## Mapping 1,000+ Language Models via the Log-Likelihood Vector

## **Anonymous ACL submission**

## Abstract

To compare autoregressive language models at scale, we propose using log-likelihood vectors computed on a predefined text set as model features. This approach has a solid theoretical basis: when treated as model coordinates, their squared Euclidean distance approximates the Kullback–Leibler divergence of text-generation probabilities. Our method is highly scalable, with computational cost growing linearly in both the number of models and text samples, and is easy to implement as the required features are derived from cross-entropy loss. Applying this method to over 1,000 language models, we constructed a "model map," providing a new perspective on large-scale model analysis.

## 1 Introduction

004

011

014

017

019

024

Language models have been evolving rapidly, and their community has grown substantially. To understand the community's structure and future directions, it is essential to systematically analyze model similarity and positioning based on language modeling principles. On the Hugging Face Hub, models are categorized based on their names and attributes, while other studies measure model similarity by outputs (Yax et al., 2024) or activations (Zhou et al., 2025). Leaderboards (Beeching et al., 2023; Fourrier et al., 2024; Chiang et al., 2024) are commonly used to assess model standings.

Since language models are probability models, we propose representing each model with coordinates that reflect the structure of the corresponding probability distribution space. Concretely, we define a language model's coordinates as its loglikelihood vector across a large collection of texts. Figure 1 shows a model map obtained through dimensionality reduction applied to the coordinates of 1,018 language models. This visualization reveals that models of the same type tend to cluster together, while text categories with the highest rel-



Figure 1: Map of 1,018 language models. Their loglikelihood vectors are visualized using t-SNE. (Top) Colors indicate model types. (Bottom) Colors indicate the model's "primary text category," the text category where the model achieves the highest relative log-likelihood among 17 categories. See Section 4 for details.

ative log-likelihoods appear as a continuous distribution across the map.

We find that the distances in our defined coordinate system accurately capture relationships among language models. On this map, each point represents a single model, with those having similar

meta-llama/Meta-Llama-3-8B Kl		mistralai/Mistral-7B-v0.3 KL		google/codegemma-2b	KL	deepseek-ai/deepseek-llm-7b-base	KL
Undi95/Meta-Llama-3-8B-hf	0.003	MaziyarPanahi/Mistral-7B-v0.3	0.000	deepseek-ai/deepseek-coder-1.3b-instruct	1654.767	deepseek-ai/deepseek-moe-16b-base	130.715
dfurman/Llama-3-8B-Orpo-v0.1	6.994	mistral-community/Mistral-7B-v0.2	7.221	bigcode/starcoderbase-1b	1812.768	deepseek-ai/deepseek-llm-7b-chat	245.408
migtissera/Tess-2.0-Llama-3-8B	28.829	unsloth/mistral-7b-v0.2	7.222	deepseek-ai/deepseek-coder-1.3b-base	2115.462	deepseek-ai/deepseek-moe-16b-chat	247.889
freewheelin/free-llama3-dpo-v0.2	67.783	mistralai/Mistral-7B-v0.1	22.038	bigcode/gpt_bigcode-santacoder	2644.348	deepseek-ai/DeepSeek-V2-Lite	306.530
jondurbin/bagel-8b-v1.0	110.365	Cartinoe5930/Llama2_init_Mistral	29.819	deepseek-ai/deepseek-coder-6.7b-instruct	2675.004	deepseek-ai/ESFT-vanilla-lite	307.487
migtissera/Llama-3-8B-Synthia-v3.5	128.141	Locutusque/Hercules-3.1-Mistral-7B	32.099	Qwen/CodeQwen1.5-7B-Chat	2759.243	deepseek-ai/DeepSeek-V2-Lite-Chat	585.090
nvidia/Llama3-ChatQA-1.5-8B	128.776	migtissera/Synthia-7B-v3.0	33.418	NTQAI/Nxcode-CQ-7B-orpo	2767.578	mistralai/Mistral-7B-Instruct-v0.1	913.708
ruslanmv/Medical-Llama3-8B	138.566	uukuguy/speechless-zephyr-code-functionary-7b	35.723	Salesforce/codegen-6B-multi	2859.701	statking/zephyr-7b-sft-full-orpo	1089.127
FairMind/Llama-3-8B-4bit-UltraChat-Ita	199.748	uukuguy/zephyr-7b-alpha-dare-0.85	35.745	bigcode/starcoderbase-7b	2990.326	Severian/ANIMA-Phi-Neptune-Mistral-7B	1140.161
NousResearch/Hermes-2-Theta-Llama-3-8B	222.357	crumb/apricot-wildflower-20	40.907	deepseek-ai/deepseek-coder-6.7b-base	3282.352	sethuiyer/Medichat-Llama3-8B	1157.599

Table 1: Top 10 nearest neighbors among the 1,018 language models for each model listed in the first row. The KL divergence is computed using formula (3).

text-generation probability distributions appearing closer together and those with more distinct distributions positioned farther apart. In Section 2, we show that the squared Euclidean distance in this coordinate system approximates the Kullback-Leibler (KL) divergence among models. Table 1 lists the nearest neighbors for each language model. For example, many of the closest neighbors of meta-llama/Meta-Llama-3-8B (AI@Meta, 2024) also contain Llama-3 in their names.

047

048

051

061

062

067

071

087

Several studies have explored methods for comparing language models (see Appendix A). In particular, prior work on comparing generated text includes approaches that construct phylogenetic trees based on model-generated text (Yax et al., 2025) and approaches that measure differences in text-generation probabilities conditioned on given prompts using KL divergence (Melamed et al., 2024). However, these methods require generating text with each model, thus incurring the cost of pairwise distance computations, which becomes prohibitively expensive at large scale. By contrast, our method does not involve actual text generation; instead, we compute generation probabilities on a predefined text corpus. This enables us to derive model coordinates without pairwise comparisons, allowing efficient large-scale comparison of many models.

A language model's log-likelihood vector can be treated as its feature vector. In Section 4, we discuss the insights gleaned from analyzing these features. Then, in Section 5, we show how this feature vector can be leveraged to predict benchmark performance.

## 2 Mapping Language Models into the Space of Text Probability Distributions

In this section, we present our proposed method. Sections 2.2 and 2.3 introduce model feature vectors derived from text-generation probabilities. Section 2.4 demonstrates that the squared Euclidean distance in the coordinate system built using these features approximates the KL divergence between models. Section 2.5 offers an interpretation of the resulting model coordinates. An extension of this method, which defines model coordinates using the sequence of conditional probabilities for generating a given token sequence, is presented in Appendix E.

## 2.1 Autoregressive language models

Let  $\mathcal{X}$  be the set of all possible texts, and let  $\mathcal{V}$  be the token vocabulary. A text  $x \in \mathcal{X}$  is represented as a sequence of tokens:

$$x = (y_1, \dots, y_n), \quad y_t \in \mathcal{V}.$$
 09

089

091

092

093

095

096

097

100

101

102

104

105

106

109

111

112

113

114

116

117

118

Denoting the maximum text length by  $n_{\text{max}}$ , we have  $\mathcal{X} = \bigcup_{n=0}^{n_{\text{max}}} \mathcal{V}^n$ . We consider a set of K language models  $\{p_i\}_{i=1}^K$ . With  $y_0$  denoting the beginning-of-sequence (BOS) token, each language model  $p_i$  predicts the next token  $y_t$  given the preceding token sequence  $y^{t-1} = (y_0, \dots, y_{t-1})$ . Thus, the conditional probability defined by  $p_i$  is given by

$$y_t \sim p_i(y_t \mid y^{t-1}), \quad t = 1, \dots, n.$$
 107

Accordingly, the probability of a text x under model  $p_i$ , denoted  $x \sim p_i$ , is

$$p_i(x) = \prod_{t=1}^n p_i(y_t \mid y^{t-1}).$$
 110

In addition to the K language models  $p_1, \ldots, p_K$ , we introduce a language model  $p_0$  that represents an underlying distribution for theoretical purposes. We assume we have a dataset (corpus)

$$D = (x_1, x_2, \dots, x_N) \in \mathcal{X}^N,$$
115

consisting of N texts, where each text is independently drawn from  $p_0$ .

## 2.2 Log-likelihood vector

For a model  $p_i$ , the probability of generating a text 119 x is denoted  $p_i(x)$ . Following the convention in 120

163

164

165

166

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

190

191

192

193

194

195

196

197

198

199

likelihood of model  $p_i$  given the text x. The loglikelihood is then defined as  $\operatorname{KL}(p_i, p_i) = \sum p_i($ 

$$\ell_i(x) = \sum_{t=1}^n \log p_i(y_t \mid y^{t-1})$$

121

122

123

124

125

126

127

128

129

130

132

133

134

135

137

138

139

140

141

142

143

144

145

146

147

In language model implementations,  $-\ell_i(x)$  corresponds to the cross-entropy loss for the text x, and  $\exp(-\ell_i(x)/n)$  is known as the perplexity.

statistical model selection, we refer to  $p_i(x)$  as the

Our approach is straightforward. Given that the dataset D consists of N texts, we use the loglikelihood vector

$$\boldsymbol{\ell}_i = (\ell_i(x_1), \dots, \ell_i(x_N))^\top \in \mathbb{R}^N$$

as the feature vector for model  $p_i$ . The first step in our model analysis is to construct the log-likelihood matrix

$$\boldsymbol{L} = (\boldsymbol{\ell}_1, \dots, \boldsymbol{\ell}_K)^\top \in \mathbb{R}^{K \times N}$$

by stacking the vectors  $\ell_i$  for the K models.

## 2.3 Double centering

As a preprocessing step for model analysis, we apply a technique called double centering (Borg and Groenen, 2005) to *L*. First, we perform rowwise centering. The mean of each row, referred to as the mean log-likelihood, is given by

$$\bar{\ell}_i = \sum_{s=1}^N \ell_i(x_s) / N$$

Subtracting this value from each component of  $\ell_i$ , we define the centered log-likelihood vector  $\boldsymbol{\xi}_i = (\xi_{i1}, \dots, \xi_{iN})^\top \in \mathbb{R}^N$ , where

$$\xi_{is} := \ell_i(x_s) - \bar{\ell}_i, \quad s = 1, \dots, N.$$

148 Next, we apply column-wise centering to the matrix 149 of centered feature vectors  $(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_K)^{\top}$ . The 150 mean vector is  $\bar{\boldsymbol{\xi}} = \frac{1}{K} \sum_{i=1}^{K} \boldsymbol{\xi}_i$ , and by subtracting 151 this vector from each  $\boldsymbol{\xi}_i$ , we define the double-152 centered log-likelihood vector

$$oldsymbol{q}_i = oldsymbol{\xi}_i - oldsymbol{\xi}$$

154 For further details, see Appendices B and C.

## 155 2.4 Kullback–Leibler divergence

The Kullback–Leibler (KL) divergence is often used to measure how far apart two models  $p_i$  and  $p_j$  are in the space of probability distributions.<sup>1</sup> It is defined as

$$\operatorname{KL}(p_i, p_j) = \sum_{x \in \mathcal{X}} p_i(x) \log \frac{p_i(x)}{p_j(x)}$$
160

$$= \mathbb{E}_{x \sim p_i} \left( \ell_i(x) - \ell_j(x) \right).$$
 (1) 161

We assume the dataset D is generated from an unknown underlying model  $p_0$  and that the models  $p_i$ and  $p_j$  provide good approximations of  $p_0$ . Under this assumption, the KL divergence can be approximated as follows:

$$2\operatorname{KL}(p_i, p_j) \approx \operatorname{Var}_{x \sim p_0} \left( \ell_i(x) - \ell_j(x) \right). \quad (2)$$

While the definition of KL divergence in (1) involves the expectation of  $\ell_i(x) - \ell_j(x)$ , the approximation in (2) takes the form of a variance. This result is somewhat surprising yet quite insightful. Notably, although KL divergence is not symmetric in the two models, the approximation in (2) is symmetric. We estimate (2) from the dataset D as

$$2 \operatorname{KL}(p_i, p_j) \approx \|\boldsymbol{q}_i - \boldsymbol{q}_j\|^2 / N.$$
 (3)

Thus, if we regard the model coordinates of  $p_i$  as  $q_i/\sqrt{N}$  by scaling with N, then the squared Euclidean distance between two points approximates  $2 \text{ KL}(p_i, p_j)$ .

The main results, namely (2) and (3), are proved in Appendix D using the theory of exponential family of distributions (Barndorff-Nielsen, 2014; Efron, 1978, 2022; Amari, 1982), similar to the discussion on the relationship between the norm of embeddings and KL divergence (Oyama et al., 2023). Although the concept of model coordinates has been discussed in statistics (Shimodaira, 2001), and there have been a few applications of model maps (Shimodaira and Hasegawa, 2005; Shimodaira and Terada, 2019), they have seldom been utilized in practice.

## 2.5 Model coordinates

We primarily use  $q_i$  as the feature vector of model  $p_i$  and refer to it as the model coordinates. As shown in (3), the squared Euclidean distance in the q-coordinate system approximates the KL divergence between language models,<sup>2</sup> indicating that  $q_i$  represents the position of  $p_i$  in the space of probability distributions. Since  $\xi_i$  differs from  $q_i$  only

 $<sup>{}^{1}\</sup>mathbb{E}(\cdot)$  denotes expectation and  $\operatorname{Var}(\cdot)$  denotes variance.

<sup>&</sup>lt;sup>2</sup>That is,  $\|\boldsymbol{q}_i - \boldsymbol{q}_j\|^2 \approx 2N \operatorname{KL}(p_i, p_j)$ . For simplicity, we omit the constant scaling factor in expressions of this type.



Figure 2: Text embeddings for 10,000 texts in dataset D, computed via simcse-roberta-large (Gao et al., 2021b) and visualized with t-SNE. Colors indicate 17 text categories.

by an offset from the origin,  $\xi_i$  also serves as a model coordinate, and  $\|q_i - q_j\|^2 = \|\xi_i - \xi_j\|^2$ . However, we prefer  $q_i$  for its more interpretable components and thus adopt it throughout this paper.

For visualization purposes, we mainly use  $\ell_i$  as the coordinates of the model map, as  $\ell_i$  can be intuitively interpreted as encoding  $\sqrt{N} \bar{\ell}_i$  in the "height" dimension and  $q_i$  in the "horizontal" dimensions. As shown in Appendix D.6,

$$\|\boldsymbol{\ell}_i - \boldsymbol{\ell}_j\|^2 = \|\boldsymbol{q}_i - \boldsymbol{q}_j\|^2 + N(\bar{\ell}_i - \bar{\ell}_j)^2,$$
 (4)

which means the squared Euclidean distance in the  $\ell$ -coordinate system can be decomposed into the sum of  $2N \operatorname{KL}(p_i, p_j)$  and  $N (\overline{\ell}_i - \overline{\ell}_j)^2$ .

## **3** Experimental Setup

We describe the key components of our experiment. In particular, Section 3.1 explains the procedure for selecting the 10,000 texts used to compute the model coordinates, and Section 3.2 discusses how we selected the 1,018 language models. Further details are given in Appendix G.

## 3.1 Selection of text data

The texts used for computing the language models' coordinates were extracted from the Pile (Gao et al., 2020), with five categories of copyrighted material removed.<sup>3</sup> This yielded a dataset D consisting of 10,000 texts, each tagged with a category label from the Pile. Figure 2 visualizes these texts.

To build the dataset, we began by dividing the first 1M texts from the Pile Uncopyrighted corpus into 1,024-byte chunks (UTF-8 encoded). In cases where decoding errors occurred, we truncated by one byte at a time. Chunks smaller than 256 bytes were discarded, resulting in about 5.7M valid chunks. From these, we randomly sampled 10,000 texts to create the final dataset used for computing model coordinates.

232

233

234

235

236

237

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

256

257

258

259

260

261

263

264

265

266

267

268

269

270

272

273

274

275

276

277

## 3.2 Selection of language models

We used K = 1,018 language models in total. Of these, 1,000 were selected from models listed on Open LLM Leaderboard v1. Specifically, we considered CausalLM models ranging from 1B to 13B parameters and ranked them by their number of downloads over the 30 days preceding February 1, 2025.<sup>4</sup> We initially selected the top 1,100 models by download count and attempted log-likelihood calculations. Among these, 1,011 successfully produced valid log-likelihood values, and we chose the 1,000 most frequently downloaded from that set. In addition, we included 18 models from the DeepSeek language model series. We obtained information on model parameter sizes and architectures from the Leaderboard. Appendix G provides basic information on the selected models and details on how model types were defined. A complete list of models used in this study is given in Appendix K.

## 3.3 Computation of the Log-Likelihood

The log-likelihood matrix L was computed in float16 precision, with the bottom 2% of values clipped. This clipping mitigates the large impact of extremely low likelihoods on (3). After computing L, we applied row-wise and column-wise centering to obtain the double-centered log-likelihood matrix Q. Figure 3 visualizes Q. Each value in the matrix can be interpreted in two ways: as the relative probability of a text for each model or as the relative likelihood of a model for each text. Both models and texts exhibit clustering patterns. Figure 4 shows a dendrogram of the top 100 models. We examined the effective dimension from the perspective of feature vector dimensionality reduction and found that the cumulative contribution ratio, based on the sum of squared singular values of Q, reached 90% at 42 dimensions and 95% at 82 dimensions.

## 3.4 Obtaining the leaderboard scores

We obtained benchmark scores for the language models used in our experiments from Open LLM

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/datasets/monology/ pile-uncopyrighted

<sup>&</sup>lt;sup>4</sup>Hugging Face's API provides the total number of downloads in the past 30 days.



Figure 3: The double-centered log-likelihood matrix Q, with rows and columns reordered by hierarchical clustering. Each row corresponds to one of the 1,018 models, color-coded by model type. Each column represents one of the 10,000 texts, color-coded by text category.



Figure 4: Hierarchical clustering of the top 100 mostdownloaded models, based on their feature vectors  $q_i$ . Model names are color-coded by model type.

Leaderboard v1.<sup>5</sup> Between April 2023 and June 2024, this leaderboard evaluated language models on six tasks: AI2 Reasoning Challenge (ARC) (Clark et al., 2018), HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021a), TruthfulQA (Lin et al., 2022a), Winogrande (Sakaguchi et al., 2019), and GSM8K (Cobbe et al., 2021). Along with individual task scores, we also use the average score across all six tasks, referred to as 6-TaskMean.

## 4 Map of Language Models

279

281

284

285

291

We applied t-SNE (van der Maaten and Hinton, 2008) to the log-likelihood matrix L for dimensionality reduction<sup>6</sup>. Using this visualization, we analyze the insights gained from the model map in this section. While this paper presents model maps us-

open-llm-leaderboard-old/open\_llm\_leaderboard

ing L, alternative maps using the double-centered log-likelihood matrix Q, as well as a model map with labels for all language models, are available in Appendix H.

294

297

299

300

301

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

321

322

323

## 4.1 Visualizing attributes on the model map

**Model type.** The top panel of Fig. 1 visualizes the distribution of model types, with each model color-coded according to its type. We observe that models belonging to the same type tend to cluster together, forming distinct regions on the map (e.g., 11ama-2, mistral, and gemma). In particular, models optimized for coding tasks<sup>7</sup> appear in a relatively compact region, suggesting that these models share notable similarities in their probability distributions.

**Text category.** The bottom panel of Fig. 1 shows the model map with each model color-coded according to the text category in which it achieves the highest relative likelihood. From this figure, we see that models exhibiting high likelihoods for the same text category are grouped together. Notably, the cluster containing codingspecialized models in the top panel aligns with the GitHub/StackExchange region in the bottom panel, suggesting that these models have relatively high likelihoods for text originating from GitHub and StackExchange.

**Model performance.** Figure 5 visualizes two evaluation metrics: mean log-likelihood and benchmark task performance. From the left and central

<sup>&</sup>lt;sup>5</sup>https://huggingface.co/spaces/

<sup>&</sup>lt;sup>6</sup>The perplexity value was set to 30, and we used the scikitlearn implementation (Varoquaux et al., 2015).

<sup>&</sup>lt;sup>7</sup>For example, starcoder2-7b (Lozhkov et al., 2024), deepseek-coder-1.3b-base (Guo et al., 2024), codegemma-2b (CodeGemma Team et al., 2024) and CodeLlama-7b-Instruct-hf (Rozière et al., 2024).



Figure 5: Model maps illustrating model performance. From left to right, the panels show each model's mean loglikelihood, 6-TaskMean score, and the "primary task," meaning the task for which each model achieves its highest standardized score among the six tasks. The color bar is clipped at the 10th percentile for mean log-likelihood and 6-TaskMean, with darker colors indicating better performance. In the primary task panel, models with standardized scores below zero on all six tasks are labeled "All Under 0."



Figure 6: A model map color-coded by (Left) model size and (Right) model creation date.

panels, we see that both metrics exhibit similar trends on the map, where models that lie close together tend to show similar metric values. Additionally, in the right panel, the GSM8K/MMLU region corresponds to the ArXiv/PubMed Central region in the bottom panel of Fig. 1, suggesting that models with high likelihoods on academic and scientific texts also tend to perform well on mathematical reasoning and academic knowledge-intensive tasks.

Model size and creation date. Figure 6 shows the distribution of models by size and creation date.
Compared to Fig. 5, newer models generally perform better, but model size does not always correlate with performance, as some smaller models perform comparably to larger ones.

## 4.2 Detection of data leakage

325

326

327

337

338

Since the text data we used was extracted from
the Pile corpus, models that were pre-trained on
the Pile are likely to exhibit higher log-likelihood
values than their actual capabilities measured by
benchmarks. We analyze this effect using the
model map in Fig. 7. The left panel highlights
models that used the Pile for pre-training. The right



Figure 7: (Left) Models tagged with "the Pile." (Right) Difference between the standardized mean log-likelihood and the standardized 6-TaskMean score.

panel shows models with high mean log-likelihood relative to their 6-TaskMean score. The alignment between these two distributions suggests that models pre-trained on the Pile tend to achieve higher likelihoods on our text data, while their benchmark performance remains comparatively lower.

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

## 5 Predicting Model Performance from Model Coordinates

As shown in the central panel of Fig. 5, the positioning of models on the map suggests that a model's benchmark performance may be inferred from its coordinates. In this section, we conduct a regression analysis using the q-coordinates to predict benchmark scores and evaluate predictive performance.

## 5.1 Benchmark scores and models

We use six benchmark scores from Open LLM Leaderboard v1, as described in Section 3.4. Our experiments are conducted on 996 models for which these benchmark scores are available<sup>8</sup>.

<sup>&</sup>lt;sup>8</sup>From the 1,018 models, we excluded four that lacked TruthfulQA scores and 18 from deepseek-ai, as their bench-

	ARC	HellaSwag	MMLU	TruthfulQA	Winogrande	GSM8K	6-TaskMean	mean log-likelihood
Pearson's $r$	0.972	0.948	0.963	0.958	0.965	0.935	0.976	0.994
Spearman's $\rho$	0.976	0.974	0.969	0.937	0.973	0.899	0.978	0.990

Table 2: Results of ridge regression for predicting benchmark scores from model coordinates. Predictions for 6-TaskMean and mean log-likelihood are also included. High correlation coefficients are observed across all settings.

## 5.2 Setting for regression analysis

370

371

374

375

376

381

394

396

400

401

402

403

404

405

For each benchmark task, the dataset is given as  $\{(q_1, f_1), \ldots, (q_K, f_K)\}$ , where  $q_i \in \mathbb{R}^N$  is the double-centered log-likelihood vector of the language model  $p_i$ , and  $f_i \in [0, 100]$  is its corresponding benchmark score. We use ridge regression to predict each benchmark score. Let  $Q \in \mathbb{R}^{K \times N}$  be the matrix of explanatory variables, and let  $f = (f_1, \ldots, f_K) \in \mathbb{R}^K$  be the vector for the target variable. The objective function with parameter  $w \in \mathbb{R}^N$  is given by:

$$L(w) = \|f - Qw\|^2 + \alpha \|w\|^2$$

where  $\alpha \in \mathbb{R}_{>0}$  is a hyperparameter that controls the strength of regularization. Since the number of variables N is much larger than the sample size  $K (N \gg K)$ , making this a high-dimensional regression setting, we carefully set  $\alpha$  using crossvalidation to avoid overfitting.

We split the set of models into five folds and perform parameter training<sup>9</sup> and benchmark score prediction. To mitigate the effect of randomness, we repeat the data splitting with five different seeds and take the average of the predictions as the final predicted score. As evaluation metrics, we compute Pearson's r and Spearman's  $\rho$  to measure the correlation between the predicted and benchmark scores. Additionally, we conduct experiments by replacing the target variable with 6-TaskMean and mean log-likelihood, leading to a total of eight experimental settings. See Appendix J.1 for details.

## 5.3 Results

Table 2 presents the results of the regression analysis. For all benchmark tasks, the correlation coefficients are high, demonstrating that ridge regression achieves strong predictive performance. Similarly, for 6-TaskMean, the correlation remains high, as also illustrated by the scatter plot in Fig. 8. Even for mean log-likelihood, the regression model attains a



Figure 8: Scatter plot of predicted scores versus 6-TaskMean scores for test sets. Pearson's r is 0.976, and Spearman's  $\rho$  is 0.978. Each point is color-coded by the mean log-likelihood, with higher values generally corresponding to higher 6-TaskMean scores. However, some models with extremely high mean log-likelihood due to data leakage (Section 4.2) appear in the lower score region, deviating from the general trend. Despite this, the regression-based predictions remain highly accurate, as reflected in the strong correlation coefficients. The color bar is clipped at the 10th percentile. Scatter plots for individual benchmark tasks are provided in Fig. 14 in Appendix J.2.



Figure 9: Relationship between the squared Euclidean distance of model coordinates and KL divergence. (Left) In the token-level experiment (Section 6.1), each point represents a text. (Right) In the text-level experiment (Section 6.2), each point represents a pair of models.

high correlation, despite the explanatory variables being double-centered.

## 6 Empirical Validation of Theory

In Section 2, we discussed model coordinates for text probability models. Appendix E extends this framework to token-sequence probability models. These theoretical results state that the squared Euclidean distance in the model-coordinate space approximates the KL divergence between models. To validate this, we first evaluate token-level condi-

- 408
- 410 411 412 413 414

mark scores were incomplete due to the models being relatively new, leaving 996 for analysis. For simplicity, we denote the number of models as K.

<sup>&</sup>lt;sup>9</sup>We used scikit-learn (Varoquaux et al., 2015) RidgeCV.



Figure 10: Visualization of 36 language models obtained by linearly interpolating pretrained model weights. Each point is color-coded according to its mean log-likelihood. (Left) Models in the weight parameter space. (Right) Models in the probability distribution space, represented by the q-coordinate system.

tional probability models (Section 6.1), since tokenlevel experiments are generally easier to conduct than text-level experiments. We then extend our analysis to text probability models (Section 6.2).Additionally, in Section 6.3, we explore the relationship between model weight parameters and model coordinates.

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

## 6.1 Validation of (26) for token-level models

**Settings.** Two models with a shared tokenizer Llama-2-7b-hf and Llama-2-7b-chat-hf (Touvron et al., 2023b), denoted as  $p_1$  and  $p_2$ , were used. Following the method described in Appendix E.1, we computed the model coordinates  $\zeta_1(x)$  and  $\zeta_2(x)$  for each text  $x = (y_1, \dots, y_n) \in D$ , where these coordinates are centered vectors with elements  $\log p_i(y_t|y^{t-1})$  as defined in (22). We then calculated the squared Euclidean distance between these coordinates,  $\|\zeta_1(x) - \zeta_2(x)\|^2$ . To obtain the exact KL divergence between models, we used the outputs of the softmax function in the language models. We computed the sum of per-token KL divergences:  $\sum_{t=1}^{n} \text{KL}(p_1(y_t|y^{t-1}), p_2(y_t|y^{t-1}))$ , which is also used in Ly et al. (2023).

**Results and discussion.** The left panel of Fig. 9 shows a scatter plot of the squared Euclidean distance and KL divergence for  $x \in D$ , with a correlation coefficient of 0.893. This result indicates that (26) provides a good approximation in actual language models.

## 6.2 Validation of (3) for text-level models

**Settings.** Among the 292 language models sharing the tokenizer with Llama-2-7b-hf (Touvron et al., 2023b), we excluded the five models with the largest values of  $\sum_{x \in D} \|\zeta_i(x)\|^2$ , leaving 287

models for our experiment. Using the approach described in Section 2, we computed the model coordinates for each model and calculated the squared Euclidean distance  $||\mathbf{q}_i - \mathbf{q}_j||^2$  between every pair of models. Because it is extremely difficult to directly compute  $\mathrm{KL}(p_i, p_j)$ , we instead used  $\frac{1}{N} \sum_{x \in D} ||\boldsymbol{\zeta}_i(x) - \boldsymbol{\zeta}_j(x)||^2$  as a proxy. For a theoretical justification that this quantity approximates  $\mathrm{KL}(p_i, p_j)$ , see (27) in Appendix E.

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

**Results and discussion.** As shown in the right panel of Fig. 9, the scatter plot of squared Euclidean distance versus KL divergence exhibits a correlation coefficient of 0.904. This finding confirms that the relationship in (3) holds approximately in practical language models.

# 6.3 Relationship between model weights and model coordinates

For language models with the same architecture, comparison can also be conducted via weight parameters. We investigated how the structure of the model-coordinate space aligns with the structure of the weight-parameter space. We generated 36 new language models by linearly interpolating the weights of Llama-2-7b-hf, Llama-2-7b-chat-hf (Touvron et al., 2023b), and vicuna-7b-v1.5 (Zheng et al., 2023). For these 36 models, we computed text-generation loglikelihoods and visualized the resulting model coordinates in two dimensions using Principal Component Analysis (PCA). Figure 10 shows these 36 models in both weight space and model coordinate space. We observe that the two-dimensional grid structure in the weight space is mapped continuously into the model coordinate space. This finding supports the validity of (6) and (14) in Appendix D, especially for models that are close to one another. Further details on this experiment can be found in Appendix I.

## 7 Conclusion

We propose a method to compare autoregressive language models using log-likelihoods from a predefined text set. By treating these as model coordinates, we show that the squared Euclidean distance approximates KL divergence, enabling efficient large-scale comparisons. Experiments with over 1,000 models confirm its effectiveness in structuring model relationships and predicting benchmark task performance, while also validating the theoretical foundations.

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

547

548

549

550

## Limitations

499

501

502

505

507

508

509

510

511

512

513

514

515

516

517

518

519

521

523

525

527

528

529

531

532

533

534

537

538

539

540

541

542

543

545

546

- Changing the text data will alter the analysis results of the model map. This is both a limitation and an advantage of the proposed method, because it allows us to choose text data according to the analysis objective. For example, if we want to investigate code-focused language models in more detail, we can increase the proportion of code data from GitHub, thereby increasing the resolution of code-focused models on the model map.
- The proof (Section 2 and Appendix D) that the squared Euclidean distance in the model coordinate system approximates the KL divergence assumes that the language model's text-generation probabilities closely match the distribution of the text data. When this assumption does not hold, the approximation accuracy decreases. However, even under such circumstances, the model coordinates should still function sufficiently as model features.
  - If the text data used for the model coordinates is contained in a language model's pre-training corpus, that model's mean log-likelihood may be overestimated. This data leakage is generally non-negligible, as shown in Fig. 7 in Section 4.2, which illustrates the effect of using the Pile corpus. However, comparing it against benchmark scores makes it possible to detect such data leakage, and one can remove models that are affected.
    - Beyond the data leakage mentioned above, other systematic errors introduced into the model coordinates can also affect the model map. As noted in Appendix C, *ℓ<sub>i</sub>* and *ℓ<sub>i</sub>* are susceptible to systematic biases. However, thanks to double centering, *q<sub>i</sub>* is less influenced by bias terms.
  - Although the calculation of model coordinates is linear time O(KN) in the number of models K and the number of texts N, it still requires a non-negligible amount of computation. In our experiments, it took about 10 minutes on a single GPU (RTX 6000 Ada) to compute coordinates ( $N = 10^4$ , float16) for a single 7B model.
- Computing the model map visualization from the model coordinates is generally not linear

in K. For example, t-SNE requires a distance matrix, incurring  $O(K^2N)$  computational cost. However, in modern computing environments, as long as K is not extremely large (e.g., in the millions), the cost of visualization is negligible compared to the cost of calculating the model coordinates.

- A sufficiently large number of text samples, N, is desirable for the model coordinates. We used  $N = 10^4$ . Since the error in the KL divergence estimate due to randomness decreases proportionally to  $N^{-1/2}$ , N must be increased according to the desired resolution of the model map.
- When using model coordinates as feature vectors,  $N = 10^4$  can be unwieldy. According to the experiment in Section 3.3, applying PCA to reduce the dimensionality of *q*-coordinates to 82 dimensions still retains 95% of the information. However, the predictive performance of such dimension-reduced features has not yet been tested.
- Since the language models used in our experiments were obtained from the Open LLM Leaderboard v1 (which ran from April 2023 to June 2024), our discussion of models released after June 2024 is limited.
- The results of the task-performance prediction in Section 5 should be interpreted conservatively. While we use a proper cross-validation setup to prevent data leakage (splitting the training set, validation set, and test set), if very similar models appear in both training and test sets (for example, models that share a base model and differ only slightly by finetuning), it might become easier to predict task performance. Also, we have not evaluated models that are not listed on the leaderboard.
- The theoretical validation experiment in Section 6 is limited. Currently, it is difficult to directly compute exact KL divergence values among text-generation probability models, so conducting more precise validation experiments remains a future challenge.
- In the method of computing model coordinates from token-sequence conditional probabilities (Appendix E and Appendix F), the proof that the squared Euclidean distance in

the model coordinate system approximates the KL divergence requires additional assumptions (Appendix F.3). In practice, Assumption 2 does not hold, and due to the variations in (33),  $\|\zeta_i - \zeta_j\|^2$  in (26) tends to overestimate the KL divergence. Nonetheless, even in such situations, token-level model coordinates will still likely function sufficiently as features.

