

## Anonymous ACL submission

## Abstract

Language models (LMs) have become ubiquitous in both NLP research and in commercial product offerings. As their commercial importance has surged, the most powerful models have become closed off, gated behind proprietary interfaces, with important details of their training data, architectures, and development undisclosed. Given the importance of these details in scientifically studying these models, including their biases and potential risks, we believe it is essential for the research community to have access to powerful, truly open LMs. To this end, we have built OLMo, a competitive, truly **Open Language Model**, to enable the scientific study of language models. Unlike most prior efforts that have only released model weights and inference code, we release OLMo alongside open training data and training and evaluation code. We hope this release will empower and strengthen the open research community and inspire a new wave of innovation.

## 1 Introduction

Language models have been at the center of NLP technologies for many years (Rosenfeld, 2000; Bengio et al., 2003; Mikolov et al., 2013; Peters et al., 2018; Brown et al., 2020). Recently, due to large-scale pretraining and human annotation for alignment, they have become commercially valuable (OpenAI, 2023). However, as their commercial value has increased, the largest models have become gated behind proprietary interfaces, with important details left undisclosed.

We believe that full access to open language models for the research community is critical to the scientific study of these models, their strengths and weaknesses, and their biases and risks. Accordingly, we introduce **OLMo**, a powerful, truly open language model alongside open training data, training and evaluation code, intermediate model checkpoints, and training logs.

Recent LM releases have varied in their degree of openness. For example, Mixtral 8x7B provided model weights and a brief report (Jiang et al., 2024), while LLaMA came with in-depth adaptation training instructions (Touvron et al., 2023b), and Mosaic Pretrained Transformer came with many details, including the dataset distribution, though not the data itself (MosaicML NLP Team, 2023). Falcon’s pretraining data was partially released (Almazrouei et al., 2023), and the most open models—the Pythia suite (Biderman et al., 2023) and BLOOM (BigScience et al., 2022)—released training code, model checkpoints, training data and more.

With OLMo, we release the whole framework from data to training to evaluation tools: multiple training checkpoints across multiple hardware types, training logs, and exact datasets used, with a permissive license. We are not the only team to do this; recent work from LLM360 targets similar goals (Liu et al., 2023). OLMo narrows the gap from their models to state-of-the-art capabilities of models like LLaMA2. This project has benefited from lessons learned from all of these previous efforts with their varying degrees of openness, and we believe that a large, diverse population of open models is the best hope for scientific progress on understanding language models and engineering progress on improving their utility.

The OLMo framework encompasses the tools and resources required for building and researching language models. For training and modeling, it includes full model weights, training code, training logs, ablations, training metrics in the form of Weights & Biases logs, and inference code. The released model includes four variants of our language model at the 7B scale corresponding to different architectures, optimizers, and training hardware, and one model at the 1B scale, all trained on at least 2T tokens. We also release hundreds of intermediate checkpoints available as revisions on HuggingFace.

Size	L	D	H	Tokens	Peak LR	Betas	Epsilon	WD	Batch size
1B	16	2048	16	2T	4.0E-4	(0.9, 0.95)	1.0E-5	0.1	~4M
7B	32	4086	32	2.46T	3.0E-4	(0.9, 0.95)	1.0E-5	0.1	~4M
65B*	80	8192	64	TBD	1.5E-4	(0.9, 0.95)	1.0E-5	0.1	~2M → ~16M

Table 1: OLMo model sizes, number of training tokens, and AdamW optimizer settings. **L** is number of layers, **D** is hidden dimension, **H** is number of attention heads, **WD** is weight decay. The batch size of the 65B model is repeatedly doubled during training,  $\sim 2\text{M} \rightarrow \sim 4\text{M} \rightarrow \sim 8\text{M} \rightarrow \sim 16\text{M}$ .

\* At the time of writing our 65B model is still training.

For dataset building and analysis, the full training data used for these models is openly available (Dolma; Anonymous, 2024), including code that produces the training data, and tools for analyzing pretraining data (Anonymous, 2023c). For evaluation, we build on Catwalk (Anonymous, 2023a) for downstream evaluation and Paloma (Anonymous, 2023b) for perplexity-based evaluation. For adaptation, we use Open Instruct (Iverson et al., 2023; Wang et al., 2023) to train with instruction and feedback data. Finally, all code and weights are released under the Apache 2.0 License.

With this release, we hope to catalyze research into as-yet poorly understood aspects of these models, for example, the relationship between pretraining data and model capabilities, the impact of design and hyperparameter choices, and various optimization methods and their impact on model training. In addition, we report on the lessons learned and important details necessary to successfully train language models at this scale.

## 2 OLMo Framework

This section describes the OLMo framework, consisting of the OLMo models (Section 2.1), our pretraining dataset, Dolma (Section 2.2), and our evaluation framework (Section 2.4).

### 2.1 OLMo Model and Architecture

We adopt a decoder-only transformer architecture based on (Vaswani et al., 2017), and deliver 1B and 7B variants as described in Table 1, with a 65B version coming soon. Our specific architecture includes several improvements over the vanilla transformer from (Vaswani et al., 2017) following other recent large language models like PaLM (Chowdhery et al., 2022), the LLaMA family (Touvron et al., 2023a,b), OpenLM (Gururangan et al., 2023), and Falcon (Almazrouei et al., 2023). Check table 5 in Appendix A for a comprehensive comparison of our 7B architecture to the similarly-sized models

from these other families.

We generally select hyperparameters by optimizing for training throughput on our hardware while minimizing the risk of loss spikes and slow divergence. We ablate choices through our in-loop evaluation setting, given available computational sources (Section 2.4). Table 5 compares our design choices with recent state-of-the-art open language models. Our main changes over the vanilla transformer architecture can be summarized as follows:

1. **No biases.** Following LLaMA, PaLM, and others, we exclude all bias terms from our architecture in order to improve training stability.
2. **Non-parametric layer norm.** We use the non-parametric formulation of layer norm (Ba et al., 2016) in which there is no affine transformation within the norm, i.e. no “adaptive gain” (or bias). We believe this was the safest option and it was also the fastest compared to the other variants we considered: parametric layer norm and RMSNorm (Zhang and Sennrich, 2019).
3. **SwiGLU activation function.** Like LLaMA, PaLM, and others we use the SwiGLU activation function (Shazeer, 2020) instead of ReLU, and following LLaMA the activation hidden size is approximately  $\frac{8}{3}d$ , but increased to the closest multiple of 128 (e.g. 11,008 for our 7B model) to improve throughput.<sup>1</sup>
4. **Rotary positional embeddings (RoPE).** Like LLaMA, PaLM, and others we replace absolute positional embeddings with rotary positional embeddings (RoPE; Su et al., 2021).
5. **Vocabulary.** We use a modified version of the BPE-based tokenizer from GPT-NeoX-20B (Black et al., 2022) with additional tokens for

<sup>1</sup> Since SwiGLU is a “gated” activation function, the output is half the size of the input. So technically our inputs to SwiGLU have a dimensionality of  $2 \times 11,008 = 22,016$  for our 7B model.

masking personal identifiable information (PII). The final vocabulary size is 50,280. However, to maximize training throughput we increase the size of the corresponding embedding matrix in our model to 50,304 so that it’s a multiple of 128.

## 2.2 Pretraining Data: Dolma

Despite progress in access to model parameters, pretraining datasets are still not as open. Pretraining data are often not released alongside open models (let alone closed models) and documentation about such data is often lacking in detail that would be needed to reproduce or fully understand the work. This has made it difficult to support certain threads of language model research, such as understanding how training data impacts model capabilities and limitations. To facilitate open research on language model pretraining, we built and released our pretraining dataset, Dolma—a diverse, multi-source corpus of 3T tokens across 5B documents acquired from 7 different data sources that are (1) commonly seen in large-scale language model pretraining and (2) accessible to the general public (Anonymous, 2024). Table 2 provides a high-level overview of the amount of data from each source.

Dolma is built using a pipeline of (1) language filtering, (2) quality filtering, (3) content filtering, (4) deduplication, (5) multi-source mixing, and (6) tokenization. We refer the reader to the Dolma report (Anonymous, 2024) for more details about its design principles, details about its construction, and a more detailed summary of its contents. The report provides additional analyses and experimental results from training language models on intermediate states of Dolma to share what we learned about important data curation practices, including the role of content or quality filters, deduplication, and mixing data from multiple sources. We keep documents from each source separate, both during curation as well as in the final release. We open-sourced our high-performance data curation tools; this toolkit can be used to further experiment on Dolma, reproduce our work, and enable fast and easy curation of pretraining corpora. Finally, we also open-sourced our WIMBD tool (Anonymous, 2023c) to help with dataset analysis.

## 2.3 Adaptation

Pretrained models are not always used as-is, but rather further fine-tuned to improve their performance, safety, and usability. Often models are first

Source	Type	UTF-8 bytes (GB)	Docs (millions)	Tokens (billions)
Common Crawl	web pages	9,022	3,370	2,006
The Stack	code	1,043	210	342
C4	web pages	790	364	174
Reddit	social media	339	377	80
peS2o	papers	268	38.8	57
Project Gutenberg	books	20.4	0.056	5.2
Wikipedia and Wikibooks	encyclopedic	16.2	6.2	3.7
Total		11,519	4,367	2,668

Table 2: Composition of Dolma. Tokens counts are based on the GPT-NeoX tokenizer

trained to follow instructions (Mishra et al., 2022; Wei et al., 2022; Sanh et al., 2022), and then further trained on human preferences (Ouyang et al., 2022) to improve the quality of their generations. We showcase the efficacy of using OLMo as a base model for further fine-tuning by training OLMo to be a general chat assistant following the TULU data and training setup (Iverson et al., 2023). This involves first performing instruction fine-tuning with a mixture of distilled and human-written instruction data and then further aligning the model with distilled preference data using Direct Preference Optimization (DPO) (Rafailov et al., 2023).

