

ABLE: Representing and Mapping LLMs via Attribution-Based Large-model Embedding

Anonymous ACL submission

Abstract

The explosive growth of large language models (LLMs) has led to an opaque ecosystem where undocumented model relationships hinder copyright protection, security auditing, and model routing. Existing representation methods struggle to address this challenge efficiently. Approaches analyzing internal parameters face scalability barriers due to structural heterogeneity across diverse architectures, while methods relying on external outputs are susceptible to behavioral mimicry, where distinct models converge to similar predictions despite differing underlying mechanisms. To bridge this gap, we propose ABLE (Attribution-Based Large-model Embedding), a novel framework that leverages the interpretability space to construct model representations. By aggregating gradient-based feature attributions via a tokenizer-agnostic word-level alignment, ABLE captures the intrinsic cognitive patterns of models rather than surface-level outputs. Beyond empirical utility, we proved that ABLE is a Lipschitz continuous mapping with finite-sample convergence guarantees, ensuring stability and reliability. Extensive experiments on 239 LLMs demonstrate that our training-free approach achieving competitive or superior performance in relation prediction, model routing and benchmark score prediction¹.

1 Introduction

The ecosystem of large language models (LLMs) (Vaswani et al., 2017; Yang et al., 2024; Hao et al., 2022) is expanding at an explosive rate, with platforms like Hugging Face hosting hundreds of thousands of models (Zhao et al., 2023; Jain, 2022; Rothman, 2022). Yet this rapid growth comes with a transparency crisis: many models are released

¹Our code for ABLE extraction and analysis are available at this link: <https://anonymous.4open.science/r/ABLE>

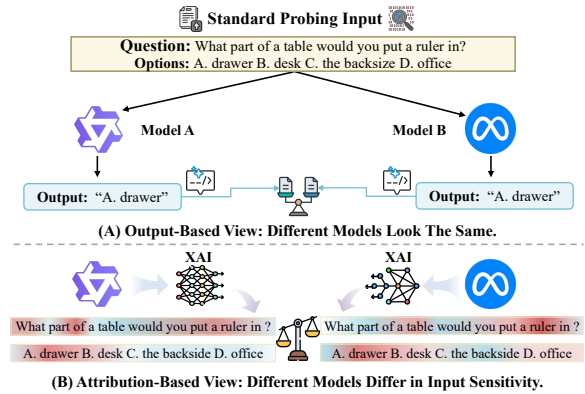


Figure 1: Motivation for attribution-based representations. (A) Given the same probing input, the output-based view cannot distinguish models that produce identical answers. (B) The attribution-based view reveals differences in input sensitivity, capturing how models rely on different tokens to reach the same prediction.

without clear documentation of their training origins, making the relationships between models opaque and invisible (Barman et al., 2024). Uncovering these hidden model relationships is crucial for copyright protection, enabling detection of unauthorized model reuse and intellectual property infringement (Menell, 1988; Shao et al., 2024; Horwitz et al., 2025); security auditing, allowing practitioners to trace backdoor propagation paths within model families (Cheng et al., 2024; Ryoo et al., 2013); and intelligent model routing, facilitating selection of the most suitable model based on capability similarity (Ong et al., 2024).

To address these challenges, researchers seek to construct compact LLM representations that enable systematic comparison and analysis across large model collections. Existing approaches generally fall into two categories: internal feature analysis and external output probing. Methods based on internal features directly analyze model parameters or activations (Yadav et al., 2023; Zhu et al., 2025). However, the structural heterogeneity in-

062	herent in diverse models, such as differences in	114
063	architecture or scale, necessitates complex weight	
064	alignment or layer mapping strategies (Mattheakis	
065	et al., 2019). These requirements impede the scalability	
066	of such methods for large-scale model representation.	
067	Conversely, methods relying on external	
068	outputs offer superior scalability by treating models	
069	as black boxes (Zhuang et al., 2024; Yax et al.,	
070	2024; Oyama et al., 2025). However, in applica-	
071	tions such as copyright protection, since distinct	
072	model architectures or training strategies may converge	
073	to similar output distributions, relying solely	
074	on surface-level behavioral overlap can lead to erroneous	
075	conclusions (Orgad et al., 2024). Therefore,	
076	accurately representing LLMs and detecting genuine	
077	relationships requires exploring the underlying	
078	mechanistic similarities between models.	
079	As illustrated in Fig. 1, much like distinct cognitive	
080	processes in humans, different models may produce	
081	identical outputs while relying on entirely different	
082	input features. Building on this insight, we propose	
083	ABLE (Attribution-Based Large-model Embedding),	
084	a framework that leverages the interpretability space	
085	to construct representations of LLMs. By computing	
086	the gradient of the model’s output with respect to	
087	its inputs via the $Gradient \times Input$ (Kim et al.,	
088	2018; Wang et al., 2024), ABLE quantifies the	
089	contribution of each input feature to a specific	
090	prediction, thereby capturing the mechanistic	
091	patterns governing how a model processes	
092	information (Sturmfels et al., 2020; Nguyen et	
093	al., 2021; Kommiya Mothilal et al., 2021; Zhou	
094	et al., 2022). To ensure scalability across diverse	
095	architectures, we implement word-level alignment	
096	to map token-specific attributions into a unified	
097	vocabulary, effectively resolving tokenizer	
098	heterogeneity. We validate this training-free	
099	framework on 239 LLMs, demonstrating its	
100	effectiveness in model relationship analysis,	
101	routing, and benchmark score prediction.	
102	Furthermore, we provide theoretical guarantees	
103	regarding the stability of the representation. Under	
104	standard regularity assumptions, we prove that	
105	ABLE constitutes a Lipschitz continuous mapping	
106	(Johnson et al., 1984) from the model parameter	
107	space to the embedding space, ensuring that	
108	models with similar parameters yield similar	
109	embeddings. We further establish finite-sample	
110	concentration bounds to guarantee that empirical	
111	representations converge to population-level	
112	embeddings. These results collectively establish	
113	ABLE as a principled and stable representation	
	grounded	
	in model parameters.	114
	In summary, our contributions are as follows:	115
	• We propose a novel paradigm that leverages	116
	the interpretability space to construct model	117
	representations. We introduce ABLE, a frame-	118
	work that utilizes gradient-based feature	119
	attribution to capture mechanistic similarities	120
	between models, thereby overcoming the	121
	limitations of surface-level output analysis.	122
	• We provide theoretical guarantees for the	123
	stability of our representation. We prove that	124
	ABLE is a Lipschitz-continuous mapping	125
	from the model parameter space to the	126
	embedding space, ensuring that models with	127
	similar parameters naturally yield similar	128
	embeddings under standard regularity	129
	assumptions.	
	• We validate ABLE on 239 LLMs across	130
	diverse tasks including lineage reconstruction,	131
	model routing, and benchmark score	132
	prediction, demonstrating its scalability	133
	and effectiveness.	134
	2 Related Work	135
	2.1 LLM Representation	136
	This section reviews related work on	137
	constructing representations for LLMs,	138
	covering both internal-feature and	139
	output-based approaches.	140
	Internal-feature methods directly analyze	141
	model parameters or activations. Representative	142
	approaches study parameter changes from	143
	fine-tuning to resolve conflicts during	144
	model merging (Yadav et al., 2023),	145
	test derivative relationships by comparing	146
	model weights (Zhu et al., 2025), and	147
	analyze activation patterns across layers	148
	to quantify cross-model similarity (Zhou	149
	et al., 2024). These methods are well-	150
	suited for scenarios where models share	151
	compatible architectures, but require	152
	careful alignment when comparing	153
	models with different structures.	154
	Output-based methods construct	155
	representations from model predictions	156
	or logits. Existing researches generate	157
	model embeddings from log-likelihoods	158
	on a standardized corpus (Oyama et	159
	al., 2025), construct phylogenetic trees	160
	by analyzing patterns in generated	161
	text (Yax et al., 2024), and learn	
	embeddings from task performance	
	vectors for model routing (Zhuang	
	et al., 2024). These methods are	
	architecture-agnostic and scalable,	
	but are limited to capturing what	
	models predict rather	

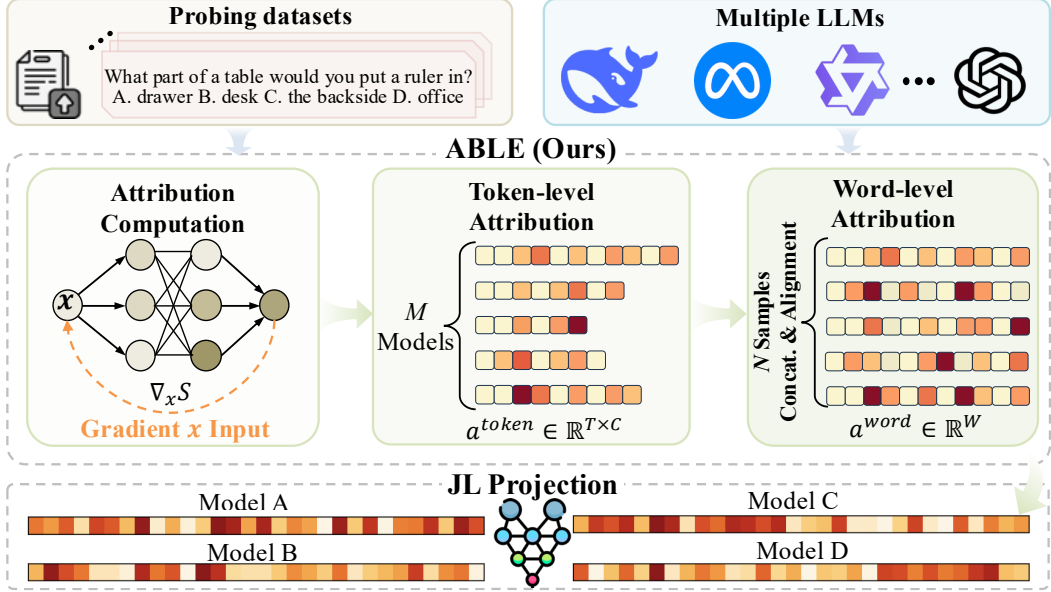


Figure 2: Overview of the ABLÉ pipeline. Given a model, we compute Gradient \times Input attributions on pre-defined probing datasets, align token-level scores to a unified word vocabulary, and project the result to a compact low-dimensional embedding. JL: Johnson–Lindenstrauss.

than how they process inputs. ABLÉ bridges these two paradigms by using feature attribution as the representation basis.

2.2 Feature Attribution Methods

Feature attribution methods quantify the contribution of each input feature to a model’s output, providing interpretable explanations for model predictions. White-box methods leverage gradient information, including Saliency Maps (Simonyan et al., 2013), Integrated Gradients (Sundararajan et al., 2017), and Layer-wise Relevance Propagation (Bach et al., 2015; Samek et al., 2021). Black-box methods probe models through input perturbations without requiring gradient access, such as LIME (Ribeiro et al., 2016) which fits local linear models, and SHAP (Lundberg and Lee, 2017) which applies game-theoretic Shapley values.

