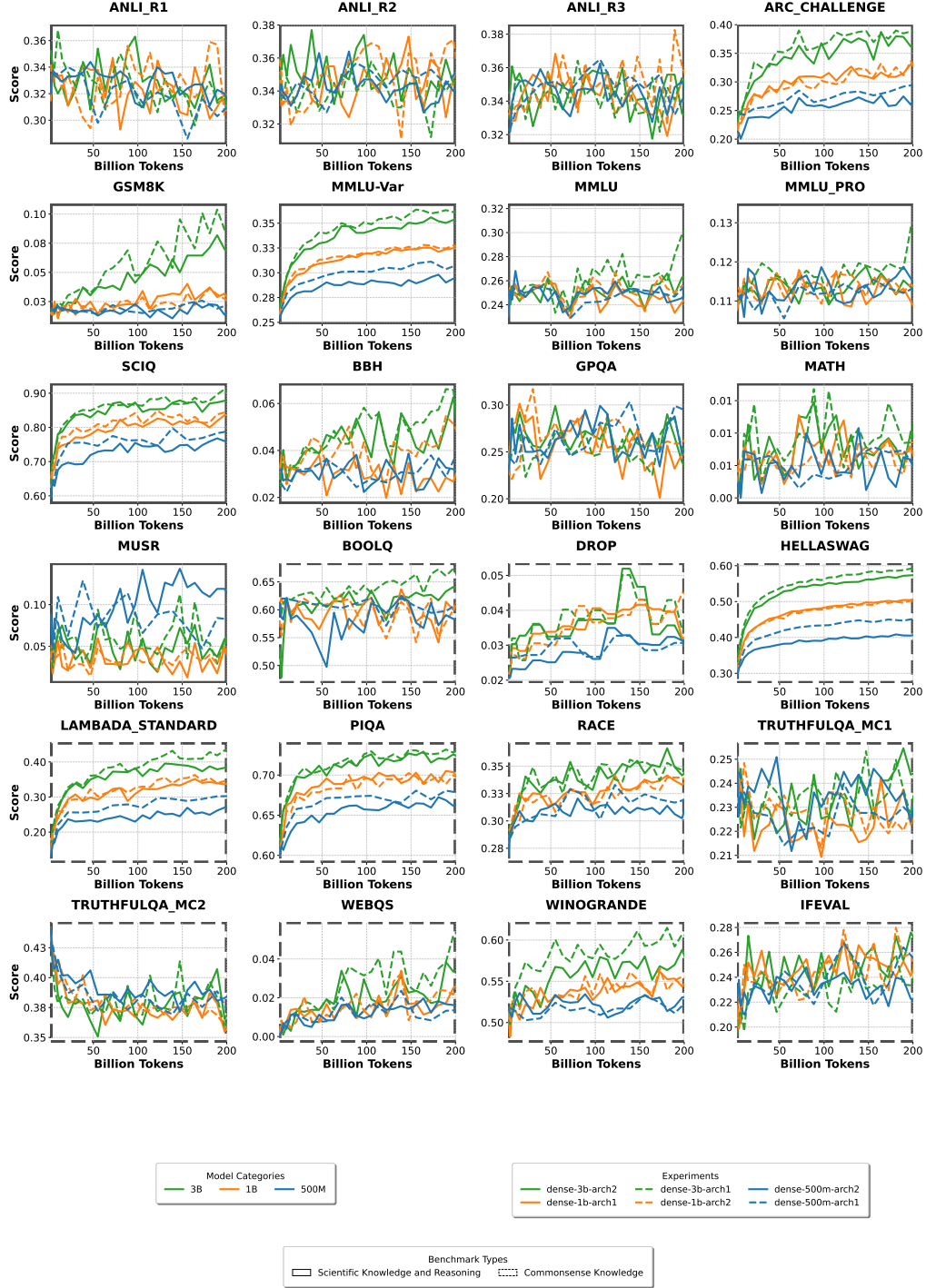


Figure 4: Results of state-of-the-art benchmarks with different models sizes.

B Example of MMLU and MMLU-var benchmarks

We present in this section an example of an evaluation task to highlight the differences between the standard version of MMLU and the MMLU-var variant.

