

# A Shared Geometry of Difficulty in Multilingual Language Models

Anonymous ACL submission

## Abstract

Predicting problem-difficulty in large language models (LLMs) refers to estimating how difficult a task is according to the model itself, typically by training linear probes on its internal representations. In this work, we study the multilingual geometry of problem-difficulty in LLMs by training linear probes using the AMC subset of the Easy2Hard benchmark, translated into 21 languages. We found that difficulty-related signals emerge at two distinct stages of the model internals, corresponding to **shallow** (early-layers) and **deep** (later-layers) internal representations, that exhibit functionally different behaviors. Probes trained on deep representations achieve high accuracy when evaluated on the same language but exhibit poor cross-lingual generalization. In contrast, probes trained on shallow representations generalize substantially better across languages, despite achieving lower within-language performance. Together, these results suggest that LLMs first form a language-agnostic representation of problem difficulty, which subsequently becomes language-specific. This closely aligns with existing findings in LLM interpretability showing that models tend to operate in an abstract conceptual space before producing language-specific outputs. We demonstrate that this two-stage representational process extends beyond semantic content to high-level meta-cognitive properties such as problem-difficulty estimation.

## 1 Introduction

Large language models (LLMs) are increasingly deployed in multilingual settings, yet our understanding of their internal reasoning remains heavily skewed toward English (Resck et al., 2025). While recent work suggests that LLMs may internally “think” in English or exhibit an English-centric representational topology (Chang et al., 2022; Kim and Lee, 2025; Li et al., 2025; Schut et al., 2025; Wendler et al., 2024), less is known about whether

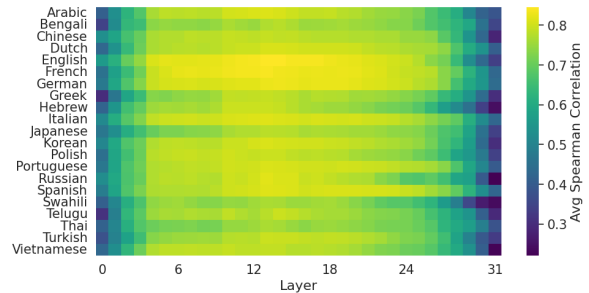


Figure 1: **Layer-wise performance of difficulty probes across languages for LLaMA-3.1-8B.** Heatmap shows, for each test language (rows) and transformer layer (columns), the average Spearman correlation between probe predictions and ground-truth difficulty, where each value is averaged over probes trained on all other languages. Performance peaks in the middle layers, indicating that difficulty representations are most consistently aligned across languages at intermediate depths of the network.

higher-level meta-cognitive attributes (e.g., the model’s internal estimate of how difficult a problem is) generalize across languages.

Lugoloobi and Russell (2025) showed that problem difficulty can be decoded linearly from model residual activations. Focusing on English inputs, they demonstrated that these difficulty signals enable effective interventions, reducing hallucinations via “difficulty vectors” steering and serving as a robust proxy for performance generalization during reinforcement learning.

Whether such internally encoded notions of difficulty persist beyond English, however, remains unexplored. If a model encounters the same mathematical problem expressed in different languages, does it construct distinct difficulty representations, or does it project all inputs onto a shared notion of “difficulty” in activation space? In this work, we address this gap by bridging difficulty probing with multilingual representation learning. We construct a multilingual version of the American Mathematics Competitions (AMC) dataset spanning 21 lan-

