

Understanding LLMs' Cross-Lingual Context Retrieval: How Good It Is And Where It Comes From

Anonymous ACL submission

Abstract

The ability of cross-lingual context retrieval is a fundamental aspect of cross-lingual alignment of large language models (LLMs), where the model extracts context information in one language based on requests in another language. Despite its importance in real-life applications, this ability has not been adequately investigated for state-of-the-art models. In this paper, we evaluate the cross-lingual context retrieval ability of over 40 LLMs across 12 languages to understand the source of this ability, using cross-lingual machine reading comprehension (xMRC) as a representative scenario. Our results show that several small, post-trained open LLMs show strong cross-lingual context retrieval ability, comparable to closed-source LLMs such as GPT-4o. Our interpretability analysis shows that the cross-lingual context retrieval process can be divided into two main phases: question encoding and answer retrieval, which are formed in pre-training and post-training, respectively. Furthermore, the bottleneck of cross-lingual context retrieval lies at the last transformer layers in the second phase, where the effect of post-training can be evidently observed. Our results also indicate that larger LLMs need further multilingual post-training to fully unlock their cross-lingual context retrieval potential.¹

1 Introduction

Since the rise of Large language models (LLMs), many models have demonstrated their strong capability in various NLP tasks (Chang et al., 2024), e.g. ChatGPT², Claude³, Gemini (Gemini Team et al., 2024), LLaMA (Grattafiori et al., 2024), Qwen (Qwen et al., 2025), DeepSeek (DeepSeek-AI et al., 2024), etc. However, due to the domi-

¹Our code and data will be publicly available after the anonymous period.

²<https://chatgpt.com/>

³<https://claude.ai/>

Context: Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager.

Question: What role does John Elway currently have in the Broncos franchise?

Question: Welche Position hat John Elway derzeit im Broncos-Franchise inne?

Question: 约翰·埃尔维目前在野马队中担任什么角色？

Answer: Executive Vice President of Football Operations and General Manager

(a). en-x

Context: John Elway, is currently Denver's Executive Vice President of Football Operations and General Manager.

Question: What role does John Elway currently have in the Broncos franchise?

Answer: Executive Vice President of Football Operations and General Manager

Context: John Elway gehalten, derzeit Denvers Executive Vice President of Football Operations und General Manager ist.

Question: Welche Position hat John Elway derzeit im Broncos-Franchise inne?

Answer: Executive Vice President of Football Operations und General Manager

Context: 过去的记录是由约翰·埃尔维保持的，他在38岁时带领野马队赢得第33届超级碗，目前担任丹佛的橄榄球运营执行副总裁兼总经理。

Question: 约翰·埃尔维目前在野马队中担任什么角色？

Answer: 橄榄球运营执行副总裁兼总经理

(b). x-x

Figure 1: Examples of our en-x and x-x testing scenarios. The figures show examples in English (en), German (de), and Chinese (zh).

nance of English training data, most of these LLMs show their best performance in English (Lai et al., 2023b; Wang et al., 2023). To improve their performance and efficiency in non-English languages, cross-lingual alignment has become a major research topic for multilingual LLMs (Qi et al., 2023; Gao et al., 2024), which encourages LLMs to share capabilities across languages. For example, given requests with the same semantics but in different languages, LLMs should give consistent answers.

One fundamental aspect of such cross-lingual alignment is cross-lingual context retrieval, where the model needs to extract context information in the source language, e.g. English, to answer requests in the target language. However, there is yet no comprehensive evaluation and exploration of this ability.

In this paper, we evaluate the cross-lingual context retrieval ability of SOTA multilingual LLMs, and analyze the source and bottleneck of such ca-

pability. We use cross-lingual machine reading comprehension (xMRC) (Cui et al., 2019) as a simplified but representative scenario, where the target knowledge to be retrieved is literally included in the context, and the model only has to copy part of the context as the answer. This ablates the need for knowledge storage and multilingual text generation, thus better focusing the research on cross-lingual context retrieval. Our main findings are:

- Several SOTA post-trained open-source LLMs, especially in their 7-9B versions show strong xMRC ability, catching up with closed-source models.
- Post-training significantly improves the estimated oracle performance to a very high level in all tested languages, setting space for improvement of real performance.
- The xMRC process can be divided into two phases within the model: question encoding and answer retrieval.
- The question encoding phase forms in pre-training, and its stability correlates with the base model capability; the answer retrieval phase forms in post-training, and serves as a bottleneck for xMRC.

2 Methodology

2.1 Evaluation methods

2.1.1 Scenarios

We use English to non-English (en-x) as a typical and common cross-lingual scenario, where the context and answer are in English, and the question is in non-English. Figure 1a shows a testing sample of en-x task across 3 example languages.

Meanwhile, we use non-English monolingual (x-x) as a comparative to the en-x scenario, in order to distinguish xMRC ability with monolingual MRC ability. Figure 1b shows a testing sample of x-x task in the same 3 example languages.

2.1.2 Metrics

Performance metrics. We use the F1 score and exact match (EM) to evaluate xMRC performance, as adopted by XQuAD. Because the models tend to include unnecessary source text in their correct answers, EM will under-estimate the correctness of the models' response. As a result, we use F1 as the main performance metric.

Cross-lingual performance metrics. To measure the level of the models' performance alignment between English and non-English languages, we calculate the average performance on all the non-English languages, as well as the ratio of non-English performances over English ones.

2.1.3 Performance bottleneck analysis

Error type ablation. We distinguish content error in the model predictions from other less important types of error. Based on preliminary observations, we choose several major types of distracting errors:

- **Language:** answering in x instead of en;
- **Gibberish:** giving irrelevant, non-sense or hallucination output;
- **Refusal:** refusing to answer the question;
- **Blank:** providing no answer at all.

The latter 3 error types can be aggregated as generation failure errors. We measure the proportion of these errors in the models' incorrect responses. For example, the language error rate for test setting en-x is calculated as:

$$E_{\text{lang}}(\text{en}, \text{x}) = \sum_{(r,a) \in W_{\text{en},\text{x}}} \frac{\mathbb{I}(\text{Lang}(r) = \text{en})\mathbb{I}(\text{Lang}(a) = \text{x})}{|W_{\text{en},\text{x}}|}$$

where r is the reference answer, a is the model prediction, and $\text{Lang}(\cdot)$ is a language detector for given text. Meanwhile, a generation failure error rate, e.g. gibberish, will be:

$$E_{\text{gib}}(\text{en}, \text{x}) = \sum_{(r,a) \in W_{\text{en},\text{x}}} \frac{\mathbb{I}(\text{Type}(r) = \text{Gibberish})}{|W_{\text{en},\text{x}}|}$$

where $\text{Type}(\cdot)$ is an LLM-based error type classifier for generation failure samples.

Oracle performance estimation Because the xMRC score can be affected by errors in the generation process, it may not reflect the full potential of the model's cross-lingual context retrieval ability. To ablate this effect, we estimate the oracle retrieval performance of the models by perturbation-based attribution⁴ on contextual sentences or spans, which calculates the importance of each of them to the output. If the sentence/span with the correct answer in it receives the largest importance,

⁴We use Captum to perform the attribution (<https://github.com/pytorch/captum>)

we consider the model’s oracle retrieval on this testing sample is correct. Then, we calculate the accuracy of such selection respectively with one generation step or the whole generation process as the attributed target.

2.2 Analysis methods

2.2.1 Layer-wise attribution to reflect forward process

To better understand the forward calculation process of models performing xMRC retrieval, we need layer-wise attribution in order to see from which input part the LLM draws information into the residual stream (Voita et al., 2024) at each layer.

Previous studies have proposed multiple layer-wise attribution methods, including attention-based (Hao et al., 2021; Ferrando et al., 2022), backpropagation-based (Voita et al., 2021), and decomposition-based (Yang et al., 2023a) etc. Here we choose AttentionLRP (Achtibat et al., 2024), because it can attribute latent representations in each model layer and is especially suitable for transformer models.

Based on this, we define the **Major Relevance Depth (MRD)** to estimate the maximum depth to which a token representation x needs to be encoded, by calculating the layer number corresponding to the 95th percentile of its attributed relevance to model m ’s output:

$$\begin{aligned} \text{MRD}(m, x) &= \min_{1 \leq n \leq N} n \\ \text{s.t. } \sum_{i=1}^n r_{\text{out}}(x, i) &\geq 0.95 \sum_{i=1}^N r_{\text{out}}(x, i) \end{aligned}$$

where $r_{\text{out}}(x, i)$ refers to the normalized relevance score of token representation x in layer i to the final output. Note that N is one less than the total layer number (e.g. 31 for LLaMA-3.1-8B) because we attribute the output of the last layer to its precedent layers. If the MRD of an input token is \hat{n} , then we estimate that the token’s information participates in the context retrieval process only in the first \hat{n} -th layers. Then, for parts of the input, i.e., task description, demonstrations, question and context, we take the maximum MRD of tokens in each part to represent them, and calculate the mean MRD for each part on a group of testing samples to reflect the general distribution.

| Name | Modes | Sizes |
|--------------------|-----------------|------------|
| LLaMA 2 | Chat | 7B |
| LLaMA 3.1 | Base / Instruct | 8B / 70B |
| Mistral V0.3 | Base / Instruct | 7B |
| Qwen 2.5 | Base / Instruct | 7B / 72B |
| DeepSeek V2 | Base / Chat | Lite (16B) |
| Gemma 2 | Base / IT | 9B |
| GPT-3.5-Turbo-0125 | - | - |
| GPT-4o | - | - |

Table 1: Selected models for main evaluation and analysis.

2.2.2 Hidden state similarity to measure cross-lingual alignment

To observe the cross-linguality of model internal representations during the context retrieval process, we collect the hidden states of the input sequence across all model layers, and calculate their cross-lingual similarity. Because our testing samples are parallel across all the testing languages, we can use the ratio of the mean of sample-wise over language-wise cosine similarities to define a cross-lingual similarity ratio $S(\text{en}, x)$ between the model calculation processes in English and language x :

$$S(\text{en}, x) = \overline{\text{Sim}}(E, X) / \overline{\text{Sim}}(X, X) = \frac{(K - 1) \sum_{e_k \in E, x_k \in X} \text{Sim}(e_k, x_k)}{2 \sum_{x_i, x_j \in X} \text{Sim}(x_i, x_j)}$$

where $\overline{\text{Sim}}$ denotes the mean cosine similarity, E denotes the hidden states from English MRC tasks, and X denotes the hidden states from en-x xMRC tasks. K is the total number of testing samples, e_k and x_k are the hidden states from the k -th parallel sample pair between English and language x , x_i, x_j are the hidden states from every two different samples in language x , and the cosine similarity $\text{Sim}(x, y) = x \cdot y / \|x\| \|y\|$.

3 Experiment Settings

3.1 Dataset

We use the XQuAD dataset to measure the xMRC performance of LLMs, because its testing samples are parallel in all the 12 included languages and thus suitable for cross-lingual transforming. The XQuAD dataset is composed of 1190 parallel data points for each of its 12 languages, ensuring a consistent evaluation setup across languages. With an average context length of 702.50, the dataset provides sufficiently rich contexts for machine reading comprehension. Furthermore, the 12 languages cover diverse language families, scripts, and resource levels, making our evaluation more trustworthy and representative.

| | F1 scores | | | | | error rates | | |
|---------------------------|-----------|-----------|----------|--------------|-------------|---------------|-----------------|------------------|
| | en-en | mean en-x | mean x-x | en-x / en-en | x-x / en-en | mean language | mean generation | en-en generation |
| LLaMA-3.1-8B | 75.97 | 49.01 | 70.28 | 0.64 | 0.93 | 0.32 | 8.88 | 5.60 |
| LLaMA-3.1-70B | 82.39 | 58.68 | 74.73 | 0.71 | 0.91 | 60.21 | 2.63 | 1.20 |
| Mistral-V0.3-7B | 79.57 | 58.74 | 64.92 | 0.74 | 0.82 | 21.24 | 14.25 | 0.49 |
| Qwen-2.5-7B | 62.42 | 57.51 | 66.11 | 0.92 | 1.06 | 0.96 | 3.29 | 1.51 |
| Qwen-2.5-72B | 86.03 | 78.92 | 81.16 | 0.92 | 0.94 | 10.90 | 2.53 | 0.00 |
| DeepSeek-V2-Lite-16B | 73.81 | 44.65 | 57.66 | 0.61 | 0.78 | 12.97 | 8.45 | 1.87 |
| Gemma-2-9B | 80.42 | 66.82 | 72.90 | 0.83 | 0.91 | 1.91 | 4.11 | 1.02 |
| LLaMA-3.1-Instruct-8B | 77.89 | 72.13 | 65.02 | 0.93 | 0.83 | 0.89 | 2.53 | 0.85 |
| LLaMA-3.1-Instruct-70B | 83.29 | 73.07 | 74.13 | 0.88 | 0.89 | 0.23 | 1.87 | 1.85 |
| Mistral-V0.3-Instruct-7B | 62.01 | 56.63 | 49.39 | 0.91 | 0.80 | 2.77 | 3.30 | 1.77 |
| Qwen-2.5-Instruct-7B | 81.83 | 76.43 | 71.61 | 0.93 | 0.88 | 0.67 | 3.21 | 2.75 |
| Qwen-2.5-Instruct-72B | 77.12 | 66.04 | 70.29 | 0.86 | 0.91 | 4.58 | 1.62 | 0.38 |
| DeepSeek-V2-Chat-Lite-16B | 70.30 | 54.03 | 49.95 | 0.77 | 0.71 | 2.36 | 5.92 | 0.58 |
| Gemma-2-IT-9B | 83.69 | 78.72 | 75.53 | 0.94 | 0.90 | 0.17 | 2.47 | 1.95 |
| GPT-3.5-Turbo-0125 | 81.74 | 68.75 | 72.04 | 0.84 | 0.88 | 0.16 | 2.80 | 0.00 |
| GPT-4o | 83.29 | 78.76 | 75.68 | 0.95 | 0.91 | 0.10 | 1.40 | 0.00 |

Table 2: 2-shot F1 scores on en-x and x-x tasks, and 2-shot language error and generation failure error rates (%) on en-x tasks.

3.2 Models and tools

We adopt a variety of SOTA open and business LLMs, including LLaMA-3.1, Qwen-2.5, GPT-4o, etc., in smaller and larger sizes. Table 1 shows our selected model list in the main evaluation and analysis results. Also, Table 4 in Appendix A.1 shows a full list of all the tested models, where some of them are of older versions and alternative sizes compared with the main-list models.

We also adopt tools for error type identification. For language error detection, we use Lingua⁵ with its high-accuracy mode, the accuracy of which is satisfactory in our tested languages. For generation failure errors detection, we use Qwen-2.5-72B-Instruct (prompt shown in Appendix A.2) to act like a classifier, distinguishing normal outputs and the three types of error.

3.3 Prompts

In our xMRC evaluation, we try our best to adopt the most suitable prompt templates for each tested model to minimize its limitation of performance. We try two prompting formats for each model and take the one with higher xMRC performance to represent that model. As shown in Appendix A.2, the v1 prompt format places the task description before the demonstrations, and the v2 prompt format places the task description after the demonstrations.

We also test each model in 2-shot and 0-shot settings. Because most of the tested LLMs perform better in the 2-shot setting, we focus our evaluation and analysis mainly on this setting. The 0-shot evaluation results are in Appendix A.3.

⁵<https://github.com/pemistahl/lingua>

4 Results

4.1 Evaluation results

The left part of Table 2 summarizes the F1 scores of our main-list models on both en-x cross-lingual MRC (xMRC) and x-x monolingual MRC tasks. Comparing the xMRC scores with the monolingual scores highlights the difference in model performance between cross-lingual and monolingual settings. Detailed F1 scores for each model, language, and task are further illustrated in Table 6 in Appendix A.3.

4.1.1 Cross-lingual performance

Generally, one can see the English MRC performance of most tested models are high (over 70 out of 100), but the xMRC scores ranges (from 45 to 78), showing the performance gap in context retrieval with English and non-English queries, even with the same English context and answer.

Down into individual models, while GPT-4o shows the highest cross-lingual performance and smallest language gap, several open-sourced post-trained model series, such as Gemma-2-IT, Qwen-2.5-Instruct and LLaMA-3.1-Instruct, also show comparable performance levels and small language gaps, indicating that newest training techniques contribute to higher xMRC performance.

An interesting observation is that, for LLaMA-3.1-8B and Gemma-2-9B, the English performances after post-training are close to those of the base versions, but the cross-lingual performances greatly improve. Also, this phenomenon is more prominent in smaller (7-9B) than larger models, bring the former a smaller performance gap between English and non-English. This shows post-training, especially on smaller LLMs, is important

293 to the xMRC task.

294 4.1.2 Comparison with Monolingual 295 performance

296 In general, the performance gaps between English and non-English monolingual MRC are much
297 smaller for most of the tested LLMs than xMRC,
298 and the Qwen models even show higher non-
299 English performance than English. This suggests
300 that non-English language fluency is not the main
301 cause of the ranging xMRC performance.
302

303 Also, for base models and post-trained large
304 models ($\sim 70B$), the monolingual performances
305 in non-English languages are always higher than
306 cross-lingual performance. A possible explanation
307 to this may be the cross-lingual task is more dif-
308 ficult and less frequent in training, and it requires
309 cross-lingual understanding.

310 However, for post-trained, smaller models (7-
311 9B), the pattern flips, where the monolingual per-
312 formances become consistently lower than cross-
313 lingual ones. This result suggests that these models
314 become capable of **using their English ability to**
315 **assist context retrieval with non-English queries**,
316 overcoming the difficulty and low-frequency of
317 the cross-lingual task. Also, since we should ex-
318 pect larger LLMs having better instruction fol-
319 lowing and understanding in general, this further
320 highlights that post-training better elicits the cross-
321 lingual context retrieval ability on smaller LLMs.

322 4.2 Performance bottleneck analysis

323 4.2.1 Error type ablation

324 **Language error.** The right part of Table 2 shows
325 the average percentage of language errors among
326 all languages. One can see an expected advan-
327 tage of post-trained models against base models,
328 because the former are better in following the cross-
329 lingual task format. However, since the error rate
330 is low for all the post-trained models, it cannot be
331 viewed as a bottleneck for the xMRC task.