In the context of LLMs, attribution has been widely used for sample-level explanations, such as identifying salient tokens in classification (Liu and Avci, 2019) or generation tasks (Zhao and Shan, 2024). We extend this by aggregating sample-level attributions to construct model representations in interpretability space.

3 Methods

This section describes the ABLÉ representation construction pipeline (see Fig. 2). Given a collection of M language models and a pre-defined

probing dataset, we produce a compact embedding $\mathbf{a}_m \in \mathbb{R}^K$ for each model m . The pipeline consists of three stages: (1) attribution computation, (2) cross-tokenizer alignment, and (3) dimensionality reduction.

3.1 Probing Data Design

We adopt multi-choice questions (MCQ) from test set of standard benchmarks as our probing format. Each sample consists of a question q and a set of candidate options $\{o_1, o_2, \dots, o_C\}$, where C is typically 4–5. This design offers two key advantages. First, MCQ provide *deterministic attribution targets*: the log-probability of each option sequence given the question, $\log P_m(o_c | q)$, serves as a well-defined scalar output for gradient computation. Second, by computing attributions for all options, we implicitly capture how the model distinguishes correct from incorrect answers, a signal that reflects the model’s cognitive strategy rather than merely its final prediction.

3.2 Attribution Computation

Among various attribution methods, we adopt Gradient \times Input (GI) (Ancona et al., 2019, 2017; Nielsen et al., 2022) for two reasons: (1) computational efficiency, as GI requires only a single forward-backward pass per sample; and (2) theoretical tractability, as we prove in Appendix A that GI yields a stable, distance-preserving embedding.

For each model m and each sample, we compute the attribution of every question token to each option’s log-probability. The attribution score for the t -th question token with respect to option o_c is:

$$\alpha_t^{(c)} = \langle \mathbf{e}_t, \nabla_{\mathbf{e}_t} S_c \rangle, \\ S_c = \sum_{j=1}^{|o_c|} \log P_m(o_c^{(j)} | q, o_c^{(<j)}) \quad (1)$$

where \mathbf{e}_t denotes the embedding of the t -th question token, and S_c is the total log-probability of option o_c .

For a question with T tokens and C options, this yields an attribution matrix $\mathbf{A}^{\text{token}} \in \mathbb{R}^{T \times C}$. To obtain a single attribution vector per sample, we flatten this matrix into $\mathbb{R}^{T \cdot C}$. Concatenating across all N samples in the probing dataset produces a model-level token attribution vector $\mathbf{a}_m^{\text{token}} \in \mathbb{R}^{D_{\text{token}}}$, where $D_{\text{token}} = \sum_{i=1}^N T_i \times C_i$.

3.3 Cross-Tokenizer Alignment

Different models employ different tokenizers, resulting in attribution vectors of incompatible dimensions (Kudo and Richardson, 2018). To enable cross-model comparison, we align all attributions to a unified word-level vocabulary through a two-step process.

Token-to-Character Mapping. For each token t in the flattened attribution vector $\mathbf{a}_m^{\text{token}}$, let a_t denote its attribution score. We identify its character span $[s_t, e_t)$ in the original text and uniformly distribute the attribution across all characters in this span:

$$a_{\text{char}}^{(i)} = \frac{a_t}{e_t - s_t}, \quad \forall i \in [s_t, e_t). \quad (2)$$

This produces a character-level attribution vector $\mathbf{a}_m^{\text{char}} \in \mathbb{R}^L$, where L is the total number of characters in the input text.

Character-to-Word Aggregation. Given the character-level attributions, we aggregate them into word-level representations. We segment the character sequence into words using whitespace as the delimiter. For a word w spanning characters $[s_w, e_w)$, its attribution is computed as the sum of its constituent characters:

$$a_{\text{word}}^{(w)} = \sum_{i=s_w}^{e_w-1} a_{\text{char}}^{(i)}. \quad (3)$$

Whitespace attributions are appended to the preceding word to ensure no attribution is lost.

Since all models process the same probing text, this procedure yields a word-level attribution vector $\mathbf{a}_m^{\text{word}} \in \mathbb{R}^W$ for each model, where W is the total number of words in the probing corpus. This alignment addresses the tokenizer heterogeneity problem.

3.4 Dimensionality Reduction via Random Projection

The word-level attribution vector $\mathbf{a}_m^{\text{word}}$ can be high-dimensional ($W \sim 10^5$). To obtain a compact representation, we apply the Johnson-Lindenstrauss (JL) random projection (Johnson et al., 1984). Specifically, we sample a random matrix $\mathbf{R} \in \mathbb{R}^{K \times W}$ with independent and identically distributed entries drawn from $\mathcal{N}(0, 1/K)$, and compute:

$$\mathbf{a}_m = \mathbf{R} \mathbf{a}_m^{\text{word}} \in \mathbb{R}^K. \quad (4)$$

By the JL lemma, pairwise distances are preserved up to a multiplicative factor $(1 \pm \epsilon)$ with high probability, provided $K = O(\epsilon^{-2} \log M)$. In practice, we determined the optimal value of K via downstream task performance (see Appendix B). The resulting \mathbf{a}_m is the final ABLE representation for model m .

4 Experiments

We evaluate ABLE through a series of experiments designed to assess both its structural validity and practical utility. We first verify that ABLE captures meaningful model relationships through lineage reconstruction and relation prediction (Sections 4.2 and 4.3). We then demonstrate ABLE’s practical value in model routing and benchmark score prediction tasks (Sections 4.4 and 4.5).

4.1 Experimental Setup

Datasets. We construct a balanced probing dataset \mathcal{D} by randomly sampling 200 instances from each of six benchmarks: ARC-Challenge (Clark et al., 2018), Winogrande (Sakaguchi et al., 2021), MMLU (Hendrycks et al., 2021b,a), Hellaswag (Zellers et al., 2019), GPQA (Rein et al., 2024), and CommonsenseQA (Talmor et al., 2019), yielding $N = 1,200$ samples in total. This combination ensures diverse coverage across reasoning and knowledge domains. For model routing experiments, we additionally use the evaluation dataset from EmbedLLM (Zhuang et al., 2024).

Models. We evaluate $M = 239$ open-source LLMs spanning parameter scales from 70M to

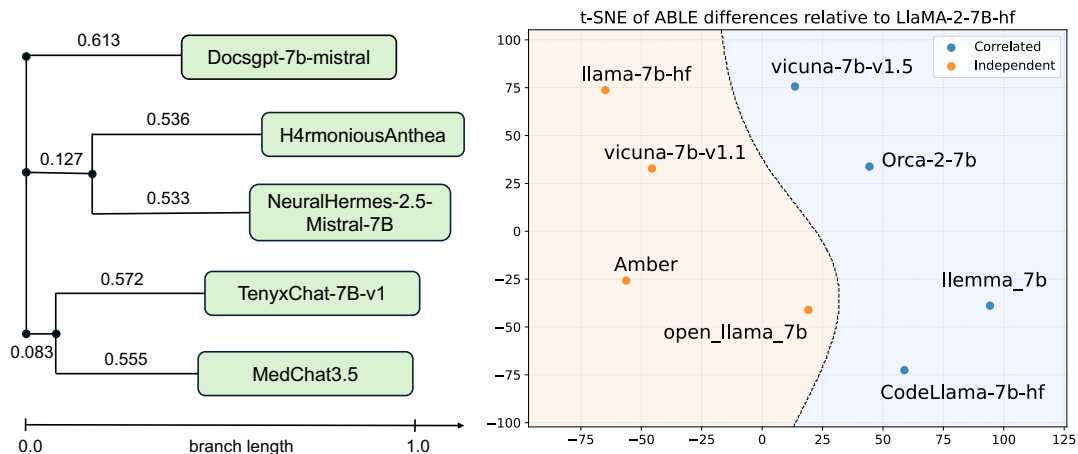


Figure 3: Lineage verification experiments. Left: Mistral family phylogenetic tree reconstructed from ABLE embeddings. Right: Llama-2 family models in ABLE space, showing clear separation between correlated and independent pairs.

70B. This collection encompasses foundation base models and instruction-tuned variants, as well as domain-specific models tailored for mathematics, coding, and medicine. More model details are provided in Appendix E.

4.2 Model Lineage and Atlas

To evaluate whether ABLE representations capture meaningful structural information and evolutionary history of LLMs, we present experiments at two scales: first, small-scale lineage verification on documented model relations; second, constructing a comprehensive atlas of the full model collection.

Lineage Verification on Mistral and Llama-2 Families. To validate that ABLE captures model lineage, we conduct two small-scale verification experiments. For the Mistral family, we select five leaf models with publicly disclosed fine-tuning histories (DocsGPT, NeuralHermes-2.5, TenyxChat-7B, MedChat3.5, and H4rmoniousAnthea) and construct an unrooted tree using the Neighbor-Joining method (Trees, 1987) on cosine distances calculated on ABLE. As shown in Fig. 3 (left), the results aligns with the ground-truth lineage disclosed by model creators (arc53, 2023; mlabonne, 2024; Tenyx, 2024; Ullah, 2024; Vallego, 2024). For the Llama-2 family, we examine eight models: four are correlated with Llama-2-7b-hf (Vicuna-7b-v1.5, CodeLlama-7b-hf, Llemma-7b, Orca-2-7b) and four are independent (Llama-7b-hf, Vicuna-7b-v1.1, Amber, Open-Llama-7b), following the relationship annotations by Zhu et al. (2025). As shown in Fig. 3 (right), correlated and independent models are clearly separated in ABLE space,

with an SVM decision boundary achieving perfect classification.

Visualizing the Global Model Map. We then project the ABLE features into a two-dimensional space using t-SNE (Maaten and Hinton, 2008) to construct a Model Map. From Fig. 4, three major insights can be inferred. First, LLMs tend to aggregate based on their respective families, such as Llama, Qwen, and Mistral family. Second, fine-tuned models tend to cluster around corresponding base models. A prime example is the Llama 3 family (red) in the central region, where multiple LoRA-tuned variants cluster tightly around the base Llama-3-8B model. Third, fine-tuned models may also deviate from their families to align with models sharing similar tasks. For instance, Llemma-34B and DeepSeek-Math-7B-Instruct cluster together due to their shared focus on mathematical reasoning, while LoRA-8B-Code and CodeGemma-2B converge due to code-related fine-tuning (names labeled in red in Fig. 4). This indicates that ABLE can capture genealogical proximity and functional affinity among heterogeneous LLMs.