066	guages and train linear probes to predict problem	probes, we connect these strands by framing diffi-	115
067	difficulty across languages and model layers. We	culty as a cross-lingual internal property.	116
068	found that problem difficulty is encoded in a depth-		
069	dependent manner that trades off cross-lingual gen-		
070	eralization and language-specific accuracy. Linear		
071	probes show that difficulty is decodable both from		
072	shallow and deep layers, but with distinct behav-		
073	ior: shallow layers support strong cross-lingual		
074	transfer (see Figure 1), while deeper layers yield		
075	higher accuracy within the training language but		
076	generalize poorly across languages, indicating an		
077	early language-agnostic representation that is later		
078	refined into language-specific form.		
079			
	<b>2 Related Work</b>		
080	<b>Difficulty Encoding in LLMs.</b> Recent work		
081	shows that LLMs encode high-level meta-		
082	properties of tasks in their internal representa-		
083	tions. Most notably, <a href="#">Lugoloobi and Russell (2025)</a>		
084	demonstrate that <i>human-perceived problem diffi-</i>		
085	<i>culty</i> is strongly linearly decodable from resid-		
086	ual activations across models and domains using		
087	Easy2Hard-Bench ( <a href="#">Ding et al., 2024</a> ). They fur-		
088	ther validate the utility of this feature, showing		
089	that manipulating the difficulty direction (steering)		
090	can suppress hallucination and improve responses.		
091	However, their analysis is restricted to English in-		
092	puts.		
093	<b>Multilingual Internal Representations.</b> Several		
094	studies suggest that multilingual LLMs do not		
095	maintain fully language-agnostic internal spaces.		
096	<a href="#">Schut et al. (2025)</a> show that multilingual mod-		
097	els perform key reasoning steps in representations		
098	closest to English, even when operating in other		
099	languages. Similarly, <a href="#">Li et al. (2025)</a> find that		
100	probing performance degrades substantially for		
101	low-resource languages and that deeper layers be-		
102	come increasingly language-specific, with reduced		
103	cross-lingual representational similarity. These re-		
104	sults connect to broader analyses of cross-lingual		
105	alignment and multilingual geometry, which find		
106	persistent language-sensitive directions alongside		
107	language-neutral structure in representation space		
108	( <a href="#">Chang et al., 2022</a> ; <a href="#">Hämmerl et al., 2024</a> ; <a href="#">Kim and</a>		
109	<a href="#">Lee, 2025</a> ).		
110	We bridge these lines of work by studying diffi-		
111	culty as a multilingual internal signal. While prior		
112	work has examined difficulty encoding in monolin-		
113	gual settings and language dependence in multilin-		
114	gual models largely through linguistic or reasoning		
		<b>3 Methodology</b>	117
		<b>3.1 Data</b>	118
		We use the AMC subset of the Easy2Hard bench-	119
		mark ( <a href="#">Ding et al., 2024</a> ), which contains approxi-	120
		mately 4,000 math problems annotated with con-	121
		tinuous difficulty scores. Difficulty is estimated via	122
		Item Response Theory (IRT) from human success	123
		rates, yielding values in $[0, 1]$ .	124
		We translated the original English problems into	125
		20 additional languages using gpt-5.1 ( <a href="#">OpenAI,</a>	126
		<a href="#">2026</a> ) resulting in a 21-language benchmark (costs	127
		in Appendix A). For each language, we use approx-	128
		imately 3,000 problems for training and 1,000 for	129
		testing. The train–test split is defined once at the	130
		level of problem, ensuring that cross-lingual eval-	131
		uation is performed on identical unseen problems	132
		(discussion on translation quality in Appendix B).	133
		<b>3.2 Models</b>	134
		We evaluate four instruction-tuned LLMs spanning	135
		different architectures and scales: LLaMA-3.1-8B	136
		( <a href="#">Meta, 2024a</a> ), LLaMA-3.2-3B ( <a href="#">Meta, 2024b</a> ),	137
		LLaMA-3.2-1B ( <a href="#">Meta, 2024b</a> ) and Qwen3-8B	138
		( <a href="#">Qwen Team, 2025</a> ). All models are prompted us-	139
		ing their standard chat templates to reflect realistic	140
		deployment conditions. To ensure basic linguistic	141
		adequacy, we qualitatively verify that each model	142
		preserves meaning when translating a small subset	143
		of prompts from each language back into English.	144
		<b>3.3 Experimental Setup</b>	145
		<b>Feature Extraction.</b> Using TransformerLens	146
		( <a href="#">Nanda and Bloom, 2022</a> ), we extract residual	147
		stream activations from every transformer layer.	148
		For each input, we record the residual vector at the	149
		final prompt token, as this token has been shown	150
		to provide the most informative representation for	151
		linear probing of problem difficulty ( <a href="#">Lugoloobi and</a>	152
		<a href="#">Russell, 2025</a> ).	153
		<b>Probing and Evaluation.</b> We train linear Ridge	154
		regression probes to predict continuous difficulty	155
		scores from residual activations. Probes are trained	156
		independently at each layer to analyze the depth-	157
		wise emergence of difficulty representations. Each	158
		probe is trained on a single language and evaluated	159
		either monolingually (training and testing on the	160
		same language $A$ ) or cross-lingually (training on	161
		language $A$ and testing on a different language	162