332 **Generation failure errors.** The right part of Ta-
333 ble 2 also shows the average summed percentages
334 of gibberish, refusal, and blank error in samples
335 with wrong model answers (detailed percentages
336 in Table 7-9c in Appendix A.4). One can see the
337 generation failure rates are minor for most models,
338 regardless of size and post-training, marking it not
339 the bottleneck of xMRC either.

| Model | Step | Sequence | | |
|------------------------|-------|----------|-------|-------|
| Language | en-en | en-x | en-en | en-x |
| LLaMA-3.1-8B | 66.55 | 67.59 | 35.14 | 36.39 |
| LLaMA-3.1-70B | 47.47 | 54.92 | 47.89 | 56.97 |
| LLaMA-3.1-Instruct-8B | 89.86 | 83.42 | 93.75 | 86.66 |
| LLaMA-3.1-Tuned-8B | 86.49 | 81.04 | 89.53 | 83.07 |
| LLaMA-3.1-Instruct-70B | 92.82 | 90.67 | 95.44 | 93.54 |

Table 3: Oracle performance estimated (F1) for LLaMA models in en-en and en-x (average) scenarios. The estimation is performed with one generation step (left) and with the whole generated sequence, respectively.

340 4.2.2 Estimated oracle performance

341 As described in §2.1.3, the oracle performances of
342 the LLaMA models estimated with one generation
343 step and the whole sequence are shown in Table 3,
344 where two phenomena stand out:

345 First, the estimated oracle performances of post-
346 trained models are significantly higher than the
347 base models, both on en-en and en-x settings, sug-
348 gesting post-training is crucial to the xMRC task.

349 Second, the estimated oracle performance for
350 en-x is close to en-en for all the LLaMA models,
351 and the oracle for post-trained models are basically
352 over 90%. This is quite different from the actual
353 performance, where en-en is better than en-x, and
354 the performances are not that high. This suggests
355 the models have actually obtained the ability to
356 locate the correct answer in the xMRC task, but the
357 ability needs to be further elicited.

358 5 Two-phased mechanism of xMRC

359 The above results show that post-training on
360 smaller LLM elicits the cross-lingual transfer of
361 contextual retrieval ability. However, the mech-
362 anism that enables such transfer is unknown. Con-
363 sidering the model forward process from the input
364 prompt to the output answer, we come up with a
365 two-phase hypothesis of the xMRC process (taking
366 en-x as an example):

1. **Question encoding.** The queries (mainly the non-English question) will be encoded into a shared semantic space, where queries in different languages are aligned and understood in a language-neutral way;
2. **Answer retrieval.** The encoded queries will be used to match the answer in the English context according to the task description and format, then the answer is generated by copying from the original context.

Figure 2 shows an illustration of the hypothe-
sized process. We test our hypothesis from at-

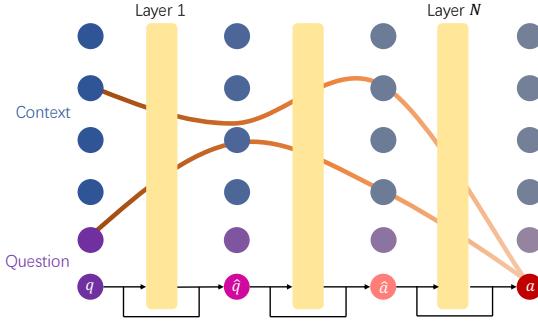


Figure 2: Illustration of the hypothesized two-phased xMRC process. As layer depths increases, the last question token will be gradually transferred to the first answer token.

tribution and hidden-state views. Since LLaMA-3.1-Instruct-8B shows high performance alignment across languages and is widely used in the community, we take its family to perform our analysis.

To further ensure the effect of post-training, we also conduct finetuning and compare the model behavior before and after it.

5.1 Attribution view

We use layer-wise attribution (§2.2.1) to identify the question encoding and answer retrieval phases on the LLaMA models. Figure 3 shows the mean MRD of the contexts and the questions for the LLaMA-3.1-Instruct-8B on testing samples that are identified as either "balanced" or "en-superior" in all tested languages.⁶

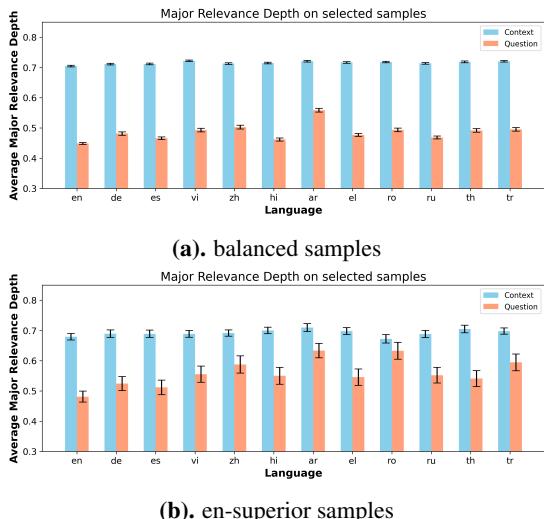
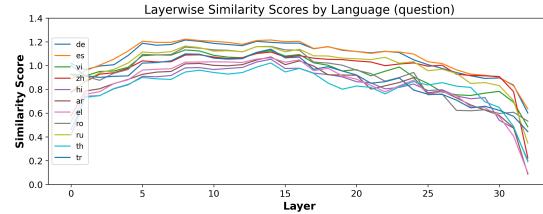


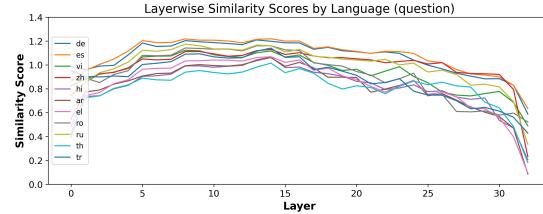
Figure 3: Mean MRD of the context and question parts for LLaMA-3.1-Instruct-8B.

One can see that the mean question MRD is

⁶We identify a sample as "balanced" if the model F1 score on it is above 0.5 in all directions; and a sample as "en-superior" if the F1 score in English is higher than the average of other languages with a margin greater than 0.5.



(a). balanced samples



(b). en-superior samples

Figure 4: Question hidden state similarity between English and other languages in each layer of the LLaMA-3.1-Instruct-8B model.

significantly and substantially lower than the mean context MRD in all tested languages and across the LLaMA models, revealing a clear phased behavior.

Also, by comparing the MRDs of "balanced" (Figures 3a) and "en-superior" samples (Figures 3b), we can see that the MRDs of "balanced" samples are more stable than those of "en-superior" samples, suggesting correlation between the high cross-lingual context retrieval ability and a clear two-phase behavior.

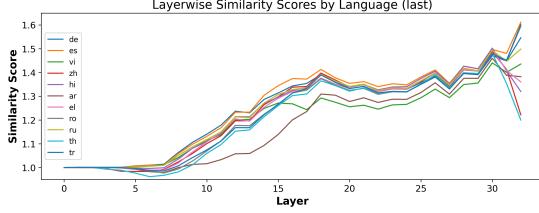
Details and More results can be found in Appendix B.1. We found this pattern consistent for different LLaMA models and not affected by the prompt format, though it is weaker in LLaMA-2-Chat-7B, which has a smaller pre-trained capacity and weaker multilingual ability.

In summary, the attribution results supports the two-phase hypothesis, and indicates that the phased behavior is already formed after pre-training, regardless of model size and prompt format. The phasing strength correlates with the pre-training capacity and the xMRC task performance.

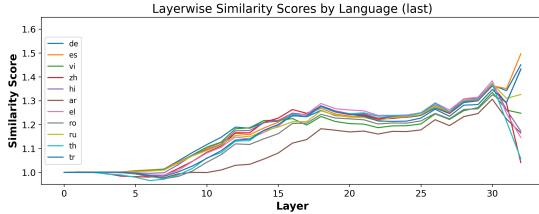
5.2 Hidden states view

The hidden state similarity results also support our hypothesis. Figures 4, 5, and 6 show the en-x hidden state similarity of the question, the last input token (predicting the start of the answer) and the context parts for the LLaMA-3.1-Instruct-8B model. Results for other LLaMA models can be found in Appendix B.2. The observed trends are consistent:

- For question representations, they all show



(a). balanced samples



(b). en-superior samples

Figure 5: Last-input-token hidden state similarity between English and other languages in each layer of the LLaMA-3.1-Instruct-8B model.

a shared arc-shaped trend, where the highest similarity to English appears at the relative depth of around 1/3;

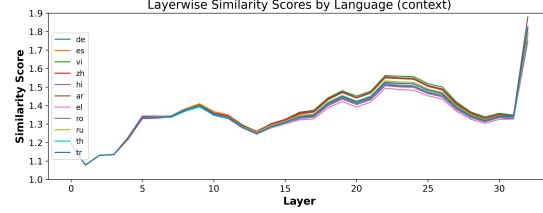
- For context representations, a consistent double-peak trend can be observed, with a "turning point" around the relative depth of 0.4 (matching the question MRD) and the second, higher peak at around 0.7 (matching the context MRD);
- For last input token representations, one can see a consistent "plateau" of similarity starting at around a relative depth of 0.5, which also matches the mean question MRD.

Again, for the less powerful model LLaMA-2-Chat-7B (in Appendix B.2 Figure 18), the trends are weaker: its question similarity to English varies much across languages, and the "plateau" of answer similarity to English starts later than other models, which is after the mean question and context MRDs.

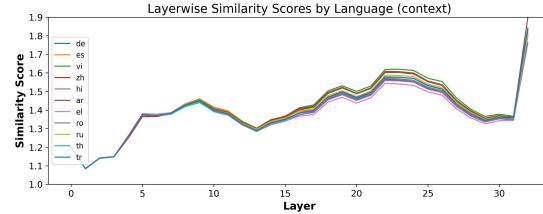
These results suggest that the hidden states similarity through the xMRC process also undergo two main phases with evident distinction. This phased behavior already exists in pre-trained LLMs, and is preserved in post-training. Also, the phasing strength correlates with the model capability built during pre-training.

5.3 Pre-training vs. post-training

Although the phasing behavior is shaped during pre-training, other results show that post-training



(a). balanced samples



(b). en-superior samples

Figure 6: Context hidden state similarity between English and other languages in each layer of the LLaMA-3.1-Instruct-8B model.

is also essential to complete the xMRC task.

We observe that (in Appendix B.2 Figure 19) the last-input-token similarity of base models experience severe and consistent decline in the last few layers, across all non-English languages. Since in our xMRC task we expect the same English output for all tested languages, this drop in cross-lingual similarity can directly affect the performance and its cross-lingual alignment. Meanwhile, post-training significantly narrows this decline, especially on the 8B model. As a result, for the Instruct-8B model, languages with Latino alphabets reveal even higher similarity than the previous layers.

In this regard, a possible reason for the larger performance gap between English and non-English for larger post-trained models could be, their post-training is insufficient to form the answer retrieving phase. As shown in Figure 17, their last-input-token similarity drops in the final layers instead of rising, which is a feature of the base models.

Another outcome of post-training is sensitivity to the demonstrations. One can see from Figure 6 and 17 (c) and (d), both post-trained models show divergence in context similarity, which is not seen in their base models. Since the context in the xMRC task is English-only, this divergence in similarity can only be the effect of cross-lingual demonstration preceding the context. Compared with the 8B model, the divergence of the 70B instruct model is too large, suggesting that it is too sensitive to the demonstrations, and may harm the cross-lingual performance.

489 5.4 Validation of post-training

490 To validate the effect of post-training, we tune the
 491 LLaMA-3.1-8B model use the TULU-v3 dataset
 492 (Lambert et al., 2025) into a model called LLaMA-
 493 3.1-Tuned-8B (Appendix B.3). Like the instruct
 494 model, the en-x xMRC performance of the tuned
 495 model become significantly higher than the base
 496 model (Table 10), and the decline of cross-lingual
 497 hidden state similarity in the last few layers also
 498 narrows (Figure 7). This is clear evidence for the
 499 essential role of post-training in completing the
 500 xMRC task.

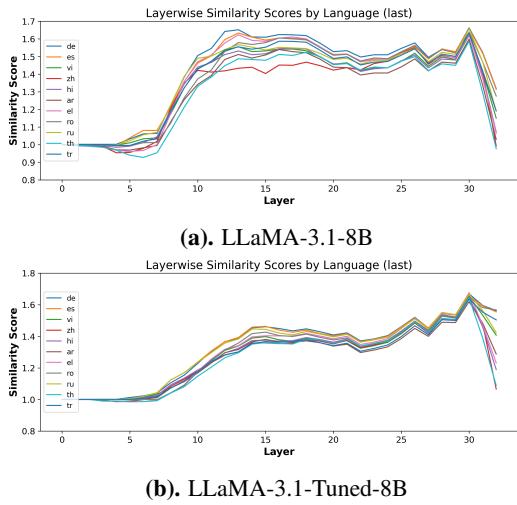


Figure 7: Change in Last-input-token hidden state similarity between English and other languages in each layer of LLaMA-3.1-8B after tuning on the same data samples as "balanced" of the base model.

501 6 Related Work

502 6.1 Cross-lingual alignment of LLMs

503 Previous studies have shown a misalignment of
 504 LLMs with English and other languages. With re-
 505 spect to performance, Lai et al. (2023b), Ahuja
 506 et al. (2024) and Wang et al. (2024) demon-
 507 strated that SOTA LLMs performed better in English,
 508 and showed inconsistency when dealing with non-
 509 English queries. Etxaniz et al. (2024) found that
 510 LLMs performed worse with non-English prompts
 511 than with self-translated English prompts. Beyond
 512 performance, Qi et al. (2023) demonstrated the
 513 low cross-lingual consistency of factual knowl-
 514 edge does of LLMs, and Gao et al. (2024) showed
 515 that multilingual pre-training and instruction tun-
 516 ing could only enhance superficial levels of cross-
 517 lingual alignment. However, the previous work
 518 mostly focused on queries with a consistent lan-
 519 guage, instead of cross-lingual queries.

520 There have also been many techniques to en-
 521 hance LLMs' cross-lingual alignment. For exam-
 522 ple, adding parallel data in the pre-training stage
 523 (Lample and Conneau, 2019; Jiang et al., 2022; Wei
 524 et al., 2023; Lu et al., 2024); and the post-training
 525 stage, including instruction tuning (Li et al., 2023b;
 526 Zhang et al., 2025; Li et al., 2023a; Cahyawijaya
 527 et al., 2023; Chai et al., 2024; Kuulmets et al., 2024;
 528 Shaham et al., 2024; Kew et al., 2024) and pref-
 529 erence tuning (Lai et al., 2023a; She et al., 2024).
 530 Especially, extra translation training is commonly
 531 used (Zhang et al., 2023; Yang et al., 2023b; Li
 532 et al., 2024; Ranaldi et al., 2024; Zhu et al., 2024;
 533 Lu et al., 2024). In this paper, we examine the
 534 effect of some of these techniques by comparing
 535 various SOTA models.

536 6.2 Cross-lingual machine reading 537 comprehension

xMRC is a relatively new task of natural language
 understanding. Cui et al. (2019) proposed the task,
 in order to improve non-English MRC performance
 by introducing English resources. There are some
 representative datasets in this area, such as XQA
 (Liu et al., 2019), BiPaR (Jing et al., 2019), MLQA
 (Lewis et al., 2020), and XQuAD (Artetxe et al.,
 541 2020). Ushio et al. (2023) also proposes a pipeline
 542 for multilingual QA generation.

Previous work on the xMRC task mainly focuses
 on enhancing the performance of task-specific mod-
 543 els using techniques such as data augmentation
 (Bornea et al., 2021; Xiang et al., 2024), knowledge
 injection (Duan et al., 2021), contrastive learning
 (Chen et al., 2022) and knowledge transfer (Cao
 544 et al., 2023; Xu et al., 2023). However, the xMRC
 545 performance of LLMs has not been studied.

555 7 Conclusion

We investigate the cross-lingual context retrieval
 556 ability of LLMs with the xMRC task. We first eval-
 557 uates the xMRC performance of existing open- and
 558 closed-sourced LLMs, and find that post-trained
 559 7-9B models show high performance and little gap
 560 across English and non-English languages. Then,
 561 we conduct analysis on the LLaMA models and
 562 identifies the two main phases (question encod-
 563 ing and answer retrieval) of the xMRC process, and
 564 find a possible explanation to the language gap for
 565 larger post-trained models. We hope our research
 566 will inspire future study to foster the cross-lingual
 567 alignment of LLMs in a broader scope.

569 8 Limitations

570 While this study provides valuable insights into the
571 cross-lingual context retrieval abilities of LLMs
572 and identifies a two-phased mechanism underlying
573 this process, it is important to acknowledge certain
574 limitations.

575 First, the scope of our empirical evaluation is
576 constrained by available resources and time. This
577 necessarily limits the breadth of our testing, pre-
578 venting us from exhaustively covering the rapidly
579 expanding landscape of LLMs. Furthermore, while
580 we test across 12 diverse languages, a more com-
581 prehensive analysis would ideally include an even
582 wider range of languages in order to ensure the
583 generalizability of our findings across linguistic
584 diversity.

585 Second, although we successfully identify a two-
586 phased feature of cross-lingual machine reading
587 comprehension and confirm its correlation with pre-
588 training and post-training stages, the precise factors
589 within these training processes that drive this out-
590 come remain unclear. Future work could delve
591 deeper into the specifics of pre-training objectives,
592 data composition, and post-training techniques to
593 pinpoint the exact elements that contribute to the
594 emergence and effectiveness of these two phases.

595 Finally, within the two major phases we discover
596 – question encoding and answer retrieval – we ob-
597 serve hints of more fine-grained changes in model
598 behavior, particularly in the hidden state similarity
599 curves. These preliminary observations suggest the
600 potential for a more nuanced understanding of the
601 xMRC process. Future studies could further in-
602 vestigate these finer-grained dynamics within each
603 phase to gain a more detailed and complete pic-
604 ture of how LLMs achieve cross-lingual context
605 retrieval.

606 Beyond these limitations, it is also important to
607 consider potential risks associated with this work.
608 While our research is foundational and not directly
609 tied to specific applications, advancements in cross-
610 lingual context retrieval, like any technology, could
611 be misused. For example, improved cross-lingual
612 capabilities might inadvertently contribute to the
613 spread of misinformation if models are used to re-
614 trieve and amplify biased or inaccurate information
615 across languages. Furthermore, if deployed with-
616 out careful consideration, these technologies could
617 exacerbate existing inequalities by favoring lan-
618 guages and knowledge systems already dominant
619 in LLM training data, potentially marginalizing

620 less-represented languages and perspectives. Fu-
621 ture work should consider these dual-use aspects
622 and explore mitigation strategies to ensure respon-
623 sible development and deployment of cross-lingual
624 NLP technologies, paying special attention to fair-
625 ness and inclusivity across diverse linguistic com-
626 munities.