## References

610

611

612

613

614

617

618

619

622

624

628

636

637

638

641

645

- 01. AI, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Guoyin Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Wen Xie, and 13 others. 2025. Yi: Open foundation models by 01.ai. *Preprint*, arXiv:2403.04652.
- 42dot Inc. 2023. 42dot llm: A series of large language model by 42dot.
- Rishiraj Acharya. 2023. Catppt.
  - AI@Meta. 2024. Llama 3 model card.
    - Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. 2023. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. *Preprint*, arXiv:2305.13245.
  - AI@Waktaverse. 2024. Waktaverse llama 3 model card.
    - Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.
      - Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. *Preprint*, arXiv:2402.17733.
    - Shun-Ichi Amari. 1982. Differential Geometry of Curved Exponential Families-Curvatures and Information Loss. *The Annals of Statistics*, 10(2):357 – 385.
    - Shun-Ichi Amari. 1998. Natural gradient works efficiently in learning. *Neural computation*, 10:251–276.
    - Yuvanesh Anand, Zach Nussbaum, Brandon Duderstadt, Benjamin Schmidt, and Andriy Mulyar. 2023.Gpt4all: Training an assistant-style chatbot with large scale data distillation from gpt-3.5-turbo.

Alex Andonian, Quentin Anthony, Stella Biderman, Sid Black, Preetham Gali, Leo Gao, Eric Hallahan, Josh Levy-Kramer, Connor Leahy, Lucas Nestler, Kip Parker, Michael Pieler, Shivanshu Purohit, Tri Songz, Wang Phil, and Samuel Weinbach. 2021. GPT-NeoX: Large Scale Autoregressive Language Modeling in PyTorch.

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *Preprint*, arXiv:2310.11511.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. Program synthesis with large language models. *Preprint*, arXiv:2108.07732.
- Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q. Jiang, Jia Deng, Stella Biderman, and Sean Welleck. 2024. Llemma: An open language model for mathematics. *Preprint*, arXiv:2310.10631.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. Layer normalization. *Preprint*, arXiv:1607.06450.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, and 29 others. 2023. Qwen technical report. *Preprint*, arXiv:2309.16609.
- Abhinand Balachandran. 2023. Tamil-llama: A new tamil language model based on llama 2. *Preprint*, arXiv:2311.05845.
- Ole Barndorff-Nielsen. 2014. Information and exponential families: in statistical theory. John Wiley & Sons.
- Pierpaolo Basile, Elio Musacchio, Marco Polignano, Lucia Siciliani, Giuseppe Fiameni, and Giovanni Semeraro. 2023. Llamantino: Llama 2 models for effective text generation in italian language. *Preprint*, arXiv:2312.09993.
- Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry Tworek, and Mark Chen. 2022. Efficient training of language models to fill in the middle. *Preprint*, arXiv:2207.14255.
- Edward Beeching, Clémentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. Open llm leaderboard.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *Preprint*, arXiv:2004.05150.

808

754

- Leonard Bereska and Stratis Gavves. 2024. Mechanistic interpretability for AI safety - a review. *Transactions on Machine Learning Research*. Survey Certification, Expert Certification.
- Stella Biderman, Kieran Bicheno, and Leo Gao. 2022. Datasheet for the pile. *Preprint*, arXiv:2201.07311.

705

711

712

714

715

716

717

718

719

720

721

723

724

725

726

727

728

733

734

735

736

737

740

741

742

744

745

746

747

748

750

751

752

753

- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. Pythia: A suite for analyzing large language models across training and scaling. *Preprint*, arXiv:2304.01373.
  - Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2019. Piqa: Reasoning about physical commonsense in natural language. *Preprint*, arXiv:1911.11641.
  - Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. 2022. Gpt-neox-20b: An opensource autoregressive language model. *Preprint*, arXiv:2204.06745.
- Sid Black, Gao Leo, Phil Wang, Connor Leahy, and Stella Biderman. 2021. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow. If you use this software, please cite it using these metadata.
- Ingwer Borg and Patrick J. F. Groenen. 2005. *Modern Multidimensional Scaling: Theory and Applications*, 2 edition. Springer Series in Statistics. Springer, New York, NY.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. *Preprint*, arXiv:2005.14165.
- CeADAR. 2023. Financeconnect-13b (revision 5f7841d).
- Sahil Chaudhary. 2023. Code alpaca: An instructionfollowing llama model for code generation.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. 2024a. Alpagasus: Training a better alpaca with fewer data. *Preprint*, arXiv:2307.08701.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg

Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. Evaluating large language models trained on code. *Preprint*, arXiv:2107.03374.

- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Yao Liu, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. 2023a. Symbolic discovery of optimization algorithms. *Preprint*, arXiv:2302.06675.
- Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2024b. Longlora: Efficient fine-tuning of long-context large language models. *Preprint*, arXiv:2309.12307.
- Yukang Chen, Shaozuo Yu, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2023b. Long alpaca: Long-context instructionfollowing models.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%\* chatgpt quality.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. Chatbot arena: An open platform for evaluating llms by human preference. *Preprint*, arXiv:2403.04132.

chiliu. 2023. Mamba-gpt-3b-v4.

- François Chollet. 2019. On the measure of intelligence. *Preprint*, arXiv:1911.01547.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *Preprint*, arXiv:1905.10044.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *Preprint*, arXiv:1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *Preprint*, arXiv:2110.14168.
- CodeGemma Team, Heri Zhao, Jeffrey Hui, Joshua Howland, Nam Nguyen, Siqi Zuo, Andrea Hu, Christopher A. Choquette-Choo, Jingyue Shen, Joe Kelley, Kshitij Bansal, Luke Vilnis, Mateo Wirth, Paul Michel, Peter Choy, Pratik Joshi, Ravin Kumar, Sarmad Hashmi, Shubham Agrawal, and 8 others.

- 810 811 812 814 815 816 817 818 822 823 824 827 829 830 831 836 837 848
- 850

- 858
- 861

2024. Codegemma: Open code models based on gemma. Preprint, arXiv:2406.11409.

- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world's first truly open instructiontuned llm.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. Ultrafeedback: Boosting language models with scaled ai feedback. Preprint, arXiv:2310.01377.
- Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. 2024. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. Preprint, arXiv:2401.06066.
- Alpin Dale, Wing Lian, Bleys Goodson, Guan Wang, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknium". 2023. Llongorca13b: Llama2-13b model instruct-tuned for long context on filtered openorcav1 gpt-4 dataset.
- Tri Dao. 2023. FlashAttention-2: Faster attention with better parallelism and work partitioning.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. Preprint, arXiv:2205.14135.
- Soham De, Samuel L. Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Albert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, Guillaume Desjardins, Arnaud Doucet, David Budden, Yee Whye Teh, Razvan Pascanu, Nando De Freitas, and Caglar Gulcehre. 2024. Griffin: Mixing gated linear recurrences with local attention for efficient language models. Preprint, arXiv:2402.19427.
- DeciAI Research Team. 2023. Decilm-7b-instruct.
  - DeepSeek-AI, Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, Huazuo Gao, Kaige Gao, Wenjun Gao, Ruiqi Ge, Kang Guan, Daya Guo, Jianzhong Guo, Guangbo Hao, and 68 others. 2024a. Deepseek llm: Scaling open-source language models with longtermism. Preprint, arXiv:2401.02954.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. Preprint, arXiv:2501.12948.

DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, and 138 others. 2024b. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. Preprint, arXiv:2405.04434.

866

867

869

870

871

872

873

874

875

876

877

878

879

880

881

883

884

885

886

887

888

889

890

891

892

893

894

895

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 2022. 8-bit optimizers via block-wise quantization. *Preprint*, arXiv:2110.02861.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. Preprint, arXiv:2305.14314.
- Peter Devine. 2024. Tagengo: A multilingual chat dataset. Preprint, arXiv:2405.12612.
- Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. 2021. Bold: Dataset and metrics for measuring biases in open-ended language generation. Preprint, arXiv:2101.11718.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. Preprint, arXiv:2305.14233.
- Duy Quang Do, Hoang Le, and Duc Thang Nguyen. 2023. Torolama: The vietnamese instructionfollowing and chat model.
- Bradley Efron. 1978. The geometry of exponential families. The Annals of Statistics, 6:362-376.
- Bradley Efron. 2022. Exponential Families in Theory and Practice. Cambridge University Press.
- Kawin Ethayarajh, Winnie Xu, Dan Jurafsky, and Douwe Kiela. 2023. Human-centered loss functions (halos).
- Clémentine Fourrier, Nathan Habib, Alina Lozovskaya, Konrad Szafer, and Thomas Wolf. 2024. Open llm leaderboard v2. https://huggingface. co/spaces/open-llm-leaderboard/open\_llm\_ leaderboard.
- Victor Gallego. 2024. Configurable safety tuning of language models with synthetic preference data. Preprint, arXiv:2404.00495.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The pile: An 800gb dataset of diverse text for language modeling. Preprint, arXiv:2101.00027.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff,

- 919 920 921
- 922
- 92 92

- 928 929 930 931 932 933 934
- 935 936 937 938
- 939 940
- 942 943

941

- 944 945
- 946
- 947 948 949
- 9

952

960

961

- 962 963 964
- 965 966

968

969 970

- 971 972
- 972 973

Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2021a. A framework for few-shot language model evaluation.

- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021b. SimCSE: Simple contrastive learning of sentence embeddings. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. *Preprint*, arXiv:2009.11462.
  - Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, and 1331 others. 2024. Gemini: A family of highly capable multimodal models. *Preprint*, arXiv:2312.11805.
  - Xinyang Geng and Hao Liu. 2023. Openllama: An open reproduction of llama.
  - Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2024. Tora: A tool-integrated reasoning agent for mathematical problem solving. *Preprint*, arXiv:2309.17452.
  - Diego Granziol, Stefan Zohren, and Stephen Roberts. 2021. Learning rates as a function of batch size: A random matrix theory approach to neural network training. *Preprint*, arXiv:2006.09092.
  - Griffin Team, Soham De, Samuel L Smith, Anushan Fernando, Alex Botev, George-Christian Muraru, Ruba Haroun, and Leonard Berrada et al. 2024. Recurrentgemma.
  - Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, and 24 others. 2024. Olmo: Accelerating the science of language models. *Preprint*, arXiv:2402.00838.
  - Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024. Deepseek-coder: When the large language model meets programming – the rise of code intelligence. *Preprint*, arXiv:2401.14196.
  - Jan Philipp Harries. 2023. orca\_mini\_v2\_ger\_7b: An explain tuned llama-7b model based on orca mini v2 and adapted to german language.
  - Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. 2022. Toxigen: A large-scale machine-generated dataset

for adversarial and implicit hate speech detection. *Preprint*, arXiv:2203.09509.

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. Measuring massive multitask language understanding. *Preprint*, arXiv:2009.03300.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the math dataset. *Preprint*, arXiv:2103.03874.
- Alex Henry, Prudhvi Raj Dachapally, Shubham Pawar, and Yuxuan Chen. 2020. Query-key normalization for transformers. *Preprint*, arXiv:2010.04245.
- Jiwoo Hong, Noah Lee, and James Thorne. 2024. Orpo: Monolithic preference optimization without reference model. *Preprint*, arXiv:2403.07691.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. *Preprint*, arXiv:1902.00751.
- Chan-Jan Hsu, Chang-Le Liu, Feng-Ting Liao, Po-Chun Hsu, Yi-Chang Chen, and Da-Shan Shiu. 2024. Breeze-7b technical report. *Preprint*, arXiv:2403.02712.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

IDEA-CCNL. 2021. Fengshenbang-lm.

- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*.
- "interstellarninja", "Teknium", "theemozilla", "karan4d", and "huemin\_art". 2024. Hermes-2-pro-mistral-7b.
- Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Daniel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, Xian Li, Brian O'Horo, Gabriel Pereyra, Jeff Wang, Christopher Dewan, Asli Celikyilmaz, Luke Zettlemoyer, and Ves Stoyanov. 2023. Opt-iml: Scaling language model instruction meta learning through the lens of generalization. *Preprint*, arXiv:2212.12017.
- Neel Jain, Ping yeh Chiang, Yuxin Wen, John Kirchenbauer, Hong-Min Chu, Gowthami Somepalli, Brian R. Bartoldson, Bhavya Kailkhura, Avi Schwarzschild, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2023. Neftune: Noisy embeddings improve instruction finetuning. *Preprint*, arXiv:2310.05914.

- 1030 1031 1033 1040 1041 1042 1043 1044 1045 1046 1047 1048 1049 1050 1051 1052 1054 1055 1056 1057 1058 1059 1060 1061 1064 1066 1068 1070 1071 1072 1073 1074 1075 1076 1077

1078 1079

- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. Preprint, arXiv:2310.06825.
- Dongfu Jiang, Yishan Li, Ge Zhang, Wenhao Huang, Bill Yuchen Lin, and Wenhu Chen. 2024. Tigerscore: Towards building explainable metric for all text generation tasks. Preprint, arXiv:2310.00752.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer, 2017. Triviaga: A large scale distantly supervised challenge dataset for reading comprehension. Preprint, arXiv:1705.03551.
- Dahyun Kim, Yungi Kim, Wonho Song, Hyeonwoo Kim, Yunsu Kim, Sanghoon Kim, and Chanjun Park. 2024a. sdpo: Don't use your data all at once. Preprint, arXiv:2403.19270.
- Dahyun Kim, Chanjun Park, Sanghoon Kim, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, Changbae Ahn, Seonghoon Yang, Sukyung Lee, Hyunbyung Park, Gyoungjin Gim, Mikyoung Cha, Hwalsuk Lee, and Sunghun Kim. 2024b. Solar 10.7b: Scaling large language models with simple yet effective depth upscaling. Preprint, arXiv:2312.15166.
- Dohyeong Kim, Myeongjun Jang, Deuk Sin Kwon, and Eric Davis. 2022. Kobest: Korean balanced evaluation of significant tasks. Preprint, arXiv:2204.04541.
- Seungduk Kim, Seungtaek Choi, and Myeongho Jeong. 2024c. Efficient and effective vocabulary expansion towards multilingual large language models. Preprint, arXiv:2402.14714.
- Hyunwoong Ko, Kichang Yang, Minho Ryu, Taekyoon Choi, Seungmu Yang, Jiwung Hyun, Sungho Park, and Kyubyong Park. 2023. A technical report for polyglot-ko: Open-source large-scale korean language models. Preprint, arXiv:2306.02254.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Large language models are zero-shot reasoners. Preprint, arXiv:2205.11916.
- Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa, Joshua Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. 2023. Sparse upcycling: Training mixture-of-experts from dense checkpoints. Preprint, arXiv:2212.05055.
- Sunjun Kweon, Junu Kim, Jiyoun Kim, Sujeong Im, Eunbyeol Cho, Seongsu Bae, Jungwoo Oh, Gyubok Lee, Jong Hak Moon, Seng Chan You, Seungjin Baek, Chang Hoon Han, Yoon Bin Jung, Yohan Jo, and Edward Choi. 2024. Publicly shareable clinical large language model built on synthetic clinical notes. Preprint, arXiv:2309.00237.

Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, 1085 Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, 1086 Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, 1088 David Glushkov, Arnav Dantuluri, Andrew Maguire, 1089 Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. 2023. Openassistant conversations - democ-1092 ratizing large language model alignment. Preprint, arXiv:2304.07327. 1093

1094

1095

1096

1098

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

- Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. Biomistral: A collection of opensource pretrained large language models for medical domains. Preprint, arXiv:2402.10373.
- Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. 2019. Quantifying the carbon emissions of machine learning. Preprint, arXiv:1910.09700.
- Ariel N. Lee, Cole J. Hunter, and Nataniel Ruiz. 2024. Platypus: Ouick, cheap, and powerful refinement of llms. Preprint, arXiv:2308.07317.
- Ariel N. Lee, Cole J. Hunter, Nataniel Ruiz, Bleys Goodson, Wing Lian, Guan Wang, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknium". 2023. Openorcaplatypus: Llama2-13b model instruct-tuned on filtered openorcav1 gpt-4 dataset and merged with divergent stem and logic dataset model.
- Junbum Lee. 2023. llama-2-ko-7b (revision 4a9993e).
- Junbum Lee. 2024a. Llama-3-koen.
- Junbum Lee. 2024b. Yi-ko-6b (revision 205083a).
- Jungwon Lee and Seungjun Ahn. 2024. Llama-3instruction-constructionsafety-layertuning.
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. 2020. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. Preprint, arXiv:2005.01643.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, and 48 others. 2023a. Starcoder: may the source be with you! Preprint, arXiv:2305.06161.
- Shenggui Li, Hongxin Liu, Zhengda Bian, Jiarui Fang, Haichen Huang, Yuliang Liu, Boxiang Wang, and Yang You. 2023b. Colossal-ai: A unified deep learning system for large-scale parallel training. Preprint, arXiv:2110.14883.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del 1134 Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023c. 1135 Textbooks are all you need ii: phi-1.5 technical report. 1136 Preprint, arXiv:2309.05463. 1137

- 1138 Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023d. Chatdoctor: A medical 1139 chat model fine-tuned on a large language model 1140 meta-ai (llama) using medical domain knowledge. 1141 1142 *Preprint*, arXiv:2303.14070. Wing Lian, Bleys Goodson, Eugene Pentland, Austin 1143 Cook, Chanvichet Vong, and "Teknium". 2023a. 1144 1145 Openorca\_preview1: A llama-13b model fine-tuned on small portion of openorcav1 dataset. 1146 Wing Lian, Bleys Goodson, Guan Wang, Eugene Pent-1147 land, Austin Cook, Chanvichet Vong, and "Teknium". 1148 2023b. Jackalope 7b: Mistral-7b model multi-turn 1149 chat tuned on filtered openorcav1 gpt-4 dataset. 1150 Wing Lian, Bleys Goodson, Guan Wang, Eugene Pent-1151 land, Austin Cook, Chanvichet Vong, and "Teknium". 1152 2023c. Llongorca7b: Llama2-7b model instruct-1153 tuned for long context on filtered openorcav1 gpt-4 1154 1155 dataset. Wing Lian, Bleys Goodson, Guan Wang, Eugene Pent-1156 land, Austin Cook, Chanvichet Vong, and "Teknium". 1157 2023d. Mistralorca: Mistral-7b model instruct-tuned 1158 on filtered openorcav1 gpt-4 dataset. 1159 1160 Wing Lian, Bleys Goodson, Guan Wang, Eugene 1161 Pentland, Austin Cook, Chanvichet Vong, and "Teknium". 2023e. Mistralslimorca: Mistral-7b 1162 model instruct-tuned on filtered, corrected, openor-1163 cav1 gpt-4 dataset. 1164 Wing Lian, Guan Wang, Bleys Goodson, Eugene Pent-1165 land, Austin Cook, Chanvichet Vong, and "Teknium". 1166 1167 2023f. Slimorca: An open dataset of gpt-4 aug-1168 mented flan reasoning traces, with verification. Wing Lian, Guan Wang, Bleys Goodson, Eugene Pent-1169 land, Austin Cook, Chanvichet Vong, "Teknium", and 1170 Nathan Hoos. 2023g. Slimorca dedup: A dedupli-1171 cated subset of slimorca. 1172 Stephanie Lin, Jacob Hilton, and Owain Evans. 2022a. 1173 Truthfulga: Measuring how models mimic human 1174 falsehoods. Preprint, arXiv:2109.07958. 1175 Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu 1176 Wang, Shuohui Chen, Daniel Simig, Myle Ott, Na-1177 man Goyal, Shruti Bhosale, Jingfei Du, Ramakanth 1178 Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav 1179 Chaudhary, Brian O'Horo, Jeff Wang, Luke Zettle-1180 1181 moyer, Zornitsa Kozareva, Mona Diab, and 2 others. 1182 2022b. Few-shot learning with multilingual language models. Preprint, arXiv:2112.10668. 1183 Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. 1184 2025. World model on million-length video and 1185 1186 language with blockwise ringattention. Preprint, arXiv:2402.08268. 1187 Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paran-1188 jape, Michele Bevilacqua, Fabio Petroni, and Percy 1189 Liang. 2023. Lost in the middle: How language mod-1190 els use long contexts. Preprint, arXiv:2307.03172. 1191
- Zihan Liu, Wei Ping, Rajarshi Roy, Peng Xu, Chankyu Lee, Mohammad Shoeybi, and Bryan Catanzaro. 2024. Chatqa: Surpassing gpt-4 on conversational qa and rag. *Preprint*, arXiv:2401.10225.

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le, Barret Zoph, Jason Wei, and Adam Roberts. 2023. The flan collection: Designing data and methods for effective instruction tuning. *Preprint*, arXiv:2301.13688.
- Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, and 47 others. 2024. Starcoder 2 and the stack v2: The next generation. *Preprint*, arXiv:2402.19173.
- Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, Yansong Tang, and Dongmei Zhang. 2025. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *Preprint*, arXiv:2308.09583.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evolinstruct. *Preprint*, arXiv:2306.08568.
- Xingtai Lv, Ning Ding, Yujia Qin, Zhiyuan Liu, and Maosong Sun. 2023. Parameter-efficient weight ensembling facilitates task-level knowledge transfer. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 270–282, Toronto, Canada. Association for Computational Linguistics.
- Manuel Romero. 2023. llama-2-coder-7b (revision d30d193).
- Pankaj Mathur. 2023a. orca\_mini\_v2\_7b: An explain tuned llama-7b model on uncensored wizardlm, alpaca, & dolly datasets.
- Pankaj Mathur. 2023b. orca\_mini\_v3\_13b: An orca style llama2-70b model.
- Pankaj Mathur. 2023c. orca\_mini\_v3\_7b: An explain tuned llama2-7b model.
- Rimon Melamed, Lucas Hurley McCabe, Tanay Wakhare, Yejin Kim, H. Howie Huang, and Enric Boix-Adserà. 2024. Prompts have evil twins. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 46–74, Miami, Florida, USA. Association for Computational Linguistics.
- Sean Meyn and Richard L. Tweedie. 2009. *Markov Chains and Stochastic Stability*, 2nd edition. Cambridge Mathematical Library. Cambridge University Press.

- 1249 1250 1251 1252 1253 1254 1255 1256 1257 1258 1259 1260 1261 1262 1263 1265 1266 1268 1269 1270 1271 1272 1273 1274 1275 1276 1277 1278 1279 1280 1281 1282 1284 1285 1286 1287 1290 1292 1293 1294 1295 1296 1297 1298 1299 1300

1301 1302

- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. Preprint, arXiv:1809.02789.
- Xu Ming. 2023. textgen: Implementation of language model finetune.
- Arindam Mitra, Luciano Del Corro, Shweti Mahajan, Andres Codas, Clarisse Simoes, Sahaj Agarwal, Xuxi Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Aggarwal, Hamid Palangi, Guoqing Zheng, Corby Rosset, Hamed Khanpour, and Ahmed Awadallah. 2023. Orca 2: Teaching small language models how to reason. Preprint, arXiv:2311.11045.
- Koh Mitsuda, Xinqi Chen, Toshiaki Wakatsuki, and Kei Sawada. 2024. rinna/llama-3-youko-8b.
- MosaicML NLP Team. 2023a. Introducing mpt-30b: Raising the bar for open-source foundation models. Accessed: 2023-06-22.
- MosaicML NLP Team. 2023b. Introducing mpt-7b: A new standard for open-source, commercially usable llms. Accessed: 2023-03-28.
- Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. 2024. Generative representational instruction tuning. Preprint, arXiv:2402.09906.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual general-Preprint, ization through multitask finetuning. arXiv:2211.01786.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4. Preprint, arXiv:2306.02707.
- Xavier Murias. 2023. Cybertron: Uniform neural alignment.
- Nexusflow.ai team. 2023. Nexusraven-v2: Surpassing gpt-4 for zero-shot function calling.
- Xuan-Phi Nguyen, Wenxuan Zhang, Xin Li, Mahani Aljunied, Zhiqiang Hu, Chenhui Shen, Yew Ken Chia, Xingxuan Li, Jianyu Wang, Qingyu Tan, Liying Cheng, Guanzheng Chen, Yue Deng, Sen Yang, Chaogun Liu, Hang Zhang, and Lidong Bing. 2024. Seallms - large language models for southeast asia. Preprint, arXiv:2312.00738.
- Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2023. Codegen: An open large language model for code with multi-turn program synthesis. Preprint, arXiv:2203.13474.

Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of gpt-4 on medical challenge problems. Preprint, arXiv:2303.13375.

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

- Krzysztof Ociepa, Łukasz Flis, Krzysztof Wróbel, Adrian Gwoździej, SpeakLeash Team, and Cyfronet Team. 2024a. Introducing bielik-7b-v0.1: Polish language model. Accessed: 2024-04-01.
- Krzysztof Ociepa, Łukasz Flis, Krzysztof Wróbel, Sebastian Kondracki, SpeakLeash Team, and Cyfronet Team. 2024b. Introducing bielik-7b-instruct-v0.1: Instruct polish language model. Accessed: 2024-04-01
- Krzysztof Ociepa, Łukasz Flis, Krzysztof Wróbel, Adrian Gwoździej, and Remigiusz Kinas. 2024c. Bielik 7b v0.1: A polish language model - development, insights, and evaluation. Preprint. arXiv:2410.18565.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. Gpt-4 technical report. Preprint, arXiv:2303.08774.
- Momose Oyama, Sho Yokoi, and Hidetoshi Shimodaira. 2023. Norm of word embedding encodes information gain. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 2108-2130, Singapore. Association for Computational Linguistics.
- Ankit Pal and Malaikannan Sankarasubbu. 2024a. Gemini goes to med school: Exploring the capabilities of multimodal large language models on medical challenge problems & hallucinations. Preprint, arXiv:2402.07023.
- Ankit Pal and Malaikannan Sankarasubbu. 2024b. Openbiollms: Advancing open-source large language models for healthcare and life sciences.
- Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White. 2024. Smaug: Fixing failure modes of preference optimisation with dpo-positive. *Preprint*, arXiv:2402.13228.
- Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel R. Bowman. 2022. Bbq: A hand-built bias benchmark for question answering. *Preprint*, arXiv:2110.08193.
- Leonid Pekelis, Michael Feil, Forrest Moret, Mark Huang, and Tiffany Peng. 2024. Llama 3 gradient: A series of long context models.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, 1355 Ruxandra Cojocaru, Alessandro Cappelli, Hamza 1356 Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, 1357

- 1359 1360 1362 1363 1364 1365 1366 1367 1369 1371 1372 1373 1374 1375 1376 1377 1378 1379 1380 1381 1382 1383 1384 1385 1388 1389 1390 1391 1392 1393 1394 1395 1396 1397 1398 1400 1401 1402 1403 1404 1405 1406 1407

1409

1410

1411

1412 1413 and Julien Launay. 2023. The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only. arXiv preprint arXiv:2306.01116.

- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023a. Instruction tuning with gpt-4. Preprint, arXiv:2304.03277.
- Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023b. Yarn: Efficient context window extension of large language models. Preprint, arXiv:2309.00071.
- Nikhil Pinnaparaju, Reshinth Adithyan, Duy Phung, Jonathan Tow, James Baicoianu, and Nathan Cooper. 2024. Stable code 3b.
- Kunat Pipatanakul, Phatrasek Jirabovonvisut, Potsawee Manakul, Sittipong Sripaisarnmongkol, Ruangsak Patomwong, Pathomporn Chokchainant, and Kasima Tharnpipitchai. 2023. Typhoon: Thai large language models. Preprint, arXiv:2312.13951.
- Marco Polignano, Pierpaolo Basile, and Giovanni Semeraro. 2024. Advanced natural-based interaction for the italian language: Llamantino-3-anita. Preprint, arXiv:2405.07101.
- Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy Lin. 2023. Rankvicuna: Zero-shot listwise document reranking with open-source large language models. Preprint, arXiv:2309.15088.
- Ofir Press, Noah A. Smith, and Mike Lewis. 2022. Train short, test long: Attention with linear biases enables input length extrapolation. Preprint, arXiv:2108.12409.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. Preprint, arXiv:2305.18290.
- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2020. Zero: Memory optimizations toward training trillion parameter models. Preprint, arXiv:1910.02054.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, and 7 others. 2024. Code llama: Open foundation models for code. Preprint, arXiv:2308.12950.
- Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution. Preprint, arXiv:1804.09301.
- Oscar Sainz, Iker García-Ferrero, Rodrigo Agerri, Oier Lopez de Lacalle, German Rigau, and Eneko Agirre. 2024. Gollie: Annotation guidelines improve zero-shot information-extraction. Preprint, arXiv:2310.03668.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Winogrande: An adversarial winograd schema challenge at scale. Preprint, arXiv:1907.10641.

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. 2019. Socialiga: Commonsense reasoning about social interactions. Preprint, arXiv:1904.09728.
- Akira Sasaki, Masato Hirakawa, Shintaro Horie, and Tomoaki Nakamura. 2023a. Elyza-japanese-llama-2-7b.
- Akira Sasaki, Masato Hirakawa, Shintaro Horie, Tomoaki Nakamura, Sam Passaglia, and Daisuke Oba. 2023b. Elyza-japanese-llama-2-13b.
- Kei Sawada, Tianyu Zhao, Makoto Shing, Kentaro Mitsui, Akio Kaga, Yukiya Hono, Toshiaki Wakatsuki, and Koh Mitsuda. 2024. Release of pre-trained models for the Japanese language. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 13898-13905. https://arxiv.org/abs/2404.01657.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. Preprint, arXiv:2402.03300.
- Pratyusha Sharma, Jordan T. Ash, and Dipendra Misra. 2023. The truth is in there: Improving reasoning in language models with layer-selective rank reduction. Preprint, arXiv:2312.13558.
- Noam Shazeer. 2019. Fast transformer decoding: One write-head is all you need. Preprint, arXiv:1911.02150.
- Noam Shazeer. 2020. Glu variants improve transformer. Preprint, arXiv:2002.05202.
- Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou, Margaret Li, Gergely Szilvasy, Rich James, Xi Victoria Lin, Noah A. Smith, Luke Zettlemoyer, Scott Yih, and Mike Lewis. 2024. In-context pretraining: Language modeling beyond document boundaries. *Preprint*, arXiv:2310.10638.
- Hidetoshi Shimodaira. 2001. Multiple comparisons of log-likelihoods and combining nonnested models with applications to phylogenetic tree selection. Communications in Statistics-Theory and Methods, 30(8-9):1751-1772.
- Hidetoshi Shimodaira and Ying Cao. 1998. A graphical technique for model selection diagnosis. The Institute of Statistical Mathematics Research Memorandum, 680.
- Hidetoshi Shimodaira and Masami Hasegawa. 2005. 1466 Assessing the Uncertainty in Phylogenetic Inference, 1467 pages 463-493. Springer New York, New York, NY. 1468

- 1469 1470
- 1471
- 1472
- 1473 1474
- 1475
- 1476
- 1477 1478
- 1479 1480 1481
- 1482
- 1483 1484
- 1485
- 1486 1487
- 1488
- 1489 1490
- 1491 1492
- 1493 1494
- 1495
- 1496

- 1499
- 1500 1501 1502
- 1503 1504
- 1505
- 1507 1508

1509

- 1510 1511
- 1512 1513 1514
- 1515

1516

1

1518 1519 1520

1522

1521

1523 1524

- Hidetoshi Shimodaira and Yoshikazu Terada. 2019. Selective inference for testing trees and edges in phylogenetics. *Frontiers in ecology and evolution*, 7:174.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2020. Megatron-Im: Training multi-billion parameter language models using model parallelism. *Preprint*, arXiv:1909.08053.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Aguera y Arcas, Dale Webster, and 11 others. 2022. Large language models encode clinical knowledge. *Preprint*, arXiv:2212.13138.
  - Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Aguera y Arcas, and 12 others. 2023. Towards expert-level medical question answering with large language models. *Preprint*, arXiv:2305.09617.
- Shashank Sonkar, Naiming Liu, Debshila Basu Mallick, and Richard G. Baraniuk. 2023. Class: A design framework for building intelligent tutoring systems based on learning science principles. *Preprint*, arXiv:2305.13272.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, and 432 others. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Preprint*, arXiv:2206.04615.
- Stability AI Language Team. 2024. Stable lm 2 1.6b.
  - Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2023a. Roformer: Enhanced transformer with rotary position embedding. *Preprint*, arXiv:2104.09864.
  - Yixuan Su, Tian Lan, and Deng Cai. 2023b. Openalpaca: A fully open-source instruction-following model based on openllama.
  - Shivchander Sudalairaj, Abhishek Bhandwaldar, Aldo Pareja, Kai Xu, David D. Cox, and Akash Srivastava. 2024. Lab: Large-scale alignment for chatbots. *Preprint*, arXiv:2403.01081.
  - Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. 2024. A simple and effective pruning approach for large language models. *Preprint*, arXiv:2306.11695.

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and<br/>Jonathan Berant. 2019. Commonsenseqa: A question<br/>answering challenge targeting commonsense knowl-<br/>edge. Preprint, arXiv:1811.00937.15251526<br/>1527<br/>15281526

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1563

1564

1566

1570

- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model.
- Yi Tay, Mostafa Dehghani, Vinh Q. Tran, Xavier Garcia, Jason Wei, Xuezhi Wang, Hyung Won Chung, Siamak Shakeri, Dara Bahri, Tal Schuster, Huaixiu Steven Zheng, Denny Zhou, Neil Houlsby, and Donald Metzler. 2023. Ul2: Unifying language learning paradigms. *Preprint*, arXiv:2205.05131.
- Yi Tay, Jason Wei, Hyung Won Chung, Vinh Q. Tran, David R. So, Siamak Shakeri, Xavier Garcia, Huaixiu Steven Zheng, Jinfeng Rao, Aakanksha Chowdhery, Denny Zhou, Donald Metzler, Slav Petrov, Neil Houlsby, Quoc V. Le, and Mostafa Dehghani. 2022. Transcending scaling laws with 0.1% extra compute. *Preprint*, arXiv:2210.11399.
- "Teknium", Charles Goddard, "interstellarninja", "theemozilla", "karan4d", and "huemin\_art". 2024a. Hermes-2-theta-llama-3-8b.
- "Teknium", "interstellarninja", "theemozilla", "karan4d", and "huemin\_art". 2024b. Hermes-2-prollama-3-8b.
- "Teknium", "theemozilla", "karan4d", and "huemin\_art". 2024c. Nous hermes 2 mistral 7b dpo.
- Migel Tissera. 2023. Synthia-7b-v1.3: Synthetic intelligent agent.
- Together Computer. 2023. Redpajama-data: An open source recipe to reproduce llama training dataset.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models. *Preprint*, arXiv:2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023b. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.
- Jonathan Tow. 2023. Stablelm alpha v2 models. 1573
- Jonathan Tow, Marco Bellagente, Dakota Mahan, and Carlos Riquelme. 2023. Stablelm 3b 4e1t. 1575

- 1577 1578 1579 1580 1581 1582 1583 1583
- 15
- 1586 1587
- 1588
- 1590 1591
- 1593 1594 1595
- 1596 1597
- 1598 1599
- 1601

- 1602 1603
- 1603
- 1605 1606
- 1607 1608
- 1609 1610 1611

1613 1614

1612

- 1615 1616 1617
- 1618
- 1619
- 1621 1622

1623 1624 1625

- 1620
- 1628 1629

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine

Fourrier, Nathan Habib, Nathan Sarrazin, Omar San-

seviero, Alexander M. Rush, and Thomas Wolf. 2023.