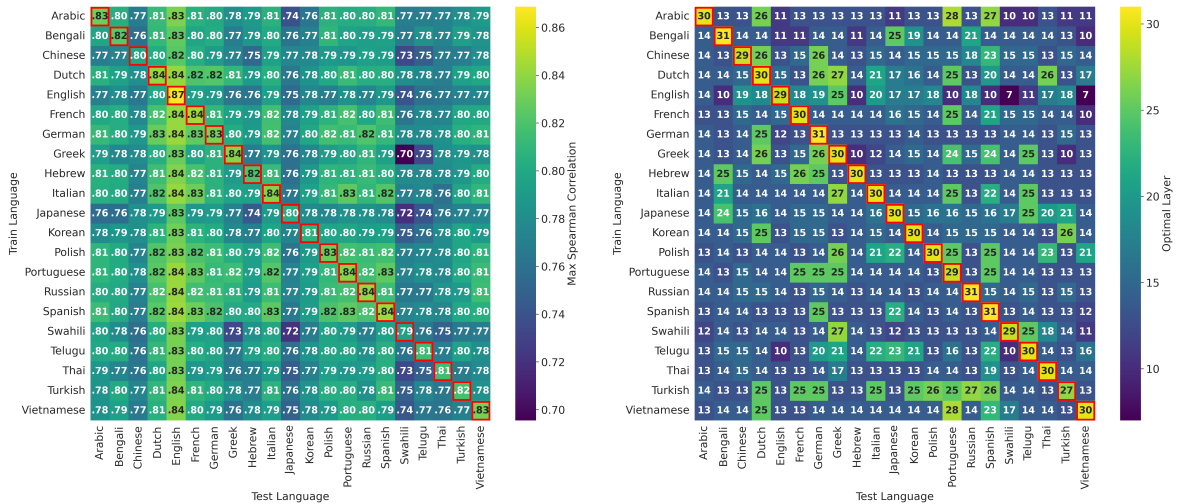


Figure 2: **Cross-lingual structure of difficulty representations in LLaMA-3.1-8B.** **Left:** Maximum Spearman  $\rho$  achieved by linear difficulty probes for each training–testing language pair, evaluated at the layer that maximizes performance for that pair. Diagonal entries correspond to same-language probing, while off-diagonal entries reflect cross-lingual transfer. **Right:** Transformer layer indices at which peak performance is attained for each language pair.

B). The regularization strength is selected from  $\{10, 100, 1000\}$  by maximizing mean Spearman  $\rho$  on the unseen split.

**Evaluation Metric.** Probe performance is measured using Spearman rank correlation ( $\rho$ ) between predicted and ground-truth difficulty scores on held-out test problems. This metric captures ordinal agreement and is robust to scale differences across probes.

## 4 Results

We present results for LLaMA-3.1-8B as a representative model; all trends described below are consistent across LLaMA-3.2 (1B, 3B) and Qwen3-8B, with full results reported in Appendix C.

### 4.1 Depth-Dependent Structure of Difficulty Representations

We first investigate whether problem difficulty is encoded in a language-specific manner or whether it corresponds to a shared internal representation across languages.

Figure 2 (left) summarizes cross-lingual probing performance. Each cell reports the *maximum* Spearman correlation achieved across all layers for a given training–testing language pair. Diagonal entries correspond to same-language probing, while off-diagonal entries reflect cross-lingual transfer. Rows therefore indicate how well a probe trained on a given language generalizes to all others.

Overall, we observe uniformly high correlations across language pairs, including transfers between typologically distant and low-resource languages, indicating that relative difficulty rankings are largely preserved across languages.

Figure 2 (right) reveals a systematic difference in *where* these correlations occur. For LLaMA-3.1-8B, probes trained and tested on the same language (i.e. diagonal scores) consistently peak in later layers (around layer 30), whereas cross-lingual transfer peaks much earlier (around layer 14). This pattern is highly stable across languages: diagonal cells concentrate tightly around a single later (deep) layer, while off-diagonal cells concentrate around earlier layers.

Table 1 quantifies this separation. For LLaMA-3.1-8B, the mean optimal layer for same-language probing is 30.05, compared to 15.67 for cross-lingual transfer. The same depth separation appears across all evaluated models, with the absolute layer indices shifting according to model depth (Appendix C).

Taken together, these results indicate a clear depth-dependent organization of difficulty representations: an earlier, shallow layer representation and a later, deeper representation that is optimized for language-specific performance.

### 4.2 Cross-Lingual Generalization

To further study the consequences of this depth divergence, we explicitly compare probe perfor-

Model	Num. Layers	Probe performance Spearman $\rho$ (mean $\pm$ std)		Optimal layer (mean peak layer)		$\Delta$	$\Delta$
		Same-lang (Diag.)	Cross-lang (Off-diag.)	Same-lang (Diag.)	Cross-lang (Off-diag.)	Transfer drop (Diag $\rightarrow$ Off-Diag)	In-lang drop (Off-Diag $\rightarrow$ Diag)
LLaMA-3.1-8B	32	0.822 $\pm$ 0.019	0.783 $\pm$ 0.024	30.05	15.67	0.177	0.014
LLaMA-3.2-3B	28	0.805 $\pm$ 0.022	0.771 $\pm$ 0.027	22.62	9.07	0.192	0.010
LLaMA-3.2-1B	16	0.797 $\pm$ 0.021	0.761 $\pm$ 0.028	12.05	6.11	0.055	0.007
Qwen3-8B	36	0.849 $\pm$ 0.019	0.783 $\pm$ 0.034	14.29	8.59	0.114	0.015