The **standard MMLU benchmark** uses a multiple choice format: The model receives a question along with a list of possible answers (A, B, C, D) and is prompted to select the correct one. For example, in a question about k-means clustering, the model sees both the question and all four answer choices in the input.

Standard MMLU prompt example:

Question: You want to cluster 7 points into 3 clusters using the k-Means Clustering algorithm. Suppose after the first iteration, clusters C1, C2 and C3 contain the following two-dimensional points...

A. C1: (3,3), C2: (4,4), C3: (6,6)
B. C1: (3,3), C2: (6,6), C3: (12,12)
C. C1: (6,6), C2: (12,12), C3: (12,12)
D. C1: (0,0), C2: (48,48), C3: (35,35)

Answer:

In contrast, **MMLU-var** uses a continuation-based approach. The model is given only the question followed by the prompt "Answer:", without any answer choices. The system then evaluates the likelihood the model assigns to each possible answer (as separate completions), and selects the one with the highest log-probability.

MMLU-var prompt example:

Question: You want to cluster 7 points into 3 clusters using the k-Means Clustering algorithm. Suppose after the first iteration, clusters C1, C2 and C3 contain the following two-dimensional points...

Answer:

Each answer candidate (e.g., C1: (3,3), C2: (4,4), C3: (6,6)) is then evaluated separately by computing its log-probability as a continuation of the prompt above.

C Models architectural details

Model size	Model name	# layers	# attention heads	Hidden size	Intermediate size
0.5B	dense-0.5B-arch1	32	16	1024	3072
	dense-0.5B-arch2	16	20	1280	3072
1B	dense-1b-arch1	32	24	1536	4096
	dense-1b-arch2	16	16	2048	7168
3B	dense-3b-arch1	32	20	2560	8960
	dense-3b-arch2	16	24	3072	15360

Table 2: Model architectures. kv heads = 4 for all the models.

D Starting kit description

Since the maximum model size we offer is around 3 Billion parameters, all starting kits can be executed on free-tier devices that can be accessible world-wide such as Google Colab or Kaggle free-tier GPU notebooks. This makes the competition free from compute constraints.

- `0-Basic_Competition_Information`: This notebook contains general information concerning the competition organization, phases, deadlines and terms. The content is the same as the one shared in the competition Codabench page.
- `1-How_to_interact_with_model`: This notebook aims to familiarize participants with the tools that are used to interact with the model and perform some easy text generation tasks.
- `2-How_to_evaluate_a_model`: Shows participants how a checkpoint of the model can be evaluated using `lm-evaluation-harness` package.
- `3-Reproduce_baseline_results`: This notebook shows how the baseline results (MMLU-Var on a single checkpoint) could be reproduced. It includes integrating the MMLU-Var benchmark within the `lm-evaluation-harness` package and running it to get its result.
- `4-How_to_Contribute`: This notebook explains how to fully integrate a new task within the `lm-evaluation-harness` package.
- `5-Scoring`: This notebook shows firstly how the score is computed by describing its different components. Next, it provides a script which can be used locally by the participants to obtain a score for their contributions. We encourage participants to evaluate their solutions via codabench (which uses the same scoring module as the one described in this notebook).
- `6-Submission_examples`: This notebook presents the composition of a submission bundle for Codabench and usable parameters

