

# Lost in Multilinguality: Dissecting Cross-lingual Factual Inconsistency in Transformer Language Models

Anonymous ACL submission

## Abstract

Multilingual language models (MLMs) store factual knowledge across languages but often struggle with cross-lingual factual consistency, i.e., with providing consistent responses to semantically equivalent prompts in different languages. While previous studies point out this issue, the underlying causes remain unexplored. In this work, we use mechanistic interpretability methods to investigate cross-lingual inconsistencies in MLMs. We find that MLMs encode knowledge in an language-independent concept space through most layers, and only transition to language-specific spaces in the final layers. Failures during this language transition process often result in incorrect predictions in the target language, even when the model correctly predicts the answer in other languages. To mitigate this inconsistency issue, we propose a linear shortcut method that bypasses computations in the final layers, enhancing both prediction accuracy and cross-lingual consistency. Overall, this study deepens the understanding of MLM mechanisms and offers insights for generating consistent factual predictions.

## 1 Introduction

Multilingual language models (MLMs) have shown remarkable capabilities in storing and retrieving factual knowledge across languages (Jiang et al., 2020; Kassner et al., 2021). However, they often exhibit inconsistencies when responding to semantically equivalent prompts in different languages. For instance, an MLM might correctly predict the capital of Canada when asked in English but fail to do so when queried in another language, e.g., Chinese. This phenomenon is known as *cross-lingual factual inconsistency* (Qi et al., 2023). It raises questions about how effectively MLMs transfer knowledge across languages, and shows limitations in their robustness and fairness.

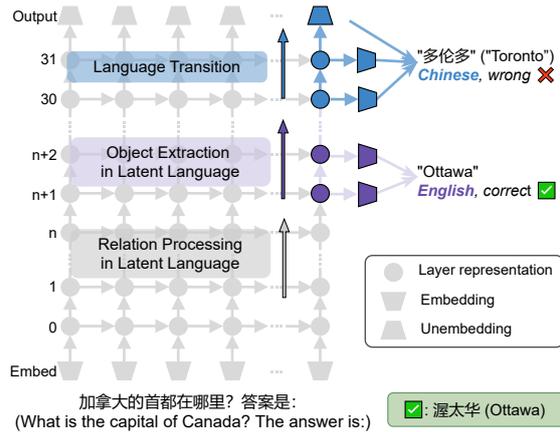


Figure 1: Illustration of language transition failure in LLaMA2 when answering the question: “加拿大的首都 在哪里? 答案是:” (“What is the capital of Canada? The answer is:”). In intermediate layers, the model processes information in its latent language, i.e., a concept space independent of the input language.<sup>1</sup> While it correctly identifies “Ottawa” in English during the concept-space object extraction, the final output “多伦多” (“Toronto”) is incorrect after transitioning to Chinese. This indicates the model’s failure to adapt knowledge from the concept space to the target language, leading to cross-lingual inconsistency.

Understanding the root causes of such inconsistencies is crucial, yet research in this area remains limited. While prior studies have explored the inner workings of MLMs (Wendler et al., 2024; Dumas et al., 2024; Fierro et al., 2024), they mainly focus on scenarios where models make correct predictions, leaving the reasons behind inconsistent predictions unexplored. Furthermore, while Qi et al. (2023) identify frequent cross-language inconsistencies in MLMs, they do not investigate the underlying causes behind them.

In this work, we address this research gap

<sup>1</sup>This concept space in LLaMA2, as seen through the Logit Lens (Nostalgebraist, 2020), exhibits a bias towards English, reflecting its English-centric nature (Wendler et al., 2024).

by analyzing cross-lingual factual inconsistency through the lens of mechanistic interpretability (Olah, 2022; Nanda et al., 2023), which aims at reverse-engineering and, thereby, understanding language models. We trace information flows within MLMs to identify where inconsistencies arise on two complementary scenarios: (1) cases where models produce correct predictions consistent with English and (2) cases where models predict correctly in English but generates incorrect answers in other languages.<sup>2</sup> This comparison aims at uncovering the causes of both success and failure in multilingual factual recall.

Our analysis reveals that MLMs process factual knowledge in a concept space largely independent of the input language through most layers, and transition to language-specific spaces in the final layers. However, even when the correct prediction is encoded in this concept space, the model can fail the language transition, leading to incorrect predictions in the target language (see Figure 1). This highlights the critical role of the language transition mechanism for cross-lingual consistency.

Overall, our contributions are as follows:

(i) **Dataset Construction (§3)**: We introduce KLAR, an enhanced KnowLedge probing dataset for Auto-Regressive models, covering 17 languages and 20 relation types. It provides a robust framework for multilingual knowledge probing, which we use to evaluate the cross-lingual consistency of two state-of-the-art MLMs (§4).

(ii) **Mechanistic Analysis (§5)**: We conduct the first interpretability-driven study of cross-lingual factual inconsistency, revealing how MLMs encode and process factual knowledge across layers.

(iii) **Failure Mode Identification (§6)**: In a detailed layer-wise analysis, we identify the language transition mechanism as main failure point that leads to cross-lingual inconsistency.

