

# MIXTUREVITAE: OPEN WEB-SCALE PRETRAINING DATASET WITH HIGH QUALITY INSTRUCTION AND REASONING DATA BUILT FROM PERMISSIVE-FIRST TEXT SOURCES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We present **MixtureVitae**, an open-access<sup>1</sup> pretraining corpus built to minimize legal risk while providing strong downstream performance. **MixtureVitae** follows a permissive-first, risk-mitigated sourcing strategy that combines public-domain and permissively licensed text (e.g., CC-BY/Apache) with carefully justified low-risk additions (e.g., government works and EU TDM-eligible sources). **MixtureVitae** adopts a simple, single-stage pretraining recipe that integrates a large proportion of permissive synthetic instruction and reasoning data—signals typically introduced during post-training and generally scarce in permissive web corpora. We categorize all sources into a three-tier scheme that reflects varying risk levels and provide shard-level provenance metadata to enable risk-aware usage. In controlled experiments using the open-sci-ref training protocol (fixed architectures and hyperparameters; 50B and 300B token budgets across 130M–1.7B parameters), models trained on **MixtureVitae** consistently outperform other permissive datasets across a suite of standard benchmarks, and at the 1.7B-parameters/300B-tokens setting, they surpass FineWeb-Edu and approach DCLM late in training. Performance is particularly strong on MMLU and on math and code benchmarks: a 1.7B model pretrained on 300B **MixtureVitae** tokens matches or exceeds a strong 1.7B instruction-tuned baseline on GSM8K, HumanEval, and MBPP, despite using over 36× fewer tokens (300B vs. ≈11T). Supported by a thorough decontamination analysis, these results show that permissive-first data with high instruction and reasoning density, tiered by licensing and provenance-related risk, can provide a practical and risk-mitigated foundation for training capable LLMs, reducing reliance on broad web scrapes without sacrificing competitiveness.

## 1 INTRODUCTION

The proliferation of large language models (LLMs) has transformed the landscape of artificial intelligence, yet their development often relies on a legally and ethically precarious foundation. The vast majority of performant models are pretrained on massive web scrapes, indiscriminately mixing public-domain content with copyrighted materials such as books, news articles, and personal websites without explicit permission (Raffel et al., 2020; Gao et al., 2020). This practice has led to a growing number of copyright infringement lawsuits, creating significant legal uncertainty for both academic researchers and commercial developers and threatening the future of the field. At the same time, practitioners who wish to avoid this risk have few alternatives, as most high-performing pretraining mixtures rely, at least in part, on opaque or non-permissive web scrapes.

Compounding this uncertainty is the prevailing assumption that state-of-the-art performance is inextricably linked to the sheer scale and diversity offered by these legally ambiguous web scrapes. The absence of a high-performance, large-scale pretraining dataset that actively mitigates these risks

---

<sup>1</sup>Dataset, source code for experiments reproduction and pre-trained models will be revealed upon acceptance.

054 has forced a difficult choice between performance and compliance. In practice, the strongest open  
055 baselines such as FineWeb-Edu (Penedo et al., 2024) and DCLM (Li et al., 2024) still rely on mixed-  
056 license or unspecified web data, whereas strictly permissive corpora tend to lag behind them on  
057 reasoning-heavy benchmarks. This raises a critical question: Can a powerful language model be  
058 trained on a dataset that provides a more legally robust foundation?

059 To this question, we answer "yes": We introduce **MixtureVitae**, a 422-billion-token, open-access  
060 pretraining dataset constructed to minimize copyright risk while explicitly demonstrating that a  
061 reasoning- and instruction-dense, permissive-first mixture can substantially close the performance  
062 gap to leading non-permissive corpora. The core of **MixtureVitae**'s "permissive-first" data com-  
063 prise (1) text with clear and permissive licenses (e.g., CC-BY-\*, Apache 2.0), public-domain text,  
064 and copyright-exempt text such as US federal works (see Appendix J) and (2) risk-mitigated text.  
065 Following Phi-4 (Abdin et al., 2024), which shows that the addition of synthetic and web-rewrite  
066 data boosts performance, we address the scarcity of real, human-written reasoning and conversa-  
067 tional dialogue in strictly permissive sources by significantly augmenting **MixtureVitae** with tar-  
068 geted synthetic data, which is derived from permissive models and sources. We call this combination  
069 of expressly licensed and risk-mitigated methods the "**permissive-first**" approach.

070 To validate our approach, we train models with **130M, 400M, 1.3B, and 1.7B parameters** on  
071 **MixtureVitae** and compare their performance against several prominent open datasets. The re-  
072 sults first confirm that **MixtureVitae significantly outperforms all other permissively licensed**  
073 **baselines**, with the performance gap widening as the model scale increases. The more critical test,  
074 however, is against popular non-permissive datasets containing higher proportions of copyrighted or  
075 ambiguously-licensed material. In this setting, our models achieve competitive performance, and on  
076 math and code benchmarks, our 1.7B base model matches or exceeds a strong 1.7B instruction-tuned  
077 baseline despite being trained on a dramatically smaller token budget.

078 In summary, our contributions are threefold:

079 **1. Permissive-first, risk-mitigated, and performant recipe for pretraining corpora.** We present  
080 **MixtureVitae**, the first highly-performant, permissive-first, and risk-mitigated pretraining corpus  
081 that deliberately front-loads high-quality reasoning and instruction data to drive capability gains in  
082 small models. It is organized into auditable provenance tiers and constructed via a positive-inclusion  
083 pipeline, avoiding the need for retroactive filtering.

084 **2. We demonstrate that reliance on indiscriminately scraped, high-risk copyrighted data is**  
085 **not a prerequisite for training capable LLMs.** Leveraging the open-sci-ref (Nezhurina  
086 et al., 2025) protocol to ensure rigorous comparison across 130M–1.7B parameter scales, we demon-  
087 strate the value of front-loading instruction and reasoning data into pre-training. Our 422B-token,  
088 permissive-first mixture closes the gap to mixed-license baselines while providing an auditable legal  
089 provenance. Furthermore, we show that our 1.7B base model, despite a limited 300B token bud-  
090 get, is comparable across multiple reasoning benchmarks to a strong 1.7B instruction-tuned base-  
091 line—trained on roughly  $36\times$  more tokens ( $\approx 11T$ ).

092 **3. Evaluation integrity and reusable artifacts.** We perform a large-scale 13-gram decontamination  
093 analysis across all benchmarks, showing that **MixtureVitae**'s gains persist on decontaminated test  
094 sets and when removing shards responsible for most detected overlap, and we release the corpus,  
095 shard-level provenance metadata, and curation code to enable compliant, reproducible pretraining  
096 in future work.

## 098 2 DATASET

099 We adopt a permissive-first, risk-mitigated strategy, combining sources with clear permissive li-  
100 censes (e.g. CC-BY, Apache, public domain) with narrowly justified inclusions (government works,  
101 EU TDM-eligible data) and targeted synthetic data. Within this framework, the **MixtureVitae**  
102 dataset is constructed from three primary categories: curated sources for domain-specific exper-  
103 tise, diverse web data for language and general knowledge and instruction-following and reasoning  
104 datasets to enhance reasoning and task-completion abilities.  
105

106 The major categories of our corpus are visualized in Figure 1a. We provide a granular breakdown  
107 showing the token count for each component (Figure 4), the license distribution (Figure 1b), and  
synthetic data usage (Figure 2a). Specific data sources are detailed in the following subsections.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

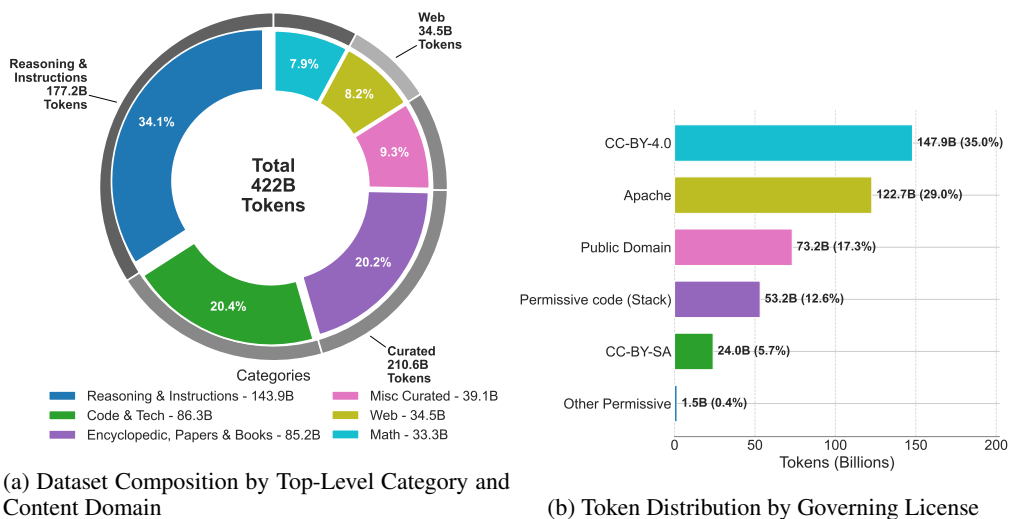


Figure 1: Composition of the **MixtureVitae** dataset (permissive-first, risk-mitigated composition).

## 2.1 DATA SOURCES

Our dataset selection process is governed by a two-layer criteria, prioritizing risk mitigation followed by quality and capability objectives:

**I. Legal & Licensing:** The primary filter is legal compliance. A dataset is considered only if it operates under a clear permissive license (e.g., CC-BY, Apache 2.0) or is in the public domain. For synthetic data, we further scrutinize the provenance of seed corpora and generator models (Appendix M). The majority of our synthetic sources satisfy full provenance transparency (classified as Tier 1), while a minority of community reasoning datasets with opaque provenance are categorized as Tier 2 to manage residual risk.

**II. Quality & Capability:** Among compliant sources, we prioritize datasets with prior evidence of high performance in community mixtures (e.g., Soldaini & Lo, 2023). Furthermore, to address the reasoning deficits typical of strictly permissive web scrapes, we target high-density instruction and reasoning data, a choice driven by the need to boost performance on tasks such as GSM8K (Cobbe et al., 2021) and MMLU(Hendrycks et al., 2021).

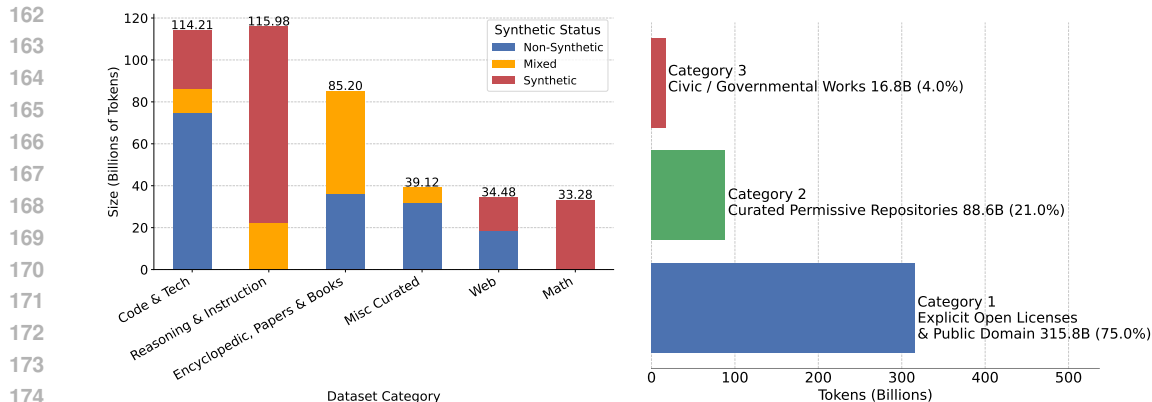
The following sections describe each of the three categories of data in **MixtureVitae**: web, curated sources, and instruction and reasoning datasets.

### 2.1.1 WEB-SCALE CORPORA

One subset of our pre-training data is derived from web-scale datasets including Nemotron-CC (Su et al., 2025), MGACorpus (Hao et al., 2025), and FineFineWeb (M-A-P et al., 2024). It also contains synthetic data generated by rephrasing web text from Nemotron-CC and MGACorpus.

### 2.1.2 CURATED DATASETS

To incorporate domain-specific knowledge and high-quality text, we curate diverse sources: public financial documents from SEC EDGAR (U.S. Securities and Exchange Commission, 2024), multi-lingual encyclopedic articles from MegaWika (Barham et al., 2023) and TxT360 (Tang et al., 2024), scientific papers from arXiv (Clement et al., 2019) and peS2o (Soldaini & Lo, 2023), medical data from PubMed (National Library of Medicine (U.S.), 1996), code from The Stack v1 (Kocetkov et al., 2023), patents from the USPTO database (United States Patent and Trademark Office, 2024) and EuroPat (Heafield et al., 2022), mathematical problems from Deepmind Math (Saxton et al., 2019), and video transcripts from both VALID (Nguyen et al., 2024) and the YouTube Commons corpus (Langlais, 2024), news and law data from the Open License Corpus (Min et al., 2024).



(a) **MixtureVitae** composition by origin (total token counts at the top in billions). Each bar represents one of the six primary content domains (see Fig. 1a), segmented by source type: **Non-Synthetic** (real human-written text and code), **Mixed** (sources with partial synthetic data), and **Synthetic** (data generated by permissive models from permissive seeds).

(b) Legal provenance and risk-mitigation tiers of the **MixtureVitae** corpus. The dataset is segmented into its three constituent legal categories, with all sources falling into a permissive-first or risk-mitigated tier. Token counts (billions) and total corpus percentages are shown for each category.

Figure 2: Composition and provenance of **MixtureVitae**: (a) Synthetic-status distribution across the six content domains, (b) licensing tiers and risk posture for the corpus.

We source 12.6% of our dataset from **The Stack v1**, a permissive-first, risk-mitigated code dataset governed by the OpenRAIL-M license. We discuss its permissiveness situation in Appendix K.

### 2.1.3 INSTRUCTION AND REASONING DATASETS

To enhance instruction-following and reasoning, we follow Abdin et al. (2024) by including considerable synthetic and web-rewrite data. We extensively use fully and partially synthetic data — all generated from permissive or public-domain seed data using models under permissive licenses.

**General Instruction Following** We include a strong instruction-following baseline with the Magpie Collection (Xu et al., 2024), its derivatives (e.g., Magpie-Phi3-Pro). This is augmented with preference data from UltraFeedback (Cui et al., 2024) and NVIDIA’s SFT data blend NVIDIA (2024), which contains a curated mixture of permissively licensed subsets from public datasets, including OASST (Köpf et al., 2023), CodeContests (Li et al., 2022), FLAN (Chung et al., 2022), OpenPlatypus (Lee et al., 2023), and the training split of GSM8K (Cobbe et al., 2021). Additionally, we augment the P3 (Sanh et al., 2022) dataset with a few-shot and multiple-choice format.

**Reasoning** To improve reasoning, we incorporate general corpora such as Glaive-AI Reasoning Dataset (Glaive AI, 2023) and OpenThoughts (Guha et al., 2025) as well as domain-specific datasets: the legal dataset CaseHOLD (Zheng et al., 2021), scientific Q&A from the OpenScience collection (NVIDIA Corporation, 2025), and agent-focused instructions from OpenManus-RL (Ulab-UIUC and MetaGPT, 2024).

**Mathematics and Coding** To strengthen quantitative reasoning, we combine our internally developed synthetic Math Word Problems dataset (Appendix I) with established datasets like Meta-MathQA (Yu et al., 2024) and DM-Math (Saxton et al., 2019), further enriched with large-scale math instruction sets, including OpenMathInstruct-2 (Toshniwal et al., 2024b), DART-MATH (Tong et al., 2024), Nemo-Math (Mahabadi et al., 2025), and Prism-Math (NVIDIA, 2025). For coding, we combine the Ling Coder collection Codefuse Team et al. (2025) with executable instructions from the StarCoder dataset Kocetkov et al. (2023) to target a wide range of software engineering tasks.

### 2.1.4 LICENSING TIERS AND RISK PROFILES

To make the provenance and legal footing of **MixtureVitae** transparent, we conceptualize all dataset components into *tiers* based on license type and expected risk profile (see Figure 2b and Table 14).<sup>2</sup>

<sup>2</sup>The high-level groupings presented in this section (e.g., “Code & Tech”, “Reasoning”) and the shard breakdowns in the Appendices are primarily organizational abstractions for visualization and provenance tracking.

**Tier 1 — Explicit Open Licenses & Public Domain.** This tier encompasses text and code under clear permissive licenses (e.g., CC0, CC-BY, Apache 2.0, MIT, BSD, a permissive subset of P3) or in the public domain, such as encyclopedic resources, scientific papers, and portions of curated math corpora. Because licenses are explicit and permissive, the legal risk of reuse is minimal. This tier also includes synthetic data generated from permissively licensed models and seed data.

**Tier 2 — Curated Permissive Corpora with Upstream Opacity.**

**(a) Permissive Corpora With Partial or Unverified Provenance.** This subset includes resources such as THE STACK V1 and Wikipedia-derived corpora. The released dataset all carries a permissive license, and curators apply filters (e.g., repository-level license heuristics). However, because provenance is only partially tracked at the file or example level, there remains some residual uncertainty about the licensing status of individual items, hence its separation from Tier 1. This Tier also includes datasets that have no license, but the underlying data is public domain or permissive and requiring the same license as the upstream data, or where the data is solely obtained synthetically from a model that is permissively licensed.