## 2.4 Evaluation

We perform base model evaluation at two stages: *online* evaluation to make decisions for model design and *offline* evaluation to evaluate model checkpoints. For the offline stage, we use the Catwalk framework (Anonymous, 2023a), a publicly available evaluation tool with access to a wide range of datasets and task formats. Using Catwalk, we perform downstream evaluation as well as intrinsic language modeling evaluation on the new perplexity benchmark, Paloma (Anonymous, 2023b).

For both downstream and perplexity evaluation, we use our fixed evaluation pipeline to compare results against publicly available models. We also report a separate evaluation of our adapted model.

**In-Loop Training Ablations** Throughout model training, we perform downstream evaluations to make decisions around model architecture, initialization, optimizers, learning rate schedule, and data mixtures. We call this our *online* evaluation as it runs in-loop every 1000 training steps (or ~4B training tokens) and provides an early and continuous signal on the quality of the model being trained.

These evaluations rely on many of the core tasks and experiment settings used for our *offline* evaluation detailed in Section 4.1, which also mirrors the task and evaluation structure of the EleutherAI eval harness (Gao et al., 2023).

**Downstream Evaluation** Following much previous work (Brown et al., 2020; Black et al., 2022; Touvron et al., 2023a,b, *inter alia*), we report zero-shot performance on a set of downstream tasks. Our evaluation suite consists of 8 core tasks corresponding closely to the commonsense reasoning task set reported by Touvron et al. (2023a) and Touvron et al. (2023b) (see Table 3 for a list of tasks). Given the scale of the models being evaluated, such tasks were selected at the beginning of model development due to their naturalness (e.g., all can be formulated as text completion scoring tasks) and ability to provide meaningful signals throughout training (see Figure 1).

**Intrinsic Language Modeling Evaluation** To measure how OLMo-7B fits distributions of language beyond held-out training data, we use Paloma (Anonymous, 2023b), a new perplexity benchmark that includes 585 different domains of text. Domains range from nytimes.com to r/depression on Reddit and are drawn from 18 separate data sources, such as C4 (Raffel et al., 2020), in stratified samples. This allows for more equal inclusion of text domains that are under-represented in their source corpora.

We aim not just to compare OLMo-7B against other models for best performance, but also to demonstrate how it enables fuller and more controlled scientific evaluations. OLMo-7B is the largest LM with explicit decontamination for perplexity evaluation. Following the approach described in Paloma, we remove any pretraining document with paragraphs leaked from Paloma evaluation data. Without decontamination, other models risk underestimating perplexity (i.e., overestimating the model’s out-of-sample fit). We also release intermediate checkpoints, allowing richer comparisons with two other models that release checkpoints, Pythia-6.9B (Biderman et al., 2023) and RPJ-INCITE-7B (Together Computer, 2023) (see Figure 2).

**Adaptation Evaluation** We also evaluate OLMo after instruction fine-tuning and DPO training using the TULU evaluation suite proposed in Wang et al. (2023); Ivison et al. (2023). We focus on eval-

uations around model chat capabilities and safety in order to showcase the efficacy of using OLMo as a base for further fine-tuning.

### 3 Training OLMo

This section describes our pretraining setup, including our distributed training framework (Section 3.1), optimizer settings (Section 3.2), data preparation (Section 3.3), and hardware (Section 3.4).

#### 3.1 Distributed Training Framework

We train our models using the *ZeRO* optimizer strategy (Rajbhandari et al., 2019) via PyTorch’s FSDP framework (Zhao et al., 2023), which reduces memory consumption by sharding the model weights and their corresponding optimizer state across GPUs. At the 7B scale, this enables training with a micro-batch size of 4096 tokens per GPU on our hardware (see Section 3.4). For OLMo-1B and -7B models, we use a constant global batch size of approximately 4M tokens (2048 instances, each with a sequence length of 2048 tokens). For OLMo-65B model (currently training), we use a batch size warmup that starts at approximately 2M tokens (1024 instances), then doubles every 100B tokens until reaching approximately 16M tokens (8192 instances).

To improve throughput, we employ mixed-precision training (Micikevicius et al., 2017) through FSDP’s built-in settings and PyTorch’s amp module. The latter ensures that certain operations like the softmax always run in full precision to improve stability, while all other operations run in half-precision with the `bf16` format. Under our specific settings, the sharded model weights and optimizer state local to each GPU are kept in full precision. The weights within each transformer block are only cast to `bf16` when the full-sized parameters are materialized on each GPU during the forward and backward passes. Gradients are reduced across GPUs in full precision.

#### 3.2 Optimizer

We use the AdamW optimizer (Loshchilov and Hutter, 2019) with the hyperparameters shown in Table 1. For all model sizes, we warm up the learning rate over 5000 steps (~21B tokens) and then decay it linearly from there down to a tenth of the peak learning rate over the remainder of training. After the warm-up period, we clip gradients such that



the total  $\ell^2$ -norm of the parameter gradients<sup>2</sup> does not exceed 1.0. Table 5 gives a comparison of our optimizer settings at the 7B scale to those of other recent LMs that also used AdamW.

### 3.3 Data

We built our training dataset out of a 2T-token sample from our open dataset, Dolma (Anonymous, 2024), which we describe in Section 2.2. The tokens from every document are concatenated together after appending a special EOS token to the end of each document, and then we group consecutive chunks of 2048 tokens to form training instances. The training instances are shuffled in the exact same way for each training run. The data order and exact composition of each training batch can be reconstructed from the artifacts we release.

All of our released models have been trained to at least 2T tokens (a single epoch over our training data), and some have been trained beyond that by starting a second epoch over the data with a different shuffling order. The impact of repeating this small amount of data should be negligible according to prior work (Muennighoff et al., 2023).

### 3.4 Hardware

In order to verify that our codebase could be used on both NVIDIA and AMD GPUs without any loss in performance, we trained models on two different clusters:

- **LUMI:** Provided by the LUMI supercomputer,<sup>3</sup> we used up to 256 nodes on this cluster, where each node consists of 4x AMD MI250X GPUs with 128GB of memory<sup>4</sup> and 800Gbps of interconnect.
- **MosaicML:** Provided by MosaicML<sup>5</sup> (Databricks), we used 27 nodes on this cluster, where each node consists of 8x NVIDIA A100 GPUs with 40GB of memory and 800Gbps interconnect.

Despite minor differences in batch size to optimize for training throughput, both runs resulted in nearly

<sup>2</sup>During gradient clipping all of the model’s parameters are treated as a single big vector (as if all parameters were flattened and concatenated together), and we take the  $\ell_2$ -norm over the corresponding single gradient vector. This is the standard way to clip gradients in PyTorch.

<sup>3</sup><https://www.lumi-supercomputer.eu>

<sup>4</sup>The MI250X is a dual-chip module, meaning in practice that each physical device consists of two logical devices, so each node has 8 logical GPU devices with 64GB of memory each.

<sup>5</sup><https://www.mosaicml.com>

identical performance on our evaluation suite by 2T tokens.

## 4 Results

The checkpoint used for evaluating OLMo-7B is trained until 2.46T tokens on the Dolma (Anonymous, 2024) dataset with a linear learning rate decay schedule mentioned in Section 3.2. In our experiments, we find that tuning this checkpoint further on the Dolma dataset for 1000 steps with the learning rate linearly decayed to 0 boosts model performance on perplexity and end-task evaluation suites described in Section 2.4. We compare OLMo with other publicly available models including LLaMA-7B (Touvron et al., 2023a), Llama-2-7B (Touvron et al., 2023b), MPT-7B (MosaicML NLP Team, 2023), Pythia-6.9B (Biderman et al., 2023), Falcon-7B (Almazrouei et al., 2023) and RPJ-INCITE-7B (Together Computer, 2023).

### 4.1 Downstream evaluation

**Setup** Our core **downstream evaluation suite** (see Table 3) consists of: arc (both arc\_easy and arc\_challenge) (Clark et al., 2018), boolq (Clark et al., 2019), openbookqa (Mihaylov et al., 2018), sciq (Welbl et al., 2017), hellaswag (Zellers et al., 2019), piqa (Bisk et al., 2020), and winogrande (Sakaguchi et al., 2021). In Appendix C, we also report results on an additional set of auxiliary tasks outside of our core evaluation set that we found to have less stable performance trends (see Figure 4).

In all cases, we perform zero-shot evaluation using the rank classification approach popularized by Brown et al. (2020). Under this approach, candidate text completions (e.g., different multiple-choice options) are ranked by likelihood (usually normalized by some normalization factor), and prediction accuracy is reported. While Catwalk implements several common likelihood normalization strategies, including normalizing by number of tokens (per-token normalization) (Brown et al., 2020; Liang et al., 2022), by number of characters (per-character normalization) (Gao et al., 2023), as well as incorporating an answer’s unconditional likelihood (Brown et al., 2020), we selected the normalization strategies for each dataset separately. Specifically, we used unconditional normalization for arc and openbookqa, per-token normalization for hellaswag, piqa, and winogrande and no normalization for boolq, and sciq (i.e., tasks formulated as single token prediction tasks).