(iv) **Approach (§7)**: We propose a shortcut method that bypasses the model’s final-layer computations, enhancing both prediction accuracy and cross-lingual consistency in MLMs.<sup>3</sup>

## 2 Related Work

**Mechanistic Interpretability (MI)** aims to understand LLMs by decomposing their computations into smaller, interpretable components. It

<sup>2</sup>English serves as the pivot language due to its central role in many multilingual language models (Held et al., 2023; Zhang et al., 2023).

<sup>3</sup>We will release our dataset and code upon publication.

has gained significant attention for studying factual knowledge recall in LLMs (Meng et al., 2022; Dai et al., 2022; Geva et al., 2023; Yu et al., 2023; Lv et al., 2024; Wang et al., 2024).

Following Olah et al. (2020) and Rai et al. (2024), MI research is categorized into the study of *features*, which capture human-interpretable properties in model representations or components like neurons and attention heads (Elhage et al., 2022; Gurnee et al., 2023), and the study of *circuits*, which refer to subgraphs of the model’s computation graph responsible for implementing specific behaviors (Wang et al., 2023; Elhage et al., 2021).

In this work, we focus on representation-level feature-based interpretability analysis to interpret the behavior of multilingual language models in the knowledge probing task. Specifically, we use Logit Lens (Nostalgebraist, 2020) to project latent state representations of LMs into the vocabulary space, enabling the analysis of intermediate representations and tracking how information evolves across layers.

### Interpreting Multilingual Language Models.

Recent studies have explored the internal workings of MLMs. Wendler et al. (2024) examine the latent language of LLaMA2 models using controlled translation, completion, and cloze tasks, finding that LLaMA2 internally relies on English as a pivot language. Building on this setup, Dumas et al. (2024) investigate the disentanglement of language and concept representations, demonstrating that LLaMA2 processes language and concept information independently. Fierro et al. (2024) analyze knowledge probing tasks to study how mechanisms identified in monolingual contexts generalize to multilingual settings, but their focus remains limited to correct prediction cases.

In contrast, our work centers on understanding the internal mechanisms responsible for cross-lingual inconsistencies. By examining both consistent and inconsistent predictions, we uncover how MLMs transition from language-independent to language-specific processing. This approach offers new insights into how MLMs encode and transfer factual knowledge across languages, addressing a key gap in prior research.

## 3 KLAR Dataset

We focus on the factual knowledge probing task, where a fact is represented as a subject-relation-object triple  $\langle s_i, r_i, o_i \rangle$  and expressed in natural

language prompts. Given a prompt constructed from the subject  $s_i$  and relation  $r_i$ , LMs are expected to predict the object  $o_i$ . For example, the fact  $\langle \text{Canada}, \text{capital}, \text{Ottawa} \rangle$  can be queried as, “What is the *capital* of *Canada*?”, and the model should predict the object *Ottawa* as the answer.

Qi et al. (2023) introduce the BMLAMA17 dataset for evaluating multilingual factual knowledge in MLMs. However, in many factual questions in BMLAMA17, the object appears in the middle of the sentence rather than at the end, which is incompatible with knowledge probing for auto-regressive models. Furthermore, BMLAMA17 includes many relations with multiple correct answers,<sup>4</sup> making it difficult to reliably evaluate the correctness of a model’s response for a given  $\langle s_i, r_i, o_i \rangle$  triple where  $o_i$  is only one of the possible answers.

To address these limitations, we construct KLAR, a KnowLedge probing dataset that ensures compatibility with Auto-Regressive models and provides clarity in factual evaluation. We extract parallel factual knowledge triples in 17 languages from BMLAMA17 and design prompts where the object consistently appears at the end. Relation-specific templates are structured as “<Question> The answer is:”, e.g.,  $\langle \text{Canada}, \text{capital}, \text{Ottawa} \rangle$  becomes: “What is the capital of Canada? The answer is:”. These templates are initially created in English and translated into 16 other languages using gpt-35-turbo. To ensure clarity, we exclude relations with multiple correct answers and inspect the semantic clarity in prompt templates manually and/or through back-translation.

The resulting KLAR dataset includes 2,621 parallel factual knowledge triples in 17 languages, covering 20 relation types. Table 1 provides an overview of the languages and sample relations. Detailed statistics are provided in Appendix A.1.

## 4 Cross-lingual Consistency Evaluation

**Models and Languages** We analyze two widely used open-source multilingual auto-regressive language models: LLaMA2-7B (Touvron et al., 2023) and BLOOM-560M (Le Scao et al., 2023). LLaMA2 is trained on a multilingual corpus dominated by English, which accounts for 89.7% of the

<sup>4</sup>For example, the relation “shares\_border\_with” (prompt: “Which country does <subject> share a border with?”) often involves multiple correct answers, as a country typically shares borders with several others.