Table 1: **Cross-lingual difficulty probing summary across models.** Performance reports Spearman  $\rho$  (mean  $\pm$  std) over language pairs. Across all models, cross-lingual (off-diagonal) performance is statistically lower than same-language (diagonal) performance under a paired significance test ( $p < 10^{-3}$ ). Optimal layer denotes the mean layer index achieving peak  $\rho$  in each regime. **Transfer drop (Diag $\rightarrow$ Off-Diag)** measures the loss in cross-lingual performance when probes are fixed at the diagonal-optimal layer. **In-lang drop (Off-Diag $\rightarrow$ Diag)** measures the loss in same-language performance when probes are fixed at the transfer-optimal layer.

mance at layers optimized for same-language versus cross-lingual evaluation. We hypothesized that *diagonal-optimal* layers are maximizing language specific performance (i.e. probes trained on those layers learn language specific features) while *off-diagonal-optimal* layers are optimized for cross-lingual representation (i.e. represent features in a language-agnostic manner). The final two columns of Table 1 report the performance impact of evaluating probes at these respective layer choices on same-language and cross-lingual scenarios (see Appendix D for details).

When probes are evaluated cross-lingually at the diagonal-optimal layer, performance drops substantially. For LLaMA-3.1-8B, fixing probes at the diagonal optimal layer leads to a mean reduction of 0.177 Spearman  $\rho$  under transfer. This sharp decline indicates that deeper layers, while highly predictive within-language, encode problem-difficulty in a manner that does not align well across languages.

In contrast, evaluating probes monolingually at the off-diagonal-optimal layer results in only a negligible loss in same-language performance (0.014 Spearman  $\rho$  for LLaMA-3.1-8B). Thus, the layer that best supports cross-lingual transfer (shallow layer) remains near-optimal for the source language itself.

Together, these findings indicate that problem difficulty is organized around a shared, language-independent direction in activation space that emerges at shallow layers. Deeper layers refine this signal in a language-specific way, improving within-language accuracy at the expense of cross-lingual alignment. This depth-dependent trade-off explains why cross-lingual generalization peaks earlier in the network while same-language perfor-

mance continues to improve at later layers.

## Conclusion

By probing residual activations across 21 languages, we find that LLMs encode problem difficulty in a depth-dependent way: an early, shared representation supports strong cross-lingual transfer, while deeper layers refine difficulty in a language-specific manner, improving monolingual accuracy but reducing alignment across languages. This shared difficulty direction emerges early and remains stable even for low-resource languages, indicating that LLMs form a language-agnostic estimate of difficulty before specializing it to individual languages.

These findings extend prior work by (Lugoloobi and Russell, 2025) in two important ways. First, they demonstrate that the difficulty signal identified in English is not an artifact of language, but instead reflects a genuinely multilingual internal property. Second, they suggest that the mechanistic role of difficulty, previously shown to support steering and hallucination reduction, originates in a shared representational subspace that precedes language-specific reasoning. This supports viewing difficulty as a high-level internal signal of the model, rather than a byproduct of surface-level language features.

From a practical perspective, our results point to new opportunities for multilingual systems, especially in low-resource settings. Lightweight difficulty predictors trained on shallow activations can be deployed without per-language tuning, enabling language-agnostic routing, curriculum design, or compute-aware inference even where labeled data are scarce.

## 293 Limitations

294 While we show that problem difficulty is encoded  
295 in a shared, cross-lingual subspace, our analysis is  
296 confined to mathematical problem solving using  
297 the AMC subset of Easy2Hard. This domain offers  
298 the advantage of well-calibrated, human-derived  
299 difficulty labels, but it represents a narrow slice  
300 of model behavior. It therefore remains unclear  
301 whether the same cross-lingual geometry of diffi-  
302 culty extends to domains where difficulty is more  
303 subjective or context-dependent, such as common-  
304 sense reasoning, programming, or open-ended gen-  
305 eration.

306 Additionally, our empirical evaluation is limited  
307 to a small set of instruction-tuned, decoder-only  
308 language models. Although the observed trends are  
309 consistent across them, it is not yet clear whether  
310 the same representational structure is present in all  
311 model.

312 Finally, our conclusions are based exclusively on  
313 probing analyses. Probing establishes the presence  
314 and cross-lingual alignment of difficulty-related  
315 signals, but does not by itself demonstrate that these  
316 signals play a causal role in shaping model behav-  
317 ior during inference. Prior work provides such  
318 causal evidence in the English setting (Lugoloobi  
319 and Russell, 2025), showing that intervening along  
320 a learned difficulty direction can steer model behav-  
321 ior and reduce hallucinations. Whether analog-  
322 ous interventions transfer across languages—for  
323 example, by steering a model using difficulty vec-  
324 tors learned in one language and applied to an-  
325 other—remains an open question, which we leave  
326 to future work.