7B Models	arc challenge	arc easy	boolq	hella- swag	open bookqa	piqa	sciq	wino- grande	avg.
<b>Falcon</b>	47.5	70.4	74.6	75.9	53.0	78.5	93.9	68.9	70.3
<b>LLaMA</b>	44.5	67.9	75.4	76.2	51.2	77.2	93.9	70.5	69.6
<b>Llama 2</b>	48.5	69.5	80.2	76.8	48.4	76.7	94.5	69.4	70.5
<b>MPT</b>	46.5	70.5	74.2	77.6	48.6	77.3	93.7	69.9	69.8
<b>Pythia</b>	44.1	61.9	61.1	63.8	45.0	75.1	91.1	62.0	63.0
<b>RPJ-INCITE</b>	42.8	68.4	68.6	70.3	49.4	76.0	92.9	64.7	66.6
<b>OLMo-7B</b>	48.5	65.4	73.4	76.4	50.4	78.4	93.8	67.9	69.3

Table 3: Zero-shot evaluation of OLMo-7B and 6 other publicly available comparable model checkpoints on 8 core tasks from the downstream evaluation suite described in Section 2.4. For OLMo-7B, we report results for the 2.46T token checkpoint.

**Results** Table 3 summarizes the result of zero-shot evaluation of OLMo-7B and compares it against 6 other publicly available models of comparable size. We report results on 8 core tasks from our evaluation suite described in Section 2.4. On aggregate, OLMo-7B is competitive against all 6 publicly available model checkpoints in our comparison table.

In Figure 1 we plot the accuracy score progression of 8 core end-tasks. All tasks, except OBQA, show an upward trend in accuracy numbers as OLMo-7B is trained on more tokens. A sharp upward tick in accuracy of many tasks between the last and the second to last step shows us the benefit of linearly reducing the LR to 0 over the final 1000 training steps. See Table 7 in Appendix C for additional evaluation results and discussion.

## 4.2 Intrinsic language modeling evaluation

**Setup** For intrinsic evaluations, Paloma proposes a range of analyses, from inspection of performance in each domain separately to more summarized results over combinations of domains. We report results at two levels of granularity: the aggregate performance over 11 of the 18 sources in Paloma as in (Anonymous, 2023b), as well as more fine-grained results over each of these sources individually. This particular subset of 11 sources from Paloma excludes sources that are not publicly available, involve fringe or toxic text, or consist of code data not supported by Paloma’s decontamination approach. This leaves C4 (Raffel et al., 2020), mC4-en (Chung et al., 2023), Wikitext 103 (Merity et al., 2016), Penn Treebank (Marcus et al., 1999; Nunes, 2020), RedPajama (Together Computer, 2023), Falcon-RefinedWeb (Penedo et al., 2023), Dolma (Anonymous, 2024), M2D2 S2ORC (Reid

et al., 2022), M2D2 Wikipedia (Reid et al., 2022), C4 100 domains (Chronopoulou et al., 2022), and Dolma 100 Subreddits (Anonymous, 2024). To allow for a fair comparison between models with different vocabularies, we report bits per byte as defined by Gao et al. (2020) over the test sets of these sources.

**Results** In the *Sources Combined* subplot of Figure 2, we show the performance of OLMo-7B against 6 comparably-sized language models on the combination of 11 data sources from Paloma. Overall we find OLMo to have a competitive fit, especially given its training data was explicitly decontaminated against Paloma. As seen through the comparison of final models (see shapes) as well intermediate checkpoints (see dashed lines), the OLMo results follow similar scaling trends of other models. Note that the performance of intermediate checkpoints is influenced by where that checkpoint occurs in the learning rate schedule. So models trained for fewer steps will tend to have steeper training curves without necessarily being more sample efficient if training duration were fixed across all models. MPT-7B, nevertheless, stands out as improving ahead of the other models in this subplot. This could be due to a number of factors, including pretraining data composition and its match to the domains in Paloma (e.g., MPT trains on 27% non-Common Crawl data rather than 18% for LLaMA, 12.2% for RedPajama, and 11.2% for OLMo) as well as various data preprocessing decisions (e.g., MPT’s use of semantic deduplication by Abbas et al., 2023, on C4).

The remaining subplots in Figure 2 provide more fine-grained analysis by reporting bits per byte separately for each of the 11 data sources that are combined in the aggregated Paloma metric. From

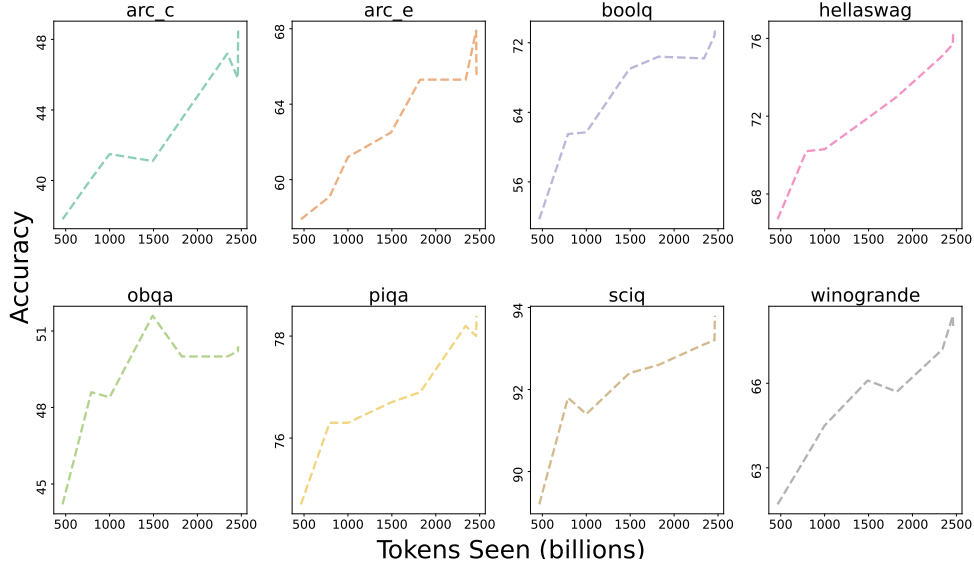


Figure 1: Accuracy score progression of OLMo-7B on 8 core end-tasks score from Catwalk evaluation suite described in Section 2.4. We can see the benefit of decaying LR to 0 in the final 1000 steps of training on most tasks.

this we see greater variation in sample efficiency, largely driven by the similarity of training and evaluation distributions. Notably, OLMo-7B fares well on evaluations predominated by Common Crawl, such as C4, though different ways of postprocessing Common Crawl are best fit by models trained with that specific data, such as Falcon-7B on Falcon RefinedWeb. Meanwhile, OLMo-7B is less sample efficient compared to other models on sources less related to scraped web text, such as WikiText-103, M2D2 S2ORC, and M2D2 Wikipedia. The RedPajama evaluation shows a similar pattern, perhaps as only 2 of its 7 domains are from Common Crawl, and Paloma weights domains within each source equally. Since heterogeneous data from curated sources like Wikipedia and ArXiv papers is much less abundant than scraped web text, maintaining sample efficiency for fit to these distributions of language will be challenging as pretraining corpora are scaled.

### 4.3 Adaptation Evaluation

**Setup** We evaluate OLMo before adaptation, and after both the supervised fine-tuning and DPO training stage, focusing on the safety and chat evaluations used by Wang et al. (2023). We additionally compare to officially released instruction-tuned variants of the models from Table 3. We finally also compare to TüLU 2 models to compare against models trained using the same post-training data

Model	MMLU 0-shot ↑	AlpacaEval % win ↑	ToxiGen % Toxic ↓	TruthfulQA % Info+True ↑
OLMo (base)	28.3	-	81.4	31.6
MPT Chat	33.8	46.8	0.1	42.7
Falcon Instruct	25.2	14.0	70.7	27.2
RPJ-INCITE Chat	27.0	38.0	46.4	53.0
Llama-2-Chat	46.8	87.3	0.0	26.3
TüLU 2	50.4	73.9	7.0	51.7
TüLU 2+DPO	50.7	85.1	0.5	- <sup>6</sup>
OLMo +SFT	47.3	57.0	14.4	41.2
OLMo +SFT+DPO	46.2	69.3	1.7	52.0

Table 4: Evaluation of various instruction-tuned 7B models, including OLMo-7B and before and after adaptation training. Lower is better for ToxiGen and higher is better for other metrics. We provide a detailed description of models and metrics in Appendix. E.

mixes and procedures.

**Results** We find that instruction tuning considerably improves the performance and safety of OLMo, increasing MMLU performance by a wide margin and improving ToxiGen and TruthfulQA scores - especially after DPO training. Additionally, we find that OLMo outperforms most other chat variants after both initial instruction tuning (OLMo +SFT) and additional preference alignment (OLMo +SFT+DPO), highlighting both the strength of OLMo as a base model and the strength of the TüLU mix used to perform adaptation training. However, we find there is still a gap with TüLU 2, which is trained by applying the TüLU

<sup>6</sup>Following Ivison et al. (2023), we do not report TüLU 2 TruthfulQA scores due to test set contamination.