**(b) Synthetic Data with Non-Permissive or Unverifiable Generators or Seeds.** This tier contains datasets that are themselves permissively licensed (e.g., Apache/MIT/CC-BY), but where either (i) the generator model used to create the synthetic data operates under a more restrictive license (e.g., Llama-3 community license, OpenAI API terms), or (ii) the seed data contains slices whose provenance cannot be fully audited (e.g., partially opaque community mixtures). These datasets constitute only  $\approx 4\%$  of [MixtureVitea](#) and are isolated for transparency so that users who require a strictly permissive generator and seed provenance can exclude them (more detail in Table 14).

**Tier 3 — Civic / Governmental Works.** This tier includes materials that are either statutory public domain (e.g., U.S. federal works) or under a strong public-purpose rationale for reuse (e.g., government websites, regulatory notices). While not always explicitly licensed, such work—typically created for dissemination—is widely recognized as low-risk for inclusion. Filtering with copyright keyword checks further reduces the possibility of inadvertently including restricted content.

## 2.2 DATA PROCESSING PIPELINE

To transform the raw data sources into a high-quality and permissively licensed pretraining corpus, we develop a multistage data processing pipeline. Our curation pipeline includes the following stages: ensuring permissive licensing, filtering for CSAM and offensive language, improving overall content quality, and reducing data redundancy. The following sections detail each component.

**Permissiveness Filtering.** In contrast to standard data pipelines that rely on the retroactive negative filtering of broad web scrapes (e.g., Fan et al., 2025), we employ a **positive inclusion** strategy for web data. Rather than ingesting broad web dumps and filtering post-hoc, we positively select sources based on auditable permissive status. Specifically, we (i) apply an explicit allowlist of governmental and international domains (Appendix L.1), (ii) curate a set of websites with known permissive licenses (Appendix L.2), and (iii) expand this set with risk-mitigated documents by searching for permissive license keywords (e.g., “CC-BY-SA”), excluding documents with restrictive terms (e.g., “all rights reserved”). This upfront design minimizes the risk of including paywalled or opted-out content (e.g., commercial news). We justify the inclusion of governmental works under a strong fair-use rationale, considering their public purpose, content type, and minimal market impact (Appendix J).

**Quality and Safety Filtering.** Per standard practices (Raffel et al., 2020), we remove documents with base64-encoded text (which can disrupt training) and duplicative headers and footers (e.g., “Home — Search”) from FineFineWeb. We remove obscene, adult and CSAM-related content with keyword-based blocklists adapted from prior work (Laurençon et al., 2022; Nakamura et al., 2025). For Wikipedia-based documents, we remove articles about films, sporting events, and biographies of living persons in English with applied targeted filtering, to minimize memorization of facts about people, in case of objection to incorrect facts about people being generated by models trained on

In practice, the actual training data construction follows a granular *domain-aware mixing strategy* (detailed in Section 2.2), where documents are clustered by base URL or provenance to preserve domain coherence per sample, rather than strictly sampling from rigid high-level partitions.

270 **MixtureVitae**. Besides dataset-level filters, we also evaluate the final model’s safety profile via  
 271 standard red-teaming (Appendix F.3).

272 **Deduplication**. Informed by recent findings in large-scale data curation, our deduplication strategy  
 273 prioritizes diversity over purity. While removing exact repetitions mitigates harmful memoriza-  
 274 tion (Lee et al., 2022), prior research finds that aggressive, global near-duplicate removal can be  
 275 detrimental. For example, the creators of the **FineWeb-Edu** dataset (Penedo et al., 2024) reported  
 276 *worsened* model performance by global fuzzy deduplication, postulating that it removed “too much  
 277 quality data.” Therefore, we adopt a local-only approach. We first apply **intra-dataset dedupli-**  
 278 **cation** using prefix-based exact matching to remove verbatim boilerplate text (Lee et al., 2022).  
 279 We **intentionally avoid full, cross-dataset fuzzy deduplication** to preserve near-duplicates (e.g.,  
 280 Wikipedia articles with different formatting across sources). We posit that doing so retains “**stylistic**  
 281 **and domain diversity**,” a factor shown to be helpful for model generalization (Chen et al., 2024).

282 **Training Example Curation**. Our process for creating training examples involves several stages:  
 283 **1. Heuristic Cleaning**. We remove boilerplate content by eliminating repetitive n-gram prefixes  
 284 and suffixes, following standard web data cleaning pipelines (Raffel et al., 2020). **2. Fine-grained**  
 285 **Deduplication**. To enhance data quality, we segment documents into sentences and remove dupli-  
 286 cate sentences within each document. Documents with high internal repetition (sentence duplication  
 287 rate > 75%) are discarded entirely, as this has been shown to improve model performance (Lee et al.,  
 288 2022). **3. Domain-Aware Mixing**. To construct the final training examples, we employ a domain-  
 289 aware data mixing strategy (Xie et al., 2023). Documents are clustered by their base URL (a proxy  
 290 for domain), and sentences are concatenated first within their original document, then packed with  
 291 other documents from the same cluster.

292 **Additional Filtering for Synthetic Datasets**. To ensure that the synthetic subsets of **MixtureVitae**  
 293 adhere to our permissive-first, risk-mitigated approach, we prioritize data originating from seeds that  
 294 are sourced from permissive sources and generated with models that are themselves permissively  
 295 licensed. A small portion ( $\approx 4\%$ ) of **MixtureVitae** originates from sources with opaque, mixed or  
 296 restrictive provenance and is isolated into Tier 2(b), as detailed in Appendix M and Table 14.

## 298 3 EXPERIMENTS

301 We empirically validate the efficacy of **MixtureVitae** through a comprehensive set of evaluations.  
 302 We begin by outlining our controlled experimental framework, model architectures, and baseline  
 303 selection in Section 3.1. We then present the primary scaling behavior and general benchmark per-  
 304 formance in Section 3.2, followed by a focused evaluation on reasoning, mathematics, and coding  
 305 tasks in Section 3.3. To ensure the integrity of these findings, we detail our decontamination pro-  
 306 tocol and leakage analysis in Section 3.4. Finally, we isolate the contributions of specific dataset  
 307 components through ablation studies in Section H, highlighting the critical impact of instruction and  
 308 reasoning data density.

### 310 3.1 EXPERIMENTAL SETUP

312 To empirically validate the quality of the **MixtureVitae** pretraining dataset, we conduct a large-  
 313 scale comparative study against a selection of prominent open pretraining datasets. We isolate the  
 314 impact of the dataset on downstream performance using the **open-sci-ref** training procedure (Nezhu-  
 315 rina et al., 2025), which enables systematic control of factors affecting benchmark scores. As in  
 316 `open-sci-ref`, we fix the model architecture (Table 4, sizes: 0.13B, 0.4B, 1.3B, 1.7B) and  
 317 training hyperparameters (Table 5), varying only the dataset. This design ensures that any perfor-  
 318 mance difference can be attributed solely to the dataset.

319 Also, following the numbers given in `open-sci-ref`, we train each model on two token budgets:  
 320 50B and 300B, to analyze scaling effects. Conducting separate training runs on each budget, rather  
 321 than using intermediate checkpoints, thus ensuring a consistent data distribution and allowing for  
 322 proper optimization of learning rate schedules for each specific token budget (Hoffmann et al., 2022).  
 323 This follows standard practice: data mixtures effective at small token budgets may not generalize to  
 larger ones (Albalak et al., 2023).

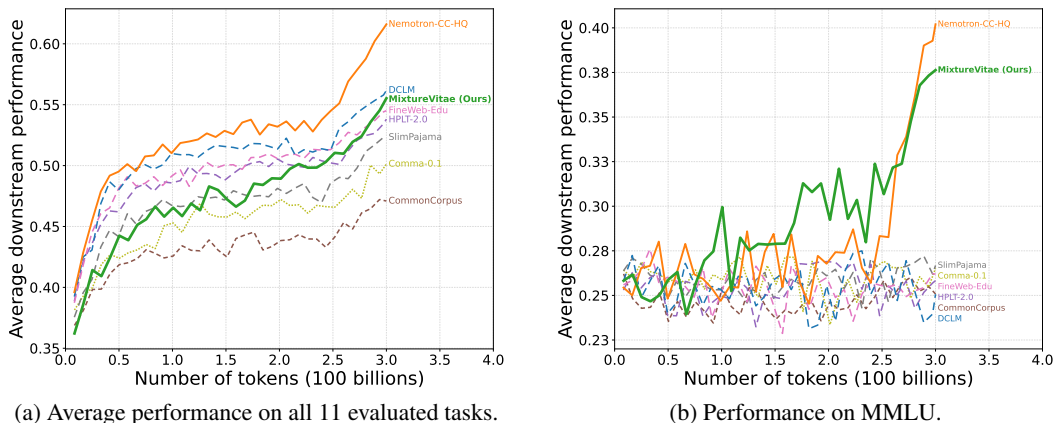


Figure 3: Performance comparison of pretraining datasets for a 1.7B-parameter model trained up to a 300B token budget, showing downstream accuracy as a function of the number of training tokens.

To guard against test-set leakage, we also perform a large-scale 13-gram decontamination analysis and re-evaluation with decontaminated dataset; Section 3.4 and Appendix G detail this procedure.

Within this controlled evaluation framework, we compare **MixtureVitae** with the set of public baselines evaluated in `open-sci-ref`, with the addition of a representative selection of permissively licensed datasets. As detailed in Table 3, the comparison set includes two groups: (i) **Non-Permissive/Mixed-License Baselines.** C4 (Raffel et al., 2020), The Pile (Gao et al., 2020), SlimPajama (Shen et al., 2024), FineWeb-Edu (Penedo et al., 2024), Nemetron-CC-HQ (Su et al., 2025), DCLM-baseline (Li et al., 2024), HPLT Monolingual Datasets v2.0 (Burchell et al., 2025); (ii) **Permissive Baselines.** CommonCorpus and its English subset (Langlais et al., 2025), as well as Comma-0.1 (Kandpal et al., 2025).

All datasets are tokenized using the GPT-NeoX-20B tokenizer (Black et al., 2022), resulting in a vocabulary size of 50,304. The models are trained using Megatron-LM (Shoeybi et al., 2020), and the evaluations are performed using LM Evaluation Harness (Gao et al., 2021). Model performance is evaluated on recognized downstream task benchmarks: MMLU (Hendrycks et al., 2021), COPA (Roemmele et al., 2011), LAMBADA (Paperno et al., 2016), OpenBookQA (Mihaylov et al., 2018), Winogrande (Sakaguchi et al., 2021), ARC (Challenge and Easy) (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), Commonsense-QA (Talmor et al., 2019) and PIQA (Bisk et al., 2020).

### 3.2 EXPERIMENT RESULTS

**Overall average performance.** At a 300B-token budget, **MixtureVitae** shows strong performance when compared to the reference permissive datasets and is almost comparable to the non-permissive datasets (Figure 3, Tab. 1). **MixtureVitae** outperforms all permissive dataset baselines by a significant margin, with gaps widening considerably for larger model sizes, in terms of average performance across all 10 tasks (see Figure 3a, Tab. 1). Non-permissive datasets, particularly Nemetron-CC-HQ and DCLM, still achieve the highest overall performance. Approaching the 300B token budget, **MixtureVitae** catches up to FineWeb-Edu and DCLM. More importantly, while the top-performing models are still trained on non-permissive datasets like Nemetron-CC-HQ and DCLM, our results demonstrate that this performance gap is no longer an inevitability. **MixtureVitae** proves that a dataset built on a fully permissive, risk-mitigated foundation can achieve highly competitive results—significantly outperforming all other permissive baselines and landing within a small, practical margin of top-tier, legally-ambiguous corpora. This finding directly challenges the prevailing assumption that reliance on high-risk, indiscriminately scraped copyrighted data is a prerequisite for training capable LLMs. **MixtureVitae** performs particularly well relative to others on reasoning related tasks like MMLU (Figure 3b, Tab. 1), where most baselines are near random chance. Among all the baselines, only Nemetron-CC-HQ catches up to **MixtureVitae** at around 260B and overtakes it past that point. Our findings also hold at the 50B token budget scale (App. Sec. F.2).

Table 1: Performance comparison of 1.7B-parameter models trained on different pretraining datasets with a 300B token budget. *Italic* denotes the best result among permissive-only datasets, while **bold** indicates the best result overall, including mixed-license datasets. **MixtureVitae** outperforms other permissive datasets across most benchmarks. On reasoning related MMLU, BoolQ, and CommonSense-QA, it also outperforms strong non-permissive baselines.

Benchmark	<b>MixtureVitae</b> (permissive)	Comma-0.1 (permissive)	CommonCorpus (permissive)	FineWeb-Edu (mixed-license)	DCLM (mixed-license)
COPA	<i>0.73</i>	0.71	0.71	0.76	<b>0.81</b>
Lambada	0.48	<i>0.54</i>	0.49	0.52	<b>0.65</b>
OpenBookQA	<i>0.35</i>	0.33	0.31	<b>0.42</b>	0.39
Winogrande	0.58	<i>0.60</i>	0.56	0.61	<b>0.62</b>
MMLU	<b>0.38</b>	0.27	0.25	0.26	0.25
ARC-Challenge	<i>0.40</i>	0.36	0.32	<b>0.44</b>	0.40
ARC-Easy	<i>0.71</i>	0.63	0.61	<b>0.75</b>	0.73
BoolQ	<b>0.75</b>	0.62	0.62	0.67	0.69
CommonSense-QA	<b>0.49</b>	0.21	0.19	0.19	0.20
HellaSwag	<i>0.54</i>	0.53	0.45	0.63	<b>0.67</b>
PIQA	0.70	<i>0.71</i>	0.66	<b>0.76</b>	<b>0.76</b>
Average	<b>0.56</b>	0.50	0.47	<b>0.55</b>	<b>0.56</b>

**Performance on single tasks.** We show performance on each single task in Tab. 1 and in the App. Sec. F.1 (App. Fig. 5). **MixtureVitae** outperforms other permissive datasets on MMLU, Arc Challenge, Arc Easy and BoolQ, while closely matching DCLM and FineWeb-Edu. On PIQA, HellaSwag, Winogrande, OpenBookQA, **MixtureVitae** is on par with Comma-0.1, while both are behind non-permissive datasets. Lambada is the only task where **MixtureVitae** falls behind Comma-0.1. We thus observe **MixtureVitae** to be particularly strong on reasoning-related tasks.

### 3.3 RESULTS ON PROBLEM SOLVING AND INSTRUCTION-BASED DOWNSTREAM TASKS

To further demonstrate the performance of the **MixtureVitae** dataset, we evaluate the model on a set of math, code, and instruction benchmarks: GSM8k (Cobbe et al., 2021), MBPP (Austin et al., 2021), IF-Eval (Zhou et al., 2023). Our evaluation uses the final 1.7B model checkpoints after training for 300B tokens using the `open-sci-ref` protocol (exact evaluation setup in Table 7).

Unlike traditional web-only baselines (e.g., C4, FineWeb, DCLM), **MixtureVitae** utilizes a *reasoning and instruction-heavy* pretraining mixture. Compared against base models with same architecture and matched training compute, this front-loading strategy shows capabilities typically associated with post-training. This pretraining composition leads to a more token-efficient and simple path to reasoning competence already after single base model pre-training stage, matching or outperforming conventional multi-stage extensive pre- and post-training procedures.

The results (Table 2) show a dramatic difference on math (GSM8K) and coding (HumanEval, MBPP). **MixtureVitae** achieves scores of **0.53**, **0.32**, and **0.38**, respectively. This performance is considerably stronger than any other dataset, all of which remain near random performance on GSM8K (0.02-0.06) and cap at 0.13 on HumanEval and 0.22 on MBPP. Most notably, our base model outperforms the post-trained SmoLLM2-1.7B-Instruct (Ben allal et al., 2025) model on GSM8K, HumanEval, and MBPP — despite the latter being trained on  $\approx 11T$  tokens (over  $36\times$  our budget). See App. Sec. H for ablation studies confirming reasoning and instruction subset being largely responsible for the strong performance on math/code tasks (App. Fig. 10)).

### 3.4 TEST LEAKAGE AND DECONTAMINATION

To rule out test-set leakage as an alternative explanation for these gains, we perform a 13-gram exact-match decontamination sweep between **MixtureVitae** and all benchmarks (Appendix G.3). Document-level overlap is negligible for most tasks (e.g., at or below 0.0003% for ARC, HellaSwag, LAMBADA, OpenBookQA, and PIQA; see Table 9); contamination rates are modest for MMLU and BoolQ; for code benchmarks such as HumanEval and MBPP, contamination rates are higher but still small. We then (i) re-evaluate all models on decontaminated test sets with all overlapping

Table 2: **Performance on math, code, and instruction-following tasks for 1.7B models.** We compare **MixtureVitae**—trained on a reasoning- and instruction-heavy, permissive-first mixture—against standard `open-sci-ref` baselines trained on predominantly web-based corpora. **MixtureVitae** shows a substantial lead in math and code tasks. Notably, the 1.7B **MixtureVitae** base model exceeds **SmolLM2-1.7B-Instruct** on GSM8K, HumanEval, and MBPP despite training on 300B rather than  $\approx 11$ T tokens.