627 9 Ethics Statements

628 This research adheres to ethical principles in its use
629 of language models and data. All language models
630 evaluated and finetuned in this study are accessed
631 and utilized in compliance with their respective
632 licenses and terms of service. Furthermore, the
633 XQuAD dataset employed for evaluation, and the
634 TULU-v3 dataset used for finetuning LLaMA-3.1-
635 8B, are both publicly available datasets intended
636 for research purposes. Based on our review and
637 the documented nature of these datasets, we have
638 determined that they are not designed to collect
639 or contain personally identifiable information or
640 offensive content. To the best of our knowledge,
641 and as indicated in their public documentation, nei-
642 ther dataset includes data that names or uniquely
643 identifies individual people, nor do they present
644 offensive content.

645 Our use of these existing artifacts, including both
646 language models and datasets, is aligned with their
647 intended use within research contexts. Specifically,
648 derivatives of data accessed for research purposes,
649 such as model outputs and analysis results, are used
650 solely within the bounds of academic inquiry and
651 are not disseminated or utilized outside of these
652 research contexts, in accordance with responsible
653 data handling practices. We have strived to conduct
654 this research responsibly and ethically, focusing
655 on understanding and improving the cross-lingual
656 capabilities of language models for the benefit of
657 the NLP community, while respecting the intended
658 use and access conditions of all resources utilized.

659 References

660 Reduan Achitbat, Sayed Mohammad Vakilzadeh Hatefi,
661 Maximilian Dreyer, Aakriti Jain, Thomas Wiegand,
662 Sebastian Lapuschkin, and Wojciech Samek. 2024.
663 AttnLRP: Attention-Aware Layer-Wise Relevance
664 Propagation for Transformers. In *Proceedings of the*
665 *41st International Conference on Machine Learning*,
666 pages 135–168. PMLR.

667 Sanchit Ahuja, Divyanshu Aggarwal, Varun Gumma,
668 Ishaan Watts, Ashutosh Sathe, Millicent Ochieng,

| | | |
|-----|---|-----|
| 669 | Rishav Hada, Prachi Jain, Mohamed Ahmed, Kalika Bali, and Sunayana Sitaram. 2024. MEGAVERSE: Benchmarking large language models across languages, modalities, models and tasks . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 2598–2637, Mexico City, Mexico. Association for Computational Linguistics. | 726 |
| 670 | | 727 |
| 671 | | 728 |
| 672 | | 729 |
| 673 | | 730 |
| 674 | | 731 |
| 675 | | 732 |
| 676 | | 733 |
| 677 | | 734 |
| 678 | Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. On the cross-lingual transferability of monolingual representations . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4623–4637, Online. Association for Computational Linguistics. | 735 |
| 679 | | 736 |
| 680 | | 737 |
| 681 | | 738 |
| 682 | | 739 |
| 683 | | 740 |
| 684 | Mihaela Bornea, Lin Pan, Sara Rosenthal, Hans Florian, and Avi Sil. 2021. Multilingual Transfer Learning for QA using Translation as Data Augmentation. In <i>AAAI Conference on Artificial Intelligence</i> . | 741 |
| 685 | | 742 |
| 686 | | 743 |
| 687 | | 744 |
| 688 | Samuel Cahyawijaya, Holy Lovenia, Tiezheng Yu, Willy Chung, and Pascale Fung. 2023. InstructAlign: High-and-low resource language alignment via continual crosslingual instruction tuning . In <i>Proceedings of the First Workshop in South East Asian Language Processing</i> , pages 55–78, Nusa Dua, Bali, Indonesia. Association for Computational Linguistics. | 745 |
| 689 | | 746 |
| 690 | | 747 |
| 691 | | 748 |
| 692 | | 749 |
| 693 | | 750 |
| 694 | | 751 |
| 695 | Tingfeng Cao, Chengyu Wang, Chuanqi Tan, Jun Huang, and Jinhui Zhu. 2023. Sharing, Teaching and Aligning: Knowledgeable Transfer Learning for Cross-Lingual Machine Reading Comprehension . <i>Preprint</i> , arXiv:2311.06758. | 752 |
| 696 | | 753 |
| 697 | | 754 |
| 698 | | 755 |
| 699 | | 756 |
| 700 | Linzheng Chai, Jian Yang, Tao Sun, Hongcheng Guo, Jiaheng Liu, Bing Wang, Xiannian Liang, Jiaqi Bai, Tongliang Li, Qiyao Peng, and Zhoujun Li. 2024. xCoT: Cross-lingual Instruction Tuning for Cross-lingual Chain-of-Thought Reasoning . <i>Preprint</i> , arXiv:2401.07037. | 757 |
| 701 | | 758 |
| 702 | | 759 |
| 703 | | 760 |
| 704 | | 761 |
| 705 | | 762 |
| 706 | Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2024. A Survey on Evaluation of Large Language Models . <i>ACM Transactions on Intelligent Systems and Technology</i> , 15(3):1–45. | 763 |
| 707 | | 764 |
| 708 | | 765 |
| 709 | | 766 |
| 710 | | 767 |
| 711 | | 768 |
| 712 | | 769 |
| 713 | Nuo Chen, Linjun Shou, Ming Gong, and Jian Pei. 2022. From Good to Best: Two-Stage Training for Cross-Lingual Machine Reading Comprehension . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 36(10):10501–10508. | 770 |
| 714 | | 771 |
| 715 | | 772 |
| 716 | | 773 |
| 717 | | 774 |
| 718 | Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shijin Wang, and Guoping Hu. 2019. Cross-lingual machine reading comprehension . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1586–1595, Hong Kong, China. Association for Computational Linguistics. | 775 |
| 719 | | 776 |
| 720 | | 777 |
| 721 | | 778 |
| 722 | | 779 |
| 723 | | 780 |
| 724 | | 781 |
| 725 | | 782 |
| 726 | DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qiiao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhua Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiying Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhi-gang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. 2024. Deepseek-v3 technical report . <i>Preprint</i> , arXiv:2412.19437. | 778 |
| 727 | | 779 |
| 728 | | 780 |
| 729 | | 781 |
| 730 | | 782 |
| 731 | | 783 |
| 732 | | 784 |
| 733 | | 785 |
| 734 | | 786 |
| 735 | | 787 |

| | | |
|-----|--|-----|
| 788 | ter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers), pages 550–564, Mexico City, Mexico. Association for Computational Linguistics. | 849 |
| 789 | | 850 |
| 790 | | 851 |
| 791 | | 852 |
| 792 | Javier Ferrando, Gerard I. Gállego, and Marta R. Costa-jussà. 2022. Measuring the mixing of contextual information in the transformer. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 8698–8714, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. | 853 |
| 793 | | 854 |
| 794 | | 855 |
| 795 | | 856 |
| 796 | | 857 |
| 797 | | 858 |
| 798 | | 859 |
| 799 | Changjiang Gao, Hongda Hu, Peng Hu, Jiajun Chen, Jixing Li, and Shujian Huang. 2024. Multilingual pre-training and instruction tuning improve cross-lingual knowledge alignment, but only shallowly. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 6101–6117, Mexico City, Mexico. Association for Computational Linguistics. | 860 |
| 800 | | 861 |
| 801 | | 862 |
| 802 | | 863 |
| 803 | | 864 |
| 804 | | 865 |
| 805 | | 866 |
| 806 | | 867 |
| 807 | | 868 |
| 808 | Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Jack Krawczyk, Cosmo Du, Ed Chi, Heng-Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty Chan, Manaal Faruqui, Aliaksei Severyn, Hanzhao Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah, Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan, Jeremiah Liu, Andras Orban, Fabian Güra, Hao Zhou, Xinying Song, Aurelien Boffy, Harish Ganapathy, Steven Zheng, HyunJeong Choe, Ágoston Weisz, Tao Zhu, Yifeng Lu, Siddharth Gopal, Jarrod Kahn, Maciej Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa, Majd Al Merey, Martin Baeuml, Zhifeng Chen, Laurent El Shafey, Yujing Zhang, Olcan Sercinoglu, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rustemi, Natalie Clay, Phil Crone, Tomas Kociský, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam Bloniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrette, Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Gaurav Singh Tomar, Evan Senter, Martin Chadwick, Ilya Korakov, Nithya Attaluri, Iñaki Iturrate, Ruibo Liu, Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse Hartman, Xavier Garcia, Thanumalayan Sankaranarayana Pillai, Jacob Devlin, Michael Laskin, Diego de Las Casas, Dasha Valter, Connie Tao, Lorenzo Blanco, Adrià Puigdomènech Badia, David Reitter, Mianna Chen, Jenny Brennan, Clara Rivera, Sergey Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski, Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yiming Gu, Kate Olszewska, Ravi Addanki, Antoine Miech, Annie Louis, Denis Teplyashin, Geoff Brown, Elliot Catt, Jan Balaguer, Jackie Xiang, Pidong Wang, Zoe Ashwood, Anton Briukhov, Albert Webson, Sanjay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-Wei Chang, Axel Stjerngren, Josip Djolonga, Yuting Sun, Ankur Bapna, Matthew Aitchison, Pedram Pejman, Henryk Michalewski, Tianhe Yu, Cindy Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich, Kehang Han, Peter Humphreys, Thibault Sellam, James Bradbury, Varun Godbole, Sina Samangooei, Bogdan Damoc, Alex Kaskasoli, Sébastien M. R. Arnold, Vijay Vasudevan, Shubham Agrawal, Jason Riesa, Dmitry Lepikhin, Richard Tanburn, Srivatsan Srinivasan, Hyeontaek Lim, Sarah Hodkinson, Pranav Shyam, Johan Ferret, Steven Hand, Ankush Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Giang, Alexander Neitz, Zaheer Abbas, Sarah York, Machel Reid, Elizabeth Cole, Aakanksha Chowdhery, Dipanjan Das, Dominika Rogozińska, Vitaliy Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas Zilka, Flavien Prost, Luheng He, Marianne Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan, Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu, Raoul de Liedekerke, Justin Gilmer, Carl Saroufim, Shruti Rijhwani, Shaobo Hou, Disha Shrivastava, Anirudh Baddepudi, Alex Goldin, Adnan Ozturk, Albin Cassirer, Yunhan Xu, Daniel Sohn, Devendra Sachan, Reinald Kim Amplayo, Craig Swanson, Dessie Petrova, Shashi Narayan, Arthur Guez, Siddhartha Brahma, Jessica Landon, Miteyan Patel, Ruizhe Zhao, Kevin Villela, Luyu Wang, Wen-hao Jia, Matthew Rahtz, Mai Giménez, Legg Yeung, James Keeling, Petko Georgiev, Diana Mincu, Boxi Wu, Salem Haykal, Rachel Saputro, Kiran Vodrahalli, James Qin, Zeynep Cankara, Abhanshu Sharma, Nick Fernando, Will Hawkins, Behnam Neyshabur, Solomon Kim, Adrian Hutter, Priyanka Agrawal, Alex Castro-Ros, George van den Driessche, Tao Wang, Fan Yang, Shuo yiin Chang, Paul Komarek, Ross McIlroy, Mario Lučić, Guodong Zhang, Wael Farhan, Michael Sharman, Paul Natsev, Paul Michel, Yamini Bansal, Siyuan Qiao, Kris Cao, Siamak Shakeri, Christina Butterfield, Justin Chung, Paul Kishan Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar Soparkar, Karel Lenc, Timothy Chung, Aedan Pope, Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo Wang, Joshua Maynez, Mary Phuong, Taylor Tobin, Andrea Tacchetti, Maja Trebacz, Kevin Robinson, Yash Katariya, Sebastian Riedel, Paige Bailey, Kefan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone, Neil Houlsby, Xuehan Xiong, Zhen Yang, Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music Li, Thais Kagohara, Jay Pavagadhi, Sophie Bridgers, Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed, | 912 |

| | | | |
|-----|--|---|---|
| 913 | Tianqi Liu, Richard Powell, Vijay Bolina, Mariko Iinuma, Polina Zablotskaia, James Besley, Da-Woon Chung, Timothy Dozat, Ramona Comanescu, Xiance Si, Jeremy Greer, Guolong Su, Martin Polacek, Raphaël Lopez Kaufman, Simon Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad Tomasev, Jin-wei Xing, Christina Greer, Helen Miller, Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Angelos Filos, Milos Besta, Rory Blewins, Ted Klimenko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir, Vered Cohen, Charline Le Lan, Krishna Haridasan, Amit Marathe, Steven Hansen, Sholto Douglas, Rajkumar Samuel, Mingqiu Wang, Sophia Austin, Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso Lorenzo, Lars Lowe Sjösund, Sébastien Cevey, Zach Gleicher, Thi Avrahami, Anudhyan Boral, Hansa Srivivasan, Vittorio Selo, Rhys May, Konstantinos Aisopos, Léonard Hussenot, Livio Baldini Soares, Kate Baumli, Michael B. Chang, Adrià Recasens, Ben Caine, Alexander Pritzel, Filip Pavetic, Fabio Pardo, Anita Gergely, Justin Frye, Vinay Ramasesh, Dan Horgan, Kartikeya Badola, Nora Kassner, Subhrajit Roy, Ethan Dyer, Víctor Campos Campos, Alex Tomala, Yunhao Tang, Dalia El Badawy, Elspeth White, Basil Mustafa, Oran Lang, Abhishek Jindal, Sharad Vikram, Zhitao Gong, Sergi Caelles, Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng, Wojciech Stokowiec, Ce Zheng, Phoebe Thacker, Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh, James Svensson, Max Bileschi, Piyush Patil, Ankesh Anand, Roman Ring, Katerina Tsihlas, Arpi Vezer, Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom Kwiatkowski, Samira Daruki, Keran Rong, Allan Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg, Mina Khan, Lisa Anne Hendricks, Marie Pellat, Vladimir Feinberg, James Cobon-Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives, Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou, Qingze Wang, Thibault Sottiaux, Michela Paganini, Jean-Baptiste Lespiau, Alexandre Moufarek, Samer Hassan, Kaushik Shivakumar, Joost van Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh Goyal, Matthew Tung, Andrew Brock, Hannah Sheahan, Vedant Misra, Cheng Li, Nemanja Rakićević, Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener, Fantine Huot, Matthew Lamm, Nicola De Cao, Charlie Chen, Sidharth Mudgal, Romina Stella, Kevin Brooks, Gautam Vasudevan, Chenxi Liu, Mainak Chain, Nivedita Melinkeri, Aaron Cohen, Venus Wang, Kristie Seymore, Sergey Zubkov, Rahul Goel, Summer Yue, Sai Krishnakumaran, Brian Albert, Nate Hurley, Motoki Sano, Anhad Mohananey, Jonah Joughin, Egor Filonov, Tomasz Kępa, Yomna Eldawy, Jiawern Lim, Rahul Rishi, Shirin Badiezadegan, Taylor Bos, Jerry Chang, Sanil Jain, Sri Gayatri Sundara Padmanabhan, Subha Puttagunta, Kalpesh Krishna, Leslie Baker, Norbert Kalb, Vamsi Bedapudi, Adam Kurzrok, Shuntong Lei, Anthony Yu, Oren Litvin, Xiang Zhou, Zhichun Wu, Sam Sobell, Andrea Siciliano, Alan Papir, Robby Neale, Jonas Bragagnolo, Tej Toor, Tina Chen, Valentin Anklin, Feiran Wang, | Richie Feng, Milad Gholami, Kevin Ling, Lijuan Liu, Jules Walter, Hamid Moghaddam, Arun Kishore, Jakub Adamek, Tyler Mercado, Jonathan Mallinson, Siddhinita Wandekar, Stephen Cagle, Eran Ofek, Guillermo Garrido, Clemens Lombriser, Maksim Mukha, Botu Sun, Hafeezul Rahman Mohammad, Josip Matač, Yadi Qian, Vikas Peswani, Paweł Janus, Quan Yuan, Leif Schelin, Oana David, Ankur Garg, Yifan He, Oleksii Duzhiyi, Anton Älgmyr, Timothée Lottaz, Qi Li, Vikas Yadav, Luyao Xu, Alex Chinien, Rakesh Shivanna, Aleksandr Chuklin, Josie Li, Carrie Spadine, Travis Wolfe, Kareem Mohamed, Subhabrata Das, Zihang Dai, Kyle He, Daniel von Dincklage, Shyam Upadhyay, Akanksha Maurya, Luyan Chi, Sebastian Krause, Khalid Salama, Pam G Rabinovitch, Pavan Kumar Reddy M, Aarush Selvan, Mikhail Dektiarev, Golnaz Ghiasi, Erdem Guven, Himanshu Gupta, Boyi Liu, Deepak Sharma, Idan Heimlich Shtacher, Shachi Paul, Oscar Akerlund, François-Xavier Aubet, Terry Huang, Chen Zhu, Eric Zhu, Elico Teixeira, Matthew Fritze, Francesco Bertolini, Liana-Eleonora Marinescu, Martin Bölle, Dominik Paulus, Khyati Gupta, Tejas Latkar, Max Chang, Jason Sanders, Roopa Wilson, Xuewei Wu, Yi-Xuan Tan, Lam Nguyen Thiet, Tulsee Doshi, Sid Lall, Swaroop Mishra, Wanming Chen, Thang Luong, Seth Benjamin, Jasmine Lee, Ewa Andrejczuk, Dominik Rabiej, Vipul Ranjan, Krzysztof Styrc, Pengcheng Yin, Jon Simon, Malcolm Rose Harriott, Mudit Bansal, Alexei Robsky, Geoff Bacon, David Greene, Daniil Mirylenka, Chen Zhou, Obaid Sarvana, Abhimanyu Goyal, Samuel Andermatt, Patrick Siegler, Ben Horn, Assaf Israel, Francesco Pongetti, Chih-Wei "Louis" Chen, Marco Selvatici, Pedro Silva, Kathie Wang, Jackson Tolins, Kelvin Guu, Roey Yogev, Xiaochen Cai, Alessandro Agostini, Maulik Shah, Hung Nguyen, Noah Ó Donnáile, Sébastien Pereira, Linda Friso, Adam Stambler, Adam Kurzrok, Chenkai Kuang, Yan Romanikhin, Mark Geller, ZJ Yan, Kane Jang, Cheng-Chun Lee, Wojciech Fica, Eric Malmi, Qijun Tan, Dan Banica, Daniel Balle, Ryan Pham, Yanping Huang, Diana Avram, Hongzhi Shi, Jasjot Singh, Chris Hidey, Niharika Ahuja, Pranab Saxena, Dan Dooley, Srividya Pranavi Potharaju, Eileen O'Neill, Anand Gokulchandran, Ryan Foley, Kai Zhao, Mike Dusenberry, Yuan Liu, Pulkit Mehta, Ragha Kotikalapudi, Chalence Safranek-Shrader, Andrew Goodman, Joshua Kessinger, Eran Globen, Prateek Kolhar, Chris Gorgolewski, Ali Ibrahim, Yang Song, Ali Eichenbaum, Thomas Brovelli, Sahitya Potluri, Preethi Lahoti, Cip Baetu, Ali Ghorbani, Charles Chen, Andy Crawford, Shalini Pal, Mukund Sridhar, Petru Gurita, Asier Mujika, Igor Petrovski, Pierre-Louis Cedoz, Chenmei Li, Shiyuan Chen, Niccolò Dal Santo, Siddharth Goyal, Jitesh Punjabi, Karthik Kappagantu, Chester Kwak, Pallavi LV, Sarmishta Velury, Himadri Choudhury, Jamie Hall, Premal Shah, Ricardo Figueira, Matt Thomas, Minjie Lu, Ting Zhou, Chintu Kumar, Thomas Jurdí, Sharat Chikkerur, Yenai Ma, Adams Yu, Soo Kwak, Victor Ähdel, Sujeevan Rajayogam, Travis Choma, Fei Liu, Aditya Barua, Colin Ji, Ji Ho | 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 |
|-----|--|---|---|