Constructing Cross-Family Phylogenetic Trees. While the Model Map provides a continuous two-dimensional projection, we further construct a hierarchical phylogenetic tree to explicitly reveal evolutionary relationships across model families. We select representative models from ten families (Mistral, Bloom, LLaMA-1, LLaMA-2, LLaMA-3, Yi, Qwen2, Qwen3, MPT, and Pythia) and compute pairwise cosine distances from their ABLE embeddings. We then apply

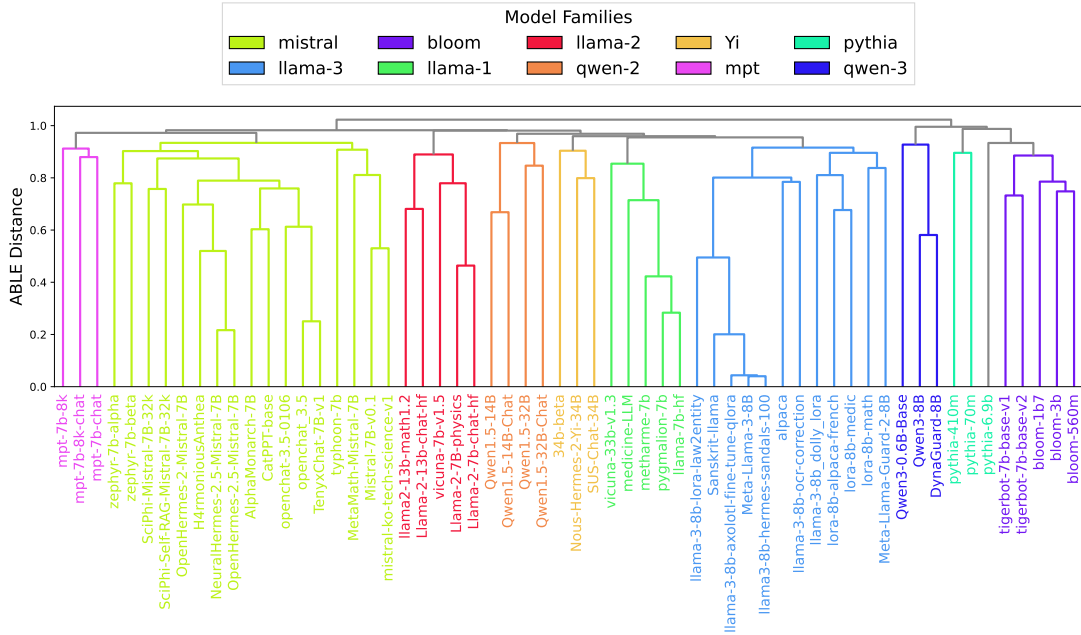


Figure 5: Phylogenetic tree constructed from ABLÉ embeddings using hierarchical clustering with cosine distance. Branches are colored by model family; gray branches indicate cross-family merging points.

Table 1: Model relation prediction performance of ABLÉ and four baselines (Values are mean \pm standard deviation over 5 random seeds). Bold: the best performance.

Method	Accuracy	Precision	Recall	F1	AUC
Random	0.489 \pm 0.124	0.367 \pm 0.093	0.733 \pm 0.186	0.489 \pm 0.124	0.550 \pm 0.139
Greedy	0.674 \pm 0.096	0.527 \pm 0.313	0.222 \pm 0.157	0.309 \pm 0.203	0.561 \pm 0.110
PhyloLM	0.704 \pm 0.141	0.596 \pm 0.182	0.844 \pm 0.169	0.667 \pm 0.090	0.893 \pm 0.059
Log-Likelihood	0.721 \pm 0.041	0.586 \pm 0.064	0.861\pm0.150	0.692 \pm 0.070	0.891 \pm 0.053
ABLE(Ours)	0.867\pm0.042	0.836\pm0.077	0.756 \pm 0.093	0.790\pm0.065	0.906\pm0.057

computes log-likelihood vectors on a fixed dataset as model features; and (4) PhyloLM (Yax et al., 2024), an output-based method uses output similarity distances to other LLMs as a signature vector.

As shown in Table 1, ABLÉ achieves the best performance on Accuracy, Precision, F1, and AUC. Notably, ABLÉ and output-based methods exhibit different precision-recall trade-offs: Log-Likelihood attains the highest recall but with lower precision, whereas ABLÉ achieves substantially higher precision at the cost of moderate recall. This reflects that ABLÉ requires stronger evidence before classifying a model pair as correlated, resulting in more confident positive predictions. This conservative behavior is valuable for intellectual property protection, where falsely accusing an independent model of infringement can lead to legal disputes and reputation damage (Van Wyk et al., 2023).

4.4 Model Routing

To rigorously evaluate the effectiveness of ABLÉ features in model routing, we adopt the matrix factorization framework, dataset, and evaluation metrics from EmbedLLM (Zhuang et al., 2024). This architecture uses a learnable linear projection layer to map frozen question embeddings into the model feature space, predicting the probability that a model answers a question correctly via element-wise interaction. Unlike EmbedLLM, which optimizes model embeddings end-to-end, we keep all ABLÉ features frozen throughout training and only learn the semantic alignment from question embeddings to the model space.

As shown in Table 2, the router based on frozen ABLÉ features achieves 67.6% accuracy, slightly higher than the fully trained EmbedLLM router. While the accuracy gain is modest, the key advantage of our approach lies in training-free scalabil-

Table 2: Model routing performance of ABLE against three baselines (Values are mean \pm standard deviation over 5 random seeds). ABLE router is competitive with EmbedLLM router.

	Random	Single-Best	EmbedLLM	ABLE
Router Accuracy	0.413 \pm 0.101	0.605 \pm 0.000	0.665 \pm 0.003	0.676 \pm 0.001

Table 3: Benchmark score prediction performance. We report the Pearson (r) and Spearman (ρ) correlation coefficients between the ground-truth and the scores predicted by ridge regression using ABLE.

	ARC	HellaSwag	MMLU	TruthfulQA	WinoGrande
Pearson r	0.8781	0.7764	0.8554	0.9127	0.8374
Spearman ρ	0.8898	0.8157	0.8637	0.8222	0.8602

ity. The ABLE router can accommodate an arbitrary number of models, and when new models are added to the pool, only their ABLE features need to be computed without retraining. In contrast, EmbedLLM learns model embeddings end-to-end, necessitating retraining whenever new models are introduced.

4.5 Benchmark Score Prediction

In this section, we investigate whether ABLE representations can predict LLM performance on standard benchmarks. Using scores from five core benchmarks on the Hugging Face Open LLM Leaderboard (Fourrier et al., 2024; Myrzakhan et al., 2024; Beeching et al., 2023) as labels and ABLE features as input, we evaluate the predictive ability of a ridge regression model (McDonald, 2009) via leave-one-out cross-validation.

As illustrated in Table 3 and Fig. 6, ABLE features demonstrate strong correlations with ground-truth benchmark scores, achieving Spearman ρ values between 0.81 and 0.89. This high consistency indicates that model hierarchies predicted solely from ABLE representations closely mirror those derived from standard evaluations. These findings suggest that ABLE effectively encodes model capabilities within a linear subspace. Specifically, for large-scale comparisons where relative ranking is prioritized over precise scoring, ABLE serves as a computationally efficient surrogate for resource-intensive benchmarks.

5 Conclusion

In this work, we aim to mitigate the growing transparency crisis in the LLM ecosystem by proposing ABLE, a framework that bridges the gap between internal mechanism analysis and external

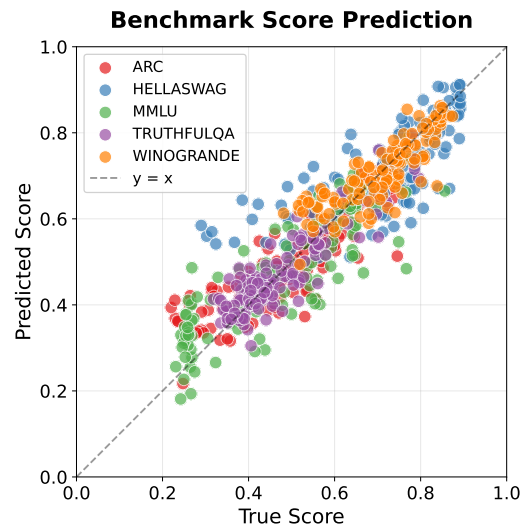


Figure 6: Scatter plots showing predicted benchmark score against ground truth scores on five benchmark using ABLE features.

behavioral probing. By leveraging the interpretability space, ABLE captures the intrinsic cognitive patterns of models, effectively distinguishing between superficial behavioral mimicry and genuine mechanistic similarity. We also provide theoretical analysis indicating that ABLE constitutes a Lipschitz-continuous mapping, thereby supporting the stability of the representation.

While extracting gradient-based features incurs a computational cost, this process represents a one-time investment that yields reusable embeddings. Once computed, ABLE facilitates efficient downstream applications, including lineage tracking, model routing, and performance estimation. We hope this work offers a valuable perspective for mapping and understanding the evolving landscape of LLMs.

520 Limitations

521 **Computational Efficiency.** Computing ABLE
522 features requires one forward and one backward
523 pass per sample, which incurs certain time costs.
524 However, ABLE embeddings are computed once
525 per model and can be reused across all downstream
526 tasks. For large-scale model comparison, this one-
527 time extraction cost is lower than running full
528 benchmark evaluations on every model.

529 **Attribution Method.** To balance accuracy and
530 efficiency, we only employ Gradient \times Input (GI)
531 as the attribution method. We have not explored
532 more sophisticated techniques such as Integrated
533 Gradient (Kapishnikov et al., 2021; Qi et al., 2020)
534 or SmoothGrad (Smilkov et al., 2017), which may
535 provide higher fidelity attributions at the cost of
536 increased computation.

537 **Closed-Source Models.** ABLE requires access
538 to model gradients, making it inapplicable to
539 closed-source APIs. Nevertheless, the concept of
540 mapping LLMs in interpretability space opens av-
541 enues for future work using perturbation-based at-
542 tribution methods (Ivanovs et al., 2021), which
543 could extend ABLE to black-box settings.

544 **Sample Sensitivity.** ABLE is computed on a fi-
545 nite set of probing samples, raising the concern
546 that the representation may be sensitive to the spe-
547 cific choice of samples. However, as shown in
548 Appendix D, ABLE embeddings exhibit high sta-
549 bility across different random subsamples of the
550 probing dataset, suggesting that ABLE captures
551 robust structural properties of models rather than
552 being sensitive to individual input samples.

553 **Alignment Granularity.** Our character-level
554 alignment uniformly distributes token attributions
555 across characters, which may lose subword-level
556 information. While this simple approach proves
557 effective for the English text benchmarks evaluated
558 in this work, alternative schemes such as frequency-
559 weighted character assignment or byte-level align-
560 ment may be more suitable for code or multilingual
561 inputs. We leave systematic exploration of align-
562 ment strategies to future work.