Zephyr: Direct distillation of lm alignment. Preprint,

Laurens van der Maaten and Geoffrey Hinton. 2008.

Bram Vanroy. 2024. Geitje 7b ultra: A conversational

Gaël Varoquaux, Lars Buitinck, Gilles Louppe, Olivier

Grisel, Fabian Pedregosa, and Andreas Mueller. 2015.

Scikit-learn: Machine learning without learning the

machinery. GetMobile Mob. Comput. Commun.,

Leandro von Werra, Alex Havrilla, Max reciprocated,

Jonathan Tow, Aman cat state, Duy V. Phung,

Louis Castricato, Shahbuland Matiana, Alan, Ayush

Thakur, Alexey Bukhtiyarov, aaronrmm, Fabrizio

Milo, Daniel, Daniel King, Dong Shin, Ethan Kim,

Justin Wei, Manuel Romero, and 10 others. 2023.

CarperAI/trlx: v0.6.0: LLaMa (Alpaca), Benchmark

Ben Wang. 2021. Mesh-Transformer-JAX: Model-

Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A

Guan Wang, Sijie Cheng, Qiying Yu, and Changling

Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li,

Sen Song, and Yang Liu. 2024a. Openchat: Advanc-

ing open-source language models with mixed-quality

Guan Wang, Bleys Goodson, Wing Lian, Eugene Pent-

land, Austin Cook, Chanvichet Vong, and "Teknium".

2023b. Openorcaxopenchatpreview2: Llama2-13b

model instruct-tuned on filtered openorcav1 gpt-4

Shenzhi Wang, Yaowei Zheng, Guoyin Wang, Shiji

Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack

Hessel, Tushar Khot, Khyathi Raghavi Chandu,

David Wadden, Kelsey MacMillan, Noah A. Smith,

Iz Beltagy, and Hannaneh Hajishirzi. 2023c. How far

can camels go? exploring the state of instruction tun-

ing on open resources. Preprint, arXiv:2306.04751.

Zhuoshu Li, and Y. Wu. 2024c. Let the expert

stick to his last: Expert-specialized fine-tuning for

Zihan Wang, Deli Chen, Damai Dai, Runxin Xu,

Song, and Gao Huang. 2024b. Llama3-8b-chinese-

guage Models with Imperfect Data.

data. Preprint, arXiv:2309.11235.

chat (revision 6622a23).

6 Billion Parameter Autoregressive Language Model.

Liu. 2023a. OpenChat: Advancing Open-source Lan-

Parallel Implementation of Transformer Language

model for dutch. *Preprint*, arXiv:2412.04092.

Learning Research, 9(86):2579–2605.

Visualizing data using t-sne. Journal of Machine

arXiv:2310.16944.

19(1):29-33.

Util, T5 ILQL, Tests.

Model with JAX.

dataset.

Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R.
Bowman. 2020. Blimp: The benchmark of linguistic minimal pairs for english. *Transactions of the Association for Computational Linguistics*, 8:377–392.

1630

1631

1632

1633

1634

1636

1637

1638

1639

1640

1641

1642

1643

1644

1645

1646

1647

1649

1650

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1668

1669

1670

1671

1672

1674

1675

1678

1679

- Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig Schmidt. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. *Preprint*, arXiv:2203.05482.
- Haoyuan Wu, Haisheng Zheng, Zhuolun He, and Bei Yu. 2024. Parameter-efficient sparsity crafting from dense to mixture-of-experts for instruction tuning on general tasks. *Preprint*, arXiv:2401.02731.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. 2024. Sheared llama: Accelerating language model pre-training via structured pruning. *Preprint*, arXiv:2310.06694.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient streaming language models with attention sinks. *Preprint*, arXiv:2309.17453.
- Shitao Xiao, Zheng Liu, Peitian Zhang, and Xingrun Xing. 2023. Lm-cocktail: Resilient tuning of language models via model merging. *Preprint*, arXiv:2311.13534.
- Huajian Xin, Daya Guo, Zhihong Shao, Zhizhou Ren, Qihao Zhu, Bo Liu, Chong Ruan, Wenda Li, and Xiaodan Liang. 2024a. Deepseek-prover: Advancing theorem proving in llms through large-scale synthetic data. *Preprint*, arXiv:2405.14333.
- Huajian Xin, Z. Z. Ren, Junxiao Song, Zhihong Shao, Wanjia Zhao, Haocheng Wang, Bo Liu, Liyue Zhang, Xuan Lu, Qiushi Du, Wenjun Gao, Qihao Zhu, Dejian Yang, Zhibin Gou, Z. F. Wu, Fuli Luo, and Chong Ruan. 2024b. Deepseek-prover-v1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search. *Preprint*, arXiv:2408.08152.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023a. Wizardlm: Empowering large language models to follow complex instructions. *Preprint*, arXiv:2304.12244.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023b. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *Preprint*, arXiv:2304.01196.
- Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024a. A paradigm shift in machine translation: Boosting translation performance of large language models. *Preprint*, arXiv:2309.11674.
- 19

- Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024b. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. *Preprint*, arXiv:2401.08417.
- Xwin-LM Team. 2023. Xwin-lm.

1687

1688

1690

1692

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

1723

1724

1725

1726

1727

1728

1729

1730

1731

1732

1733

1734

1735

1736

1737

1738

- Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. 2023a. Ties-merging: Resolving interference when merging models. *Preprint*, arXiv:2306.01708.
  - Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. 2023b. Ties-merging: Resolving interference when merging models. In *Advances in Neural Information Processing Systems*, volume 36, pages 7093–7115. Curran Associates, Inc.
  - An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 43 others. 2024a. Qwen2 technical report. *Preprint*, arXiv:2407.10671.
  - Kailai Yang, Tianlin Zhang, Ziyan Kuang, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2024b. Mentallama: Interpretable mental health analysis on social media with large language models. *Preprint*, arXiv:2309.13567.
  - Ping Yang, Junjie Wang, Ruyi Gan, Xinyu Zhu, Lin Zhang, Ziwei Wu, Xinyu Gao, Jiaxing Zhang, and Tetsuya Sakai. 2022. Zero-shot learners for natural language understanding via a unified multiple choice perspective. *Preprint*, arXiv:2210.08590.
  - Nicolas Yax, Pierre-Yves Oudeyer, and Stefano Palminteri. 2024. Phylolm : Inferring the phylogeny of large language models and predicting their performances in benchmarks. *Preprint*, arXiv:2404.04671.
- Nicolas Yax, Pierre-Yves Oudeyer, and Stefano Palminteri. 2025. PhyloLM: Inferring the phylogeny of large language models and predicting their performances in benchmarks. In *The Thirteenth International Conference on Learning Representations*.
- Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024a. Language models are super mario: Absorbing abilities from homologous models as a free lunch. *Preprint*, arXiv:2311.03099.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024b. Metamath: Bootstrap your own mathematical questions for large language models. *Preprint*, arXiv:2309.12284.

- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, and 3 others. 2024a. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *Preprint*, arXiv:2311.16502.
  Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhu Chen.
- Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhu Chen. 2024b. Mammoth2: Scaling instructions from the web. *Preprint*, arXiv:2405.03548.

1749

1750

1751

1752

1753

1754

1755

1756

1757

1758

1759

1760

1761

1762

1763

1764

1765

1766

1767

1768

1769

1770

1771

1772

1773

1774

1775

1776

1777

1778

1779

1780

1781

1782

1783

1784

1785

1786

1787

1788

1789

1790

1791

1792

- YuLan-Team. 2023. Yulan-chat: An open-source bilingual chatbot.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? *Preprint*, arXiv:1905.07830.
- Biao Zhang and Rico Sennrich. 2019. Root mean square layer normalization. *Preprint*, arXiv:1910.07467.
- Chen Zhang, Dawei Song, Zheyu Ye, and Yan Gao. 2024a. Towards the law of capacity gap in distilling language models. *Preprint*, arXiv:2311.07052.
- Ge Zhang, Xinrun Du, Bei Chen, Yiming Liang, Tongxu Luo, Tianyu Zheng, Kang Zhu, Yuyang Cheng, Chunpu Xu, Shuyue Guo, Haoran Zhang, Xingwei Qu, Junjie Wang, Ruibin Yuan, Yizhi Li, Zekun Wang, Yudong Liu, Yu-Hsuan Tsai, Fengji Zhang, and 3 others. 2024b. Cmmmu: A chinese massive multi-discipline multimodal understanding benchmark. *Preprint*, arXiv:2401.11944.
- Jiaxing Zhang, Ruyi Gan, Junjie Wang, Yuxiang Zhang, Lin Zhang, Ping Yang, Xinyu Gao, Ziwei Wu, Xiaoqun Dong, Junqing He, Jianheng Zhuo, Qi Yang, Yongfeng Huang, Xiayu Li, Yanghan Wu, Junyu Lu, Xinyu Zhu, Weifeng Chen, Ting Han, and 6 others. 2022a. Fengshenbang 1.0: Being the foundation of chinese cognitive intelligence. *CoRR*, abs/2209.02970.
- Kaiyan Zhang, Ning Ding, Biqing Qi, Sihang Zeng, Haoxin Li, Xuekai Zhu, Zhang-Ren Chen, and Bowen Zhou. 2024c. Ultramedical: Building specialized generalists in biomedicine.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022b. Opt: Open pre-trained transformer language models. *Preprint*, arXiv:2205.01068.
- Wenxuan Zhang, Sharifah Mahani Aljunied, Chang Gao, Yew Ken Chia, and Lidong Bing. 2023. M3exam: A multilingual, multimodal, multilevel benchmark for examining large language models. *Preprint*, arXiv:2306.05179.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Or-

Tianyu Zhao, Akio Kaga, and Kei Sawada. 2023a.

Tianyu Zhao and Kei Sawada. 2023. rinna/japanese-gpt-

Tianyu Zhao, Toshiaki Wakatsuki, Akio Kaga, Koh Mit-

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan

Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,

Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,

Joseph E. Gonzalez, and Ion Stoica. 2023. Judg-

ing llm-as-a-judge with mt-bench and chatbot arena.

Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo

Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu

Chen, and Nan Duan. 2023. Agieval: A human-

centric benchmark for evaluating foundation models.

Xinyu Zhou, Delong Chen, Samuel Cahyawijaya,

Xufeng Duan, and Zhenguang Cai. 2025. Linguis-

tic minimal pairs elicit linguistic similarity in large

language models. In Proceedings of the 31st Inter-

national Conference on Computational Linguistics,

pages 6866-6888, Abu Dhabi, UAE. Association for

Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu,

Banghua Zhu, Hiteshi Sharma, Felipe Vieira Frujeri,

Shi Dong, Chenguang Zhu, Michael I. Jordan, and

Jiantao Jiao. 2023b. Fine-tuning language models

with advantage-induced policy alignment. Preprint,

Sally Zhu, Ahmed M Ahmed, Rohith Kuditipudi, and

Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B.

Brown, Alec Radford, Dario Amodei, Paul Chris-

tiano, and Geoffrey Irving. 2020. Fine-tuning lan-

guage models from human preferences. Preprint,

In recent years, research on comparing large lan-

guage models (LLMs) has gained attention. This

section provides an overview of existing studies

Percy Liang. 2025. Independence tests for language

Karthik Ganesan, Wei-Lin Chiang, Jian Zhang, and

Jiantao Jiao. 2023a. Starling-7b: Improving llm help-

Preprint, arXiv:2306.05685.

Preprint, arXiv:2304.06364.

Computational Linguistics.

arXiv:2306.02231.

arXiv:1909.08593.

**Related Work** 

models.

Δ

fulness & harmlessness with rlaif.

suda, and Kei Sawada. 2023b. rinna/bilingual-gpt-

methods. Preprint, arXiv:1804.06876.

rinna/youri-7b.

neox-3.6b.

neox-4b.

donez, and Kai-Wei Chang. 2018. Gender bias in

coreference resolution: Evaluation and debiasing

- 1801
- 1802 1803
- 1804
- 1805 1806
- 1807 1808 1809
- 1810
- 1811 1812
- 1813 1814 1815
- 1816
- 1817
- 1818 1819
- 1820 1821

1822

1824 1825

1826

1827 1828 1829

1830 1831

1832 1833

1834

1835 1836

1837 1838

1839

1840

1841 1842 1843 from three perspectives: model parameters, activations<sup>10</sup>, and probability distributions. 1844

1846

1847

1848

1849

1851

1852

1853

1854

1855

1856

1857

1860

1861

1862

1864

1865

1866

1870

1871

1872

1874

1876

1877

1878

1879

1882

1883

1884

1886

1887

**Comparison of model parameters.** One approach to comparing LLMs is to analyze their parameters. Zhu et al. (2025) proposed a statistical framework for evaluating parameter similarity between different models and introduced a method for determining whether these models were trained independently. Additionally, Yadav et al. (2023b) focused on parameter changes due to task adaptation, specifically analyzing task vectors<sup>11</sup>. They proposed a method to mitigate interference when integrating task vectors from different models. Specifically, by reducing redundant numerical components and adjusting for conflicting signs, their approach enables effective model merging.

**Comparison of activations.** Comparisons of LLMs based on activations have also been studied. Zhou et al. (2025) quantified the similarity between LLMs by measuring the cosine similarity of activations for linguistic minimal pairs. In particular, they used datasets such as BLiMP (Warstadt et al., 2020) and showed that model similarity is significantly influenced by the pre-training dataset.

Comparison of probability distributions. Several approaches compare LLMs using probability distributions. Lv et al. (2023) proposed a method for computing coefficients in parameter ensembling by providing the same input text to two models and comparing the softmax probability distributions at each token. Specifically, they used KL divergence and summed the results to derive appropriate coefficients. Furthermore, Yax et al. (2025) proposed a similarity metric based on the conditional probabilities of LLMs and introduced a method for calculating the phylogenetic distance between different models. Additionally, a method has been proposed for measuring differences in conditional probabilities based on prompts using KL divergence (Melamed et al., 2024).

## **B** Double Centering

We confirm the notation and computational operations. The matrices L,  $\Xi$ , and Q are all of size  $K \times N$ , and their elements are denoted by

<sup>&</sup>lt;sup>10</sup>In general, "activations" refer to the intermediate outputs of Transformer models (e.g., the residual stream or neurons) (Bereska and Gavves, 2024).

<sup>&</sup>lt;sup>11</sup>The difference in parameters before and after fine-tuning is referred to as a task vector, and arithmetic operations on these vectors function effectively (Ilharco et al., 2023).

 $\ell_{is}, \xi_{is}, q_{is}$ , respectively. In particular, we have 1888  $\ell_{is} = \ell_i(x_s)$ . First, row-wise centering of L is 1889 performed by subtracting the mean log-likelihood 1890  $\overline{\ell}_i$  of each model from each row  $(\ell_{i1}, \ldots, \ell_{iN})$ , 1891 resulting in  $\boldsymbol{\Xi} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_K)^{\top}$ . Next, column-1892 wise centering of  $\Xi$  is performed by subtracting 1893 the coordinate component  $\bar{\xi}_s$  of the mean vec-1894 tor  $\bar{\boldsymbol{\xi}}$  from each column  $(\xi_{1s}, \ldots, \xi_{Ks})^{\top}$ , yielding 1895  $\boldsymbol{Q} = (\boldsymbol{q}_1, \dots, \boldsymbol{q}_K)^{\top}$ . Thus, this process involves 1896 double centering, where column-wise centering 1897 follows row-wise centering. Notably, even after 1898 column-wise centering, the row-wise mean of Q1899 remains zero: 1900

1901  

$$\frac{1}{N} \sum_{s=1}^{N} q_{is} = \frac{1}{N} \sum_{s=1}^{N} (\xi_{is} - \bar{\xi}_s)$$
1902  

$$= \frac{1}{N} \sum_{s=1}^{N} \xi_{is} - \frac{1}{NK} \sum_{i=1}^{K} \sum_{s=1}^{N} \xi_{is}$$
1903  

$$= 0 - 0 = 0.$$
(5)

The column-wise centering can be interpreted as follows. In the  $\boldsymbol{\xi}$ -coordinate system, the mean vector  $\bar{\boldsymbol{\xi}}$  of the *K* model coordinates  $\boldsymbol{\xi}_1, \ldots, \boldsymbol{\xi}_K$ can be regarded as representing an "average model." By redefining this average model as the new origin, we obtain the *q*-coordinate system. From the definition  $\boldsymbol{q}_i = \boldsymbol{\xi}_i - \bar{\boldsymbol{\xi}}$ , its mean satisfies

1911 
$$\sum_{i=1}^{K} q_i/K = \mathbf{0}.$$

1904 1905

1906

1907

1908

1909

1910

1912

1913

1914

1915

1917

1918

1920

1921

1922

1923

1924

1925

1926

1927

1928

The row-wise centering can be interpreted as follows. Let  $\mathbf{1}_N = (1, ..., 1)^\top \in \mathbb{R}^N$ . Since  $\bar{\ell}_i = \mathbf{1}_N^\top \boldsymbol{\ell}_i / N$ , we have  $\boldsymbol{\xi}_i = \boldsymbol{\ell}_i - \bar{\ell}_i \mathbf{1}_N$ . Thus,

$$\mathbf{1}_N^{\top} \boldsymbol{\xi}_i = \mathbf{1}_N^{\top} \boldsymbol{\ell}_i - \bar{\ell}_i N = 0,$$

and furthermore,  $\mathbf{1}_N^{\top} \bar{\boldsymbol{\xi}} = 0$ . From this, equation (5) is actually trivial, as

$$\mathbf{1}_N^{\top} \boldsymbol{q}_i = \mathbf{1}^{\top} (\boldsymbol{\xi}_i - \bar{\boldsymbol{\xi}}) = 0 - 0 = 0.$$

The row-wise centering implies that  $\xi_1, \ldots, \xi_K$ and  $q_1, \ldots, q_K$  lie in the subspace orthogonal to  $\mathbf{1}_N$ .

## C Effect of errors in model coordinates

We analyze the impact of additive errors  $\epsilon_{is}$  in the log-likelihood vector components  $\ell_{is}$  for  $i = 1, \ldots, K$  and  $s = 1, \ldots, N$ . Denoting the true values with an asterisk as  $\ell_{is}^*$ , the observed values can be expressed as

$$\ell_{is} = \ell_{is}^* + \epsilon_{is}$$

We decompose the error as follows:

6

$$a_{is} = a + b_i + c_s + d_{is},$$
 1930

1929

1931

1932

1945

1951

1952

1953

1954

1955

1957

1958

where we assume, without loss of generality, the constraints

$$\sum_{i=1}^{K} b_i = \sum_{s=1}^{N} c_s = \sum_{i=1}^{K} d_{is} = \sum_{s=1}^{N} d_{is} = 0.$$
 1933

Here, a,  $b_i$ , and  $c_s$  are bias terms, while  $d_{is}$  represents interaction terms. Using simple calculations, we obtain: 1936

$$\bar{\ell}_i = \frac{1}{N} \sum_{s=1}^N \ell_{is} = \bar{\ell}_i^* + a + b_i,$$
1937

$$\xi_{is} = \ell_{is} - \bar{\ell}_i = \xi_{is}^* + c_s + d_{is},$$
1930

$$\bar{\xi}_s = \frac{1}{K} \sum_{i=1}^{K} \xi_{is} = \bar{\xi}_s^* + c_s,$$
1941

$$q_{is} = \xi_{is} - \bar{\xi}_s = q_{is}^* + d_{is}.$$
 1943

Additionally, for the differences between two models, which are crucial for the model map, we obtain:

$$\ell_{is} - \ell_{js} = \ell_{is}^* - \ell_{js}^* + b_i - b_j + d_{is} - d_{js},$$

$$\xi_{is} - \xi_{js} = \xi_{is}^* - \xi_{js}^* + d_{is} - d_{js},$$

$$q_{is} - q_{js} = q_{is}^* - q_{js}^* + d_{is} - d_{js}.$$

$$1946 \\
1947 \\
1948 \\
1949 \\
1949 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\
1950 \\$$

Thus, the terms affected by the error are  $\bar{\ell}_i$ , which is influenced by  $a + b_i$ , and  $\ell_{is} - \ell_{js}$ , which is affected by  $b_i + d_{is}$ , meaning it is influenced by the bias terms. However, in the centered values  $\xi_{is} - \xi_{js}$  and  $q_{is} - q_{js}$ , only the interaction term  $d_{is}$  contributes to the error.

## D Theory of Model Coordinates for Text Probability Distributions

In this section, we prove the main results of Sec-1959 tion 2, namely (2) and (3). Our discussion applies 1960 not only to text probability distributions but also 1961 more generally to any setting where i.i.d. obser-1962 vations  $x_1, \ldots, x_N \sim p_i(x)$  are available. Com-1963 pared to the previous study that proposed model 1964 maps (Shimodaira and Cao, 1998), we conduct a 1965 more precise analysis in this paper. Specifically, the previous study provided only a rough evalu-1967 ation of the approximation and omitted the cen-1968 tering argument in the proof of (2). In the next 1969 section, as the starting point of our discussion, we 1970 construct a *super model* that includes the K mod-1971 els  $p_i$ ,  $i = 1, \ldots, K$ , as submodels. This model 1972 is introduced as a mathematical tool to rigorously 1973 prove the main theorem of this paper, and we do not compute it numerically in practice. 1975

## 1978

- 979

1980

1981

1982

1983

1984

1985

1986

1987

1989

1990

1991

1992

1993

1995

1996

1997

1998

1999

2004

2007

## **D.1** Exponential family of distributions

We first consider a model in the exponential family of distributions parameterized by a K-dimensional parameter  $\boldsymbol{\theta} \in \mathbb{R}^{K}$ :

$$p(x; \boldsymbol{\theta}) = p_0(x) \exp(\boldsymbol{\theta}^\top \boldsymbol{b}(x) - \psi(\boldsymbol{\theta})). \quad (6)$$

Here, the function  $\boldsymbol{b}(x) = (b_1(x), \dots, b_K(x))^\top$ will be defined later using the K models. The normalization constant is given by

$$Z(\boldsymbol{\theta}) = \sum_{x \in \mathcal{X}} p_0(x) \exp(\boldsymbol{\theta}^\top \boldsymbol{b}(x)),$$
  
$$\psi(\boldsymbol{\theta}) = \log Z(\boldsymbol{\theta}),$$

which ensures that  $\sum_{x \in \mathcal{X}} p(x; \theta) = 1$ . For  $\theta = 0$ , where  $\mathbf{0} = (0, \dots, 0)^{\top}$ , we obtain

$$p(x;\mathbf{0}) = p_0(x).$$

To associate the K models with (6), we define  $\boldsymbol{b}(x)$ . For a constant  $\lambda > 0$ , we set

$$\lambda b_i(x) := \ell_i(x) - \ell_0(x). \tag{7}$$

The constant  $\lambda$  is an order parameter introduced for theoretical convenience, and in our theoretical framework, we assume that  $b_i(x)$  is of constant order and that  $\lambda$  is sufficiently small<sup>12</sup>. Thus, we essentially assume that  $|\ell_i(x) - \ell_0(x)| = O_p(\lambda)$  is sufficiently small, implying that each model  $p_i$  provides a good approximation of the true generative model  $p_0$ . In the proof of the main theorem, we consider the asymptotic theory as  $\lambda \to 0$ , retaining terms up to  $O(\lambda^2)$  while ignoring those of  $O(\lambda^3)$ . A one-hot vector is defined as  $e_i$  =

 $(0,\ldots,0,1,0,\ldots,0)^{\top} \in \mathbb{R}^{K}$  for  $i = 1,\ldots,K$ , where only the *i*-th element is 1. Then, setting  $\boldsymbol{\theta} = \lambda \boldsymbol{e}_i$  gives

$$p(x; \lambda \boldsymbol{e}_i) = p_i(x). \tag{8}$$

Indeed, substituting (7) into (6) yields

2008  

$$p(x; \lambda e_i)$$
2009  

$$=p_0(x) \exp(\lambda e_i^\top b(x) - \psi(\lambda e_i))$$
2010  

$$=p_0(x) \exp(\ell_i(x) - \ell_0(x) - \psi(\lambda e_i))$$
2011  

$$=p_0(x)(p_i(x)/p_0(x)) \exp(-\psi(\lambda e_i))$$
2012  

$$=p_i(x) \exp(-\psi(\lambda e_i))$$
2013  

$$=p_i(x),$$

where  $\psi(\lambda e_i) = 0$ . 2014

#### **D.2** Properties of the exponential family of distributions

This section outlines some well-known basic properties of the exponential family of distributions, which have been established in the literature (Barndorff-Nielsen, 2014; Efron, 1978, 2022; Amari, 1982). We define the expectation and covariance matrix of  $\boldsymbol{b}(x)$  as follows:

$$\boldsymbol{\eta}(\boldsymbol{\theta}) := \mathop{\mathbb{E}}_{x \sim p(\boldsymbol{\theta})} \left( \boldsymbol{b}(x) \right) = \sum_{x \in \mathcal{X}} \boldsymbol{b}(x) p(x; \boldsymbol{\theta}),$$
 2023

$$G(\boldsymbol{\theta}) := \mathop{\mathbb{E}}_{x \sim p(\boldsymbol{\theta})} \left\{ (\boldsymbol{b}(x) - \boldsymbol{\eta}(\boldsymbol{\theta})) (\boldsymbol{b}(x) - \boldsymbol{\eta}(\boldsymbol{\theta}))^{\top} \right\}$$

$$= \mathop{\mathrm{Var}}_{x \sim p(\boldsymbol{\theta})} (\boldsymbol{b}(x)).$$
2024

These quantities can be expressed in terms of  $\psi(\theta)$ as follows:

$$\boldsymbol{\eta}(\boldsymbol{\theta}) = \frac{\partial \psi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \quad (9) \quad 202$$

2016

2017

2019

2026

2039

2045

$$G(\boldsymbol{\theta}) = \frac{\partial^2 \psi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^{\top}}.$$
 (10) 20

We now derive these two equations. First, from

$$\frac{\partial Z(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{x \in \mathcal{X}} \boldsymbol{b}(x) p_0(x) e^{\boldsymbol{\theta}^\top \boldsymbol{b}(x)},$$
 2031

we obtain

$$\frac{\partial \psi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{\partial \log Z}{\partial \boldsymbol{\theta}} = \frac{1}{Z(\boldsymbol{\theta})} \frac{\partial Z}{\partial \boldsymbol{\theta}}$$
 2033

$$=\frac{1}{Z(\boldsymbol{\theta})}\sum_{x\in\mathcal{X}}\boldsymbol{b}(x)p_0(x)e^{\boldsymbol{\theta}^{\top}\boldsymbol{b}(x)}$$
2034

$$=\sum_{x\in\mathcal{X}}\boldsymbol{b}(x)p(x;\boldsymbol{\theta})=\boldsymbol{\eta}(\boldsymbol{\theta}).$$
 2035

Thus, we have established (9). Next, using 2036

$$\frac{\partial p(x;\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \left(\boldsymbol{b}(x) - \frac{\partial \psi}{\partial \boldsymbol{\theta}}\right) p(x;\boldsymbol{\theta})$$
2037

$$= (\boldsymbol{b}(x) - \boldsymbol{\eta}(\boldsymbol{\theta}))p(x; \boldsymbol{\theta}),$$
 2038

we obtain

$$\frac{\partial^2 \psi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^{\top}} = \frac{\partial \boldsymbol{\eta}(\boldsymbol{\theta})^{\top}}{\partial \boldsymbol{\theta}}$$
 2040

$$= \frac{\partial}{\partial \boldsymbol{\theta}} \sum_{x \in \mathcal{X}} \boldsymbol{b}(x)^\top \boldsymbol{p}(x; \boldsymbol{\theta})$$
 2041

$$= \sum_{x \in \mathcal{X}} (\boldsymbol{b}(x) - \boldsymbol{\eta}(\boldsymbol{\theta})) \boldsymbol{b}(x)^{\top} p(x; \boldsymbol{\theta})$$
 2042

$$= \sum_{x \in \mathcal{X}} (\boldsymbol{b}(x) - \boldsymbol{\eta}(\boldsymbol{\theta})) (\boldsymbol{b}(x) - \boldsymbol{\eta}(\boldsymbol{\theta}))^\top p(x; \boldsymbol{\theta})$$
 2043

$$=G(\boldsymbol{\theta}).$$
 2044

Thus, we have established (10).

 $<sup>^{12}</sup>$  If numerical computation were to be performed,  $\lambda$  could be set to any arbitrary value (e.g.,  $\lambda = 1$ ).

**D.3** 

 $\mathbb{R}^{K}$  is given by

 $\psi(\theta')$ 

 $\mathrm{KL}(p(\boldsymbol{\theta}), p(\boldsymbol{\theta}'))$ 

 $= \sum_{\boldsymbol{x}} p(\boldsymbol{x}; \boldsymbol{\theta}) \log \frac{p(\boldsymbol{x}; \boldsymbol{\theta})}{p(\boldsymbol{x}; \boldsymbol{\theta}')}$ 

 $= (\boldsymbol{\theta} - \boldsymbol{\theta}')^{\top} \boldsymbol{\eta}(\boldsymbol{\theta}) - \psi(\boldsymbol{\theta}) + \psi(\boldsymbol{\theta}').$ 

## 2047 2048

- 2050
- 2051

2054

2056

2058

2067

2070

2071

2072

2073

2074

2075

2077

2078

2079

2062

$$=\psi(\boldsymbol{\theta}) + \boldsymbol{\eta}(\boldsymbol{\theta})^{\top}(\boldsymbol{\theta}' - \boldsymbol{\theta}) \\ + \frac{1}{2}(\boldsymbol{\theta}' - \boldsymbol{\theta})^{\top}G(\boldsymbol{\theta})(\boldsymbol{\theta}' - \boldsymbol{\theta}) + O(\lambda^3).$$
(12)

 $+\frac{1}{2}(\boldsymbol{\theta}'-\boldsymbol{\theta})^{\top}\frac{\partial^{2}\psi(\boldsymbol{\theta})}{\partial\boldsymbol{\theta}\partial\boldsymbol{\theta}^{\top}}(\boldsymbol{\theta}'-\boldsymbol{\theta})$ 

Substituting (12) into (11) gives

 $=\psi(\theta)+rac{\partial\psi}{\partial\theta^{ op}}(\theta'-\theta)$ 

 $+ O(\|\boldsymbol{\theta}' - \boldsymbol{\theta}\|^3)$ 

$$\operatorname{KL}(p(\boldsymbol{\theta}), p(\boldsymbol{\theta}')) = \frac{1}{2} (\boldsymbol{\theta}' - \boldsymbol{\theta})^{\top} G(\boldsymbol{\theta}) (\boldsymbol{\theta}' - \boldsymbol{\theta}) + O(\lambda^3). \quad (13)$$

Approximation of the KL divergence

The Kullback-Leibler (KL) divergence between the

models  $p(\theta)$  and  $p(\theta')$  at parameter values  $\theta, \theta' \in$ 

 $= \sum_{x \in \mathcal{X}} p(x; \boldsymbol{\theta}) \big\{ (\boldsymbol{\theta} - \boldsymbol{\theta}')^\top \boldsymbol{b}(x) - \psi(\boldsymbol{\theta}) + \psi(\boldsymbol{\theta}') \big\}$ 

Here, we assume that the parameter values  $\theta$  and  $\theta'$ 

are sufficiently close to 0. In particular, we assume  $\|\boldsymbol{\theta}\| = O(\lambda)$  and  $\|\boldsymbol{\theta}'\| = O(\lambda)$ . Substituting (9)

and (10) into the Taylor expansion of  $\psi(\theta)$  gives

This corresponds to eq. (9) of Oyama et al. (2023). Here, the equation holds approximately by ignoring higher-order terms of  $O(\lambda^3)$ . For more details, refer to Amari (1982, p. 369) and Efron (2022, p. 35). More generally,  $G(\theta)$  represents the Fisher information metric, and (14) holds for a wide class of probability models (Amari, 1998). Furthermore, since  $G(\boldsymbol{\theta}) = G(\mathbf{0}) + O(\|\boldsymbol{\theta}\|) = G(\mathbf{0}) + O(\lambda)$ , we obtain

$$\operatorname{KL}(p(\boldsymbol{\theta}), p(\boldsymbol{\theta}')) = \frac{1}{2} (\boldsymbol{\theta}' - \boldsymbol{\theta})^{\top} G(\mathbf{0}) (\boldsymbol{\theta}' - \boldsymbol{\theta}) + O(\lambda^3). \quad (14)$$

#### **D.4** The variance representation of the KL divergence

Substituting  $p_i = p(\lambda e_i)$  into (14) gives 2080

2081  

$$2\mathrm{KL}(p_i, p_j) = 2\mathrm{KL}(p(\lambda \boldsymbol{e}_i), p(\lambda \boldsymbol{e}_j))$$

$$= \lambda^2 (\boldsymbol{e}_i - \boldsymbol{e}_j)^\top G(\mathbf{0}) (\boldsymbol{e}_i - \boldsymbol{e}_j) + O(\lambda^3). \quad (15)$$

Here, we have

(11)

$$\boldsymbol{e}_i^\top G(\mathbf{0}) \boldsymbol{e}_j = G_{ij}(\mathbf{0})$$
 2084

$$= \mathop{\mathbb{E}}_{x \sim p_0} \Big\{ (b_i(x) - \eta_i(\mathbf{0})) (b_j(x) - \eta_j(\mathbf{0})) \Big\}.$$
(16) 208

Next, we substitute (16) and (7) into the right-hand side of (15) to derive an alternative expression for 2087 the KL divergence: 2088

$$\lambda^2 (\boldsymbol{e}_i - \boldsymbol{e}_j)^\top G(\boldsymbol{0}) (\boldsymbol{e}_i - \boldsymbol{e}_j)$$
 2089

$$=\lambda^2 \mathop{\mathbb{E}}_{x \sim p_0} \left[ \left\{ (b_i(x) - \eta_i(\mathbf{0})) \right\} \right]$$
 2090

$$-\left(b_j(x) - \eta_j(\mathbf{0})\right)\Big\}^2\Big]$$
 209

$$= \mathop{\mathbb{E}}_{x \sim p_0} \left[ \left\{ (\ell_i(x) - \ell_0(x)) - (\ell_j(x) - \ell_0(x)) \right\} \right]$$
2092

$$- \mathop{\mathbb{E}}_{x' \sim p_0} \left( \left( \ell_i(x') - \ell_0(x') \right) \right)$$
 2093

$$-\left(\ell_{j}(x')-\ell_{0}(x')\right)\Big\}^{2}\Big]$$
 2094

$$= \mathop{\mathbb{E}}_{x \sim p_0} \left[ \left\{ \ell_i(x) - \ell_j(x) \right\}^2 \right]$$

$$= \left[ \left\{ \ell_i(x) - \ell_j(x) \right\}^2 \right]$$

$$- \mathop{\mathbb{E}}_{x' \sim p_0} \left( \ell_i(x') - \ell_j(x') \right) \Big\} \qquad \qquad 209$$
$$= \mathop{\mathrm{Var}}_{x \sim p_0} \left( \ell_i(x) - \ell_j(x) \right). \qquad \qquad 209$$

Finally, substituting this result into (15) yields

$$2\mathrm{KL}(p_i, p_j) = \operatorname{Var}_{x \sim p_0} \left( \ell_i(x) - \ell_j(x) \right) + O(\lambda^3).$$
(17)

This establishes (2). Furthermore, since  $|\ell_i(x) - \ell_i(x)| \leq |\ell_i(x)|$ 2100  $|\ell_i(x)| = O_p(\lambda)$ , the magnitude of (17) is  $O(\lambda^2)$ . 2101

## D.5 Estimation of the KL divergence

If the expected value  $\mathbb{E}_{x \sim p_0}(f(x))$  of a function f(x) exists and is bounded, then by the law of large numbers, the sample mean<sup>13</sup> converges to the expected value as  $N \to \infty$ , and we have

$$\mathop{\mathbb{E}}_{x \sim D}(f(x)) = \mathop{\mathbb{E}}_{x \sim p_0}(f(x)) + O_p(N^{-1/2}).$$
 2107

Applying this to (17) and ignoring the terms of order  $O_p(\lambda^3 + \lambda^2 N^{-1/2})$ , we obtain the following approximation:

$$2\mathrm{KL}(p_i, p_j) \approx \operatorname{Var}_{x \sim D} \left( \ell_i(x) - \ell_j(x) \right).$$
(18) 2111

 ${}^{13}\mathbb{E}_{x \sim D}(f(x)) = \frac{1}{N} \sum_{x \in D} f(x) = \frac{1}{N} \sum_{s=1}^{N} f(x_s)$  represents the sample mean, and  $\operatorname{Var}_{x \sim D}(\cdot)$  represents the sample variance.