## 327 References

328 Tyler Chang, Zhuowen Tu, and Benjamin Bergen. 2022.  
329 The geometry of multilingual language model repre-  
330 sentations. In *Proceedings of the 2022 Conference on*  
331 *Empirical Methods in Natural Language Processing*,  
332 pages 119–136.

333 Mucong Ding, Chenghao Deng, Jocelyn Choo, Zichu  
334 Wu, Aakriti Agrawal, Avi Schwarzschild, Tianyi  
335 Zhou, Tom Goldstein, John Langford, Animashree  
336 Anandkumar, and 1 others. 2024. Easy2hard-bench:  
337 Standardized difficulty labels for profiling llm per-  
338 formance and generalization. *Advances in Neural*  
339 *Information Processing Systems*, 37:44323–44365.

340 Katharina Hämmerl, Jindřich Libovický, and Alexan-  
341 der Fraser. 2024. [Understanding cross-lingual](#)  
342 [Alignment—A survey](#). In *Findings of the Associa-*  
343 *tion for Computational Linguistics: ACL 2024*, pages

10922–10943, Bangkok, Thailand. Association for  
344 Computational Linguistics. 345

Tianjie Ju, Weiwei Sun, Wei Du, Xinwei Yuan,  
346 Zhaochun Ren, and Gongshen Liu. 2024. How large  
347 language models encode context knowledge? a layer-  
348 wise probing study. In *Proceedings of the 2024 Joint*  
349 *International Conference on Computational Linguis-*  
350 *tics, Language Resources and Evaluation (LREC-*  
351 *COLING 2024)*, pages 8235–8246. 352

JaeSeong Kim and Suan Lee. 2025. How language  
353 directions align with token geometry in multilingual  
354 llms. *arXiv preprint arXiv:2511.16693*. 355

Daoyang Li, Haiyan Zhao, Qingcheng Zeng, and Meng-  
356 nan Du. 2025. Exploring multilingual probing in  
357 large language models: A cross-language analysis.  
358 In *Proceedings of the 1st Joint Workshop on Large*  
359 *Language Models and Structure Modeling (XLLM*  
360 *2025)*, pages 61–70. 361

William Lugoloobi and Chris Russell. 2025. Llms  
362 encode how difficult problems are. *arXiv preprint*  
363 *arXiv:2510.18147*. 364

Meta. 2024a. [The llama 3 herd of models](#). 365

Meta. 2024b. Llama 3.2: Model cards and  
366 prompt formats. [https://www.llama.com/docs/  
367 model-cards-and-prompt-formats/llama3\\_2/](https://www.llama.com/docs/model-cards-and-prompt-formats/llama3_2/). 368

Neel Nanda and Joseph Bloom. 2022. Transformerlens.  
369 [https://github.com/TransformerLensOrg/  
370 TransformerLens](https://github.com/TransformerLensOrg/TransformerLens). 371

OpenAI. 2026. [GPT-5.1 \(version date\) \[large language  
372 model\]](#). Generated text from prompt: "Insert your  
373 specific prompt here.". 374

Alibaba Group Qwen Team. 2025. [Qwen3 technical  
375 report](#). 376

Ricardo Rei, Marcos Treviso, Nuno M Guerreiro,  
377 Chrysoula Zerva, Ana C Farinha, Christine Maroti,  
378 José GC De Souza, Taisiya Glushkova, Duarte Alves,  
379 Luisa Coheur, and 1 others. 2022. Cometkiwi: Ist-  
380 unbabel 2022 submission for the quality estimation  
381 shared task. In *Proceedings of the Seventh Confer-*  
382 *ence on Machine Translation (WMT)*, pages 634–  
383 645. 384

Lucas Resck, Isabelle Augenstein, and Anna Korhonen.  
385 2025. [Explainability and interpretability of multilin-](#)  
386 [gual large language models: A survey](#). In *Proceed-*  
387 *ings of the 2025 Conference on Empirical Methods in*  
388 *Natural Language Processing*, pages 20454–20486,  
389 Suzhou, China. Association for Computational Lin-  
390 guistics. 391

Lisa Schut, Yarin Gal, and Sebastian Farquhar. 2025.  
392 Do multilingual llms think in english? In *ICLR 2025*  
393 *Workshop on Building Trust in Language Models and*  
394 *Applications*. 395

Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2024. Do llamas work in english? on the latent language of multilingual transformers. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15366–15394.

## A API Cost Analysis

To construct the multilingual version of the Easy2Hard AMC benchmark, we translated all English problems into 20 additional languages using the gpt-5.1 API. Each translation call used a fixed prompt in which, for each problem–language pair, the placeholders [TARGET LANGUAGE] and [PROBLEM\_TEXT] were replaced with the corresponding language identifier and problem content, respectively (see prompt box). The benchmark contains 3,975 unique problems, yielding a total of 79,500 translations.