Languages (17)	
Arabic ( <i>ar</i> ), Catalan ( <i>ca</i> ), Greek ( <i>el</i> ), English ( <i>en</i> ), Spanish ( <i>es</i> ), Persian ( <i>fa</i> ), French ( <i>fr</i> ), Hebrew ( <i>he</i> ), Hungarian ( <i>hu</i> ), Japanese ( <i>ja</i> ), Korean ( <i>ko</i> ), Dutch ( <i>nl</i> ), Russian ( <i>ru</i> ), Turkish ( <i>tr</i> ), Ukrainian ( <i>uk</i> ), Vietnamese ( <i>vi</i> ), Chinese ( <i>zh</i> )	
Relations (4/20)	Prompt example
capital	What is the capital of <subject>? The answer is:
continent	Which continent is <subject> located in? The answer is:
field_of_work	What field does <subject> work in? The answer is:
religion	What is the religious belief of <subject>? The answer is:

Table 1: Overview of the languages and 4 sample relations (out of 20 relations in total) in KLAR.

data, whereas BLOOM’s training data is more balanced, with English comprising 31.3% of the corpus. Our analysis considers the languages shared between each model and our dataset, covering 12 languages for LLaMA2 and 7 for BLOOM. Details on the selected languages are provided in Table 4 in Appendix A.1.

**Evaluation** Many prior studies (Geva et al., 2023; Qi et al., 2023; Fierro et al., 2024) assess correctness based on the model’s first predicted token. However, this approach is problematic, especially in multilingual settings with complex tokenization. In many cases, even if the model predicts the correct first token, its complete output can still be incorrect.<sup>5</sup> To address this issue, we evaluate correctness based on the model’s full answer to each factual question rather than relying solely on the first token. Following Jiang et al. (2020), we evaluate cross-lingual consistency using the overlap ratio of correct predictions for parallel facts between language pairs.<sup>6</sup>

**Results** Figure 2 shows the cross-lingual consistency results for LLaMA2 and BLOOM. While LLaMA2 generally performs better than BLOOM, both models face challenges in achieving high consistency across languages, particularly between

<sup>5</sup>For example, given the Chinese prompt “文森山位于哪个大陆? 答案是: ” (“Which continent is Vinson Massif located in? The answer is:”), the BLOOM model outputs “南美洲” (“South America”) instead of the correct answer “南极洲” (“Antarctica”). Although both responses share the same first token, the final prediction is incorrect.

<sup>6</sup>We do not adopt the candidate-based consistency metric proposed by Qi et al. (2023), as it relies on the next-token prediction, which, as discussed in Section 4, is unreliable in a multilingual setup.

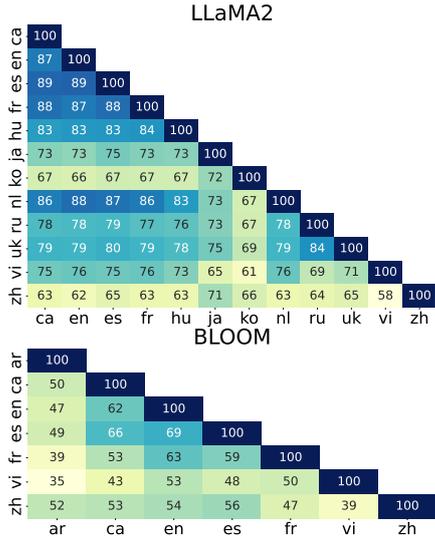


Figure 2: Cross-lingual consistency results across language pairs. The heatmaps show the overlap ratio of correct predictions between language pairs.

linguistically diverse pairs. The impact of language scripts is especially evident: Non-Latin scripts, such as Arabic (*ar*), Chinese (*zh*), and Korean (*ko*), consistently show lower consistency scores. This underscores that cross-lingual consistency remains a key limitation for both models, emphasizing the need for more robust approaches to effectively analyze and address this issue.

## 5 Analyzing Multilingual Factual Recall

To understand how multilingual language models recall factual knowledge across languages, we analyze their internal mechanisms from multiple perspectives: the layer-wise evolution of prediction ranks (§5.1), latent state similarities across language pairs (§5.2), information flow within the model (§5.3), and the composition of the latent concept space (§5.4).

### 5.1 From the Perspective of Rankings

First, we use Logit Lens (Nostalgebraist, 2020) to project latent states at each layer to the vocabulary (unembedding) and measure the rank (the lower, the better) of the target object at each layer. Specifically, we compare the rank of the correct object in its target language (`rank_target_correct`) and its English equivalent (`rank_en_correct`). This approach allows us to trace how the model processes factual knowledge across layers and transitions between different representation modes.

Figure 3a shows distinct phases of knowledge processing in both models. In the early lay-

ers, both ranks remain high, indicating that the models have not begun extracting the target object. Around layer 15 in BLOOM and layer 12 in LLaMA2, both (`rank_target_correct`) and (`rank_en_correct`) drop significantly, marking the beginning of the object extraction phase.

This phase continues until layer 28 in LLaMA2 and layer 19 in BLOOM, where a notable divergence occurs. The English rank (`rank_en_correct`) begins to increase, while the target-language rank (`rank_target_correct`) continues to decrease. This divergence reflects a transition from language-independent object extraction to target language-specific object extraction, where the models adapt the representations to align with the target language.