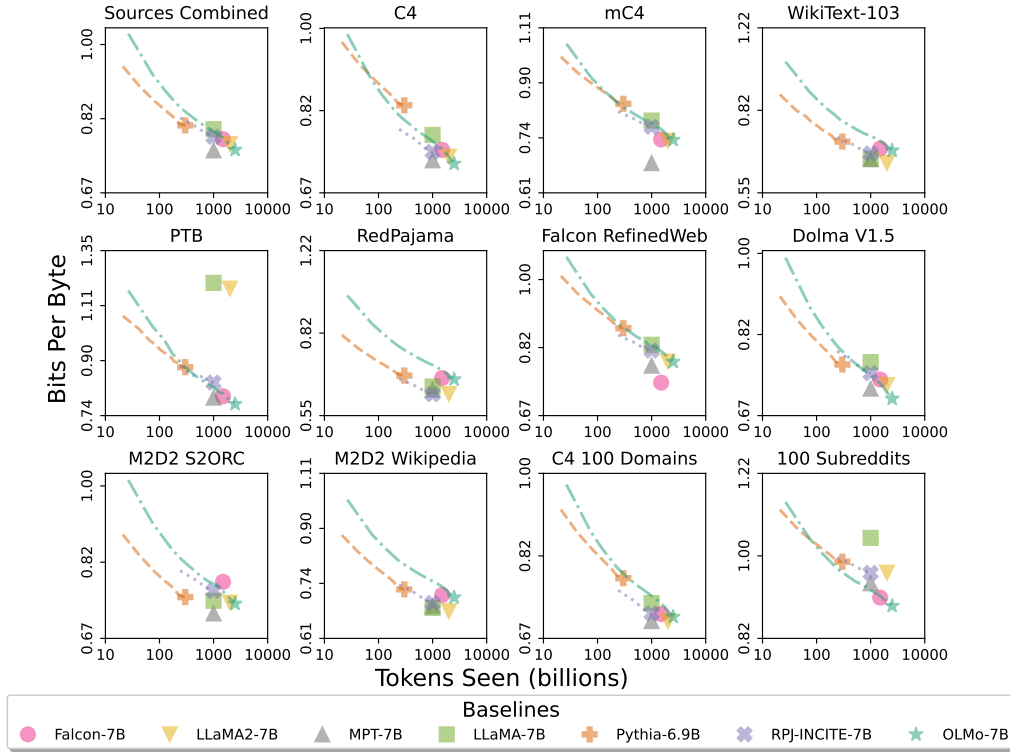


Figure 2: Bits per byte on 11 evaluation data sources from Paloma and their combination (Anonymous, 2023b), decontaminated from OLMo’s pretraining data. While models follow a general data scaling trend, sample efficiency is most favorable on in-distribution data. For example, OLMo-7B overtakes all other models on C4, perhaps from having 88.8% Common Crawl pretraining data.

mix on Llama 2. This gap may be due to test set contamination in Llama 2<sup>7</sup> and because the TULU mix was primarily designed for Llama models - we will investigate the cause of this gap in future work. Overall, we see that OLMo greatly benefits from additional tuning and serves as a strong base model for downstream applications.

## 5 Artifacts Released

By sharing artifacts from all pipeline stages, we aim to encourage open research and reduce duplicated, often costly efforts, by academics and practitioners. We release the following:

1. The training and modeling code.
2. The trained model weights for the 7B model, 7B-twin-2T, and the 1B model. For all the models, we release not only the final model weights but also 500+ intermediate checkpoints at intervals of 1000 steps.
3. The evaluation code and framework.
4. The complete set of metrics logged to Weights & Biases during training.

<sup>7</sup>Touvron et al. (2023b) report that Llama 2 was pretrained on data contaminated with MMLU test data.

5. Training logs, ablations, and findings.

6. Adapted OLMo including its training and evaluation code and data.

## 6 Conclusion and Future Work

This paper presents our first release of OLMo, a state-of-the-art, truly open language model and its framework to build and study the science of language modeling. Unlike most prior efforts that have only released model weights and inference code, we release OLMo and the whole framework, including training data and training and evaluation code. Soon, we will also release training logs, ablations, findings and Weights & Biases logs. We will additionally release the adapted models as well as all of our model adaptation code and data.

We intend to continuously support and extend OLMo and its framework, and continue to push the boundaries of open LMs to empower the open research community. To that end, we look forward to bringing different model sizes, modalities, datasets, safety measures, and evaluations into the OLMo family. We hope this and future releases will empower and strengthen the open research community and inspire a new wave of innovation.



## Limitations

We recognize building a large language model has many limitations. In fact, each step of the process of creating a language model, from the data to training to adaptation to evaluation each have their own limitations, and so we’ve added sections for each below. Of course we recognize that AI systems today can have broad societal reach, and therefore there are significant limitations beyond what we are able to fit into this section.

**Data** Our work focuses on pretraining data in English. We hope that our open framework enables the development of future models in more languages as well as multi-lingual models. The data that models are trained on is what gives models their capabilities, and at the scale of training a large language model we recognize that the data likely contains problematic content like toxic language, personal information, and copyrighted text. We mitigated this to the best of our ability but recognize there are no perfect approaches today that can completely remove such content.

**Training** Training a large language model is currently a challenging endeavor which is missing significant support from the open source community, and with our limited page count we did not provide extensive training logs documenting, for example, training runs that diverged or failed to learn. In addition, we plan to release documentation of our ablations and hyperparameter tuning upon publication.

**Adaptation** Our pretrained models face the same issues as existing pretrained LLMs, such as bias, toxicity and, hallucinations. Our adapted models are better at avoiding these generations but they are still not perfect. Additionally, we note that we largely adopt an existing data mixture designed for a different model family (TÜLU, designed for Llama models), and OLMo may require different data mixing to adjust for its unique strengths and weaknesses. The TÜLU mix itself also relies on data distilled from a variety of models, and we hope to reduce our reliance on such data in the future.

**Evaluation** While we’ve included comparisons on a variety of datasets to other current language models, many of the downstream tasks are not actually representative of how users interact with language models (i.e., as a chatbot). In addition, language model evaluations are currently very noisy;

we aimed to include only evaluations on datasets that provided some signal as to which model performs best, but recognize that there is no perfect automatic evaluation, and thus comparisons should be taken with a grain of salt.

## Ethics Statement

Through this work, we take the position that increased openness of language models is essential for scientific understanding of their abilities and limitations and for broad participation in the continued development of such models. Training on open data further enhances these benefits. In addition, our open release enables practitioners to take our models and build on them instead of having to train their own from scratch, in which case they would be repeating our work while consuming more resources and leading to an increased environmental impact. Of course, openness is not without risk; the possibility remains that these models will be used in unintended ways that cause harm. We believe that research and development efforts to understand and mitigate those potential harms will also be accelerated by the openness of the models, allowing a diversity of approaches and analyses. Over the past year there have been a number of comparable models released with very permissive licenses, so using a more strict license for our work will not remove the overall risk in the field. We believe this tradeoff on the side of being more open is the best option.

## References

- Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S Morcos. 2023. [Semdedup: Data-efficient learning at web-scale through semantic deduplication](#). *arXiv preprint arXiv:2303.09540*.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra-Aimée Cojocaru, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. [The falcon series of open language models](#). *ArXiv*, abs/2311.16867.
- Yuvanesh Anand, Zach Nussbaum, Brandon Duderstadt, Benjamin Schmidt, and Andriy Mulyar. 2023. Gpt4all: Training an assistant-style chatbot with large scale data distillation from gpt-3.5-turbo. <https://github.com/nomic-ai/gpt4all>.
- Anonymous. 2023a. Catwalk: A unified language model evaluation framework for many datasets.

689	Anonymous. 2023b. Paloma: A benchmark for evaluating language model fit.	
690		
691	Anonymous. 2023c. What’s in my big data?	
692	Anonymous. 2024. Dolma: an Open Corpus of Three Trillion Tokens for Language Model Pretraining Research.	
693		
694		
695	Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. <a href="#">Layer normalization</a> . <i>ArXiv</i> , abs/1607.06450.	
696		
697	Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. 2022. <a href="#">Training a helpful and harmless assistant with reinforcement learning from human feedback</a> .	
698		
699		
700		
701		
702		
703		
704		
705		
706		
707		
708		
709	Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. 2003. <a href="#">A neural probabilistic language model</a> . <i>J. Mach. Learn. Res.</i> , 3:1137–1155.	
710		
711		
712	Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, Usven Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar Van Der Wal. 2023. <a href="#">Pythia: A suite for analyzing large language models across training and scaling</a> . In <i>Proceedings of the 40th International Conference on Machine Learning</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 2397–2430. PMLR.	
713		
714		
715		
716		
717		
718		
719		
720		
721		
722	BigScience, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. <i>arXiv preprint arXiv:2211.05100</i> .	
723		
724		
725		
726		
727		
728	Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. <a href="#">Piqa: Reasoning about physical commonsense in natural language</a> . In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pages 7432–7439.	
729		
730		
731		
732		
733	Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. 2022. <a href="#">GPT-NeoX-20B: An open-source autoregressive language model</a> . In <i>Proceedings of the ACL Workshop on Challenges &amp; Perspectives in Creating Large Language Models</i> .	
734		
735		
736		
737		
738		
739		
740		
741		
742	Su Lin Blodgett, Lisa Green, and Brendan O’Connor. 2016. <a href="#">Demographic dialectal variation in social media: A case study of African-American English</a> . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1119–1130, Austin, Texas. Association for Computational Linguistics.	745
743		746
744		747
		748
	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. <a href="#">Language models are few-shot learners</a> . <i>ArXiv</i> , abs/2005.14165.	749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. <a href="#">Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality</a> .	761
		762
		763
		764
		765
		766
	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. <a href="#">Palm: Scaling language modeling with pathways</a> .	767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
	Alexandra Chronopoulou, Matthew Peters, and Jesse Dodge. 2022. <a href="#">Efficient hierarchical domain adaptation for pretrained language models</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1336–1351, Seattle, United States. Association for Computational Linguistics.	790
		791
		792
		793
		794
		795
		796
		797
	Hyung Won Chung, Noah Constant, Xavier García, Adam Roberts, Yi Tay, Sharan Narang, and Orhan Firat. 2023. <a href="#">Unimax: Fairer and more effective language sampling for large-scale multilingual pretraining</a> . <i>ArXiv</i> , abs/2304.09151.	798
		799
		800
		801
		802