Using the official API pricing at the time of experimentation (input: \$0.625/M tokens, output: \$5.00/M tokens, cached input: \$0.0625/M tokens). Overall, the full multilingual benchmark was produced for less than \$50 USD.

### Prompt for Math Problem Translation

#### <SYSTEM PROMPT>

You are a helpful assistant that translates math problems accurately.

#### <USER PROMPT>

Translate the following math problem into [TARGET LANGUAGE].

Problem: [PROBLEM\_TEXT]

Translation:

## B Translation Quality Assessment

To assess the semantic adequacy of the automatically generated translations used in our multilingual benchmark, we employ COMET-Kiwi, a reference-free machine translation quality estimation metric introduced by [Rei et al. \(2022\)](#). Unlike traditional n-gram-based metrics (e.g., BLEU), COMET-Kiwi estimates translation quality by predicting human adequacy judgments from multilingual neural representations, without requiring reference translations ([Ju et al., 2024](#)).

Table 2 reports the average COMET-Kiwi score for each target language. Scores are consistently high across the majority of languages, with most values exceeding 0.75, indicating strong relative semantic adequacy with respect to the English source. High-resource European languages (e.g.,

Italian, French, Spanish, German) achieve the highest scores, while typologically distant and lower-resource languages exhibit more modest degradation, most notably Swahili.

Importantly, COMET-Kiwi scores are not calibrated to absolute quality thresholds and are intended to be interpreted comparatively rather than as guarantees of human-level translation quality. Accordingly, we treat these results as evidence against severe semantic distortion rather than as a definitive certification of translation correctness.

Crucially, translation adequacy is further validated indirectly through downstream probing behavior. If translation quality were poor or systematically distorted problem semantics, difficulty probes trained on one language would fail to transfer to others. Instead, we observe strong cross-lingual probe transfer across all languages considered, including those with lower COMET-Kiwi scores, indicating that the translated problems preserve the underlying difficulty signal required for our analysis. This task-level invariance provides complementary evidence that residual translation noise does not materially affect our main findings.

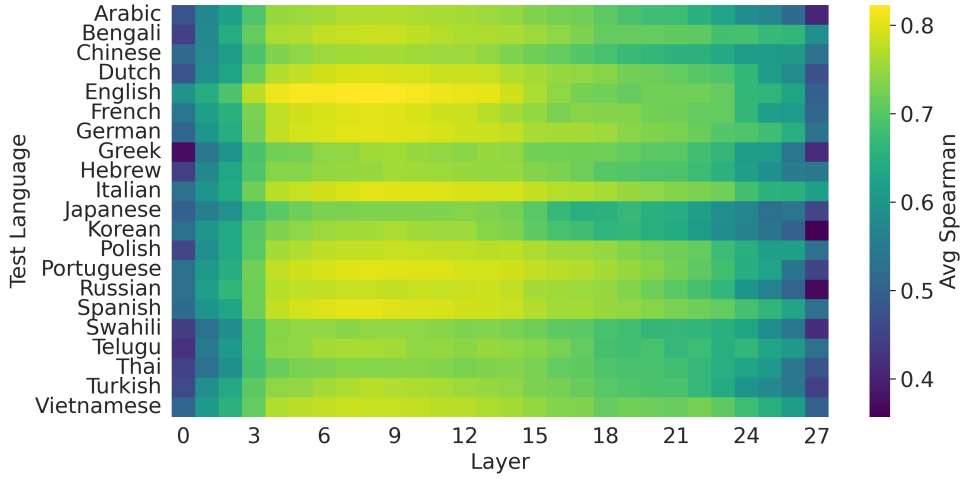
## C Additional Results

This appendix reports supplementary analyses for LLaMA-3.2-3B, LLaMA-3.2-1B, and Qwen3-8B, extending the main results presented