| | | | |
|------|---|---|------|
| 1040 | Park, Vincent Hellendoorn, Alex Bailey, Taylan Bi- | Felix Fischer, Jun Xu, Christina Sorokin, Chris Al- | 1104 |
| 1041 | ial, Huanjie Zhou, Mehrdad Khatir, Charles Sutton, Wojciech Rzadkowski, Fiona Macintosh, Kon- | berti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev, | 1105 |
| 1042 | stantin Shagin, Paul Medina, Chen Liang, Jinjing | Hannah Forbes, Dylan Banarse, Zora Tung, Mark | 1106 |
| 1043 | Zhou, Pararth Shah, Yingying Bi, Attila Dankovics, | Omernick, Colton Bishop, Rachel Sterneck, Rohan | 1107 |
| 1044 | Shipra Banga, Sabine Lehmann, Marissa Bredesen, | Jain, Jiawei Xia, Ehsan Amid, Francesco Piccinno, | 1108 |
| 1045 | Zifan Lin, John Eric Hoffmann, Jonathan Lai, Ray- | Xingyu Wang, Praseem Banzal, Daniel J. Mankowitz, | 1109 |
| 1046 | nald Chung, Kai Yang, Nihal Balani, Arthur Braži- | Alex Polozov, Victoria Krakovna, Sasha Brown, | 1110 |
| 1047 | skas, Andrei Sozanschi, Matthew Hayes, Héctor Fer- | MohammadHossein Bateni, Dennis Duan, Vlad Firoiu, | 1111 |
| 1048 | nández Alcalde, Peter Makarov, Will Chen, Anto- | Meghana Thotakuri, Tom Natan, Matthieu Geist, | 1112 |
| 1049 | nio Stella, Liselotte Snijders, Michael Mandl, Ante | Ser tan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko | 1113 |
| 1050 | Kärrman, Paweł Nowak, Xinyi Wu, Alex Dyck, Kr- | Tojo, Michael Kwong, James Lee-Thorp, Christo- | 1114 |
| 1051 | ishnan Vaidyanathan, Raghavender R, Jessica Mal- | pher Yew, Danila Sinopalnikov, Sabela Ramos, John | 1115 |
| 1052 | ilet, Mitch Rudominer, Eric Johnston, Sushil Mit- | Mellor, Abhishek Sharma, Kathy Wu, David Miller, | 1116 |
| 1053 | tal, Akhil Udathu, Janara Christensen, Vishal Verma, | Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jen- | 1117 |
| 1054 | Zach Irving, Andreas Santucci, Gamaleldin Elsayed, | ifer Beattie, Emily Caveness, Libin Bai, Julian | 1118 |
| 1055 | Elnaz Davoodi, Marin Georgiev, Ian Tenney, Nan | Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi | 1119 |
| 1056 | Hua, Geoffrey Cideron, Edouard Leurent, Mah- | Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, | 1120 |
| 1057 | moud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy | Frederick Liu, Fan Yang, Rui Zhu, Tian Huey Teh, | 1121 |
| 1058 | Zheng, Dylan Scandinaro, Heinrich Jiang, Jasper | Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, | 1122 |
| 1059 | Snoek, Mukund Sundararajan, Xuezhi Wang, Zack | Daniel Toyama, Evan Rosen, Sasan Tavakkol, Lint- | 1123 |
| 1060 | Ontiveros, Itay Karo, Jeremy Cole, Vinu Rajashekhar, | ing Xue, Chen Elkind, Oliver Woodman, John Car- | 1124 |
| 1061 | Lara Tumeh, Eyal Ben-David, Rishub Jain, Jonathan | penter, George Papamakarios, Rupert Kemp, Sushant | 1125 |
| 1062 | Uesato, Romina Datta, Oskar Bunyan, Shimu Wu, | Kafle, Tanya Grunina, Rishika Sinha, Alice Tal- | 1126 |
| 1063 | John Zhang, Piotr Stanczyk, Ye Zhang, David Steiner, | bert, Diane Wu, Denese Owusu-Afriyie, Cosmo | 1127 |
| 1064 | Subhajit Naskar, Michael Azzam, Matthew Johnson, | Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna | 1128 |
| 1065 | Adam Paszke, Chung-Cheng Chiu, Jaume Sanchez | Narayana, Jing Li, Saaber Fatehi, John Wieting, | 1129 |
| 1066 | Elias, Afroz Mohiuddin, Faizan Muhammad, Jin | Omar Ajmeri, Benigno Uria, Yeongil Ko, Laura | 1130 |
| 1067 | Miao, Andrew Lee, Nino Vieillard, Jane Park, Ji- | Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi | 1131 |
| 1068 | ageng Zhang, Jeff Stanway, Drew Garmon, Abhijit | Pang, Yeqing Li, Nir Levine, Ariel Stolovich, Re- | 1132 |
| 1069 | Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Lu- | becca Santamaría-Fernandez, Sonam Goenka, Wenny | 1133 |
| 1070 | owei Zhou, Jonathan Evens, William Isaac, Geoffrey | Yustalim, Robin Strudel, Ali Elqursh, Charlie Deck, | 1134 |
| 1071 | Irving, Edward Loper, Michael Fink, Isha Arkatkar, | Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoff- | 1135 |
| 1072 | Nanxin Chen, Izhak Shafran, Ivan Petrychenko, | mann, Dan Holtmann-Rice, Olivier Bachem, Sho | 1136 |
| 1073 | Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai | Arora, Christy Koh, Soheil Hassas Yeganeh, Siim | 1137 |
| 1074 | Zhu, Peter Grabowski, Yu Mao, Alberto Magni, | Pöder, Mukarram Tariq, Yanhua Sun, Lucian Ionita, | 1138 |
| 1075 | Kaisheng Yao, Javier Snaider, Norman Casagrande, | Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu Liu, An- | 1139 |
| 1076 | Evan Palmer, Paul Suganthan, Alfonso Castaño, | mol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz, | 1140 |
| 1077 | Irene Giannoumis, Wooyeon Kim, Mikołaj Rybiński, | Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, | 1141 |
| 1078 | Ashwin Sreevatsa, Jennifer Prendki, David Soergel, | Shreya Singh, Wei Fan, Aaron Parisi, Joe Stan- | 1142 |
| 1079 | Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, | ton, Vinod Koverkathu, Christopher A. Choquette- | 1143 |
| 1080 | Meenu Gaba, Jeremy Wiesner, Diana Gage Wright, | Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash | 1144 |
| 1081 | Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay | Shroff, Mani Varadarajan, Sanaz Bahargam, Rob | 1145 |
| 1082 | Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu, | Willoughby, David Gaddy, Guillaume Desjardins, | 1146 |
| 1083 | Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert | Marco Cornero, Brona Robenek, Bhavishya Mit- | 1147 |
| 1084 | Cui, Tian LIN, Marcus Wu, Ricardo Aguililar, Keith | tal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, | 1148 |
| 1085 | Pallo, Abhishek Chakladar, Ginger Perng, Elena Al- | Henrik Jacobsson, Alireza Ghaffarkhah, Morgane | 1149 |
| 1086 | llica Abellan, Mingyang Zhang, Ishita Dasgupta, | Rivièvre, Alanna Walton, Clément Crepy, Alicia | 1150 |
| 1087 | Nate Kushman, Ivo Penchev, Alena Repina, Xihui | Parish, Zongwei Zhou, Clement Farabet, Carey Rade- | 1151 |
| 1088 | Wu, Tom van der Weide, Priya Ponnappalli, Car- | baugh, Praveen Srinivasan, Claudia van der Salm, | 1152 |
| 1089 | oline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier | Andreas Fidjeland, Salvatore Scellato, Eri Latorre- | 1153 |
| 1090 | Dousse, Fan Yang, Jeff Piper, Nathan Ie, Rama Pa- | Chimoto, Hanna Klimczak-Plucińska, David Bridson, | 1154 |
| 1091 | sumarthi, Nathan Lintz, Anitha Vijayakumar, Daniel | Dario de Cesare, Tom Hudson, Piermaria Mendolic- | 1155 |
| 1092 | Andor, Pedro Valenzuela, Minnie Lui, Cosmin Padu- | chio, Lexi Walker, Alex Morris, Matthew Mauger, | 1156 |
| 1093 | raru, Daiyi Peng, Katherine Lee, Shuyuan Zhang, | Alexey Guseynov, Alison Reid, Seth Odoom, Lu- | 1157 |
| 1094 | Somer Greene, Duc Dung Nguyen, Paula Kurylew- | cia Loher, Victor Cotruta, Madhavi Yenugula, Do- | 1158 |
| 1095 | icz, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam | minik Grewe, Anastasia Petrushkina, Tom Duerig, | 1159 |
| 1096 | Choo, Ziqiang Feng, Biao Zhang, Achintya Sing- | Antonio Sanchez, Steve Yadlowsky, Amy Shen, | 1160 |
| 1097 | hal, Dayou Du, Dan McKinnon, Natasha Antropova, | Amir Globerson, Lynette Webb, Sahil Dua, Dong | 1161 |
| 1098 | Tolga Bolukbasi, Orgad Keller, David Reid, Daniel | Li, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, | 1162 |
| 1099 | Finchelstein, Maria Abi Raad, Remi Crocker, Pe- | Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj | 1163 |
| 1100 | ter Hawkins, Robert Dadashi, Colin Gaffney, Ken | Khare, Shreyas Rammohan Belle, Lei Wang, Chetan | 1164 |
| 1101 | Franko, Anna Bulanova, Rémi Leblond, Shirley | Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin | 1165 |
| 1102 | Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, | Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao | 1166 |
| 1103 | | Zhao, Stephan Lee, Pandu Nayak, Doug Fritz, Man- | 1167 |

| | | | | |
|------|---|------|--|------|
| 1168 | ish Reddy Vuyyuru, John Aslanides, Nidhi Vyas, Martin Wicke, Xiao Ma, Evgenii Eltyshev, Nina Martin, Hardie Cate, James Manyika, Keyvan Amiri, Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier, Nilesh Tripuraneni, David Madras, Mandy Guo, Austin Waters, Oliver Wang, Joshua Ainslie, Jason Baldridge, Han Zhang, Garima Pruthi, Jakob Bauer, Feng Yang, Riham Mansour, Jason Gelman, Yang Xu, George Polovets, Ji Liu, Honglong Cai, Warren Chen, XiangHai Sheng, Emily Xue, Sherjil Ozair, Christof Angermueller, Xiaowei Li, Anoop Sinha, Weiren Wang, Julia Wiesinger, Emmanouil Koukoumidis, Yuan Tian, Anand Iyer, Madhu Gurumurthy, Mark Goldenson, Parashar Shah, MK Blake, Hongkun Yu, Anthony Urbanowicz, Jennimaria Palomaki, Chrisantha Fernando, Ken Durden, Harsh Mehta, Nikola Momchev, Elahe Rahimtoroghi, Maria Georgaki, Amit Raul, Sebastian Ruder, Morgan Redshaw, Jinhyuk Lee, Denny Zhou, Komal Jalan, Dinghua Li, Blake Hechtman, Parker Schuh, Milad Nasr, Kieran Milan, Vladimir Mikulik, Juliana Franco, Tim Green, Nam Nguyen, Joe Kelley, Aroma Mahendru, Andrea Hu, Joshua Howland, Ben Vargas, Jeffrey Hui, Kshitij Bansal, Vikram Rao, Rakesh Ghiya, Emma Wang, Ke Ye, Jean Michel Sarr, Melanie Moranski Preston, Madeleine Elish, Steve Li, Aakash Kaku, Jigar Gupta, Ice Pasupat, Da-Cheng Juan, Milan Someswar, Tejvi M., Xinyun Chen, Aida Amini, Alex Fabrikant, Eric Chu, Xuanyi Dong, Amruta Muthal, Senaka Butthepitiya, Sarthak Jauhari, Nan Hua, Urvashi Khan-delwal, Ayal Hitron, Jie Ren, Larissa Rinaldi, Sharhar Drath, Avigail Dabush, Nan-Jiang Jiang, Harshal Godhia, Uli Sachs, Anthony Chen, Yicheng Fan, Hagai Taitelbaum, Hila Noga, Zhuyun Dai, James Wang, Chen Liang, Jenny Hamer, Chun-Sung Ferng, Chenel Elkind, Aviel Atias, Paulina Lee, Vít Listík, Mathias Carlen, Jan van de Kerkhof, Marcin Pilus, Krunoslav Zaher, Paul Müller, Sasha Zykova, Richard Stefanec, Vitaly Gatsko, Christoph Hirn-schall, Ashwin Sethi, Xingyu Federico Xu, Chetan Ahuja, Beth Tsai, Anca Stefanou, Bo Feng, Keshav Dhandhana, Manish Katyal, Akshay Gupta, Atharva Parulekar, Divya Pitta, Jing Zhao, Vivaan Bhatia, Yashodha Bhavnani, Omar Alhadlaq, Xiaolin Li, Peter Danenberg, Dennis Tu, Alex Pine, Vera Filippova, Abhipso Ghosh, Ben Limonchik, Bhar-gava Urala, Chaitanya Krishna Lanka, Derik Clive, Yi Sun, Edward Li, Hao Wu, Kevin Hongtongsak, Ianna Li, Kalind Thakkar, Kuanysh Omarov, Kushal Majmundar, Michael Alverson, Michael Kucharski, Mohak Patel, Mudit Jain, Maksim Zabelin, Paolo Pelagatti, Rohan Kohli, Saurabh Kumar, Joseph Kim, Swetha Sankar, Vineet Shah, Lakshmi Ramachandruni, Xiangkai Zeng, Ben Bariach, Laura Weidinger, Tu Vu, Alek Andreev, Antoine He, Kevin Hui, Sheleem Kashem, Amar Subramanya, Sissie Hsiao, Demis Hassabis, Koray Kavukcuoglu, Adam Sadovsky, Quoc Le, Trevor Strohman, Yonghui Wu, Slav Petrov, Jeffrey Dean, and Oriol Vinyals. 2024. <i>Gemini: A family of highly capable multimodal models.</i> Preprint, arXiv:2312.11805. | 1228 | Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-han Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-teng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhota, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-badur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Va-sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Shan-ran Narang, Sharath Raparth, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Van-denhenne, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Syd-ney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir- | 1294 |
| 1169 | | | 1293 | |
| 1170 | | | 1292 | |
| 1171 | | | 1291 | |
| 1172 | | | 1290 | |
| 1173 | | | 1289 | |
| 1174 | | | 1288 | |
| 1175 | | | 1287 | |
| 1176 | | | 1286 | |
| 1177 | | | 1285 | |
| 1178 | | | 1284 | |
| 1179 | | | 1283 | |
| 1180 | | | 1282 | |
| 1181 | | | 1281 | |
| 1182 | | | 1280 | |
| 1183 | | | 1279 | |
| 1184 | | | 1278 | |
| 1185 | | | 1277 | |
| 1186 | | | 1276 | |
| 1187 | | | 1275 | |
| 1188 | | | 1274 | |
| 1189 | | | 1273 | |
| 1190 | | | 1272 | |
| 1191 | | | 1271 | |
| 1192 | | | 1270 | |
| 1193 | | | 1269 | |
| 1194 | | | 1268 | |
| 1195 | | | 1267 | |
| 1196 | | | 1266 | |
| 1197 | | | 1265 | |
| 1198 | | | 1264 | |
| 1199 | | | 1263 | |
| 1200 | | | 1262 | |
| 1201 | | | 1261 | |
| 1202 | | | 1260 | |
| 1203 | | | 1259 | |
| 1204 | | | 1258 | |
| 1205 | | | 1257 | |
| 1206 | | | 1256 | |
| 1207 | | | 1255 | |
| 1208 | | | 1254 | |
| 1209 | | | 1253 | |
| 1210 | | | 1252 | |
| 1211 | | | 1251 | |
| 1212 | | | 1250 | |
| 1213 | | | 1249 | |
| 1214 | | | 1248 | |
| 1215 | | | 1247 | |
| 1216 | | | 1246 | |
| 1217 | | | 1245 | |
| 1218 | | | 1244 | |
| 1219 | | | 1243 | |
| 1220 | | | 1242 | |
| 1221 | | | 1241 | |
| 1222 | | | 1240 | |
| 1223 | | | 1239 | |
| 1224 | | | 1238 | |
| 1225 | | | 1237 | |
| 1226 | | | 1236 | |
| 1227 | | | 1235 | |
| 1228 | | | 1234 | |

| | | | |
|------|---|--|------|
| 1295 | ginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenber, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Chingsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazari, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-delwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madijan Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa- | tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghatham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaodu Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. <i>The llama 3 herd of models</i> . Preprint, arXiv:2407.21783. | 1359 |
| 1296 | | | 1360 |
| 1297 | | | 1361 |
| 1298 | | | 1362 |
| 1299 | | | 1363 |
| 1300 | | | 1364 |
| 1301 | | | 1365 |
| 1302 | | | 1366 |
| 1303 | | | 1367 |
| 1304 | | | 1368 |
| 1305 | | | 1369 |
| 1306 | | | 1370 |
| 1307 | | | 1371 |
| 1308 | | | 1372 |
| 1309 | | | 1373 |
| 1310 | | | 1374 |
| 1311 | | | 1375 |
| 1312 | | | 1376 |
| 1313 | | | 1377 |
| 1314 | | | 1378 |
| 1315 | | | 1379 |
| 1316 | | | 1380 |
| 1317 | | | 1381 |
| 1318 | | | 1382 |
| 1319 | | | 1383 |
| 1320 | | | 1384 |
| 1321 | | | 1385 |
| 1322 | | | 1386 |
| 1323 | | | 1387 |
| 1324 | | | 1388 |
| 1325 | | | 1389 |
| 1326 | | | 1390 |
| 1327 | | | 1391 |
| 1328 | | | 1392 |
| 1329 | | | 1393 |
| 1330 | | | 1394 |
| 1331 | | | 1395 |
| 1332 | | | 1396 |
| 1333 | | | 1397 |
| 1334 | | | 1398 |
| 1335 | | | 1399 |
| 1336 | | | 1400 |
| 1337 | | | 1401 |
| 1338 | | | 1402 |
| 1339 | | | 1403 |
| 1340 | | | 1404 |
| 1341 | | | 1405 |
| 1342 | | | 1406 |
| 1343 | | | 1407 |
| 1344 | | | 1408 |
| 1345 | | | 1409 |
| 1346 | | | 1410 |
| 1347 | | | 1411 |
| 1348 | | | 1412 |
| 1349 | | | 1413 |
| 1350 | | | 1414 |
| 1351 | | | 1415 |
| 1352 | | | 1416 |
| 1353 | | | 1417 |
| 1354 | | | 1418 |
| 1355 | | | 1419 |
| 1356 | | | 1419 |
| 1357 | | | 1419 |
| 1358 | | | 1419 |