803	Christopher Clark, Kenton Lee, Ming-Wei Chang,	Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi,	860
804	Tom Kwiatkowski, Michael Collins, and Kristina	Maarten Sap, Dipankar Ray, and Ece Kamar. 2022.	861
805	Toutanova. 2019. Boolq: Exploring the surprising	<a href="#">TOXIGEN: Controlling Language Models to Gener-</a>	862
806	difficulty of natural yes/no questions. <i>arXiv preprint</i>	<a href="#">ate Implied and Adversarial Toxicity</a> . In <i>ACL</i> .	863
807	<i>arXiv:1905.10044</i> .		
808	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,	Dan Hendrycks, Collin Burns, Steven Basart, Andy	864
809	Ashish Sabharwal, Carissa Schoenick, and Oyvind	Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-	865
810	Tafjord. 2018. <a href="#">Think you have solved question an-</a>	hardt. 2021. Measuring massive multitask language	866
811	<a href="#">swering? try arc, the ai2 reasoning challenge</a> . <i>arXiv</i>	understanding. <i>Proceedings of the International Con-</i>	867
812	<i>preprint arXiv:1803.05457</i> .	<i>ference on Learning Representations (ICLR)</i> .	868
813	Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie,	Hamish Ivison, Yizhong Wang, Valentina Pyatkin,	869
814	Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell,	Nathan Lambert, Matthew Peters, Pradeep Dasigi,	870
815	Matei Zaharia, and Reynold Xin. 2023. <a href="#">Free dolly:</a>	Joel Jang, David Wadden, Noah A. Smith, Iz Belt-	871
816	<a href="#">Introducing the world’s first truly open instruction-</a>	agy, and Hannaneh Hajishirzi. 2023. <a href="#">Camels in a</a>	872
817	<a href="#">tuned llm</a> .	<a href="#">changing climate: Enhancing lm adaptation with tulu</a>	873
818	Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao,	2.	874
819	Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and	Albert Q Jiang, Alexandre Sablayrolles, Antoine	875
820	Maosong Sun. 2023. <a href="#">Ultrafeedback: Boosting lan-</a>	Roux, Arthur Mensch, Blanche Savary, Chris Bam-	876
821	<a href="#">guage models with high-quality feedback</a> .	ford, Devendra Singh Chaplot, Diego de las Casas,	877
822	Jesse Dodge, Taylor Prewitt, Remi Tachet Des Combes,	Emma Bou Hanna, Florian Bressand, et al. 2024.	878
823	Erika Odmark, Roy Schwartz, Emma Strubell,	<a href="#">Mixtral of experts</a> . <i>arXiv preprint arXiv:2401.04088</i> .	879
824	Alexandra Sasha Luccioni, Noah A. Smith, Nicole	Andreas Köpf, Yannic Kilcher, Dimitri von Rütte,	880
825	DeCario, and Will Buchanan. 2022. <a href="#">Measuring the</a>	Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens,	881
826	<a href="#">carbon intensity of ai in cloud instances</a> .	Abdullah Barhoum, Duc Minh Nguyen, Oliver	882
827	William B. Dolan and Chris Brockett. 2005. <a href="#">Automati-</a>	Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri,	883
828	<a href="#">cally constructing a corpus of sentential paraphrases</a> .	David Alexandrovich Glushkov, Arnav Varma Dan-	884
829	In <i>International Joint Conference on Natural Lan-</i>	tuluri, Andrew Maguire, Christoph Schuhmann, Huu	885
830	<i>guage Processing</i> .	Nguyen, and Alexander Julian Mattick. 2023. <a href="#">Ope-</a>	886
831	Leo Gao, Stella Biderman, Sid Black, Laurence Gold-	<a href="#">nassistant conversations - democratizing large lan-</a>	887
832	ing, Travis Hoppe, Charles Foster, Jason Phang, Ho-	<a href="#">guage model alignment</a> . In <i>Thirty-seventh Con-</i>	888
833	race He, Anish Thite, Noa Nabeshima, et al. 2020.	<i>ference on Neural Information Processing Systems</i>	889
834	<a href="#">The pile: An 800gb dataset of diverse text for lan-</a>	<i>Datasets and Benchmarks Track</i> .	890
835	<a href="#">guage modeling</a> . <i>arXiv preprint arXiv:2101.00027</i> .	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori,	891
836	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman,	Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and	892
837	Sid Black, Anthony DiPofi, Charles Foster, Laurence	Tatsunori B. Hashimoto. 2023. <a href="#">AlpacaEval: An au-</a>	893
838	Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li,	<a href="#">tomatic evaluator of instruction-following models</a> .	894
839	Kyle McDonell, Niklas Muennighoff, Chris Ociepa,	Github repository.	895
840	Jason Phang, Laria Reynolds, Hailey Schoelkopf,	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris	896
841	Aviya Skowron, Lintang Sutawika, Eric Tang, An-	Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian	897
842	ish Thite, Ben Wang, Kevin Wang, and Andy Zou.	Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Ku-	898
843	2023. <a href="#">A framework for few-shot language model</a>	mar, et al. 2022. <a href="#">Holistic evaluation of language</a>	899
844	<a href="#">evaluation</a> .	<a href="#">models</a> . <i>arXiv preprint arXiv:2211.09110</i> .	900
845	Sidney Greenbaum and Gerald Nelson. 1996. <a href="#">The in-</a>	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	901
846	<a href="#">ternational corpus of english (ICE) project</a> . <i>World</i>	<a href="#">Truthfulqa: Measuring how models mimic human</a>	902
847	<i>Englishes</i> , 15(1):3–15.	<a href="#">falsehoods</a> . In <i>Proceedings of the 60th Annual Meet-</i>	903
848	Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang,	<i>ing of the Association for Computational Linguistics</i>	904
849	Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng	<i>(Volume 1: Long Papers)</i> , pages 3214–3252.	905
850	Wu. 2023. How close is chatgpt to human experts?	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang,	906
851	comparison corpus, evaluation, and detection. <i>arXiv</i>	Yile Wang, and Yue Zhang. 2020. <a href="#">Logiqa: A chal-</a>	907
852	<i>preprint arxiv:2301.07597</i> .	<a href="#">lenge dataset for machine reading comprehension</a>	908
853	Suchin Gururangan, Mitchell Wortsman, Samir Yitzhak	<a href="#">with logical reasoning</a> . <i>CoRR</i> , abs/2007.08124.	909
854	Gadre, Achal Dave, Maciej Kilian, Weijia Shi,	Zhengzhong Liu, Aurick Qiao, Willie Neiswanger,	910
855	Jean Mercat, Georgios Smyrnis, Gabriel Ilharco,	Hongyi Wang, Bowen Tan, Tianhua Tao, Junbo Li,	911
856	Matt Jordan, Reinhard Heckel, Alex Dimakis, Ali	Yuqi Wang, Suqi Sun, Omkar Pangarkar, et al. 2023.	912
857	Farhadi, Vaishaal Shankar, and Ludwig Schmidt.	<a href="#">Llm360: Towards fully transparent open-source llms</a> .	913
858	2023. <a href="#">OpenLM: a minimal but performative lan-</a>	<i>arXiv preprint arXiv:2312.06550</i> .	914
859	<a href="#">guage modeling (lm) repository</a> . GitHub repository.		



915	Ilya Loshchilov and Frank Hutter. 2019. <a href="#">Decoupled weight decay regularization</a> . In <i>International Conference on Learning Representations</i> .	968
916		969
917		970
918	Alexandra Sasha Luccioni, Sylvain Viguier, and Anne-Laure Ligozat. 2022. <a href="#">Estimating the carbon footprint of bloom, a 176b parameter language model</a> .	971
919		972
920		973
921	Mitchell P. Marcus, Beatrice Santorini, Mary Ann Marcinkiewicz, and Ann Taylor. 1999. <a href="#">Treebank-3</a> .	974
922		975
923		976
924	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. <a href="#">Pointer sentinel mixture models</a> . <i>ArXiv</i> , abs/1609.07843.	977
925		978
926		979
927	Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Frederick Diamos, Erich Elsen, David García, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. 2017. <a href="#">Mixed precision training</a> . <i>ArXiv</i> , abs/1710.03740.	980
928		981
929		982
930		983
931		984
932	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. <a href="#">Can a suit of armor conduct electricity? a new dataset for open book question answering</a> . <i>arXiv preprint arXiv:1809.02789</i> .	985
933		986
934		987
935		988
936	Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. <a href="#">Distributed representations of words and phrases and their compositionality</a> . In <i>Neural Information Processing Systems</i> .	989
937		990
938		991
939		992
940	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. <a href="#">Cross-task generalization via natural language crowdsourcing instructions</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3470–3487, Dublin, Ireland. Association for Computational Linguistics.	993
941		994
942		995
943		996
944		997
945		998
946		999
947	MosaicML NLP Team. 2023. <a href="#">Introducing mpt-7b: A new standard for open-source, commercially usable llms</a> . Accessed: 2023-05-05.	1000
948		1001
949		1002
950	Niklas Muennighoff, Alexander M Rush, Boaz Barak, Teven Le Scao, Aleksandra Piktus, Nouamane Tazi, Sampo Pyysalo, Thomas Wolf, and Colin Raffel. 2023. <a href="#">Scaling data-constrained language models</a> . <i>arXiv preprint arXiv:2305.16264</i> .	1003
951		1004
952		1005
953		1006
954		1007
955	Daivde Nunes. 2020. <a href="#">Preprocessed penn tree bank</a> .	1008
956		1009
957		1010
958	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> . <i>ArXiv</i> , abs/2303.08774.	1011
959		1012
960		1013
961		1014
962		1015
963		1016
964		1017
965		1018
966		1019
967		1020
		1021
		1022
		1023
		1024
		1025
	Antonios Papasavva, Savvas Zannettou, Emiliano De Cristofaro, Gianluca Stringhini, and Jeremy Blackburn. 2020. <a href="#">Raiders of the lost kek: 3.5 years of augmented 4chan posts from the politically incorrect board</a> . <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , 14:885–894.	
	David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel Rothchild, David So, Maud Texier, and Jeff Dean. 2021. <a href="#">Carbon emissions and large neural network training</a> .	
	Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra-Aimée Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. <a href="#">The refined-web dataset for falcon llm: Outperforming curated corpora with web data, and web data only</a> . <i>ArXiv</i> , abs/2306.01116.	
	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. <a href="#">Deep contextualized word representations</a> . <i>ArXiv</i> , abs/1802.05365.	
	Mohammad Taher Pilehvar and José Camacho-Collados. 2018. <a href="#">Wic: 10, 000 example pairs for evaluating context-sensitive representations</a> . <i>CoRR</i> , abs/1808.09121.	
	Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Buden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsim-poukelli, Nikolai Grigorev, Doug Fritz, Thibault Sotiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorrayne Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. 2022. <a href="#">Scaling language models: Methods, analysis &amp; insights from training gopher</a> .	
	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. <a href="#">Direct preference optimization: Your language model is secretly a reward model</a> . In <i>Thirty-seventh</i>	