Training Dataset	Tokens	IF-Eval	GSM8K	HumanEval	MBPP	Average
<i>Models Trained with open-sci-ref for 300B Tokens</i>						
<b>MixtureVitae</b>	300B	0.19	<b>0.53</b>	<b>0.32</b>	<b>0.38</b>	<b>0.36</b>
Comma-0.1	300B	0.19	0.06	0.13	0.22	0.15
CommonCorpus	300B	0.13	0.02	0.05	0.05	0.06
C4	300B	0.20	0.02	0.00	0.00	0.06
SlimPajama	300B	0.14	0.02	0.05	0.00	0.05
HPLT-2.0	300B	0.17	0.02	0.00	0.00	0.05
DCLM	300B	0.13	0.02	0.01	0.01	0.04
Nemotron-CC-HQ	300B	0.09	0.03	0.02	0.00	0.03
<i>Models Trained with open-sci-ref for 1T Tokens</i>						
FineWeb-Edu	1T	0.20	0.03	0.00	0.00	0.06
Nemotron-CC-HQ	1T	0.13	0.03	0.01	0.04	0.05
DCLM	1T	0.15	0.03	0.00	0.01	0.05
<i>Other Models</i>						
SmolLM2-1.7B	11T	0.18	0.31	0.01	0.35	0.21
SmolLM2-1.7B-Instruct	11T	<b>0.28</b>	0.37	0.28	0.37	0.33

items removed and (ii) retrain a 1.7B model, removing the shards responsible for the majority of the contamination signal. In both cases, **MixtureVitae**’s performance remains essentially unchanged (Table 10 and Figure 8), confirming that our results are not an artifact of test-set leakage.

## 4 RELATED WORK

LLM development is intrinsically linked to the scale and quality of pretraining datasets, which have become larger, more diverse, with a growing emphasis on provenance and licensing recently.

**Pioneering Large-Scale Datasets** Early large-scale text corpora for language modeling often rely on web-crawled data for scale. C4 (Raffel et al., 2020), derived from Common Crawl, is instrumental in training the T5 model, setting standards for large-scale data cleaning and deduplication. Gao et al. (2020) then introduce The Pile, demonstrating the benefit of a more varied data mixture on model generalization and downstream performance. Similarly, ROOTS (Laurençon et al., 2022) supports the training of the BLOOM model with its 498 Common Crawl multilingual scrapes. While foundational, these datasets often have complex or unspecified licenses, mixing permissive data with content of unknown or non-commercial licensing, creating potential legal risks for commercial applications.

**Open and Reproducible Datasets** Amidst many proprietary “black box” datasets, the community has pushed for more openness and reproducibility, moving toward permissive datasets that are also performant, e.g., RedPajama-1T (Weber et al., 2024) and its processing recipes (Touvron et al., 2023), Dolma (Soldaini et al., 2024) and its open-source curation toolkits, SILO (Min et al., 2024). Our work joins this effort, contributing a new risk-mitigated dataset featuring explicit consideration for the underlying copyright.

**Permissively Licensed and Synthetic Data** Growing awareness of copyright and data ownership has spurred interest in datasets built solely from permissively licensed materials. The Stack (Kotchetkov et al., 2023) curates such data for code-generation models, but creating a large, diverse, and high-quality corpus for natural language from exclusively permissive sources remains a challenge. Recent efforts like Common Corpus (Langlais et al., 2025) and The Common Pile (Kandpal et al., 2025) advance the creation of large-scale corpora of permissively licensed and public-domain text. While foundational, our experiments (Section 3) show that models trained on them can lag in com-

plex reasoning, math, and instruction following, suggesting that strictly permissive human text alone is insufficient to instill these advanced skills.

With this scarcity of high-quality reasoning and instruction data, researchers have turned to synthetic data. Alpaca (Taori et al., 2023) and OpenMathInstruct-1 (Toshniwal et al., 2024a) use instructional data for fine-tuning. Phi4 proposes using synthetic data for reasoning tasks (Abdin et al., 2024). Our work, **MixtureVitae**, extends these trends with a meticulously curated, permissive-first, risk-mitigated dataset augmented with targeted synthetic data, providing a strong, legally considered foundation for LLM training to mitigate copyright risks in many existing corpora.

The concurrent work of Apertus (Hernández-Cano et al., 2025) provides large-scale multilingual coverage via retroactive opt-out filtering on web-scale sources. In contrast, **MixtureVitae** focuses on data efficiency through a *positive inclusion strategy*, curating sources known to be permissive (e.g., government works, The Stack) and prioritizing English-centric reasoning density. While Apertus so far released only dataset composition recipes, **MixtureVitae** provides an already composed, ready-to-use pre-training dataset which facilitates reproduction and experimentation by broad community. Moreover, the dataset provides shard-level provenance, enabling risk-tiered remixes that can be flexibly adapted to dataset usage scenario and risk posture.

**Mixing Reasoning Data into Pre-Training** Concurrent with our work, Akter et al. (2025) systematically investigate the “front-loading” of reasoning data, finding that injecting reasoning data into the pretraining phase establishes foundational capabilities that cannot be replicated by scaling supervised fine-tuning (SFT) alone. Similarly, Wang et al. (2025) augment pre-training text data with synthetically generated thinking trajectories. They observe that pre-training augmented with thinking traces strongly outperforms vanilla pretraining using matched compute and token budget (8B model, 100BT) on reasoning/math/language understanding evals. Our findings with **MixtureVitae** align with and extend this observation to the permissive dataset landscape: we show that by front-loading a diverse, risk-mitigated mixture of reasoning and instruction data, we can achieve competitive performance against non-permissive baselines even with a constrained token budget. For a comparison of **MixtureVitae** composition to other permissive and non-permissive baselines, see Table 3.

## 5 DISCUSSION & CONCLUSION

We have introduced **MixtureVitae**, a pretraining corpus serving as a proof-of-concept: **Permissively licensed and permissively-sourced** real and synthetic data can achieve high performance. Our results suggest a shift in the **compliance–performance frontier**. **MixtureVitae** demonstrates that capabilities previously associated with mixed-license corpora are reachable with a permissive first, risk-mitigated approach. In our controlled 300B-token experiments, not only does **MixtureVitae** catch up to leading non-permissive baselines like DCLM and FineWeb-Edu, but our 1.7B base model also outperforms the *post-trained* SmolLM2-1.7B-Instruct—a model trained on  $\approx 11$ T tokens—on GSM8K, HumanEval and MBPP.

**Dominant fraction of reasoning & instruction data in the pre-training.** **MixtureVitae**’s performance is enhanced by the large proportion of reasoning and instruction data, as demonstrated in the ablation study in Section H. Removing this subset (“w/o Instructions” in Figure 10) causes a substantial degradation across tasks—far larger than the impact of removing the web component. This observation validates and extends the findings of Phi-4 (Abdin et al., 2024). Importantly, we show that heavily increasing reasoning and instruction data fraction on expense of generic web text creates overall boost in performance, specifically on math/code tasks (Tab. 2), *without* hurting core language understanding capabilities (Fig. 3, Tab. 1).

Beyond this specific corpus, the three-tier licensing scheme and its shard-level annotations provide a **concrete template for structuring risk-mitigated mixtures in future work**, and **MixtureVitae** as a whole serves as a reusable blueprint for compliant pretraining. We demonstrate a pipeline built on positive-inclusion “pseudo-crawling,” tiered provenance tracking, targeted synthetic generation with audited seeds and careful decontamination against test set leakage. As detailed in our scaling outlook (Appendix N), this recipe provides a path to extend compliant pretraining to the multi-trillion-token regime—via subset upsampling, multilingual expansion, and synthetic growth—providing the community with a sustainable alternative to the legal uncertainty of broad web scrapes.

## REFERENCES

- 540  
541  
542 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar,  
543 Michael Harrison, Russell J. Hewett, Mojan Javaheripi, Piero Kauffmann, James R. Lee, Yin Tat  
544 Lee, Yuanzhi Li, Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric Price, Gustavo de Rosa,  
545 Olli Saarikivi, Adil Salim, Shital Shah, Xin Wang, Rachel Ward, Yue Wu, Dingli Yu, Cyril  
546 Zhang, and Yi Zhang. Phi-4 technical report, 2024. URL <https://arxiv.org/abs/2412.08905>.
- 547  
548 Syeda Nahida Akter, Shrimai Prabhume, Eric Nyberg, Mostofa Patwary, Mohammad Shoeybi,  
549 Yejin Choi, and Bryan Catanzaro. Front-loading reasoning: The synergy between pretraining and  
550 post-training data, 2025. URL <https://arxiv.org/abs/2510.03264>.
- 551  
552 Alon Albalak, Liangming Pan, Colin Raffel, and William Yang Wang. Efficient online data mix-  
553 ing for language model pre-training. In *RO-FoMo: Robustness of Few-shot and Zero-shot Learn-  
554 ing in Large Foundation Models*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=9Tze4oy41w)  
555 9Tze4oy41w.
- 556  
557 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
558 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large  
559 language models, 2021. URL <https://arxiv.org/abs/2108.07732>.
- 560  
561 Samuel Barham, Orion Weller, Michelle Yuan, Kenton Murray, Mahsa Yarmohammadi, Zhengping  
562 Jiang, Siddharth Vashishtha, Alexander Martin, Anqi Liu, Aaron Steven White, Jordan Boyd-  
563 Graber, and Benjamin Van Durme. Megawika: Millions of reports and their sources across 50  
564 diverse languages, 2023. URL <https://arxiv.org/abs/2307.07049>.
- 565  
566 Loubna Ben allal, Anton Lozhkov, Elie Bakouch, Gabriel Martin Blazquez, Guilherme Penedo,  
567 Lewis Tunstall, Andrés Marafioti, Agustín Piqueres Lajarín, Hynek Kydlíček, Vaibhav Srivastav,  
568 Joshua Lochner, et al. Smollm2: When smol goes big—data-centric training of a fully open small  
569 language model. In *Second Conference on Language Modeling*, 2025.
- 570  
571 Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil,  
572 Zach Moshe, Tomer Ronen, Najeeb Nabwani, Ido Shahaf, Oren Tropp, Ehud Karpas, Ran Zil-  
573 berstein, Jiaqi Zeng, Soumye Singhal, Alexander Bukharin, Yian Zhang, Tugrul Konuk, Gerald  
574 Shen, Ameya Sunil Mahabaleshwarkar, Bilal Kartal, Yoshi Suhara, Olivier Delalleau, Zijia Chen,  
575 Zhilin Wang, David Mosallanezhad, Adi Renduchintala, Haifeng Qian, Dima Rekeshe, Fei Jia,  
576 Somshubra Majumdar, Wahid Noroozi, Wasi Uddin Ahmad, Sean Narenthiran, Aleksander Ficek,  
577 Mehrzad Samadi, Jocelyn Huang, Siddhartha Jain, Igor Gitman, Ivan Moshkov, Wei Du, Shub-  
578 ham Toshniwal, George Armstrong, Branislav Kisacanin, Matvei Novikov, Daria Gitman, Evelina  
579 Bakhturina, Prasoon Varshney, Makesh Narsimhan, Jane Polak Scowcroft, John Kamalu, Dan Su,  
580 Kezhi Kong, Markus Kliegl, Rabeeh Karimi Mahabadi, Ying Lin, Sanjeev Satheesh, Jupinder Par-  
581 mar, Pritam Gundecha, Brandon Norick, Joseph Jennings, Shrimai Prabhume, Syeda Nahida  
582 Akter, Mostofa Patwary, Abhinav Khattar, Deepak Narayanan, Roger Waleffe, Jimmy Zhang,  
583 Bor-Yiing Su, Guyue Huang, Terry Kong, Parth Chadha, Sahil Jain, Christine Harvey, Elad  
584 Segal, Jining Huang, Sergey Kashirsky, Robert McQueen, Izzy Putterman, George Lam, Arun  
585 Venkatesan, Sherry Wu, Vinh Nguyen, Manoj Kilaru, Andrew Wang, Anna Warno, Abhilash  
586 Somasamudramath, Sandip Bhaskar, Maka Dong, Nave Assaf, Shahar Mor, Omer Ullman Ar-  
587 gov, Scot Junkin, Oleksandr Romanenko, Pedro Larroy, Monika Katariya, Marco Rovinelli, Viji  
588 Balas, Nicholas Edelman, Anahita Bhiwandiwala, Muthu Subramaniam, Smita Ithape, Karthik  
589 Ramamoorthy, Yuting Wu, Suguna Varshini Velury, Omri Almog, Joyjit Daw, Denys Fridman,  
590 Erick Galinkin, Michael Evans, Shaona Ghosh, Katherine Luna, Leon Derczynski, Nikki Pope,  
591 Eileen Long, Seth Schneider, Guillermo Siman, Tomasz Grzegorzec, Pablo Ribalta, Monika  
592 Katariya, Chris Alexiuk, Joey Conway, Trisha Saar, Ann Guan, Krzysztof Pawelec, Shyamala  
593 Prayaga, Oleksii Kuchaiev, Boris Ginsburg, Oluwatobi Olabiyi, Kari Briski, Jonathan Cohen,  
594 Bryan Catanzaro, Jonah Alben, Yonatan Geifman, and Eric Chung. Llama-nemotron: Efficient  
595 reasoning models, 2025. URL <https://arxiv.org/abs/2505.00949>.
- 596  
597 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical com-  
598 monsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*,  
599 volume 34, pp. 7432–7439, 2020.

- 594 Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace  
595 He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivan-  
596 shu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. GPT-NeoX-20B:  
597 An open-source autoregressive language model, 2022. URL [https://arxiv.org/abs/  
598 2204.06745](https://arxiv.org/abs/2204.06745).
- 599 Michael J. Bommarito, II, Jillian Bommarito, and Daniel Martin Katz. The k13m data project:  
600 Copyright-clean training resources for large language models, 2025. URL [https://arxiv.  
601 org/abs/2504.07854](https://arxiv.org/abs/2504.07854).
- 602 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
603 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
604 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 605 Laurie Burchell, Ona De Gibert Bonet, Nikolay Arefyev, Mikko Aulamo, Marta Bañón, Pinzhen  
606 Chen, Mariia Fedorova, Liane Guillou, Barry Haddow, Jan Hajic, et al. An expanded massive  
607 multilingual dataset for high-performance language technologies (hplt). In *Proceedings of the  
608 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,  
609 pp. 17452–17485, 2025.
- 610 Hao Chen, Abdul Waheed, Xiang Li, Yidong Wang, Jindong Wang, Bhiksha Raj, and Marah I  
611 Abdin. On the diversity of synthetic data and its impact on training large language models. *arXiv  
612 preprint arXiv:2410.15226*, 2024.
- 613 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan  
614 Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu,  
615 Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pel-  
616 lat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao,  
617 Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin,  
618 Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language  
619 models, 2022. URL <https://arxiv.org/abs/2210.11416>.
- 620 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina  
621 Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Jill  
622 Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of  
623 the North American Chapter of the Association for Computational Linguistics: Human Lan-  
624 guage Technologies, Volume 1 (Long and Short Papers)*, pp. 2924–2936, Minneapolis, Min-  
625 nesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1300. URL  
626 <https://aclanthology.org/N19-1300/>.
- 627 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
628 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
629 *arXiv preprint arXiv:1803.05457*, 2018.
- 630 Colin B Clement, Matthew Bierbaum, Kevin P O’Keeffe, and Alexander A Alemi. On the use of  
631 arxiv as a dataset. *arXiv preprint arXiv:1905.00075*, 2019.
- 632 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
633 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
634 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 635 Codefuse Team, Ling Team, Wenting Cai, Yuchen Cao, Chaoyu Chen, Chen Chen, Siba Chen, Qing  
636 Cui, Peng Di, Junpeng Fang, Zi Gong, Ting Guo, Zhengyu He, Yang Huang, Cong Li, Jianguo Li,  
637 Zheng Li, Shijie Lian, BingChang Liu, Songshan Luo, Shuo Mao, Min Shen, Jian Wu, Jialong  
638 Yang, Wenjie Yang, Tong Ye, Hang Yu, Wei Zhang, Zhenduo Zhang, Hailin Zhao, Xunjin Zheng,  
639 and Jun Zhou. Every sample matters: Leveraging mixture-of-experts and high-quality data for  
640 efficient and accurate code llm, 2025. URL <https://arxiv.org/abs/2503.17793>.
- 641 Common Crawl Foundation. Common Crawl. <https://commoncrawl.org/>, 2025. Ac-  
642 cessed: 2025-08-25.
- 643 United States Congress. Copyright act of 1976, 1976. URL [https://www.copyright.gov/  
644 title17/](https://www.copyright.gov/title17/). Public Law 94-553, Enacted October 19, 1976.