| | | |
|------|--|------|
| 1420 | <i>of the Association for Computational Linguistics: EMNLP 2024</i> , pages 13097–13124, Miami, Florida, USA. Association for Computational Linguistics. | 1477 |
| 1421 | | 1478 |
| 1422 | | 1479 |
| 1423 | Hele-Andra Kuulmets, Taido Purason, Agnes Luhtaru, and Mark Fishel. 2024. <i>Teaching Llama a new language through cross-lingual knowledge transfer</i> . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 3309–3325, Mexico City, Mexico. Association for Computational Linguistics. | 1480 |
| 1424 | | |
| 1425 | | |
| 1426 | | |
| 1427 | | |
| 1428 | | |
| 1429 | | |
| 1430 | Viet Lai, Chien Nguyen, Nghia Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan Rossi, and Thien Nguyen. 2023a. <i>Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback</i> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 318–327, Singapore. Association for Computational Linguistics. | 1481 |
| 1431 | | 1482 |
| 1432 | | 1483 |
| 1433 | | 1484 |
| 1434 | | 1485 |
| 1435 | | 1486 |
| 1436 | | |
| 1437 | | |
| 1438 | | |
| 1439 | Viet Dac Lai, Nghia Trung Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Huu Nguyen. 2023b. <i>ChatGPT Beyond English: Towards a Comprehensive Evaluation of Large Language Models in Multilingual Learning</i> . Preprint, arXiv:2304.05613. | 1487 |
| 1440 | | 1488 |
| 1441 | | 1489 |
| 1442 | | 1490 |
| 1443 | | 1491 |
| 1444 | | 1492 |
| 1445 | Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. 2025. <i>Tulu 3: Pushing frontiers in open language model post-training</i> . Preprint, arXiv:2411.15124. | 1493 |
| 1446 | | |
| 1447 | | |
| 1448 | | |
| 1449 | | |
| 1450 | | |
| 1451 | | |
| 1452 | | |
| 1453 | | |
| 1454 | | |
| 1455 | Guillaume Lample and Alexis Conneau. 2019. <i>Cross-lingual Language Model Pretraining</i> . Preprint, arXiv:1901.07291. | 1498 |
| 1456 | | 1499 |
| 1457 | | 1500 |
| 1458 | Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020. <i>MLQA: Evaluating cross-lingual extractive question answering</i> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7315–7330, Online. Association for Computational Linguistics. | 1501 |
| 1459 | | 1502 |
| 1460 | | 1503 |
| 1461 | | 1504 |
| 1462 | | 1505 |
| 1463 | | 1506 |
| 1464 | | 1507 |
| 1465 | Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. 2023a. <i>Bactrian-X: Multilingual Replicable Instruction-Following Models with Low-Rank Adaptation</i> . Preprint, arXiv:2305.15011. | 1508 |
| 1466 | | 1509 |
| 1467 | | |
| 1468 | | |
| 1469 | Jiahuan Li, Hao Zhou, Shujian Huang, Shanbo Cheng, and Jiajun Chen. 2024. <i>Eliciting the translation ability of large language models via multilingual finetuning with translation instructions</i> . <i>Transactions of the Association for Computational Linguistics</i> , 12:576–592. | 1510 |
| 1470 | | 1511 |
| 1471 | | 1512 |
| 1472 | | 1513 |
| 1473 | | 1514 |
| 1474 | | 1515 |
| 1475 | Lei Li, Yuwei Yin, Shicheng Li, Liang Chen, Peiyi Wang, Shuhuai Ren, Mukai Li, Yazheng Yang, | 1516 |
| 1476 | | |
| 1477 | Jingjing Xu, Xu Sun, Lingpeng Kong, and Qi Liu. 2023b. <i>M\$^3IT: A Large-Scale Dataset towards Multi-Modal Multilingual Instruction Tuning</i> . Preprint, arXiv:2306.04387. | 1477 |
| 1478 | | 1478 |
| 1479 | | 1479 |
| 1480 | | 1480 |
| 1481 | Jiahua Liu, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2019. <i>XQA: A cross-lingual open-domain question answering dataset</i> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2358–2368, Florence, Italy. Association for Computational Linguistics. | 1481 |
| 1482 | | 1482 |
| 1483 | | 1483 |
| 1484 | | 1484 |
| 1485 | | 1485 |
| 1486 | | 1486 |
| 1487 | Yinquan Lu, Wenhao Zhu, Lei Li, Yu Qiao, and Fei Yuan. 2024. <i>LLaMAX: Scaling linguistic horizons of LLM by enhancing translation capabilities beyond 100 languages</i> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 10748–10772, Miami, Florida, USA. Association for Computational Linguistics. | 1487 |
| 1488 | | 1488 |
| 1489 | | 1489 |
| 1490 | | 1490 |
| 1491 | | 1491 |
| 1492 | | 1492 |
| 1493 | | 1493 |
| 1494 | Jirui Qi, Raquel Fernández, and Arianna Bisazza. 2023. <i>Cross-Lingual Consistency of Factual Knowledge in Multilingual Language Models</i> . Preprint, arXiv:2310.10378. | 1494 |
| 1495 | | 1495 |
| 1496 | | 1496 |
| 1497 | | 1497 |
| 1498 | Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. <i>Qwen2.5 technical report</i> . Preprint, arXiv:2412.15115. | 1498 |
| 1499 | | 1499 |
| 1500 | | 1500 |
| 1501 | | 1501 |
| 1502 | | 1502 |
| 1503 | | 1503 |
| 1504 | | 1504 |
| 1505 | | 1505 |
| 1506 | | 1506 |
| 1507 | | 1507 |
| 1508 | | 1508 |
| 1509 | | 1509 |
| 1510 | Leonardo Ranaldi, Giulia Pucci, and Andre Freitas. 2024. <i>Empowering cross-lingual abilities of instruction-tuned large language models by translation-following demonstrations</i> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 7961–7973, Bangkok, Thailand. Association for Computational Linguistics. | 1510 |
| 1511 | | 1511 |
| 1512 | | 1512 |
| 1513 | | 1513 |
| 1514 | | 1514 |
| 1515 | | 1515 |
| 1516 | | 1516 |
| 1517 | Uri Shaham, Jonathan Herzig, Roee Aharoni, Idan Szpektor, Reut Tsarfaty, and Matan Eyal. 2024. <i>Multilingual instruction tuning with just a pinch of multilinguality</i> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 2304–2317, Bangkok, Thailand. Association for Computational Linguistics. | 1517 |
| 1518 | | 1518 |
| 1519 | | 1519 |
| 1520 | | 1520 |
| 1521 | | 1521 |
| 1522 | | 1522 |
| 1523 | | 1523 |
| 1524 | Shuaijie She, Wei Zou, Shujian Huang, Wenhao Zhu, Xiang Liu, Xiang Geng, and Jiajun Chen. 2024. <i>MAPO: Advancing multilingual reasoning through multilingual-alignment-as-preference optimization</i> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10015–10027, Bangkok, Thailand. Association for Computational Linguistics. | 1524 |
| 1525 | | 1525 |
| 1526 | | 1526 |
| 1527 | | 1527 |
| 1528 | | 1528 |
| 1529 | | 1529 |
| 1530 | | 1530 |
| 1531 | | 1531 |
| 1532 | Asahi Ushio, Fernando Alva-Manchego, and Jose Camacho-Collados. 2023. <i>A Practical Toolkit</i> | 1532 |
| 1533 | | 1533 |

- 1534 for Multilingual Question and Answer Generation.
 1535 *Preprint*, arXiv:2305.17416.
- 1536 Elena Voita, Javier Ferrando, and Christoforos Nalm-
 1537 pantis. 2024. *Neurons in large language models: Dead, n-gram, positional*. In *Findings of the Asso-*
 1538 *ciation for Computational Linguistics: ACL 2024*,
 1539 pages 1288–1301, Bangkok, Thailand. Association
 1540 for Computational Linguistics.
- 1541 Elena Voita, Rico Sennrich, and Ivan Titov. 2021. *Ana-*
 1542 *lyzing the source and target contributions to predic-*
 1543 *tions in neural machine translation*. In *Proceedings*
 1544 *of the 59th Annual Meeting of the Association for*
 1545 *Computational Linguistics and the 11th International*
 1546 *Joint Conference on Natural Language Processing*
 1547 *(Volume 1: Long Papers)*, pages 1126–1140, Online.
 1548 Association for Computational Linguistics.
- 1549 Bin Wang, Zhengyuan Liu, Xin Huang, Fangkai Jiao,
 1550 Yang Ding, Ai Ti Aw, and Nancy F. Chen. 2023.
 1551 *SeaEval for Multilingual Foundation Models: From*
 1552 *Cross-Lingual Alignment to Cultural Reasoning*.
 1553 *Preprint*, arXiv:2309.04766.
- 1554 Bin Wang, Zhengyuan Liu, Xin Huang, Fangkai Jiao,
 1555 Yang Ding, AiTi Aw, and Nancy Chen. 2024. *SeaE-*
 1556 *val for multilingual foundation models: From cross-*
 1557 *lingual alignment to cultural reasoning*. In *Proce-*
 1558 *edings of the 2024 Conference of the North Ameri-*
 1559 *can Chapter of the Association for Computational Lin-*
 1560 *guistics: Human Language Technologies (Volume 1:*
 1561 *Long Papers)*, pages 370–390, Mexico City, Mexico.
 1562 Association for Computational Linguistics.
- 1563 Xiangpeng Wei, Haoran Wei, Huan Lin, Tianhao Li,
 1564 Pei Zhang, Xingzhang Ren, Mei Li, Yu Wan, Zhi-
 1565 wei Cao, Binbin Xie, Tianxiang Hu, Shangjie Li,
 1566 Binyuan Hui, Bowen Yu, Dayiheng Liu, Baosong
 1567 Yang, Fei Huang, and Jun Xie. 2023. *PolyLM:*
 1568 *An Open Source Polyglot Large Language Model*.
 1569 *Preprint*, arXiv:2307.06018.
- 1570 Junyi Xiang, Maofu Liu, Qiyuan Li, Chen Qiu, and
 1571 Huijun Hu. 2024. *A cross-guidance cross-lingual*
 1572 *model on generated parallel corpus for classical Chi-*
 1573 *nese machine reading comprehension*. *Information*
 1574 *Processing & Management*, 61(2):103607.
- 1575 Weiwen Xu, Xin Li, Wai Lam, and Lidong Bing. 2023.
 1576 *mPMR: A Multilingual Pre-trained Machine Reader*
 1577 *at Scale*. *Preprint*, arXiv:2305.13645.
- 1578 Sen Yang, Shujian Huang, Wei Zou, Jianbing Zhang,
 1579 Xinyu Dai, and Jiajun Chen. 2023a. *Local inter-
 1580 pretation of transformer based on linear decomposi-*
 1581 *tion*. In *Proceedings of the 61st Annual Meeting of the*
 1582 *Association for Computational Linguistics (Volume 1:*
 1583 *Long Papers)*, pages 10270–10287, Toronto, Canada.
 1584 Association for Computational Linguistics.
- 1585 Wen Yang, Chong Li, Jiajun Zhang, and Chengqing
 1586 Zong. 2023b. *BigTranslate: Augmenting Large Lan-*
 1587 *guage Models with Multilingual Translation Capabil-*
 1588 *ity over 100 Languages*. *Preprint*, arXiv:2305.18098.

| Name | Modes | Sizes |
|--------------------|-----------------|----------------|
| LLaMA 2 | Base / Chat | 7B / 13B / 70B |
| LLaMA 3 | Base / Instruct | 8B / 70B |
| LLaMA 3.1 | Base / Instruct | 8B / 70B |
| LLaMAX-2-Alpaca | - | 7B |
| LLaMAX-3-Alpaca | - | 8B |
| Mistral V0.1 | Base / Instruct | 7B |
| Mistral V0.3 | Base / Instruct | 7B |
| Qwen 1.5 | Base / Chat | 7B / 14B / 72B |
| Qwen 2 | Base / Instruct | 7B / 72B |
| Qwen 2.5 | Base / Instruct | 7B / 72B |
| DeepSeek V2 | Base / Chat | Lite (16B) |
| Gemma 2 | Base / IT | 9B |
| GPT-3.5-Turbo-0125 | - | - |
| GPT-4o | - | - |

Table 4: Full list of models evaluated. This table presents a complete list of all models tested in this study, encompassing older versions and alternative sizes.

Kexun Zhang, Jane Dwivedi-Yu, Zhaojiang Lin, Yuning
 1590 Mao, William Yang Wang, Lei Li, and Yi-Chia Wang.
 1591 2025. *Extrapolating to unknown opinions using*
 1592 *LLMs*. In *Proceedings of the 31st International Con-*
 1593 *ference on Computational Linguistics*, pages 7819–
 1594 7830, Abu Dhabi, UAE. Association for Computa-
 1595 tional Linguistics.

Shaolei Zhang, Qingkai Fang, Zhuocheng Zhang, Zhen-
 1596 grui Ma, Yan Zhou, Langlin Huang, Mengyu Bu,
 1597 Shangtong Gui, Yunji Chen, Xilin Chen, and Yang
 1598 Feng. 2023. *BayLing: Bridging Cross-lingual Align-
 1599 ment and Instruction Following through Interactive*
 1600 *Translation for Large Language Models*. *Preprint*,
 1601 arXiv:2306.10968.

Wenhao Zhu, Shujian Huang, Fei Yuan, Shuaijie She,
 1602 Jiajun Chen, and Alexandra Birch. 2024. *Question*
 1603 *translation training for better multilingual reasoning*.
 1604 In *Findings of the Association for Computational Lin-*
 1605 *guistics: ACL 2024*, pages 8411–8423, Bangkok,
 1606 Thailand. Association for Computational Linguistics.

A Evaluation

A.1 Full list of models evaluated

A comprehensive list of all models evaluated during
 1612 this study, including older versions and alternative
 1613 sizes, is available in Table 4 to supplement the main
 1614 list.

A.2 Prompt used in error type detection

The prompt for error type detection is:

<|im_start|>system
 You are Qwen, created by Alibaba Cloud. You are a helpful assistant.<|im_end|>
 <|im_start|>user
 You are tasked with identifying the type of a given raw answer. You will be provided with a question and a raw answer. Your job is to determine whether the raw answer falls into one of the following categories based on the given question:

0. Reasonable Answer: The answer seems like some attempt to answer the question, regardless of whether it is correct or not.
1. Blank Answer: No response is provided.
2. Gibberish: Incoherent text with no clear meaning or cannot be seen as some kind of answer to the question, e.g. “{Your Answer}”.
3. Denial of Answer: A statement indicating inability to answer, such as “I apologize, but I cannot answer this question because...”.

You must provide your response as a SINGLE number representing the category (0, 1, 2, or 3) without extra output.

1619 The v1 prompt format is:

{system prompt}

Below is a reading comprehension task. There will be paragraphs of context, each followed by a question related to its content. You should only present your answer to the last question by strictly copying the corresponding part of the context. Please provide a direct answer in English without extra output. Your answer should be in the form of “Answer: {Your Answer}”

Context: {demo context 1}

Question: {demo question 1}

Answer: {demo answer 1}

Context: {demo context 2}

Question: {demo question 2}

Answer: {demo answer 2}

Your task starts here:

Context: {text context}

Question: {text question}

1621 The v2 prompt format is:

{system prompt}

Context: {demo context 1}

Question: {demo question 1}

Answer: {demo answer 1}

Context: {demo context 2}

Question: {demo question 2}

Answer: {demo answer 2}

Your task starts here:

Context: {text context}

Question: {text question}

You should only present your answer to the last question by strictly copying the corresponding part of the context. Please provide a direct answer in English without extra output. Your answer should be in the form of “Answer: {Your Answer}”

A.3 Detailed evaluation results

Table 5 presents the complete evaluation results from our 0-shot experiments, encompassing both the English-to-NonEnglish (en-x) and non-English monolingual (x-x) tasks in all models and languages tested. Table 6 further presents detailed F1 scores for en-x and x-x tasks in the 2-shot setting.

A.4 Detailed language error and generation failure error rates

Tables 7 and 8 show detailed language and generation failure error rates across all tested languages. Meanwhile, Table 9 provides a more granular view of the generation failure errors discussed in the main text. While Table 8 presents the aggregated rate of these errors, Table 9 is further subdivided into three separate tables: Table 9a, Table 9b, and Table 9c. These tables individually display the error rates for gibberish errors, refusal errors, and blank errors, respectively, across all tested models and languages in the 2-shot en-x xMRC task setting.

B Two-phased xMRC Analysis

B.1 Further analysis on MRD

B.1.1 Example of Attribution

Figure 8 shows an example of the attribution outcome for LLaMA-3.1-Instruct-8B.

B.1.2 MRD for other LLaMA models

Figures 9 and 10 provide further illustrative examples of the mean MRD for context and question

1622

1623

1624

1625

1626

1627

1628

1629

1630

1631

1632

1633

1634

1635

1636

1637

1638

1639

1640

1641

1642

1643

1644

1645

1646

1647

1648

1649

1650

1651

1652

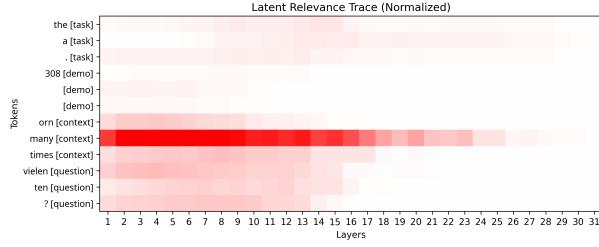


Figure 8: An example output of layer-wise attribution with LLaMA-3.1-Instruct-8B, where only the top 3 tokens from each input part are shown.

components, specifically for LLaMA-3.1-Instruct-70B and LLaMA-2-Chat-7B. These figures complement the MRD analysis presented in the main body of this paper.