These findings show that MLMs recall knowledge through an initial concept-space object extraction phase (marked by significant rank drops for both English and target language answers) before transitioning to language-specific object extraction and producing the final output.

### 5.2 From the Perspective of Latent States

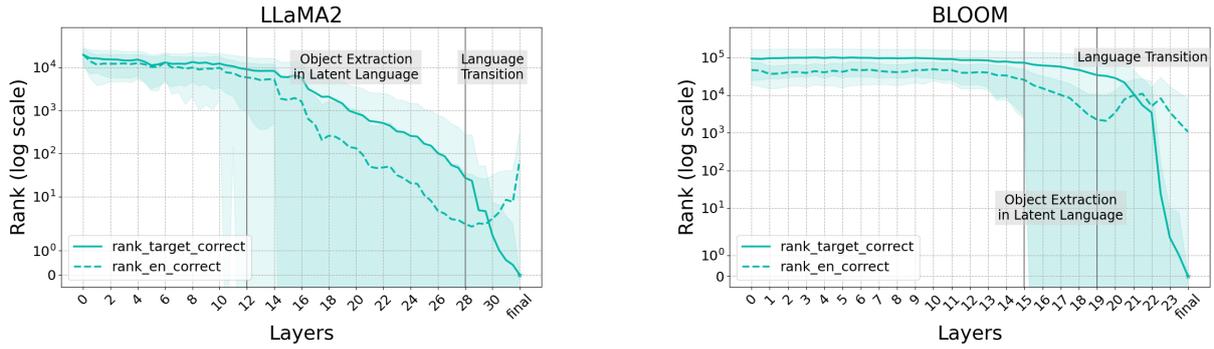
Moreover, we measure the cosine similarity of latent states between language pairs across layers.

Figure 3b shows the average cosine similarity of latent states between English and individual target languages for LLaMA2 and BLOOM.<sup>7</sup> As information propagates through the layers, similarity increases, peaking around 0.8 in the middle layers for both models. This trend holds even for linguistically diverse pairs, such as English and Arabic, suggesting the formation of a shared concept space where factual knowledge is encoded in the model’s latent language which is generic and independent of the input language. In the final layers, similarity decreases, reflecting a transition to language-specific processing. This aligns with the divergence observed in Section 5.1, where the rank changes of the target language object and its English equivalent begin to differ. These observations confirm the model’s transition from concept-space object extraction to language-specific adaptations in the final layers.

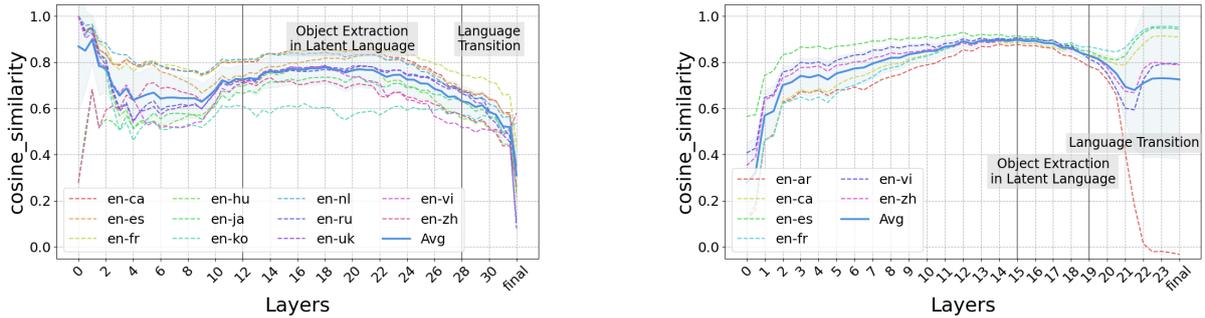
### 5.3 Information Flow Dissection

While Sections 5.1 and 5.2 demonstrate the presence of a concept space in the middle layers, they

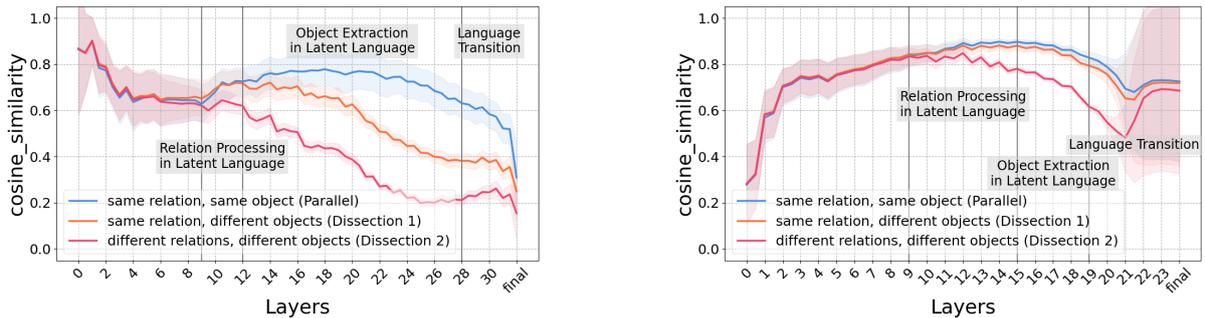
<sup>7</sup>For clarity, only language pairs involving English are shown here. Complete results for all language pairs are provided in Appendix A.2.1.