563 References

564 Marco Ancona, Enea Ceolini, Cengiz Öztireli, and
565 Markus Gross. 2017. Towards better understand-
566 ing of gradient-based attribution methods for deep
567 neural networks. *arXiv preprint arXiv:1711.06104*.

Marco Ancona, Enea Ceolini, Cengiz Öztireli, and 568
Markus Gross. 2019. Gradient-based attribution 569
methods. In *Explainable AI: Interpreting, explain- 570*
ing and visualizing deep learning, pages 169–191. 571
Springer. 572

arc53. 2023. Docsgpt-7b-mistral. [https:// 573](https://huggingface.co/Arc53/docsgpt-7b-mistral)
huggingface.co/Arc53/docsgpt-7b-mistral. 574

Sebastian Bach, Alexander Binder, Grégoire Montavon, 575
Frederick Klauschen, Klaus-Robert Müller, and Wo- 576
jciech Samek. 2015. On pixel-wise explanations 577
for non-linear classifier decisions by layer-wise rele- 578
vance propagation. *PLoS one*, 10(7):e0130140. 579

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, 580
Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei 581
Huang, and 1 others. 2023. Qwen technical report. 582
arXiv preprint arXiv:2309.16609. 583

Kristian González Barman, Nathan Wood, and Pawel 584
Pawlowski. 2024. Beyond transparency and explain- 585
ability: on the need for adequate and contextualized 586
user guidelines for llm use. *Ethics and Information 587*
Technology, 26(3):47. 588

Edward Beeching, Clémentine Fourier, Nathan 589
Habib, Sheon Han, Nathan Lambert, Nazneen 590
Rajani, Omar Sanseviero, Lewis Tunstall, and 591
Thomas Wolf. 2023. Open llm leaderboard 592
(2023-2024). [https://huggingface.co/ 593](https://huggingface.co/spaces/open-llm-leaderboard-old/open_llm_leaderboard)
[spaces/open-llm-leaderboard-old/open_ 594](https://huggingface.co/spaces/open-llm-leaderboard-old/open_llm_leaderboard)
[llm_leaderboard](https://huggingface.co/spaces/open-llm-leaderboard-old/open_llm_leaderboard). 595

Pengzhou Cheng, Zongru Wu, Tianjie Ju, Wei Du, and 596
Zhuosheng Zhang Gongshen Liu. 2024. Transferring 597
backdoors between large language models by knowl- 598
edge distillation. *arXiv preprint arXiv:2408.09878*. 599

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, 600
Ashish Sabharwal, Carissa Schoenick, and Oyvind 601
Tafjord. 2018. Think you have solved question an- 602
swering? try arc, the ai2 reasoning challenge. *arXiv 603*
preprint arXiv:1803.05457. 604

Clémentine Fourier, Nathan Habib, Alina Lozovskaya, 605
Konrad Szafer, and Thomas Wolf. 2024. Open 606
llm leaderboard v2. [https://huggingface. 607](https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard)
[co/spaces/open-llm-leaderboard/open_llm_ 608](https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard)
[leaderboard](https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard). 609

Yaru Hao, Haoyu Song, Li Dong, Shaohan Huang, 610
Zewen Chi, Wenhui Wang, Shuming Ma, and Furu 611
Wei. 2022. Language models are general-purpose 612
interfaces. *arXiv preprint arXiv:2206.06336*. 613

Dan Hendrycks, Collin Burns, Steven Basart, Andrew 614
Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 615
2021a. Aligning ai with shared human values. *Pro- 616*
ceedings of the International Conference on Learning 617
Representations (ICLR). 618

Dan Hendrycks, Collin Burns, Steven Basart, Andy 619
Zou, Mantas Mazeika, Dawn Song, and Jacob Stein- 620
hardt. 2021b. Measuring massive multitask language 621
understanding. *Proceedings of the International Con- 622*
ference on Learning Representations (ICLR). 623

624	Eliahu Horwitz, Nitzan Kurer, Jonathan Kahana, Liel Amar, and Yedid Hoshen. 2025. We should chart an atlas of all the world’s models. <i>arXiv preprint arXiv:2503.10633</i> .	Gary C McDonald. 2009. Ridge regression. <i>Wiley Interdisciplinary Reviews: Computational Statistics</i> , 1(1):93–100.	677
625			678
626			679
627			
628	Maksims Ivanovs, Roberts Kadikis, and Kaspars Ozols. 2021. Perturbation-based methods for explaining deep neural networks: A survey. <i>Pattern Recognition Letters</i> , 150:228–234.	Peter S Menell. 1988. An analysis of the scope of copyright protection for application programs. <i>Stan. L. Rev.</i> , 41:1045.	680
629			681
630			682
631		mlabonne. 2024. Neuralhermes-2.5-mistral-7b. https://huggingface.co/mlabonne/neuralHermes-2.5-mistral-7B .	683
632	Shashank Mohan Jain. 2022. Hugging face. In <i>Introduction to transformers for NLP: With the hugging face library and models to solve problems</i> , pages 51–67. Springer.		684
633			685
634		Aidar Myrzakhan, Sondos Mahmoud Bsharat, and Zhiqiang Shen. 2024. Open-llm-leaderboard: From multi-choice to open-style questions for llms evaluation, benchmark, and arena. <i>arXiv preprint arXiv:2406.07545</i> .	686
635			687
636	William B Johnson, Joram Lindenstrauss, and 1 others. 1984. Extensions of lipschitz mappings into a hilbert space. <i>Contemporary mathematics</i> , 26(189-206):1.		688
637			689
638			690
639	Andrei Kapishnikov, Subhashini Venugopalan, Besim Avci, Ben Wedin, Michael Terry, and Tolga Bolukbasi. 2021. Guided integrated gradients: An adaptive path method for removing noise. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 5050–5058.	Giang Nguyen, Daeyoung Kim, and Anh Nguyen. 2021. The effectiveness of feature attribution methods and its correlation with automatic evaluation scores. <i>Advances in Neural Information Processing Systems</i> , 34:26422–26436.	691
640			692
641			693
642			694
643			695
644		Ian E Nielsen, Dimah Dera, Ghulam Rasool, Ravi P Ramachandran, and Nidhal Carla Bouaynaya. 2022. Robust explainability: A tutorial on gradient-based attribution methods for deep neural networks. <i>IEEE Signal Processing Magazine</i> , 39(4):73–84.	696
645	Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, and 1 others. 2018. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In <i>International conference on machine learning</i> , pages 2668–2677. PMLR.		697
646			698
647			699
648			700
649		Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez, M Waleed Kadous, and Ion Stoica. 2024. Routellm: Learning to route llms with preference data. <i>arXiv preprint arXiv:2406.18665</i> .	701
650			702
651	Ramaravind Kommiya Mothilal, Divyat Mahajan, Chenhao Tan, and Amit Sharma. 2021. Towards unifying feature attribution and counterfactual explanations: Different means to the same end. In <i>Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society</i> , pages 652–663.		703
652			704
653			705
654		Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Reichart, Idan Szpektor, Hadas Kotek, and Yonatan Belinkov. 2024. Llms know more than they show: On the intrinsic representation of llm hallucinations. <i>arXiv preprint arXiv:2410.02707</i> .	706
655			707
656			708
657	Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. <i>arXiv preprint arXiv:1808.06226</i> .		709
658			710
659		Momose Oyama, Hiroaki Yamagiwa, Yusuke Takase, and Hidetoshi Shimodaira. 2025. Mapping 1,000+ language models via the log-likelihood vector. <i>arXiv preprint arXiv:2502.16173</i> .	711
660			712
661	Frederick Liu and Besim Avci. 2019. Incorporating priors with feature attribution on text classification . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6274–6283, Florence, Italy. Association for Computational Linguistics.		713
662			714
663		Zhongang Qi, Saeed Khorram, and Fuxin Li. 2020. Visualizing deep networks by optimizing with integrated gradients. In <i>AAAI</i> , volume 34, pages 11890–11898.	715
664			716
665			717
666			718
667	Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. <i>Advances in neural information processing systems</i> , 30.	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2024. GPQA: A graduate-level google-proof q&a benchmark . In <i>First Conference on Language Modeling</i> .	719
668			720
669			721
670	Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. <i>Journal of machine learning research</i> , 9(Nov):2579–2605.		722
671			723
672		Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?" explaining the predictions of any classifier. In <i>Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining</i> , pages 1135–1144.	724
673	Marios Mattheakis, Pavlos Protopapas, David Sondak, Marco Di Giovanni, and Efthimios Kaxiras. 2019. Physical symmetries embedded in neural networks. <i>arXiv preprint arXiv:1904.08991</i> .		725
674			726
675			727
676			728
			729

730	Denis Rothman. 2022. <i>Transformers for Natural Language Processing: Build, train, and fine-tune deep neural network architectures for NLP with Python, Hugging Face, and OpenAI's GPT-3, ChatGPT, and GPT-4</i> . Packt Publishing Ltd.	784
731		785
732		
733		
734		
735	Jungwoo Ryoo, Syed Rizvi, William Aiken, and John Kissell. 2013. Cloud security auditing: challenges and emerging approaches. <i>IEEE Security & Privacy</i> , 12(6):68–74.	788
736		789
737		790
738		791
739	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavathula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106.	792
740		793
741		794
742		795
743	Wojciech Samek, Grégoire Montavon, Sebastian Lapuschkin, Christopher J Anders, and Klaus-Robert Müller. 2021. Explaining deep neural networks and beyond: A review of methods and applications. <i>Proceedings of the IEEE</i> , 109(3):247–278.	796
744		
745		
746		
747		
748	Shuo Shao, Yiming Li, Hongwei Yao, Yiling He, Zhan Qin, and Kui Ren. 2024. Explanation as a watermark: Towards harmless and multi-bit model ownership verification via watermarking feature attribution. <i>arXiv preprint arXiv:2405.04825</i> .	797
749		798
750		799
751		800
752		
753	Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. <i>arXiv preprint arXiv:1312.6034</i> .	801
754		802
755		803
756		804
757		805
758	Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. 2017. Smoothgrad: removing noise by adding noise. <i>arXiv preprint arXiv:1706.03825</i> .	806
759		807
760		808
761		809
762		810
763		
764	Pascal Sturmfels, Scott Lundberg, and Su-In Lee. 2020. Visualizing the impact of feature attribution baselines. <i>Distill</i> , 5(1):e22.	811
765		812
766		813
767		814
768	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In <i>International conference on machine learning</i> , pages 3319–3328. PMLR.	815
769		816
770		817
771		818
772		819
773		
774		
775		
776		
777	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	820
778		821
779		822
780		823
781		824
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		
810		
811		
812		
813		
814		
815		
816		
817		
818		
819		
820		
821		
822		
823		
824		
825		
826		
827		
828		
829		
830		
831		
832		
833		
834		
835		
836		
837		

838 Yilun Zhou, Serena Booth, Marco Tulio Ribeiro, and
 839 Julie Shah. 2022. Do feature attribution methods
 840 correctly attribute features? In *Proceedings of the*
 841 *AAAI conference on artificial intelligence*, volume 36,
 842 pages 9623–9633.