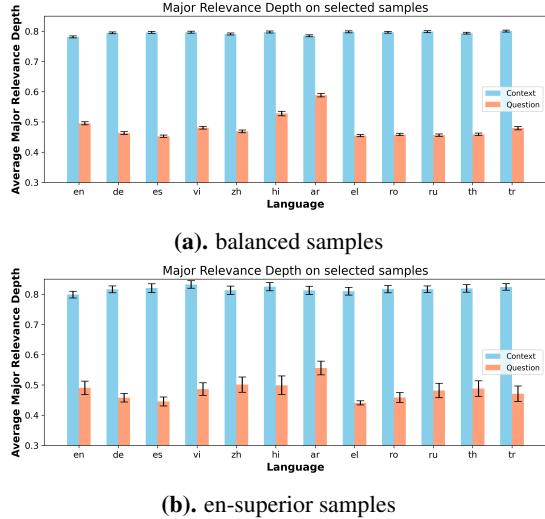


Figure 9: Mean MRD of the context and question parts for LLaMA-3.1-Instruct-70B.

B.1.3 Analysis of task descriptions and demonstrations

Analyzing the MRD of task descriptions and demonstrations in our 2-shot setting (Figures 11–13) reveals a general trend where demonstrations tend to exhibit a comparable or slightly higher MRD than task descriptions across the LLaMA model family, suggesting demonstrations are at least as important as, if not slightly more impactful than, task descriptions in guiding the models. This could indicate that providing concrete examples is a particularly effective way to communicate the desired behavior for cross-lingual context retrieval to these models.

However, the precise relationship is not uniform and varies across models. For example, while LLaMA-3.1-Instruct-8B shows a relatively bal-

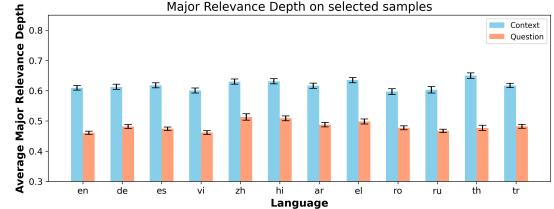


Figure 10: Mean MRD of the context and question parts for LLaMA-2-Chat-7B.

anced MRD between task descriptions and demonstrations, LLaMA-2-Chat-7B consistently demonstrates a higher MRD for demonstrations, which implies that older or smaller models might lean more heavily on the provided in-context examples. In contrast, LLaMA-3.1-Instruct-70B exhibits the most pronounced difference, with a significantly elevated MRD for task descriptions across all languages and sample types, suggesting that larger models can become highly attuned to and reliant on user-specified task commands.

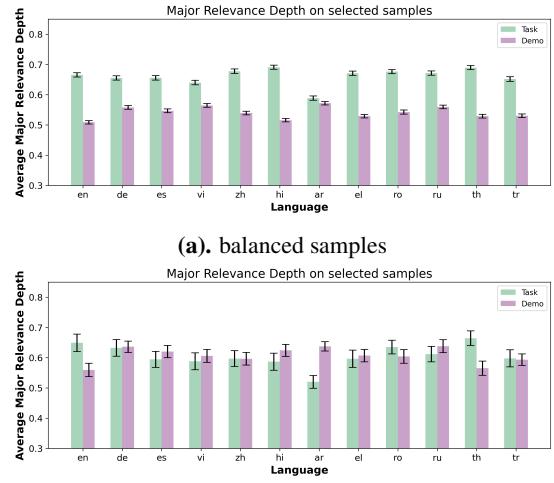


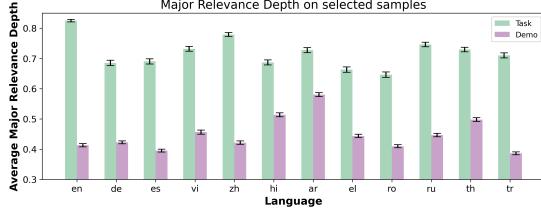
Figure 11: Mean MRD of the task descriptions and demonstrations parts for LLaMA-3.1-Instruct-8B.

B.1.4 Influence of prompt format on MRD pattern

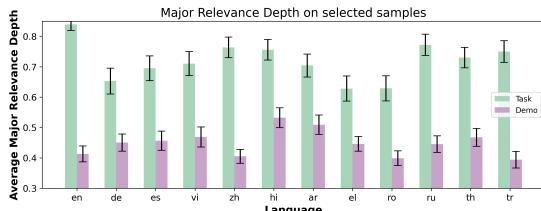
We test the influence of different prompt formats (v1, v2) on LLaMA-3.1-Instruct-8B, and by com-

1653
1654
1655
1656

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684

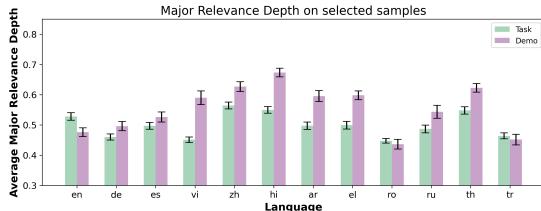


(a). balanced samples

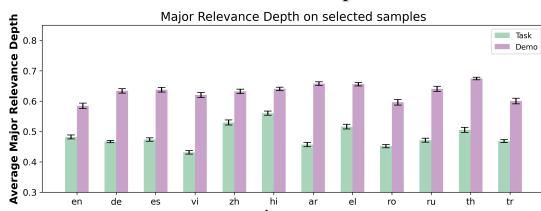


(b). en-superior samples

Figure 12: Mean MRD of the task descriptions and demonstrations parts for LLaMA-3.1-Instruct-70B.



(a). balanced samples



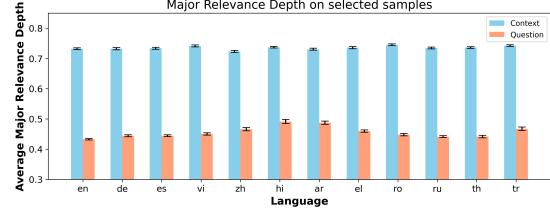
(b). en-superior samples

Figure 13: Mean MRD of the task descriptions and demonstrations parts for LLaMA-2-Chat-7B.

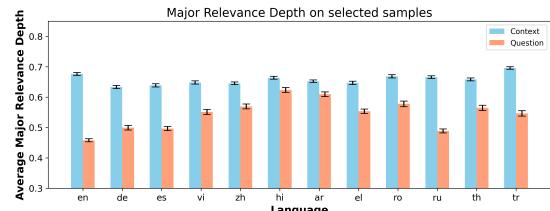
paring the results in Figure 14 and Figure 3, which present results obtained using the prompt format v1 and v2, respectively, it is clear that the fundamental pattern observed in the mean MRD is consistent across both formats. Therefore, the trend of the mean question MRD being consistently and substantially lower than the mean context MRD is maintained regardless of the prompt format employed.

B.2 Hidden state similarity results for other LLaMA models

Figures 15–18 present the hidden state similarity results for additional LLaMA models, complementing the analysis of the LLaMA-3.1-Instruct-8B model discussed in the main body of the paper.



(a). balanced samples



(b). en-superior samples

Figure 14: Mean MRD for LLaMA-3.1-Instruct-8B on both "balanced" and "en-superior" samples in v1 prompting format. Only the results of context and question parts of the prompt are displayed.

Meanwhile, Figure 19 provides an overview of the last-input-token hidden state similarity for base LLaMA models, as discussed in §5.3.

B.3 Training details and evaluation results of our finetuned LLaMA-3.1-8B

We tune the LLaMA-3.1-8B base model on TULU-V3 for 1 epoch with 8 * H800 GPUs for 15 hours using the LLaMA-Factory repository. The data cut-off length is 2048, batch size per device is 8, learning rate is 1.0e-5, and the warm-up ratio is 0.1 with cosine learning rate scheduling.

Regarding evaluation, Table 10 summarizes the performance of our finetuned model on both en-x cross-lingual and x-x monolingual MRC tasks. Furthermore, Figure 20 illustrates the hidden state similarity between English and other tested languages across layers, focusing on question, context, and last-input-token representations derived from balanced samples.

1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722

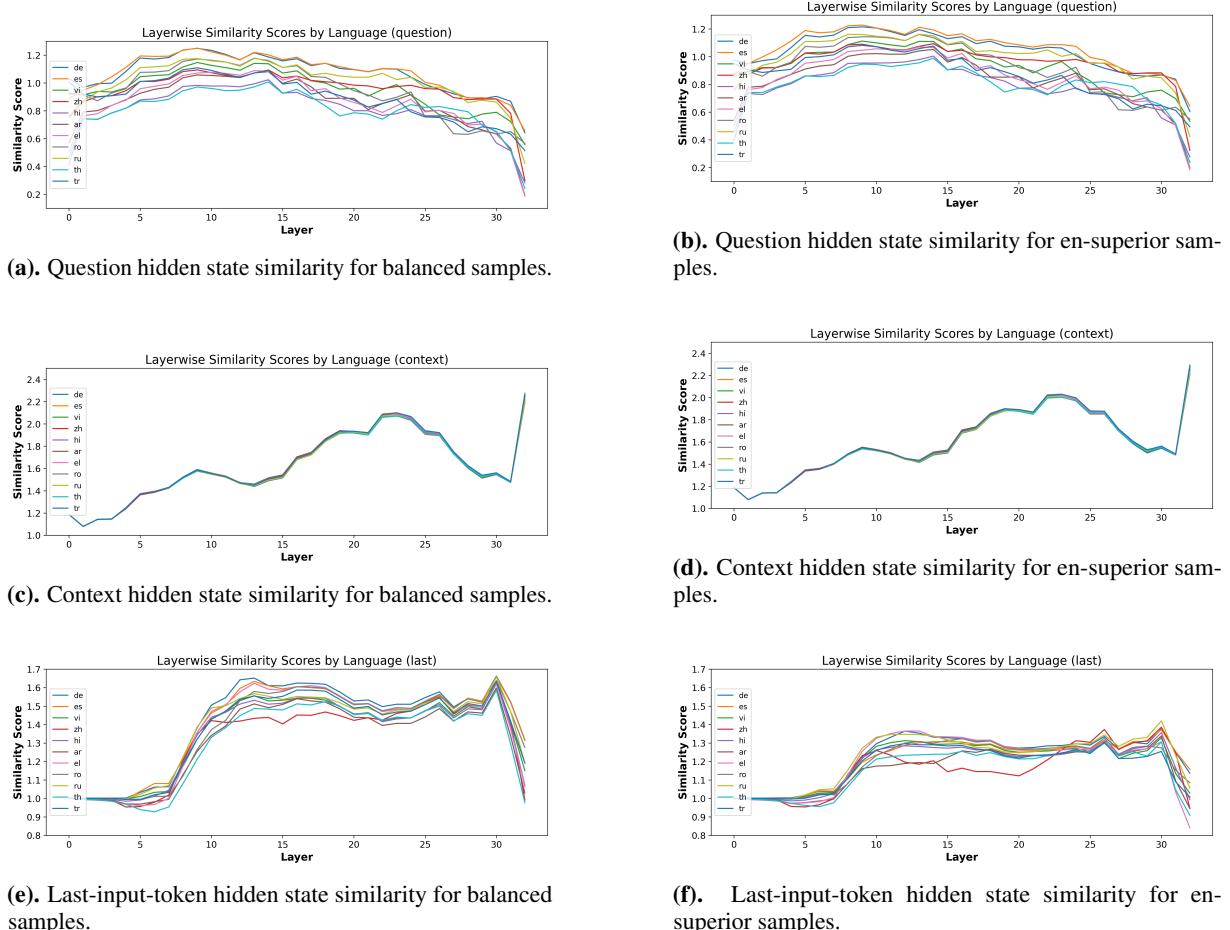


Figure 15: Hidden state similarity between English and other languages on different parts of the selected samples in each layer of the LLaMA-3.1-8B model.

| | en-en | en-de | en-es | en-vi | en-zh | en-hi | en-ar | en-el | en-ro | en-ru | en-th | en-tr |
|----------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| LLaMA-3.1-8B | 28.49 | 26.39 | 23.70 | 16.97 | 21.76 | 21.68 | 19.66 | 24.96 | 22.27 | 22.69 | 27.73 | 15.71 |
| LLaMA-3.1-70B | 49.92 | 38.32 | 36.22 | 41.34 | 45.13 | 24.62 | 27.40 | 37.73 | 38.49 | 34.87 | 47.90 | 33.95 |
| Mistral-V0.3-7B | 20.50 | 21.43 | 19.50 | 10.84 | 8.66 | 11.60 | 14.87 | 14.20 | 12.44 | 18.32 | 16.13 | 10.67 |
| Qwen-2.5-7B | 38.66 | 36.05 | 37.31 | 48.32 | 42.02 | 38.91 | 42.44 | 34.87 | 35.55 | 38.57 | 45.63 | 41.60 |
| Qwen-2.5-72B | 63.95 | 56.55 | 56.72 | 53.19 | 52.86 | 55.97 | 56.30 | 52.27 | 54.03 | 55.55 | 55.29 | 51.68 |
| DeepSeek-V2-Lite-16B | 12.35 | 11.26 | 4.12 | 4.20 | 9.58 | 9.58 | 4.79 | 7.73 | 7.06 | 10.84 | 5.80 | 4.96 |
| Gemma-2-9B | 15.21 | 5.97 | 14.87 | 14.87 | 26.47 | 7.40 | 11.43 | 4.62 | 13.70 | 20.42 | 32.27 | 17.31 |
| LLaMA-2-Chat-7B | 34.96 | 26.81 | 24.54 | 19.50 | 22.52 | 19.58 | 20.00 | 20.50 | 22.10 | 25.04 | 13.95 | 16.55 |
| LLaMA-3.1-Instruct-8B | 40.92 | 25.55 | 28.82 | 28.24 | 33.95 | 27.65 | 24.87 | 20.42 | 19.24 | 29.66 | 29.75 | 30.42 |
| LLaMA-3.1-Instruct-70B | 57.82 | 47.23 | 47.23 | 45.88 | 49.16 | 39.33 | 42.86 | 46.72 | 43.19 | 41.51 | 51.93 | 44.20 |
| Mistral-V0.3-Instruct-7B | 3.19 | 2.61 | 2.69 | 5.13 | 2.86 | 2.02 | 4.12 | 3.11 | 2.35 | 3.87 | 4.54 | 3.36 |
| Qwen-2.5-Instruct-7B | 53.28 | 40.17 | 35.97 | 36.13 | 40.67 | 36.22 | 35.46 | 33.70 | 38.99 | 35.71 | 36.22 | 38.24 |
| Qwen-2.5-Instruct-72B | 36.89 | 27.98 | 26.81 | 23.53 | 23.28 | 22.94 | 23.45 | 23.19 | 31.01 | 24.37 | 23.03 | 23.95 |
| DeepSeek-V2-Chat-Lite-16B | 16.30 | 18.24 | 13.87 | 8.24 | 15.13 | 11.01 | 11.09 | 13.95 | 13.87 | 12.61 | 9.75 | 12.61 |
| Gemma-2-IT-9B | 57.06 | 46.30 | 42.77 | 42.86 | 42.35 | 44.37 | 40.25 | 41.85 | 39.24 | 42.69 | 44.96 | 43.87 |
| GPT-3.5-Turbo-0125 | 32.18 | 24.12 | 22.61 | 21.34 | 21.93 | 17.48 | 22.44 | 23.70 | 23.87 | 21.43 | 20.34 | 20.42 |
| GPT-4o | 51.01 | 39.58 | 36.97 | 40.50 | 41.18 | 37.48 | 36.05 | 37.82 | 36.97 | 39.33 | 39.16 | 39.66 |

(a). 0-shot Exact Match (EM) scores (%) on en-x tasks.

| | en-en | en-de | en-es | en-vi | en-zh | en-hi | en-ar | en-el | en-ro | en-ru | en-th | en-tr |
|----------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| LLaMA-3.1-8B | 37.60 | 34.11 | 33.87 | 21.86 | 35.12 | 36.97 | 30.11 | 36.66 | 30.82 | 30.37 | 39.43 | 21.79 |
| LLaMA-3.1-70B | 68.03 | 58.29 | 56.22 | 60.33 | 55.20 | 50.63 | 41.23 | 56.41 | 57.13 | 55.87 | 62.87 | 54.32 |
| Mistral-V0.3-7B | 40.07 | 27.17 | 31.26 | 19.26 | 13.19 | 17.94 | 23.98 | 19.16 | 18.23 | 28.26 | 21.86 | 18.63 |
| Qwen-2.5-7B | 59.02 | 55.58 | 55.61 | 64.37 | 60.26 | 56.32 | 60.61 | 53.14 | 51.21 | 56.72 | 63.23 | 55.52 |
| Qwen-2.5-72B | 80.50 | 73.73 | 74.48 | 72.48 | 71.37 | 73.60 | 73.23 | 70.59 | 71.44 | 73.36 | 72.78 | 68.18 |
| DeepSeek-V2-Lite-16B | 24.88 | 25.14 | 11.07 | 13.25 | 18.55 | 14.01 | 13.69 | 15.32 | 14.72 | 23.10 | 9.28 | 13.73 |
| Gemma-2-9B | 24.08 | 11.12 | 24.96 | 20.17 | 36.04 | 10.03 | 16.04 | 8.68 | 21.19 | 32.67 | 43.05 | 24.34 |
| LLaMA-2-Chat-7B | 56.83 | 45.97 | 44.45 | 37.78 | 36.50 | 33.65 | 32.13 | 34.03 | 39.14 | 45.51 | 25.18 | 30.64 |
| LLaMA-3.1-Instruct-8B | 64.47 | 50.69 | 53.41 | 52.33 | 57.10 | 50.09 | 48.61 | 45.36 | 43.46 | 54.14 | 53.88 | 52.72 |
| LLaMA-3.1-Instruct-70B | 78.28 | 70.13 | 68.39 | 64.43 | 68.91 | 56.20 | 60.19 | 67.79 | 65.78 | 64.72 | 67.69 | 62.13 |
| Mistral-V0.3-Instruct-7B | 35.19 | 32.41 | 32.68 | 30.23 | 32.15 | 28.49 | 30.21 | 29.42 | 28.93 | 32.51 | 31.48 | 29.23 |
| Qwen-2.5-Instruct-7B | 73.03 | 59.69 | 56.82 | 56.99 | 60.64 | 56.99 | 55.54 | 53.93 | 58.27 | 56.62 | 57.08 | 59.22 |
| Qwen-2.5-Instruct-72B | 59.84 | 46.79 | 46.86 | 43.46 | 36.20 | 43.82 | 40.01 | 44.46 | 50.91 | 43.98 | 44.68 | 44.28 |
| DeepSeek-V2-Chat-Lite-16B | 43.24 | 35.99 | 32.46 | 26.34 | 38.20 | 27.90 | 29.16 | 35.33 | 29.12 | 32.53 | 30.33 | 29.19 |
| Gemma-2-IT-9B | 76.80 | 67.91 | 64.91 | 64.74 | 64.83 | 65.31 | 62.35 | 63.89 | 61.29 | 64.94 | 65.84 | 64.10 |
| GPT-3.5-Turbo-0125 | 60.08 | 51.06 | 50.45 | 47.23 | 49.06 | 41.58 | 48.90 | 50.57 | 50.90 | 49.32 | 46.66 | 46.22 |
| GPT-4o | 74.24 | 64.38 | 58.94 | 65.71 | 65.40 | 62.48 | 58.81 | 64.00 | 62.15 | 64.82 | 64.60 | 64.95 |