(a) Layer-wise rank of correct predictions averaged across all languages and relations (§5.1). “rank\_target\_correct” denotes the rank of correct predictions in the target language, while “rank\_en\_correct” represents the rank of their English equivalents.



(b) Cosine similarity of latent state similarity between each language pair averaged across all relations (§5.2).



(c) Comparative study of latent state similarity across language pairs (§5.3). We compare the latent state similarity for parallel facts, non-parallel facts sharing the same relation, and non-parallel facts belonging to different relations, respectively.

Figure 3: Analysis of multilingual knowledge probing of LLaMA2 and BLOOM, including (3a) layer-wise evolution of correct prediction ranks, (3b) latent state similarities across languages, and (3c) the development of latent state similarities in different settings.

do not clarify the type of information contributing to the observed high similarity between language pairs. To disentangle whether this similarity arises from relational information, object information, or both, we perform comparative experiments under three conditions: (1) **Same relation, same object (Parallel)**, as in Section 5.2): Latent state similarity is calculated using parallel facts between each language pair (e.g., "the capital of Canada" in both English and another language); (2) **Same relation, different objects (Dissection 1)**: Similarity is calculated using non-parallel facts sharing the same relation (e.g., "the capital of Canada" in one language versus "the capital of Spain" in another);

(3) **Different relation, different objects (Dissection 2)**: Similarity is calculated using non-parallel facts from different relations (e.g., "the capital of Canada" versus "the official language of Spain").

Figure 3c shows distinct processing phases. Around layer 9, the *Dissection 2* curve drops significantly in both models, while *Parallel* and *Dissection 1* curves remain close, indicating that models process relational information specific to the current fact’s relation. The high similarity during this stage suggests that such relation processing happens in a language-independent concept space.

From layer 12 in LLaMA2 and layer 15 in BLOOM, the *Dissection 1* curve begins to drop,

marking a transition to object-specific processing. During layers 12–28 in LLaMA2 and layers 15–19 in BLOOM, the *Parallel* curve remains high, indicating that object information is processed in the model’s latent language.

At layer 28 in LLaMA2 and layer 19 in BLOOM, the *Parallel* curve drops significantly, signaling the language transition phase, where the concept-space object representations are adapted to the target language.

Together, the progression shows the models’ transitions from relation processing to object extraction and to language-specific adaptation.

#### 5.4 Concept Space Language Composition

To further explore how the concept space encodes information in MLMs, we analyze the language composition of their latent states. Using Logit Lens, we project intermediate layer representations onto the vocabulary space and identify the language of the top-10 predicted tokens at each layer using fasttext (Joulin et al., 2017).<sup>8</sup>

Figure 4 shows the language composition for LLaMA2 and BLOOM with Chinese (zh) as the input language, averaged across factual queries spanning all relations. Results for other input languages are provided in Appendix A.2.3.

In LLaMA2, English dominates the middle-to-upper layers, suggesting that factual knowledge is processed in an English-centric concept space. This is consistent with prior findings that “LLaMA2 models think in English” (Wendler et al., 2024). In contrast, BLOOM exhibits a more diverse composition in the middle-to-upper layers, comprising primarily Latin-based languages like English, French, Spanish, German, etc. Across different input languages (see Appendix Figures 10 and 11), both models show similar middle-to-upper layer compositions, regardless of the input language. This demonstrates that MLMs encode knowledge in a concept space independent of the input language.

#### 5.5 Summary

Our analysis reveals a three-stage knowledge recall process in MLMs (as illustrated in Figure 1): first relation processing, then object extraction in the model’s latent language, and finally the transition to language-specific processing to adapt the object to the target language. These findings pro-

<sup>8</sup>We filter out tokens with confidence scores below 0.5.

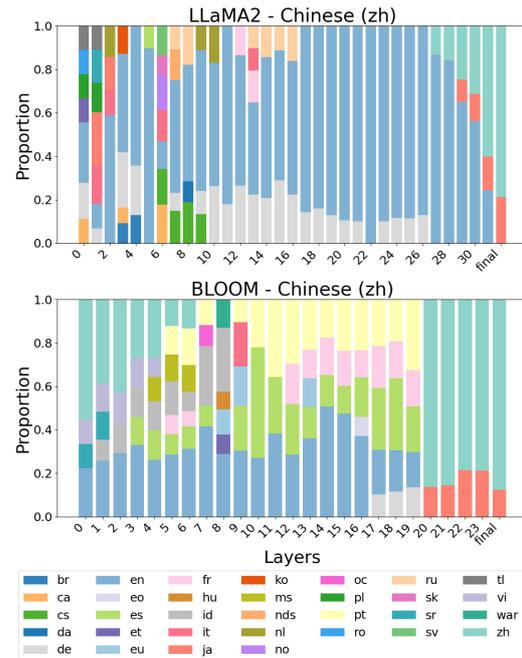


Figure 4: Language composition of latent representations with Chinese as the input language. In LLaMA2, English dominates the middle-to-upper layers, whereas BLOOM has a more diverse language composition.

vide a comprehensive view on the mechanisms of multilingual factual recall.