E Metrics Calculation Examples

E.1 $Score_{SQ}$ Calculation examples

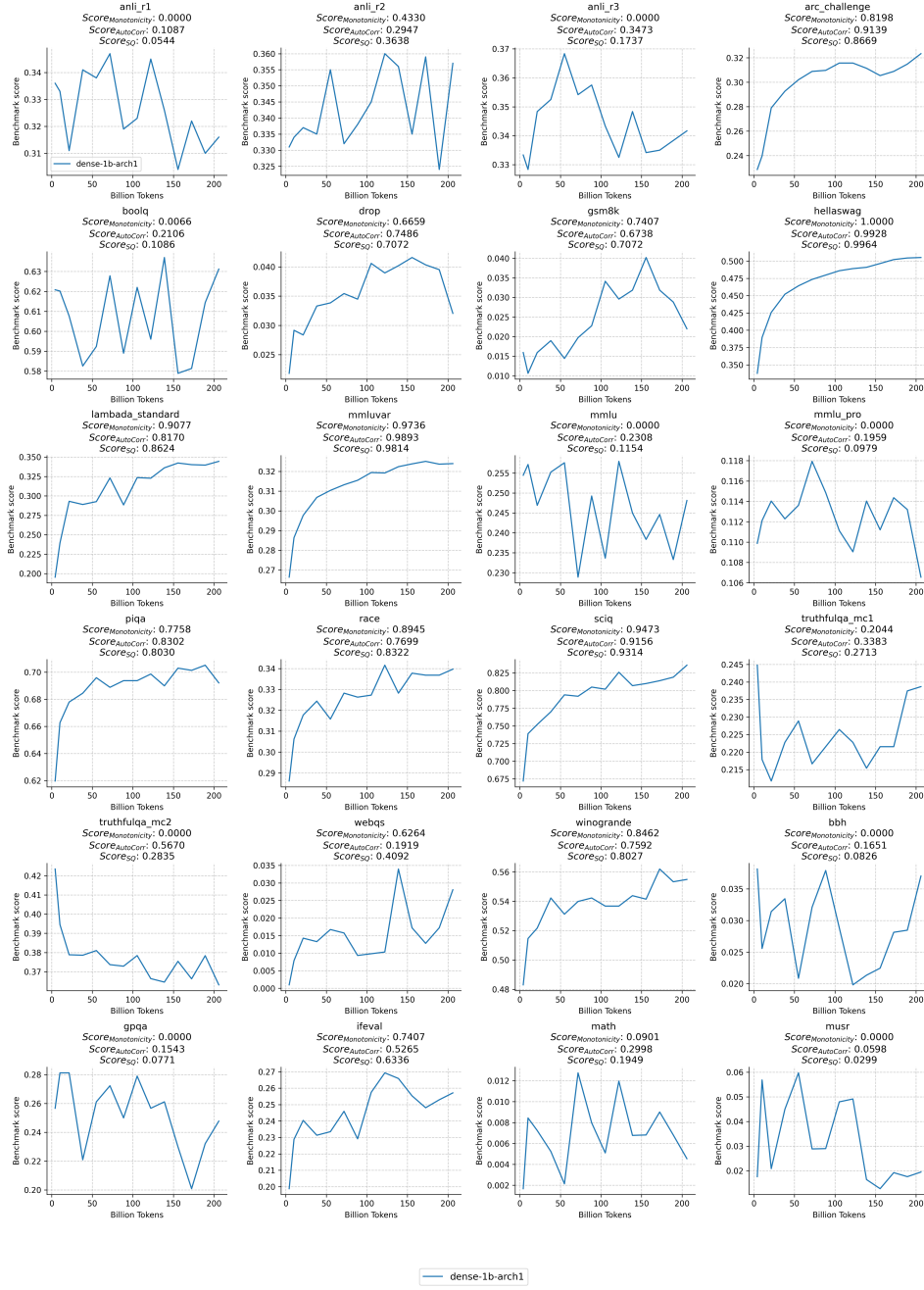


Figure 5: Signal quality metric calculated for different benchmarks using 1B model. (SC corresponds to $Score_{Monotonicity}$ in equation 2, AC corresponds to $Score_{AutoCorr}$ in equation 4 and SQ corresponds to $Score_{SQ}$ in equation 5 with $\beta_{Mono} = \beta_{AutoCorr} = 0.5$).