843 Sally Zhu, Ahmed Ahmed, Rohith Kudithipudi, and
 844 Percy Liang. 2025. Independence tests for language
 845 models. *arXiv preprint arXiv:2502.12292*.

846 Richard Zhuang, Tianhao Wu, Zhaojin Wen, Andrew Li,
 847 Jiantao Jiao, and Kannan Ramchandran. 2024. Em-
 848 bedllm: Learning compact representations of large
 849 language models. *arXiv preprint arXiv:2410.02223*.

850 A Theoretical Analysis

851 In this section, we provide a theoretical justifica-
 852 tion that ABLE constitutes a valid and mechanism-
 853 aware embedding of large language models. Our
 854 analysis abstracts away implementation details and
 855 focuses on the relationship between model param-
 856 eters and attribution-based representations. Sec-
 857 tion A.1 defines ABLE. Section A.2 establishes
 858 that ABLE is a stable feature map induced by
 859 model parameters. Section A.3 proves that the
 860 low-dimensional embedding preserves inter-model
 861 distances via random projection. Section A.4 pro-
 862 vides finite-sample concentration guarantees.

863 A.1 Preliminaries and Definitions

864 Let a language model be parameterized by $\theta \in \mathbb{R}^P$.
 865 For an input representation $x \in \mathbb{R}^d$ and an output
 866 class y , denote the log-probability function by

$$867 s_{\theta}^{(y)}(x) \triangleq \log p_{\theta}(y | x),$$

868 which we assume to be differentiable with respect
 869 to x almost everywhere.

870 **Attribution-based sensitivity.** For each input x ,
 871 we define the attribution score as the inner product
 872 between the input and its gradient:

$$873 a_{\theta}(x) \triangleq \langle x, \nabla_x s_{\theta}(x) \rangle \in \mathbb{R},$$

874 where $s_{\theta}(x)$ denotes the log-probability function
 875 stacking $s_{\theta}^{(y)}(x)$ over output classes.

876 **Model representation.** Let $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$ be a
 877 deterministic aggregation function. Given a probe
 878 distribution \mathcal{D} over inputs, we define the (pre-
 879 projection) ABLE representation as

$$880 \Phi(\theta) \triangleq \mathbb{E}_{x \sim \mathcal{D}}[\phi(a_{\theta}(x))] \in \mathbb{R}^D.$$

881 Finally, the ABLE embedding is obtained via a
 882 random projection

$$883 \Psi(\theta) \triangleq \mathbf{R} \Phi(\theta) \in \mathbb{R}^K,$$

884 where $\mathbf{R} \in \mathbb{R}^{K \times D}$ is a Johnson–Lindenstrauss
 885 random matrix.

886 A.2 The Parameter Stability of ABLE

887 The objective of this section is to demonstrate the
 888 parameter stability of ABLE, meaning that the em-
 889 bedding varies continuously with changes in model
 890 parameters. Our main line of reasoning is to first
 891 prove that ABLE is a Lipschitz continuous map-
 892 ping under standard regularity assumptions, and
 893 subsequently demonstrate that Transformer-based
 894 LLMs satisfy these assumptions.

895 **Theorem A.1** (Parameter Stability of ABLE). *As-*
 896 *sume the following regularity conditions hold:*

- 897 1. *The gradient field of the log-probability func-*
 898 *tion, $\nabla_x s_{\theta}(x)$, is M -Lipschitz continuous in*
 899 *θ : $\|\nabla_x s_{\theta}(x) - \nabla_x s_{\theta'}(x)\|_2 \leq M\|\theta - \theta'\|_2$.*
- 900 2. *For all x in the support of \mathcal{D} , $s_{\theta}(x)$ is differ-*
 901 *entiable with respect to x almost everywhere.*
- 902 3. *The aggregation function ϕ is L_{ϕ} -Lipschitz,*
 903 *and inputs are bounded on average:*
 904 $\mathbb{E}_{x \sim \mathcal{D}}\|x\|_2 \leq B$.

905 *Then the ABLE mapping $\Phi : \theta \mapsto \Phi(\theta)$ is Lipschitz*
 906 *continuous:*

$$907 \|\Phi(\theta) - \Phi(\theta')\|_2 \leq L\|\theta - \theta'\|_2,$$

908 where $L = BML_{\phi}$.

909 *Proof.* Using the Cauchy-Schwarz inequality and
 910 the assumption on $\nabla_x s_{\theta}(x)$:

$$911 |a_{\theta}(x) - a_{\theta'}(x)| = |\langle x, \nabla_x s_{\theta}(x) - \nabla_x s_{\theta'}(x) \rangle|$$

$$912 \leq \|x\|_2 \|\nabla_x s_{\theta}(x) - \nabla_x s_{\theta'}(x)\|_2$$

$$913 \leq \|x\|_2 M \|\theta - \theta'\|_2.$$

914 Applying the Lipschitz continuity of ϕ and taking
 915 expectation over $x \sim \mathcal{D}$:

$$916 \|\Phi(\theta) - \Phi(\theta')\|_2$$

$$917 \leq \mathbb{E}_{x \sim \mathcal{D}}[L_{\phi}|a_{\theta}(x) - a_{\theta'}(x)|]$$

$$918 \leq L_{\phi} \mathbb{E}_{x \sim \mathcal{D}}[\|x\|_2] \cdot M \cdot \|\theta - \theta'\|_2$$

$$919 \leq BML_{\phi} \|\theta - \theta'\|_2.$$

920 **Theorem A.1** provides the formal guarantee for
 921 ABLE’s parameter stability. The Lipschitz inequal-
 922 ity $\|\Phi(\theta) - \Phi(\theta')\|_2 \leq L\|\theta - \theta'\|_2$ ensures that
 923 the mapping does not diverge: finite differences
 924

in model parameters result in strictly bounded differences in the embedding space. This property is crucial for a reliable representation, as it guarantees that models with similar internal mechanisms (proximal parameters) are mapped to nearby points in the vector space, preventing chaotic behavior where minor weight variations could lead to disparate embeddings.

The validity of Theorem A.1 relies on three assumptions. Assumptions 2 and 3 are satisfied by design in our setting:

- **Assumption 2 (Differentiability):** While text is discrete, the model operates on continuous word embeddings. Standard embedding lookup layers and subsequent smooth transformations ensure differentiability with respect to the input representation x almost everywhere.
- **Assumption 3 (Boundedness):** The aggregation function ϕ is linear (averaging) and thus 1-Lipschitz. Input embeddings belong to a finite vocabulary set \mathcal{V} , so $\|x\|_2$ is strictly bounded by $B = \max_{v \in \mathcal{V}} \|v\|_2$.

The remaining Assumption 1, which requires the model’s gradient field to be Lipschitz continuous with respect to its parameters, is the most non-trivial condition. We now demonstrate that this assumption holds for Transformer-based Large Language Models under standard architectural constraints.

Proposition A.1 (Regularity of Transformer Gradient Fields). *Consider a Transformer-based language model with the following properties:*

1. *Activation functions (e.g., GELU, Swish) have bounded first and second derivatives: $|\sigma'(z)| \leq L_\sigma$ and $|\sigma''(z)| \leq L_{\sigma'}$ for all z .*
2. *Token embeddings are drawn from a finite vocabulary with bounded norms: $\|x\|_2 \leq B_x$.*
3. *The network has finite depth L and bounded weight matrices: $\|W_l\|_2 \leq B_W$ for all layers l .*

Then the gradient field $\nabla_x s_\theta(x)$ is M -Lipschitz continuous in θ for some constant $M > 0$.

Proof. We proceed by analyzing the Lipschitz dependence layer by layer.

Step 1: Linear layer. For $y = Wx$, the input gradient is $\nabla_x y = W^T$. For two parameter configurations W, W' :

$$\begin{aligned} \|\nabla_x y_W - \nabla_x y_{W'}\|_2 &= \|W^T - W'^T\|_2 \\ &= \|W - W'\|_2, \end{aligned}$$

which is 1-Lipschitz in W .

Step 2: Layer with activation. For $y = \sigma(Wx)$, the chain rule gives $\nabla_x y = W^T \cdot \text{diag}(\sigma'(Wx))$. When W changes to W' :

$$\begin{aligned} \|\nabla_x y_W - \nabla_x y_{W'}\|_2 &\leq \|W^T - W'^T\|_2 \cdot \|\text{diag}(\sigma'(Wx))\|_2 \\ &\quad + \|W'^T\|_2 \cdot \|\sigma'(Wx) - \sigma'(W'x)\|_2 \\ &\leq L_\sigma \|W - W'\|_2 \\ &\quad + B_W L_{\sigma'} \|x\|_2 \|W - W'\|_2 \\ &\leq (L_\sigma + B_W L_{\sigma'} B_x) \|W - W'\|_2. \end{aligned}$$

Step 3: Multi-layer composition. For an L -layer network $f = f_L \circ \dots \circ f_1$, the chain rule yields $\nabla_x f = \prod_{l=1}^L J_l$, where J_l is the Jacobian of layer l . Since each J_l is Lipschitz in θ_l with bounded spectral norm (due to bounded weights and the smoothness of softmax attention), perturbation analysis shows that the product $\nabla_x f$ remains Lipschitz in θ .

For Transformer architectures specifically, self-attention layers use softmax, whose Jacobian satisfies $\|\partial \text{softmax} / \partial z\|_2 \leq 1$. Combined with bounded query, key, and value projections, each attention layer contributes a bounded Lipschitz factor. Aggregating across all L layers yields the global constant M .

Consequently, ABLE constitutes a stable feature map as implied by its Lipschitz continuity, where small perturbations in model parameters result in bounded changes in the representation space. It is worth noting that this guarantee fundamentally relies on differentiability; architectures employing non-differentiable components (e.g., hard attention) may require alternative theoretical analysis.

A.3 Distance Preservation via Random Projection

We next show that the final ABLE embedding preserves inter-model geometry up to a small distortion.

Theorem A.2 (Johnson–Lindenstrauss Property). *Let $\mathcal{F} = \{\theta_1, \dots, \theta_m\}$ be a finite set of models and*

let $0 < \epsilon < 1$. If the projection dimension satisfies $K = O(\epsilon^{-2} \log m)$, then with high probability, for all $i, j \in \{1, \dots, m\}$, the projected distance $\|\Psi(\theta_i) - \Psi(\theta_j)\|_2$ lies within $(1 \pm \epsilon)$ of the original distance $\|\Phi(\theta_i) - \Phi(\theta_j)\|_2$.

Proof. Let $\Delta_{ij} = \Phi(\theta_i) - \Phi(\theta_j)$ be the difference vector between any two models. Only the linear projection $\Psi(\theta) = \mathbf{R}\Phi(\theta)$ comprises random variables, where entries of $\mathbf{R} \in \mathbb{R}^{K \times D}$ are i.i.d. Gaussian $\mathcal{N}(0, 1/K)$.