(b). 0-shot F1 Scores on en-x tasks.

| | de-de | es-es | vi-vi | zh-zh | hi-hi | ar-ar | el-el | ro-ro | ru-ru | th-th | tr-tr |
|----------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| LLaMA-3.1-8B | 24.87 | 19.16 | 18.66 | 33.95 | 16.05 | 17.31 | 10.00 | 18.15 | 17.90 | 33.87 | 14.37 |
| LLaMA-3.1-70B | 42.94 | 37.23 | 42.86 | 50.08 | 30.00 | 40.17 | 36.30 | 40.76 | 35.13 | 57.06 | 37.40 |
| Mistral-V0.3-7B | 27.31 | 25.29 | 20.17 | 34.62 | 7.98 | 21.09 | 17.73 | 22.69 | 11.93 | 23.45 | 12.44 |
| Qwen-2.5-7B | 29.50 | 35.04 | 38.57 | 57.23 | 23.11 | 42.94 | 21.09 | 31.85 | 32.18 | 53.95 | 31.51 |
| Qwen-2.5-72B | 46.89 | 46.39 | 51.93 | 78.32 | 41.18 | 50.67 | 36.64 | 50.67 | 43.03 | 63.78 | 41.34 |
| DeepSeek-V2-Lite-16B | 4.87 | 2.44 | 1.01 | 7.56 | 2.02 | 1.43 | 4.03 | 3.19 | 2.02 | 5.13 | 2.94 |
| Gemma-2-9B | 2.02 | 14.03 | 8.91 | 12.10 | 8.66 | 1.43 | 1.43 | 4.62 | 10.92 | 9.66 | 13.28 |
| LLaMA-2-Chat-7B | 22.86 | 18.82 | 15.55 | 4.54 | 1.68 | 4.79 | 6.05 | 22.86 | 11.93 | 3.78 | 10.59 |
| LLaMA-3.1-Instruct-8B | 25.71 | 26.97 | 36.64 | 48.40 | 27.06 | 33.03 | 13.70 | 28.32 | 26.22 | 35.63 | 29.83 |
| LLaMA-3.1-Instruct-70B | 40.25 | 35.29 | 47.31 | 57.65 | 35.71 | 44.96 | 30.92 | 40.76 | 38.66 | 59.50 | 40.84 |
| Mistral-V0.3-Instruct-7B | 2.69 | 2.44 | 4.29 | 0.92 | 0.84 | 3.70 | 1.51 | 3.11 | 1.43 | 4.29 | 2.61 |
| Qwen-2.5-Instruct-7B | 32.86 | 30.92 | 30.42 | 47.40 | 24.71 | 27.90 | 21.68 | 36.05 | 27.82 | 42.44 | 30.67 |
| Qwen-2.5-Instruct-72B | 29.41 | 24.45 | 26.97 | 45.71 | 20.42 | 32.18 | 15.88 | 32.02 | 25.29 | 36.30 | 21.51 |
| DeepSeek-V2-Chat-Lite-16B | 12.18 | 7.48 | 10.08 | 7.82 | 7.65 | 9.24 | 8.40 | 7.73 | 9.41 | 11.93 | 7.31 |
| Gemma-2-IT-9B | 42.94 | 40.92 | 46.97 | 51.09 | 41.93 | 43.03 | 37.98 | 45.97 | 42.52 | 56.13 | 36.30 |
| GPT-3.5-Turbo-0125 | 23.53 | 23.19 | 30.50 | 38.32 | 25.71 | 28.74 | 24.03 | 27.31 | 25.80 | 39.41 | 26.47 |
| GPT-4o | 40.34 | 34.87 | 44.37 | 54.29 | 32.10 | 42.10 | 26.22 | 38.94 | 37.82 | 52.79 | 30.59 |

(c). 0-shot Exact Match (EM) scores (%) on x-x tasks

| | de-de | es-es | vi-vi | zh-zh | hi-hi | ar-ar | el-el | ro-ro | ru-ru | th-th | tr-tr |
|----------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| LLaMA-3.1-8B | 37.62 | 34.52 | 36.04 | 43.66 | 41.18 | 34.54 | 32.58 | 34.61 | 35.68 | 47.38 | 28.02 |
| LLaMA-3.1-70B | 62.66 | 59.02 | 61.54 | 55.26 | 56.81 | 63.43 | 56.52 | 59.81 | 55.99 | 70.20 | 57.39 |
| Mistral-V0.3-7B | 41.31 | 47.37 | 39.40 | 41.73 | 24.04 | 42.28 | 35.43 | 37.91 | 29.51 | 37.70 | 28.32 |
| Qwen-2.5-7B | 49.22 | 58.04 | 62.82 | 68.32 | 46.19 | 63.63 | 44.78 | 52.10 | 52.96 | 68.59 | 54.30 |
| Qwen-2.5-72B | 70.77 | 72.81 | 74.95 | 84.33 | 68.09 | 72.94 | 64.72 | 73.17 | 68.56 | 78.50 | 67.79 |
| DeepSeek-V2-Lite-16B | 12.56 | 8.24 | 7.64 | 9.73 | 8.35 | 6.98 | 9.59 | 9.39 | 7.79 | 9.03 | 8.89 |
| Gemma-2-9B | 4.16 | 20.57 | 13.66 | 14.74 | 15.58 | 2.53 | 3.07 | 6.96 | 18.84 | 14.14 | 23.02 |
| LLaMA-2-Chat-7B | 40.12 | 40.37 | 32.22 | 10.68 | 13.62 | 17.59 | 18.83 | 38.52 | 26.51 | 10.64 | 23.93 |
| LLaMA-3.1-Instruct-8B | 51.95 | 58.67 | 62.55 | 54.15 | 52.86 | 56.45 | 45.21 | 55.02 | 52.99 | 56.89 | 55.40 |
| LLaMA-3.1-Instruct-70B | 68.60 | 68.79 | 73.54 | 66.32 | 62.66 | 70.97 | 66.97 | 68.90 | 65.53 | 73.91 | 66.35 |
| Mistral-V0.3-Instruct-7B | 24.81 | 26.71 | 26.83 | 11.22 | 15.57 | 20.21 | 22.81 | 25.11 | 20.37 | 30.99 | 21.98 |
| Qwen-2.5-Instruct-7B | 57.65 | 57.89 | 56.71 | 57.29 | 50.26 | 53.38 | 50.18 | 59.65 | 53.25 | 63.81 | 55.85 |
| Qwen-2.5-Instruct-72B | 54.08 | 51.89 | 52.80 | 57.27 | 45.64 | 57.66 | 45.27 | 56.80 | 50.74 | 62.11 | 48.48 |
| DeepSeek-V2-Chat-Lite-16B | 35.78 | 34.43 | 34.37 | 15.75 | 25.36 | 29.23 | 33.21 | 32.04 | 34.46 | 28.91 | 28.62 |
| Gemma-2-IT-9B | 68.38 | 68.84 | 71.81 | 65.03 | 67.91 | 66.82 | 67.78 | 69.48 | 66.19 | 73.71 | 64.21 |
| GPT-3.5-Turbo-0125 | 53.58 | 56.26 | 58.50 | 52.01 | 50.36 | 57.77 | 58.48 | 56.60 | 55.95 | 57.45 | 54.36 |
| GPT-4o | 67.90 | 67.79 | 71.45 | 65.72 | 61.33 | 69.67 | 62.12 | 66.61 | 66.59 | 74.43 | 62.68 |

(d). 0-shot F1 Scores on x-x tasks.

Table 5: 0-shot evaluation results on en-x and x-x tasks.

| | en-en | en-de | en-es | en-vi | en-zh | en-hi | en-ar | en-el | en-ro | en-ru | en-th | en-tr |
|----------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| LLaMA-3.1-8B | 75.97 | 45.97 | 43.90 | 50.71 | 47.70 | 48.41 | 50.98 | 50.40 | 47.22 | 49.12 | 59.77 | 44.94 |
| LLaMA-3.1-70B | 82.39 | 60.46 | 60.00 | 59.57 | 62.68 | 57.38 | 51.50 | 59.43 | 54.28 | 56.15 | 66.54 | 57.48 |
| Mistral-V0.3-7B | 79.57 | 66.94 | 67.57 | 59.82 | 53.66 | 48.36 | 52.56 | 47.83 | 67.30 | 70.62 | 48.56 | 62.94 |
| Qwen-2.5-7B | 62.42 | 56.68 | 56.45 | 59.15 | 58.84 | 55.85 | 62.94 | 50.62 | 54.43 | 58.37 | 61.61 | 57.72 |
| Qwen-2.5-72B | 86.03 | 77.24 | 79.22 | 80.16 | 80.14 | 79.09 | 78.70 | 76.72 | 80.41 | 80.77 | 78.48 | 77.16 |
| DeepSeek-V2-Lite-16B | 73.81 | 46.34 | 47.65 | 52.75 | 41.77 | 34.45 | 44.61 | 46.18 | 50.39 | 51.91 | 34.58 | 40.57 |
| Gemma-2-9B | 80.42 | 60.13 | 61.74 | 64.76 | 71.33 | 64.55 | 68.84 | 75.81 | 65.49 | 72.09 | 70.59 | 59.72 |
| LLaMA-3.1-Instruct-8B | 77.89 | 74.81 | 73.50 | 73.29 | 72.78 | 68.59 | 69.66 | 70.41 | 72.24 | 73.40 | 73.60 | 71.20 |
| LLaMA-3.1-Instruct-70B | 83.29 | 73.58 | 72.98 | 73.97 | 73.79 | 70.87 | 75.96 | 72.91 | 71.69 | 72.73 | 75.30 | 70.00 |
| Mistral-V0.3-Instruct-7B | 62.01 | 59.06 | 60.98 | 54.81 | 58.84 | 47.16 | 60.51 | 57.49 | 59.80 | 60.81 | 55.63 | 47.86 |
| Qwen-2.5-Instruct-7B | 81.83 | 77.37 | 77.02 | 76.06 | 78.82 | 73.70 | 76.11 | 74.70 | 76.80 | 77.44 | 77.05 | 75.69 |
| Qwen-2.5-Instruct-72B | 77.12 | 67.65 | 68.89 | 64.34 | 52.67 | 70.35 | 59.59 | 69.52 | 69.81 | 70.19 | 67.01 | 66.40 |
| DeepSeek-V2-Chat-Lite-16B | 70.30 | 55.96 | 58.96 | 51.21 | 62.05 | 48.39 | 52.04 | 54.80 | 52.86 | 57.04 | 50.01 | 51.01 |
| Gemma-2-IT-9B | 83.69 | 78.72 | 78.13 | 79.38 | 79.17 | 77.86 | 76.53 | 79.82 | 79.96 | 79.80 | 79.28 | 77.24 |
| GPT-3.5-Turbo-0125 | 81.74 | 71.98 | 72.81 | 71.53 | 68.63 | 63.20 | 65.77 | 63.05 | 70.86 | 70.70 | 65.21 | 72.54 |
| GPT-4o | 83.29 | 78.31 | 74.51 | 80.29 | 79.40 | 77.64 | 78.29 | 80.03 | 78.10 | 80.23 | 79.56 | 80.00 |

(a). 2-shot F1 scores on en-x tasks.

| | de-de | es-es | vi-vi | zh-zh | hi-hi | ar-ar | el-el | ro-ro | ru-ru | th-th | tr-tr |
|----------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| LLaMA-3.1-8B | 71.67 | 74.17 | 73.19 | 64.96 | 71.91 | 68.20 | 69.35 | 74.65 | 67.17 | 70.89 | 66.89 |
| LLaMA-3.1-70B | 76.24 | 78.68 | 76.90 | 71.67 | 75.85 | 76.43 | 72.41 | 78.06 | 68.43 | 76.70 | 70.67 |
| Mistral-V0.3-7B | 71.02 | 72.91 | 69.91 | 66.65 | 56.73 | 58.25 | 62.81 | 71.73 | 63.57 | 62.08 | 58.44 |
| Qwen-2.5-7B | 58.74 | 57.36 | 72.29 | 73.38 | 63.74 | 72.74 | 67.24 | 58.82 | 64.12 | 77.89 | 60.84 |
| Qwen-2.5-72B | 81.29 | 81.15 | 82.82 | 89.08 | 79.12 | 79.53 | 77.83 | 82.18 | 77.49 | 85.00 | 77.29 |
| DeepSeek-V2-Lite-16B | 65.26 | 67.33 | 64.81 | 63.68 | 48.64 | 46.98 | 51.04 | 65.73 | 56.19 | 49.43 | 55.20 |
| Gemma-2-9B | 75.62 | 76.34 | 73.60 | 66.92 | 73.90 | 72.51 | 71.87 | 78.14 | 67.38 | 74.20 | 71.43 |
| LLaMA-3.1-Instruct-8B | 66.22 | 69.74 | 69.38 | 61.99 | 66.00 | 66.19 | 58.81 | 68.18 | 61.08 | 66.18 | 61.43 |
| LLaMA-3.1-Instruct-70B | 75.26 | 76.36 | 78.83 | 71.09 | 74.05 | 72.34 | 72.10 | 77.23 | 70.01 | 76.23 | 71.98 |
| Mistral-V0.3-Instruct-7B | 55.84 | 53.27 | 57.29 | 39.27 | 34.57 | 45.94 | 44.52 | 59.13 | 52.33 | 55.19 | 45.97 |
| Qwen-2.5-Instruct-7B | 73.75 | 75.24 | 78.26 | 70.21 | 67.08 | 70.82 | 67.65 | 75.61 | 67.99 | 73.89 | 67.23 |
| Qwen-2.5-Instruct-72B | 73.09 | 71.26 | 73.36 | 71.12 | 64.14 | 69.76 | 64.70 | 75.14 | 69.13 | 73.92 | 67.60 |
| DeepSeek-V2-Chat-Lite-16B | 56.24 | 59.33 | 56.18 | 50.06 | 41.19 | 42.46 | 44.03 | 54.63 | 56.10 | 40.71 | 48.52 |
| Gemma-2-IT-9B | 76.22 | 77.12 | 79.86 | 72.25 | 75.47 | 74.44 | 74.89 | 77.32 | 72.64 | 78.57 | 72.04 |
| GPT-3.5-Turbo-0125 | 75.68 | 77.58 | 73.09 | 70.49 | 67.64 | 70.09 | 71.03 | 77.17 | 71.55 | 67.00 | 71.12 |
| GPT-4o | 76.94 | 78.78 | 77.25 | 71.37 | 72.02 | 76.17 | 73.40 | 77.66 | 77.34 | 80.00 | 71.58 |

(b). 2-shot F1 scores on x-x tasks

Table 6: Detailed 2-shot F1 scores on en-x and x-x tasks in each language.

| | en-de | en-es | en-vi | en-zh | en-hi | en-ar | en-el | en-ro | en-ru | en-th | en-tr |
|----------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| LLaMA-3.1-8B | 0.68 | 0.00 | 1.28 | 0.00 | 0.10 | 0.10 | 0.00 | 0.00 | 0.19 | 0.10 | 1.10 |
| LLaMA-3.1-70B | 60.20 | 54.59 | 63.73 | 48.90 | 69.09 | 66.86 | 56.83 | 67.88 | 56.69 | 56.09 | 61.50 |
| Mistral-V0.3-7B | 2.11 | 1.78 | 16.39 | 23.80 | 61.82 | 54.17 | 6.96 | 2.65 | 1.02 | 46.35 | 16.63 |
| Qwen-2.5-7B | 2.42 | 1.01 | 0.43 | 0.43 | 0.00 | 0.00 | 0.00 | 1.15 | 0.00 | 0.00 | 5.17 |
| Qwen-2.5-72B | 7.39 | 6.35 | 11.22 | 0.98 | 5.74 | 2.82 | 19.78 | 11.36 | 3.59 | 27.01 | 23.64 |
| DeepSeek-V2-Lite-16B | 8.90 | 1.87 | 1.69 | 26.75 | 46.80 | 4.76 | 3.68 | 2.27 | 0.47 | 35.17 | 10.31 |
| Gemma-2-9B | 0.00 | 0.10 | 9.78 | 2.37 | 0.50 | 0.20 | 0.34 | 3.15 | 0.00 | 2.23 | 2.37 |
| LLaMA-2-Chat-7B | 3.07 | 0.45 | 0.72 | 2.25 | 0.27 | 0.00 | 0.32 | 1.55 | 0.00 | 0.12 | 1.62 |
| LLaMA-3.1-Instruct-8B | 1.15 | 1.87 | 0.37 | 0.36 | 0.59 | 0.00 | 0.00 | 1.75 | 0.36 | 0.37 | 3.01 |
| LLaMA-3.1-Instruct-70B | 0.79 | 0.76 | 0.00 | 0.00 | 0.68 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.34 |
| Mistral-V0.3-Instruct-7B | 4.30 | 1.73 | 12.79 | 0.00 | 0.35 | 0.00 | 0.44 | 3.12 | 0.00 | 0.67 | 7.03 |
| Qwen-2.5-Instruct-7B | 1.89 | 0.71 | 0.76 | 0.38 | 0.00 | 0.00 | 0.00 | 0.68 | 0.00 | 0.00 | 2.91 |
| Qwen-2.5-Instruct-72B | 6.31 | 1.13 | 5.68 | 9.73 | 2.46 | 7.51 | 2.61 | 4.41 | 1.96 | 0.35 | 8.18 |
| DeepSeek-V2-Chat-Lite-16B | 5.43 | 2.35 | 0.90 | 0.41 | 4.47 | 1.33 | 0.94 | 4.09 | 0.65 | 0.66 | 4.74 |
| Gemma-2-IT-9B | 0.96 | 0.00 | 0.31 | 0.00 | 0.00 | 0.00 | 0.00 | 0.28 | 0.00 | 0.00 | 0.28 |
| GPT-3.5-Turbo-0125 | 1.12 | 0.36 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.32 |
| GPT-4o | 0.39 | 0.33 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.38 | 0.00 | 0.00 |

Table 7: 2-shot language error rates (%) on en-x xMRC tasks.