1139	<i>the Association for Computational Linguistics</i> , pages	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	1195
1140	960–966, Florence, Italy. Association for Computa-	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	1196
1141	tional Linguistics.	wan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mi-	1197
1142	Alex Wang, Amanpreet Singh, Julian Michael, Felix	haylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel	1198
1143	Hill, Omer Levy, and Samuel R. Bowman. 2018.	Simig, Punit Singh Koura, Anjali Sridhar, Tianlu	1199
1144	<a href="#">Glue: A multi-task benchmark and analysis plat-</a>	Wang, and Luke Zettlemoyer. 2022. <a href="#">Opt: Open pre-</a>	1200
1145	<a href="#">form for natural language understanding.</a> <i>ArXiv</i> ,	<a href="#">trained transformer language models.</a>	1201
1146	abs/1804.07461.		
1147	Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack	Yanli Zhao, Andrew Gu, Rohan Varma, Liangchen Luo,	1202
1148	Hessel, Tushar Khot, Khyathi Raghavi Chandu,	Chien chin Huang, Min Xu, Less Wright, Hamid	1203
1149	David Wadden, Kelsey MacMillan, Noah A. Smith,	Shojanazeri, Myle Ott, Sam Shleifer, Alban Desmai-	1204
1150	Iz Beltagy, and Hannaneh Hajishirzi. 2023. <a href="#">How</a>	son, Can Balioglu, Bernard Nguyen, Geeta Chauhan,	1205
1151	<a href="#">far can camels go? exploring the state of instruction</a>	Yuchen Hao, and Shen Li. 2023. <a href="#">Pytorch fsdp: Expe-</a>	1206
1152	<a href="#">tuning on open resources.</a>	<a href="#">riences on scaling fully sharded data parallel.</a> <i>Proc.</i>	1207
		<i>VLDB Endow.</i> , 16:3848–3860.	1208
1153	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,		
1154	Adams Wei Yu, Brian Lester, Nan Du, Andrew M.		
1155	Dai, and Quoc V Le. 2022. <a href="#">Finetuned language mod-</a>		
1156	<a href="#">els are zero-shot learners.</a> In <i>International Confer-</i>		
1157	<i>ence on Learning Representations.</i>		
1158	Johannes Welbl, Nelson F Liu, and Matt Gardner. 2017.		
1159	<a href="#">Crowdsourcing multiple choice science questions.</a>		
1160	<i>arXiv preprint arXiv:1707.06209.</i>		
1161	Carole-Jean Wu, Ramya Raghavendra, Udit Gupta,		
1162	Bilge Acun, Newsha Ardalani, Kiwan Maeng, Glo-		
1163	ria Chang, Fiona Aga Behram, James Huang,		
1164	Charles Bai, Michael Gschwind, Anurag Gupta,		
1165	Myle Ott, Anastasia Melnikov, Salvatore Candido,		
1166	David Brooks, Geeta Chauhan, Benjamin Lee, Hsien-		
1167	Hsin S. Lee, Bugra Akyildiz, Maximilian Balandat,		
1168	Joe Spisak, Ravi Jain, Mike Rabbat, and Kim Hazel-		
1169	wood. 2022. <a href="#">Sustainable ai: Environmental implica-</a>		
1170	<a href="#">tions, challenges and opportunities.</a>		
1171	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,		
1172	Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei		
1173	Lin, and Daxin Jiang. 2024. <a href="#">WizardLM: Empow-</a>		
1174	<a href="#">ering large pre-trained language models to follow</a>		
1175	<a href="#">complex instructions.</a> In <i>The Twelfth International</i>		
1176	<i>Conference on Learning Representations.</i>		
1177	Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley.		
1178	2023. Baize: An open-source chat model with		
1179	parameter-efficient tuning on self-chat data. <i>arXiv</i>		
1180	<i>preprint arXiv:2304.01196.</i>		
1181	Savvas Zannettou, Barry Bradlyn, Emiliano De Cristo-		
1182	faro, Haewoon Kwak, Michael Sirivianos, Gianluca		
1183	Stringini, and Jeremy Blackburn. 2018. <a href="#">What is gab:</a>		
1184	<a href="#">A bastion of free speech or an alt-right echo chamber.</a>		
1185	In <i>Companion Proceedings of the The Web Confer-</i>		
1186	<i>ence 2018</i> , WWW ’18, page 1007–1014, Republic		
1187	and Canton of Geneva, CHE. International World		
1188	Wide Web Conferences Steering Committee.		
1189	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali		
1190	Farhadi, and Yejin Choi. 2019. <a href="#">Hellaswag: Can a</a>		
1191	<a href="#">machine really finish your sentence?</a> <i>arXiv preprint</i>		
1192	<i>arXiv:1905.07830.</i>		
1193	Biao Zhang and Rico Sennrich. 2019. <a href="#">Root mean square</a>		
1194	<a href="#">layer normalization.</a> <i>ArXiv</i> , abs/1910.07467.		

## A Training Settings

Table 5 summarizes the model architecture and the optimizer parameters of OLMo-7B as well as recent similar-sized models.

## B Power Consumption and Carbon Footprint

Following previous literature (Strubell et al., 2019; Patterson et al., 2021; Wu et al., 2022; Dodge et al., 2022), we estimate the total energy consumed and carbon released while pretraining our models by calculating the total power consumption required for training, and then multiplying it by the carbon emission intensity of the power grid where the model was trained. While reporting these operational emissions is standard practice, it does not account for other sources of emissions such as the embodied emissions due to the manufacturing, transportation and disposal of hardware and datacenter infrastructure, lifetime operational emissions due to use, rebound effects, or other environmental impacts such as water consumption or mining. Thus our estimates should be viewed as lower bounds.

We calculate the total power consumption for our models by measuring the power consumption of a single node every 25ms, calculating an average across the entire training run, and multiplying by the total number of nodes. We then account for the energy efficiency of the data center by multiplying the previous total by a power usage effectiveness (PUE) factor, which we set to 1.1, representing a conservative 10% energy consumption overhead typical of energy efficient datacenters.<sup>89</sup> We estimate that pretraining our 7B models consumed **239 MWh** of energy.

To calculate carbon emissions, we multiply the total power consumption by a carbon intensity factor, measured in kg CO<sub>2</sub> emitted per KWh, based on the physical location of the data center where each model was trained. The model trained on A100-40GB GPUs was trained in Australia, so we assume a carbon intensity factor of 0.610, the national average for Australia in 2022.<sup>10</sup> The model trained on MI250X GPUs was trained in the LUMI

supercomputer, which runs on 100% renewable, carbon-neutral energy, so we assume a carbon intensity factor of 0. LUMI is powered entirely by hydroelectric power and some sources (Ubierna et al., 2022) measure the carbon intensity factor of hydroelectric power to be 0.024, which would imply total carbon emissions of 3.54 tCO<sub>2</sub>eq.<sup>11</sup> However, we rely on the official LUMI data for our calculations, and thus we estimate total pretraining emissions of **69.78 tCO<sub>2</sub>eq.**<sup>12</sup> In Table 6 we compare our models with other previously released models based on publicly available information.

We hope that openly releasing our models can reduce future emissions by allowing others to avoid the need to pretrain models from scratch, and give insights into the true cost of developing state of the art models. We also highlight that our estimates are lower bounds, because they do not include other critical pieces of development such as debugging, hyperparameter tuning, and downtime.

## C Additional Evaluation

**Additional perplexity results** In Figure 3 we provide results for each of the 7 data sources in Paloma (Anonymous, 2023b) that are excluded from the combined metric in Figure 2. Some of these sources such as Pile (Gao et al., 2020) and ICE (Greenbaum and Nelson, 1996) are not publicly available at this time. Dolma 100 Programming Languages (Anonymous, 2024) consists of code data that is not supported by the decontamination approach used in Paloma. TwitterAAE (Blodgett et al., 2016), along with ICE, are datasets for targeted analyses of disparities in performance between different dialects and as such should be evaluated separately. And finally, the Manosphere, Gab, and 4chan corpora (Ribeiro et al., 2021; Zannettou et al., 2018; Papasavva et al., 2020) are intended to examine model fit to language from fringe online communities that are studied for prevalent hate speech and toxicity. Thus minimizing perplexity on these fringe corpora is not always desirable.