### 6 Examining the Cause of Cross-Lingual Inconsistency

Next, we analyze incorrect predictions across languages to investigate the causes of cross-lingual inconsistencies in MLMs.

Figure 5 shows the rank evolution for incorrect predictions in LLaMA2 and BLOOM. While the rank of the correct answer decreases significantly in the middle layers (both in the target language and in English) — consistent with the behavior observed in correct predictions (Figure 3a) — the rank of the incorrect answer surpasses that of the correct answer during language transition in the final layers. This suggests that factual knowledge is processed in the concept space in the middle layers as in correct predictions, but errors arise during the transition to language-specific processing.

To further investigate this phenomenon, we examine individual examples of LLaMA2.<sup>9</sup> Figure 6 presents cases in Spanish and Chinese, with additional examples provided in Appendix A.2.3. A

<sup>9</sup>LLaMA2’s English-biased latent space provides clearer insights into the switch from English to the target language, while BLOOM’s latent space is less interpretable, as shown in Figure 4.

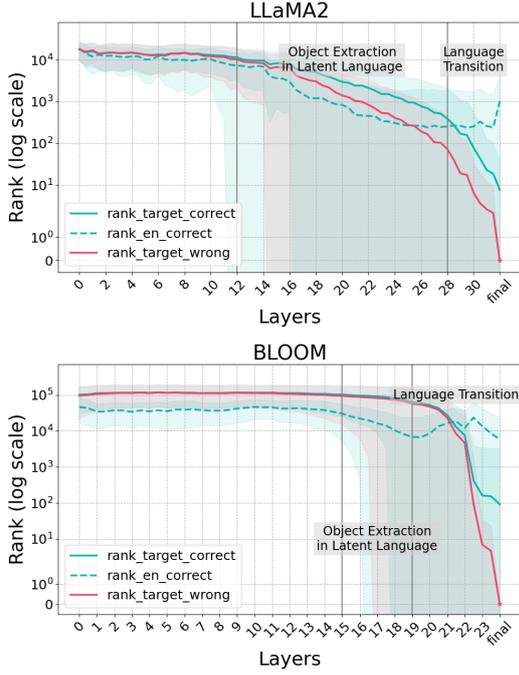


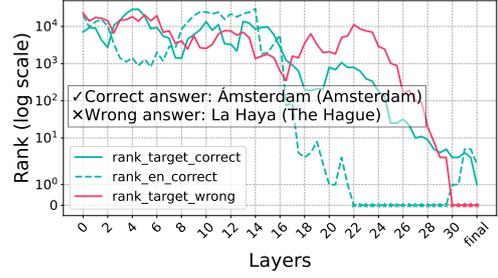
Figure 5: Layer-wise rank of incorrect predictions averaged across all languages and relations. The rank\_target\_wrong curve represents the rank of the model’s final incorrect prediction across layers, while rank\_target\_correct and rank\_en\_correct denote the ranks of the correct answer in the target language and the English equivalent, respectively.

consistent pattern emerges: in the middle-to-upper layers, the correct answer in English often ranks lowest (rank\_en\_correct=0), indicating accurate recall during the concept space processing stage. However, in the final layers, the rank of the incorrect target-language answer decreases, surpassing the correct answer during language transition.

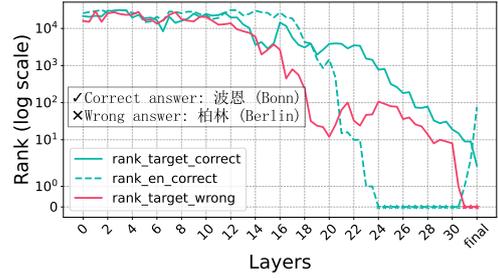
This observation underscores the critical role of language transition in cross-lingual inconsistencies. Although MLMs encode correct factual knowledge in the middle-layer concept space, the transition to language-specific processing introduces errors, causing incorrect predictions. Addressing this issue is crucial for improving cross-lingual consistency and robustness of MLMs.

## 7 Linear Shortcut for Improving Cross-Lingual Consistency

In this section, we propose a linear shortcut method to address language transition errors. Our approach bypasses final-layer computations, directly adapting concept-space representations to the target language, enhancing both prediction accuracy and cross-lingual consistency of MLMs.



(a) Prompt in Spanish: “¿Dónde se encuentra la capital de Reino de los Países Bajos? La respuesta es:” (“What is the capital of the Kingdom of Netherlands? The answer is:”).



(b) Prompt in Chinese: “西德的首都在哪里? 答案是:” (“What was the capital of West Germany? The answer is:”).