2102

2103 2104 2105

2106

2108

2122

2123

2124

2125

2126

We define the coordinates  $\boldsymbol{\xi}_i \in \mathbb{R}^N$  of model  $p_i$  as

2113 
$$\boldsymbol{\xi}_i = (\xi_{i1}, \dots, \xi_{iN})^{\top}$$

with

$$\xi_{is} := \ell_i(x_s) - \mathop{\mathbb{E}}_{x \sim D}(\ell_i(x))$$

for  $s = 1, \ldots, N$ . From (18), we obtain 2116

2117 
$$2\text{KL}(p_i, p_j) \approx \frac{1}{N} \sum_{s=1}^{N} (\xi_{is} - \xi_{js})^2$$
$$= \frac{1}{N} \|\boldsymbol{\xi}_i - \boldsymbol{\xi}_j\|^2. \tag{19}$$

Since 2119

2120 
$$\|\boldsymbol{\xi}_i - \boldsymbol{\xi}_j\|^2 = \|(\boldsymbol{q}_i + \bar{\boldsymbol{\xi}}) - (\boldsymbol{q}_j + \bar{\boldsymbol{\xi}})\|^2 = \|\boldsymbol{q}_i - \boldsymbol{q}_j\|^2$$

this establishes (3). 2121

#### **Relationships among the three types of D.6** model coordinates

Let  $\mathbf{1}_N = (1, \dots, 1)^\top \in \mathbb{R}^N$ . From the definitions of the  $\xi$ -coordinate system and the q-coordinate system, we have

2127 
$$\xi_i = q_i + \bar{\xi},$$
  
2128  $\ell_i = \xi_i + \bar{\ell}_i \mathbf{1}_N$   
2129  $= q_i + \bar{\ell}_i \mathbf{1}_N + \bar{\xi}.$  (20)

Additionally, equation (5) in Appendix B can be 2130 rewritten as 2131

$$\mathbf{1}_N^\top \boldsymbol{q}_i = 0. \tag{21}$$

Thus, we obtain 2133

2136

2137 2138

2132

2135 
$$= \|(q_i - q_j) + (\bar{\ell}_i - \bar{\ell}_i)\|$$

$$\begin{split} \|\boldsymbol{\ell}_{i} - \boldsymbol{\ell}_{j}\|^{2} \\ = \|(\boldsymbol{q}_{i} - \boldsymbol{q}_{j}) + (\bar{\ell}_{i} - \bar{\ell}_{j})\mathbf{1}_{N}\|^{2} \\ = \|\boldsymbol{q}_{i} - \boldsymbol{q}_{j}\|^{2} + N(\bar{\ell}_{i} - \bar{\ell}_{j})^{2} \end{split}$$

$$+2(\bar{\ell}_i-\bar{\ell}_j)\mathbf{1}_N^\top(\boldsymbol{q}_i-\boldsymbol{q}_j)$$

$$= \|\boldsymbol{q}_i - \boldsymbol{q}_j\|^2 + N(\ell_i - \ell_j)^2,$$

where (20) and (21) are used in the first and last 2139 equations, respectively. This establishes (4). 2140

Moreover, since  $\bar{\ell}_i = \mathbf{1}_N^{\top} \boldsymbol{\ell}_i / N$ , it is straightfor-2141 ward that the component of the  $\ell$ -coordinate system 2142 in the  $\mathbf{1}_N$  direction is given by 2143

2144 
$$(\mathbf{1}_N/\sqrt{N})^{\top} \boldsymbol{\ell}_i = \sqrt{N} \bar{\ell}_i.$$

### Ε Mapping Language Models into the **Space of Token Probability Distributions**

2145

2146

2147

2168

2169

2172

2174

2175

2176

2177

2178

2179 2180

2182

2183

2184

2185

In Section 2, we discussed model maps based on 2148 the probability distributions  $p_i(x)$  of texts gener-2149 ated by language models. This approach requires 2150 computing probabilities for a large number of texts 2151 in the dataset  $D = (x_1, \ldots, x_N)$ , leading to high 2152 computational costs. To mitigate this issue, we 2153 focus on the fact that a text  $x = (y_1, \ldots, y_n)$  is a 2154 sequence of tokens. Instead of using text probabil-2155 ities, we discuss model maps based on the condi-2156 tional probability distributions of token generation, 2157  $p_i(y_t|y^{t-1})$ . In this approach, model coordinates 2158 are computed using only a single text x. A limita-2159 tion of this approach is that it can only be used for 2160 comparing models that share the same tokenizer. 2161 Furthermore, the current estimation method ignores 2162 the variance of the expected log-likelihood ratio of 2163 conditional probabilities, resulting in a rough ap-2164 proximation. Thus, the estimated values should 2165 be regarded only as reference values rather than 2166 precise measurements. 2167

## E.1 Model coordinates

For a text  $x = (y_1, \ldots, y_n)$ , the coordinates of model  $p_i$ 

$$\boldsymbol{\zeta}_i = (\zeta_{i1}, \dots, \zeta_{in})^\top \in \mathbb{R}^n$$
 2171

are defined as

$$\zeta_{it} := \log p_i(y_t | y^{t-1}) - \ell_i(x) / n \qquad (22) \qquad 2173$$

for t = 1, ..., n. This is centered for each *i* and for each text, satisfying  $\sum_{t=1}^{n} \zeta_{it} = 0$ .

## E.2 Kullback-Leiber divergence

The KL divergence for next-token generation in language models, where  $y_t \sim p_i(y_t|y^{t-1})$ , is given by

$$\mathrm{KL}(p_i(y_t|y^{t-1}), p_j(y_t|y^{t-1})) =$$
(23)

$$\sum_{y_t \in \mathcal{V}} p_i(y_t | y^{t-1}) \log \frac{p_i(y_t | y^{t-1})}{p_j(y_t | y^{t-1})}.$$
 (24) 2181

We apply the results for text probability distributions from Section 2 and Appendix D to the conditional probability distributions of token generation. The equation corresponding to (2) is

$$2\mathrm{KL}(p_i(y_t|y^{t-1}), p_j(y_t|y^{t-1})) \approx 2180$$

$$\operatorname{Var}_{y_t \sim p_0(y_t|y^{t-1})} \left\{ \log \frac{p_i(y_t|y^{t-1})}{p_j(y_t|y^{t-1})} \right\}.$$
 (25) 2187

2228

2231

2236

2244

2252

2253

2260

2263

2264

2267

The squared Euclidean distance in the  $\zeta$ -coordinate system provides an estimate of the sum of (25) over all tokens in the text x:

2191 
$$\|\boldsymbol{\zeta}_i - \boldsymbol{\zeta}_j\|^2$$

2188

2189

2190

2192

2193

2195

2196

2197

2198

2199

2200

2201

2202

2205

2206

2207

2208

2212

2214

2216

2218

2220

2221

2223

2224

2226

$$\approx 2 \sum_{t=1} \mathrm{KL}(p_i(y_t|y^{t-1}), p_j(y_t|y^{t-1})) \qquad (26)$$

$$\approx 2 \mathrm{KL}(p_i, p_j).$$
 (27)

The proof is provided in Appendix F. To justify the estimation in (26), we assume the following:

$$\mathbb{E}_{y_t \sim p_0(y_t | y^{t-1})} \left\{ \log \frac{p_i(y_t | y^{t-1})}{p_j(y_t | y^{t-1})} \right\}$$
(28)

takes a constant value independent of t. In reality, this assumption is not entirely correct, and the degree of variation affects the accuracy of the approximation in  $(26)^{14}$ . On the other hand, the approximation in (27) holds more generally and is demonstrated in Appendix F.5.

#### F **Theory of Model Coordinates for Token Probability Distributions**

In this section, we provide a more detailed explanation of the content discussed in Appendix E. We extend the discussion of text probability distributions in Appendix D to the case of conditional probability distributions for token generation.

## F.1 Exponential family of distributions

We apply the same setting as for  $p_i(x)$  in Section D to the conditional probability distributions of tokens:

$$y_t \sim p_i(y_t | y^{t-1}), \quad t = 1, \dots, n.$$

The exponential family of distributions incorporating K models, corresponding to (6), is given here as

$$p(y_t|y^{t-1};\boldsymbol{\theta}) := p_0(y_t|y^{t-1})$$
$$\exp(\boldsymbol{\theta}^{\top}\boldsymbol{b}(y_t|y^{t-1}) - \psi(\boldsymbol{\theta}|y^{t-1})). \quad (29)$$

The setting (7), which associates the K models with (29), is given here as

$$\lambda b_i(y_t|y^{t-1}) := \log p_i(y_t|y^{t-1}) - \log p_0(y_t|y^{t-1}).$$
(30)

Thus, we have

$$p_i(y_t|y^{t-1}) = p(y_t|y^{t-1}; \lambda \boldsymbol{e}_i)$$

for i = 1, ..., K.

#### **F.2** The variance representation of the KL divergence

The KL divergence is given by (24). Applying the result for the model  $p(x; \theta)$  in (17) to the tokenlevel conditional distribution model  $p(y_t|y^{t-1}; \boldsymbol{\theta})$ , we obtain

$$2\text{KL}(p_i(y_t|y^{t-1}), p_j(y_t|y^{t-1})) = 223$$

$$\operatorname{Var}_{y_t \sim p_0(y_t|y^{t-1})} \left\{ \log \frac{p_i(y_t|y^{t-1})}{p_j(y_t|y^{t-1})} \right\} + O(\lambda^3).$$
(31)

#### Two additional assumptions **F.3**

To estimate the KL divergence from a single text x, two additional assumptions are required, as described below. Such assumptions were not necessary when estimating the KL divergence from the dataset D in Appendix D. In reality, these two assumptions are not strictly satisfied, and the discrepancy between these assumptions and reality affects the accuracy of the KL divergence approximation.

**Assumption 1:** We assume that the probability distribution of  $y_t$  depends only on the past k tokens, denoted as  $y_{t-k}^{t-1} = (y_{t-k}, y_{t-k+1}, \dots, y_{t-1})$ . That is,

$$p_i(y_t|y^{t-1}) = p_i(y_t|y^{t-1}_{t-k}),$$

which allows us to regard  $y_{t-k}^t$  as the state of a Markov chain. More generally, we use the notation  $y^k$  to represent a state variable. We consider a function f of the state variable  $y^k$ . Furthermore, we assume that this Markov chain is positive Harris recurrent, has a stationary distribution  $\pi$ , and that f is absolutely integrable, i.e.,

$$\mathbb{E}_{k \sim \pi}(|f(y^k)|) < \infty.$$
225

Then, by the strong law of large numbers for Markov chains (Meyn and Tweedie, 2009, Theorem 17.0.1 (i)), in the limit  $n \to \infty$ ,

y

For simplicity in notation and discussion, we assume that  $y_{-k+1}, \ldots, y_0$  are appropriately defined. Since the Markov chain converges to  $\pi$ , we also have

$$\frac{1}{n}\sum_{t=1}^{n}f(y_{t-k}^{t}) \to \frac{1}{n}\sum_{t=1}^{n}\sum_{y^{t}\sim p_{0}}^{\mathbb{E}}(f(y_{t-k}^{t})) \quad (32)$$

almost surely as  $n \to \infty$ .

<sup>&</sup>lt;sup>14</sup>Since this assumption does not affect (3), there is no concern regarding the use of model maps based on text probabilities.

## **Assumption 2:**

$$\mathbb{E}_{y_t \sim p_0(y_t|y^{t-1})} \log \frac{p_i(y_t|y^{t-1})}{p_j(y_t|y^{t-1})} = c \qquad (33)$$

for some  $c \in \mathbb{R}$  that can depend on the indices iand j but not on t. In other words, (33) takes a constant value independent of t.

#### F.4 Estimation of the KL divergence 2272

Define 2273

2268

2269

2270

2271

2275 2276

2277

2278

2280

2284

2286

2287

2289

2290

2291

2292

2293

2295 2296

# $h(y^{t}) := \log \frac{p_{i}(y_{t}|y^{t-1})}{p_{i}(y_{t}|y^{t-1})}.$

From Assumption 1,  $h(y^t)$  can be written in the form  $h(y^t) = f_1(y_{t-k}^t)$  for some  $f_1$ , so applying (32), for sufficiently large n, we obtain

$$\frac{1}{n}\sum_{t=1}^n h(y^t) \approx \frac{1}{n}\sum_{t=1}^n \mathop{\mathbb{E}}_{y^t \sim p_0}(h(y^t)).$$

Applying (33) to the right-hand side gives 2279

$$\frac{1}{n}\sum_{t=1}^n h(y^t) \approx c$$

Next, since  $(h(y^t) - c)^2$  can be written in the form 2281 of  $f_2(y_{t-k}^t)$  for some function  $f_2$ , applying (32) again yields 2283

$$\begin{split} &\frac{1}{n} \sum_{t=1}^{n} (h(y^{t}) - c)^{2} \\ &\approx \frac{1}{n} \sum_{t=1}^{n} \sum_{y^{t-1} \sim p_{0}}^{n} \left\{ \sum_{y_{t} \sim p_{0}(y_{t}|y^{t-1})}^{N} (h(y^{t})) \right\} \\ &\approx \frac{1}{n} \sum_{t=1}^{n} \sum_{y_{t} \sim p_{0}(y_{t}|y^{t-1})}^{N} (h(y^{t})). \end{split}$$

In the final equation, we applied (32) using the fact that  $\operatorname{Var}_{y_t \sim p_0(y_t|y^{t-1})}(h(y^t)) = f_3(y_{t-k}^t)$  for some  $f_3$ . Using (31), we obtain

$$\sum_{t=1}^{n} (h(y^{t}) - c)^{2} \approx 2\sum_{t=1}^{n} \operatorname{KL}(p_{i}(y_{t}|y^{t-1}), p_{j}(y_{t}|y^{t-1})). \quad (34)$$

Meanwhile, the components of the model coordinate  $\zeta_i$  are given by

2294 
$$\zeta_{it} = \log p_i(y_t|y^{t-1}) -$$

where

$$c_i = \frac{1}{n} \sum_{t=1}^n \log p_i(y_t | y^{t-1}).$$

 $c_i$ 

Since

$$\zeta_{it} - \zeta_{jt} = h(y^t) - (c_i - c_j)$$
<sup>2298</sup>

with  $c_i - c_j \approx c$ , equation (34) can be rewritten as

$$\|oldsymbol{\zeta}_i-oldsymbol{\zeta}_j\|^2pprox$$
2300

$$2\sum_{t=1}^{n} \mathrm{KL}(p_i(y_t|y^{t-1}), p_j(y_t|y^{t-1})).$$
2301

Thus, (26) is established.

#### Connecting the KL divergence of token F.5 and text probability distributions

Here, we fix the sequence length of the text x = $(y_1,\ldots,y_n)$  as n, i.e., we set  $\mathcal{X} = \mathcal{V}^n$ . For notational simplicity, we define

$$g_i(y^t) = \log p_i(y_t|y^{t-1}).$$
 2308

Noting that

$$p_i(x) = \prod_{t=1}^n p_i(y_t | y^{t-1}) = \prod_{t=1}^n e^{g_i(y^t)},$$
2310

we obtain

=

$$\operatorname{KL}(p_i, p_j)$$
 2312

$$= \sum_{x \in \mathcal{X}} \prod_{t'=1}^{n} e^{g_i(y^{t'})} \sum_{t=1}^{n} (g_i(y^t) - g_j(y^t))$$
 2313

$$=\sum_{t=1}^{n}\sum_{y^{t}\in\mathcal{V}^{t}}\prod_{t'=1}^{t}e^{g_{i}(y^{t'})}(g_{i}(y^{t})-g_{j}(y^{t}))$$
2314

$$= \sum_{t=1}^{n} \sum_{y^{t-1} \in \mathcal{V}^{t-1}} \prod_{t'=1}^{t-1} e^{g_i(y^{t'})}$$
231

$$\sum_{y_t \in \mathcal{V}} e^{g_i(y^t)} (g_i(y^t) - g_j(y^t))$$
2316

$$=\sum_{t=1}^{n}\sum_{y^{t-1}\in\mathcal{V}^{t-1}}p_{i}(y^{t-1})$$
2317

$$\mathrm{KL}(p_i(y_t|y^{t-1}), p_j(y_t|y^{t-1}))$$
2318

$$= \sum_{t=1}^{n} \mathbb{E}_{y^{t-1} \sim p_i} \Big\{ \mathrm{KL}(p_i(y_t|y^{t-1}), p_j(y_t|y^{t-1})) \Big\}$$
 2319

$$= \mathop{\mathbb{E}}_{x \sim p_i} \left\{ \sum_{t=1}^n \operatorname{KL}(p_i(y_t | y^{t-1}), p_j(y_t | y^{t-1})) \right\}.$$
 2320

Thus, for sufficiently large n, by the strong law of large numbers for Markov chains, we obtain

$$\mathrm{KL}(p_i, p_j) \approx \sum_{t=1}^{n} \mathrm{KL}(p_i(y_t | y^{t-1}), p_j(y_t | y^{t-1})).$$
(35)

2297

2299

2303

2304

2305

2307

2309

G

2324

## 2330

2331 2333

2334

2337

2340

2342

2347

2349

2350

2358

2359

2361 2362

2363

2367

2368

2371

G.1 Information obtained via the Hugging Face Hub API We used the Hugging Face Hub API to retrieve

**Details of Experiments** 

information about each language model's tags, the date the model was created, the number of downloads over the past 30 days, and the model's configuration details. All of this information is current as of February 1, 2025.

This corresponds to (27). Assumption 1 from Ap-

pendix F.3 is used in (35), but Assumption 2 is not

needed in the discussion of this subsection.

Among the model tags, we specifically used llama2, llama-2, license:llama2, llama3, llama-3, and license:llama3 to determine the model type (llama-1, llama-2, or llama-3). Furthermore, to identify language models that were pre-trained on the Pile in Section 4, we employed tags such as dataset:eleutherai/pile, dataset:eleutherai/the\_pile,

dataset:eleutherai/the\_pile\_deduplicated, and arxiv:2101.00027.

## G.2 How the model type was determined

In principle, we used the value of model\_type in the config retrieved from the Hugging Face Hub API as the model type. However, out of the 1,018 language models we examined, there were 587 whose config model\_type was listed as llama. For these, we used the following procedure to determine whether they were the original Llama (llama-1), Llama-2, or Llama-3; if we were able to identify which version they were, we reclassified them as llama-1, llama-2, or llama-3 accordingly.

- 1. We checked the tags assigned to each model. Of these, 136 models that included any of llama2, llama-2, or license:llama2 were classified as 11ama-2. Similarly, 39 models that included any of llama3, llama-3, or license:llama3 were classified as llama-3.
- 2. For the remaining 412 models whose classification was not determined by tags alone, we used the creation date and the model name (converted to lowercase) to make a decision. First, 69 models that were created prior to July 18, 2023 (the Llama-2 release date) were classified as 11ama-1. Next, 88 models whose lowercase model name contained either llama2 or llama-2 were classified as

Model type	Models
llama-1	69
llama-2	223
llama-3	62
llama	217
mistral	232
gpt_neox	54
deepseek	26
gptj	19
gemma	18
opt	15
bloom	12
falcon	11
qwen2	10
mixtral	9
mpt	6
stablelm	6
gpt_neo	3
phi	3
gpt_bigcode	3
phi3	3
xglm	3
rwkv	3
starcoder2	2
olmo	2
camelidae	2
codegen	2
deci	1
recurrent_gemma	1
<pre>stablelm_alpha</pre>	1
Total	1,018

Table 3: Number of models by model type.

llama-2. Among those whose lowercase model name contained llama3 or llama-3, 22 models whose creation date was after April 18, 2024 (the Llama-3 release date) were classified as 11ama-3.

2372

2373

2374

2378

2380

2381

2382

2383

2386

3. After following the steps above, the 217 models that could not be classified were left as llama.

Furthermore, any model whose name prior to the slash (/) was deepseek-ai was defined as deepseek. In addition, even though abacusai/Llama-3-Smaug-8B was tagged with license:llama2, we manually reclassified it as 11ama-3.

Table 3 shows the number of models classified

Text category	Texts
Pile-CC	2,353
PubMed Central	1,763
ArXiv	1,172
Github	925
FreeLaw	837
StackExchange	712
Wikipedia (en)	567
USPTO Backgrounds	487
PubMed Abstracts	464
Gutenberg (PG-19)	251
DM Mathematics	151
EuroParl	83
HackerNews	67
Ubuntu IRC	54
PhilPapers	51
NIH ExPorter	41
Enron Emails	22
Total	10,000

Table 4: Number of texts in each text category.

into each model type.

2390

2393

2394

2397

2400

2401

2402

2403

2404

2405

2407

2408

G.3 Basic information on the dataset

The dataset used in our experiments consists of a total of 10,000 texts, which are divided into 17 text categories. Table 4 shows the number of texts in each category.

To assign colors to the text categories, we first compute the average text embedding for each category in the Pile using simcse-roberta-large (Gao et al., 2021b). Next, we calculate a tour over the 17 average embedding vectors by solving the traveling salesman problem. Based on the adjacency relationships along this tour, we segment the hue circle at equal intervals and color each category accordingly.

G.4 Standard scores

In the three experiments described below, we use standardized values<sup>15</sup>, or Z-score normalization, of both the log-likelihood and the benchmark scores calculated for the K language models.

• In Figure 1, where we define each language model's primary text category, we use the av-

erage log-likelihood for each category, stan-2409 dardized across all models. 2410 • In Section 4, to determine each language 2411 model's primary task, we standardize each 2412 task's score across the K models. 2413 • Furthermore, in the data leakage detection de-2414 scribed in Section 4, we use the difference 2415 between the standardized mean log-likelihood 2416 and the standardized 6-TaskMean score as the 2417 indicator. 2418

2419

2420

2421

2422

2423

2424

2425

2426

2427

2428

2430

2431

2432

2433

2434

2436

2437

2438

2439

2440

2441

2442

2443

2444

2447

2448

## G.5 Hierarchical clustering settings

In Fig. 3, which visualizes the double-centered loglikelihood matrix Q, hierarchical clustering is applied to both the rows and the columns of Q. We use correlation as the metric and average as the linkage method. For the hierarchical clustering shown in Fig. 4, which presents a dendrogram of 100 language models, we use euclidean as the metric and ward as the linkage method.

## H Additional Model Maps

In this section, we present additional model maps, including a figure that lists all the model names and a map obtained through dimensionality reduction of the double-centered log-likelihood matrix Q.

## H.1 Model map with model names

Figure 11 shows the model map defined in Section 4, with each point labeled with its corresponding model name.

# H.2 Model map via the double centered log-likelihood matrix

In the main text, we use a model map generated by dimensionality reduction of the log-likelihood matrix L. Here, in Figs. 12 and 13, we present model maps obtained by dimensionality reduction of the double-centered log-likelihood matrix Q.

## I Details of Weight Interpolation

In this section, we describe the experimental details concerning the relationship between model coordinates and model weights, as introduced in Section 6.3.

Constructing the weight grid.Let  $p_0$  denote the2449base model, and  $p_1$  and  $p_2$  denote the fine-tuned2450models derived from  $p_0$ . We denote the weight2451parameter vectors of these models as  $W_0, W_1$ , and2452

<sup>&</sup>lt;sup>15</sup>By subtracting the mean from each value and dividing by the standard deviation, the data is transformed to have a mean of 0 and a variance of 1.



Figure 11: Each point on the model map is labeled with the corresponding model name.



Figure 12: Model maps obtained by double centered loglikelihood matrix Q. These maps correspond to Figure 1. (Top) Colors indicate model types. (Bottom) Colors indicate the text category in which each model attains the highest relative log-likelihood among 17 categories.

 $W_2$ , respectively. To construct the weight grid, we merged the model weights using the following linear operation:

2453 2454

2455

2456

2457

2459

2460

2461

2464

2465

2466

$$W_{\alpha,\beta} = W_0 + \alpha (W_1 - W_0) + \beta (W_2 - W_0),$$
(36)

where the merge ratios  $\alpha, \beta \in \mathbb{R}$  were chosen from 36 evenly spaced combinations within the interval [0,1]:  $\{0.0, 0.2, 0.6, 0.8, 1.0\}$ . The original models  $W_0, W_1$ , and  $W_2$  correspond to  $W_{0,0}, W_{1,0}$ , and  $W_{0,1}$ , respectively. When  $\alpha + \beta \leq 1$ , the operation corresponds to linear interpolation between models. Even among models with the same architecture, the sizes of the embedding/unembedding matrices may differ. In such cases, we truncated or reshaped the weight parameters to match the base model.

**2467 Computing model coordinates.** For each composed model  $p_{\alpha,\beta}$  with weight  $W_{\alpha,\beta}$ , we computed the model coordinates  $q_{\alpha,\beta}$  following the method described in Section 2. The tokenizer of the base model was used to ensure consistency. The text data consists of 1,000 randomly sampled sentences from the dataset of 10,000 sentences created in Section 3.1.

2469

2470

2471

2472

2474

2475

2476

2477

2478

2479

2480

2481

2482

2483

2484

2486

2487

2489

2491

2492

2494

2495

2496

2497

2499

2504

2505

2507

2509

2510

2511

2513

**Selection of models.** We set the base model  $p_0$  as Llama-2-7b-hf. Among the fine-tuned models based on Llama-2-7b-hf, we selected the two most downloaded models: vicuna-7b-v1.5 and Llama-2-7b-chat-hf, denoted as  $p_1$  and  $p_2$ , respectively.

**Visualization.** Figure 10 shows the linearly merged models, visualized in both weight space and log-likelihood space. The corners of the grid are labeled with their corresponding  $(\alpha, \beta)$  values.

In the left panel, we visualized  $W_{\alpha,\beta}$  in the weight space. Since the dimensionality of  $W_{\alpha,\beta}$ , i.e., the number of model parameters, is extremely high, we employed a 2D projection method using the norms of the difference vectors  $r_1 = ||W_1 - W_0||_2$ ,  $r_2 = ||W_2 - W_0||_2$ , and the angle between them,  $\phi = \arccos((W_1 - W_0)^\top (W_2 - W_0)/r_1r_2)$ . Each point was placed at  $(\alpha r_1 + \beta r_2 \cos \phi, \beta r_2 \sin \phi)$ .

In the right panel, we visualized the model coordinates  $q_{\alpha,\beta}$  by Principal Component Analysis (PCA). Each model  $p_{\alpha,\beta}$  was mapped onto the *q*coordinate system to analyze the structure of the interpolated models.

## J Details of Model Performance Prediction

This section provides additional details on the prediction of benchmark scores using model coordinates, as discussed in Section 5.

## J.1 Details of ridge regression

As described in Section 5.2, ridge regression requires setting a regularization strength parameter,  $\alpha$ . To determine  $\alpha$  from  $\{10^1, \ldots, 10^9\}$ , we performed a five-fold cross-validation within each training dataset<sup>16</sup>. As a post-processing step, we clipped the predicted scores to the range [0, 100].

For the setting where the target variable  $\mathbf{f}$  was replaced with the mean log-likelihood  $(\bar{\ell}_1, \ldots, \bar{\ell}_K) \in \mathbb{R}^K$ , we searched for  $\alpha$  within

<sup>&</sup>lt;sup>16</sup>The training dataset consisted of four out of the five folds obtained by splitting the dataset for each benchmark task.



Figure 13: Model maps obtained by double centered log-likelihood matrix Q. These maps correspond to Figure 5. These maps are illustrating model performance. From left to right, the panels show each model's mean log-likelihood, 6-TaskMean score, and the "primary task," which refers to the task where each model achieves the highest standardized score among the six tasks, color-coded accordingly. The color bar is clipped at the 10th percentile for mean log-likelihood and 6-TaskMean, with darker colors indicating better performance. In the primary task panel, models with standardized scores below zero across all six tasks are labeled as "All Under 0."



Figure 14: Scatter plots of predicted scores and benchmark scores for six benchmark tasks. Additionally, results for predicting 6-TaskMean (identical to Fig. 8) and the mean log-likelihood are also shown. Each point is color-coded by the mean log-likelihood, with higher mean log-likelihood values generally corresponding to higher task scores. For better visualization, the color bar range is clipped to the 10th–100th percentile.

	ARC	HellaSwag	MMLU	TruthfulQA	Winogrande	GSM8K	Average	mean log-likelihood
Pearson's $r$ Spearman's $\rho$	$\substack{0.969 \pm 0.002 \\ 0.975 \pm 0.001}$	$_{0.943\pm0.004}^{0.943\pm0.004}_{0.971\pm0.002}$	$\substack{0.959 \pm 0.003 \\ 0.966 \pm 0.003}$	$\substack{0.954 \pm 0.001 \\ 0.933 \pm 0.001}$	$_{0.961\pm0.002}^{0.961\pm0.002}$	$\substack{0.930 \pm 0.002 \\ 0.892 \pm 0.001}$	$\substack{0.973 \pm 0.001 \\ 0.976 \pm 0.001}$	$0.993 \pm 0.001$ $0.989 \pm 0.001$

Table 5: This table presents the mean and standard deviation of correlation coefficients between predicted and actual benchmark scores. These coefficients were computed using ridge regression across five different data splits. Results for predicting 6-TaskMean and the mean log-likelihood are also included.

 $\{10^{-4}, \ldots, 10^4\}$  and did not apply clipping as a post-processing step.

## J.2 Details of prediction results

2514

2516

2517

2518

2521

Figure 14 shows scatter plots of predicted scores and actual benchmark scores for each benchmark task, as well as for 6-TaskMean and the mean loglikelihood. As in Fig. 8, the scatter plots show strong correlations for all six benchmark tasks and for the mean log-likelihood.

To account for randomness, we ran five different2523data splits when predicting each benchmark score.2524The final predicted score was the average of these2525five runs. For each split, we computed the correla-2526tion coefficients between the predicted scores and2527the actual benchmark scores. Table 5 presents their2528mean and standard deviation, and shows a similar2529trend as Table 2.2530

## K Model List

2531

2533

2534

2535

2537

2539

2540

2542

2543

2544

2546

2547

2548

2551 2552

2554

2556

2557 2558

2559

2564

2570

2571

2572

Table 6 lists the 1,018 models used in this study. The BibTeX entries cited for each model were determined through the following procedure.

First, we extracted the BibTeX entries available in each model's Hugging Face model card<sup>17</sup>. If the BibTeX entry was missing a year of publication, we filled it in with the model's creation date<sup>18</sup>. Additionally, we generated BibTeX entries using the arXiv IDs found in the model card tags by querying the arXiv API<sup>19</sup>. This process resulted in a set of BibTeX entries for each model.

Next, we manually checked pairs of different BibTeX entries where the title similarity<sup>20</sup> was high, or the authors matched, to determine whether they corresponded to the same source. This step allowed us to create groups of BibTeX entries that were considered identical.

Then, for each BibTeX group, we selected a representative entry as follows. Within each group, the entry most frequently cited by the models was chosen as the representative. If multiple candidates met this criterion, we prioritized BibTeX entries generated from arXiv IDs when available. If no such entry existed, we selected the one with the longest string.

Finally, we replaced each model's BibTeX entry with the representative entry from its corresponding group. Any selected BibTeX entry that contained typos or formatting errors was manually corrected based on compilation errors. If the author information was incomplete, we corrected it manually by checking the source.

Note that for mistralai/Mistral-7B-v0.3, google/codegemma-2b,

deepseek-ai/deepseek-llm-7b-base

in Table 1, as well as for the example deepseek-ai/deepseek-coder-1. 3b-base, we manually prepared the BibTeX entries for citation based on their respective sources. We also used the same BibTeX entries for all other models that were considered to be of the same type.

<sup>19</sup>https://github.com/lukasschwab/arxiv.py.

<sup>&</sup>lt;sup>17</sup>https://github.com/huggingface/huggingface\_ hub.

<sup>&</sup>lt;sup>18</sup>https://github.com/sciunto-org/ python-bibtexparser.

<sup>&</sup>lt;sup>20</sup>We used Python's difflib.SequenceMatcher.