E.2 $Score_{RC}$ Calculation examples

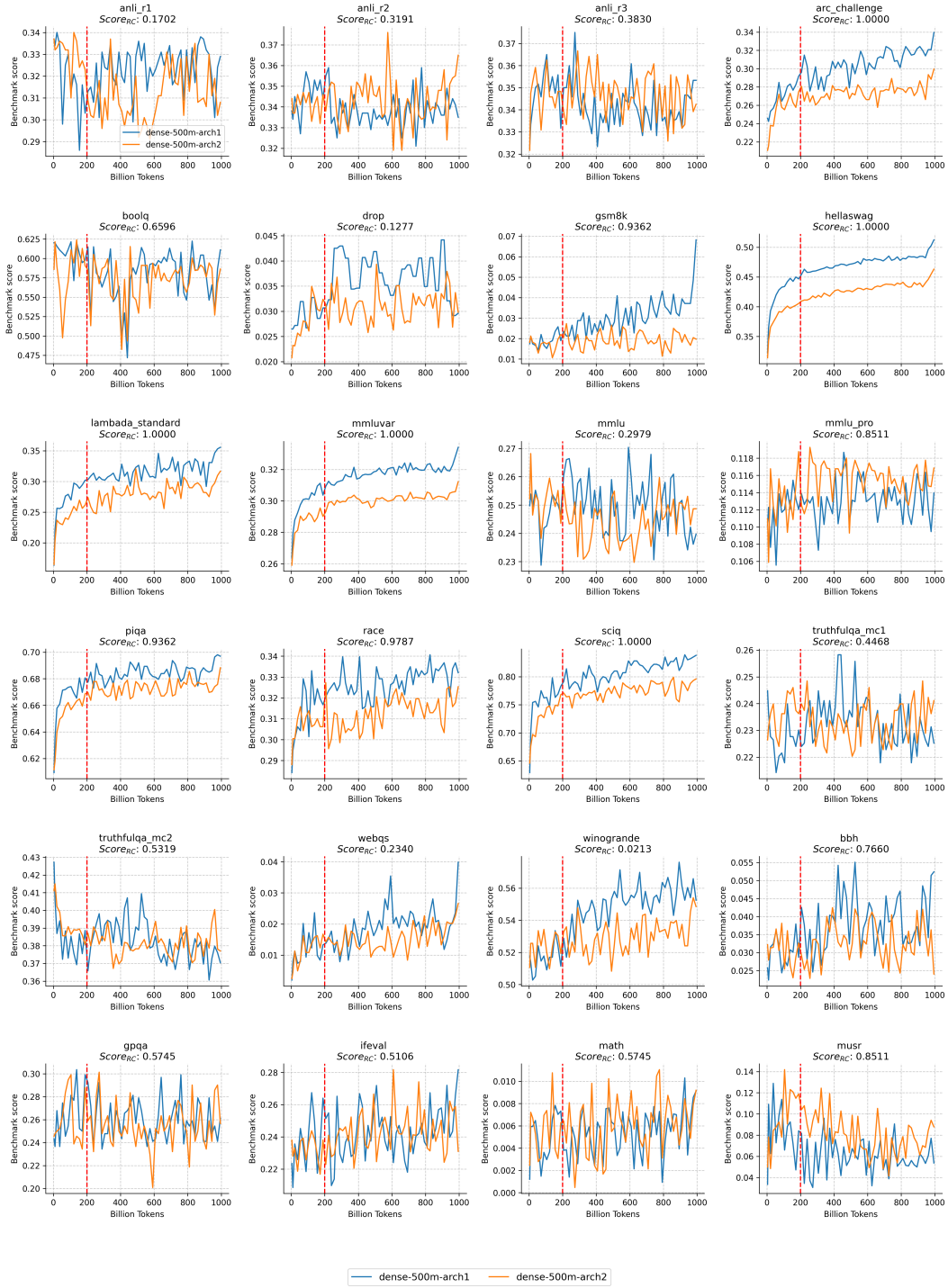


Figure 6: Ranking Consistency metric calculated for different benchmarks using 0.5B model.