Figure 6: Rank evolution for prompts in Spanish (6a) and Chinese (6b). rank\_target\_wrong represents the rank of the model’s final incorrect prediction across layers, while rank\_target\_correct and rank\_en\_correct denote the ranks of the correct answer in the target language and the English equivalent, respectively. The plots show the impact of errors during language transition, where the rank of the incorrect answer surpasses the correct answer in the final layers.

### 7.1 Shortcut with Linear Approximation

The proposed method involves a two-step process (as illustrated in Figure 7): (a) **Deriving the linear shortcut:** Inspired by Hernandez et al. (2023), we hypothesize that the mapping from the model’s latent state at layer  $n$  to the final layer  $N$ , i.e.,  $h_n \rightarrow h_N$  can be well-approximated by a linear function  $f(h_n) = Wh_n + b \approx h_N$ . Using  $m$  correctly predicted samples, we use first-order approximation to estimate  $W$  and  $b$ , approximating the adaptation of concept-space representations to the target language.<sup>10</sup> For further details on the derivation and hyperparameters, please refer to Appendix A.3. (b) **Applying the linear shortcut:** During inference, the shortcut  $f(\cdot)$  is applied to by-

<sup>10</sup>The selection of layer  $n$  and training size  $m$  for approximating the linear transformation are treated as hyperparameters, set to  $n = 30$  for LLaMA2,  $n = 20$  for BLOOM, and  $m = 25$  for both models. Details are provided in Appendix A.3.2.

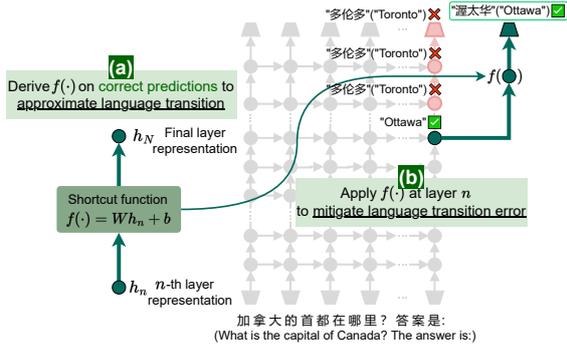


Figure 7: Illustration of the proposed shortcut method for mitigating cross-lingual inconsistency. (a) The shortcut function is learned on correct predictions to approximate language transition; (b) The learned function is then applied to bypass the error-prone final layers. In the example, the shortcut successfully recovers the correct answer, “渥太华” (“Ottawa”), in Chinese.

pass the original final-layer computations, mitigating errors introduced during language transition.

## 7.2 Results and Discussion

We evaluate the prediction accuracy and cross-lingual consistency of LLaMA2 and BLOOM and the shortcut on all KLAR samples.

**Baselines.** We compare our shortcut method to two translation-based baselines: (1) *translation-en*: We translate all input queries from each language to English using Google Translate, obtain model predictions in English, and then translate them back to the target language. (2) *translation-early-exit*: We use Logit Lens to extract top-predicted tokens from the same layers as the shortcut method, translate them into the target language and evaluate their accuracy.

**Results.** Figure 8 shows the effectiveness of the shortcut mapping: It improves prediction accuracy and cross-lingual consistency across models and languages. This demonstrates its ability to adapt concept-space knowledge to target languages for more reliable predictions.

	original	shortcut	trans-en	trans-exit
LLaMA2	71.47	<b>76.08</b>	53.88	13.93
BLOOM	43.24	<b>51.67</b>	28.03	15.68

Table 2: Average accuracy across languages..

As shown in Table 2, both baseline methods perform poorly (see Table 7 and 8 in Appendix for more details), indicating that existing transla-

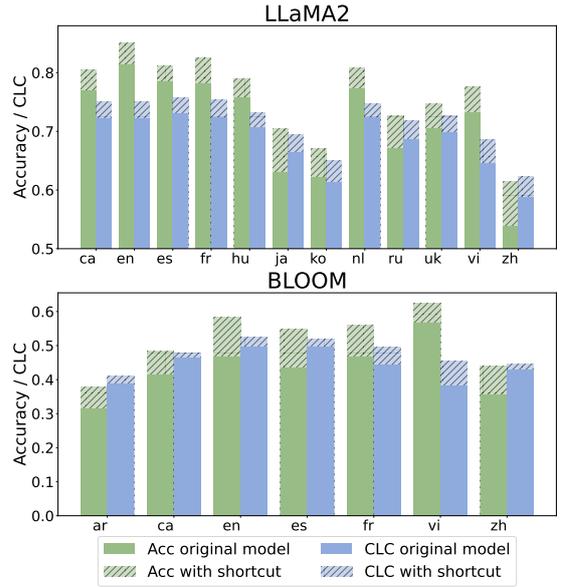


Figure 8: Accuracy (ACC) and cross-lingual consistency (CLC) per language for LLaMA2 and BLOOM, with and without the shortcut method.

tors are insufficient for cross-lingual factual prediction. In contrast, our shortcut method directly adapts latent representations from earlier layers, preserving richer contextual information and thus achieving higher prediction accuracy. Moreover, it is lightweight and efficient, relying only on linear operations, making it easily adaptable to existing MLMs.