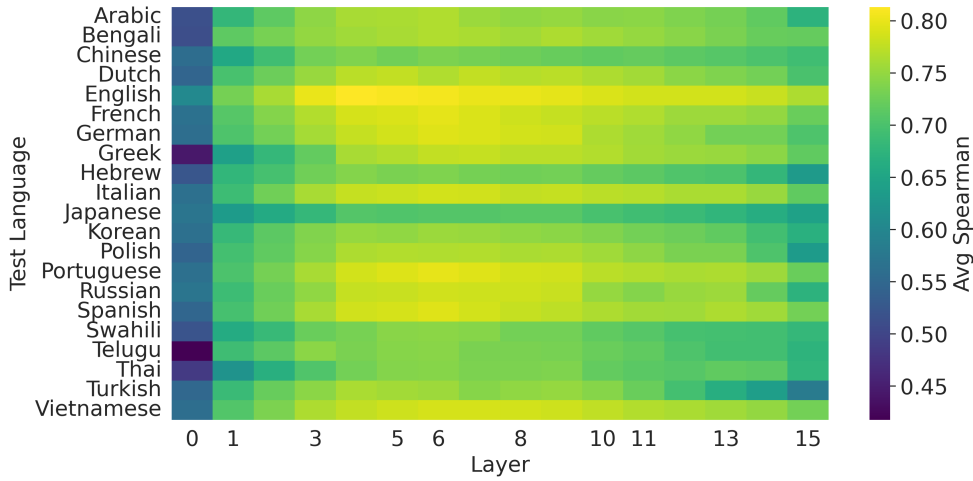
Table 2: **Average COMET-Kiwi scores for translated problems by target language.** Higher scores indicate stronger semantic adequacy with respect to the source language (English).

Language	COMET-Kiwi Score
Italian	0.8250
Japanese	0.8232
Dutch	0.8232
French	0.8218
Spanish	0.8183
Turkish	0.8106
Vietnamese	0.8101
Russian	0.8069
Portuguese	0.8052
Korean	0.8050
Greek	0.8045
Chinese	0.8026
Bengali	0.8019
German	0.8002
Thai	0.7906
Polish	0.7894
Hebrew	0.7681
Arabic	0.7576
Telugu	0.7553
Swahili	0.6134

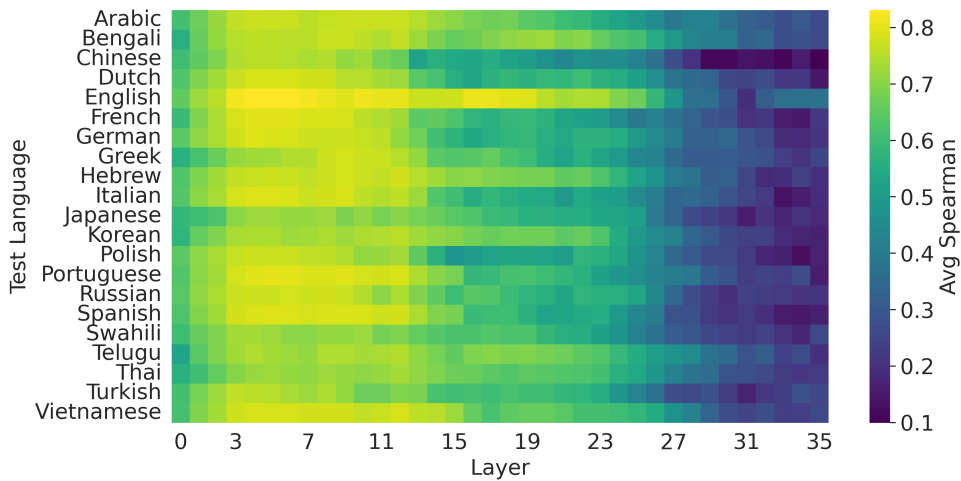
465	for LLaMA-3.1-8B. As summarized in Table 1, all	at which same-language difficulty encoding is	513
466	three models exhibit the same qualitative trends	strongest.	514
467	discussed in Section 4; the figures here provide a		
468	more fine-grained, model-specific view.		
469	<b>Layer-wise cross-lingual performance.</b> Fig-		
470	ure 3 reports the average cross-lingual Spearman		
471	correlation as a function of layer and test lan-		
472	guage. For both LLaMA-3.2 variants, probe per-		
473	formance peaks in early-to-middle layers and de-		
474	grades toward the top of the network, closely mir-		
475	roring the depth-dependent divergence observed for		
476	LLaMA-3.1-8B. Qwen3-8B shows a similar early		
477	peak, followed by a sharper decline in later lay-		
478	ers, consistent with its larger separation between		
479	same-language and cross-lingual optimal layers re-		
480	ported in Table 1. Overall, these results reinforce		
481	the conclusion that the most transferable difficulty		
482	signal emerges before language-specific processing		
483	dominates deeper layers.		
484	<b>Cross-lingual consistency across languages.</b>		
485	Figure 4 presents full Spearman correlation ma-		
486	trices for each model. The problem-wise matrices		
487	(left column) are uniformly high across language		
488	pairs, indicating that relative difficulty rankings		
489	are preserved across translations. The layer-wise		
490	matrices (right column) show that cross-lingual		
491	alignment concentrates within a narrow band of		
492	earlier layers, while optimal layers diverge in		
493	deeper regions. This effect is most pronounced		
494	for LLaMA-3.2-1B, aligning with Table 1, which		
495	shows that reduced model capacity shifts transfer-		
496	optimal layers earlier in the network.		
497	<b>D Cross-Lingual Generalization Analysis</b>		
498	This appendix details the procedure used to com-		
499	pute the <i>Transfer drop</i> and <i>In-language drop</i> re-		
500	ported in the final two columns of Table 1.		
501	1. <b>Layer-wise probing.</b> For each train-		
502	ing-testing language pair $(A, B)$ , linear dif-		
503	ficulty probes are evaluated at every trans-		
504	former layer. Performance is measured using		
505	Spearman $\rho$ on held-out test problems, yield-		
506	ing a layer-wise performance profile for each		
507	pair.		
508	2. <b>Diagonal-optimal layers.</b> For each language		
509	$A$ , the <i>diagonal-optimal layer</i> is defined as		
510	the layer that maximizes performance when		
511	training and testing on the same language		
512	$(A, A)$ . This layer corresponds to the depth		
		3. <b>Transfer drop (Diag→Off-Diag).</b> For a fixed	515
		training language $A$ and each target language	516
		$B \neq A$ :	517
		• compute the best achievable cross-	518
		lingual performance for $(A, B)$ across	519
		all layers;	520
		• compute the cross-lingual performance	521
		obtained when evaluating at $A$ 's	522
		diagonal-optimal layer;	523
		• take the difference between the two.	524
		These differences are averaged across all $B \neq$	525
		$A$ and then across all $A$ to obtain the mean	526
		<i>Transfer drop</i> , measuring how much cross-	527
		lingual performance is lost when using same-	528
		language-optimal layers.	529
		4. <b>Transfer-optimal layers.</b> For each training	530
		language $A$ :	531
		• identify, for each target language $B \neq A$ ,	532
		the layer that maximizes performance for	533
		$(A, B)$ ;	534
		• take the statistical mode of these layers	535
		to obtain a single <i>transfer-optimal layer</i> .	536
		5. <b>In-language drop (Off-Diag→Diag).</b> For	537
		each language $A$ :	538
		• compute the best same-language perfor-	539
		mance across all layers;	540
		• compute the same-language performance	541
		at the transfer-optimal layer;	542
		• take the difference between the two.	543
		These differences are averaged across lan-	544
		guages to yield the mean <i>In-language drop</i> ,	545
		quantifying the cost of fixing probes at	546
		transfer-optimal layers for monolingual evalu-	547
		ation.	548