Consider a single pair with difference vector $v = \Delta_{ij}$. We are interested in the distribution of $\|\mathbf{R}v\|^2$. Let R_k be the k -th row of \mathbf{R} . Then the k -th component of the projected vector is $y_k = \langle R_k, v \rangle$. Since $R_k \sim \mathcal{N}(0, \frac{1}{K}I)$, the linear combination y_k is essentially a univariate Gaussian:

$$y_k \sim \mathcal{N}\left(0, \frac{\|v\|^2}{K}\right).$$

Consequently, the squared norm of the projection is:

$$\|\mathbf{R}v\|^2 = \sum_{k=1}^K y_k^2 = \frac{\|v\|^2}{K} \sum_{k=1}^K Z_k^2,$$

where $Z_k \sim \mathcal{N}(0, 1)$ are standard Gaussian variables. The sum $X = \sum_{k=1}^K Z_k^2$ follows a Chi-squared distribution with K degrees of freedom, denoted as χ_K^2 .

We use standard concentration bounds for the Chi-squared distribution (derived from the Moment Generating Function bound). For any $\epsilon \in (0, 1)$:

$$\mathbb{P}(X \geq K(1 + \epsilon)) \leq \exp\left(-\frac{K}{4}\epsilon^2\right),$$

$$\mathbb{P}(X \leq K(1 - \epsilon)) \leq \exp\left(-\frac{K}{4}\epsilon^2\right).$$

Combining these, the probability that the squared length is distorted by more than ϵ is:

$$\begin{aligned} \mathbb{P}\left(\left|\|\mathbf{R}v\|^2 - \|v\|^2\right| > \epsilon\|v\|^2\right) \\ &= \mathbb{P}\left(\left|\frac{X}{K} - 1\right| > \epsilon\right) \\ &\leq 2 \exp\left(-\frac{K\epsilon^2}{4}\right). \end{aligned}$$

We apply this to the set of all pairwise differences $\mathcal{V} = \{\Delta_{ij} \mid 1 \leq i < j \leq m\}$, which has cardinal-

ity $|\mathcal{V}| = \binom{m}{2} < m^2$. Using the union bound:

$$\begin{aligned} P_{\text{fail}} &= \mathbb{P}\left(\exists i, j : \left|\|\Psi(\theta_i) - \Psi(\theta_j)\|^2 - \|\Delta_{ij}\|^2\right| \right. \\ &\quad \left. > \epsilon\|\Delta_{ij}\|^2\right) \\ &\leq \sum_{i < j} \mathbb{P}\left(\left|\|\mathbf{R}\Delta_{ij}\|^2 - \|\Delta_{ij}\|^2\right| > \epsilon\|\Delta_{ij}\|^2\right) \\ &\leq m^2 \cdot 2 \exp\left(-\frac{K\epsilon^2}{4}\right). \end{aligned}$$

To ensure $P_{\text{fail}} \leq \delta$, we require:

$$2m^2 \exp\left(-\frac{K\epsilon^2}{4}\right) \leq \delta \implies K \geq \frac{4}{\epsilon^2} \ln \frac{2m^2}{\delta}.$$

Thus, choosing $K = O(\epsilon^{-2} \log m)$ suffices to preserve all pairwise distances with high probability.

We next show that the final ABLE embedding preserves inter-model geometry up to a small distortion.

Theorem A.3 (Johnson–Lindenstrauss Property).

Let $\mathcal{F} = \{\theta_1, \dots, \theta_m\}$ be a finite set of models and let $0 < \epsilon < 1$. If the projection dimension satisfies $K = O(\epsilon^{-2} \log m)$, then with high probability, for all $i, j \in \{1, \dots, m\}$, the projected distance $\|\Psi(\theta_i) - \Psi(\theta_j)\|_2$ lies within $(1 \pm \epsilon)$ of the original distance $\|\Phi(\theta_i) - \Phi(\theta_j)\|_2$.

Consequence. Theorem A.3 implies that ABLE preserves pairwise distances between models up to a controlled distortion, enabling reliable comparison, clustering, and retrieval in a low-dimensional embedding space.

A.4 Finite-Sample Concentration

In practice, the ABLE representation is computed from N i.i.d. samples rather than the population expectation. We now establish that the empirical estimate concentrates around its expectation, bridging the gap between population-level theory and finite-sample implementation.

Theorem A.4 (Finite-Sample Guarantee). Let $\hat{\Phi}_N(\theta) = \frac{1}{N} \sum_{i=1}^N \phi(a_\theta(x_i))$ be the empirical ABLE representation based on N i.i.d. samples from \mathcal{D} . Assume that $\|\phi(a_\theta(x))\|_\infty \leq B_\phi$ for all x in the support of \mathcal{D} . Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$:

$$\|\hat{\Phi}_N(\theta) - \Phi(\theta)\|_2 \leq B_\phi \sqrt{\frac{2D \log(2D/\delta)}{N}},$$

where D is the dimension of the pre-projection representation.

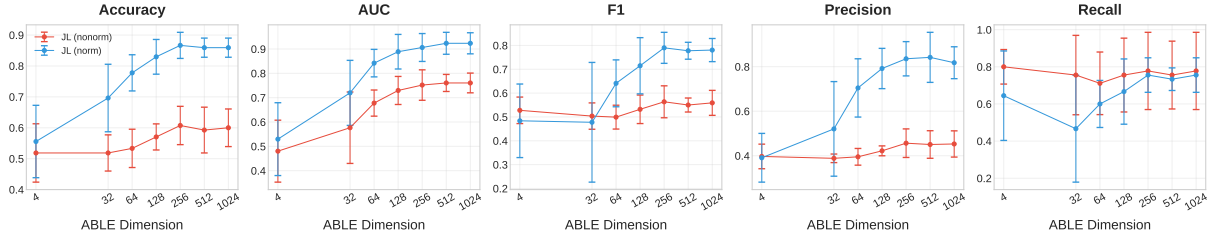


Figure 7: Effect of ABL embedding dimension on relation prediction performance. The blue curve corresponds to applying L2 normalization to each model’s feature vector prior to Johnson-Lindenstrauss projection, while the red curve represents the unnormalized variant. Error bars indicate standard deviation computed over 5 independent runs.

Proof. For each coordinate $j \in \{1, \dots, D\}$, the empirical mean $[\hat{\Phi}_N(\theta)]_j = \frac{1}{N} \sum_{i=1}^N [\phi(a_\theta(x_i))]_j$ is an average of N i.i.d. bounded random variables with $|\phi(a_\theta(x))_j| \leq B_\phi$. By Hoeffding’s inequality:

$$\mathbb{P}(|[\hat{\Phi}_N(\theta)]_j - [\Phi(\theta)]_j| > \epsilon) \leq 2 \exp\left(-\frac{N\epsilon^2}{2B_\phi^2}\right).$$

Setting the right-hand side to δ/D and applying a union bound over all D coordinates, we obtain that with probability at least $1 - \delta$, $|\hat{\Phi}_N(\theta)_j - [\Phi(\theta)]_j| \leq B_\phi \sqrt{\frac{2 \log(2D/\delta)}{N}}$ for all j . The ℓ_2 bound follows from $\|\cdot\|_2 \leq \sqrt{D} \|\cdot\|_\infty$.

Consequence. Theorem A.4 quantifies the gap between the theoretical population-level representation $\Phi(\theta)$ and its finite-sample estimate $\hat{\Phi}_N(\theta)$. The bound scales as $O(\sqrt{D \log D/N})$, indicating that the estimation error diminishes with more samples. This theoretical guarantee is empirically substantiated in Appendix D. There, we calculate pairwise ABL distances between models using two disjoint subsets of the probing data. The strong correlation observed between these two sets of distances indicates that the finite sample size is sufficient to achieve stable representations, thereby validating the convergence predicted by the concentration bound.

B Ablation Study of ABL Representation Dimensionality

This ablation study examines how the ABL embedding dimension d affects representation quality. We vary d across $\{4, 32, 64, 128, 256, 512, 1024\}$ and evaluate each setting on the relation prediction task described in §4.3. We report Accuracy, AUC, F1, Precision, and Recall as evaluation metrics.

As shown in Fig. 7, representation quality improves as d increases and converges around $d =$

256. Very low dimensions (e.g., $d = 4$) yield poor performance because excessive compression discards structural information from the original high-dimensional space. Conversely, dimensions beyond 256 offer only marginal gains while incurring higher computational and storage costs. This trade-off between quality and efficiency motivates our choice of $d = 256$ for all experiments.

We also compare two variants: with and without L2 normalization of each model’s feature vector prior to JL projection. Normalization rescales each feature vector to unit length, i.e., $\mathbf{a}'_m = \mathbf{a}_m / \|\mathbf{a}_m\|_2$. As shown in Fig. 7, the normalized variant (blue) generally outperforms the unnormalized variant (red) across most metrics. This suggests that normalization balances the scale across dimensions, enabling the distance metric to better capture directional differences between models.

C Empirical Validation via Model Interpolation

The theoretical analysis in Appendix A establishes that ABL is a Lipschitz-continuous mapping from model parameters to the embedding space: models with similar parameters should yield similar ABL embeddings. To empirically validate this correspondence on real nonlinear Transformers, we conduct a model interpolation experiment.

We select three Llama-2-7B variants sharing the same architecture as anchor models: the base model (W_0), Llama-2-7b-chat-hf (W_1), and Vicuna-7b-v1.5 (W_2). We construct 25 merged models via linear weight interpolation:

$$W_{\alpha,\beta} = (1 - \alpha - \beta)W_0 + \alpha W_1 + \beta W_2, \quad 1159$$

where $\alpha, \beta \in \{0, 0.25, 0.5, 0.75, 1.0\}$. In weight space, we reconstruct the 2D geometric structure using inner products among the three anchor models; in ABL space, we apply Principal Component

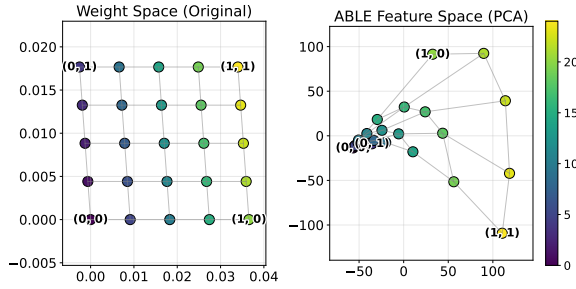


Figure 8: Comparison of weight space and ABLÉ feature space. Left: weight space reconstructed via inner products. Right: ABLÉ feature space after PCA reduction. ABLÉ space preserves the topological structure of weight space but exhibits non-linear distortion. Notably, (0,0) and (0,1) are very close in ABLÉ space, indicating that Vicuna and the base model exhibit similar behavior, consistent with Vicuna’s relatively small fine-tuning magnitude.