One notable result here is that OLMo-7B is much farther ahead of the other models on Dolma 100 Programming Languages (100 PLs). Note that this effect may be due in part to underestimation from contamination, as decontaminating code data is beyond the scope of the method in Paloma. At the

<sup>89</sup><https://www.nrel.gov/computational-science/measuring-efficiency-pue.html>

<sup>90</sup><https://www.google.com/about/datacenters/efficiency/>

<sup>10</sup><https://www.cleanenergyregulator.gov.au/Infohub/Markets/Pages/qcmr/december-quarter-2022/Emissions-Reduction.aspx>

<sup>11</sup><https://www.lumi-supercomputer.eu>

<sup>12</sup>These metrics were in part collected using Carbonara’s AI agent and monitoring platform. Learn more at: <https://trycarbonara.com>

	<b>OLMo-7B</b>	<b>LLaMA2-7B</b>	<b>OpenLM-7B</b>	<b>Falcon-7B</b>	<b>PaLM-8B</b>
Dimension	4096	4096	4096	4544	4096
Num heads	32	32	32	71	16
Num layers	32	32	32	32	32
MLP ratio	~8/3	~8/3	~8/3	4	4
Layer norm type	non-parametric	RMSNorm	parametric	parametric	parametric
Positional embeddings	RoPE	RoPE	RoPE	RoPE	RoPE
Attention variant	full	GQA	full	MQA	MQA
Biases	none	none	in LN only	in LN only	none
Block type	sequential	sequential	sequential	parallel	parallel
Activation	SwiGLU	SwiGLU	SwiGLU	GeLU	SwiGLU
Sequence length	2048	4096	2048	2048	2048
Batch size (instances)	2160	1024	2048	2304	512
Batch size (tokens)	~4M	~4M	~4M	~4M	~1M
Weight tying	no	no	no	no	yes
Warmup steps	5000	2000	2000	1000	
Peak LR	3.0E-04	3.0E-04	3.0E-04	6.0E-04	
Minimum LR	3.0E-05	3.0E-05	3.0E-05	1.2E-05	
Weight decay	0.1	0.1	0.1	0.1	
Beta1	0.9	0.9	0.9	0.99	
Beta2	0.95	0.95	0.95	0.999	
Epsilon	1.0E-05	1.0E-05	1.0E-05	1.0E-05	
LR schedule	linear	cosine	cosine	cosine	
Gradient clipping	global 1.0	global 1.0	global 1.0	global 1.0	
Gradient reduce dtype	FP32	FP32	FP32	BF16	
Optimizer state dtype	FP32	most likely FP32	FP32	FP32	

Table 5: LM architecture and optimizer comparison at the 7–8B scale. In the “layer norm type” row, “parametric” and “non-parametric” refer to the usual layer norm implementation with and without adaptive gain and bias, respectively. All models are trained using AdamW.

same time other models that are trained on code data from GitHub such as RPJ-INCITE-7B, that are just as likely to have contamination, fair much worse. Another factor then is that OLMo-7B trains on code data with exactly the same post-processing as that in 100 PLs while the code data in other models will have been processed differently. Similarly, Pile evaluation demonstrates these in-distribution and potential contamination effects as Pythia-6.9B achieves top performance despite being trained on almost an order of magnitude fewer tokens than OLMo-7B.

The results on the remaining 5 targeted sources should be interpreted with care, as Paloma often finds that perplexity on these sources is dominated by superficial features such as low average document length rather than fit to that which would actually be salient to members of these speech communities. TwitterAAE and Gab have among the shortest documents in Paloma contributing to unusually high bits per byte in this figure. Other than these two, the models are notably very closely grouped in a data scaling trend in ICE, Manosphere, and 4chan.

**Additional end-task results** Next, in Table 7, we provide results from zero-shot evaluation of

OLMo-7B on 6 additional end-tasks apart from the 8 in our core evaluation suite. These tasks are headqa\_en (Vilares and Gómez-Rodríguez, 2019), logiqa (Liu et al., 2020), mrpc (Dolan and Brockett, 2005), qnli (Wang et al., 2018), wic (Pilehvar and Camacho-Collados, 2018), and wnli (Wang et al., 2018).

We note, however, that in contrast to our core evaluation set described in Section 4.1, we found these additional end-tasks to have less stable performance during model development, and to provide a limited signal. This is illustrated in Figure 4, where we see the progress of task performance throughout training to be more random (compare with the more stable upward trends in Figure 1). While tasks such as mrpc and wic appear more stable, they offered additional difficulties related to performance being tied to random chance (e.g., wic) or the tendency of models to make spurious predictions (e.g., always predicting a single label) that either inflate or deflate performance due to dataset class imbalances (e.g., mrpc). We therefore caution against relying too heavily on these tasks when measuring model performance throughout training and comparing models.



	GPU Type	GPU Power Consumption (MWh)	Power Usage Effectiveness	Carbon Intensity (kg CO <sub>2</sub> e/KWh)	Carbon Emissions (tCO <sub>2</sub> eq)
<b>Gopher-280B</b>	TPU v3	1,066	1.08	0.330	380
<b>BLOOM-176B</b>	A100-80GB	433	1.2	0.057	30
<b>OPT-175B</b>	A100-80GB	324	1.1	0.231	82
<b>T5-11B</b>	TPU v3	77	1.12	0.545	47
<b>LLaMA-7B</b>	A100-80GB	33	1.1	0.385	14
<b>LLaMA2-7B</b>	A100-80GB	74	1.1	0.385	31
<b>OLMo-7B</b>	MI250X	135	1.1	0.000*	0*
<b>OLMo-7B</b>	A100-40GB	104	1.1	0.610	70

Table 6: CO<sub>2</sub> emissions during pretraining. We estimate the total carbon emissions for various models using publicly available data on PUE, carbon intensity of local power grid, and reported power consumption. Numbers for Gopher-280B (Rae et al., 2022), BLOOM-176B (Luccioni et al., 2022), OPT-175B (Zhang et al., 2022), T5-11B (Patterson et al., 2021), LLaMA (Touvron et al., 2023a), and LLaMA2 (Touvron et al., 2023b) are taken from their respective papers. See Section B for details on how tCO<sub>2</sub>eq was calculated.

\* LUMI runs entirely on hydroelectric power<sup>11</sup> and some estimates (Ubierna et al., 2022) measure the intensity factor of hydroelectric power to be 0.024, implying total emissions of 3.54 tCO<sub>2</sub>eq.

	headqa_en	logiqa	mrpc	qnli	wic	wnli	avg.
<b>Falcon-7B</b>	38.6	23.7	62.8	49.8	49.5	47.9	45.4
<b>LLaMA-7B</b>	38.7	19.5	68.6	50.1	49.1	52.1	46.4
<b>LLaMA2-7B</b>	39.5	26.1	69.1	49.4	49.8	45.1	46.5
<b>MPT-7B</b>	37.4	22.9	67.7	52.1	48.1	47.9	46.0
<b>Pythia-6.9B</b>	40.1	21.5	65.4	53.8	55.0	38.0	45.6
<b>RPJ-INCITE-7B</b>	36.9	27.8	58.8	53.8	48.9	57.8	47.3
<b>OLMo-7B</b>	37.3	23.4	68.4	49.1	50.2	56.3	47.5

Table 7: Zero-shot evaluation of OLMo-7B on 6 additional end-tasks apart from the 8 present in our core evaluation suite. Once again, we compare OLMo-7B to 6 other model checkpoints which are publicly available. We find that OLMo-7B outperforms the other models on aggregate taken over 6 additional end-tasks from this table, however these tasks were also found to provide limited signal during training (see Figure 4).

## D Adaptation Training Details

We use the following hyperparameters when instruction tuning OLMo. These were chosen through small pilot experiments.

- Learning Rate:  $2 \times 10^{-6}$
- Epochs: 3
- Warmup: Linear warmup for the first 3% of total training time, and then linear cooldown to a learning rate of 0 over the remaining steps.
- Weight Decay: 0
- Gradient clipping: 0
- Maximum sequence length: 2048

After instruction finetuning, we then use the following hyperparameters for DPO training, following Ivison et al. (2023):

- Learning Rate:  $5 \times 10^{-7}$

- $\beta$ : 0.1
- Epochs: 3
- Warmup: Linear warmup for the first 10% of total training time, and then linear cooldown to a learning rate of 0 over the remaining steps.
- Weight Decay: 0
- Gradient clipping: 0
- Maximum sequence length: 2048

## E Adaptation Evaluation and Model details

We choose the models in Table 4 by choosing the ‘canonical’ best versions (that is, the best instruction-tuned or otherwise adapted models released by the same organisation) of the base models

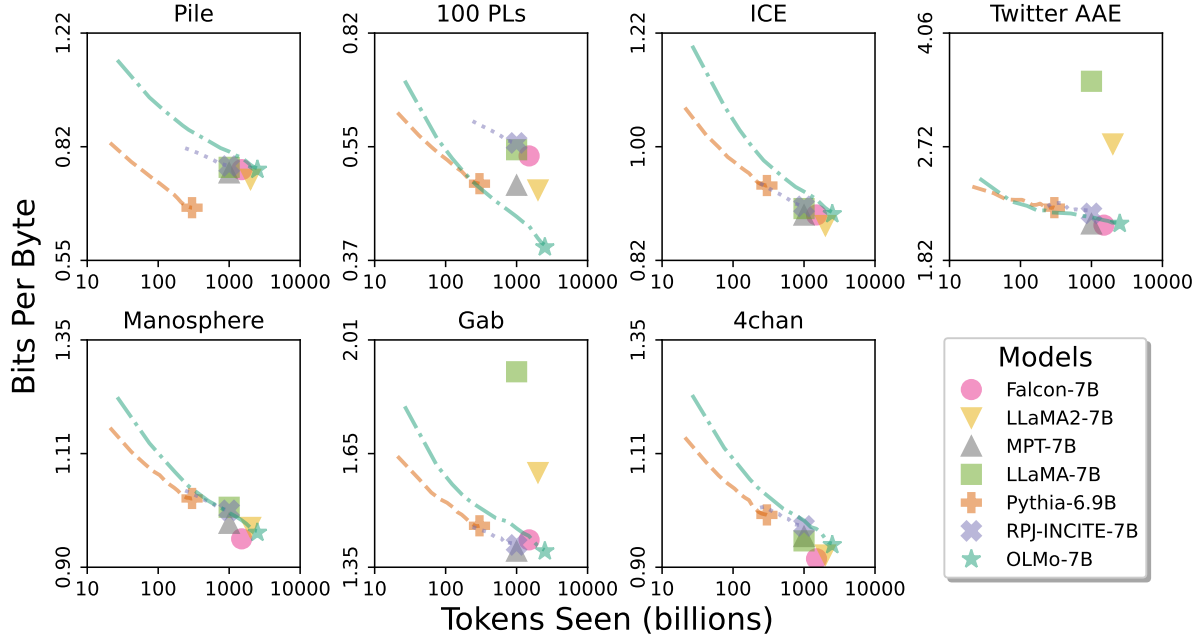


Figure 3: Bits per byte for each of the 7 remaining Paloma data sources not aggregated in Figure 2.