- 648 Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu,  
649 and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback, 2024.  
650 URL <https://openreview.net/forum?id=pNkOx3IVWI>.  
651
- 652 European Union. Directive (EU) 2019/790 of the European Parliament and of the Council of 17  
653 April 2019 on copyright and related rights in the Digital Single Market, 2019. L 130/92.
- 654 Dongyang Fan, Vinko Sabolčec, Matin Ansari-pour, Ayush Kumar Tarun, Martin Jaggi, Antoine  
655 Bosselut, and Imanol Schlag. Can performant LLMs be ethical? quantifying the impact of  
656 web crawling opt-outs. In *Second Conference on Language Modeling*, 2025. URL <https://openreview.net/forum?id=a6QsOjr3wo>.  
657
- 658 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason  
659 Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile:  
660 An 800gb dataset of diverse text for language modeling, 2020. URL <https://arxiv.org/abs/2101.00027>.  
661
- 662 Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence  
663 Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric  
664 Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot lan-  
665 guage model evaluation, September 2021. URL [https://doi.org/10.5281/zenodo.](https://doi.org/10.5281/zenodo.5371628)  
666 5371628.  
667
- 668 Glaive AI. Glaive-ai reasoning dataset. [https://huggingface.co/datasets/](https://huggingface.co/datasets/glaiveai/reasoning-v1-20m)  
669 [glaiveai/reasoning-v1-20m](https://huggingface.co/datasets/glaiveai/reasoning-v1-20m), 2023.  
670
- 671 Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna  
672 Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, et al. Opendthoughts: Data recipes for reason-  
673 ing models. *arXiv preprint arXiv:2506.04178*, 2025.
- 674 Xintong Hao, Ruijie Zhu, Ge Zhang, Ke Shen, and Chenggang Li. Reformulation for pretraining  
675 data augmentation, 2025. URL <https://arxiv.org/abs/2502.04235>.  
676
- 677 Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Ka-  
678 mar. ToxiGen: A large-scale machine-generated dataset for adversarial and implicit hate speech  
679 detection. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings*  
680 *of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
681 *Papers)*, pp. 3309–3326, Dublin, Ireland, May 2022. Association for Computational Linguis-  
682 tics. doi: 10.18653/v1/2022.acl-long.234. URL [https://aclanthology.org/2022.](https://aclanthology.org/2022.acl-long.234/)  
683 [acl-long.234/](https://aclanthology.org/2022.acl-long.234/).
- 684 Kenneth Heafield, Elaine Farrow, Jelmer van der Linde, Gema Ramírez-Sánchez, and Dion Wig-  
685 gins. The EuroPat corpus: A parallel corpus of European patent data. In Nicoletta Calzolari,  
686 Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara  
687 Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Jan Odijk, and Stelios  
688 Piperidis (eds.), *Proceedings of the Thirteenth Language Resources and Evaluation Conference*,  
689 pp. 732–740, Marseille, France, June 2022. European Language Resources Association. URL  
690 <https://aclanthology.org/2022.lrec-1.78/>.
- 691 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja-  
692 cob Steinhardt. Measuring massive multitask language understanding. In *International Confer-*  
693 *ence on Learning Representations*, 2021. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=d7KBjmI3GmQ)  
694 [d7KBjmI3GmQ](https://openreview.net/forum?id=d7KBjmI3GmQ).
- 695 Alejandro Hernández-Cano, Alexander Hägele, Allen Hao Huang, Angelika Romanou, Antoni-Joan  
696 Solergibert, Barna Pasztor, Bettina Messmer, Dhia Garbaya, Eduard Frank Ďurech, Ido Hakimi,  
697 et al. Apertus: Democratizing open and compliant llms for global language environments. *arXiv*  
698 *preprint arXiv:2509.14233*, 2025.  
699
- 700 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
701 Rutherford, Diego de las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hen-  
nigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia

- 702 Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack William Rae, and  
703 Laurent Sifre. An empirical analysis of compute-optimal large language model training. In  
704 Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural  
705 Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=iBBcRU1OAPR>.  
706  
707
- 708 Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael  
709 Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabza. Llama guard: Llm-  
710 based input-output safeguard for human-ai conversations, 2023. URL [https://arxiv.org/  
711 abs/2312.06674](https://arxiv.org/abs/2312.06674).
- 712 Nikhil Kandpal, Brian Lester, Colin Raffel, Sebastian Majstorovic, Stella Biderman, Baber Abbasi,  
713 Luca Soldaini, Enrico Shippole, A. Feder Cooper, Aviya Skowron, John Kirchenbauer, Shayne  
714 Longpre, Lintang Sutawika, Alon Albalak, Zhenlin Xu, Guilherme Penedo, Loubna Ben Al-  
715 lal, Elie Bakouch, John David Pressman, Honglu Fan, Dashiell Stander, Guangyu Song, Aaron  
716 Gokaslan, Tom Goldstein, Brian R. Bartoldson, Bhavya Kailkhura, and Tyler Murray. The  
717 common pile v0.1: An 8tb dataset of public domain and openly licensed text, 2025. URL  
718 <https://arxiv.org/abs/2506.05209>.
- 719 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,  
720 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
721 models, 2020. URL <https://arxiv.org/abs/2001.08361>.  
722
- 723 Denis Kocetkov, Raymond Li, Chenghao Mou, Yacine Jernite, Margaret Mitchell, Carlos Muñoz  
724 Ferrandis, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro Von Werra, et al. The stack:  
725 3 tb of permissively licensed source code. *Transactions on Machine Learning Research*, 2023.  
726
- 727 Andreas Köpf, Yannic Kilcher, Dimitri Von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith  
728 Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, et al. Openassistant  
729 conversations-democratizing large language model alignment. *Advances in neural information  
730 processing systems*, 36:47669–47681, 2023.
- 731 Pierre-Carl Langlais. Releasing youtube-commons: a massive open corpus for conversational and  
732 multimodal data. *Hugging Face blog*, April 2024. URL [https://huggingface.co/blog/  
733 Pclanglais/youtube-commons](https://huggingface.co/blog/Pclanglais/youtube-commons).
- 734
- 735 Pierre-Carl Langlais, Carlos Rosas Hinostroza, Mattia Nee, Catherine Arnett, Pavel Chizhov,  
736 Eliot Krzystof Jones, Irène Girard, David Mach, Anastasia Stasenko, and Ivan P. Yamshchikov.  
737 Common corpus: The largest collection of ethical data for llm pre-training, 2025. URL  
738 <https://arxiv.org/abs/2506.01732>.
- 739 Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral,  
740 Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu  
741 Nguyen, Jörg Froberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gerard  
742 Dupont, Stella Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, Giada Pistilli,  
743 Olivier Nguyen, Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa,  
744 Paulo Villegas, Tristan Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel  
745 Muñoz, Jian Zhu, Daniel Van Strien, Zaid Alyafeai, Khalid Almubarak, Minh Chien Vu,  
746 Itziar Gonzalez-Dios, Aitor Soroa, Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan,  
747 Shamik Bose, David Adelani, Long Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim,  
748 Violette Lepercq, Suzana Ilic, Margaret Mitchell, Sasha Alexandra Luccioni, and Yacine  
749 Jernite. The bigscience roots corpus: A 1.6tb composite multilingual dataset. In S. Koyejo,  
750 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural In-  
751 formation Processing Systems*, volume 35, pp. 31809–31826. Curran Associates, Inc., 2022.  
752 URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/  
753 ce9e92e3de2372a4b93353eb7f3dc0bd-Paper-Datasets\\_and\\_Benchmarks.  
754 pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/ce9e92e3de2372a4b93353eb7f3dc0bd-Paper-Datasets_and_Benchmarks.pdf).
- 755 Ariel N. Lee, Cole J. Hunter, and Nataniel Ruiz. Platypus: Quick, cheap, and powerful refinement  
of llms. *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*, 2023.

- 756 Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-  
757 Burch, and Nicholas Carlini. Deduplicating training data makes language models better. In  
758 *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume*  
759 *1: Long Papers)*, pp. 8424–8445, 2022.
- 760  
761 Mark A Lemley and Bryan Casey. Fair learning. *Texas Law Review*, 95:1, 2017.
- 762  
763 Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Gadre, Hritik Bansal, Etash  
764 Guha, Sedrick Keh, Kushal Arora, Saurabh Garg, Rui Xin, Niklas Muennighoff, Reinhard Heckel,  
765 Jean Mercat, Mayee Chen, Suchin Gururangan, Mitchell Wortsman, Alon Albalak, Yonatan Bit-  
766 ton, Marianna Nezhurina, Amro Abbas, Cheng-Yu Hsieh, Dhruva Ghosh, Josh Gardner, Maciej  
767 Kilian, Hanlin Zhang, Rulin Shao, Sarah Pratt, Sunny Sanyal, Gabriel Ilharco, Giannis Daras,  
768 Kalyani Marathe, Aaron Gokaslan, Jieyu Zhang, Khyathi Chandu, Thao Nguyen, Igor Vasiljevic,  
769 Sham Kakade, Shuran Song, Sujay Sanghavi, Fartash Faghri, Sewoong Oh, Luke Zettlemoyer,  
770 Kyle Lo, Alaaeldin El-Nouby, Hadi Pouransari, Alexander Toshev, Stephanie Wang, Dirk Groen-  
771 eveld, Luca Soldaini, Pang Wei Koh, Jenia Jitsev, Thomas Kollar, Alexandros G. Dimakis, Yair  
772 Carmon, Achal Dave, Ludwig Schmidt, and Vaishaal Shankar. Datacomp-lm: In search of the  
773 next generation of training sets for language models. *Advances in Neural Information Processing*  
774 *Systems*, 37:14200–14282, 2024.
- 775  
776 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom  
777 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien  
778 de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven  
779 Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson,  
780 Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level  
781 code generation with alphacode. *Science*, 378(6624):1092–1097, December 2022. ISSN 1095-  
782 9203. doi: 10.1126/science.abq1158. URL <http://dx.doi.org/10.1126/science.abq1158>.
- 783  
784 Varvara Logacheva, Daryna Dementieva, Sergey Ustyantsev, Daniil Moskovskiy, David Dale, Irina  
785 Krotova, Nikita Semenov, and Alexander Panchenko. ParaDetox: Detoxification with paral-  
786 lel data. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings*  
787 *of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
788 *Papers)*, pp. 6804–6818, Dublin, Ireland, May 2022. Association for Computational Linguis-  
789 tics. doi: 10.18653/v1/2022.acl-long.469. URL <https://aclanthology.org/2022.acl-long.469/>.
- 790  
791 M-A-P, Ge Zhang, Xinrun Du, Zhimiao Yu, Zili Wang, Zekun Wang, Shuyue Guo, Tianyu Zheng,  
792 Kang Zhu, Jerry Liu, Shawn Yue, Binbin Liu, Zhongyuan Peng, Yifan Yao, Jack Yang, Ziming  
793 Li, Bingni Zhang, Minghao Liu, Tianyu Liu, Yang Gao, Wenhui Chen, Xiaohuan Zhou, Qian  
794 Liu, Taifeng Wang, and Wenhao Huang. Finefineweb: A comprehensive study on fine-grained  
795 domain web corpus. <https://huggingface.co/datasets/m-a-p/FineFineWeb>,  
796 December 2024.
- 797  
798 Rabeeh Karimi Mahabadi, Sanjeev Satheesh, Shrimai Prabhumoye, Mostofa Patwary, Mohammad  
799 Shoeybi, and Bryan Catanzaro. Nemotron-cc-math: A 133 billion-token-scale high quality math  
800 pretraining dataset. *arXiv preprint arXiv:2508.15096*, 2025.
- 801  
802 Pratyush Maini, Skyler Seto, Richard Bai, David Grangier, Yizhe Zhang, and Navdeep Jaitly.  
803 Rephrasing the web: A recipe for compute and data-efficient language modeling. In *Proceed-*  
804 *ings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1:*  
805 *Long Papers)*, pp. 14044–14072, 2024.
- 806  
807 Thomas Margoni and Martin Kretschmer. A deeper look into the eu text and data mining exceptions:  
808 Harmonisation, data ownership, and the future of technology. *GRUR International*, 71(8):685–  
809 701, 07 2022. ISSN 2632-8623. doi: 10.1093/grurint/ikac054. URL <https://doi.org/10.1093/grurint/ikac054>.
- 809  
810 Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024. URL  
<https://ai.meta.com/blog/meta-llama-3/>.

- 810 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct elec-  
811 tricity? a new dataset for open book question answering. In *Proceedings of the 2018 Conference*  
812 *on Empirical Methods in Natural Language Processing*, pp. 2381–2391, 2018.
- 813  
814 Sewon Min, Suchin Gururangan, Eric Wallace, Weijia Shi, Hannaneh Hajishirzi, Noah A. Smith,  
815 and Luke Zettlemoyer. SILO language models: Isolating legal risk in a nonparametric datastore.  
816 In *The Twelfth International Conference on Learning Representations*, 2024. URL [https://](https://openreview.net/forum?id=rुक0nyQPec)  
817 [openreview.net/forum?id=rुक0nyQPec](https://openreview.net/forum?id=rुक0nyQPec).
- 818 Taishi Nakamura, Mayank Mishra, Simone Tedeschi, Yekun Chai, Jason T Stillerman, Felix  
819 Friedrich, Prateek Yadav, Tanmay Laud, Vu Minh Chien, Terry Yue Zhuo, et al. Aurora-m:  
820 Open source continual pre-training for multilingual language and code. In *Proceedings of the*  
821 *31st International Conference on Computational Linguistics: Industry Track*, pp. 656–678, 2025.
- 822 National Library of Medicine (U.S.). Pubmed. <https://pubmed.ncbi.nlm.nih.gov/>,  
823 1996.
- 824  
825 Marianna Nezhurina, Jörg Franke, Taishi Nakamura, Timur Carstensen, Niccolò Ajroldi, Ville Ko-  
826 mulainen, David Salinas, and Jenia Jitsev. Open-sci-ref-0.01: open and reproducible reference  
827 baselines for language model and dataset comparison, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2509.09009)  
828 [abs/2509.09009](https://arxiv.org/abs/2509.09009).
- 829 Huu Nguyen, Ken Tsui, Andrej Radonjic, and Christoph Schuhmann. Valid (video-audio large  
830 interleaved dataset), 2024. URL [https://huggingface.co/datasets/ontocord/](https://huggingface.co/datasets/ontocord/VALID)  
831 [VALID](https://huggingface.co/datasets/ontocord/VALID).
- 832 NVIDIA. SFT DataBlend v1. [https://huggingface.co/datasets/nvidia/sft\\_](https://huggingface.co/datasets/nvidia/sft_datablend_v1)  
833 [datablend\\_v1](https://huggingface.co/datasets/nvidia/sft_datablend_v1), 2024.
- 834  
835 NVIDIA. Nemotron-PrismMath Dataset. [https://huggingface.co/datasets/](https://huggingface.co/datasets/nvidia/Nemotron-PrismMath)  
836 [nvidia/Nemotron-PrismMath](https://huggingface.co/datasets/nvidia/Nemotron-PrismMath), 2025.
- 837 NVIDIA Corporation. OpenScience Dataset, 2025. URL [https://huggingface.co/](https://huggingface.co/datasets/nvidia/OpenScience)  
838 [datasets/nvidia/OpenScience](https://huggingface.co/datasets/nvidia/OpenScience).
- 839  
840 Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi,  
841 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The LAMBADA dataset:  
842 Word prediction requiring a broad discourse context. In Katrin Erk and Noah A. Smith (eds.),  
843 *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Vol-*  
844 *ume 1: Long Papers)*, pp. 1525–1534, Berlin, Germany, August 2016. Association for Com-  
845 putational Linguistics. doi: 10.18653/v1/P16-1144. URL [https://aclanthology.org/](https://aclanthology.org/P16-1144/)  
846 [P16-1144/](https://aclanthology.org/P16-1144/).
- 847  
848 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin  
849 Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the  
850 finest text data at scale, 2024. URL <https://arxiv.org/abs/2406.17557>.
- 851  
852 Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John  
853 Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan,  
854 Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks,  
855 Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron  
856 Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu,  
857 Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen  
858 Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kun-  
859 coro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Men-  
860 sch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux,  
861 Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yu-  
862 jia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Au-  
863 relia 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, Llorayne Bennett, Demis Hassabis, Koray Kavukcuoglu,  
and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training go-  
pher, 2022. URL <https://arxiv.org/abs/2112.11446>.