Analysis (PCA) to the low-dimensional ABLÉ features for 2D visualization.

As shown in Fig. 8, a notable observation is that (0,0) and (0,1) are very close in ABLÉ space, while (0,0) and (1,0) are much farther apart. This indicates that Vicuna (W_2) exhibits higher behavioral similarity to the base model than Chat (W_1) does. This observation is consistent with the weight space structure: inner products computed directly from model weights show $\|W_2 - W_0\|^2 \approx 0.00032$, whereas $\|W_1 - W_0\|^2 \approx 0.00134$. Vicuna’s fine-tuning magnitude is approximately one quarter that of Chat, demonstrating that ABLÉ distance reflects the similarity structure in parameter space.

Examining the overall geometry of both spaces, weight space exhibits a regular parallelogram grid reflecting the geometry of linear interpolation. Although ABLÉ feature space displays irregular distortion, it crucially preserves the topological structure consistent with weight space: (1) adjacent models remain neighbors in both spaces; (2) grid edges do not intersect, indicating that local neighborhood structure is preserved; (3) color gradients of model indices show similar transition patterns across both spaces.

This result validates the theoretical prediction at two levels. The topological preservation confirms the continuity of the ABLÉ mapping: models close in parameter space remain close in ABLÉ space. The geometric distortion reveals non-linear effects present in Transformers, indicating that ABLÉ sensitivity to parameter changes varies across different regions. This is consistent with the Lipschitz conti-

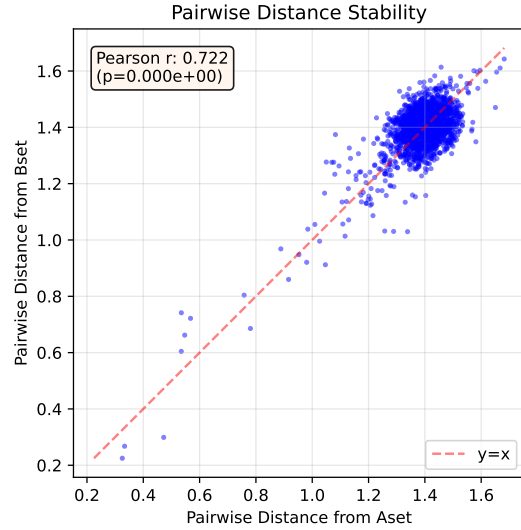


Figure 9: Robustness test for ABLÉ embeddings. Each point corresponds to a model pair, with x-axis showing the pairwise distance computed from Subset A and y-axis from Subset B.

nity established in Appendix A: ABLÉ guarantees bounded propagation of model similarity without requiring strict linear correspondence.

D Robustness to Sample Selection

To evaluate the robustness of ABLÉ to probing sample selection, we conduct a Mantel test. Specifically, we randomly partition the probing dataset into two non-overlapping subsets of equal size, and independently compute ABLÉ embeddings for the same set of models on each subset. We then compute the Euclidean distance between all pairs of models based on embeddings from each subset, yielding two sets of pairwise distances.

As shown in Fig. 9, the two sets of distances exhibit a strong positive correlation (Pearson $r = 0.722$, $p < 0.001$), with points clustering tightly around the $y = x$ diagonal. This result indicates that the inter-model distance structure captured by ABLÉ embeddings remains consistent across different probing samples, demonstrating that ABLÉ captures intrinsic behavioral characteristics of models rather than patterns dependent on specific probing inputs.

E Model List

Table 4 lists all $M = 239$ models included in this work, sorted alphabetically.

Table 4: Complete list of models included in this study.

Model Name	Size	Model Type	Downloads
01-ai/Yi-34B-Chat	34B	Llama	30,397
01-ai/Yi-6B	6B	Llama	16,054
01-ai/Yi-6B-200K	6B	Llama	18,155
abocide/Qwen2.5-7B-Instruct-R1-forfinance	7B	Qwen2	529
AdaptLLM/medicine-chat	6.57B	Llama	1,244
AdaptLLM/medicine-LLM	6.57B	Llama	168
AdaptLLM/medicine-LLM-13B	13B	Llama	34
allenai/tulu-2-dpo-70b	70B	Llama	2,543
Arc53/docsgpt-7b-mistral	7B	Mistral	89
augmnt/shisa-base-7b-v1	7B	Mistral	1,041
bardsai/jaskier-7b-dpo-v5.6	7B	Mistral	159
berkeley-nest/Starling-LM-7B-alpha	7B	Mistral	1,493
bigcode/octocoder	15.5B	Other	148
bigscience/bloom-1b1	1.1B	Bloom	6,749
bigscience/bloom-1b7	1.7B	Bloom	28,490
bigscience/bloom-3b	3B	Bloom	8,904
bigscience/bloom-560m	0.56B	Bloom	112,394
bigscience/bloom-7b1	7B	Bloom	11,136
bigscience/bloomz-3b	3B	Bloom	4,235
bigscience/bloomz-560m	0.56B	Bloom	846,321
bigscience/bloomz-7b1	7B	Bloom	5,299
Biomimicry-AI/ANIMA-Nectar-v2	6.57B	Mistral	975
BioMistral/BioMistral-7B	7B	Mistral	120,331
BioMistral/BioMistral-7B-DARE	7B	Mistral	1,904
bxod/Llama-3.2-1B-Instruct-uz	1B	Llama	35
CausalLM/34b-beta	34B	Llama	8,349
cerebras/Cerebras-GPT-1.3B	1.3B	GPT-2	1,266
cerebras/Cerebras-GPT-111M	0.11B	GPT-2	4,694
cerebras/Cerebras-GPT-2.7B	2.7B	GPT-2	977
cerebras/Cerebras-GPT-256M	0.26B	GPT-2	1,185
cerebras/Cerebras-GPT-590M	0.59B	GPT-2	1,088
cerebras/Cerebras-GPT-6.7B	6.7B	GPT-2	943
cloudyu/Mixtral_11Bx2_MoE_19B	19B	Mixtral	920
codefuse-ai/CodeFuse-DeepSeek-33B	33B	Llama	127
codellama/CodeLlama-13b-Instruct-hf	13B	Llama	19,108
codellama/CodeLlama-34b-Instruct-hf	34B	Llama	20,472
codellama/CodeLlama-7b-hf	7B	Llama	53,172
cognitivecomputations/yayi2-30b-llama	30B	Llama	52
continuedev/instinct	4.86B	Qwen2	241
ConvexAI/Luminex-34B-v0.1	34B	Llama	7,565
ConvexAI/Luminex-34B-v0.2	34B	Llama	7,640
CorticalStack/pastiche-crown-clown-7b-dare-dpo	7B	Mistral	71
CultriX/NeuralTrix-bf16	6.57B	Mistral	74
databricks/dolly-v2-12b	11.58B	GPT-NeoX	3,047
databricks/dolly-v2-3b	3B	GPT-NeoX	-
databricks/dolly-v2-7b	7B	GPT-NeoX	-
davanstrien/query-gen	8B	Llama	56
DeepHat/DeepHat-V1-7B	7B	Qwen2	1,599
deepseek-ai/deepseek-coder-1.3b-base	1.3B	Llama	21,029
deepseek-ai/deepseek-coder-6.7b-instruct	6.7B	Llama	44,962
deepseek-ai/deepseek-llm-67b-chat	67B	Llama	1,815
deepseek-ai/deepseek-math-7b-instruct	7B	Llama	5,002
dfurman/HermesBagel-34B-v0.1	34B	Llama	96
EleutherAI/gpt-neo-1.3B	1.3B	GPT-Neo	28,899
EleutherAI/gpt-neo-125m	0.12B	GPT-Neo	111,607
EleutherAI/gpt-neo-2.7B	2.7B	GPT-Neo	16,759
EleutherAI/llemma_34b	34B	Llama	232
EleutherAI/llemma_7b	7B	Llama	820
EleutherAI/pythia-1.4b	1.4B	GPT-NeoX	18,921
EleutherAI/pythia-1.4b-deduped	1.4B	GPT-NeoX	8,043
EleutherAI/pythia-12b	12B	GPT-NeoX	15,377
EleutherAI/pythia-160m	0.16B	GPT-NeoX	89,343
EleutherAI/pythia-1b-deduped	1B	GPT-NeoX	9,688

continued on next page

continued from previous page

Model Name	Size	Model Type	Downloads
EleutherAI/pythia-2.8b	2.8B	GPT-NeoX	23,704
EleutherAI/pythia-2.8b-deduped	2.8B	GPT-NeoX	7,199
EleutherAI/pythia-410m	0.41B	GPT-NeoX	43,070
EleutherAI/pythia-6.9b	6.9B	GPT-NeoX	17,578
EleutherAI/pythia-70m	0.07B	GPT-NeoX	137,323
eren23/ogno-monarch-jaskier-merge-7b-OH-PREF-DPO	7B	Mistral	71
facebook/opt-1.3b	1.3B	OPT	351,633
facebook/opt-125m	0.12B	OPT	4,242,012
facebook/opt-13b	13B	OPT	9,457
facebook/opt-2.7b	2.7B	OPT	14,711
facebook/opt-350m	0.35B	OPT	100,491
facebook/opt-6.7b	6.7B	OPT	16,057
fblgit/UNA-SimpleSmaug-34b-v1beta	34B	Llama	7,579
FelixChao/llama2-13b-math1.2	13B	Llama	1,093
FelixChao/Scorpio-7B	7B	Mistral	58
FelixChao/vicuna-7B-chemical	7B	Llama	1,086
FelixChao/vicuna-7B-physics	7B	Llama	1,089
galaxy/gowizardlm	6.72B	Llama	950
google/codegemma-1.1-7b-it	7B	Gemma	150
google/codegemma-2b	2B	Gemma	2,143
google/codegemma-7b	7B	Gemma	2,006
google/gemma-2b	2B	Gemma	180,941
google/gemma-2b-it	2B	Gemma	59,882
google/gemma-7b	7B	Gemma	54,233
google/gemma-7b-it	7B	Gemma	108,191
Harshvir/Llama-2-7B-physics	7B	Llama	1,024
HuggingFaceH4/zephyr-7b-alpha	7B	Mistral	1,817
HuggingFaceH4/zephyr-7b-beta	7B	Mistral	75,649
ibivibiv/alpaca-dragon-72b-v1	72B	Llama	877
Imran1/MedChat3.5	7B	Mistral	37
inflatebot/MN-12B-Mag-Mell-R1	12B	Mistral	429
Intel/neural-chat-7b-v3-3	7B	Mistral	31,866
janhq/Jan-v1-4B	4B	Qwen3	1,350
janhq/Jan-v1-edge	1.72B	Qwen3	51
jiawei-ucas/Qwen-2.5-7B-ConsistentChat	7B	Qwen2	16
kamrr/llama-3-8b_dolly_lora	8B	Llama	44
kevin009/llamaRAGdrama	6.57B	Mistral	809
klodia/alpaca	8B	Llama	82
klodia/lora-8b-alpaca-french	8B	Llama	87
klodia/lora-8b-bio	8B	Llama	101
klodia/lora-8b-code	8B	Llama	133
klodia/lora-8b-math	8B	Llama	93
klodia/lora-8b-medic	8B	Llama	109
klodia/lora-8b-physic	8B	Llama	123
kyujinpy/Sakura-SOLRCA-Math-Instruct-DPO-v1	9.79B	Llama	1,008
LLM360/Amber	7B	Llama	3,443
lmsys/vicuna-13b-v1.5	13B	Llama	116,184
lmsys/vicuna-13b-v1.5-16k	13B	Llama	10,289
lmsys/vicuna-33b-v1.3	33B	Llama	1,403
lmsys/vicuna-7b-v1.1	7B	Llama	2,310
lmsys/vicuna-7b-v1.5	7B	Llama	162,278
lmsys/vicuna-7b-v1.5-16k	7B	Llama	4,445
MaziyarPanahi/WizardLM-Math-70B-v0.1	70B	Llama	76
meta-llama/Llama-2-13b-chat-hf	13B	Llama	209,634
meta-llama/Llama-2-70b-chat-hf	70B	Llama	8,026
meta-llama/Llama-2-7b-chat-hf	7B	Llama	325,422
meta-llama/Llama-2-7b-hf	7B	Llama	599,490
meta-llama/Llama-3.1-8B-Instruct	8B	Llama	10,639,638
meta-llama/LlamaGuard-7b	7B	Llama	1,720
meta-llama/Meta-Llama-3-70B	70B	Llama	480,747
meta-llama/Meta-Llama-3-70B-Instruct	70B	Llama	56,883
meta-llama/Meta-Llama-3-8B	8B	Llama	2,236,547
meta-llama/Meta-Llama-3-8B-Instruct	8B	Llama	1,531,726
meta-llama/Meta-Llama-Guard-2-8B	8B	Llama	15,469