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
1	01-ai/Yi-1.5-6B (01. AI et al., 2025)	llama	6B	2024-05-11	5102	-525.57	61.57
2	01-ai/Yi-1.5-6B-Chat (01. AI et al., 2025)	llama	6B	2024-05-11	19443	-645.33	66.17
3	01-ai/Yi-1.5-9B (01. AI et al., 2025)	llama	8B	2024-05-11	20252	-513.39	66.73
4	01-ai/Yi-1.5-9B-32K (01. AI et al., 2025) 01 $ai/Yi-1.5-9B-52K$ (01. AI et al., 2025)	llama	8B 9D	2024-05-15	9287	-510.38	55.22 60.56
5	$01-ai/Yi-1.5-9B-Chat_16K (01 ALet al. 2025)$	llama	8B	2024-03-10	14262	-620.90	66.98
7	01-ai/Yi-6B (Zhang et al., 2024b; Yue et al., 2024a; 01. AI et al., 2025)	llama	6B	2023-11-01	6997	-515.62	54.02
8	01-ai/Yi-6B-200K (Zhang et al., 2024b; Yue et al., 2024a; 01. AI et al., 2025)	llama	6B	2023-11-06	8370	-529.12	56.76
9	01-ai/Yi-9B-200K (Zhang et al., 2024b; Yue et al., 2024a; 01. AI et al., 2025)	llama	8B	2024-03-15	8915	-517.82	61.94
10	42dot/42dot_LLM-PLM-1.3B (42dot Inc., 2023)	llama	1B	2023-09-04	1268	-600.03	35.70
11	42d0t/42d0t_LLM-SF1-1.3B (42d0t Inc., 2023) 922-CA/monika-ddlc-7b-y1	llama-2	1B 7B	2023-09-04	1455	-605.50	50.01 50.49
13	AIChenKai/TinyLlama-1.1B-Chat-v1.0-x2-MoE	mixtral	1B	2024-01-03	1217	-619.71	36.98
14	AIJUUD/juud-Mistral-7B	mistral	7B	2024-01-31	1312	-564.68	61.72
15	AIJUUD/juud-Mistral-7B-dpo (Lacoste et al., 2019)	mistral	7B	2024-02-07	3217	-567.00	60.89
16	AlekseyKorshuk/pygmalion-6b-vicuna-chatml	gptj	6B	2023-06-22	1169	-564.90	42.08
17	Andronooe/ retAnother_Open-Liania-5B-LokA	mistral	3B 7B	2023-07-21	9634	-040.40	69 57
19	Artples/L-MChat-Small	phi	2B	2024-04-11	2793	-640.49	63.14
20	Aspik101/StableBeluga-13B-instruct-PL-lora_unload	llama-2	13B	2023-08-04	1244	-544.64	56.24
21	Aspik101/WizardVicuna-Uncensored-3B-instruct-PL-lora_unload	llama-2	3B	2023-08-07	1161	-626.59	39.95
22	Aspik101/vicuna-13b-v1.5-PL-lora_unload	llama-2	13B	2023-08-03	1263	-566.80	55.24
23	AtAndDev/Capybaraiviarcoroni-/B Austism/chronos_hermes_13b_v2	mistrai	/B 13B	2024-01-03	1102	-541.55	70.32 56.10
25	Azazelle/Argetsu	mistral	7B	2023-12-30	1223	-560.37	69.64
26	Azazelle/Dumb-Maidlet	mistral	7B	2023-12-30	1201	-548.86	68.34
27	Azazelle/Half-NSFW_Noromaid-7b	mistral	7B	2023-12-29	1208	-547.63	62.32
28	Azazelle/Maylin-7b	mistral	7B	2024-01-04	1198	-575.94	70.26
29	Azazelle/Silicon-Medley	mistral	7B 7D	2023-12-29	1200	-566.66	69.49 60.70
30	Azazelle/Tippy-Toppy-7b	mistral	7B 7B	2023-12-30	1205	-560.94	69.58
32	Azazelle/Yuna-7b-Merge	mistral	7B	2024-01-05	1201	-580.65	71.46
33	Azazelle/smol_bruin-7b	mistral	7B	2023-12-29	1207	-572.80	71.05
34	Azazelle/xDAN-SlimOrca	mistral	7B	2023-12-29	1206	-575.20	68.04
35	Azure99/blossom-v1-3b	bloom	3B 2P	2023-07-29	1228	-641.52	36.90
30	Azure99/blossom-v2-llama2-7b	llama-2	5B 7B	2023-08-08	1237	-572.17	51 71
38	Azure99/blossom-v3-mistral-7b	mistral	7B	2023-11-20	1319	-565.31	62.95
39	Azure99/blossom-v3_1-mistral-7b	mistral	7B	2023-11-27	1320	-567.17	62.53
40	Azure99/blossom-v4-mistral-7b	mistral	7B	2023-12-26	1315	-550.37	63.61
41	BEE-spoke-data/Mixtral-GQA-400m-v2	mixtral	2B 7B	2023-12-20	1182	-794.25	28.45
42	BioMistral/BioMistral-7B (Labrak et al. 2024)	mistral	7В	2024-01-22	11349	-592.94	74.85 52.33
44	BioMistral/BioMistral-7B-DARE (Yadav et al., 2023a; Labrak et al., 2024;	mistral	7B	2024-02-05	1209	-586.37	57.03
	Yu et al., 2024a)						
45	BlueNipples/TimeCrystal-12-13B	llama-2	13B	2023-11-11	1282	-567.72	59.26
40 47	Bram vanroy/GEITje-/B-uitra (vanroy, 2024) Brouz/Sierpeno	llama	/B 13B	2024-01-27	1397	-0/4.95	56.59
48	CHIH-HUNG/Ilama-2-13b-FINETUNE2 TEST 2.2w	llama-2	13B 13B	2023-09-08	1174	-530.90	53.20
49	CHIH-HUNG/llama-2-13b-FINETUNE3_3.3w-r16-gate_up_down	llama-2	13B	2023-09-20	1171	-533.48	54.32
50	CHIH-HUNG/llama-2-13b-FINETUNE3_3.3w-r4-gate_up_down	llama-2	13B	2023-09-19	1166	-535.81	53.35
51	CHIH-HUNG/Ilama-2-13b-FINETUNE3_3.3w-r4-q_k_v_o	llama-2	13B	2023-09-19	1170	-536.50	53.62
52 53	CHIH-HUNG/Ilama-2-13b-FINETUNE3_3.3w-r8-gate_up_down	llama-2	13B 13B	2023-09-20	1165	-533.51	53.71
54	CHIH-HUNG/Ilama-2-13b-FINETUNE3 3.3w-r8-q k v o gate up down	llama-2	13B	2023-09-20	1165	-536.24	53.43
55	CHIH-HUNG/llama-2-13b-FINETUNE4_3.8w-r16-gate_up_down	llama-2	13B	2023-09-22	1165	-533.56	53.52
56	CHIH-HUNG/llama-2-13b-FINETUNE4_3.8w-r16-gate_up_down-test1	llama-2	13B	2023-10-07	1165	-531.27	53.66
57	CHIH-HUNG/Ilama-2-13b-FINETUNE4_3.8w-r16-q_k_v_o	llama-2	13B	2023-09-22	1164	-534.36	53.68
58	CHIH-HUNG/Ilama-2-13b-FINETUNE4_3.8w-r4-gate_up_down	llama-2	13B 12D	2023-09-21	11/4	-537.05	53.48
59 60	CHIH-HUNG/llama-2-13b-FINETUNE4_3.8w-r8-gate_up_down	llama-2	13B 13B	2023-09-23	1105	-540.47	53.25 53.58
61	CHIH-HUNG/Ilama-2-13b-FINETUNE4 3.8w-r8-g k v o	llama-2	13B	2023-09-21	1171	-533.89	54.06
62	CHIH-HUNG/Ilama-2-13b-FINETUNE4_3.8w-r8-q_k_v_o_gate_up_down	llama-2	13B	2023-09-25	1165	-537.12	52.88
63	CHIH-HUNG/llama-2-13b-FINETUNE4_addto15k_4.5w-r16-	llama-2	13B	2023-10-08	1161	-531.46	54.88
64	gate_up_down CHIH-HUNG/llama-2-13b-FINETUNE4_compare15k_4.5w-r16-	llama-2	13B	2023-10-08	1169	-532.00	53.94
65	gate_up_down	11 0	120	2022 10 04	11/7	524.00	52.44
65	CHIH-HUNG/Ilama-2-13b-FINETUNE5_4w-r16-gate_up_down	llama-2	13B 12P	2023-10-04	1167	-534.88	53.44
67	CHIH-HUNG/llama-2-13b-FINETUNE5_4w-r10-q_k_v_0	llama-2	13B 13B	2023-10-04	1172	-555.84	53 32
68	CHIH-HUNG/Ilama-2-13b-FINETUNE5 4w-r4-q k v o gate up down	llama-2	13B	2023-10-02	1163	-541.64	53.38
69	CHIH-HUNG/llama-2-13b-FINETUNE5_4w-r8-gate_up_down	llama-2	13B	2023-10-03	1163	-534.71	54.02
70	CHIH-HUNG/Ilama-2-13b-FINETUNE5_4w-r8-q_k_v_o	llama-2	13B	2023-10-02	1167	-534.88	54.64
71	CHIH-HUNG/Ilama-2-13b-FINETUNE5_4w-r8-q_k_v_o_gate_up_down	llama-2	13B	2023-10-04	1164	-540.68	53.69
12	CHIH-HUNG/Ilama-2-13b-OpenOrca_5W CHIH-HUNG/Ilama-2-13h-Open Platynus and cop 2 fw-3 epoch	nama-2 llama-2	13B 13R	2023-08-24	11/8	-555.15	53.80 53.80
74	CHIH-HUNG/Ilama-2-13b-dolphin_20w	llama-2	13B	2023-08-29	1187	-533.73	55.06
75	CHIH-HUNG/Ilama-2-13b-dolphin_5w	llama-2	13B	2023-08-25	1175	-533.03	55.53
76	CalderaAI/13B-BlueMethod	llama-1	13B	2023-07-07	1193	-569.51	54.12
77	CalderaAI/13B-Legerdemain-L2	Ilama-2	13B	2023-08-03	1200	-551.69	55.13
/8 70	Caldera AI/13B-Out000108 Caldera AI/13B-Thorns-12	nama Ilama	13B 13B	2023-07-20	1197	-1133.30	49.54 54 72
80	Cartinoe5930/Llama2 init Mistral	llama-2	7B	2023-09-00	1207	-528.82	60.98
81	Changgil/K2S3-SOLAR-11b-v1.0	llama	10B	2024-03-03	2290	-811.23	36.67
82	Changgil/k2s3_test_24001	llama-2	13B	2024-02-14	2346	-561.61	56.67

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
83	CobraMamba/mamba-gpt-3b-v4 (chiliu, 2023)	llama	3B	2023-09-05	1183	-605.83	41.24
84	ContextualAI/archangel_sft-kto_llama13b (Ethayarajh et al., 2023)	llama	13B	2023-12-03	1214	-542.49	52.87
85	Corianas/gpt-j-6B-Dolly	gptj	6B	2023-03-28	1173	-502.52	39.60
80	2024: Lee, 2024a; Lee and Ahn, 2024)	nama-3	<b>6</b> D	2024-03-22	2/12	-038.38	30.32
87	Dampish/Dante-2.8B	gpt_neox	2B	2023-05-09	1237	-671.96	-
88	Dampish/StellarX-4B-V0 (Black et al., 2022)	gpt_neox	4B	2023-05-27	1268	-722.64	37.31
89	Dampish/StellarX-4B-V0.2 (Black et al., 2022)	gpt_neox llama_2	4B 13B	2023-06-03	1245	-866.96	30.15 55.83
91	Danielbrdz/Barcenas-35	llama-2	3B	2023-11-15	1167	-554.05	41.74
92	Danielbrdz/Barcenas-7b	llama	7B	2023-08-25	1175	-603.96	50.87
93	Danielbrdz/Barcenas-Llama3-8b-ORPO	llama-3	8B	2024-04-29	12814	-557.50	72.50
94	Danielbrdz/CodeBarcenas-7b	llama-2	7B	2023-09-03	1174	-656.68	40.09
95 96	Deathsquad10/TinyLlama-Remix	llama	1B 1B	2023-11-26	1213	-812.12	34.00
97	Deathsquad10/TinyLlama-repeat	llama	1B	2024-01-06	1191	-622.61	37.09
98	Deci/DeciLM-7B-instruct (DeciAI Research Team, 2023)	deci	7B	2023-12-10	9769	-570.53	63.19
99 100	DeepMount00/Llama-3-8b-Ita	llama-3	8B 7P	2024-05-01	179401	-578.66	73.65
100	Delcos/Mistral-Pygmalion-7b	llama-2	7В	2023-11-08	1185	-579.30	51.02
102	Doctor-Shotgun/CalliopeDS-L2-13B (Yadav et al., 2023a)	llama-2	13B	2023-09-16	1428	-568.72	56.34
103	Doctor-Shotgun/CalliopeDS-v2-L2-13B	llama-2	13B	2023-09-28	1250	-572.65	57.12
104	DopeorNope/You_can_cry_Snowman-13B	llama	13B	2023-12-27	1203	-577.68	69.46
105	EleutherAl/gpt-j-bb (Gao et al., 2020; Wang, 2021; Wang and Komatsuzaki, 2021: Su et al. 2023a)	gptj	6B	2022-03-02	241435	-4/9.14	40.10
106	EleutherAI/gpt-neo-1.3B (Gao et al., 2020; Black et al., 2021)	gpt_neo	1B	2022-03-02	243332	-545.00	33.58
107	EleutherAI/gpt-neo-2.7B (Gao et al., 2020; Black et al., 2021)	gpt_neo	2B	2022-03-02	205503	-520.30	36.20
108	Eleuther AI/Ilemma_7b (Azerbayev et al., 2024)	llama-2	7B	2023-09-12	4671	-554.43	48.75
109	EleutherAl/polyglot-ko-12.8b (Kim et al., 2022; Ko et al., 2023; Su et al., 2023a)	gpt_neox	13B	2022-10-14	2941	-967.65	33.33
110	EleutherAI/pythia-1.4b (Gao et al., 2020; Biderman et al., 2022, 2023)	gpt_neox	1B	2023-02-09	22152	-531.68	34.75
111	EleutherAI/pythia-1.4b-deduped (Gao et al., 2020; Biderman et al., 2022,	gpt_neox	1B	2023-02-09	12730	-534.52	35.00
112	2023) Elauthor A Unythia 12b (Gao at al. 2020; Bidarman at al. 2022, 2023)	ant noor	120	2022 02 28	0102	175 12	20.07
112	EleutherAl/pythia-12b (Gao et al., 2020; Biderman et al., 2022, 2025) EleutherAl/pythia-12b-deduped (Gao et al., 2020; Biderman et al., 2022)	gpt_neox	12B	2023-02-28	9192	-477.81	38.82 39.70
110	2023)	Spi_neon	120	2020 02 27	,,,,,,		27110
114	EleutherAI/pythia-1b-deduped (Gao et al., 2020; Biderman et al., 2022, 2023)	gpt_neox	1B	2023-02-14	17983	-551.07	32.78
115	EleutherAl/pythia-2.8b-deduped (Gao et al., 2020; Biderman et al., 2022, 2022)	gpt_neox	2B	2023-02-10	11649	-507.42	36.72
116	EleutherAI/pythia-6.9b-deduped (Gao et al., 2020; Biderman et al., 2022,	gpt_neox	6B	2023-02-25	9824	-490.70	39.30
	2023)					~~~ ~ <b>~</b>	
117	Enno-Ai/ennodata-raw-pankajmathur-13b-pett	llama mixtrol	13B	2023-09-29	1215	-613.07	55.40
118	Expert68/Ilama2, 13b instructed version2	llama-2	13B 13B	2023-10-14	9078 1184	-561 35	71.75 55.41
120	FairMind/Llama-3-8B-4bit-UltraChat-Ita	llama-3	8B	2024-05-03	5014	-543.07	61.54
121	FairMind/Phi-3-mini-4k-instruct-bnb-4bit-Ita	mistral	4B	2024-05-02	2746	-634.71	66.61
122	FelixChao/CodeLlama13B-Finetune-v1	llama	13B	2023-09-13	1188	-647.40	47.19
125	FelixChao/Ilama2-13b-math1.1	llama-2	13B 13B	2023-08-12	1101	-602.88	54.18 54.19
125	FinancialSupport/saiga-7b	mistral	7B	2023-12-28	3902	-558.71	64.51
126	FlagAlpha/Llama2-Chinese-7b-Chat	llama-2	7B	2023-07-23	1383	-627.19	51.13
127	FlagAlpha/Llama3-Chinese-8B-Instruct	llama-3	8B	2024-04-23	1978	-557.17	63.50
128	Fredimensh/Guanaco-3B-Uncensored-V2 FreedomIntelligence/AceGPT-7B	gpt_neox	2B 7B	2023-08-27	2189	-032.40	38.98 49.47
129	FreedomIntelligence/phoenix-inst-chat-7b	bloom	7B 7B	2023-04-11	12962	-693.44	43.03
131	GeneZC/MiniChat-1.5-3B (Jain et al., 2023; Zhang et al., 2024a; Rafailov	llama	3B	2023-11-26	1601	-599.51	50.23
122	et al., 2024) Cons7C(MiniChot 2, 2P. (Join et al., 2022). There at al., 2024a, Bafailay et al.	11.0000	20	2022 12 27	4167	607 10	51.40
132	GenezC/MiniChat-2-3B (Jain et al., 2023; Zhang et al., 2024a; Rafallov et al., 2024)	llama	38	2023-12-27	4107	-007.10	51.49
133	GeneZC/MiniChat-3B (Zhang et al., 2024a)	llama	3B	2023-11-11	1633	-585.38	45.31
134	GeneZC/MiniMA-2-3B (Zhang et al., 2024a)	llama	3B	2023-12-27	1492	-526.37	44.75
135	GeneZC/MiniMA-3B (Zhang et al., 2024a) Gritt M/Gritt M 7P (Muonnichoff et al., 2024)	llama mistrol	3B 7P	2023-11-11	1494	-537.57	41.44
130	Gryphe/MythoLogic-L2-13b	llama	13B	2023-08-03	1198	-570.07	56.19
138	Gryphe/MythoMax-L2-13b	llama	13B	2023-08-10	14308	-583.80	56.00
139	Gryphe/MythoMix-L2-13b	llama	13B	2023-08-08	1174	-566.44	56.31
140	HWERI/Llama2-7b-sharegpt4	llama-2	7B	2023-08-04	1200	-660.22	51.05
141	HwEKI/pytnia-1.40-deduped-snaregpt Henry II/Instruct Mistral 7B v0.1 Dolly 15K	gpt_neox mistral	1B 7B	2023-08-10	1200	-557.95	55.11 60.45
142	HenryJJ/Instruct Yi-6B Dollv15K	llama	6B	2024-01-02	1170	-524.87	56.85
144	HenryJJ/Instruct_Yi-6B_Dolly_CodeAlpaca	llama	6B	2024-01-07	1177	-527.06	56.11
145	HiTZ/GoLLIE-7B (Sainz et al., 2024)	llama-2	7B	2023-09-25	1394	-632.52	37.48
146	HuggingFaceFW/ablation-model-fineweb-v1 (Lacoste et al., 2019)	Ilama mistral	1B 7D	2024-04-20	2238	-771.42	36.76 50 79
147	HuggingFaceH4/zephyr-7b-alpha (Ding et al., 2023; Tunstall et al., 2023; Cui	mistral	7B 7B	2023-10-20	12911	-559.73	59.78 59.50
	et al., 2024; Rafailov et al., 2024)						
149	HuggingFaceH4/zephyr-7b-beta (Ding et al., 2023; Tunstall et al., 2023; Cui	mistral	7B	2023-10-26	294019	-570.40	59.08
150	et al., 2024; Ratailov et al., 2024) IDEA-CCNL/Ziva-LLaMA-13B-Pretrain-v1 (IDEA-CCNL 2021; Zhang	llama-1	13B	2023-06-01	1171	-1397 30	29.96
150	et al., 2022a; Yang et al., 2022)		1.50	2022 00 01		1071.00	0
151	IDEA-CCNL/Ziya-LLaMA-13B-v1 (IDEA-CCNL, 2021; Zhang et al., 2022a;	llama-1	13B	2023-05-16	1244	-1397.34	29.82
152	Yang et al., 2022) INSAIT-Institute/BgGPT-7B-Instruct-v0.2	mistral	7 <b>P</b>	2024-03-03	3265	-594 56	63.08
152	IkariDev/Athena-v1	llama	13B	2023-08-30	1225	-566.12	54.11
154	IkariDev/Athena-v4	llama	13B	2023-10-07	1214	-555.80	57.23
155	Intel/neural-chat-7b-v3-1 (Mukherjee et al., 2023)	mistral	7B	2023-11-14	3976	-563.71	59.90

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
156	Intel/neural-chat-7b-v3-2 (Yu et al., 2024b)	mistral	7B	2023-11-21	2102	-554.92	68.29
157	Intel/neural-chat-7b-v3-3 (Yu et al., 2024b)	mistral	7B	2023-12-09	166226	-575.02	69.83
158	Jayant9928/orpo_med_v3	llama	8B	2024-05-01	2705	-557.53	62.21
159	Jiayi-Pan/Tiny-Vicuna-1B	llama	1B	2023-11-22	3094	-647.57	34.76
160	Josephgflowers/TinyLlama-3T-Cinder-v1.2	llama	1B 7D	2023-12-31	1413	-777.50	35.26
161	JosephusCheung/Qwen-LLaMAfied-/B-Chat	llama-2	/B 7D	2023-08-04	1221	-040.//	51.99
163	Kabster/Bio-Mistraly2-Squared	mistral	7B 7B	2023-08-30	2762	-603 34	43.00 57.73
164	Kabster/Bio/Mistral-Zephyr-Beta-SLERP	mistral	7B	2024-03-09	2758	-597.30	56.35
165	KnutJaegersberg/LLongMA-3b-LIMA	llama	3B	2023-09-03	2010	-617.15	38.51
166	KnutJaegersberg/MistralInstructLongish	mistral	7B	2023-11-15	1970	-543.74	53.62
167	KnutJaegersberg/Qwen-1_8B-Llamafied	llama	1B	2024-01-03	2839	-760.12	44.75
168	KnutJaegersberg/Walter-Falcon-1B	falcon	1B	2023-12-09	2014	-1010.10	34.07
169	KnutJaegersberg/deacon-13b	llama	13B	2023-09-20	2018	-533.27	53.63
170	KnutJaegersberg/deacon-3b	llama	3B	2023-09-18	2016	-621.85	39.05
171	KnutJaegersberg/falcon-1b-t-sft	falcon	1B 7D	2023-12-04	2029	-9/8.54	35.02
172	Kiluljaegersberg/webivilstrai- / D	misuai	/ D 6B	2023-11-17	1207	-554.02	35.97
173	KoboldAI/GPT-I-6B-Janeway (Gao et al. 2020: Wang and Komatsuzaki	gptj onti	6B	2022-03-02	4382	-496 37	39.54
17.	2021)	SPJ	02	2022 00 02	1002	190107	0,00
175	KoboldAI/GPT-J-6B-Shinen (Gao et al., 2020; Wang and Komatsuzaki, 2021)	gptj	6B	2022-03-02	1498	-498.02	39.60
176	KoboldAI/GPT-J-6B-Skein (Lacoste et al., 2019; Wang, 2021)	gptj	6B	2022-03-02	1236	-509.27	40.02
177	KoboldAI/LLaMA2-13B-Holomax	llama-2	13B	2023-08-14	1255	-563.47	54.52
178	KoboldAI/LLaMA2-13B-Tiefighter	llama-2	13B	2023-10-18	2321	-599.88	54.51
179	KoboldAI/OPT-13B-Erebus (Zhang et al., 2022b)	opt	13B	2022-09-09	6998	-652.81	39.61
180	KoboldAI/OPT-13B-Neryous-Wilx (Zhang et al., 2022b)	opt	13B 12D	2023-02-13	1057	-040.08	39.01
181	KoboldAI/OPT-27B-Frebus (Zhang et al. 2022b)	opt	13D 2B	2022-09-19	4239	-692.01	39.33
183	Kobold AI/OPT-2 7B-Nervbus-Mix	opt	2D 2B	2022-09-19	1464	-686 59	36.88
184	KoboldAI/OPT-6.7B-Erebus (Zhang et al., 2022b)	opt	6B	2022-09-15	5587	-641.11	39.09
185	KoboldAI/OPT-6.7B-Nerybus-Mix	opt	6B	2023-02-13	1585	-642.68	38.83
186	KoboldAI/OPT-6B-nerys-v2 (Zhang et al., 2022b)	opt	6B	2022-06-26	4974	-656.62	38.72
187	Korabbit/Llama-2-7b-chat-hf-afr-100step-flan	llama-2	7B	2023-11-30	1208	-650.87	52.88
188	Korabbit/Llama-2-7b-chat-hf-afr-100step-flan-v2	llama-2	7B	2023-12-03	1208	-650.96	52.92
189	Korabbit/Llama-2-7b-chat-hf-afr-100step-v2	llama-2	7B	2023-11-22	1205	-655.71	50.89
190	Korabbit/Llama-2-7b-chat-hf-afr-200step-flan	llama-2	7B	2023-11-30	1201	-647.72	52.62
191	Korabbit/Llama-2-/b-chat-hf-afr-200step-flan-v2	llama-2	7B 7D	2023-12-03	1211	-648.67	52.75
192	Korabbit/Llama-2-7b-chat-hf-afr-200step-merged	llama-2	/B 7D	2023-11-21	1222	-055.05	52.20
193	Korabbit/Llama-2-7b-chat-hf-afr-300step-v2	llama-2	7B 7B	2023-11-22	1203	-647.99	52 41
195	Korabbit/Llama-2-7b-chat-hf-afr-441sten-flan-v2	llama-2	7B	2023-12-03	1200	-647.85	52.28
196	Kukedlc/NeuralExperiment-7b-MagicCoder-v7.5	mistral	7B	2024-03-07	4044	-569.72	74.28
197	Kukedlc/NeuralLLaMa-3-8b-DT-v0.1	llama-3	8B	2024-05-11	5707	-571.65	72.52
198	Kukedlc/NeuralLLaMa-3-8b-ORPO-v0.3	llama-3	8B	2024-05-14	7985	-562.41	72.66
199	Kukedlc/NeuralSynthesis-7B-v0.1	mistral	7B	2024-04-06	8910	-596.75	76.80
200	Kukedlc/NeuralSynthesis-7B-v0.3	mistral	7B	2024-04-07	3114	-597.52	76.70
201	Kukedlc/NeuralSynthesis-7b-v0.4-slerp	mistral	7B	2024-04-12	3067	-597.78	76.76
202	LDCC/LDCC-SOLAR-10.7B (Kim et al., 2024b)	llama	10B	2024-01-03	3267	-549.97	71.40
203	LTC ALLabs/L2-7b Baluga WVG-Uncensored	llama	/B 7D	2023-09-23	1241	-330.11	52.04
204	LTC-AL-Labs/L2-7b-Bernes-Synthia	llama_2	7B 7B	2023-10-03	1222	-562.42	52.04
205	LTC-ALLabs/L2-7b-Hermes-WVG-Test	llama	7B	2023-09-27	1225	-570.23	51 35
200	LTC-AI-Labs/L2-7b-Synthia-WVG-Test	llama	7B	2023-09-28	1223	-576.22	51.25
208	Lazycuber/L2-7b-Base-Guanaco-Uncensored	llama	7B	2023-09-19	1172	-555.04	50.45
209	LeoLM/leo-hessianai-7b	llama	7B	2023-08-22	4502	-592.31	47.72
210	LeoLM/leo-hessianai-7b-chat	llama	7B	2023-09-10	3924	-735.02	49.29
211	LinkSoul/Chinese-Llama-2-7b	llama-2	7B	2023-07-20	40799	-584.42	52.59
212	Locutusque/Hercules-2.5-Mistral-7B	mistral	7B	2024-02-10	1953	-533.73	63.59
213	Locutusque/Hercules-3.1-Mistral-/B	mistral	7B	2024-02-19	2780	-528.16	62.09
214	Locutusque/Orca-2-13D-SF1-V4	nama Ilama	13B 12D	2023-11-25	2545	-01/.05	59.75 56.15
215	Locutusque/Orca 2 13b SET v5	llama	13D 13B	2023-12-22	2519	-074.55	56.77
210	MTSAIR/multi verse model	mistral	7B	2024-03-07	6357	-598.96	76 74
218	MayaPH/FinOPT-Franklin	opt	1B	2023-05-26	1207	-1380.66	29.78
219	MayaPH/opt-flan-iml-6.7b (Iyer et al., 2023)	opt	6B	2023-08-15	1193	-729.78	35.84
220	MaziyarPanahi/Llama-3-8B-Instruct-v0.4	llama-3	8B	2024-05-01	1407	-569.34	70.30
221	MaziyarPanahi/Llama-3-8B-Instruct-v0.8	llama-3	8B	2024-05-01	7281	-576.95	73.17
222	MaziyarPanahi/Llama-3-8B-Instruct-v0.9	llama-3	8B	2024-05-30	6241	-575.60	73.29
223	MaziyarPanahi/Mistral-7B-Instruct-v0.3	mistral	7B	2024-05-22	5230	-549.67	65.21
224	MaziyarPanahi/Mistral-7B-v0.3	mistral	7B	2024-05-22	5758	-531.44	60.40
225	Mihaiii/Matia 0.2	mistral	7B	2024-03-05	3065	-359.49	65.40
220	Miniman/Mini DPO test02	mistral	/B 7P	2023-12-10	3470	-507.04	61 23
228	MoaData/Myrrh solar 10.7b 3.0	llama	10R	2023-11-30	9855	-673 46	67.61
220	NExtNewChattingAl/shark tank ai 7 b	mistral	7R	2023-12-17	1224	-546.93	71 10
230	NExtNewChattingAI/shark tank ai 7b v2	mistral	7B	2023-12-23	1208	-618.14	66.54
231	NTQAI/Nxcode-CQ-7B-orpo (Hong et al., 2024)	qwen2	7B	2024-04-24	12700	-685.73	42.98
232	Neko-Institute-of-Science/metharme-7b	llama-1	6B	2023-04-30	1161	-562.36	47.48
233	Neko-Institute-of-Science/pygmalion-7b	llama-1	6B	2023-04-30	1188	-562.13	46.04
234	NekoPunchBBB/Llama-2-13b-hf_Open-Platypus-QLoRA-multigpu	llama-2	13B	2023-09-15	1208	-536.60	54.40
235	NeverSteep/Llama-3-Lumimaid-8B-v0.1	IIama-3	8B	2024-04-30	1499	-545.48	66.55
236	INeverSieep/Noromaid-/b-vU.2	mistral	7B	2023-12-21	1179	-552.34	01./8 50.27
237	NewstaR/Starlight-7B (Clark et al. 2018: Zellere et al. 2010: Hondersele	llama-7	/B 7R	2023-10-01	1180	-000.75	20.57 20.72
200	The sum of the second of the s	manna-2	7.0	2020-07-11	1100	272.07	12.13