- 864 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
865 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
866 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 867
- 868 Melissa Roemmele, Cosmin Bejan, and Andrew Gordon. Choice of plausible alternatives: An  
869 evaluation of commonsense causal reasoning. In *AAAI spring symposium: logical formalizations*  
870 *of commonsense reasoning*, pp. 90–95, 2011.
- 871
- 872 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-  
873 sarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 874 Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, An-  
875 toine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M Saiful Bari, Canwen  
876 Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chh-  
877 ablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo  
878 Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala  
879 Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan  
880 Teehan, Stella Biderman, Leo Gao, Tali Bers, Thomas Wolf, and Alexander M. Rush. Multitask  
881 prompted training enables zero-shot task generalization. In *International Conference on Learning*  
882 *Representations*, 2022. URL <https://openreview.net/forum?id=9Vrb9D0WI4>.
- 883 David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. Analysing mathematical rea-  
884 soning abilities of neural models. In *International Conference on Learning Representations*,  
885 2019.
- 886
- 887 Zhiqiang Shen, Tianhua Tao, Liqun Ma, Willie Neiswanger, Zhengzhong Liu, Hongyi Wang, Bowen  
888 Tan, Joel Hestness, Natalia Vassilieva, Daria Soboleva, and Eric Xing. Slimpajama-dc: Under-  
889 standing data combinations for llm training, 2024. URL [https://arxiv.org/abs/2309.](https://arxiv.org/abs/2309.10818)  
890 10818.
- 891 Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan  
892 Catanzaro. Megatron-lm: Training multi-billion parameter language models using model par-  
893 allelism, 2020. URL <https://arxiv.org/abs/1909.08053>.
- 894
- 895 Luca Soldaini and Kyle Lo. peS2o (Pretraining Efficiently on S2ORC) Dataset. Technical report,  
896 Allen Institute for AI, 2023. ODC-By, <https://github.com/allenai/pes2o>.
- 897
- 898 Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur,  
899 Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. Dolma: An open corpus of  
900 three trillion tokens for language model pretraining research. In *Proceedings of the 62nd annual*  
901 *meeting of the association for computational linguistics (volume 1: long papers)*, pp. 15725–  
902 15788, 2024.
- 903
- 904 Dan Su, Kezhi Kong, Ying Lin, Joseph Jennings, Brandon Norrick, Markus Kliegl, Mostofa Patwary,  
905 Mohammad Shoeybi, and Bryan Catanzaro. Nemotron-cc: Transforming common crawl into  
906 a refined long-horizon pretraining dataset. In *Proceedings of the 63rd Annual Meeting of the*  
907 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2459–2475, 2025.
- 908
- 909 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A ques-  
910 tion answering challenge targeting commonsense knowledge. In Jill Burstein, Christy Doran, and  
911 Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of*  
912 *the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long*  
913 *and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Com-  
914 putational Linguistics. doi: 10.18653/v1/N19-1421. URL [https://aclanthology.org/](https://aclanthology.org/N19-1421/)  
915 N19-1421/.
- 916
- 917 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. URL [https://huggingface.co/spaces/LLM360/](https://huggingface.co/spaces/LLM360/TxT360)  
TxT360.

- 918 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
919 Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.  
920
- 921 Yuxuan Tong, Xiwen Zhang, Rui Wang, Ruidong Wu, and Junxian He. Dart-math: Difficulty-aware  
922 rejection tuning for mathematical problem-solving. *Advances in Neural Information Processing*  
923 *Systems*, 37:7821–7846, 2024.
- 924 Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Gitman.  
925 Openmathinstruct-1: A 1.8 million math instruction tuning dataset. *Advances in Neural Informa-*  
926 *tion Processing Systems*, 37:34737–34774, 2024a.
- 927 Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Gitman.  
928 Openmathinstruct-1: A 1.8 million math instruction tuning dataset, 2024b. URL [https://](https://arxiv.org/abs/2402.10176)  
929 [arxiv.org/abs/2402.10176](https://arxiv.org/abs/2402.10176).  
930
- 931 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
932 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Ar-  
933 mand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation  
934 language models, 2023. URL <https://arxiv.org/abs/2302.13971>.
- 935 Ulab-UIUC and MetaGPT. OpenManus-RL Dataset. [https://huggingface.co/](https://huggingface.co/datasets/CharlieDreemur/OpenManus-RL)  
936 [datasets/CharlieDreemur/OpenManus-RL](https://huggingface.co/datasets/CharlieDreemur/OpenManus-RL), 2024.  
937
- 938 United States Patent and Trademark Office. USPTO Patent Public Data Sets. [https://](https://developer.uspto.gov/product/patent-public-data-sets)  
939 [developer.uspto.gov/product/patent-public-data-sets](https://developer.uspto.gov/product/patent-public-data-sets), 2024.
- 940 U.S. Securities and Exchange Commission. EDGAR: Electronic Data Gathering, Analysis, and  
941 Retrieval System. <https://www.sec.gov/edgar>, 2024.  
942
- 943 Liang Wang, Nan Yang, Shaohan Huang, Li Dong, and Furu Wei. Thinking augmented pre-training.  
944 *arXiv preprint arXiv:2509.20186*, 2025.
- 945 Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. Do-not-answer:  
946 Evaluating safeguards in llms. In *Findings of the Association for Computational Linguistics:*  
947 *EACL 2024*, pp. 896–911, 2024.  
948
- 949 Maurice Weber, Dan Fu, Quentin Anthony, Yonatan Oren, Shane Adams, Anton Alexandrov, Xi-  
950 aozhong Lyu, Huu Nguyen, Xiaozhe Yao, Virginia Adams, et al. Redpajama: an open dataset for  
951 training large language models. *Advances in neural information processing systems*, 37:116462–  
952 116492, 2024.
- 953 Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang,  
954 Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up  
955 language model pretraining. In *Advances in Neural Information Processing Systems*, volume 36,  
956 pp. 69798–69818, 2023.  
957
- 958 Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and  
959 Bill Yuchen Lin. Magpie: Alignment data synthesis from scratch by prompting aligned llms with  
960 nothing, 2024. URL <https://arxiv.org/abs/2406.08464>.
- 961 Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhen-  
962 guo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions  
963 for large language models. In *The Twelfth International Conference on Learning Representations*,  
964 2024.
- 965 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-  
966 chine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association*  
967 *for Computational Linguistics*, 2019. URL <https://arxiv.org/abs/1905.07830>.  
968
- 969 Lucia Zheng, Neel Guha, Brandon R Anderson, Peter Henderson, and Daniel E Ho. When does  
970 pretraining help? assessing self-supervised learning for law and the casehold dataset of 53,000+  
971 legal holdings. In *Proceedings of the eighteenth international conference on artificial intelligence*  
*and law*, pp. 159–168, 2021.

972 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny  
 973 Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv preprint*  
 974 *arXiv:2311.07911*, 2023.

975 Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson.  
 976 Universal and transferable adversarial attacks on aligned language models. *arXiv preprint*  
 977 *arXiv:2307.15043*, 2023.

## 979 A APPENDIX

### 980 B REPRODUCIBILITY STATEMENT

981 We release our code at <https://anonymous.4open.science/r/mixturevitae-FEFE>  
 982 , with a frozen snapshot at commit 6785991a corresponding to this submission.

#### 983 B.1 DATASET AND CURATION RECIPES

- 984 • **Public Release:** The full **422B** token dataset, along with the 100B and 50B subsets used  
 985 for scaling ablations experiments, will be made publicly available upon acceptance of this  
 986 paper.
- 987 • **Curation Methodology:**
  - 988 – **Dataset Composition** The detailed list of sources and their composition are shown in  
 989 Figure 4.
  - 990 – **Code:** We are including our data curation and math word problem generation scripts  
 991 with the submission.

#### 992 B.2 TRAINING PROCEDURE

993 To ensure our experiments are directly comparable and reproducible, we adhered to a controlled,  
 994 public framework.

- 1000 • **Framework:** All experiments were conducted using the **open-sci-ref** training proce-  
 1001 dure (Nezhurina et al., 2025), which standardizes key factors affecting performance.
- 1002 • **Architectures:** The exact model architectures for all four scales (0.13B, 0.4B, 1.3B, 1.7B)  
 1003 are detailed in Table 4.
- 1004 • **Hyperparameters:** The complete training schedules and hyperparameters (learning rate,  
 1005 batch size, warmup, etc.) for both the 50B and 300B token budgets are specified in Ta-  
 1006 ble E.1.
- 1007 • **Software:** Models were trained using Megatron-LM (Shoeybi et al., 2020) with the GPT-  
 1008 NeoX-20B tokenizer(Black et al., 2022).
- 1009 • **Code:** We are including our training script with the submission.

#### 1010 B.3 EVALUATION AND ANALYSIS

1011 Our evaluation protocol is fully specified to allow for independent verification of our results.

- 1012 • **Framework:** All general and reasoning task evaluations were performed using the public  
 1013 LM Evaluation Harness (Gao et al., 2021).
- 1014 • **Settings:** The exact settings for each benchmark, including the number of few-shot exam-  
 1015 ples, are provided in Table 6 and Table 7.
- 1016 • **Decontamination:** Our 13-gram decontamination protocol is detailed in Appendix G.
- 1017 • **Code:** We are including our evaluation and decontamination scripts with the submission.

1018 While model checkpoints and training logs are not included in the initial submission due to size and  
 1019 anonymity constraints, we plan to release these upon publication to facilitate future research.

## C LIMITATIONS AND BROADER IMPACT STATEMENT

**Limitations.** While the dataset improves the current state-of-the-art in the legal risk mitigation of high-performing pretraining data, upstream provenance may still contain errors with respect to licensing. We mitigate by tiering sources, providing explicit shard-level audit metadata, and applying filtering and decontamination; we encourage downstream users to select tiers consistent with their risk posture. Further automation of licensing check procedures is a subject of future work. The dataset has a scale of 422B tokens, which is not sufficient for larger-scale pre-training, and future work should investigate scaling up the presented dataset composition recipe.

**Broader Impact Statement.** The dataset improves transparency and reduces legal uncertainty for open pre-training, providing a safe ground for research, experimentation and development for the open-source community. It also can boost the trust of the general public into open-source machine-learning research that can be executed on well-validated, transparent artifacts with clear origins and widely accepted licensing schemes for broad re-use.

## D DATASET COMPOSITION AND COMPARISON

This appendix provides a detailed view of the **MixtureVitae** corpus, both in relation to other datasets and in its internal construction.

Table 3: Comparison of large-scale pretraining datasets, grouped by their licensing philosophy to provide context for our performance results. **MixtureVitae** is unique in its combination of a risk-mitigated licensing approach and the inclusion of a large subset of reasoning, coding and instruction synthetic data.

Dataset	Size (Tokens)	Primary Data Types	Licensing Philosophy
<i>Non-Permissive / Mixed-License Baselines</i>			
Nemotron-CC-HQ (Su et al., 2025)	≈ 1.1T	Web, Synthetic	Unspecified
DCLM-baseline (Li et al., 2024)	≈ 3.8T	Web, Code, Academic	Mixed / Unspecified
FineWeb-Edu (Penedo et al., 2024)	≈ 1.3T	Web (Educational)	Unspecified
The Pile (Gao et al., 2020)	≈ 183.28B	Web, Books, Code	Mixed / Unspecified
SlimPajama (Shen et al., 2024)	≈ 627B	Web, Books, Code	Mixed / Unspecified
C4 (Raffel et al., 2020)	≈ 156B	Web	ODC-BY
HPLT-2.0 (eng.) (Burchell et al., 2025)	≈ 2.86T	Web, Books, News	Mixed / Unspecified
<i>Permissive Baselines</i>			
CommonCorpus (Langlais et al., 2025)	≈ 2T	Web, Curated	Strictly Permissive
Comma-0.1 (Kandpal et al., 2025)	≈ 1T	Web, Curated	Strictly Permissive
KL3M (Bommarito et al., 2025)	≈ 580B	Web, Curated	Strictly Permissive
OLC (Min et al., 2024)	≈ 228B	Web, Curated	Strictly Permissive
<i>Our Contribution</i>			
<b>MixtureVitae</b>	≈ <b>422B</b>	<b>Web, Curated, Synthetic</b>	<b>Permissive-First, Risk-Mitigated</b>

**Shard Definitions and Mixing.** It is important to note that the dataset shards and categories listed in this appendix serve as logical groupings for transparency, licensing audits, and ablation analysis. They do not dictate a rigid sequential training order. As noted in the main text, the physical construction of training batches utilizes domain-aware packing to maximize local coherence, prioritizing the density of reasoning and factual tokens over these high-level taxonomic boundaries.

Table 3 presents a high-level comparison of **MixtureVitae** against the other prominent pretraining datasets evaluated in our experiments, detailing their respective sizes, primary data types, and licensing philosophies.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

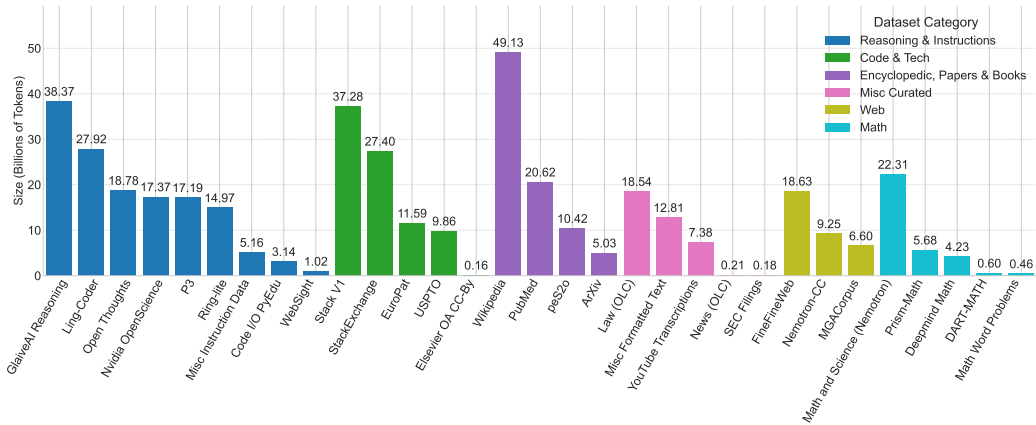


Figure 4: Detailed composition of the **MixtureVitae** dataset.

Figure 4 presents the detailed composition of the **MixtureVitae** dataset. The individual components are color-coded by their primary dataset category, as presented in the main text.

- **Code & Tech (Blue):** This domain is anchored by our largest code sources, Stack V1 and Lingo-Coder, and supplemented by StackExchange.
- **Reasoning & Instruction (Green):** The largest contributor to this category is Open Thoughts, followed by P3 and NVIDIA OpenScience.
- **Encyclopedic, Papers & Books (Purple):** This category is dominated by Wikipedia, the single largest component in the dataset. It is complemented by large-scale text from PubMed and arXiv.
- **Math (Cyan):** The math component is a diverse mixture of sources, led by the Math and Science (Nemotron) corpus and Prism-Math.
- **Web (Yellow):** Our web data is primarily sourced from corpora such as SEC Filings, MGA-Corpus, and FineFineWeb.
- **Misc Curated (Pink):** This category includes a variety of high-quality curated sources, notably Law (Open License Corpus) and YouTube Transcriptions.

## E EXPERIMENT SETUP DETAILS

To ensure full reproducibility, this appendix details the complete experimental setup. This includes the model architectures for all scales, the training hyperparameters for both 50B and 300B token budgets, and the specific settings used for all general evaluation benchmarks.

Table 4: **open-sci-ref** (Nezhurina et al., 2025) model architecture and scales. We used tied embedding weights in all experiments.

Parameters (B) (Non-Emb + Emb)	Layers	Hidden	Heads	FFN Hidden	Memory	FLOPs
0.1 + 0.03 = 0.13	22	512	8	2256	0.89 GB	$7.8 \times 10^8$
0.35 + 0.05 = 0.40	22	1024	16	3840	2.88 GB	$2.4 \times 10^9$
1.21 + 0.10 = 1.31	24	2048	32	5440	7.544 GB	$7.9 \times 10^9$
1.61 + 0.10 = 1.71	24	2048	32	8192	9.884 GB	$1.0 \times 10^{10}$

Table 5: The training schedules used in our experiments.

Tokens	Global Batch Size (tokens)	Iterations	Learning Rate	Warmup	Cooldown (20%)
50B	4.12M	11,921	$4 \times 10^{-3}$	1,000	2,384
300B	4.12M	72,661	$4 \times 10^{-3}$	25,000	14,532

## E.1 TRAINING SETUP PARAMETERS

This appendix details the exact model architectures and training hyperparameters used for all experiments, ensuring full reproducibility.

We adopt the standard architectures and scales from the **open-sci-ref** framework to allow for a fair and direct comparison against other published baselines. All models were trained with tied embedding weights.

**Model Architecture** Table 4 defines the four model scales used in our study. The columns are defined as follows:

**Parameters (B) (Non-Emb + Emb)** The total model parameters in billions, separated into **Non-Embedding** (Non-Emb) parameters (the core transformer blocks) and **Embedding** (Emb) parameters (the token lookup tables). As noted in the caption, we used tied embedding weights.

**Layers** The total number of transformer blocks stacked in the model.

**Hidden** The hidden size (or embedding dimension,  $d_{\text{model}}$ ) of the model.

**Heads** The number of attention heads in the multi-head attention mechanism.

**FFN Hidden** The inner dimension of the Feed-Forward Network (FFN) layer within each transformer block.

**Memory** The approximate VRAM required to store the model weights, in bfloat16.

**FLOPs** An approximation of the training compute cost using the **6N** rule: a standard estimate for a transformer’s forward-and-backward pass, where **N** is the number of *non-embedding* parameters (Kaplan et al., 2020).

**Training Schedules** Table 5 defines the training hyperparameters for our two main experimental runs (50B and 300B tokens). We use a single stage training with no post-training.

**Tokens** The total number of tokens in the training run.

**Global Batch Size (tokens)** The total number of tokens processed in a single training step (i.e., one gradient update) across all GPUs.

**Iterations** The total number of training steps.

**Learning Rate** The peak learning rate used during training.

**Warmup** The number of initial *iterations* (steps) over which the learning rate linearly increases from 0 to its peak value.

**Cooldown (20%)** The number of final *iterations* (the last 20% of training) over which the learning rate decays to zero.