(a) Llama3.2\_3B

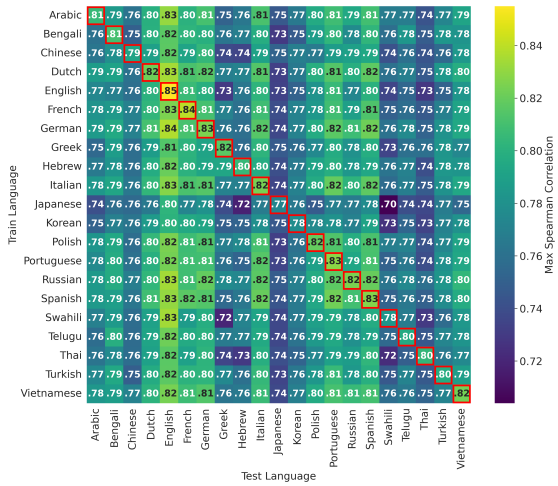


(b) Llama3.2\_1B

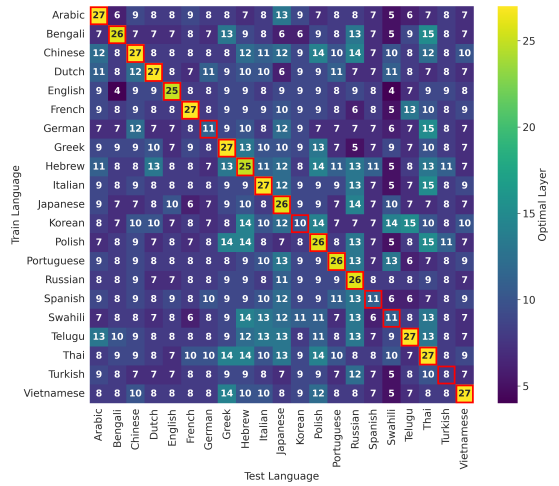


(c) Qwen3

Figure 3: Layer-wise cross-lingual difficulty probing for additional models. Same setup as Figure 1, but for LLaMA-3.2-3B, LLaMA-3.2-1B, and Qwen3-8B. Heatmaps report, for each test language and transformer layer, the average Spearman correlation between predicted and ground-truth difficulty, averaged over probes trained on all other languages. As in Figure 1, cross-lingual performance peaks in early-to-middle layers across models.



(a) Llama3.2\_3B (problem-wise)



(b) Llama3.2\_3B (layer-wise)