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
239	Nexusflow/NexusRaven-V2-13B (Nexusflow.ai team, 2023; Rozière et al., 2024)	llama	13B	2023-12-04	3966	-593.55	48.21
240 241	Nexusflow/Starling-LM-7B-beta (Ziegler et al., 2020; Zhu et al., 2023a) NickyNicky/Mistral-7B-OpenOrca-oasst_top1_2023-08-25-v2 (Xiao et al.,	mistral mistral	7B 7B	2024-03-19 2023-10-11	4847 1236	-559.22 -560.11	69.88 61.65
242	2024) NickyMistral-7B-OpenOrca-oasst_top1_2023-08-25-v3 (Dao et al., 2022; Xiao et al. 2024)	mistral	7B	2023-10-13	1227	-568.73	61.26
243	Norquinal/Mistral-7B-claude-instruct	mistral	7B	2023-09-28	1248	-534.14	59.27
244	Norquinal/Ilama-2-7b-claude-chat	llama-2	7B	2023-08-11	1221	-557.34	50.98
245	Norquinal/llama-2-7b-claude-chat-rp	llama-2	7B	2023-08-14	1218	-557.09	51.25
246	NousResearch/CodeLlama-13b-hf	llama	13B	2023-08-24	6299	-582.02	43.35
247	NousResearch/CodeLlama-/b-ni NousPasaarah/Harmas 2 Pro Llama 2 8P ("Taknium" at al. 2024b)	llama	/B 9D	2023-08-24	9428	-597.05	39.81 69.72
248	NousResearch/Hermes-2-Pro-Mistral-7B ("interstellarninia" et al. 2024)	mistral	6B 7B	2024-04-30	14586	-620.90	67.35
250	NousResearch/Hermes-2-Theta-Llama-3-8B ("Teknium" et al., 2024a)	llama-3	8B	2024-05-05	12392	-565.47	69.21
251	NousResearch/Meta-Llama-3-8B-Instruct (AI@Meta, 2024)	llama-3	8B	2024-04-18	104388	-573.95	67.10
252	NousResearch/Nous-Hermes-13b	llama-1	13B	2023-06-03	1536	-605.76	54.04
253	NousResearch/Nous-Hermes-2-Mistral-7B-DPO ("Teknium" et al., 2024c)	mistral	7B	2024-02-18	8725	-565.26	68.10
254	NousResearch/Nous-Hermes-2-SOLAR-10.7B	llama	10B	2024-01-01	9918	-542.37	71.00
255	NousResearch/Nous-Hermes-Llama2-13b	llama-2	13B	2023-07-20	39412	-586.34	55.75
256	NousResearch/Nous-Hermes-Ilama-2-7b	Ilama-2	6B 7D	2023-07-25	12019	-569.93	51.87
257	NousResearch/Yarn-Mistral-7b-128k (Peng et al., 2023b)	mistral	/B 7B	2023-10-31	20727	-532.22	59.42 59.63
259	OFvortex/FMO-2B	gemma	2B	2023-10-31	4137	-976 17	44.26
260	Open-Orca/LlongOrca-13B-16k (Mukherjee et al., 2023; Wang et al., 2023a;	llama-2	13B	2023-08-16	1217	-551.07	56.59
261	Dale et al., 2023; Longpre et al., 2023; Touvron et al., 2023b) Open-Orca/LlongOrca-7B-16k (Mukherjee et al., 2023; Wang et al., 2023a;	llama-2	7B	2023-08-05	1223	-583.06	53.02
262	Longpre et al., 2023; Touvron et al., 2023b; Lian et al., 2023c) Open-Orca/Mistral-7B-OpenOrca (Mukherjee et al., 2023; Longpre et al., 2022; Ling et al., 2023; d)	mistral	7B	2023-09-29	18070	-575.75	60.17
263	Open-Orca/Mistral-7B-SlimOrca (Mukherjee et al., 2023; Lian et al., 2023f,e; Longre et al., 2023)	mistral	7B	2023-10-08	3696	-569.42	60.37
264	Open-Orca/OpenOrca-Platypus2-13B (Hu et al., 2022; Mukherjee et al., 2023; Wang et al., 2023a,b; Lee et al., 2023; Longpre et al., 2023; Touvron et al., 2024).	llama	13B	2023-08-11	6853	-557.72	57.28
265	Open-Orca/OpenOrca-Preview1-13B (Mukherjee et al., 2023; Lian et al., 2023a; Longpre et al., 2023; Touvron et al., 2023a)	llama-1	13B	2023-07-12	1240	-679.87	51.38
266	Open-Orca/OpenOrcaxOpenChat-Preview2-13B (Mukherjee et al., 2023; Wang et al., 2023a,b; Longpre et al., 2023; Touvron et al., 2023b)	llama-2	13B	2023-07-31	2172	-547.85	56.70
267	OpenAssistant/codellama-13b-oasst-stt-v10	llama-2	13B	2023-08-26	2122	-591.56	44.85
268	OpenAssistant/liama2-13b-orca-8k-5319 (Mukherjee et al., 2023)	llama-2	13B 12P	2023-07-24	1270	-531.8/	55.09 40.77
209	OpenAssistant/oasst-sit-1-pytilia-120 OpenAssistant/oasst-sit-4-pytilia-12b-enoch-3 5	gpt_neox	12B	2023-03-09	458718	-013.23	40.77
270	OpenAssistant/outsis sit 1 pythia 125 open 5.5	gpt_neox	12B	2023-05-07	1323	-524.99	42.21
272	OpenAssistant/stablelm-7b-sft-v7-epoch-3	gpt_neox	7B	2023-04-20	1215	-764.93	34.85
273	OpenBuddy/openbuddy-atom-13b-v9-bf16	llama	13B	2023-08-05	1206	-625.44	52.31
274	OpenBuddy/openbuddy-llama2-13b-v11-bf16	llama-2	13B	2023-08-23	1199	-606.17	52.93
275	OpenBuddy/openbuddy-llama2-13b-v11.1-bf16	llama-2	13B	2023-08-24	1210	-595.93	55.28
276	OpenBuddy/openbuddy-llama2-13b-v8.1-tp16	llama-2	13B	2023-07-25	7339	-587.40	57.76
277	OpenBuddy/openbuddy-mistral-7b-y13	mistral	6D 7B	2024-04-20	1335	-578.14	53 50
278	OpenBuddy/openbuddy-mistral-7b-v13-base	mistral	7B 7B	2023-10-11	1197	-677.88	51.99
280	OpenBuddy/openbuddy-mistral-7b-v13.1	mistral	7B	2023-10-11	1217	-638.68	52.62
281	OpenBuddy/openbuddy-mistral2-7b-v20.3-32k	mistral	7B	2024-03-27	2302	-642.78	62.73
282	OpenBuddy/openbuddy-openllama-13b-v7-fp16	llama-1	13B	2023-07-03	1209	-654.65	49.31
283	OpenBuddy/openbuddy-openllama-3b-v10-bf16	llama	3B	2023-08-10	1206	-794.41	36.87
284	OpenBuddy/openbuddy-openllama-7b-v12-bf16	llama	7B	2023-09-19	3153	-854.52	45.28
285	OpenBuddy/openbuddy-zephyr-7b-v14.1	mistral	7B	2023-11-06	3392	-657.03	51.86
280	by the second start germina - 1, -7, -10, (1) so that a second start, 2015; Mi- haylov et al., 2018; Clark et al., 2019; Talmor et al., 2019; Sap et al., 2019; Bisk et al., 2019; Chollet, 2019; Sakaguchi et al., 2019; Zellers et al., 2019; Austin et al., 2021; Cobbe et al., 2021; Hendrycks et al., 2021a; Chen et al.,	gemma	ŏБ	2024-04-08	5128	-1555.00	39.78
	2021; Parrish et al., 2022; Srivastava et al., 2023; Zhong et al., 2023; Gemini						
787	Ieam et al., 2024) OpenPine/mistral-ft-optimized-1218	mistral	78	2023-12 17	1405	-546 73	71.04
287	OpenPipe/mistral-ft-optimized-1227	mistral	7B 7B	2023-12-27	6273	-564.28	70.54
289	Orenguteng/Llama-3-8B-Lexi-Uncensored	llama-3	8B	2024-04-23	9653	-564.93	66.18
290	PathFinderKR/Waktaverse-Llama-3-KO-8B-Instruct (AI@Meta, 2024; AI@Waktaverse, 2024)	llama-3	8B	2024-04-19	2305	-552.92	66.77
291	Pirr/pythia-13b-deduped-green_devil	gpt_neox	13B	2023-02-09	1351	-505.96	40.31
292	PistachioAlt/Synatra-MCS-7B-v0.3-RP-Slerp	mistral	7B	2023-12-11	1162	-549.18	69.18
293	PracticeLLM/SOLAK-tall-10./B-Merge-v1.0	ilama	10B 1B	2023-12-20	1235	-530.70	71.08
294	PysmalionAl/mythalion-13b	llama-2	13R	2023-09-02	2383	-552.46	56.48
296	PygmalionAI/pygmalion-1.3b	gpt_neox	1B	2022-12-25	1519	-963.99	31.14
297	PygmalionAI/pygmalion-2-13b	llama-2	13B	2023-09-04	2083	-540.27	55.12
298	PygmalionAI/pygmalion-2-7b	llama-2	6B	2023-09-04	2063	-559.12	51.11
299	PygmalionAI/pygmalion-2.7b	gpt_neo	2B	2023-01-05	1986	-659.52	33.98
300	PygmalionAl/pygmalion-6b	gptj	6B	2023-01-07	4312	-555.36	38.47
301	Q-bert/Dumblebee-/B	mistral	7B 7D	2023-12-03	1226	-549.82	60.00
302	Q-bert/Terminis-7B	mistral	/B 7R	2023-12-03	1237	-550.18	70 73
304	Owen/CodeOwen1.5-7B-Chat (Bai et al., 2023)	awen2	7B	2024-04-15	77169	-691.05	43.26
305	Qwen/Qwen1.5-1.8B (Bai et al., 2023)	qwen2	1B	2024-01-22	137256	-610.58	46.55
306	Qwen/Qwen1.5-1.8B-Chat (Bai et al., 2023)	qwen2	1B	2024-01-30	10856	-673.92	43.99
307	Qwen/Qwen1.5-4B (Bai et al., 2023)	qwen2	3B	2024-01-22	6431	-553.86	57.05

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
308	Qwen/Qwen1.5-4B-Chat (Bai et al., 2023)	qwen2	3B	2024-01-30	5525	-601.85	46.79
309	Qwen/Qwen1.5-7B (Bai et al., 2023)	qwen2	7B	2024-01-22	122768	-533.48	61.76
310	Qwen/Qwen1.5-7B-Chat (Bai et al., 2023)	qwen2	7B	2024-01-30	26143	-608.59	55.15
311	Qwen/Qwen2-1.5B (Yang et al. 2024a)	qwen2	1B 7B	2024-05-31	46510	-581.07	55.80 68.40
313	RWKV/rwkv-raven-1b5	rwky	1B	2023-05-04	1570	-526.30	33.56
314	RoversX/llama-2-7b-hf-small-shards-Samantha-V1-SFT	llama-2	7B	2023-08-11	1164	-558.46	49.96
315	RubielLabarta/LogoS-7Bx2-MoE-13B-v0.2	mixtral	12B	2024-01-21	3055	-589.76	77.15
316	S4sch/zephyr-neural-chat-frankenmerge11b	mistral	11B	2023-11-28	1181	-618.27	58.57
317	SJ-Donald/SJ-SOLAR-10. / D-DPO SJ-Donald/SOLAR-10. 7B-slerp	llama	10B 10B	2024-01-25	2306	-533.85	72.67
319	SJ-Donald/llama3-passthrough-chat	llama-3	11B	2024-05-17	2241	-607.16	60.15
320	Salesforce/codegen-6B-multi (Nijkamp et al., 2023)	codegen	6B	2022-04-13	1651	-733.14	32.43
321	Salesforce/codegen-6B-nl (Nijkamp et al., 2023)	codegen	6B	2022-04-13	1176	-478.78	40.00
322	Sanji Watsuki/Kunoichi-DPO-y2-7B	mistral	7B 7B	2024-01-04	1230	-579.55	72.13
323	SanjiWatsuki/Kulotell-Di O-v2-7B SanjiWatsuki/Loyal-Macaroni-Maid-7B	mistral	7B	2023-12-24	1210	-575.31	71.68
325	SanjiWatsuki/Silicon-Maid-7B	mistral	7B	2023-12-27	1578	-580.90	70.31
326	SanjiWatsuki/Sonya-7B	mistral	7B	2023-12-31	5399	-594.55	68.48
327	et al. 2023: Asai et al. 2023)	mistral	7B	2023-10-27	1214	-664.82	56.46
328	SealLLM-7B-v2 (Kojima et al., 2023; Zheng et al., 2023; Zhang et al., 2023; Nguyen et al., 2024)	mistral	7B	2024-01-29	6294	-548.56	67.57
329	SeaLLMs/SeaLLM-7B-v2.5 (Zhang et al., 2023; Nguyen et al., 2024)	gemma	8B	2024-04-03	13928	-565.35	69.07
330	Severian/ANIMA-Phi-Neptune-Mistral-7B	mistral	7B	2023-10-11	1191	-631.73	55.54
332	SuperAGI/SAM	nama mistral	13B 7B	2023-09-02	1205	-555.28	54.51 59.30
333	TIGER-Lab/MAmmoTH2-7B-Plus (Yue et al., 2024b)	mistral	7B	2024-05-06	10155	-673.56	67.75
334	TIGER-Lab/MAmmoTH2-8B-Plus (Yue et al., 2024b)	llama	8B	2024-05-06	12819	-595.71	67.49
335	TIGER-Lab/TIGERScore-13B (Jiang et al., 2024)	llama	13B	2023-11-26	1429	-548.56	56.79
336	TaylorAI/Flash-Llama-13B	llama	13B	2023-08-19	1181	-530.82	53.67
338	TaylorAl/Flash-Llama-7B	llama	эв 7В	2023-08-13	1177	-578.15	40.13
339	TeeZee/Bielik-SOLAR-LIKE-10.7B-Instruct-v0.1	mistral	10B	2024-04-10	1566	-854.79	53.50
340	TehVenom/Dolly_Malion-6b	gptj	6B	2023-03-27	1201	-486.77	39.77
341	TehVenom/Dolly_Shygmalion-6b	gptj	6B	2023-03-29	1195	-491.76	39.89
342	Teh Venom/GPT L Dvg, PPO 6B	gptj	6B	2023-05-23	1197	-487.65	40.11
343	TehVenom/GPT-J-Pvg PPO-6B-Dev-V8p4	gptj	6B	2023-03-05	1203	-493.17	39.61
345	TehVenom/Metharme-13b-Merged	llama-1	13B	2023-05-18	1199	-561.32	54.15
346	TehVenom/Moderator-Chan_GPT-JT-6b	gptj	6B	2023-03-19	1185	-502.43	42.17
347	Teh Venom/PPO_Pygway-V 8p4_Dev-6b (Wang and Komatsuzaki, 2021)	gptj	6B	2023-03-17	1191	-489.49	39.85
348	TehVenom/PPO Shygmalion-V8p4 Dev-6b	gptj	6B	2023-03-23	1199	-490.13	39.85
350	TehVenom/Pygmalion-13b-Merged	llama-1	13B	2023-05-18	1207	-588.35	48.49
351	TehVenom/Pygmalion-Vicuna-1.1-7b	llama-1	6B	2023-05-02	1246	-572.52	49.25
352	TehVenom/Pygmalion_AlpacaLora-7b	llama-1	7B	2023-04-30	1195	-605.84	46.49
353	TencentARC/LLaMA-Pro-8B	llama-2	8B	2024-03-17	1319	-557.58	51.67
355	TencentARC/LLaMA-Pro-8B-Instruct	llama-2	8B	2024-01-06	1423	-634.71	58.06
356	TheBloke/CodeLlama-13B-Instruct-fp16	llama-2	13B	2023-08-24	2193	-579.86	45.82
357	TheBloke/CodeLlama-13B-Python-fp16	llama-2	13B	2023-08-24	2114	-587.68	37.52
358	TheBloke/Liama-2-13B-Tp10 TheBloke/Nous-Hermes-13B-SuperHOT-8K-fp16	llama-2	13B 13B	2023-07-18	0470 1225	-530.82	53.07 52.18
360	TheBloke/Planner-7B-fp16	llama-1	7B	2023-06-05	1219	-562.04	45.65
361	TheBloke/UltraLM-13B-fp16 (Ding et al., 2023)	llama-1	13B	2023-06-29	1211	-567.38	54.62
362	TheBloke/Vicuna-13B-CoT-fp16 (Lacoste et al., 2019)	llama-1	13B	2023-06-08	1219	-646.61	53.28
363	TheBloke/Wizard-Vicuna-13B-Uncensored-HF	llama-1	13B 7P	2023-05-13	1648	-579.62	54.14
365	TheBloke/WizardLM-13B-V1-1-SuperHOT-8K-fp16 (Xu et al., 2023a)	llama-1	13B	2023-07-07	1970	-626.40	53.16
366	TheBloke/airoboros-13B-HF	llama-1	13B	2023-05-23	1214	-581.17	54.05
367	TheBloke/airoboros-7b-gpt4-fp16	llama-1	7B	2023-06-04	1208	-603.39	47.70
368	TheBloke/gpt4-alpaca-lora-13B-HF	llama-1	13B	2023-04-17	1181	-547.42	53.98
309	TheBloke/guanaco-15B-HF	llama-1	13B 7B	2023-05-25	1222	-586.92	55.54 47 34
371	TheBloke/koala-13B-HF	llama-1	13B	2023-04-07	2501	-594.23	51.16
372	TheBloke/koala-7B-HF	llama-1	7B	2023-04-07	1238	-617.95	44.29
373	TheBloke/stable-vicuna-13B-HF (von Werra et al., 2023; Anand et al., 2023; Chiang et al., 2023; Taori et al., 2023; Touvron et al., 2023a)	llama-1	13B	2023-04-28	1337	-590.21	51.64
374	TheBloke/tulu-13B-fp16 (Köpf et al., 2023; Conover et al., 2023; Wang et al., 2023c; Peng et al., 2023a; Longpre et al., 2023; Chaudhary, 2023; Touvron	llama-1	13B	2023-06-10	1231	-583.06	53.58
375	et al., 2023a) TheBloke/tulu-7B-fp16 (Köpf et al., 2023; Conover et al., 2023; Wang et al., 2023c; Peng et al., 2023a; Longpre et al., 2023; Chaudhary, 2023; Touvron et al. 2023a)	llama-1	7B	2023-06-10	4079	-616.51	50.24
376	TheBloke/wizard-vicuna-13B-HF	llama-1	13B	2023-05-04	1219	-610.98	52.75
377	TheTravellingEngineer/llama2-7b-chat-hf-dpo	llama-2	7B	2023-08-14	1201	-659.39	50.38
378	The TravellingEngineer/Ilama2-7b-chat-hf-guanaco	llama-2	6B	2023-08-02	1209	-597.96	50.02
380	The TravellingEngineer/llama2-7b-chat-hf-v3	llama-2	6B	2023-08-08	1208	-549.87	48.81
381	TheTravellingEngineer/llama2-7b-chat-hf-v4	llama-2	6B	2023-08-10	1213	-549.87	49.78
382	TheTravellingEngineer/Ilama2-7b-hf-guanaco	llama-2	6B	2023-07-25	1206	-554.41	50.12
383	TigerResearch/tigerbot-7b-base	Ilama Ilama	7B	2023-08-19	1228	-587.17	47.93
385	TinyLiama/TinyLiama-1.1B-Chat-v0.0	llama	1B 1R	2023-11-20	13743	-042.34 -619 71	34.94 37.17
200	,			12 00			

ID	Model Name	Model Type	Size	Date	DLs	$\overline{\ell}_i$	Task
386	TinyLlama/TinyLlama-1.1B-intermediate-step-1195k-token-2.5T	llama	1B	2023-12-11	1365	-613.87	36.26
387	TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-3T	llama	1B	2023-12-28	798206	-608.46	36.42
388	TinyLlama/TinyLlama-1.1B-intermediate-step-955k-token-2T TomGra/EusionNet	llama	1B 10B	2023-11-19	8119	-646.34	34.56
390	TomGrc/FusionNet_linear	llama	10B	2023-12-31	1167	-561.26	74.43
391	Toten5/Marcoroni-neural-chat-7B-v2	mistral	7B	2023-12-12	1201	-557.34	72.50
392	TsinghuaC3I/Llama-3-8B-UltraMedical (Zhang et al., 2024c)	llama-3	8B	2024-04-27	3949	-558.79	63.73
393 394	Unbabel/TowerInstruct-7B-v0.1 (Alves et al., 2024)	llama	6B	2024-01-03	2230	-582.00	49.11 52.39
395	Undi95/MLewd-Chat-v2-13B	llama	13B	2023-09-26	1187	-569.88	57.23
396	Undi95/MLewd-L2-13B	llama	13B	2023-09-04	1172	-666.07	53.12
397	Undi95/MLewd-L2-Chat-13B Undi95/MLewd-y2 4_13B	llama	13B 13B	2023-09-16	1177	-551.09	57.75
399	Undi95/MLewdBoros-L2-13B	llama	13B	2023-09-20	1200	-542.84	56.51
400	Undi95/Meta-Llama-3-8B-hf (AI@Meta, 2024)	llama-3	8B	2024-04-18	11376	-514.33	62.35
401	Undi95/Mistral-11B-v0.1	mistral	10B	2023-10-09	1196	-556.96	58.05
402	Undi95/Nous-Hermes-13B-Code Undi95/OpenRP-13B	llama	13B 13B	2023-09-02	1207	-602.47	55.95 56.57
404	Undi95/ReMM-SLERP-L2-13B	llama	13B	2023-09-04	1511	-585.33	56.03
405	Undi95/ReMM-v2-L2-13B	llama	13B	2023-09-09	1207	-565.06	56.99
406	Undi95/ReMM-v2.1-L2-13B Undi95/ReMM-v2.2.1.2.13B	llama	13B 13B	2023-09-12	1217	-564.83	56.71 57.10
407	Undi95/UndiMix-v1-13b	llama	13B 13B	2023-09-21	1220	-649.19	55.50
409	Undi95/UndiMix-v4-13B	llama	13B	2023-09-12	1212	-569.06	56.93
410	Undi95/Unholy-v1-12L-13B	llama	13B	2023-09-10	1210	-541.04	57.47
411 412	Undi95/X-MythoChronos-13B VAGOsolutions/Llama-3-SauerkrautLM-8b-Instruct	llama-3	13B 8B	2023-11-18	1203 55400	-604.69	58.43 73.74
413	VAGOsolutions/SauerkrautLM-7b-HerO	mistral	7B	2023-11-24	1226	-553.88	64.49
414	VAGOsolutions/SauerkrautLM-Gemma-7b	gemma	8B	2024-02-27	5691	-561.27	67.83
415	VAGOsolutions/SauerkrautLM-SOLAR-Instruct	llama llama 1	10B 7D	2023-12-20	1175	-560.76	74.21
410	VMware/open-llama-7b-open-instruct	llama-1	7Б 7В	2023-05-51	6470	-642.10	41.11
418	Voicelab/trurl-2-13b-academic	llama-2	13B	2023-09-18	2774	-568.46	53.94
419	Voicelab/trurl-2-7b	llama-2	7B	2023-08-16	2861	-602.32	50.58
420	Weyaxi/ChatAY1-Lora-Assamble-Marcoroni Weyaxi/Dolphin2.1-OpenOrca-7B	llama mistral	13B 7B	2023-09-14	1203	-569.45	57.76
422	Weyaxi/Instruct-v0.2-Seraph-7B	mistral	7B	2023-10-11	1203	-574.63	68.48
423	Weyaxi/Luban-Marcoroni-13B-v2	llama	13B	2023-09-13	1212	-566.43	57.92
424	Weyaxi/Luban-Marcoroni-13B-v3	llama	13B	2023-09-13	1214	-566.44	57.94
425	weyaxi/MetaMath-Chupacabra-7B-V2.01-Sterp Weyaxi/MetaMath-NeuralHermes-2 5-Mistral-7B-Linear	mistral	7B 7B	2023-12-08	1210	-546.54	70.26 67.60
427	Weyaxi/MetaMath-NeuralHermes-2.5-Mistral-7B-Ties	mistral	7B	2023-12-05	1203	-604.36	67.03
428	Weyaxi/MetaMath-OpenHermes-2.5-neural-chat-v3-3-Slerp	mistral	7B	2023-12-10	1224	-548.15	69.92
429	Weyaxi/MetaMath-Tulpar-7b-v2-Slerp	mistral	7B 7P	2023-12-08	1217	-555.63	70.07
430	Weyaxi/MetaMath-neural-chat-7b-v3-2-Sterp	mistral	7В	2023-12-08	1208	-590.97	67.54
432	Weyaxi/MetaMath-una-cybertron-v2-bf16-Ties	mistral	7B	2023-12-06	1218	-602.74	68.88
433	Weyaxi/OpenHermes-2.5-neural-chat-7b-v3-1-7B	mistral	7B	2023-11-24	1230	-567.74	67.84
434	Weyaxi/OpenOrca_Zenbyr-7B	mistral	7B 7B	2023-12-03	1221	-576.02	68.71 64.97
436	Weyaxi/Samantha-Nebula-7B	mistral	7B	2023-10-11	1166	-584.51	54.58
437	Weyaxi/SauerkrautLM-UNA-SOLAR-Instruct	llama	10B	2023-12-21	1225	-561.98	74.26
438	Weyaxi/Seraph-7B	mistral	7B	2023-12-11	1214	-546.73	71.86
439	Weyaxi/Seraph-openchat-5.5-1210-Sterp Weyaxi/SlimOnenOrca-Mistral-7B	mistral	7B 7B	2023-12-27	1210	-567.06	70.89 60.84
441	Weyaxi/neural-chat-7b-v3-1-OpenHermes-2.5-7B	mistral	7B	2023-12-01	1205	-570.29	67.19
442	Weyaxi/openchat-3.5-1210-Seraph-Slerp	mistral	7B	2023-12-27	1210	-553.81	71.82
443	WhiteRabbitNeo/WhiteRabbitNeo-13B-v1	llama-2	13B	2023-12-17	2090	-615.73	49.11
444	Xwin-LM/Xwin-LM-7B-V0.1 (Xwin-LM Team, 2023)	llama-2	7B	2023-09-15	1401	-592.93	52.08
446	Yhyu13/LMCocktail-10.7B-v1 (Xiao et al., 2023)	llama-2	10B	2023-12-20	3332	-556.38	74.06
447	Yhyu13/chimera-inst-chat-13b-hf	llama-1	13B	2023-05-11	1295	-582.39	52.86
448	Yukang/Llama-2-/b-longlora-32k-ft (Chen et al., 2024b) Yukang/Long Alpaca-13B (Chen et al., 2023b, 2024b)	llama-2 llama	7B 13B	2023-09-12	2423	-1387.05	29.20 41.74
450	Yukang/LongAlpaca-7B (Chen et al., 2023b, 2024b)	llama	6B	2023-10-03	2773	-794.46	39.36
451	aaditya/Llama3-OpenBioLLM-8B (Singhal et al., 2022; Nori et al., 2023;	llama-3	8B	2024-04-20	9308	-590.14	54.06
452	Singhal et al., 2023; Pal and Sankarasubbu, 2024a,b; Rafailov et al., 2024)	mistral	70	2024 01 00	2251	696 19	62.82
452	abacusai/Giraffe-13b-32k-v3	llama-2	13B	2024-01-09	1214	-546.12	57.24
454	abacusai/Llama-3-Smaug-8B (Pal et al., 2024)	llama-3	8B	2024-04-19	12873	-565.82	64.61
455	abacusai/Slerp-CM-mist-dpo	mistral	7B	2024-01-03	4030	-553.48	73.10
456 457	abacusai/bigstrai-12b-32K abhinand/Llama-3-OnenBioMed-8B-slern-v0-3	mistral llama-3	12B 8B	2024-03-06	5216 2713	-632.41 -558 51	62.17 62.08
458	abhinand/tamil-llama-7b-base-v0.1 (Balachandran, 2023)	llama-2	7B	2023-11-08	1382	-877.00	44.52
459	abhinand/tamil-llama-7b-instruct-v0.1 (Balachandran, 2023)	llama-2	7B	2023-11-08	3507	-708.97	45.52
460	abhishekchohan/Yi-9B-Forest-DPO-v1.0	llama	9B	2024-03-18	2755	-528.97	64.11
461	aomsnekchonan/mistrai-/B-forest-apo acrastt/Bean-3B	inistrat llama	/B 3R	2024-01-21 2023-09-02	1961	-301.51 -592 51	03.28 40.18
463	acrastt/Griffin-3B	llama	3B	2023-08-18	1198	-584.63	41.13
464	acrastt/Marx-3B	llama	3B	2023-08-15	2006	-603.32	41.71
465	acrastt/Marx-3B-V2 acrastt/OmegLLaMA_3B	llama llama	3B 2P	2023-08-22	1234	-609.79	42.08
467	acrastt/Puma-3B	llama	3B 3B	2023-08-20	1198	-584.54	41.02
468	acrastt/RedPajama-INCITE-Chat-Instruct-3B-V1	gpt_neox	2B	2023-07-27	1202	-565.14	39.23
469	adamo1139/Mistral-7B-AEZAKMI-v1	mistral	7B	2023-11-27	1199	-596.56	54.92

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
470	adonlee/LLaMA_2_13B_SFT_v0	llama	13B	2023-10-03	1268	-556.59	57.31
471	adonlee/LLaMA_2_13B_SFT_v1	llama	13B	2023-11-06	1265	-546.75	63.04
472	aerdincdal/CBDDO-LLM-8B-Instruct-v1	llama	8B	2024-05-02	4114	-613.44	56.94
473	ahnyeonchan/OpenOrca-AYT-13B	llama-2	13B	2023-09-07	1184	-569.10	54.91
474	ajibawa-2023/Pothon-Code-13B	llama	13B	2023-12-08	1201	-582.71	53.61
476	ajibawa-2023/SlimOrca-13B	llama	13B	2023-11-27	1177	-608.08	60.39
477	ajibawa-2023/Uncensored-Frank-13B	llama	13B	2023-09-14	1192	-577.46	55.64
478	ajibawa-2023/Uncensored-Frank-7B	llama	7B	2023-09-14	1178	-654.45	47.90
479	ajibawa-2023/Uncensored-Jordan-13B ajibawa-2023/Uncensored-Jordan-7B	llama	13B 7B	2023-10-23	1174 1174	-5/8.48	56.27 49.95
481	ajibawa-2023/carl-7b	llama	7B	2023-07-22	1212	-599.09	46.16
482	akjindal53244/Mistral-7B-v0.1-Open-Platypus	mistral	7B	2023-10-05	1286	-543.69	58.92
483	alignment-handbook/zephyr-7b-sft-full	mistral	7B	2023-11-09	11605	-569.88	57.52
484	allbyai/ToRoLaMa-/b-v1.0 (Do et al., 2023) allanai/OLMo 1B hf (Touvron et al., 2023a; Groeneveld et al., 2024)	llama-2	7B 1B	2023-12-19	1185	-/32.57	47.87
486	allenai/OLMo-7B-hf (Touvron et al., 2023a; Groeneveld et al., 2024)	olmo	6B	2024-04-12	5407	-554.02	43.36
487	allknowingroger/MultiverseEx26-7B-slerp	mistral	7B	2024-03-30	3034	-599.87	76.80
488	alnrg2arg/blockchainlabs_7B_merged_test2_4	mistral	7B	2024-01-17	1455	-577.39	75.28
489	alorg2arg/blockchainlabs_7B_merged_test2_4_prune	mistral	7B 7D	2024-01-18	1939	-673.45	57.91
490	aloobun/falcon-1b-cot-t2	falcon	7B 1B	2023-12-20	2161	-735 39	28 56
492	aloobun/open-llama-3b-v2-elmv3	llama	3B	2023-12-08	1207	-604.18	41.14
493	andreaskoepf/llama2-13b-megacode2_min100	llama-2	13B	2023-08-14	1162	-571.36	56.92
494	argilla/notus-7b-v1	mistral	7B 7D	2023-11-16	7660	-564.14	60.22
495	athirdpath/NSFW_DPO_Noromaid-/b augmynt/shisa-7b-y1 (Jain et al. 2023: Rafailoy et al. 2024)	mistral	/B 7B	2023-12-12	1280	-546.50	61.59 55.01
497	augmxnt/shisa-76-v1 (sain et al., 2023, Karanov et al., 2024)	mistral	7B	2023-11-27	1195	-648.62	51.64
498	augmxnt/shisa-gamma-7b-v1	mistral	7B	2023-12-23	151426	-590.25	55.50
499	automerger/YamshadowExperiment28-7B	mistral	7B	2024-03-18	3152	-595.65	76.86
500	beaugogh/Llama2-7b-openorca-mc-v1	llama-2	7B 7D	2023-08-20	1177	-613.56	52.24
502	beaugogh/Llama2-7b-sharegnt4	llama-2	7В	2023-10-00	11/5	-660.74	51.05
502	beaugogh/pythia-1.4b-deduped-sharegpt	gpt_neox	1B	2023-07-25	1293	-557.93	35.11
504	beomi/KoAlpaca-KoRWKV-6B	rwkv	6B	2023-06-02	2288	-1099.15	28.57
505	beomi/KoAlpaca-Polyglot-5.8B	gpt_neox	6B	2023-03-16	4198	-1309.35	29.46
506 507	beomi/KoRWKV-6B beomi/Yi-Ko-6B (Lee 2024b)	rwkv llama	6B 6B	2023-05-26	2127 4261	-1148.78	28.19
508	beomi/lama-2-ko-7b (Lee, 2024)	llama-2	6B	2023-07-20	5226	-907.04	45.32
509	beowolx/CodeNinja-1.0-OpenChat-7B	mistral	7B	2023-12-20	6300	-555.93	67.40
510	berkeley-nest/Starling-LM-7B-alpha (Zhu et al., 2023b,a)	mistral	7B	2023-11-25	18366	-571.63	67.05
511	bhavinjawade/SOLAR-10B-OrcaDPO-Jawade	llama	10B	2024-01-06	1224	-563.19	74.27
512	higcode/gpt_bigcode-samacoder higcode/starcoder2-3h (Beltagy et al. 2020: Bayarian et al. 2022: Dao et al.	starcoder2	3B	2023-04-00	40198	-620.70	20.49
010	2022; Ainslie et al., 2023; Lozhkov et al., 2024)	stareoderz	02	2020 11 22	110010	020170	07.20
514	bigcode/starcoder2-7b (Beltagy et al., 2020; Bavarian et al., 2022; Dao et al.,	starcoder2	7B	2024-02-20	11834	-584.03	42.95
515	2022; Ainshe et al., 2023; Lozhkov et al., 2024) bigcode/starcoderbase-1b (Shazeer, 2019; Dao et al., 2022; Bavarian et al.,	gpt bigcode	1B	2023-07-03	5247	-734.99	30.06
	2022; Li et al., 2023a)	818					
516	bigcode/starcoderbase-7b (Shazeer, 2019; Dao et al., 2022; Bavarian et al.,	gpt_bigcode	7B	2023-07-26	3566	-652.23	33.75
517	2022; Li et al., 2023a) higscience/bloom-1b1 (Shoeyhi et al. 2020: Press et al. 2022: Dettmers et al.	bloom	1B	2022-05-19	11054	-687 97	32.47
517	2022)	bioom	15	2022 05 17	11001	007.97	52.17
518	bigscience/bloom-1b7 (Shoeybi et al., 2020; Press et al., 2022; Dettmers et al.,	bloom	1B	2022-05-19	40840	-657.05	33.98
519	2022) higscience/bloom-3h (Shoeyhi et al. 2020: Press et al. 2022: Dettmers et al.	bloom	3B	2022-05-19	13014	-633 79	36.07
517	2022)	biobili	50	2022-05-17	15714	-055.17	50.07
520	bigscience/bloom-7b1 (Shoeybi et al., 2020; Press et al., 2022; Dettmers et al.,	bloom	7B	2022-05-19	21428	-604.20	39.18
521	2022) higsgience/bloomz 3b (Muennighoff et al. 2023)	bloom	3B	2022 10.08	7075	701 42	37.03
521	bigscience/bloomz-7b1 (Muennighoff et al., 2023)	bloom	5B 7B	2022-09-27	12152	-652.34	42.21
523	bigscience/bloomz-7b1-mt (Muennighoff et al., 2023)	bloom	7B	2022-09-28	2427	-653.08	42.14
524	bofenghuang/vigogne-2-13b-instruct	llama-2	13B	2023-07-26	1203	-534.88	55.14
525	bofenghuang/vigogne-2-7b-chat	llama-2	7B 7D	2023-07-29	1170	-558.60	52.45
526 527	bofenghuang/vigogne-2-7b-instruct	llama-2	7В 7В	2023-07-20	1243	-555.11	52.02 47.76
528	bofenghuang/vigostral-7b-chat	mistral	7B	2023-09-29	4118	-535.15	59.18
529	budecosystem/boomer-1b	llama	1B	2023-10-03	1267	-846.13	28.44
530	castorini/rank_vicuna_7b_v1_fp16 (Touvron et al., 2023b; Pradeep et al.,	llama-2	7B	2023-09-27	1165	-720.69	44.36
531	2023) ceadar-ie/FinanceConnect-13B (CeADAR 2023)	llama	13B	2023-11-28	1270	-661 62	49 34
532	chargoddard/storytime-13b	llama-2	13B	2023-09-22	1166	-599.45	56.64
533	chickencaesar/llama2-platypus-llama2-chat-13B-hf	llama-2	13B	2023-09-28	1216	-536.92	54.11
534	chinoll/Yi-6b-200k-dpo	llama	6B	2023-12-01	1187	-573.79	51.93
535 536	circulus/Llama-2-13b-orca-v1	llama-2	13B 7B	2023-08-01	1228	-540.82	57.05 53.56
537	clibrain/Llama-2-7b-ft-instruct-es	llama-2	7B	2023-08-09	2062	-559.81	49.63
538	codellama/CodeLlama-13b-Instruct-hf (Rozière et al., 2024)	llama-2	13B	2023-08-24	25248	-579.86	45.82
539	codellama/CodeLlama-13b-Python-hf (Rozière et al., 2024)	llama-2	13B	2023-08-24	2537	-587.68	37.00
540	codellama/CodeLlama-13b-ht (Rozière et al., 2024)	llama-2	13B	2023-08-24	8432	-582.07	43.35
542	codellama/CodeLlama-7b-Python-hf (Rozière et al., 2024)	llama-2	6B	2023-08-24	5138	-605.60	+0.05 36.42
543	codellama/CodeLlama-7b-hf (Rozière et al., 2024)	llama-2	6B	2023-08-24	66040	-597.05	39.81
544	cognitivecomputations/Llama-3-8B-Instruct-abliterated-v2	llama-3	8B	2024-05-09	8728	-581.93	66.00
545	cognitivecomputations/TinyDolphin-2.8-1.1b	llama mistrol	1B	2024-01-21	1678	-683.06	36.34
546	cognitivecomputations/ westLake-/B-V2-laser	mistrai	/B	2024-01-26	3980	-391.19	/4./8