## E.2 EVALUATION SETTINGS

We used the `lm-evaluation-harness` (Gao et al., 2021) for all general evaluations. The specific tasks and few-shot counts are detailed in Table 6. The settings for the reasoning tasks (e.g., GSM8K, IFEval) are listed separately in Table 7.

Table 6: General evaluation benchmark settings. All tasks use Accuracy as the primary metric.

Task	Citation	# of Shots
MMLU	Hendrycks et al. (2021)	5
HellaSwag	Zellers et al. (2019)	10
CommonSenseQA	Talmor et al. (2019)	10
ARC-Challenge	Clark et al. (2018)	10
ARC-Easy	Clark et al. (2018)	10
PIQA	Bisk et al. (2020)	10
BoolQ	Clark et al. (2019)	10
Winogrande	Sakaguchi et al. (2021)	0
OpenBookQA	Mihaylov et al. (2018)	0
COPA	Roemmele et al. (2011)	0
LAMBADA	Paperno et al. (2016)	0

Table 7: Evaluation settings for reasoning tasks. All tasks use Accuracy as the primary metric. To execute the evaluation, we used LM Evaluation Harness Gao et al. (2021).

Task	Citation	# of Shots
GSM8k	Cobbe et al. (2021)	4
IFEval	Zhou et al. (2023)	0
MBPP	Austin et al. (2021)	4

## F ADDITIONAL EXPERIMENTS

This appendix provides additional experimental results to supplement the findings presented in the main paper. We offer a more granular breakdown of the 300B token experiment, analyze performance at a smaller 50B token scale to assess the generalization of our results, and report the results of a model red-teaming analysis to evaluate the model’s safety profile.

### F.1 300B EXPERIMENT - DETAILED RESULTS

The detailed results for each evaluated task (contributing to the average over 10 tasks as shown in Figure 3 ) are given in Figure 5.

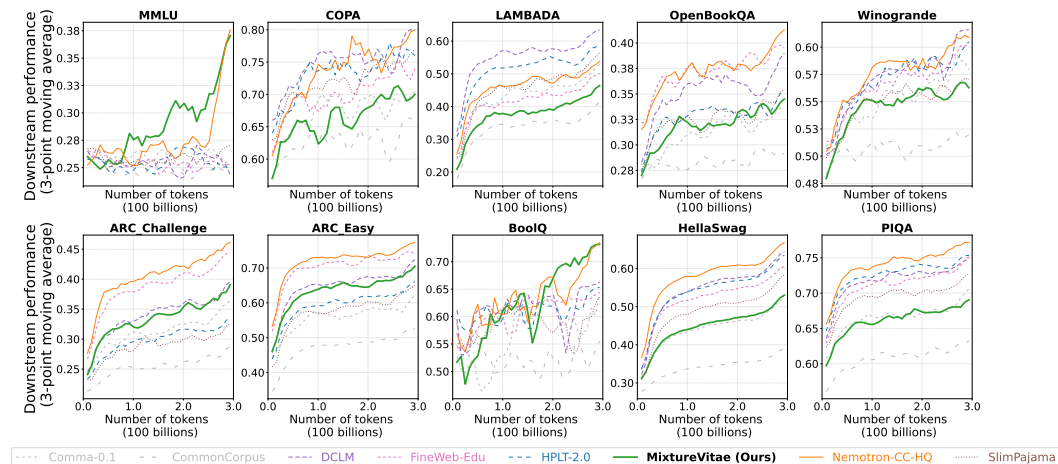


Figure 5: Comparing performance of 1.7B model trained on **MixtureVitae** and baseline datasets for a 300B token budget. While some evaluations provide clear dataset rankings (e.g. ARC, HellaSwag, Lambada), others do not provide a good signal for dataset comparison, on an individual basis.

## F.2 PERFORMANCE AT 50B TOKENS SCALE.

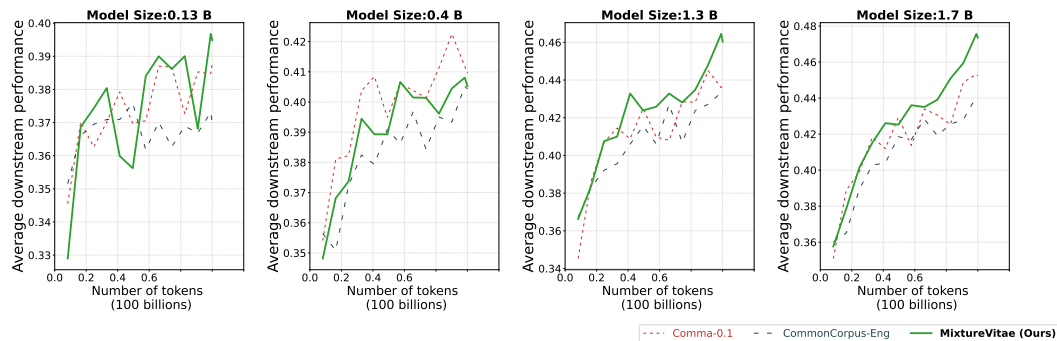


Figure 6: Average performance of permissive datasets after 50B training tokens. **MixtureVitae** shows an early and consistent lead at larger model scales.

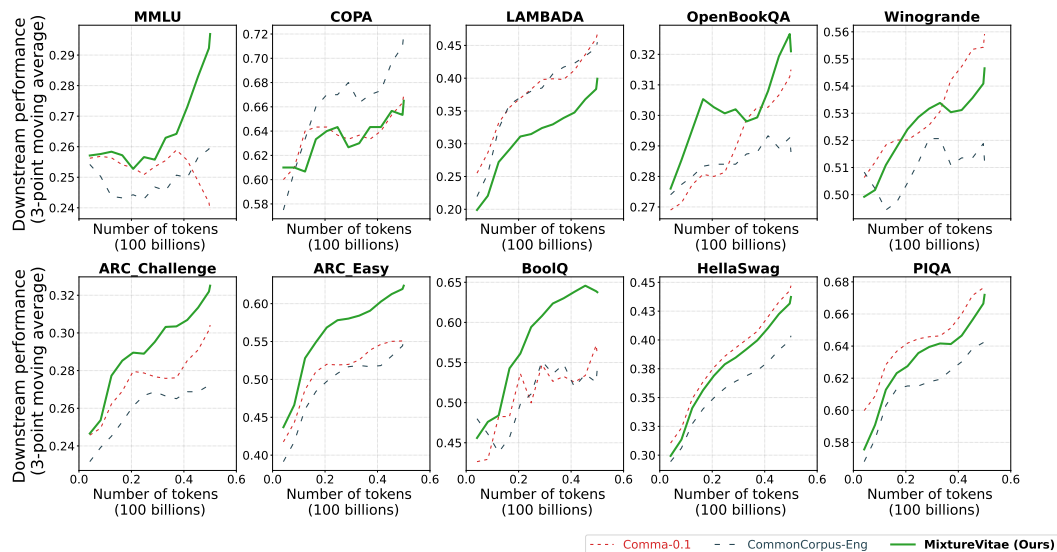


Figure 7: Per-benchmark performance of permissive datasets after 50B training tokens. **MixtureVitae**'s advantage on MMLU is apparent even at this early stage.

To assess performance on a smaller reference tokens scale, we also evaluated models trained on a 50B token subset of each dataset. The results, shown in Figure 6 and Figure 7, indicate that the advantages of **MixtureVitae** manifest already at the smaller token scales. Figure 6 shows that **MixtureVitae** establishes a consistent performance lead over other permissive datasets within the first 50B tokens, especially at the 1.3B and 1.7B model scales. The per-benchmark analysis further reinforces this finding (see Figure 7). On MMLU, **MixtureVitae** is the only permissive dataset to show a significant learning signal early in training, demonstrating that its composition provides immediate benefits, which might be both due to knowledge rich and instruction like content. Arguably, this suggests that the reasoning capability shown by **MixtureVitae** is not a late-stage phenomenon but rather an indication of efficient instillation from the early stages of training. This strong initial performance underscores the learning efficiency of **MixtureVitae**, making it a compelling choice for achieving high performance with less computational cost.

## F.3 MODEL RED TEAMING

To evaluate the safety of the model trained on **MixtureVitae** for 300B tokens, we performed a red-teaming analysis to measure the Attack Success Rate (ASR) against three standard benchmarks:

**ToxiGen** (Hartvigsen et al., 2022), **Do-Not-Answer** (Wang et al., 2024), and **AdvBench** (Zou et al., 2023). The results (Table 8) shows that our model is competitive with the baselines.

The model responses were evaluated using two safety classifiers: (i) **Llama Guard-8B** (Inan et al., 2023), used to evaluate the **Do-Not-Answer** and **AdvBench** datasets, while (ii) the **toxigen\_roberta** classifier (Logacheva et al., 2022) was used for the **ToxiGen** benchmark.

Table 8: Attack Success Rate in %, lower is better. All models are trained with the same **open-sci-ref** procedure (300B-token budget) while varying only the pretraining dataset.

Benchmark	MixtureVitae	Comma	CommonCorpus-Eng	Nemotron-HQ-CC
ToxiGen	8.07	9.04	12.77	10.21
Do-Not-Answer	28.22	24.71	21.62	20.98
AdvBench	86.92	92.12	70.58	85.77

## G CONTAMINATION ANALYSIS

### G.1 CONTAMINATION DETECTION PROTOCOL

To ensure the integrity of our evaluation, we implemented a comprehensive decontamination protocol to measure the overlap between our training dataset and all evaluation benchmarks we report results on. This protocol consists of three main stages: Index Construction, Dataset Scanning, and Leakage Reporting.

#### G.1.1 INDEX CONSTRUCTION

The first stage creates a compact, indexed set of unique n-grams from all benchmark evaluation data.

- Text Normalization:** All text from the benchmarks is processed through a normalization pipeline, similar to Laurençon et al. (2022): (1) Unicode normalization (NFKC), (2) conversion to lowercase, (3) tokenization, and (4) removal of a predefined list of common English stop words. This procedure focuses the resulting n-grams on substantive content.
- N-gramming and Filtering:** We generate 13-grams, a common n-gram size for this task Brown et al. (2020); Gao et al. (2020) from the normalized token lists. As in Laurençon et al. (2022), a set of regular expressions is used to filter out common boilerplate, exam instructions, and formatting artifacts.
- Train/Test De-duplication:** as in Gao et al. (2020), we compute the set of all 13-gram hashes from the `train` split and subtract this set from the 13-gram hashes generated from the `test` split. This ensures our index only contains n-grams that are unique to the evaluation set.

#### G.1.2 DATASET SCANNING

The second stage analyzes the target training dataset against the generated index.

- Document Processing:** Each document in the training dataset is processed using the *exact same* normalization, 13-gramming, and hashing pipeline used for index construction.
- Contamination Criteria:** A document is flagged as "contaminated" if it meets two criteria, based on the set intersection of its n-gram hashes with the benchmark index:
  - Minimum Hits:** The number of distinct matching n-grams is  $\geq 3$ .
  - Minimum Coverage:** As proposed in Rae et al. (2022), the coverage of matching n-grams is  $\geq 0.1\%$ . Coverage is defined as:

$$\text{Coverage} = \frac{\text{distinct\_hits}}{\text{total\_unique\_13grams\_in\_doc}}$$

### 1350 G.1.3 LEAKAGE REPORTING

1351

1352 The final stage aggregates the scan results into a summary report.

1353

1354 1. **Numerator (Leaked N-grams):** The procedure aggregates the reports from all scanned  
 1355 partitions. It performs a global *set union* to find all unique n-gram hashes that were found  
 1356 *at least once* in the target dataset, aggregated by benchmark source. This provides the  
 1357 `unique_ngrams_leaked` count for each benchmark.

1358 2. **Denominator (Total N-grams):** The procedure retrieves the pre-computed metadata to  
 1359 obtain the total unique n-gram count for each benchmark.

1360 3. **Final Metric:** As proposed in Touvron et al. (2023), the **Leak Percentage** for each bench-  
 1361 mark is then calculated as:

1362

$$1363 \text{ Leak Percentage} = \frac{\text{unique\_ngrams\_leaked}_{\text{benchmark}}}{\text{total\_unique\_ngrams\_in\_index}_{\text{benchmark}}} \times 100$$

1365

### 1366 G.2 CONTAMINATION REPORT

1367

1368 We executed our 13-gram contamination scan across the entire 345 697 271 documents of the  
 1369 **MixtureVitae** dataset. The global summary of contaminated documents per benchmark is presented  
 1370 in Table 9.

1371 The results confirm that for the vast majority of benchmarks—including ARC, HellaSwag, LAM-  
 1372 BADA, OpenBookQA, and PIQA—the document-level contamination rate is negligible (at or below  
 1373 0.0003%), strongly validating the integrity of our evaluation on these tasks.

1374

1375 The scan did, however, flag a minor overlap for MMLU (0.0098%) and BoolQ (0.0087%), and  
 1376 a more significant overlap for our key code benchmarks: HumanEval (0.0988%) and MBPP  
 1377 (0.0878%). This overlap in code benchmarks is a known challenge when including large-scale  
 1378 permissive code corpora like The Stack, which may naturally contain snippets of common coding  
 1379 problems (a “source overlap” rather than a direct “test-set leak”).

1380 To ensure this overlap did not artificially inflate our model’s strong performance on these key tasks,  
 1381 we conducted case studies for the benchmarks with the highest overlap. This analysis is detailed in  
 1382 the following section (Appendix G.3).

1383 Table 9: Global contamination summary by document count, based on a 13-gram overlap scan. This  
 1384 table shows the total number of documents in **MixtureVitae** that contained at least one overlapping  
 1385 n-gram from each benchmark’s test set. The total documents in **MixtureVitae** is 345 697 271 and  
 1386 the overall contamination rate is 0.1420%.

1387

Benchmark	Contaminated Docs	Contamination Rate (%)
ALERT	12	0.0000%
ARC	17	0.0000%
BoolQ	30 144	0.0087%
CommonSenseQA	0	0.0000%
GPQA	1077	0.0003%
GSM8K	230	0.0001%
HellaSwag	186	0.0001%
HumanEval	341 554	0.0988%
IfEval	756	0.0002%
LAMBADA	23	0.0000%
MBPP	303 558	0.0878%
MMLU	33 922	0.0098%
OpenBookQA	60	0.0000%
PIQA	5	0.0000%
SimpleQA	98	0.0000%

1403

### G.3 DECONTAMINATED TEST SET PERFORMANCE

To understand how test data leakage affects final performance on downstream tasks, we conducted the following experiment on all models and benchmarks reported in Section 3.3:

1. Identify problems from the test set that have at least one 13-gram match in the training dataset.
2. Evaluate the model on a decontaminated benchmark version obtained by removing problems that were identified.

Table 10: Validating math, code, and instruction performance by comparing original (Orig) vs. decontaminated (Decont) test sets for 1.7B models trained for 300B tokens. **MixtureVitae**’s high scores are shown to be genuine, as performance is maintained after removing all overlapping test items. This confirms the model’s capabilities are not an artifact of test set leakage.

Training Dataset	GSM8K		GSM8K-CoT		MBPP		MBPP+		IFEval	
	Orig	Decont	Orig	Decont	Orig	Decont	Orig	Decont	Orig	Decont
<b>MixtureVitae</b>	0.53	0.54	0.50	0.50	0.38	0.38	0.55	0.59	0.19	0.23
SmolLM2	0.30	0.30	0.28	0.29	0.35	0.35	0.48	0.48	0.17	0.20
Comma-0.1	0.06	0.06	0.09	0.09	0.21	0.23	0.28	0.28	0.18	0.20
CommonCorpus-Eng	0.02	0.01	0.01	0.01	0.02	0.02	0.04	0.05	0.12	0.16
C4	0.01	0.01	0.01	0.02	0.00	0.00	0.00	0.00	0.20	0.21
DCLM	0.01	0.02	0.02	0.02	0.01	0.00	0.02	0.02	0.12	0.13
FineWeb	0.02	0.01	0.03	0.03	0.00	0.00	0.00	0.00	0.18	0.20
HPLT	0.02	0.02	0.02	0.02	0.00	0.00	0.00	0.00	0.17	0.21
Nemotron-CC-HQ	0.03	0.02	0.03	0.03	0.00	0.00	0.00	0.00	0.09	0.10
SlimPajama	0.02	0.02	0.02	0.02	0.00	0.00	0.00	0.00	0.14	0.15

Table 11: Benchmark test set sizes (number of examples) for the original benchmarks versus the final decontaminated versions. The ‘Decontaminated’ column shows the reduced set size after removing all examples with detected 13-gram training data overlap.

Dataset	Original	Decontaminated
MBPP	500	331
IFEval	541	429
GSM8K	1319	1235
MBPP+	378	339

As we can see from Table 10, the performance of the evaluated models is consistent between the original and decontaminated versions, aside from some upward bias in the decontaminated versions of MBPP+ and IFEval. Crucially, the result for the model trained on **MixtureVitae** was not affected by the strict decontamination procedure applied to the benchmarks. This rules out the possibility that the strong performance of **MixtureVitae** on math and coding is due to benchmark leakage.

### G.4 DECONTAMINATION CASE STUDY

To further alleviate concerns about contamination issues, we performed an experiment where we trained a 1.7B model on a version of **MixtureVitae** that excludes three dataset shards that contribute 27% of the total contaminated docs— which in particular showed MMLU contamination rates that are high relative to the rest of the dataset, as shown in Table 12. The shards we removed were **Misc-Instruct**, **DART-Math** and **Nemotron Science & Math**. The results are shown in Figure 8 and demonstrate that removing these shards had no effect on MMLU performance.