E.3 $Score_{CS}$ Calculation examples

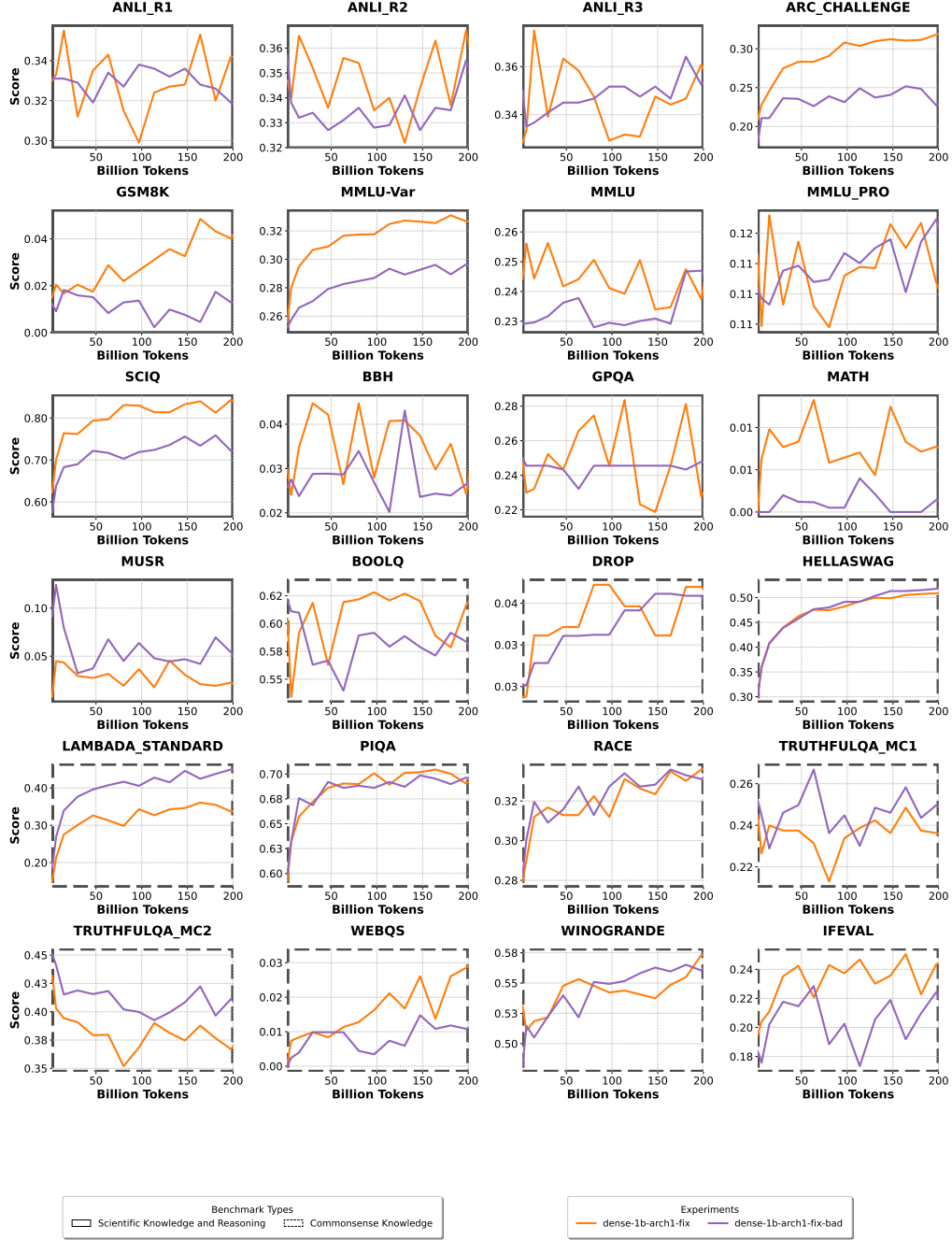


Figure 7: Compliance to scientific knowledge domains metric calculated for different benchmarks using two 1B models, one trained with a web-only data mixture and the other one trained with a rich scientific data mixture. For all the benchmarks, we use acc_norm as a metric except for TruthfulQA, we use MC1 and MC2.

E.4 MMLU-var Score Calculation Example

Calculating Signal Quality After Calculating the monotonicity score using Equation (2) and the Autocorrelation Strength using Equation (4), we compute the Signal quality score using Equation (5), where $\beta_{Monotonicity} = 0.5$ and $\beta_{AutoCorr} = 0.5$. The results for the 3 models are given in Table 3:

Model	Score _{Monotonicity}	Score _{AutoCorr}	Score _{SQ}
dense-500m-arch1	0.9387	0.9324	0.9354
dense-1b-arch1	0.9736	0.9893	0.9814
dense-3b-arch1	0.9516	0.9685	0.9601

Table 3: Detailed Signal Quality scores across the different models

Then, we average over the different model sizes to get the final Signal Quality score:

$$\text{Score}_{SQ} = \frac{0.9354 + 0.9814 + 0.9601}{3} = 0.9590$$

Calculating Ranking Consistency We compute the ranking at 200 billion tokens using Equation (6) then, the ranking consistency until 1 trillion tokens using Equation (7) for each model size $s \in \{0.5B, 1B, 3B\}$. The results are reported in Table 4:

Models	Score _{RC}
(dense-0.5b-arch1, dense-0.5b-arch2)	1.0
(dense-1b-arch1, dense-1b-arch2)	0.5106
(dense-3b-arch1, dense-3b-arch2)	1.0

Table 4: Detailed Ranking Consistency scores across different model sizes

After that we average on model sizes:

$$\text{Score}_{RC} = \frac{1.0 + 0.5106 + 1.0}{3} = 0.8368$$

Calculating Compliance Score We compute the compliance score using Equation (8) between the two 1B models with the same architecture (Arch-1), but one was trained with a rich scientific knowledge-focused datamixture and the other one was trained using a web-only datamixture. We obtained for MMLUvar the result: 0.384

Calculating the total score Using Equation 1, we compute our metric’s total score by setting $\alpha_{SQ} = 0.4$, $\alpha_{RC} = 0.2$, and $\alpha_{CS} = 0.4$.

$$\text{Score} = 0.5 \times 0.9590 + 0.1 \times 0.8369 + 0.4 \times 0.384 = 0.7167$$

F Check procedures

F.1 Scientific alignment check

To ensure that submitted benchmarks adhere to the expected *scientific knowledge domains*, we implement an initial compliance check using *gpt-4o-2024-08-06* as a zero-shot classifier. This step is designed to automatically filter out benchmarks that rely solely on general language understanding or common-sense reasoning, and to promote benchmarks that demand substantive academic or scientific knowledge.

Classification Protocol: Each submission is evaluated by prompting *gpt-4o-2024-08-06* with a structured instruction that defines two categories:

- **Accept:** for tasks requiring domain-specific, scientific, or expert-level reasoning.
- **Reject:** for tasks relying only on general language skills, common sense, or trivia.

The complete classification prompt is provided below:

Scientific scoring prompt

Classification Instructions

You are tasked with classifying each of the following question-answer pairs into one of two groups:

Accept

Classify as **Accept** if the question requires **domain-specific knowledge, scientific understanding, academic expertise, or professional reasoning**. This includes but is not limited to:

Scientific & Technical Fields

- Biology
- Chemistry
- Physics
- Mathematics
- Computer Science
- Engineering (all disciplines)
- Medicine and Health Sciences
- Neuroscience
- Pharmacology
- Veterinary Science
- Environmental Science
- Earth Science / Astronomy
- Statistics & Data Science

Professional & Applied Domains

- Law
- Business (e.g., Finance, Accounting, Marketing)
- Economics
- Political Science
- Education
- Sociology
- Anthropology
- Linguistics
- Communications / Media Studies
- Library & Information Science
- Social Work
- Public Policy
- Nursing / Allied Health
- Architecture / Urban Planning
- Agriculture / Food Science

Humanities with Reasoning Requirements

- History
- Philosophy (e.g., logic, ethics, epistemology)
- Theology / Religious Studies
- Art History
- Literary Theory

Any question that falls under these or similar knowledge-intensive domains should be marked **Accept**, even if it's not explicitly listed.

Reject

Classify as **Reject** only if the question is based on:

- General language understanding
- Common sense or cultural knowledge
- Vocabulary, grammar, spelling, or idioms
- Word analogies or word associations
- Trivia or factoids that don't require reasoning
- Sentiment, emotion, or tone recognition
- Simple reading comprehension without technical content
- NLP-specific tasks (e.g., joke detection, paraphrasing, summarization)

These do not require deep subject-matter knowledge and should be marked **Reject**.

Thresholding and Manual Review: For each benchmark, we compute the proportion of questions classified as **Accept**. If this classification accuracy exceeds **80%**, the benchmark is considered **automatically compliant** and advances to the next evaluation stage.

If the **Accept accuracy is below 80%**, the submission is flagged for **manual review** by domain experts. This ensures that potentially valuable but misclassified benchmarks are not prematurely excluded.

Observed Performance on Example Benchmarks: To calibrate the prompt, we tested the classification protocol on existing benchmarks:

- **HellaSwag** (focused on common-sense reasoning): **10% Accept**, indicating low compliance with the scientific knowledge requirement.
- **MMLU**: **100% Accept**, demonstrating full alignment with domain expectations.

As shown in Table 5, the benchmarks are well classified, where the scientific ones get a high accuracy whereas the rest obtains a lower score (HellaSwag, TruthfulQA, WinoGrande). This automated classifier serves as an efficient preliminary filter while still allowing judgment in edge cases.