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
547	cognitivecomputations/dolphin-2.2.1-mistral-7b	mistral	7B	2023-10-30	7292	-546.22	65.01
548	cognitivecomputations/dolphin-2.6-mistral-7b	mistral	7B	2023-12-27	1418	-559.17	64.92
549	cognitivecomputations/dolphin-2.6-mistral-7b-dpo	mistral	7B	2023-12-31	1236	-565.34	67.20
550	cognitivecomputations/dolphin-2.0-mistrai-/b-dpo-laser (Snarma et al., 2023)	llama-3	7B 8B	2024-01-01	130977	-558.72	65.92
552	cognitivecomputations/dolphin-2.9.1-llama-3-8b	llama-3	8B	2024-04-20	7575	-647.01	66.23
553	cognitivecomputations/dolphin-2.9.1-yi-1.5-9b	llama	8B	2024-05-18	4880	-598.79	68.92
554	cognitivecomputations/openchat-3.5-0106-laser	mistral	7B	2024-01-27	6155	-554.46	69.46
555	cookinai/BruinHermes	mistral	7B 7D	2023-12-17	1193	-551.19	73.42
557	crumb/apricot-wildflower-20	mistral	7В	2023-12-31	1179	-532.22	72.08 59.74
558	cyberagent/open-calm-7b (Andonian et al., 2021)	gpt_neox	7B	2023-05-15	24168	-1002.50	28.21
559	cypienai/cymist-2-v02-SFT (Lacoste et al., 2019)	mistral	7B	2024-05-22	2717	-537.73	62.57
560	cypienai/cymist2-v01-SFT (Lacoste et al., 2019)	mistral	7B	2024-05-12	2754	-628.53	51.71
562	danielpark/gorani-100k-flama2-13b-instruct databricks/dolly-y2-12b (Conover et al. 2023)	nama-2	13B 12B	2023-10-04	3860	-1410.79	29.69 39.46
563	databricks/dolly-v2-3b (Conover et al., 2023)	gpt_neox	3B	2023-04-13	21443	-556.97	-
564	databricks/dolly-v2-7b (Conover et al., 2023)	gpt_neox	7B	2023-04-13	10069	-555.60	39.24
565	deepseek-ai/DeepSeek-Coder-V2-Lite-Base (Dai et al., 2024)	deepseek	16B	2024-06-14	13642	-527.73	-
566 567	deepseek-ai/DeepSeek-Coder-V2-Lite-Instruct (Dai et al., 2024) deepseek-ai/DeepSeek-Prover-V1 (Xin et al., 2024a)	deepseek	16B 7B	2024-06-14	142758	-569.70	_
568	deepseek-ai/DeepSeek-Prover-V1.5-Base (Xin et al., 2024a)	deepseek	7B 7B	2024-08-15	230	-554.13	_
569	deepseek-ai/DeepSeek-Prover-V1.5-RL (Xin et al., 2024b)	deepseek	7B	2024-08-15	12223	-672.91	_
570	deepseek-ai/DeepSeek-Prover-V1.5-SFT (Xin et al., 2024b)	deepseek	7B	2024-08-15	6459	-670.56	-
571	deepseek-ai/DeepSeek-R1-Distill-Llama-8B (DeepSeek-AI et al., 2025)	deepseek	8B	2025-01-20	199313	-693.69	-
572	deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B (DeepSeek-AI et al., 2025)	deepseek	2B 14B	2025-01-20	290094	-854.97	-
574	deepseek-ai/DeepSeek-R1-Distill-Owen-7B (DeepSeek-AI et al., 2025)	deepseek	7B	2025-01-20	201118	-758.13	_
575	deepseek-ai/DeepSeek-V2-Lite (DeepSeek-AI et al., 2024b)	deepseek	16B	2024-05-15	25474	-533.54	_
576	deepseek-ai/DeepSeek-V2-Lite-Chat (DeepSeek-AI et al., 2024b)	deepseek	16B	2024-05-15	18056	-588.92	-
577	deepseek-ai/ESFT-vanilla-lite (Wang et al., 2024c)	deepseek	16B	2024-07-04	270	-550.46	-
578 579	deepseek-ai/deepseek-coder-1.3b-base (Guo et al., 2024)	deepseek	1B 1B	2023-10-28	48960	-/18.82	32 40
580	deepseek-ai/deepseek-coder-6.7b-base (Guo et al., 2024)	deepseek	6B	2023-10-23	38348	-647.51	32.40 40.87
581	deepseek-ai/deepseek-coder-6.7b-instruct (Guo et al., 2024)	deepseek	6B	2023-10-29	28828	-686.56	43.57
582	deepseek-ai/deepseek-coder-7b-base-v1.5 (Guo et al., 2024)	deepseek	7B	2024-01-25	895	-570.98	-
583	deepseek-ai/deepseek-coder-7b-instruct-v1.5 (Guo et al., 2024)	deepseek	6B 7D	2024-01-25	28928	-604.47	50.89
584 585	deepseek-ai/deepseek-lim-7b-chat (DeepSeek-AI et al., 2024a) deepseek-ai/deepseek-lim-7b-chat (DeepSeek-AI et al., 2024a)	deepseek	7B 7B	2023-11-29	34271	-552.59	59 38
586	deepseek-ai/deepseek-math-7b-base (Shao et al., 2024)	deepseek	7B	2024-02-05	37614	-566.93	57.61
587	deepseek-ai/deepseek-math-7b-instruct (Shao et al., 2024)	deepseek	7B	2024-02-05	8352	-593.77	51.48
588	deepseek-ai/deepseek-math-7b-rl (Shao et al., 2024)	deepseek	6B	2024-02-05	1633	-607.25	49.54
589	deepseek-ai/deepseek-moe-16b-base (Dai et al., 2024)	deepseek	16B	2024-01-08	13924	-544.49	-
590 591	deepseek-al/deepseek-moe-10b-chat (Dal et al., 2024) dfurman/Llama-3-8B-Orno-v0 1	llama-3	10B 8B	2024-01-09	8852 5146	-5/8.25	64 67
592	dhmeltzer/Llama-2-13b-hf-ds_eli5_1024_r_64_alpha_16_merged	llama-2	13B	2023-09-14	1188	-542.35	54.16
593	dhmeltzer/Llama-2-13b-hf-ds_wiki_1024_full_r_64_alpha_16_merged	llama-2	13B	2023-09-14	1198	-536.93	52.94
594	dhmeltzer/Llama-2-13b-hf-eli5-wiki-1024_r_64_alpha_16_merged	llama-2	13B	2023-09-14	1192	-538.52	53.57
595 596	dhmeltzer/llama-/b-SFT_eli5_wiki65k_1024_r_64_alpha_16_merged	llama	6B 13B	2023-08-25	1257	-568.38	50.00 50.47
590	digitous/13B-Chimera	llama-1	13B	2023-05-23	1359	-562.36	54.92
598	digitous/Alpacino13b	llama-1	13B	2023-04-13	1182	-547.35	52.39
599	dotvignesh/perry-7b	llama	7B	2023-09-28	1171	-618.90	49.55
600	dreamgen/WizardLM-2-7B (Luo et al., 2023; Xu et al., 2023a; Luo et al., 2025)	mistral	7B	2024-04-16	2533	-641.93	63.75
601	dvruette/oasst-pythia-12b-flash-attn-5000-steps	gpt neox	12B	2023-03-12	1175	-602.37	40.73
602	dvruette/oasst-pythia-12b-pretrained-sft	gpt_neox	12B	2023-04-03	1170	-552.58	41.48
603	dvruette/oasst-pythia-12b-reference	gpt_neox	12B	2023-04-03	1174	-557.88	40.33
604	eachadea/vicuna-13b-1.1	llama-1	13B	2023-04-13	1336	-646.61	53.29
605	eldoghbhed/Peagle-9h	mistral	7 D 8 B	2023-04-13	10180	-015.88	73 30
607	elinas/chronos-13b-v2	llama-2	13B	2023-08-02	1944	-595.38	55.25
608	elliotthwang/elliott_Llama-2-7b-hf	llama-2	6B	2023-10-09	1181	-553.66	50.20
609	elyza/ELYZA-japanese-Llama-2-13b (Touvron et al., 2023b; Sasaki et al.,	llama-2	13B	2023-12-25	1200	-587.25	56.14
610	2023b) elyza/ELYZA_iapapese_Llama_2_13b_instruct (Touvron et al. 2023b; Sasaki	llama_2	13B	2023-12-25	1710	-594.98	54 72
010	et al., 2023b)	nama-2	150	2025-12-25	1710	-574.70	54.72
611	elyza/ELYZA-japanese-Llama-2-7b (Touvron et al., 2023b; Sasaki et al.,	llama-2	7B	2023-08-28	2356	-630.66	48.70
612	2023a) alvza/ELVZA japanese Llama 2 7h fast (Touvron et al. 2023h: Sasaki et al.	llama 2	7B	2023 08 28	1766	640.71	17 67
012	2023a)	nama-2	/ <b>D</b>	2023-08-28	1700	-040.71	47.07
613	elyza/ELYZA-japanese-Llama-2-7b-fast-instruct (Touvron et al., 2023b;	llama-2	7B	2023-08-28	2387	-629.72	49.15
(14	Sasaki et al., 2023a)	11	70	2022 08 28	(20)	(25.45	40.79
014	et al 2023a)	nama-2	/B	2023-08-28	6296	-025.45	49.78
615	ericzzz/falcon-rw-1b-chat	falcon	1B	2023-12-05	1319	-722.18	37.37
616	ericzzz/falcon-rw-1b-instruct-openorca	falcon	1B	2023-11-24	1585	-752.98	37.63
617	euclaise/Ferret_7B	mistral	7B	2023-10-28	1172	-634.24	53.87
018 610	euclaise/falcon_lb_stage1	falcon	1B 1R	2023-09-15	2119 4016	-734.28	37.23 37.59
620	facebook/opt-1.3b (Brown et al., 2020: Zhang et al., 2022b)	opt	1B 1B	2022-05-11	139758	-700.64	34.60
621	facebook/opt-13b (Brown et al., 2020; Zhang et al., 2022b)	opt	13B	2022-05-11	19130	-625.38	40.06
622	facebook/opt-2.7b (Brown et al., 2020; Zhang et al., 2022b)	opt	2B	2022-05-11	61450	-673.03	36.74
623	tacebook/opt-6.7b (Brown et al., 2020; Zhang et al., 2022b)	opt	6B	2022-05-11	44909	-641.13	39.08
625	facebook/xglm-1.7B (Lin et al., 2022b)	xglm	1B	2022-03-02	2388	-723.91	31.42

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
626	facebook/xglm-4.5B (Lin et al., 2022b)	xglm	5B	2022-03-02	1436	-669.92	34.31
627	facebook/xglm-7.5B (Lin et al., 2022b)	xglm	7B	2022-03-02	4654	-663.87	36.38
628	failspy/Meta-Llama-3-8B-Instruct-abliterated-v3	llama-3	8B	2024-05-20	9975	-572.33	67.27
629	failspy/Phi-3-medium-4k-instruct-abliterated-v3	phi3	13B	2024-05-22	5430	-569.05	70.12
630	IDIgit/LUNA-SOLARKrautLM-Instruct	llama	10B	2023-12-22	1180	-582.98	73.79
632	fblgit/una-cybertron-7b-y2-bf16 (Murias 2023)	mistral	10B 7B	2024-01-02	1440	-568 27	74.07 69.67
633	feidfoe/Metamath-reproduce-7b	llama-2	7B	2023-11-24	1198	-724.03	55.81
634	fireballoon/baichuan-vicuna-chinese-7b	llama-1	7B	2023-06-18	1210	-702.08	46.06
635	freewheelin/free-llama3-dpo-v0.2 (Kim et al., 2024b)	llama-3	8B	2024-05-09	3896	-499.17	62.69
636	gagan3012/MetaModelv2	llama	10B	2024-01-03	1203	-560.16	74.24
637	gagan3012/MetaModelv3	llama	10B	2024-01-05	1207	-561.14	74.39
038	et al. 2024)	nama	130	2023-08-03	5/10	-355.92	54.69
639	garage-bAInd/Platypus2-7B (Hu et al., 2022; Touvron et al., 2023b; Lee et al.,	llama	6B	2023-08-22	6198	-554.08	49.97
640	2024)	llama 2	7B	2023 07 20	2602	568 34	10.67
641	ghost-x/ghost-7h-alnha	mistral	7B	2023-07-20	4854	-662 59	57.65
642	glaiveai/glaive-coder-7b	llama-2	7B	2023-09-17	1222	-692.31	41.56
643	google/codegemma-2b (CodeGemma Team et al., 2024)	gemma	2B	2024-03-21	8798	-793.92	32.19
644	google/codegemma-7b (CodeGemma Team et al., 2024)	gemma	8B	2024-03-21	3779	-556.79	56.73
645	google/codegemma-7b-it (CodeGemma Team et al., 2024)	gemma	8B	2024-03-21	11227	-1040.98	58.28
646	google/gemma-1.1-7b-it (Joshi et al., 2017; Zhao et al., 2018; Mihaylov et al., 2019; Clark et al., 2010; Talmor et al., 2010; San et al., 2010; Pick et al., 2010; San et al.	gemma	8B	2024-03-26	19835	-1355.60	60.09
	2016; Clark et al., 2019; Tallior et al., 2019; Sap et al., 2019; Bisk et al., 2019; Chollet 2019: Sakaguchi et al. 2019; Zellers et al. 2019; Austin et al. 2021;						
	Cobbe et al., 2021: Hendrycks et al., 2021a: Chen et al., 2021: Parrish et al.,						
	2022; Srivastava et al., 2023; Zhong et al., 2023; Gemini Team et al., 2024)						
647	google/gemma-2b (Joshi et al., 2017; Mihaylov et al., 2018; Rudinger et al.,	gemma	2B	2024-02-08	256798	-581.35	46.51
	2018; Zhao et al., 2018; Sap et al., 2019; Bisk et al., 2019; Zellers et al., 2019;						
	Clark et al., 2019; Talmor et al., 2019; Chollet, 2019; Sakaguchi et al., 2019; Cohmon et al. 2020; Cohho et al. 2021; Hendrycke et al. 2021a; Dhemale						
	et al. 2021: Chen et al. 2021: Austin et al. 2021: Hartyigsen et al. 2022:						
	Parrish et al., 2022; Lin et al., 2022a; Srivastava et al., 2023; Zhong et al.,						
	2023; Gemini Team et al., 2024)						
648	google/gemma-2b-it (Joshi et al., 2017; Mihaylov et al., 2018; Rudinger et al.,	gemma	2B	2024-02-08	103851	-886.70	42.75
	2018; Zhao et al., 2018; Sap et al., 2019; Bisk et al., 2019; Zellers et al., 2019; Clade et al., 2010; Telmer et al., 2010; Challet, 2010; Schermehi et al., 2010;						
	Gehman et al. 2019; Talmor et al. 2019; Chonel, 2019; Sakaguchi et al., 2019; Gehman et al. 2020; Cobbe et al. 2021; Hendrycks et al. 2021a; Dhamala						
	et al., 2021: Chen et al., 2021: Austin et al., 2021: Hartyigsen et al., 2022:						
	Parrish et al., 2022; Lin et al., 2022a; Srivastava et al., 2023; Zhong et al.,						
~	2023; Gemini Team et al., 2024)						
649	google/gemma-7b (Joshi et al., 2017; Mihaylov et al., 2018; Rudinger et al.,	gemma	8B	2024-02-08	72047	-539.12	63.75
	2016; Znao et al., 2018; Sap et al., 2019; Bisk et al., 2019; Zellers et al., 2019; Clark et al. 2010: Talmor et al. 2010: Challet 2010: Sakaguchi et al. 2019;						
	Gehman et al., 2020; Cobbe et al., 2017; Chohet, 2017; Sakaguen et al., 2019; Gehman et al., 2020; Cobbe et al., 2021; Hendrycks et al., 2021a; Dhamala						
	et al., 2021; Chen et al., 2021; Austin et al., 2021; Hartvigsen et al., 2022;						
	Parrish et al., 2022; Lin et al., 2022a; Srivastava et al., 2023; Zhong et al.,						
(50	2023; Dettmers et al., 2023; Gemini Team et al., 2024)		on	2024 02 12	((77)	1227 42	52 56
650	2018; Zhao et al. 2019; Sap et al. 2010; Bick et al. 2010; Zellers et al. 2019;	gemma	ðВ	2024-02-13	00//0	-1327.43	55.50
	Clark et al., 2019; Talmor et al., 2019; Chollet, 2019; Sakaguchi et al., 2019;						
	Gehman et al., 2020; Cobbe et al., 2021; Hendrycks et al., 2021a; Dhamala						
	et al., 2021; Chen et al., 2021; Austin et al., 2021; Hartvigsen et al., 2022;						
	Parrish et al., 2022; Lin et al., 2022a; Srivastava et al., 2023; Zhong et al.,						
651	2023; Gemini Team et al., 2024) google/recurrent comma 2b it (Joshi et al. 2017: Mihaylov et al. 2018:	recurrent commo	2B	2024 04 08	4046	1033.00	40.86
051	Budinger et al. 2018: Zhao et al. 2018: San et al. 2019: Bisk et al. 2019:	recurrent_gennna	20	2024-04-08	4040	-1055.09	40.80
	Zellers et al., 2019; Clark et al., 2019; Talmor et al., 2019; Chollet, 2019;						
	Sakaguchi et al., 2019; Gehman et al., 2020; Cobbe et al., 2021; Hendrycks						
	et al., 2021a; Dhamala et al., 2021; Chen et al., 2021; Hendrycks et al., 2021b;						
	Austin et al., 2021; Hartvigsen et al., 2022; Parrish et al., 2022; Lin et al., 2022; Chinet al., 2022						
	2022a; Srivastava et al., $2023$ ; Zhong et al., $2023$ ; De et al., $2024$ ; Griffin Team et al. $2024$ )						
652	gradientai/Llama-3-8B-Instruct-262k (Peng et al., 2023b; Ding et al., 2023;	llama-3	8B	2024-04-25	12536	-555.17	60.26
	Pekelis et al., 2024; AI@Meta, 2024; Liu et al., 2025)						
653	gradientai/Llama-3-8B-Instruct-Gradient-1048k (Peng et al., 2023b; Ding	llama-3	8B	2024-04-29	6854	-564.98	59.84
654	guardrail/llama-2-7b-guanaco-instruct-sharded	llama-2	6B	2023-07-21	1320	-647.09	50.58
655	gywy/llama2-13b-chinese-v2	llama-2	13B	2023-08-22	1163	-719.40	49.58
656	h2oai/h2ogpt-gm-oasst1-en-1024-12b	gpt_neox	12B	2023-05-02	1185	-495.39	40.65
657	h2oai/h2ogpt-gm-oasst1-en-2048-open-llama-7b-preview-300bt	llama-1	7B	2023-05-04	1185	-1042.98	34.32
658	h2oai/h2ogpt-gm-oasst1-en-2048-open-Ilama-7b-preview-300bt-v2	llama-1	7B	2023-05-10	1190	-746.41	37.55
039 660	n20an/n20gpt-0ass(1-312-12D h20ai/h20gnt-0ig-0ass(1-356-6-9h	gpt_neox	12B 0R	2023-04-17	1319	-484.68 _492.04	40.48 38.62
661	h2oai/h2ogpt-oig-oasst1-20-0_90	gpt_neox	9D 9R	2023-04-17	1749	-492.94	38.52
662	habanoz/TinyLlama-1.1B-intermediate-step-715k-1.5T-lr-5-1epch-	llama	1B	2023-11-22	1279	-658.65	34.98
672	airoboros3.1-1k-instruct-V1	llama	110	2022 11 20	1000	610 00	25 15
005	nabano2/11hyLiana-1.1B-intermediate-step-/15k-1.51-in-5-2.2epochs-	nama	ID	2023-11-20	1265	-046.66	55.45
664	habanoz/TinyLlama-1.1B-intermediate-step-715k-1.5T-lr-5-3epochs-oasst1-	llama	1B	2023-11-21	1277	-651.34	35.42
	top1-instruct-V1			2022 11 1			25.5-
665	habanoz/TinyLlama-1.1B-intermediate-step-715k-1.5T-lr-5-4epochs-oasst1- top1 instruct V1	llama	1B	2023-11-21	1288	-646.55	35.28
666	habanoz/tinvllama-oasst1-top1-instruct-full-lr1-5-v0 1	llama	1B	2023-11-19	1178	-679 96	35.58
667	hakurei/instruct-12b	gpt_neox	12B	2023-04-09	1177	-602.56	38.63
668	hakurei/mommygpt-3B	llama	3B	2023-11-12	1177	-598.84	41.36
669	haoranxu/ALMA-13B (Xu et al., 2024b,a)	llama	13B	2023-09-17	2012	-563.34	50.16

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
670	haoranxu/ALMA-13B-Pretrain (Xu et al., 2024b,a)	llama	13B	2023-09-17	5510	-558.03	51.68
671	haoranxu/ALMA-7B (Xu et al., 2024b,a)	llama	7B	2023-09-17	1292	-599.26	45.32
672	harborwater/open-llama-3b-claude-30k	llama	3B	2023-11-21	1200	-603.17	40.93
673	health 360/Healix-1.1B-V1-Chat-dDPO	llama	1B 13B	2023-11-05	2790	-796.10	33.00 52.61
675	heegyu/LIMA2-13b-hf (Touvron et al., 2023b)	llama-2	13B	2023-08-07	3330	-596.33	52.98
676	heegyu/LIMA2-7b-hf (Touvron et al., 2023b)	llama-2	7B	2023-08-04	3463	-651.31	49.27
677	heegyu/RedTulu-Uncensored-3B-0719	gpt_neox	3B	2023-07-23	1187	-703.96	39.19
678	heegyu/WizardVicuna-3B-0719	llama	3B	2023-07-23	3356	-655.31	39.48
680	heegyu/WizardVicuna-Uncensored-3B-0/19	llama	3B 3B	2023-07-23	0364	-038.40	39.73
681	heegyu/WizardVicuna2-13b-hf (Touvron et al., 2023b)	llama-2	13B	2023-08-07	6481	-612.15	51.05
682	hfl/chinese-alpaca-2-13b	llama	13B	2023-08-14	1308	-592.75	57.41
683	hfl/chinese-llama-2-1.3b	llama-2	1B	2023-10-08	2269	-1132.44	28.59
684 685	htl/chinese-llama-2-13b	llama-2	13B 7P	2023-08-11	1176	-665.69	52.00
686	hiyouga/Baichuan2-7B-Chat-LLaMAfied	llama-2	7B 7B	2023-09-08	1218	-609.32	40.99 51.42
687	hoskinson-center/proofGPT-v0.1-6.7B	gpt_neox	6B	2023-02-04	1187	-974.98	29.72
688	hpcai-tech/Colossal-LLaMA-2-7b-base (Li et al., 2023b; Touvron et al.,	llama-2	7B	2023-09-18	1192	-690.99	51.39
680	2023b; Dao, 2023) huggyllama/llama 13b	llama 1	13B	2023 04 03	6345	541.00	51 33
690	huggyllama/llama-150	llama-1	6B	2023-04-03	192383	-562.04	46 37
691	hyunseoki/ko-en-llama2-13b	llama-2	13B	2023-10-02	3351	-550.60	51.27
692	hyunseoki/ko-ref-llama2-13b	llama-2	13B	2023-10-04	3364	-775.26	43.62
693	hyunseoki/ko-ref-llama2-7b	llama-2	7B	2023-10-04	3312	-851.38	40.75
694	hywu/Camelidae-8x13B (Houlsby et al., 2019; Komatsuzaki et al., 2023; Dettmers et al. $2023$ ; Wu et al. $2024$ )	camelidae	13B	2024-01-10	1886	-535.46	59.40
695	hywu/Camelidae-8x7B (Houlsby et al., 2019; Komatsuzaki et al., 2023;	camelidae	7B	2024-01-10	1901	-565.23	54.47
	Dettmers et al., 2023; Wu et al., 2024)						
696	ibivibiv/llama-3-nectar-dpo-8B (Lacoste et al., 2019)	llama-3	8B 7D	2024-05-14	6085	-576.43	67.92
697 698	ibndias/NeuralHermes-MoE-2x7B	mistral	7В 12В	2024-03-02	11156	-541.12	64.00 64.08
699	ibranze/araproje-llama2-7b-hf	llama-2	7B	2023-10-06	1161	-549.87	49.73
700	ignos/LeoScorpius-GreenNode-Alpaca-7B-v1	mistral	7B	2023-12-15	1206	-558.75	74.74
701	ignos/LeoScorpius-GreenNode-Platypus-7B-v1	mistral	7B	2023-12-15	1188	-540.87	68.96
702	ignos/Mistral-T5-7B-v1	mistral	7B 2D	2023-12-18	1266	-557.24	72.47
703	invalid-coder/Sakura-SOLAR-Instruct-CarbonVillain-en-10 7B-v2-slern	llama	10B	2023-03-07	12810	-722.30	57.58 74.45
705	itsliupeng/llama2_7b_code	llama-2	7B	2023-09-28	1236	-538.58	49.05
706	itsliupeng/openllama-7b-base	llama	7B	2023-12-08	1200	-536.60	47.09
707	itsliupeng/openllama-7b-icl (Shi et al., 2024)	llama	7B	2023-12-08	1186	-535.25	47.93
708	jae24/openhermes_dpo_norobot_0201	mistral	7B 7B	2024-01-02	1161	-562.15	63.78 74.80
710	ieonsworld/CarbonVillain-en-10.7B-v1	llama	10B	2023-12-14	1192	-561.74	74.28
711	jeonsworld/CarbonVillain-en-10.7B-v4	llama	10B	2023-12-30	12898	-561.44	74.52
712	jerryjalapeno/nart-100k-7b	llama-1	7B	2023-07-14	1194	-585.01	46.39
713	jingyeom/freeze_KoSoLAR-10.7B-v0.2_1.4_dedup	llama	10B	2024-01-29	2287	-594.17	60.06
714	johnsnowlabs/BioLing-/B-Dare	mistrai	/B 8B	2024-04-08	2082	-591.07	61.93
716	jondurbin/airoboros-gpt-3.5-turbo-100k-7b	llama-1	7B	2023-05-12	1473	-614.18	47.05
717	jondurbin/airoboros-12-13b-2.1	llama-2	13B	2023-08-28	2322	-565.56	53.34
718	jondurbin/airoboros-12-13b-gpt4-2.0	llama	13B	2023-07-27	1395	-575.42	52.49
719	jondurbin/airoboros-12-13b-gpt4-m2.0	llama	13B	2023-07-28	1378	-614.41	52.66
720	iphme/Llama-2-13b-chat-german (Touvron et al., 2023b)	llama-2	13B	2023-07-21	1265	-572.28	55.07
722	jphme/em_german_leo_mistral	mistral	7B	2023-10-07	1900	-625.44	51.69
723	jphme/orca_mini_v2_ger_7b (Conover et al., 2023; Mathur, 2023a; Su et al.,	llama-1	7B	2023-07-04	1181	-603.44	47.65
	2023b; Harries, 2023; Taori et al., 2023; Touvron et al., 2023a; Xu et al.,						
724	2023a) junelee/wizard-vicuna-13b	llama-1	13B	2023-05-03	2196	-610.98	52.73
725	kaist-ai/mistral-orpo-capybara-7k (Hong et al., 2024)	mistral	7B	2024-03-23	4821	-540.67	63.36
726	kekmodel/StopCarbon-10.7B-v5	llama	10B	2023-12-30	13910	-560.12	74.41
727	kevin009/IlamaRAGdrama	mistral	7B 7D	2024-02-04	4248	-610.94	74.65
728	kings/Liama-2-ko-70-Chat (100Vron et al., 2023b)	llama-2	/B 13B	2023-07-25	5450 1174	-707.03	40.27
730	kingbri/chronolima-airo-grad-12-13B	llama-2	13B	2023-08-04	1182	-582.73	55.50
731	klyang/MentaLLaMA-chat-7B (Yang et al., 2024b)	llama	7B	2023-09-26	2634	-640.26	51.17
732	kodonho/Solar-OrcaDPO-Solar-Instruct-SLERP	llama	10B	2024-01-12	3266	-560.31	74.35
733	kodonho/SolarM-SakuraSolar-SLERP	llama	10B	2024-01-12	3254	-562.65	74.29
735	kyujinpy/Sakura-SOLAR-Instruct-DPO-v2	llama	10B	2023-12-24	3298	-560.08	74.40
736	kyujinpy/Sakura-SOLRCA-Math-Instruct-DPO-v1	llama	10B	2023-12-25	3284	-562.45	74.13
737	l3utterfly/llama2-7b-layla	llama-2	7B	2023-08-07	1182	-561.12	52.05
738	13utterfly/minima-3b-layla-v1	llama-2	3B	2023-12-12	1178	-548.62	43.21
739	13utterfly/minima-30-tayla-v2 13utterfly/mistral-7h-v() 1-lavla-v1	nama-2 mistral	3B 7R	2023-12-19	1176	-547.03	43.39 57 56
740	13utterfly/mistral-7b-v0.1-layla-v1	mistral	7В	2023-10-31	1199	-567.42	57.60
742	13utterfly/open-llama-3b-v2-layla	llama	3B	2023-08-18	1172	-832.18	40.25
743	lcw99/llama-3-10b-it-ko-2024-0527	llama-3	9B	2024-05-27	2236	-564.03	63.70
744	lcw99/Ilama-3-10b-it-kor-extented-chang	llama-3	9B	2024-05-15	2230	-558.58	54.76
745 746	ICw99/Ilama-3-100-II-Kor-extented-chang-pro8 Icw99/Ilama-3-8b-it-kor-extented-chang	llama-3	9B 8R	2024-05-21	2234	-301.88 -527.86	03.70 66.27
747	lemon-mint/gemma-7b-openhermes-v0.80	gemma	8B	2024-04-09	4867	-691.34	56.91
748	lemon-mint/gemma-ko-1.1-2b-it	gemma	2B	2024-04-26	2337	-1058.64	30.92
749	lemon-mint/gemma-ko-7b-instruct-v0.62	gemma	8B	2024-04-03	7543	-589.67	69.25