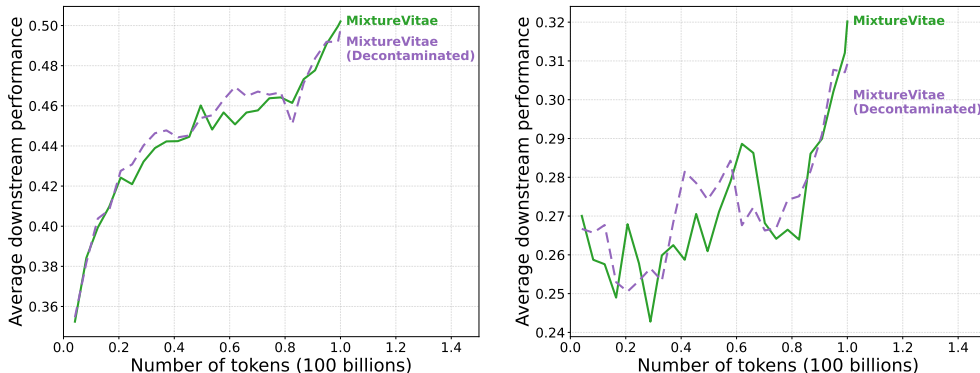
### D.5 Full Decontamination Experiment

We also performed a stronger decontamination experiment in which we removed *every* document in **MixtureVitae** that was flagged as contaminated by our 13-gram procedure (Appendix G.2) and retrained a 1.7B model for 300B tokens under the `open-sci-ref` setup. The results—shown in Figure 9—indicate that the fully decontaminated variant performs slightly *better* than the original

1458 Table 12: Contamination sources for the MMLU benchmark in **MixtureVitae**, sorted by the number  
 1459 of contaminated documents, high to low.

Dataset Shard	Contaminated Docs
Misc-Instruct	14 649
DART-Math (Tong et al., 2024)	11 102
Nemotron Science & Math (Bercovich et al., 2025)	4793
MGACorpus (Hao et al., 2025)	241
(All Remaining)	3137

1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480



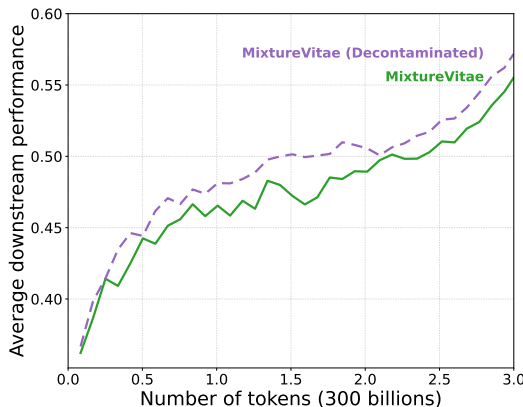
1481 (a) Average accuracy across all tasks (as listed in Table 6) as a function of number of training steps.  
 1482 (b) Accuracy on MMLU as a function of number of training steps.

1483  
 1484 **Figure 8: Validation of 1.7B model performance.** The **MixtureVitae (Decontaminated)** model  
 1485 (purple, dashed), trained with dataset shards responsible for benchmark leakage removed, performs  
 1486 closely to the full **MixtureVitae** (green, solid) model. This confirms our results are not an artifact  
 1487 of test set leakage.

1488  
1489

1490 **MixtureVitae** model, further addressing concerns that our benchmark results might be inflated by  
 1491 data leakage.

1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506



1507 **Figure 9: 1.7B model performance on a fully decontaminated dataset.** The model trained on the  
 1508 fully decontaminated **MixtureVitae** corpus (purple, dashed) performs slightly better than the model  
 1509 trained on the full **MixtureVitae** dataset (green, solid), further indicating that benchmark gains are  
 1510 not driven by contaminated examples.

1511

1512 G.5 DISCUSSION ON DECONTAMINATION METHODOLOGY

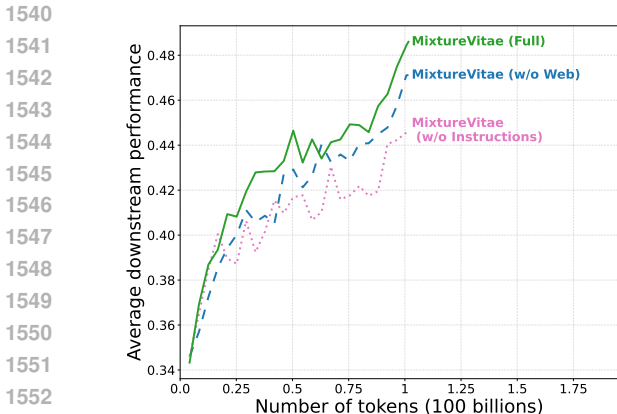
1513  
 1514 Our decontamination pipeline employs the standard 13-gram exact-match procedure (Abdin et al.,  
 1515 2024) to ensure high precision, scalability and comparability with prior baselines. While we ac-  
 1516 knowledge that exact matching overlooks paraphrased content, we avoided approximate methods  
 1517 (e.g., LSH, embedding-based) due to their tendency to produce false positives on common factual  
 1518 or algorithmic templates (Lee et al., 2022). As noted in Penedo et al. (2024), aggressive removal  
 1519 of semantically similar content risks distorting the training distribution by discarding high-value  
 1520 instructional data.

1521 H ABLATION STUDIES

1522 To isolate the impact of primary data components in **MixtureVitae**, we define **Web** and **Instructions**  
 1523 subsets (see Figure 1; **Instructions** encompasses Reasoning & Instruction and Math parts of the full  
 1524 mixture) and conduct an ablation study on a 100B-token scale. We train three separate models:  
 1525 (1) **MixtureVitae (full)**, the complete dataset; (2) **MixtureVitae (w/o Web)**, removing the **Web**  
 1526 component; (3) **MixtureVitae (w/o Instructions)**, removing the **Instructions** component.

1527 The average downstream performance of these models (Figure 10a) shows varying contributions by  
 1528 each component: The **Instructions** data is the most critical driver of performance, as its removal  
 1529 results in the largest, consistent drop of average performance compared to other configurations.  
 1530 Removing **Instructions** particularly leads to severe drop on GSM8k (from 0.47 to 0.03) and MBPP,  
 1531 as shown in Figure 10b. Absent the **Instructions** data, the model fails to match the gains of the full  
 1532 mix, underscoring the essential role of instruction-following data in generalization.

1533 Removing the **Web** component (**w/o Web**, blue dashed line) also results in a performance drop  
 1534 below the full dataset, albeit less dramatically. Figure 10b shows a drop from 0.47 to 0.41 on  
 1535 GSM8k, far less severe than the drop close to 0 for **w/o Instructions** and only slight changes on  
 1536 code evals. The comparison of ablation effects again highlights the **crucial role of instruction and**  
 1537 **reasoning data in achieving high performance.**



1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552 (a) Ablation on full **MixtureVitae** against two ver-  
 1553 sions, each excluding a data subset as indicated by  
 1554 w/o. Average performance on 10 downstream tasks.  
 1555  
 1556

(b) A performance breakdown on math, coding and instruction following tasks for the ablated dataset variants. Best results are in **bold**. Numbers in **red** indicate strong performance drop.

Training Dataset	IF-Eval	GSM8K	MBPP	Average
<b>MixtureVitae</b>	0.14	<b>0.47</b>	<b>0.34</b>	<b>0.25</b>
<b>MixtureVitae</b> (w/o Web)	0.18	0.41	0.33	<b>0.25</b>
<b>MixtureVitae</b> (w/o Instructions)	<b>0.19</b>	<b>0.03</b>	<b>0.14</b>	0.14

1557 Figure 10: An ablation study on components of the **MixtureVitae** dataset. Figure 10a shows perfor-  
 1558 mance average on 10 downstream evals during training, while Figure 10b shows scores on further  
 1559 separate math, code and instruction benchmarks. The evaluation setup is given in Table 7.  
 1560

1561 I SYNTHETIC MATH DATA GENERATION

1562 The synthetic math dataset was programmatically generated to produce a diverse range of mathe-  
 1563 matical problems and their solutions. The generation process covers a wide array of mathematical  
 1564 domains, including fundamental arithmetic operations, multi-term fractional expressions, and the  
 1565

step-by-step solution of algebraic linear equations. A key component of the dataset consists of word problems, where numerical challenges are embedded in narrative scenarios.

A significant feature of this generation pipeline is the creation of detailed, step-by-step solutions formatted as a chain-of-thought. For many problem categories, the scripts produce a human-readable explanation of the entire solution process. This is achieved by using a variety of randomized natural language templates to describe each logical step, such as carrying a digit in addition or isolating a variable in an equation.

Following the initial generation, a final post-processing step is applied to format the dataset for model training. This stage programmatically identifies data entries containing human-like, explanatory text by searching for common instructional words. For these selected entries, a descriptive header (e.g., "Here are examples of addition, division exercises") is dynamically generated. The content and phrasing of this header are randomized and based on the mathematical operations present within the text, adding significant linguistic diversity.

For example, in the generated math problem below, a model may be able to generalize to new numbers, but if the problem were to add three students instead of two, the model may not be robust enough to generalize. We leave this analysis for future work.

The age difference between Sarah and Asaf's age is half the total number of pencils Sarah has. The sum of their ages is 132, and Sarah is 27 years old. If Asaf has 60 more pencils than Sarah, calculate the total number of pencils they have together. Solution: If the sum of their ages is 132, and Sarah is 27 years old, Asaf is 105 years old. The age difference between Sarah and Asaf's age is  $105 - 27 = 78$ . Since the age difference between Sarah and Asaf's age is half the total number of pencils Sarah has, Sarah has  $2 * 78 = 156$  pencils. If Asaf has 60 more pencils than Sarah, Asaf has  $156 + 60 = 216$  pencils. Together, they have  $156 + 216 = 372$  pencils.

## J OUR POSITION ON USING GOVERNMENTAL AND OTHER WORKS UNDER FAIR USE AND RELATED ETHICAL AND LEGAL BASIS

To contextualize our licensing tiers and clarify the rationale behind including certain higher-risk but legally supportable sources, we outline here the ethical and legal considerations underlying [MixtureVitae](#)'s construction. Our goal is not to offer legal advice or definitive interpretations of copyright law, but rather to articulate the principles—fair use, permissive upstream licensing, government-works doctrine, and the EU text-and-data-mining (TDM) (European Union, 2019; Margoni & Kretschmer, 2022) exception—that inform our "permissive-first, risk-mitigated" design philosophy. We provide this discussion so downstream users can understand how specific dataset subsets were evaluated and what residual risks remain despite our filtering and provenance-tracking efforts.

### J.1 FAIR USE OF GOVERNMENT WORKS

In order to increase the diversity of our dataset, we included  $\approx 11\text{B}$  tokens of governmental website data from US federal, US non-federal, and non-US government sources. While works created by the US federal government are generally not copyrightable, other governmental website content may neither be expressly in the public domain nor explicitly licensed. For those sources, we rely on fair use principles (Congress (1976); Lemley & Casey (2017)) and the EU text and data mining exceptions (European Union (2019); Margoni & Kretschmer (2022)), which together mitigate the risk associated with using this subset.

Our ethical and legal reasoning for using this government web content—sourced from Common Crawl-related datasets (Common Crawl Foundation, 2025) that respect *robots.txt* opt-out—is as follows:

- 1620 • **Public Purpose Alignment:** The content created by governments is normally meant to be  
1621 shared with the public, and by using the data for training we are assisting this purpose.
- 1622 • **Purpose of Use:** From a legal perspective, the government works are being redistributed  
1623 as part of an open source, no-fee dataset, used to create models are less likely to violate  
1624 copyright. This purpose is clearly not to compete with the government’s own usage.
- 1625 • **Effect on Potential Market:** Our use of government website content is unlikely to affect  
1626 any potential market for that content, as governments typically do not exploit these mate-  
1627 rials commercially in ways that would compete with our dataset or downstream models.  
1628 This factor favors a finding of fair use.
- 1629 • **Nature of the Content:** The nature of the content is mostly public announcements, content  
1630 of public interest, governmental functions or the like. Again, we believe there is strong  
1631 public policy interest for fair use of this type of information.
- 1632 • **Amount Used:** While we use all or almost all of the content of the government website,  
1633 the amount of usage is not determinative of fair-use or not fair-use.
- 1634 • **Federal vs. Non-Federal Works:** Lastly, US works created by the federal governments  
1635 are generally not copyrightable. However, we recognize that this is not the case for other  
1636 foreign governmental works, or non-federal works.

1638  
1639 For these reasons, we believe using government website data presents relatively lower copyright  
1640 risk. To further minimize risk—for example, the potential inclusion of third-party copyrighted works  
1641 embedded in government web pages—we apply keyword filters such as “All Rights Reserved” and  
1642 “Copyright ©” to exclude pages that contain such terms.

1643 Recent court cases, as of the writing of this paper, include:

- 1644 • **Bartz v. Anthropic PBC:** district court ruling that use of purchased copies of books for AI  
1645 training is fair use.
- 1646 • **Kadrey v. Meta Platforms, Inc.:** district court ruling that training on authors’ books was  
1647 transformative fair use.

1648  
1649 These developments lend some support to the argument that AI training on web-text data—including  
1650 our relatively small, public-facing government subset—can fall within fair use, though the case law  
1651 is still evolving.

## 1652 J.2 OTHER TIER-2 DATA WITH OPAQUE OR MIXED PROVENANCE

1653  
1654 Similarly, our dataset includes data whose provenance is not entirely transparent even though the  
1655 license on the upstream dataset appears permissive, such as `The Stack V1` and other Tier-2  
1656 sources identified in Appendix M. In the case of `The Stack V1`, a line-by-line audit to remove  
1657 copyrighted content has not been performed, and therefore some risk remains in its usage. Nonethe-  
1658 less, we rely on fair use to justify this usage because the data are used to train models, rather than  
1659 to provide a substitutive or competing software product. For a more detailed discussion of `The`  
1660 `Stack V1`, see Section K.

1661  
1662 For other Tier-2 data, some upstream generator models impose conditions on downstream use—such  
1663 as the Llama license, which requires model users to adhere to certain limitations. We do not believe  
1664 we are bound by terms that were not contractually passed through to us by our direct licensor,  
1665 although this issue is subject to debate. We therefore classify this small portion ( $\approx 4\%$ ) of the dataset  
1666 as **Tier 2(b)**.

1667 There are additional Tier-2 data where the provenance is partially opaque. For example, a small por-  
1668 tion of our P3 dataset, when converted into a few-shot format, may pose higher risk than other Tier-1  
1669 data. While the ultimate source datasets that constitute P3 are well-known academic benchmarks,  
1670 some of those component datasets do not provide explicit licenses. Nonetheless, we consider the  
1671 resulting few-shot datasets to be highly transformative and unlikely to compete with the underlying  
1672 works: they are mixed and reformatted multiple times for the specific purpose of training classifi-  
1673 cation and few-shot models, rather than, for example, serving as standalone product reviews. We  
classify these higher-risk works as **Tier 2(b)** and include them in our dataset with that caveat.

1674 J.3 RELIANCE ON EU TEXT AND DATA MINING EXCEPTIONS  
1675

1676 We also rely, to some extent, on the EU text and data mining exception (European Union, 2019)  
1677 for our inclusion of web-crawled data. This regime is complementary to US fair-use doctrine, and  
1678 we mention it here for completeness. In particular, we depend on Common Crawl’s practice of  
1679 respecting *robots.txt* at the time of crawling. We do not believe retroactive recrawling is legally  
1680 necessary to determine whether a work was subsequently opted out, but we nonetheless commend  
1681 efforts towards doing so, such as Apertus (Hernández-Cano et al., 2025).

1682  
1683 J.4 RESIDUAL COPYRIGHT AND TRADEMARK RISKS  
1684

1685 The copyright risks in machine learning are complex. For example, copyrighted materials may  
1686 appear as limited fair-use quotations in Wikipedia articles<sup>3</sup>. A model trained on such materials in  
1687 the aggregate could, in principle, generate more substantial and potentially infringing text than the  
1688 short quotations present in the dataset. Future work should address this risk, including (i) copyright  
1689 evaluation audits of datasets, and (ii) model-level mitigations that encourage limited direct quotation  
1690 and discourage reproduction of substantial protected passages.

1691 As with other large, permissively licensed datasets, additional legal risks remain, including trade-  
1692 mark risks. For instance, while training on a Wikipedia article about “Spiderman” may be relatively  
1693 low risk (given its CC-BY-SA license and the educational, summarizing nature of the article), a  
1694 model that subsequently generates new stories featuring the character name “Spiderman”—even if  
1695 the plots themselves are not derived from existing human-created stories—may still implicate trade-  
1696 mark rights. Addressing those issues thoroughly is beyond the scope of this work and is left for  
1697 future research.

1698 We do not and cannot guarantee that, even with rigorous provenance tracking and standard filtering,  
1699 the dataset is free of legal risk. Nothing in this section constitutes legal advice. We recommend that  
1700 anyone who uses our datasets consult their own legal counsel in their jurisdiction before deploying  
1701 models trained on this data in commercial settings.