Benchmark	Gpt4-o classification accuracy
MMLU	100
HellaSwag	10
ARC-Easy	93
SciQ	87
TruthfulQA	27
WinoGrande	7
MuSR	34

Table 5: Gpt4-o classification of the common benchmarks using the detailed prompt

F.2 Leakage check

We provide an additional leakage check to make sure the expected answer is not semantically present within the context, by doing a simple string match between elements in the answer with the context. We calculate the leakage rate for SCiQ[17] benchmark using the snippet below and obtained a score of 0.82, which might explain the quality of signal of that benchmark for small models.

Below is an example of the test set of SCiQ benchmark which contains answer leakage within the passed context.

Question sample from SCiQ

Oxidants and Reductants

Compounds that are capable of accepting electrons, such as O_2 or F_2 , are called **oxidants** (or oxidizing agents) because they can oxidize other compounds. In the process of accepting electrons, an oxidant is reduced. Compounds that are capable of donating electrons, such as sodium metal or cyclohexane (C_6H_{12}), are called reductants (or reducing agents) because they can cause the reduction of another compound. In the process of donating electrons, a reductant is oxidized. These relationships are summarized in Equation 3.30: Equation 3.30 Saylor URL: <http://www.saylor.org/books>.

Question: Compounds that are capable of accepting electrons, such as O_2 or F_2 , are called what?

Answer:

Answers sample from previous sample

- residues
- antioxidants
- Oxygen
- **oxidants**

We will conduct similar checks for all proposed benchmarks and provide a score for each of them when relevant. For other benchmarks such as MMLU, this score is not relevant as answers are formatted in MCQ format:

Question and Answer sample from MMLU

Find the degree for the given field extension $\mathbb{Q}(\sqrt{2}, \sqrt{3}, \sqrt{18})$ over \mathbb{Q} .

- A. 0
- B. 4
- C. 2
- D. 6

Answer:

Since the expected answer from the model (choice from A, B, C and D) is not within the context, this check is not relevant for MCQ-style benchmarks.

G Global scores

The table presents a comprehensive ranking of 25 benchmarks based on their performance across three metrics that we proposed. MMLU-var ranks highest with a total score of 0.717, followed by ARC-Easy and SciQ in the second and third positions with total scores of 0.655 and 0.627 respectively. Interestingly, HellaSwag achieved very high scores on the first two metrics, as it is well adapted to early training assessment. However, it received a score of 0 on the third metric as it is not designed to assess scientific knowledge. Note that several tasks known to be relevant to scientific knowledge, such as Math and MMLU-Pro, received relatively poor scores. This is principally due to the fact that their low score in the SC metric is correlated with their low score in the SQ metric, as they are noisy and fail to provide a clear signal at this stage of training, making them unable to discriminate effectively between the web-only model and the knowledge-oriented model.

Benchmark	SQ (50%)	RC (10%)	SC (40%)	Total Score
MMLU-var[15]	0.959	0.837	0.384	0.717
ARC-Easy[16]	0.832	0.822	0.394	0.655
SciQ[17]	0.846	0.772	0.316	0.627
HellaSwag[5]	0.992	0.936	0.000	0.590
GSM8K[10]	0.655	0.915	0.369	0.566
LAMBADA_Standard[39]	0.894	0.752	0.000	0.522
PIQA[40]	0.842	0.690	0.013	0.495
IFEval[41]	0.532	0.562	0.353	0.463
RACE[42]	0.677	0.877	0.000	0.426
WebQuestions[43]	0.477	0.573	0.256	0.398
MATH[11]	0.266	0.672	0.494	0.398
WinoGrande[6]	0.590	0.588	0.022	0.363
MMLU[7]	0.285	0.574	0.400	0.360
BBH[44]	0.352	0.675	0.283	0.357
DROP[45]	0.573	0.386	0.059	0.349
BoolQ[46]	0.327	0.771	0.188	0.316
ANLI_r2[47]	0.296	0.433	0.241	0.288
MMLU_Pro[8]	0.342	0.716	0.000	0.243
GPQA[9]	0.227	0.402	0.049	0.173
ANLI_r3[47]	0.224	0.425	0.000	0.154
TruthfulQA_mc2[48]	0.180	0.625	0.000	0.152
TruthfulQA_mc1[48]	0.226	0.355	0.000	0.149
MuSR[49]	0.116	0.713	0.000	0.129
ANLI_r1[47]	0.148	0.319	0.000	0.106