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
750	lemon-mint/gemma-ko-7b-instruct-v0.71	gemma	8B	2024-04-09	2232	-711.63	59.23
751	leveldevai/TurdusBeagle-7B	mistral	7B	2024-01-18	1926	-577.38	75.15
752	lex-hue/Delexa-7b	mistral	7B	2024-04-05	11797	-575.31	70.86
753	lgaalves/llama-2-13b-chat-platypus	llama-2	13B	2023-09-06	1193	-570.00	53.92
754	Igaalves/nama-2-/b-nf_open-platypus	nama-2	0B 7B	2023-08-30	1204	-304.27	49.75
756	lgaalves/mistral-7b-pratypus ik	mistral	7B 7B	2023-10-10	1259	-557.93	56 29
757	lgaalves/tinyllama-1.1b-chat-v0.3 platypus	llama	1B	2023-10-09	1193	-665.44	34.50
758	lightblue/suzume-llama-3-8B-multilingual (Devine, 2024)	llama-3	8B	2024-04-23	16442	-557.50	65.55
759	lighteternal/Llama3-merge-biomed-8b (Yadav et al., 2023a; Yu et al., 2024a)	llama-3	8B	2024-05-28	2731	-575.20	66.30
760	liminerity/M7-7b	mistral	7B	2024-03-07	4186	-599.03	76.82
761	lim-agents/tora-/b-v1.0 (Gou et al., 2024)	llama-2	7B 7D	2023-10-08	1177	-628.41	48.50
763	Imi-agents/tota-code-70-v1.0 (Couret al., 2024)	llama-1	13B	2023-10-08	1242	-621.23	40.21
764	Imsys/longchat-7b-v1.5-32k	llama	7B	2023-08-01	5114	-650.92	47.95
765	lmsys/vicuna-13b-delta-v1.1 (Zheng et al., 2023; Touvron et al., 2023a)	llama-1	13B	2023-04-12	1216	-1396.42	53.28
766	lmsys/vicuna-13b-v1.1 (Zheng et al., 2023; Touvron et al., 2023a)	llama-1	13B	2023-04-12	2662	-646.61	53.28
767	lmsys/vicuna-13b-v1.3 (Zheng et al., 2023; Touvron et al., 2023a)	llama-1	13B	2023-06-18	11951	-604.80	54.27
768	Imsys/vicuna-13b-v1.5 (Touvron et al., 2023b; Zheng et al., 2023)	llama-2	13B 12P	2023-07-29	48653	-585.07	55.41 54.07
709	lmsys/vicuna-150-v1.5-10k (100v10il et al., 2025b; Zhelig et al., 2025) lmsys/vicuna-7b-delta-v1.1 (Zheng et al., 2023; Touvron et al., 2023a)	llama-1	13B 7B	2023-08-01	1330	-395.75	50.37
771	Imsys/vicuna-7b-v1.3 (Zheng et al., 2023; Touvron et al., 2023a)	llama-1	7B	2023-04-12	31436	-621.79	49.78
772	Imsys/vicuna-7b-v1.5 (Touvron et al., 2023b; Zheng et al., 2023)	llama-2	7B	2023-07-29	368410	-608.63	52.06
773	lmsys/vicuna-7b-v1.5-16k (Touvron et al., 2023b; Zheng et al., 2023)	llama-2	7B	2023-07-31	3362	-622.49	51.42
774	lu-vae/llama2-13B-sharegpt4-orca-openplatypus-8w	llama-2	13B	2023-09-14	1190	-544.97	55.75
775	lu-vae/llama2-13b-sharegpt4-test	llama-2	13B	2023-09-07	1195	-557.27	55.69
776	luffycodes/nash-vicuna-13b-v1dot5-ep2-w-rag-w-simple (Sonkar et al., 2023)	llama-2	13B 12D	2023-08-21	1165	-603.30	55.40
778	luffycodes/vicuna-class-shishya-75-ep3 (Sonkar et al., 2023)	llama-2	13D 7B	2023-12-21	2092	-635.45	46.52
779	luffycodes/vicuna-class-sinsitya-76-cp3 (Sonkar et al., 2023)	llama-2	13B	2023-12-14	1514	-605.75	55.88
780	luffycodes/vicuna-class-tutor-7b-ep3 (Sonkar et al., 2023)	llama-2	7B	2023-12-15	3735	-643.38	51.45
781	lyogavin/Anima-7B-100K	llama-2	7B	2023-09-14	1207	-611.20	42.98
782	macadeliccc/WestLake-7B-v2-laser-truthy-dpo	mistral	7B	2024-01-27	5064	-593.71	75.37
783	macadeliccc/laser-dolphin-mixtral-2x7b-dpo (Gao et al., 2021a; Sharma et al.,	mixtral	12B	2024-01-08	1275	-594.20	67.16
784	2023) malhaiar/Mistral-7B-v0.2-meditron-turkish	mistral	7B	2024-01-05	3865	-641.84	63.34
785	martyn/llama2-megamerge-dare-13b-v2	llama-2	13B	2023-12-17	1213	-594.59	57.94
786	martyn/mistral-megamerge-dare-7b	mistral	7B	2023-12-14	1203	-638.88	48.93
787	martyn/solar-megamerge-dare-10.7b-v1	llama	10B	2023-12-31	1190	-533.45	68.79
788	matsuo-lab/weblab-10b	gpt_neox	10B	2023-08-04	1596	-485.56	38.59
789	matsuo-lab/weblab-10b-instruction-sft	gpt_neox	10B 7P	2023-08-04	1235	-515.77	39.13
790	maywell/PiVoT-10.7B-Mistral-v0.2	mistral	10B	2023-12-13	3207	-578.51	64.25
792	maywell/Synatra-10.7B-v0.4	llama	10B	2023-12-27	3282	-532.85	65.48
793	maywell/Synatra-7B-v0.3-RP	mistral	7B	2023-10-29	8360	-582.37	59.26
794	maywell/Synatra-7B-v0.3-dpo	mistral	7B	2023-11-08	4302	-575.15	60.55
795	maywell/Synatra-RP-Orca-2-7b-v0.1	llama	6B	2023-11-21	3278	-632.02	59.65
796	maywell/Synatra-V0.1-/B-Instruct medalpace/medalpace 7b (Li et al. 2023d)	mistral	/B 7B	2023-10-09	3381 8136	-623.00	55.86 48.45
798	meta-llama/Llama-2-13b-chat-hf (Touvron et al. 2023b)	llama-2	13B	2023-03-29	266090	-616.08	40.43 54 91
799	meta-llama/Llama-2-13b-hf (Touvron et al., 2023b)	llama-2	13B	2023-07-13	120871	-530.82	55.69
800	meta-llama/Llama-2-7b-chat-hf (Touvron et al., 2023b)	llama-2	6B	2023-07-13	1402244	-656.64	50.74
801	meta-llama/Llama-2-7b-hf (Touvron et al., 2023b)	llama-2	6B	2023-07-13	1294737	-549.87	50.97
802	meta-llama/Meta-Llama-3-8B (AI@Meta, 2024)	llama-3	8B	2024-04-17	675850	-514.33	62.62
803	meta-Ilama/Meta-Liama-3-8B-Instruct (Al@Meta, 2024) meta-meth/MetaMeth Liamma 7B (Azerbayov et al. 2024; Vu et al. 2024b)	llama-3	8B 7D	2024-04-17	2058011	-5/3.95	52 10
804 805	meta-math/MetaMath-Mistral-7B (Jiang et al. 2023; Yu et al. 2024b)	mistral	7В	2023-11-19	3022	-571.33	55.19 65.78
806	microsoft/Orca-2-13b (Mitra et al., 2023)	llama	13B	2023-11-14	13616	-651.29	58.64
807	microsoft/Orca-2-7b (Mitra et al., 2023)	llama	7B	2023-11-14	11084	-684.33	54.55
808	microsoft/Phi-3-medium-128k-instruct	phi3	13B	2024-05-07	34054	-544.42	73.00
809	microsoft/Phi-3-medium-4k-instruct	phi3	13B	2024-05-07	35987	-552.43	73.45
810	microsoft/phi 2	phi	1B 2B	2023-09-10	112470	-724.15	47.09 61.33
812	migtissera/Llama-3-8B-Synthia-v3.5	llama-3	2B 8B	2023-12-13	3310	-533 92	67.15
813	miglissera/SynthIA-7B-v1.3 (Mukherjee et al., 2023; Tissera, 2023)	mistral	7B	2023-09-28	3336	-541.46	59.34
814	migtissera/SynthIA-7B-v1.5	mistral	7B	2023-10-07	1228	-537.68	59.59
815	migtissera/Synthia-7B-v3.0	mistral	7B	2023-12-08	1204	-531.17	61.99
816	migtissera/Tess-2.0-Llama-3-8B	llama-3	8B	2024-05-05	3314	-523.68	64.81
817	migtissera/Tess-/B-v1.4	mistral	7B 7D	2023-12-04	1234	-599.60	62.19 50.20
819	mindy-labs/mindy-7b-v2	mistral	7B 7B	2023-11-22	1213	-550.38	72 11
820	mistral-community/Mistral-7B-v0.2	mistral	7B	2024-03-23	30074	-532.63	60.41
821	mistralai/Mistral-7B-Instruct-v0.1 (Jiang et al., 2023)	mistral	7B	2023-09-27	1370245	-593.06	54.96
822	mistralai/Mistral-7B-Instruct-v0.2 (Jiang et al., 2023)	mistral	7B	2023-12-11	3565248	-574.91	65.71
823	mistralai/Mistral-7B-v0.1 (Jiang et al., 2023)	mistral	7B	2023-09-20	1750089	-527.60	60.97
824	mistralai/Mistral-7B-v0.3 (Jiang et al., 2023)	mistral	7B	2024-05-22	1443065	-531.44	60.28
825 826	maoonne/Appiavionator-75	mistral	/В 7R	2024-02-14	12004	-595.40 -563.97	73.99 74 76
827	mlabonne/ChimeraLlama-3-8B-v2	llama-3	8B	2024-04-22	2806	-565.81	69.69
828	mlabonne/ChimeraLlama-3-8B-v3	llama-3	8B	2024-05-01	5718	-568.78	70.06
829	mlabonne/Daredevil-8B-abliterated	llama	8B	2024-05-26	9010	-565.64	71.82
830	mlabonne/GML-Mistral-merged-v1	mistral	8B	2023-12-27	1166	-1384.76	48.54
831	miabonne/Marcoro14-7B-slerp	mistral	7B	2023-12-29	3792	-554.34	73.01
832 833	mlabonne/NeuralMonarch-7B	mistral	/В 7R	2024-01-00	13486	-500.09	76.15
000	······································		7.0	202102-14	10,000	271.01	

ID	Model Name	Model Type	Size	Date	DLs	$\overline{\ell}_i$	Task
834	mncai/agiin-13.6B-v0.1	mistral	13B	2023-12-15	3317	-595.75	68.40
835	monology/openinstruct-mistral-7b	mistral	7B	2023-11-20	1191	-537.83	63.64
836	mosaicml/mpt-7b (Henry et al., 2020; Shoeybi et al., 2020; Press et al., 2022; Dec et al. 2022; Chen et al. 2023a; MagaieML NL P. Team. 2023b; Teauren	mpt	7B	2023-05-05	31480	-548.53	44.28
	et al., 2023a)						
837	mosaicml/mpt-7b-8k (Henry et al., 2020; Shoeybi et al., 2020; Press et al.,	mpt	7B	2023-06-30	2106	-546.57	47.24
	2022; Dao et al., 2022; Chen et al., 2023a; MosaicML NLP Team, 2023b;						
838	mosaicml/mpt-7b-8k-chat (Henry et al., 2020; Dao et al., 2022; Press et al.,	mpt	7B	2023-06-22	1304	-586.09	47.78
	2022; MosaicML NLP Team, 2023a)	•	-	2022 05 04	00070	< 1 <b>7</b> 10	15.00
839	mosaicml/mpt-/b-chat (Henry et al., 2020; Dao et al., 2022; Press et al., 2022; MosaicML NLP Team 2023b)	mpt	7B	2023-05-04	88069	-647.13	45.39
840	mosaicml/mpt-7b-instruct (Henry et al., 2020; Dao et al., 2022; Press et al.,	mpt	7B	2023-05-05	8599	-569.62	44.83
941	2022; MosaicML NLP Team, 2023b)	mat	70	2022 05 04	1020	650.05	20.21
641	2023a: MosaicML NLP Team, 2023b)	mpt	/ D	2023-03-04	1920	-030.03	39.31
842	mrm8488/Ilama-2-coder-7b (Manuel Romero, 2023)	llama-2	7B	2023-07-26	1313	-572.92	49.95
843	mwitiderrick/open_llama_3b_code_instruct_0.1	llama	3B 7D	2023-12-11	1252	-598.99	39.72
845	nlpguy/ColorShadow-7B-v2	mistral	7B 7B	2023-12-30	1214	-569.63	66.88
846	nlpguy/ColorShadow-7B-v3	mistral	7B	2023-12-30	1222	-561.80	67.29
847	nvidia/Llama3-ChatQA-1.5-8B (Liu et al., 2024)	llama-3	8B	2024-04-28	10608	-515.82	56.71
848	occultml/CatMarcoro14-7B-slerp	mistral	7B	2024-01-06	1234	-551.33	73.25
849	occultml/Helios-10./B	llama	7B 7D	2023-12-31	1215	-1381.38	42.19
851	oh-veontaek/llama-2-13B-LoRA-assemble	llama-2	13B	2023-09-13	3303	-568.38	57.91
852	oh-yeontaek/llama-2-7B-LoRA-assemble	llama-2	7B	2023-09-13	1337	-614.33	52.26
853	openaccess-ai-collective/DPOpenHermes-11B	mistral	10B	2023-12-03	1185	-591.33	66.83
854	openaccess-ai-collective/jackalope-7b (Mukherjee et al., 2023; Lian et al.,	mistral	7B	2023-10-07	1198	-550.33	61.16
855	2023b; Longpre et al., 2023)	llama 1	13B	2023 05 17	1224	553.22	54.86
856	openaccess-ai-collective/manticore-13b-chat-pvg	llama-1	13B 13B	2023-05-17	214	-553.85	54.00
857	openaccess-ai-collective/manacol-130-chat-pyg	llama-1	13B	2023-06-06	1202	-577.38	53.97
858	openaccess-ai-collective/minotaur-13b-fixed	llama-1	13B	2023-06-12	1194	-568.52	55.19
859	openaccess-ai-collective/mistral-7b-slimorcaboros	mistral	7B	2023-10-13	1186	-569.20	61.18
860	openaccess-ai-collective/wizard-mega-13b	llama-1	13B	2023-05-14	2168	-565.46	54.27
861	openbmb/UltraLM-13b-v2.0	llama	13B	2023-09-22	1174	-554.24	58.72
862	openchat/openchat-3.5-0106 (Wang et al., 2024a; OpenAI et al., 2024)	mistral	7B	2024-01-07	26858	-554.42	69.30
803 864	openchat/openchat-3.5-0106-gemma (wang et al., 2024a) openchat/openchat-3.5-1210 (Wang et al., 2024a; OpenAL et al., 2024)	mistral	8B 7B	2024-03-09	2123	-937.52	69.42 68.89
865	openchat/openchat-3.6-8b-20240522 (Wang et al., 2024a)	llama-3	7B 8B	2023-12-12	10464	-561.80	68.14
866	openIm-research/open_llama_13b (Geng and Liu, 2023; Together Computer,	llama-1	13B	2023-06-15	2373	-582.20	47.26
867	2023; Touvron et al., 2023a) openIm-research/open_llama_3b (Geng and Liu, 2023; Together Computer,	llama-1	3B	2023-06-07	157405	-612.46	38.26
868	2023; Touvron et al., 2023a) openIm-research/open Ilama 3b v2 (Geng and Liu, 2023; Together Com-	llama-1	3B	2023-07-16	21275	-578.08	40.28
869	puter, 2023; Touvron et al., 2023a) openIm-research/open llama 7b (Geng and Liu, 2023; Together Computer,	llama-1	7B	2023-06-07	44968	-610.20	42.31
870	2023; Touvron et al., 2023a) openIm-research/open Ilama 7b v2 (Geng and Liu, 2023: Together Com-	llama-1	7B	2023-07-06	2792	-556.22	44 26
071	puter, 2023; Touvron et al., 2023a)	11	120	2022 12 10	1005	710 (1	50.45
871	openthaigpt/openthaigpt-1.0.0-beta-13b-chat-hf	llama	13B 7D	2023-12-18	1225	-712.61	50.45
873	pankaimathur/orca mini v3 13b (Mukheriee et al. 2023: Touvron et al.	llama	13B	2023-08-09	4717	-546.48	43.33 57.24
874	2023b; Mathur, 2023b) pankaimathur/orca mini v3 7b (Mukheriee et al. 2023; Touvron et al.	llama	7B	2023-08-07	2338	-572.18	53 47
875	2023b; Mathur, 2023c; Touvron et al., 2023a) pe-nh/llama-2-13b-vicuna-wizard	llama-2	13B	2023-08-11	1176	-549 39	51 94
876	pillowtalks-ai/delta13b	llama-1	13B	2023-04-14	1168	-646.61	53.29
877	princeton-nlp/Sheared-LLaMA-1.3B (Xia et al., 2024)	llama	1B	2023-10-10	28905	-646.81	35.95
878	princeton-nlp/Sheared-LLaMA-1.3B-ShareGPT (Xia et al., 2024)	llama	1B	2023-11-22	1737	-721.70	37.14
879	princeton-nlp/Sheared-LLaMA-2.7B (Xia et al., 2024)	llama-2	2B	2023-10-10	2590	-610.74	40.84
880	princeton-nlp/Sheared-LLaMA-2.7B-ShareGPT (Xia et al., 2024)	llama-2	2B	2023-11-22	1866	-681.47	42.11
882	project-baize/baize-v2-150 (Xu et al., 2025b) project-baize/baize-v2-7b (Xu et al., 2025b)	llama-1	13D 7B	2023-05-23	1212	-578.09	32.94 46.72
883	anguven3/Master-Yi-9B	llama	8B	2024-05-18	8810	-522.22	67.44
884	quantumaikr/QuantumLM-7B	llama	7B	2023-07-22	1192	-625.18	49.51
885	quantumaikr/llama-2-7b-hf-guanaco-1k	llama-2	7B	2023-08-06	1186	-605.93	50.13
886	quantumaikr/quantum-v0.01	mistral	7B	2023-12-17	1192	-559.84	74.68
887	refuelai/Llama-3-Refueled	llama-3	8B	2024-05-03	1246	-582.52	63.62
888	revolutionarybukhari/Llama-2-/b-chat-finetune-AUTOMATE	llama-2	/B 2P	2023-10-14	2267	-613.72	22.14
890	rinna/oninguai-gpt-neox-40 (Zhao et al., 2023), Sawada et al., 2024)	gpt_neox	3B 3B	2023-07-31	6042	-1117 39	29.28
891	rinna/lama-3-youko-8b (Andonian et al., 2021; AI@Meta, 2024; Sawada et al. 2024). Mituda et al. 2024	llama-3	8B	2023-05-01	1468	-502.57	57.55
892	rinna/youri-7b (Andonian et al., 2024) 2023; Sawada et al. 2024; Touvron et al., 2023b; Zhao et al., 2023a; Sawada et al. 2024)	llama-2	7B	2023-10-30	2512	-545.03	47.11
893	rishiraj/CatPPT-base (Acharya, 2023)	mistral	7B	2023-12-17	4392	-558.36	72.25
894	ruslanmv/Medical-Llama3-8B	llama-3	8B	2024-04-21	5457	-514.00	60.61
895	ruslanmv/ai-medical-model-32bit	llama	8B	2024-05-13	2810	-557.43	67.67
896	rwitz2/go-bruins-v2.1 (Murias, 2023)	mistral	7B	2023-12-14	1193	-561.68	74.50
897	rwitz2/go-pruins-v2.1.1 (Murias, 2023)	mistral	7B 2B	2023-12-14	1205	-562.22	74.95 38.02
090 899	samir-fama/FernandoGPT-v1	gpt_neox mistral	2D 7R	2023-12-23	1200	-551 17	72.87
900	samir-fama/SamirGPT-v1	mistral	7B	2023-12-28	1215	-552.06	73.11
901	sarvamai/OpenHathi-7B-Hi-v0.1-Base	llama-2	6B	2023-12-13	1766	-673.13	46.64
902	scaledown/ScaleDown-7B-slerp-v0.1	mistral	7B	2024-01-01	1206	-538.58	71.57

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
903	scb10x/llama-3-typhoon-v1.5-8b-instruct (Pipatanakul et al., 2023)	llama-3	8B	2024-05-06	6088	-590.46	65.62
904	scb10x/typhoon-7b (Pipatanakul et al., 2023)	mistral	7B	2023-12-20	1908	-617.60	58.05
905	selfrag/selfrag_llama2_7b (Asai et al., 2023)	llama-2	7B	2023-10-18	4388	-605.61	51.30
906	senseable/WestLake-7B-v2	mistral	7B	2024-01-22	1189	-594.85	74.68
907	setnuiyer/Medichat-Liama5-8B shadowml/BeagSake 7B	nama-5	8B 7B	2024-04-22	4542	-555.77	00.03 75.38
908	shaudwill/BeagSake-7B	llama-3	7B 8B	2024-01-31	2708	-502.99	63.00
910	shanchen/llama3-8B-slerp-med-262k	llama-3	8B	2024-04-30	2697	-599.35	53.65
911	shanchen/llama3-8B-slerp-med-chinese	llama-3	8B	2024-04-30	8028	-593.26	58.99
912	shenzhi-wang/Llama3-8B-Chinese-Chat (Wang et al., 2024b)	llama-3	8B	2024-04-21	55718	-550.92	67.10
913	shibing624/chinese-alpaca-plus-7b-hf (Ming, 2023)	llama-1	7B	2023-05-01	1527	-741.70	44.77
914	shitshow123/tinylamma-20000	llama	1B 7D	2024-01-09	1195	-1213.36	27.95
915	speakleash/Biolik 7B Instruct v0.1 (Levine et al. 2020; Granziol et al. 2021;	mistral	7B 7B	2024-04-21	4013	-000.14	57.44
910	Ociepa et al., 2024b; Wang et al., 2024a; Ociepa et al., 2024c)	mistrai	7 <b>D</b>	2024-03-30	5575	-0/1.00	51.24
917	speakleash/Bielik-7B-v0.1 (Ociepa et al., 2024a,c)	mistral	7B	2024-03-30	2822	-721.18	50.01
918	stabilityai/StableBeluga-13B (Touvron et al., 2023b; Mukherjee et al., 2023)	llama	13B	2023-07-27	6154	-540.85	57.05
919	stabilityai/StableBeluga-7B (Touvron et al., 2023b; Mukherjee et al., 2023)	llama	6B	2023-07-27	6691	-565.19	53.56
920	stabilityai/japanese-stablelm-base-gamma-7b (Jiang et al., 2023)	mistral	/B 7D	2023-10-16	2072	-581.58	52.59
921	stabilityai/japanese-stablem-instruct-gamma-70 (Jiang et al., 2025) stabilityai/stable-code-3b (Raibhandari et al., 2020; Black et al., 2022; Li	stablelm	7 D 2 B	2023-10-10	1412 5670	-564.50	52.82 41.53
/22	et al., 2023a: Su et al., 2023a: Touvron et al., 2023b: Pinnaparaju et al., 2024:	studienn	20	20210109	5070	010.20	11.55
	Azerbayev et al., 2024; Yu et al., 2024b)						
923	stabilityai/stablelm-2-12b-chat (Stability AI Language Team, 2024; Rafailov	stablelm	12B	2024-04-04	4843	-588.55	68.38
024	et al., 2024) stabilitya/stablelm 2.1.6b.shat (Stability ALL anguage Team, 2024; Bafailoy	stablalm	1D	2024 04 08	4156	692.01	50.71
924	et al., 2024)	stablemi	ID	2024-04-08	4150	-005.91	50.71
925	stabilityai/stablelm-2-zephyr-1_6b (Stability AI Language Team, 2024;	stablelm	1B	2024-01-19	18239	-660.56	49.99
026	Rafailov et al., 2024)	atablalm	20	2022 00 20	11112	510.56	16 50
926	stabilityal/stableim-50-4eft (Ba et al., 2010; Zhang and Sennfich, 2019; Kajb- bandari et al. 2020; Gao et al. 2020; Black et al. 2022; Li et al. 2023a; Su	stableim	2 <b>B</b>	2023-09-29	11112	-510.56	40.58
	et al. 2023a: Tow et al. 2023; Touvron et al. 2023b)						
927	stabilityai/stablelm-base-alpha-3b (Andonian et al., 2021)	gpt_neox	3B	2023-04-17	1728	-677.47	31.50
928	stabilityai/stablelm-base-alpha-7b (Andonian et al., 2021)	gpt_neox	7B	2023-04-11	1750	-623.09	34.37
929	stabilityai/stablelm-base-alpha-7b-v2 (Shazeer, 2020; Rajbhandari et al., 2020;	stablelm_alpha	6B	2023-08-04	2209	-505.45	46.18
930	Gao et al., 2020; Li et al., 2023a; Su et al., 2023a; Tow, 2023) stabilityai/stablelm-tuned-alpha-3b (Taori et al., 2023; Apand et al., 2023;	ant neox	3B	2023-04-19	2145	-736.44	32.14
250	Chiang et al., 2023)	spt_neox	50	2025 01 19	2115	750.11	52.11
931	stabilityai/stablelm-tuned-alpha-7b (Taori et al., 2023; Anand et al., 2023;	gpt_neox	7B	2023-04-19	3765	-694.37	34.04
	Chiang et al., 2023)						
932	stabilityai/stablelm-zephyr-3b (Zheng et al., 2023; Rafailov et al., 2024)	stablelm	2B	2023-11-21	8261	-777.98	53.43
955	starmpcc/Asclepius-Liama2-15B (Kweon et al., 2024)	llama 2	15D 7B	2023-09-19	1231	-045.45	30.23 47.15
935	starking/zephyr-7b-sft-full-orpo	mistral	7B	2024-05-18	2278	-548.89	53.16
936	swap-uniba/LLaMAntino-3-ANITA-8B-Inst-DPO-ITA (Basile et al., 2023;	llama-3	8B	2024-04-29	5995	-635.92	75.12
	AI@Meta, 2024; Polignano et al., 2024)						
937	teilomillet/MiniMerlin-3B	llama	3B	2023-12-15	1168	-681.62	47.63
938	teknium/OpenHermes-13B	llama-2	13B 7D	2023-09-06	1594	-543.02	55.24
939	time/falcon_7b (Shazeer 2019; Brown et al. 2020; Gao et al. 2020; Dao	falcon	7B	2023-10-29	100997	-549.44	01.45 44 17
240	et al., 2022: Su et al., 2023a: Penedo et al., 2023; Almazrouei et al., 2023)	Taleon	/10	2025-04-24	104010	-547.44	44.17
941	tiiuae/falcon-7b-instruct (Shazeer, 2019; Brown et al., 2020; Dao et al., 2022;	falcon	7B	2023-04-25	179952	-621.07	43.16
0.42	Su et al., 2023a; Penedo et al., 2023; Almazrouei et al., 2023)	6.1	10	2022 04 26	22021	(00.26	27.07
942	tituae/falcon-rw-1b (Brown et al., 2020; Dao et al., 2022; Press et al., 2022; Penedo et al. 2023)	falcon	IB	2023-04-26	22921	-688.36	37.07
943	timpal0l/Mistral-7B-v0.1-flashback-v2	mistral	7B	2023-12-04	1275	-578.05	57.53
944	togethercomputer/GPT-JT-6B-v0	gptj	6B	2022-11-22	1422	-484.70	44.05
945	togethercomputer/GPT-JT-6B-v1 (Tay et al., 2022, 2023)	gptj	6B	2022-11-24	5765	-503.02	43.13
946	togethercomputer/LLaMA-2-7B-32K	llama-2	7B	2023-07-26	8884	-530.20	47.07
947	togethercomputer/Llama-2-7B-32K-Instruct (Liu et al., 2023)	llama-2	7B	2023-08-08	5596	-564.83	50.02
948	togethercomputer/Pythia-Chat-Base-/B	gpt_neox	/B 7B	2023-03-22	/040 1487	-522.45	39.81
950	togethercomputer/RedPajama-INCITE-7B-Chat	gpt_neox	7B	2023-05-04	1583	-931.82	39 37
951	togethercomputer/RedPajama-INCITE-7B-Instruct	gpt_neox	7B	2023-05-05	1274	-525.00	42.38
952	togethercomputer/RedPajama-INCITE-Base-3B-v1	gpt_neox	3B	2023-05-04	2635	-590.97	38.54
953	togethercomputer/RedPajama-INCITE-Chat-3B-v1	gpt_neox	3B	2023-05-05	1647	-610.76	39.53
954	togethercomputer/RedPajama-INCITE-Instruct-3B-v1	gpt_neox	3B	2023-05-05	2014	-552.27	39.06
955	totally-not-an-llm/EverythingLM-13b-16k	llama-2	13B	2023-08-12	2132	-557.79	52.33
956	totally-not-an-lim/PuddleJumper-13b-V2	llama mistral	13B 13B	2023-09-21	3004	-033.45	54.19 73.57
958	unsloth/Phi-3-mini-4k-instruct	mistral	3B	2024-03-23	12224	-575 74	69.86
959	unsloth/mistral-7b-v0.2	mistral	7B	2024-03-24	4119	-532.63	60.34
960	unsloth/tinyllama-chat	llama	1B	2024-02-14	5959	-619.71	37.24
961	upstage/SOLAR-10.7B-Instruct-v1.0 (Kim et al., 2024b,a)	llama	10B	2023-12-12	67725	-557.67	74.20
962	upstage/SOLAR-10.7B-v1.0 (Kim et al., 2024b)	llama	10B	2023-12-12	24478	-525.51	66.04
963	uukuguy/speechless-code-mistral-/b-v1.0	mistral	/B	2023-10-10	4393	-540.50	58.85 57.17
904	2023b)	nama-2	13B	2023-09-01	1308	-301.91	57.17
965	uukuguy/speechless-llama2-hermes-orca-platypus-wizardlm-13b (Touvron	llama-2	13B	2023-09-01	2101	-609.55	57.52
0.11	et al., 2023b)	• . •		2024 01 22	1000	<b>53</b> 0 50	(2.02
966	uukuguy/speechless-zephyr-code-tunctionary-7b	mistral	7B	2024-01-23	4098	-529.18	62.93
907	uuxuguy/zepiiyi-/b-aipiia-uare-0.85 uvgarkurt/llama-3-merged-linear	misual Ilama 3	/В 90	2023-11-23	0141	-529.44	02.33 73.02
969	v1olet/v1olet marcoroni-go-bruins-merge-7B	mistral	7B	2023-12-11	1225	-559.92	72.81
970	v1olet/v1olet_merged_dpo_7B	mistral	7B	2023-12-12	1210	-592.71	70.26
971	varox34/Bio-Saul-Dolphin-Beagle-Breadcrumbs	mistral	7B	2024-05-01	2679	-610.01	48.72
972	vibhorag101/llama-2-13b-chat-hf-phr_mental_therapy	llama-2	13B	2023-09-17	1269	-630.34	42.50

ID	Model Name	Model Type	Size	Date	DLs	$\bar{\ell}_i$	Task
973	vicgalle/CarbonBeagle-11B (Wortsman et al., 2022)	mistral	10B	2024-01-21	6696	-549.83	74.64
974	vicgalle/CarbonBeagle-11B-truthy	mistral	10B	2024-02-10	13887	-554.84	76.10
975	vicgalle/Configurable-Hermes-2-Pro-Llama-3-8B (Gallego, 2024)	llama-3	8B	2024-05-02	10273	-580.31	70.10
976	vicgalle/Configurable-Llama-3-8B-v0.3 (Gallego, 2024)	llama-3	8B	2024-04-20	6030	-555.07	68.79
977	vicgalle/Configurable-Yi-1.5-9B-Chat (Gallego, 2024)	llama	8B	2024-05-12	6537	-633.71	70.50
978	vicgalle/ConfigurableBeagle-11B (Gallego, 2024)	mistral	10B	2024-02-17	6045	-548.90	75.40
979	vicgalle/ConfigurableHermes-7B (Gallego, 2024)	mistral	7B	2024-02-17	6085	-572.68	68.89
980	vicgalle/ConfigurableSOLAR-10.7B (Gallego, 2024)	llama	10B	2024-03-10	5135	-559.06	73.94
981	viethq188/LeoScorpius-7B	mistral	7B	2023-12-12	1194	-552.34	72.21
982	viethq188/Rabbit-7B-v2-DPO-Chat	mistral	7B	2023-12-12	1181	-579.23	69.36
983	vihangd/dopeyplats-1.1b-2T-v1	llama	1B	2023-11-26	1191	-682.92	35.28
984	vihangd/dopeyshearedplats-1.3b-v1	llama-2	1B	2023-12-12	1168	-766.02	36.74
985	vihangd/dopeyshearedplats-2.7b-v1	llama-2	2B	2023-12-16	1173	-709.66	42.90
986	vihangd/neuralfalcon-1b-v1	falcon	1B	2023-12-17	1185	-895.18	29.72
987	vihangd/shearedplats-1.3b-v1	llama-2	1B	2023-11-16	1184	-706.16	35.97
988	vihangd/shearedplats-2.7b-v2	llama-2	2B	2023-11-18	2081	-647.66	41.61
989	vihangd/smartyplats-3b-v1	llama	3B	2023-09-11	1180	-600.04	40.00
990	vihangd/smartyplats-3b-v2	llama	3B	2023-09-14	1179	-597.07	40.29
991	vikash06/llama-2-7b-small-model-new	llama-2	6B	2023-12-22	1174	-925.37	46.62
992	vmajor/Orca2-13B-selfmerge-26B	llama	13B	2023-12-01	2041	-653.10	62.24
993	vmajor/Orca2-13B-selfmerge-39B	llama	13B	2023-12-01	1208	-653.10	62.24
994	vonjack/Qwen-LLaMAfied-HFTok-7B-Chat	llama-2	7B	2023-08-09	1189	-839.96	50.64
995	w601sxs/b1ade-1b	gpt_neox	1B	2023-07-17	1204	-929.29	32.59
996	wang7776/Llama-2-7b-chat-hf-10-sparsity (Sun et al., 2024)	llama-2	6B	2023-12-11	1178	-660.87	52.48
997	wang7776/Llama-2-7b-chat-hf-20-sparsity (Sun et al., 2024)	llama-2	7B	2023-12-13	1175	-668.34	52.01
998	wang7776/Llama-2-7b-chat-hf-30-sparsity (Sun et al., 2024)	llama-2	6B	2023-12-11	1179	-676.50	51.02
999	webbigdata/ALMA-7B-Ja-V2 (Xu et al., 2024a)	llama-2	7B	2023-10-21	1177	-599.69	47.85
1000	wei123602/Llama-2-13b-FINETUNE4_TEST	llama-2	13B	2023-09-18	1174	-538.02	53.62
1001	wenbopan/Faro-Yi-9B (OpenAI et al., 2024)	llama	8B	2024-03-27	6127	-594.07	66.37
1002	wenbopan/Faro-Yi-9B-DPO (OpenAI et al., 2024)	llama	8B	2024-04-07	6121	-597.40	68.77
1003	wenge-research/yayi-7b	bloom	7B	2023-06-02	1192	-653.40	41.88
1004	wenge-research/yayi-7b-llama2	llama-2	7B	2023-07-21	1195	-562.60	49.88
1005	winglian/Llama-2-3b-hf	llama-2	3B	2023-09-19	1565	-1376.16	29.53
1006	winglian/llama-2-4b	llama-2	4B	2023-09-19	1193	-676.31	34.23
1007	xDAN-AI/xDAN-L1-Chat-RL-v1	mistral	7B	2023-12-20	1178	-574.92	68.38
1008	yam-peleg/Hebrew-Mistral-7B	mistral	7B	2024-04-26	5993	-625.28	58.76
1009	yanolja/Bookworm-10.7B-v0.4-DPO (Mukherjee et al., 2023; Lian et al., 2023g: Cui et al. 2024)	llama	10B	2024-01-18	2239	-586.77	66.59
1010	yanolja/EEVE-Korean-Instruct-10.8B-v1.0 (Lian et al., 2023g; Mukherjee	llama	10B	2024-02-22	13241	-555.03	66.48
	et al., 2023; Kim et al., 2024c; Cui et al., 2024)		-	2022 10 02	2200	1 402 04	20.56
1011	yeen214/llama2_7b_small_tuning_v1	llama-2	7B	2023-10-02	3289	-1402.04	28.56
1012	yeen214/test_llama2_7b	llama-2	7B	2023-09-30	3289	-549.87	49.73
1013	yeen214/test_llama2_ko_7b	llama-2	7B	2023-10-02	3285	-1405.43	29.99
1014	yhyhy3/open_llama_/b_v2_med_instruct	llama-l	7B	2023-07-09	1197	-561.62	46.24
1015	yulan-team/YuLan-Chat-2-13b-tp16 (YuLan-Team, 2023)	IIama	13B	2023-08-04	1165	-645.56	57.01
1016	yunconglong/DARE_TIES_13B (Yadav et al., 2023a; Yu et al., 2024a)	mixtral	12B	2024-01-30	7029	-611.23	77.10
1017	yunconglong/MoE_13B_DPO	mixtral	12B	2024-01-28	3939	-603.88	77.05
1018	yunconglong/Truthful_DPO_TomGrc_FusionNet_7Bx2_MoE_13B	mıxtral	12 <b>B</b>	2024-01-21	7965	-598.84	77.44

Table 6: List of 1018 models. "ID" denotes the alphabetical index; "Model Name" denotes the name of the model; "Model Type" denotes the classification defined in this paper; "Size" denotes the size of the model (B: billion); "Date" denotes the date of model creation; "DLs" denotes the total number of downloads;  $\bar{\ell}_i$  denotes the mean log-likelihood; "Task" denotes the mean of the 6 benchmark scores (i.e., 6-TaskMean).