1702  
1703 K PROVENANCE AND RATIONALE FOR THE STACK v1 (OPENRAIL-M AND  
1704 TERMS OF USE)  
1705

1706 Our inclusion of 53.2B tokens sourced from **The Stack v1** (Kocetkov et al., 2023), which we catego-  
1707 rize by its governing dataset card terms of use and which subsequent model uses the **OpenRAIL-M**  
1708 license, warrants this specific note on provenance. The data was included based on the following  
1709 rationale:

- 1710 • **Source and Filtering Methodology:** The dataset originates from a large-scale scrape of  
1711 GitHub. The BigCode project curated this data by applying a filter to include only those  
1712 repositories that contained a clear permissive license file (e.g., MIT, Apache 2.0, BSD) at  
1713 the root level.
- 1714 • **Acknowledged Heuristic:** This repository-level filtering is a *heuristic* and not a file-level  
1715 guarantee. As acknowledged by the dataset’s creators, this process cannot perfectly re-  
1716 solve complex cases of multi-licensing within a single repository, such as the inclusion of  
1717 non-permissively licensed vendor libraries or mixed-license assets alongside permissively-  
1718 licensed code.
- 1719 • **Inclusion Justification:** Despite this caveat, The Stack v1 represents the largest-available  
1720 public corpus curated with the *explicit goal* of permissive filtering. Excluding it would  
1721 make training a high-performance, open, and risk-mitigated code model nearly impossi-  
1722 ble. Its “best-effort” permissive curation philosophy directly aligns with our dataset’s core  
1723 principle of risk-mitigation.

1724  
1725 Thus we include it in our dataset with the classification of Tier-2, as defined in Section 2.1.4.  
1726

1727 <sup>3</sup>[https://en.wikipedia.org/wiki/Wikipedia:Quotations#Copyrighted\\_](https://en.wikipedia.org/wiki/Wikipedia:Quotations#Copyrighted_material_and_fair_use)  
material\_and\_fair\_use

## 1728 L DATA FILTERING REASONING AND PROTOCOL

1729  
1730 To promote transparency, we describe our protocol for defining and checking the lists and content  
1731 of the pseudo-crawled portion of [MixtureVitae](#).

### 1733 L.1 GOVERNMENTAL AND NGO DOMAIN PATTERNS

1735 The following list of URL patterns was used to filter for governmental, non-governmental, and in-  
1736 ternational organization websites from the web datasets. We gathered the list by examining public  
1737 records, Wikipedia lists, and the like. The list is not as simple as **gov**, because international govern-  
1738 ments use different TLDs. Moreover, some spam websites masquerades as **.gov** websites. Two of the  
1739 authors examined each domain either online or through the *Internet Archives' Wayback Machine* to  
1740 confirm they belonged to a government website. After performing a pseudo-crawl on FineFineWeb,  
1741 Nemotron-CC and MGACorpus, the authors manually audited the data for quality, and filtered out  
1742 spam websites with similar website names, which were added to blocklists.

1743 The **.gov**, **.gov/**, and **.mil/** websites are US Federal governmental works. To the extent we could, we  
1744 filtered any sites that had keywords indicating reservations of rights. We believe this lowers the risk  
1745 of inadvertent third party copyrighted works appearing on US Federal works, and is in the spirit of  
1746 the EU text data mining opt-out conventions. We also note that the ultimate source of these websites  
1747 is from Common Crawl which already also respects the *robots.txt* opt-out.

- 1748 • gov (as a suffix)
- 1749 • gov/
- 1750 • mil/

1752 All other websites in this category are specifically international governments or NGOs.

1753 vlada.mk, vlada.cz, kormany.hu, regeringen.\*, rijksoverheid.nl, government.nl,  
1754 bund.de, bundesregierung.de, government.ru, gc.ca, admin.ch, www.gob.cl/  
1755 www.gob.ec/, guatemala.gob.gt/, presidencia.gob.hn/, www.gob.mx/  
1756 presidencia.gob.pa/, www.gob.pe/, gob.es/, argentina.gob.ar/, tanzania.go.tz/  
1757 indonesia.go.id/, go.kr/, go.jp/, thailand.go.th/, europa.eu/, un/, int/, govt.,  
1758 www.gub.uy, gov., gouv.

### 1760 L.2 CURATED PERMISSIVE DOMAIN LIST

1762 The following list of approximately 50 domains was curated based on their known public domain  
1763 or CC-BY-SA\* license status or a permissive status. The websites were chosen for their diversity of  
1764 content. Two of the authors—one of which has a legal background—examined the websites' terms  
1765 of use, or relevant sections online or on the Way Back Machine to confirm licensing and permission  
1766 status. After performing a pseudo-crawl on FineFineWeb, Nemotron-CC and MGACorpus, the  
1767 authors manually reviewed the data for quality, and filtered out spam websites with similar website  
1768 names as the below. These spam sites were added to blocklists.

1769 free.law, europeana.eu, publicdomainreview.org, wisdomcommons.org,  
1770 intratext.com, mediawiki.org, wikimedia.org, wikidata.org, wikipedia.org,  
1771 wikisource.org, wikifunctions.org, wikiquote.org, wikinews.org,  
1772 wikivoyage.org, wiktionary.org, wikibooks.org, courtlistener.com/<sup>4</sup>, case.law,  
1773 pressbooks.oer.hawaii.edu, huggingface.co/docs, opencourselibrary.org,  
1774 medbig.org, doabooks.org, bccampus.ca, open.umn.edu/opentextbooks,  
1775 www.gutenberg.org, mozilla.org, www.eclipse.org, apache.org, python.org,  
1776 pytorch.org, numpy.org, scipy.org, opencv.org, scikit-learn.org, pydata.org,  
1777 matplotlib.org, palletsprojects.com, sqlalchemy.org, pypi.org, sympy.org,  
1778 nltk.org, scrapy.org, owasp.org, creativecommons.org, wikia.com, foodista.com,  
1779 fandom.com, attack.mitre.org

1778 The vast majority of these sites are CC-BY licensed. However, there are some that have other open  
1779 licenses as shown in Table 13.

1781 <sup>4</sup>For [courtlistener.com](#), the terms of use says it is CC-BY-ND, but the underlying court cases are public domain, and the content from this website is merely 176KB and is de minimis.

Table 13: Software Licenses and Associated Websites

License	Websites
BSD 3-Clause	scipy.org, sympy.org, matplotlib.org, scrapy.org, scikit-learn.org, pydata.org, pytorch.org, palletsprojects.com
Mozilla Public License	mozilla.org
Python Software Foundation License 2.0	python.org
Apache 2.0	apache.org, nltk.org, opencv.org
MIT License	sqlalchemy.org
Eclipse Public License	www.eclipse.org
MedBiquitous Standards Public License	medbiq.org

## M SYNTHETIC DATA SOURCE PROVENANCE

To ensure full transparency regarding the “permissive-first” nature of **MixtureVitae**, we provide a detailed provenance audit of our synthetic data components in Table 14, including classification to tiers as defined in Section 2.1.4.

To validate the robustness of our permissive-first strategy, we further analyze the contribution of synthetic components categorized as **Tier 2(b)**. This small part of **MixtureVitae** is comprised of subsets which are permissively licensed (e.g., Apache 2.0) but are derived from generator models with restrictive community licenses (such as Llama-3) or seed data with partially opaque origins. As illustrated in Figure 11, removing these Tier 2(b) components yields a training trajectory indistinguishable from the full **MixtureVitae** baseline. This result confirms that our model’s strong performance is driven by its core, fully verifiable permissive sources, ensuring that users with strict compliance requirements can safely exclude Tier 2(b) data without compromising downstream quality.

See Section J for a further discussion on our justification for including Tier 2 and in particular Tier 2(b) data.

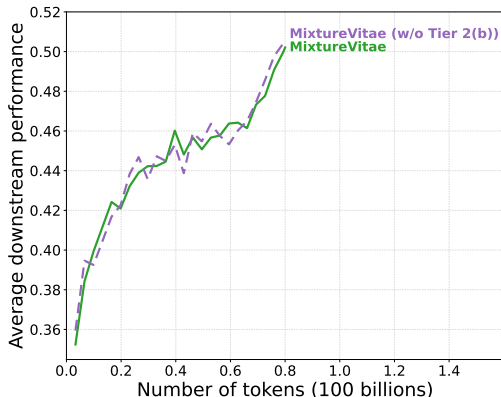


Figure 11: **Ablation of Tier 2(b) components.** We compare the training trajectory of the full **MixtureVitae** dataset (green) against a version excluding Tier 2(b) (purple dashed). Tier 2(b) consists of synthetic data derived from non-permissive generators (e.g., Llama-3) or seeds with opaque provenance. The nearly identical performance curves demonstrate that users requiring strict permissive compliance can exclude these components with negligible impact on downstream model quality.

## N SCALING OUTLOOK AND FUTURE DIRECTIONS

While **MixtureVitae** currently comprises **422** billion tokens—a scale smaller than frontier runs which often exceed 10 trillion tokens—our primary objective in this work was to establish a proof-

Table 14: Detailed provenance of synthetic data sources in **MixtureVitae**.

	Dataset Name	Dataset License	Model	Seed Data Provenance	Token Count(B)	Notes
1836						
1837	<b>Tier 1: Fully Permissive (≈ 161B Tokens)</b>					
1838	GlaiveAI Reasoning	Apache 2.0	Permissive	N/A	38.366	Fully synthetic
1839	Nemotron (Science & Math)	CC-BY-4.0	Permissive	Permissive (StackOverflow, WildChat)	22.310	Science & Math subset of Llama-Nemotron-Post-Training-Dataset
1840	Ling-Coder/SyntheticQA	Apache 2.0	Permissive	N/A	19.852	
1841	Open Thoughts	Apache 2.0	Permissive	Permissive (OpenMath-2-Math, CodeGolf, OpenCode, etc)	18.786	Excludes Organic Chemistry subset
1842	EuroPat	Public Domain	Permissive	Permissive	11.586	Synthetic image captions created from patents
1843	P3 (Permissive Subset)	Apache 2.0	N/A	Permissive (ARC, PIQA, BoolQ, etc)	10.130	
1844						
1845	Nemotron-CC	Common Crawl ToS	Permissive	Permissive (Common Crawl)	6.230	Using a Permissive-only subset
1846	YouTube	CC-BY-4.0	Permissive	Permissive (VALID, CommonCorpus)	7.386	Derived from CC-BY licensed YouTube content
1847	Prism-Math	CC-BY-4.0	Permissive	Permissive (NuminaMath-1.5)	5.682	
1848	DeepMind Math	Apache 2.0	N/A	Permissive (Procedurally Generated)	4.232	
1849	Misc. Instruct. / NVidia	NVIDIA license	Permissive	Permissive (GSM8K, MATH)	2.440	
1850	OpenMathInstruct-1	CC-BY-4.0	Permissive	N/A	1.018	Fully synthetic
1851	Websights	CC-BY-4.0	Permissive	N/A	1.018	Fully synthetic
1852	Misc Instruct. / MetaMathQA-R1 (responses)	MIT	Permissive	Permissive (GSM8K, MATH)	0.672	
1853	Math Word Problems	Apache 2.0	N/A	Permissive	0.456	Procedurally Generated
1854	Ling-Coder/DPO	Apache 2.0	Permissive	Unknown (Common-Crawl)	0.398	
1855	Misc. Instruct. / OpenR1-Math-220k	Apache 2.0	Permissive	Permissive (NuminaMath-1.5)	0.320	
1856	Misc. Instruct. / NVIDIA SFT Datablend	CC-BY-4.0	Permissive	Permissive (MNLI, COPA, PIQA, etc)	0.286	
1857	Misc. Instruct. / OpenThoughts-114k-Code (decontaminated)	Apache 2.0	Permissive	Permissive (TACO, Apps, CodeContests, etc)	0.150	
1858						
1859	Misc. Instruct. / Synthetic Code Generations	Apache 2.0	Permissive	N/A	0.104	Fully synthetic
1860						
1861	Misc. Instruct. / PrimeIntellect	Apache 2.0	N/A	N/A	0.076	
1862	StackExchange QnA	Apache 2.0	Permissive	Permissive (CommitPack)	0.006	
1863	Misc. Instruct. / PrimeIntellect Real World SWE Problems	Apache 2.0	Permissive	N/A	0.004	Fully synthetic
1864	Misc. Instruct. / PrimeIntellect Synthetic Code Understanding	Apache 2.0	Permissive	N/A	0.004	Fully synthetic
1865	Misc. Instruct. / GSM8K (train)	MIT	Permissive	N/A	0.004	Fully synthetic
1866						
1867	<b>Tier 2(a): Permissive with Upstream Opacity (≈ 35B Tokens)</b>					
1868	OS-Q2 (OpenScience)	CC-BY-4.0	Permissive	N/A	17.366	Fully synthetic
1869	Ring-lite SFT Data	Apache 2.0	Permissive	Permissive (CodeContest, APPS, TACO, etc)	14.968	
1870	PyEdu Reasoning	Stack V1, ODC-BY	Permissive	Permissive (The Stack V1)	3.138	
1871	Misc. Instruct. / Magpie-Phi3-Pro-1M-v0.1	N/A	Permissive	N/A	0.386	Fully synthetic
1872	MegaWika	CC-By-SA/4.0	Permissive	Permissive (Wikipedia)	0.356	
1873	Misc. Instruct. / Magpie-Qwen2.5-Coder-Pro-300K-v0.1	N/A	Permissive	N/A	0.120	Fully synthetic
1874	Misc. Instruct. / NovaSky-AI Sky-T1	Apache 2.0	Permissive	Permissive (AIME, MATH, etc)	0.080	
1875	Misc. Instruct. / BigCode	Stack v1	Permissive	Permissive	0.018	
1876	Self-OSS-Instruct	MIT	Permissive	Permissive (UltraChat, TruthfulQA, etc)	0.014	
1877	Misc. Instruct. / UltraFeedback	N/A	Permissive	Permissive (CaseHOLD)	0.002	Case law is public domain
1878	Misc. Instruct. / CaseHOLD (Phi4 Reasoning Traces)					
1879	<b>Tier 2(b): Restricted, Mixed or Opaque Provenance (≈ 17B Tokens)</b>					
1880	Ling-Coder/SFT	Apache 2.0	Permissive	Partially Unknown (Github, CommonCrawl, The Stack, etc)	7.668	Unknown provenance of CommonCrawl subset
1881	P3 (Commercial Subset)	Apache 2.0	N/A	Ambiguous-license (Amazon, Rotten Tomatoes, IMDb, etc)	5.730	User-generated content on commercial platforms with no verifiable licenses
1882	Misc. Instruct. / OpenMathInstruct-2	CC-By-4.0	Restricted	Permissive (GSM8K and MATH)	3.202	Llama-generated responses
1883	Misc. Instruct. / OpenManus-RL	Apache 2.0	Mixed	Permissive (AgentTraj-L, Agent-FLAN, etc)	0.024	GPT-4 used for generating traces in the AgentTraj-L subset
1884	<b>Tier 3: Civic / Governmental Works (≈ 15B Tokens)</b>					
1885	Nemotron-CC (Gov. Portion)	Common Crawl ToS	Permissive	Permissive (Common Crawl)	8.531	Using a Permissive-only subset
1886	MGACorpus	ODC-By	Permissive	Permissive (FinewebEdu-dedup)	6.596	Using only permissive subsets

1890 of-concept for data efficiency and strong downstream performance within a strict permissive-first,  
1891 risk-mitigated licensing framework. We identify several concrete avenues to scale this approach to  
1892 the multi-trillion token regime required for larger foundation models:  
1893

1894 **Subset Upsampling.** Standard industry recipes for large-scale training often heavily upsample  
1895 high-quality data. For instance, Llama 3 (Meta, 2024) employs upsampling factors of 4–10× for its  
1896 highest-quality subsets to reach its training budget. In contrast, the current iteration of **MixtureVitae**  
1897 does not assign aggressive upsampling factors to individual shards. Applying standard upsampling  
1898 techniques to our highest-value subsets (such as curated reasoning) would immediately scale their  
1899 contribution to the total token count.

1900 **Multilingual Expansion.** The current release of **MixtureVitae** is primarily English-centric. Ex-  
1901 panding the sourcing strategy to include multilingual data represents an order-of-magnitude oppor-  
1902 tunity for scaling. This can be achieved through two primary methods: (1) identifying and allowing  
1903 international permissively licensed sources, and (2) using machine translation to expand the existing  
1904 data in **MixtureVitae**.  
1905

1906 **Synthetic Expansion.** Our Math and Reasoning synthetic subsets are generated procedurally or  
1907 via LLMs. This generation process is horizontally scalable. By increasing the compute budget for  
1908 generation, these high-density subsets can be expanded significantly without incurring the legal risks  
1909 associated with scraping organic web data.  
1910

1911 **Web Data Rephrasing.** Recent work has demonstrated the utility of rephrasing web data to im-  
1912 prove quality and standardize style (Maini et al., 2024). Applying a similar rephrasing pipeline on  
1913 top of the **MixtureVitae** web data processing pipeline can further expand the corpus volume while  
1914 maintaining the strict safety and licensing standards defined in our framework.  
1915

1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943