Table 6: Benchmark Performance Comparison by Rank. **SQ** = Signal Quality (50%), **RC** = Ranking Consistency (10%), **SC** = Scientific Knowledge Compliance (40%). **Gray-shaded rows** indicate benchmarks classified as targeting scientific knowledge & reasoning. Other tasks are related to commonsense knowledge

H Tracking Validation loss

Tracking validation loss during model training is a standard practice in machine learning, with LLM scaling laws using predicted loss as a proxy for model capability. Figure 8 shows validation loss across training iterations for two model variants trained on different data mixtures: Web-Only and Scientific-Mix.

Our analysis of these loss curves reveals several important findings:

1. **Domain transfer gap:** While FineWebEdu (a subset of the training data for Web-Only) shows a small gap between variants, all other datasets exhibit substantially larger gaps, often exceeding the total reduction in loss throughout training, indicating significant out-of-distribution effects for the WebOnly model.

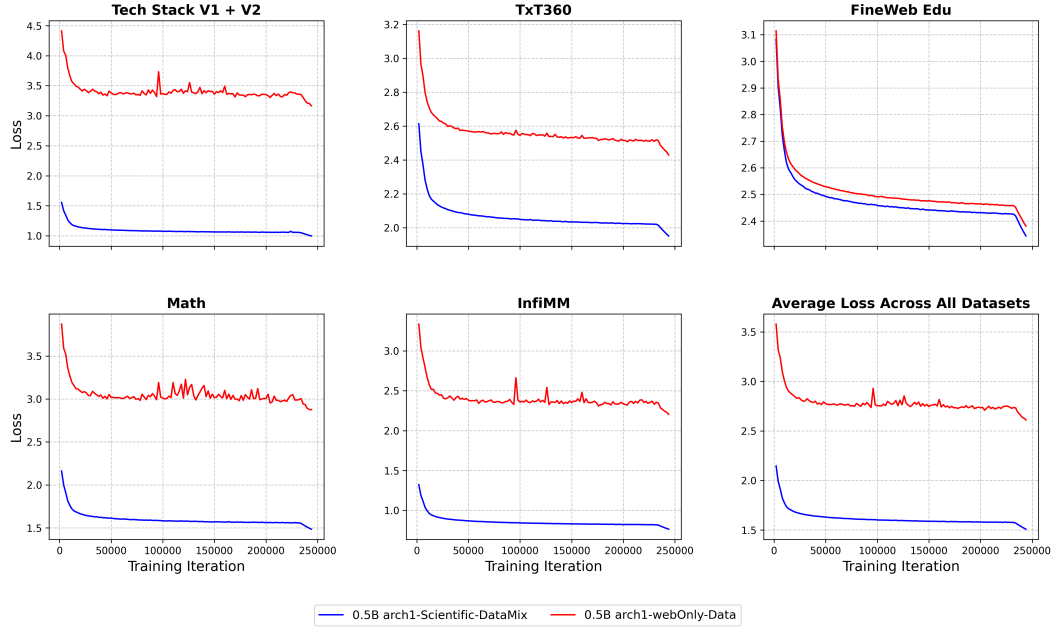


Figure 8: Validation loss across different datasets for two model variants: Scientific-DataMix (blue) and webOnly-Data (red).

2. **Cross-dataset comparison challenges:** Each dataset produces distinct numerical loss ranges due to inherent characteristics, making direct comparison across domains difficult. Lower loss on one dataset does not necessarily translate to higher capability on corresponding benchmarks.
3. **Loss interpretation limitations:** The small numerical differences in loss values throughout training are difficult to translate into meaningful capability improvements, particularly when averaging across datasets with different baseline values.
4. **Early training ambiguity:** Despite Scientific-Mix eventually showing lower loss on FineWebEdu compared to WebOnly, both variants exhibit similar loss trajectories in early training stages, making it challenging to distinguish performance differences during this critical period.

These observations highlight the limitations of using raw validation loss as the sole metric for evaluating model capability across different data mixtures and domains, especially during the early phase of training.