continued on next page

continued from previous page

Model Name	Size	Model Type	Downloads
meta-math/MetaMath-Llemma-7B	7B	Llama	1,033
meta-math/MetaMath-Mistral-7B	7B	Mistral	2,201
microsoft/MediPhi-Instruct	3.72B	Phi-3	2,007
microsoft/Orca-2-13b	13B	Llama	10,627
microsoft/Orca-2-7b	7B	Llama	9,125
microsoft/phi-1_5	1.5B	Phi	46,747
microsoft/phi-2	2.65B	Phi	1,071,489
microsoft/Phi-3.5-mini-instruct	3.72B	Phi-3	335,634
mistralai/Mistral-7B-Instruct-v0.1	7B	Mistral	491,519
mistralai/Mistral-7B-v0.1	7B	Mistral	355,407
mlabonne/AlphaMonarch-7B	7B	Mistral	12,378
mlabonne/NeuralHermes-2.5-Mistral-7B	7B	Mistral	105
mosaicml/mpt-30b-instruct	30B	MPT	2,646
mosaicml/mpt-7b	7B	MPT	-
mosaicml/mpt-7b-8k	7B	MPT	-
mosaicml/mpt-7b-8k-chat	7B	MPT	-
mosaicml/mpt-7b-chat	7B	MPT	81,083
mosaicml/mpt-7b-instruct	7B	MPT	-
Neko-Institute-of-Science/metharme-7b	7B	Llama	1,061
Neko-Institute-of-Science/pygmalion-7b	7B	Llama	1,102
neovalle/H4rmoniousAnthea	7B	Mistral	87
Nexusflow/Starling-LM-7B-beta	7B	Mistral	1,434
nomie-ai/gpt4all-13b-snoozy	13B	Llama	1,012
NousResearch/Hermes-4-14B	14B	Qwen3	5,100
NousResearch/Nous-Hermes-13b	13B	Llama	1,330
NousResearch/Nous-Hermes-2-Yi-34B	34B	Llama	7,907
NousResearch/Nous-Hermes-llama-2-7b	7B	Llama	1,253
OpenAssistant/oasst-sft-4-pythia-12b-epoch-3.5	12B	GPT-NeoX	1,519
OpenBuddy/openbuddy-codellama2-34b-v11.1.1-bf16	34B	Llama	1,143
openchat/openchat-3.5-0106	6.57B	Mistral	11,448
openchat/openchat_3.5	6.57B	Mistral	2,622
openlm-research/open_llama_7b	7B	Llama	18,239
pbevan11/llama-3-8b-ocr-correction	8B	Llama	64
PharMolix/BioMedGPT-LM-7B	7B	Llama	578
Plaban81/Moe-4x7b-math-reason-code	7B	Mixtral	67
prithivMLmods/rStar-Coder-Qwen3-0.6B	0.6B	Qwen3	21
project-baize/baize-v2-13b	13B	Llama	1,379
Q-bert/Optimus-7B	7B	Mistral	1,106
Qwen/Qwen1.5-0.5B	0.5B	Qwen2	46,541
Qwen/Qwen1.5-0.5B-Chat	0.5B	Qwen2	57,690
Qwen/Qwen1.5-1.8B	1.8B	Qwen2	17,565
Qwen/Qwen1.5-1.8B-Chat	1.8B	Qwen2	89,028
Qwen/Qwen1.5-14B	14B	Qwen2	50,314
Qwen/Qwen1.5-14B-Chat	14B	Qwen2	15,409
Qwen/Qwen1.5-32B	32B	Qwen2	11,415
Qwen/Qwen1.5-32B-Chat	32B	Qwen2	38,786
Qwen/Qwen1.5-4B	4B	Qwen2	13,977
Qwen/Qwen1.5-4B-Chat	4B	Qwen2	17,665
Qwen/Qwen1.5-7B	7B	Qwen2	70,483
Qwen/Qwen1.5-7B-Chat	7B	Qwen2	16,435
Qwen/Qwen2-0.5B	0.5B	Qwen2	329,782
Qwen/Qwen2-1.5B	1.5B	Qwen2	212,957
Qwen/Qwen2.5-0.5B-Instruct	0.5B	Qwen2	2,456,001
Qwen/Qwen2.5-1.5B	1.5B	Qwen2	504,513
Qwen/Qwen2.5-1.5B-Instruct	1.5B	Qwen2	4,990,149
Qwen/Qwen2.5-7B	7B	Qwen2	821,385
Qwen/Qwen2.5-7B-Instruct	7B	Qwen2	6,379,230
Qwen/Qwen2.5-Coder-7B	7B	Qwen2	91,843
Qwen/Qwen2.5-Coder-7B-Instruct	7B	Qwen2	553,942
Qwen/Qwen3-0.6B	0.6B	Qwen3	8,292,645
Qwen/Qwen3-0.6B-Base	0.6B	Qwen3	161,046
Qwen/Qwen3-1.7B	1.7B	Qwen3	5,473,707
Qwen/Qwen3-14B	14B	Qwen3	686,547
Qwen/Qwen3-4B	4B	Qwen3	3,966,909
Qwen/Qwen3-4B-Instruct-2507	4B	Qwen3	4,166,985

continued on next page

1225

continued from previous page

Model Name	Size	Model Type	Downloads
Qwen/Qwen3-4B-Thinking-2507	4B	Qwen3	501,267
Qwen/Qwen3-8B	8B	Qwen3	4,376,885
rishiraj/CatPPT-base	6.57B	Mistral	3,740
rubenamtz0/llama-3-8b-lora-law2entity	8B	Llama	47
sail/Sailor-7B	7B	Qwen2	100
scb10x/typhoon-7b	7B	Mistral	36,854
SciPhi/SciPhi-Mistral-7B-32k	7B	Mistral	1,098
SciPhi/SciPhi-Self-RAG-Mistral-7B-32k	7B	Mistral	1,241
shleeeee/mistral-ko-tech-science-v1	6.57B	Mistral	13
stabilityai/stable-code-3b	3B	StableLM	3,617
stabilityai/stablelm-2-1_6b	6B	StableLM	1,829
stabilityai/stablelm-3b-4e1t	3B	StableLM	18,150
stabilityai/stablelm-base-alpha-3b	3B	GPT-NeoX	2,084
stabilityai/stablelm-base-alpha-7b	7B	GPT-NeoX	1,980
stabilityai/stablelm-tuned-alpha-3b	3B	GPT-NeoX	2,124
stabilityai/stablelm-tuned-alpha-7b	7B	GPT-NeoX	2,393
stabilityai/stablelm-zephyr-3b	3B	StableLM	9,110
SUSTech/SUS-Chat-34B	34B	Llama	998
teknium/OpenHermes-2-Mistral-7B	7B	Mistral	402
teknium/OpenHermes-2.5-Mistral-7B	7B	Mistral	202,048
tenyx/Tenychat-7B-v1	7B	Mistral	889
Tesslate/UIGEN-X-4B-0729	4B	Qwen3	30
Tesslate/WEBGEN-4B-Preview	4B	Qwen3	55
TheBloke/koala-13B-HF	13B	Llama	1,469
TheBloke/tulu-30B-fp16	30B	Llama	1,174
theprint/TiTan-Qwen2.5-0.5B	0.5B	Qwen2	-
TigerResearch/tigerbot-13b-base	13B	Llama	54
TigerResearch/tigerbot-7b-base-v1	7B	Bloom	18
TigerResearch/tigerbot-7b-base-v2	7B	Bloom	39
TigerResearch/tigerbot-7b-sft-v1	7B	Bloom	47
TigerResearch/tigerbot-7b-sft-v2	7B	Bloom	32
tiiuae/falcon-40b-instruct	40B	Falcon	50,635
tiiuae/falcon-7b	7B	Falcon	64,169
tiiuae/falcon-7b-instruct	7B	Falcon	52,800
tiiuae/falcon-rw-1b	1B	Falcon	6,477
tomg-group-umd/DynaGuard-8B	8B	Qwen3	124
upstage/SOLAR-10.7B-Instruct-v1.0	10.7B	Llama	28,257
venetis/llama3-8b-hermes-sandals-100	8B	Llama	52
VinitT/Sanskrit-llama	8B	Llama	64
Vortex5/Lunar-Nexus-12B	12B	Mistral	19
Vortex5/Moonlit-Shadow-12B	12B	Mistral	14
WizardLM/WizardLM-70B-V1.0	70B	Llama	17,901
worldboss/llama-3-8b-axolotl-fine-tune-qlora	8B	Llama	70
Writer/palmyra-med-20b	20B	GPT-2	1,153
yahma/llama-7b-hf	7B	Llama	4,983
yam-peleg/Experiment26-7B	7B	Mistral	173
zhengr/MixTAO-7Bx2-MoE-v8.1	6.57B	Mixtral	19,428

1226