| | en-en | en-de | en-es | en-vi | en-zh | en-hi | en-ar | en-el | en-ro | en-ru | en-th | en-tr |
|----------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| LLaMA-3.1-8B | 5.60 | 6.77 | 8.55 | 10.97 | 5.00 | 9.59 | 8.37 | 10.56 | 9.64 | 7.94 | 9.74 | 10.52 |
| LLaMA-3.1-70B | 1.20 | 1.98 | 2.41 | 2.19 | 2.62 | 4.12 | 2.33 | 2.35 | 1.71 | 1.57 | 2.70 | 4.94 |
| Mistral-V0.3-7B | 0.49 | 4.85 | 4.93 | 17.66 | 10.00 | 32.00 | 13.97 | 18.50 | 4.86 | 5.25 | 26.43 | 18.33 |
| Qwen-2.5-7B | 1.51 | 2.03 | 1.30 | 5.72 | 1.96 | 2.26 | 3.11 | 6.09 | 2.10 | 3.99 | 2.85 | 4.77 |
| Qwen-2.5-72B | 0.00 | 1.19 | 1.35 | 0.97 | 1.82 | 3.69 | 2.62 | 3.16 | 1.91 | 2.94 | 4.93 | 3.29 |
| DeepSeek-V2-Lite-16B | 1.87 | 3.26 | 2.33 | 11.97 | 2.39 | 20.61 | 10.10 | 8.04 | 4.27 | 2.76 | 15.97 | 11.23 |
| Gemma-2-9B | 1.02 | 1.34 | 1.38 | 5.39 | 2.60 | 6.51 | 4.69 | 4.98 | 3.68 | 2.64 | 5.05 | 6.98 |
| LLaMA-2-Chat-7B | 6.18 | 9.59 | 16.79 | 22.45 | 17.69 | 60.33 | 55.46 | 46.73 | 10.01 | 8.93 | 62.89 | 21.38 |
| LLaMA-3.1-Instruct-8B | 0.85 | 2.51 | 1.73 | 1.36 | 1.99 | 1.66 | 2.58 | 2.37 | 6.48 | 2.67 | 1.03 | 3.43 |
| LLaMA-3.1-Instruct-70B | 1.85 | 1.06 | 0.68 | 1.78 | 2.94 | 1.53 | 1.54 | 1.75 | 2.91 | 2.07 | 2.26 | 2.10 |
| Mistral-V0.3-Instruct-7B | 1.77 | 2.03 | 2.18 | 7.79 | 2.25 | 1.81 | 4.00 | 1.80 | 2.71 | 2.78 | 5.61 | 3.37 |
| Qwen-2.5-Instruct-7B | 2.75 | 2.47 | 3.20 | 3.42 | 4.04 | 2.78 | 5.38 | 3.25 | 2.36 | 2.56 | 2.04 | 3.80 |
| Qwen-2.5-Instruct-72B | 0.38 | 1.58 | 1.09 | 1.19 | 0.54 | 1.71 | 1.68 | 1.35 | 2.77 | 1.65 | 1.30 | 3.00 |
| DeepSeek-V2-Chat-Lite-16B | 0.58 | 4.28 | 2.93 | 9.77 | 3.27 | 6.80 | 4.98 | 6.31 | 6.69 | 8.60 | 6.15 | 5.39 |
| Gemma-2-IT-9B | 1.95 | 2.26 | 1.76 | 1.92 | 3.30 | 3.03 | 3.73 | 1.47 | 3.43 | 2.84 | 0.94 | 2.52 |
| GPT-3.5-Turbo-0125 | 0.00 | 0.65 | 2.35 | 1.61 | 3.49 | 2.15 | 2.60 | 1.44 | 2.89 | 3.70 | 5.17 | 4.73 |
| GPT-4o | 0.00 | 0.45 | 0.86 | 1.56 | 0.50 | 0.00 | 0.92 | 1.00 | 3.57 | 2.56 | 0.00 | 4.00 |

Table 8: 2-shot generation failure error rates (%) on en-x xMRC tasks.

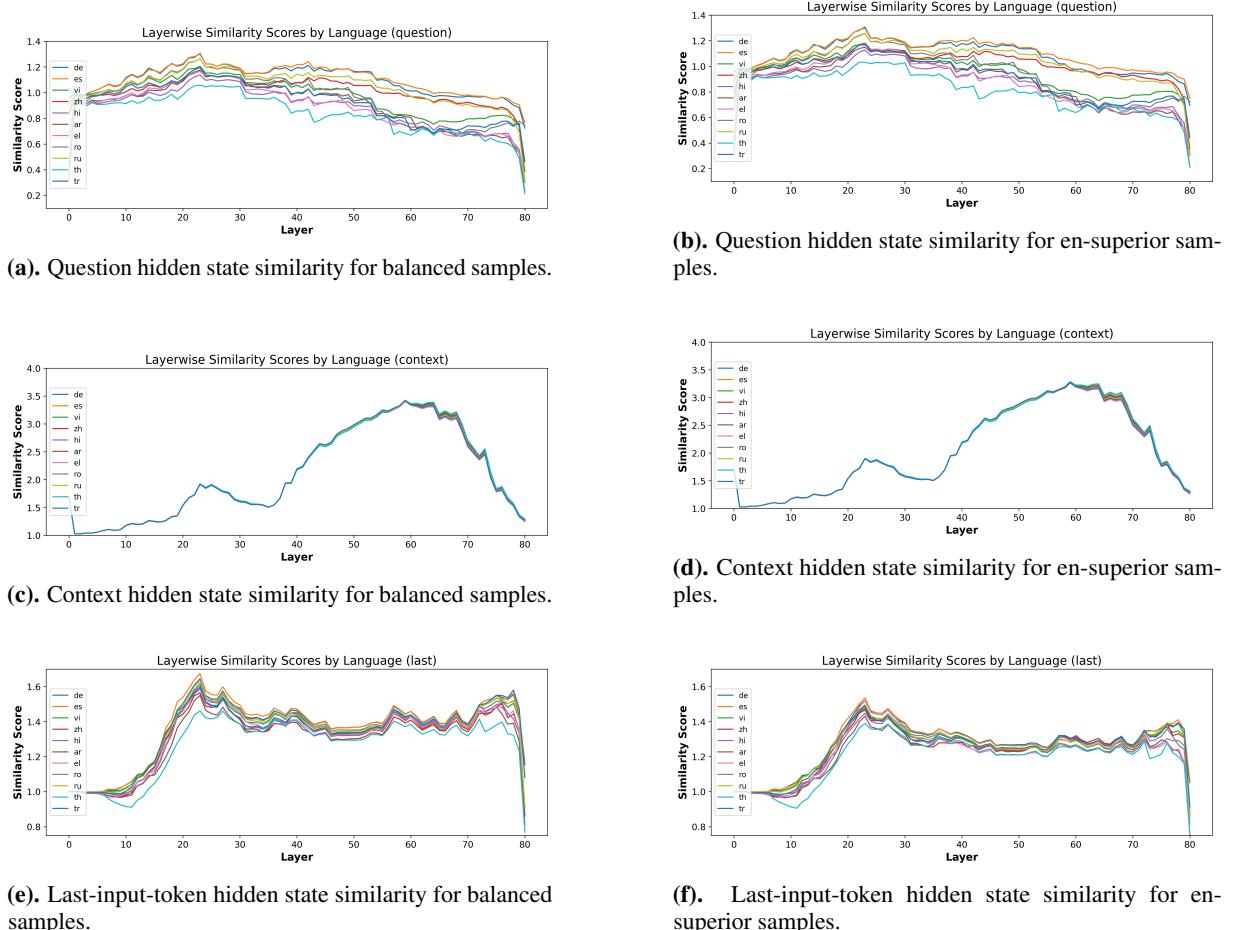


Figure 16: Hidden state similarity between English and other languages on different parts of the selected samples in each layer of the LLaMA-3.1-70B model.

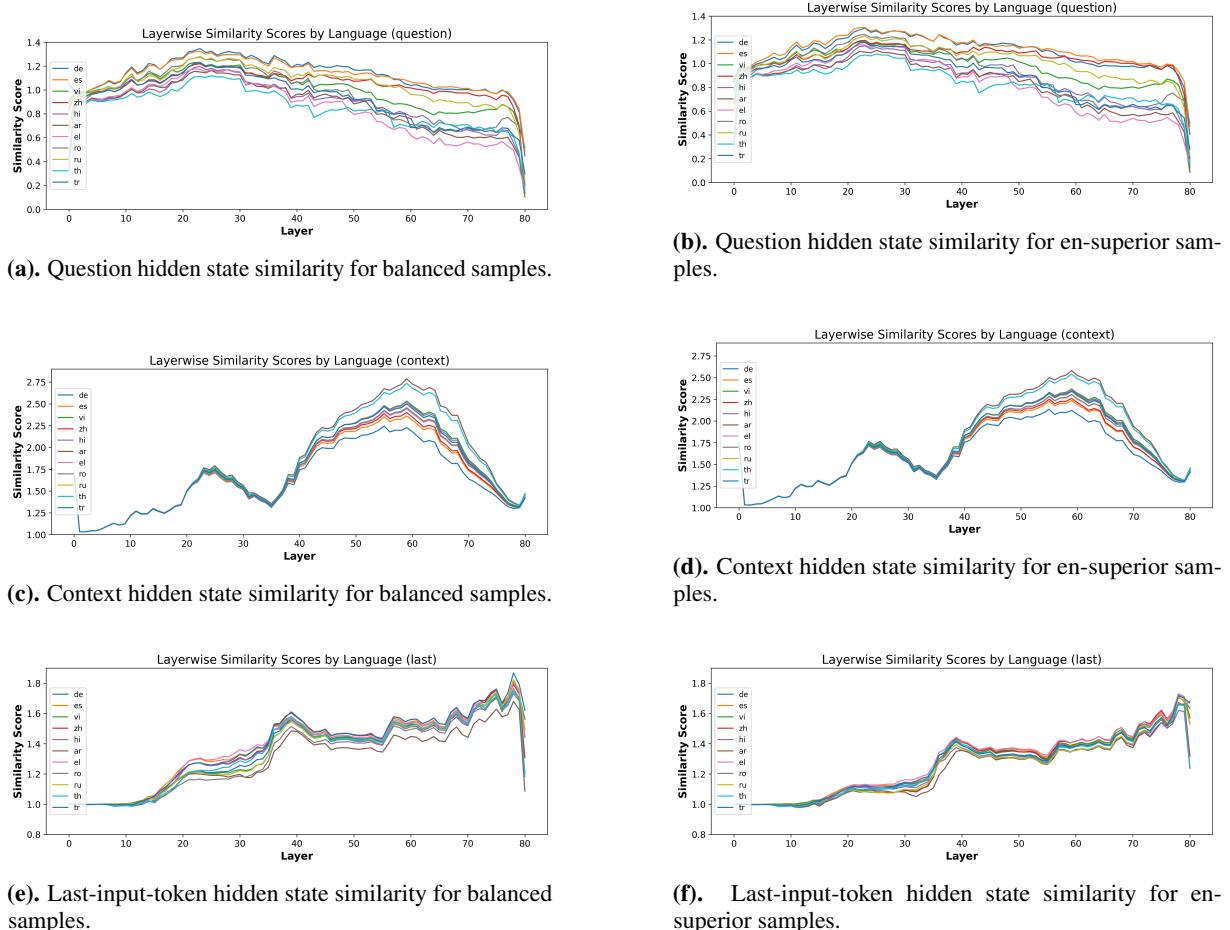
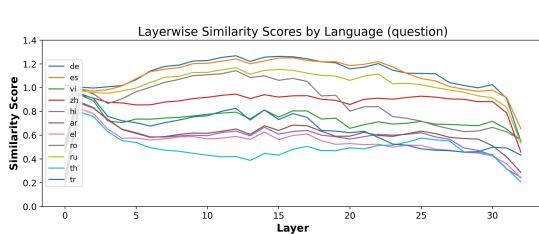
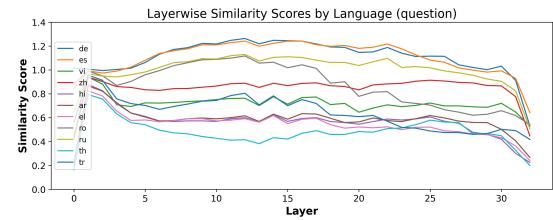


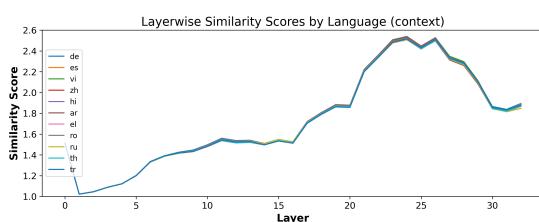
Figure 17: Hidden state similarity between English and other languages on different parts of the selected samples in each layer of the LLaMA-3.1-Instruct-70B model.



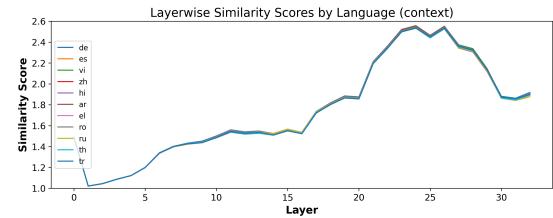
(a). Question hidden state similarity for balanced samples.



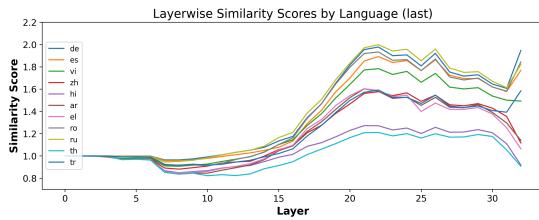
(b). Question hidden state similarity for en-superior samples.



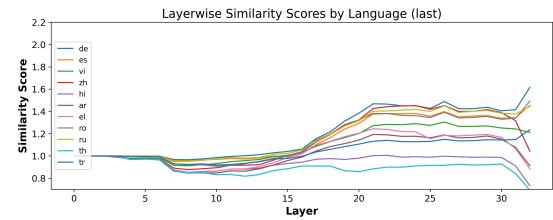
(c). Context hidden state similarity for balanced samples.



(d). Context hidden state similarity for en-superior samples.

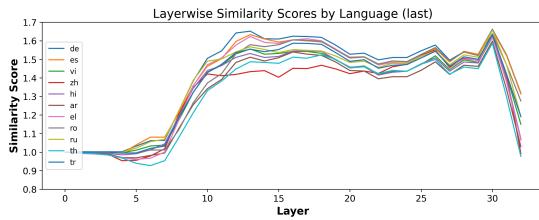


(e). Last-input-token hidden state similarity for balanced samples.

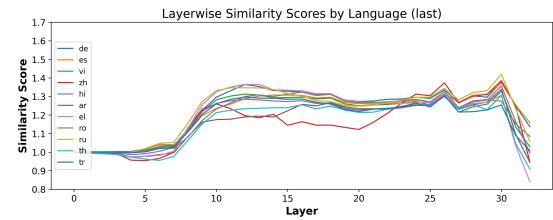


(f). Last-input-token hidden state similarity for en-superior samples.

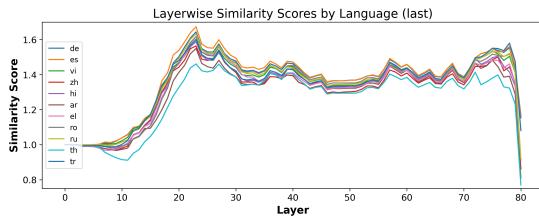
Figure 18: Hidden state similarity between English and other languages on different parts of the selected samples in each layer of the LLaMA-2-Chat-7B model.



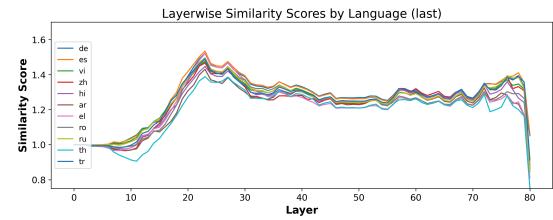
(a). "Balanced" samples for LLaMA-3.1-8B.



(b). "En-superior" samples for LLaMA-3.1-8B.



(c). "Balanced" samples for LLaMA-3.1-70B.



(d). "En-superior" samples for LLaMA-3.1-70B.

Figure 19: Last-input-token hidden state similarity between English and other languages in each layer of the base LLaMA models.

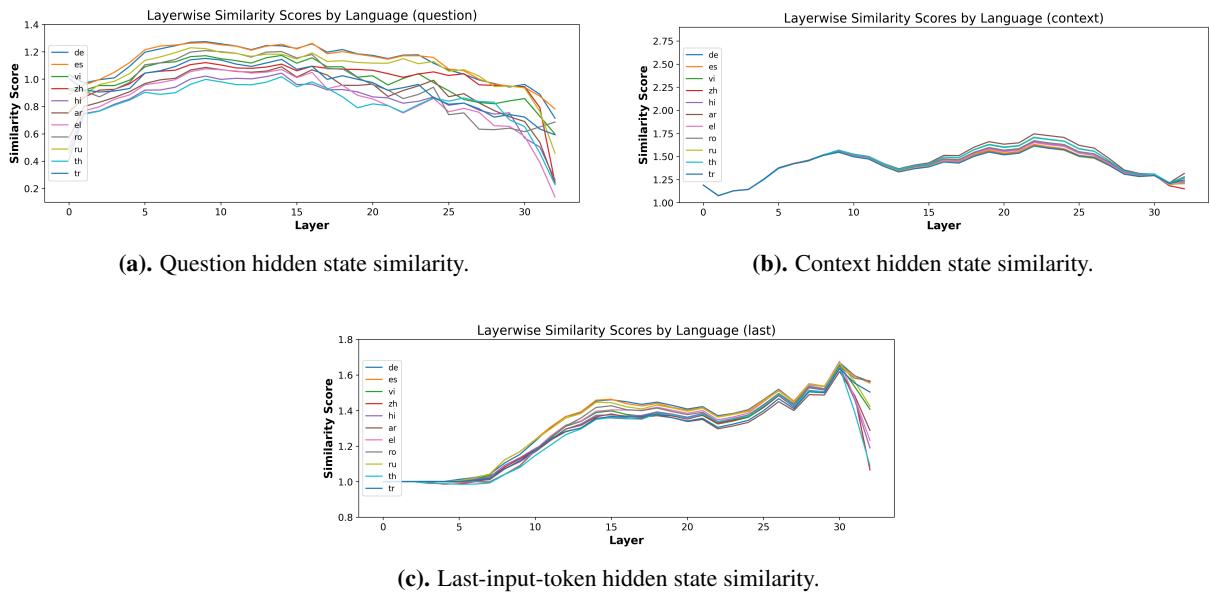


Figure 20: Hidden state similarity between English and other languages on different parts of the balanced samples in each layer for our finetuned LLaMA-3.1-8B model.