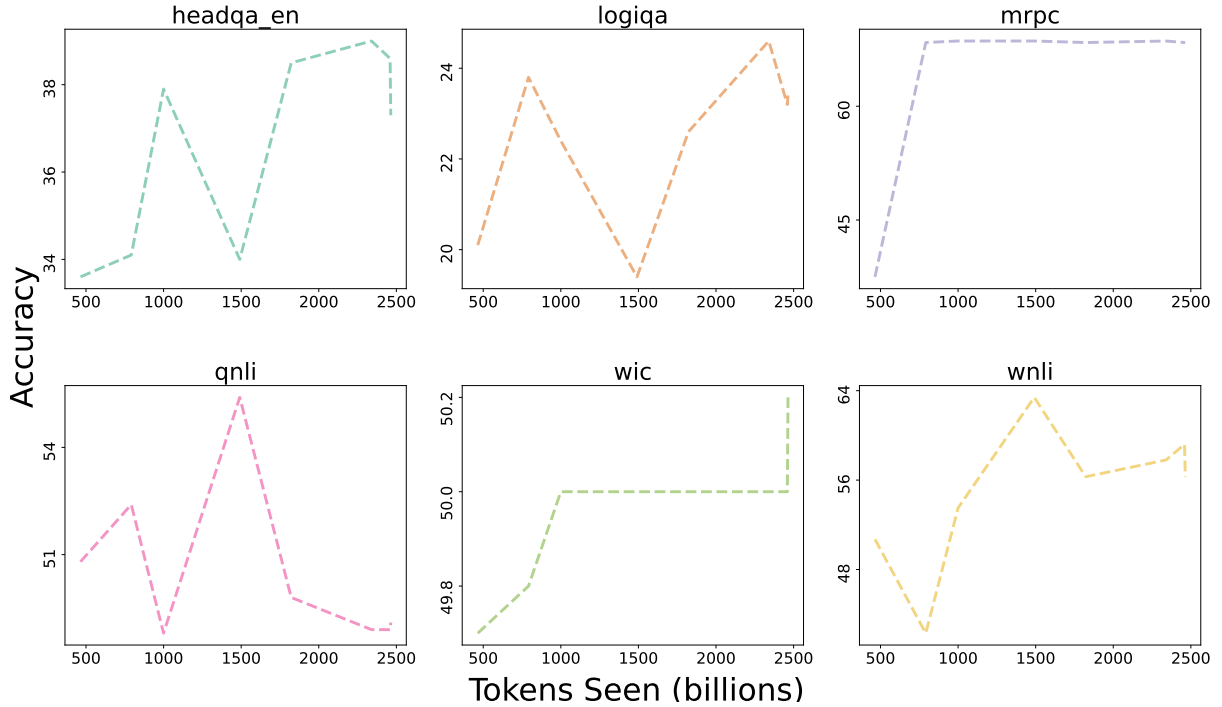


Figure 4: Accuracy score progression of OLMo-7B on 6 additional end-tasks. The performance of these additional end-tasks was unstable and provided limited signal during model development.

we compare against in Table 3. We additionally compare to TüLU 2 to show the current best models trained using the TüLU mix used to finetune OLMo. We display evaluations on MMLU, AlpacaEval, ToxiGen, and Truthfulness to focus on displaying how instruction tuning can generally

help capabilities (MMLU), how the models perform in an open-ended chat setting (AlpacaEval), and to test how instruction tuning aids in model safety and truthfulness (AlpacaEval, ToxiGen). We additionally report OLMo’s performance over the entire TüLU evaluation suite in Table 8.

Model	MMLU 0-shot	GSM8k 8-shot CoT	BBH 3-shot CoT	TyDiQA 1-shot	Codex-Eval Pass@10	AlpacaEval % win	ToxiGen % Toxic	TruthfulQA % Info + True
OLMo-7B	28.3	8.5	31.7	32.3	21.4	-	81.4	31.6
+SFT	47.3	15.5	36.9	35.2	28.6	57.0	14.4	41.2
+SFT+DPO	46.1	11.0	35.8	21.7	27.8	69.3	1.7	52.0

Table 8: Evaluation of OLMo-7B models before and after instruction finetuning and DPO training on the full TüLU evaluation suite. Lower is better for ToxiGen and higher is better for other metrics.

We provide a brief description of each model evaluated in Table 4 below. For all models, we use the provided chat template for prompt formatting when available.

- **MPT Chat:** A version of MPT 7B finetuned on the ShareGPT-Vicuna (Chiang et al., 2023), HC3 (Guo et al., 2023), Alpaca (Taori et al., 2023), HH-RLHF (Bai et al., 2022), and Evol-Instruct (Xu et al., 2024) datasets. Retrieved from <https://huggingface.co/mosaicml/mpt-7b-chat>.
- **Falcon Instruct:** A version of Falcon 7B finetuned on the Baize (Xu et al., 2023), GPT4All (Anand et al., 2023), GPTEacher (Teknum1, 2023), and Refined-Web English (Penedo et al., 2023) datasets. Retrieved from <https://huggingface.co/tiiuae/falcon-7b-instruct>.
- **RPJ-INCITE Chat:** A version of RPJ-INCITE 7B finetuned on the OASST1 (Köpf et al., 2023) and Dolly V2 (Conover et al., 2023) datasets. Retrieved from <https://huggingface.co/togethercomputer/RedPajama-INCITE-7B-Chat>.
- **Llama-2 Chat:** A version of Llama 2 7B finetuned on a mixture of instruction datasets and further trained with RLHF. We refer the reader to Touvron et al. (2023b) for further details.
- **TüLU 2:** A version of Llama 2 7B finetuned on a mixture of instruction datasets (the TüLU 2 mix). We refer the reader to Ivison et al. (2023) for further details.
- **TüLU 2+DPO:** TüLU 2 further trained with DPO on the UltraFeedback dataset (Cui et al., 2023). We refer the reader to Ivison et al. (2023) for further details.
- **OLMo +SFT:** A version of OLMo 7B finetuned on the same data as TüLU 2.

- **OLMo +SFT+DPO:** OLMo +SFT further trained with DPO on the UltraFeedback dataset (Cui et al., 2023).

We additionally provide a brief description of each evaluation setting from Table 4:

- **MMLU:** We use the official MMLU (Hendrycks et al., 2021) evaluation script and prompts available at <https://github.com/hendrycks/test>, with modifications to allow for batch processing. We evaluate using 0 few-shot examples, following the original setup of MMLU. We report average accuracy across test examples.
- **ToxiGen:** We follow the setup in Touvron et al. (2023b), but use the original set of prompts from Hartvigsen et al. (2022), which are designed to elicit toxic generations for certain groups. We take only the prompts designed to produce toxic language (‘hateful’ prompts) and use 500 prompts per group to reduce evaluation costs. For base language models, we pass in the original ToxiGen prompts unchanged and greedily decode up to the first new line (or a maximum of 512 tokens). For instruction-tuned models, we place the prompt in the corresponding template, and ask the model to complete the prompt, until the model generates a stop token (or a maximum of 512 tokens). We pass the generated text into a roberta-large model trained to detect toxic content finetuned as part of Hartvigsen et al. (2022)<sup>12</sup>. We then report the percentage of generations deemed toxic by the classifier.
- **TruthfulQA:** Following Touvron et al. (2023b), we mainly use the generation setting of TruthfulQA (Lin et al., 2022). The TruthfulQA dataset contains 818 questions, which are used to prompt the tested model to generate answers. We use the default QA prompt format with 6 in-context QA examples. We follow the official script in their of-

<sup>12</sup>[https://huggingface.co/tomh/toxigen\\_roberta](https://huggingface.co/tomh/toxigen_roberta)

1469 ficial implementation<sup>13</sup> to do greedy decoding and  
1470 answer postprocessing. We train two LLaMA 2-  
1471 based classifiers for judging the truthfulness and  
1472 informativeness of the model response, due to the  
1473 deprecation of GPT-3 making exact replication  
1474 of the original TruthfulQA evaluation infeasible.  
1475 We find that the LLaMA 2 judges are generally  
1476 able to match the performance of the original  
1477 GPT-3-based judges used by Lin et al. (2022).  
1478 We report the rate of the responses being truth-  
1479 ful and informative (% Informative and Truthful)  
1480 following Touvron et al. (2023b). We only report  
1481 the % Informative and Truthful as our primary  
1482 metric.

- 1483 • **AlpacaEval:** We use the package provided by Li  
1484 et al. (2023), following the default setup which  
1485 asks the evaluated model to generate responses  
1486 for 805 prompts and employ GPT-4 to compare  
1487 the response with Davinci-003. We employ the  
1488 “alpaca\_eval\_gpt4” annotator. We allow the eval-  
1489 uated model to generate up to 2048 tokens, with-  
1490 out specifying special stop sequences. The re-  
1491 ported win-rate is the percentage of model gener-  
1492 ations that GPT-4 reports as being preferred over  
1493 the generations from Davinci-003.

---

<sup>13</sup><https://github.com/sylinr1/TruthfulQA/>