## 8 Conclusion

This study investigates cross-lingual factual inconsistency in multilingual language models, revealing a three-stage knowledge recall process: language-independent relation processing, object extraction, and a final transition to language-specific adaptation. Errors in this transition often lead to incorrect predictions despite accurate object extraction. To address this, we propose a shortcut method that bypasses final-layer computations, improving prediction accuracy and cross-lingual consistency. Our findings enhance understanding of multilingual knowledge processing and introduce an efficient, interpretable solution for mitigating language transition errors.

Future work could expand the investigation to more languages and additional language models to assess broader applicability. Additionally, developing non-linear shortcut methods could better mitigate language transition errors, offering more robust solutions for cross-lingual consistency.

## 491 Limitations

492 First, our cross-lingual consistency analysis as-  
493 sumes English as the pivot language, reflecting the  
494 English-centric nature of most multilingual mod-  
495 els. While this aligns with prior studies (Wendler  
496 et al., 2024; Dumas et al., 2024; Fierro et al.,  
497 2024), it may limit applicability to language pairs  
498 that do not involve English.

499 Second, although the KLAR dataset covers 17  
500 languages, it does not fully capture the diversity of  
501 world languages. Expanding the analysis to more  
502 languages and exploring models with different ar-  
503 chitectures and sizes could provide deeper insights  
504 into cross-lingual inconsistencies.

505 Additionally, our shortcut method relies on lin-  
506 ear approximation for simplicity. Investigating  
507 non-linear approaches could better capture com-  
508 plex transformations during language switching  
509 and further enhance performance.

510 Finally, our analysis provides insights relevant  
511 to downstream tasks, such as multilingual knowl-  
512 edge localization (Chen et al., 2024; Kojima et al.,  
513 2024; Tang et al., 2024) and cross-lingual knowl-  
514 edge editing (Xu et al., 2023; Nie et al., 2024).  
515 However, these applications fall beyond the scope  
516 of this study and are left for future work.

## 517 References

518 Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and  
519 Jun Zhao. 2024. Journey to the center of the knowl-  
520 edge neurons: Discoveries of language-independent  
521 knowledge neurons and degenerate knowledge neu-  
522 rons. In *Proceedings of the AAAI Conference on Ar-  
523 tificial Intelligence*, volume 38, pages 17817–17825.

524 Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao  
525 Chang, and Furu Wei. 2022. Knowledge neurons  
526 in pretrained transformers. In *Proceedings of the  
527 60th Annual Meeting of the Association for Compu-  
528 tational Linguistics (Volume 1: Long Papers)*, pages  
529 8493–8502, Dublin, Ireland. Association for Com-  
530 putational Linguistics.

531 Clément Dumas, Veniamin Veselovsky, Giovanni  
532 Monea, Robert West, and Chris Wendler. 2024. How  
533 do llamas process multilingual text? a latent explo-  
534 ration through activation patching. In *ICML 2024  
535 Workshop on Mechanistic Interpretability*.

536 Nelson Elhage, Tristan Hume, Catherine Olsson,  
537 Neel Nanda, Tom Henighan, Scott Johnston, Sheer  
538 ElShowk, Nicholas Joseph, Nova DasSarma, Ben  
539 Mann, Danny Hernandez, Amanda Askell, Ka-  
540 mal Ndousse, Andy Jones, Dawn Drain, Anna  
541 Chen, Yuntao Bai, Deep Ganguli, Liane Lovitt,

Zac Hatfield-Dodds, Jackson Kernion, Tom Con-  
erly, Shauna Kravec, Stanislav Fort, Saurav Kada-  
vath, Josh Jacobson, Eli Tran-Johnson, Jared Kaplan,  
Jack Clark, Tom Brown, Sam McCandlish, Dario  
Amodei, and Christopher Olah. 2022. Softmax lin-  
ear units. *Transformer Circuits Thread*. 542  
543  
544  
545  
546  
547

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom  
Henighan, Nicholas Joseph, Ben Mann, Amanda  
Askell, Yuntao Bai, Anna Chen, Tom Conerly,  
Nova DasSarma, Dawn Drain, Deep Ganguli, Zac  
Hatfield-Dodds, Danny Hernandez, Andy Jones,  
Jackson Kernion, Liane Lovitt, Kamal Ndousse,  
Dario Amodei, Tom Brown, Jack Clark, Jared Ka-  
plan, Sam McCandlish, and Chris Olah. 2021. A  
mathematical framework for transformer circuits.  
*Transformer Circuits Thread*. 548  
549  
550  
551  
552  
553  
554  
555  
556  
557

Constanza Fierro, Negar Foroutan, Desmond Elliott,  
and Anders Søgaard. 2024. How do multilin-  
gual models remember? investigating multilin-  
gual factual recall mechanisms. *arXiv preprint  
arXiv:2410.14387*. 558  
559  
560  
561  
562

Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir  
Globerson. 2023. Dissecting recall of factual asso-  
ciations in auto-regressive language models. In *Pro-  
ceedings of the 2023 Conference on Empirical Meth-  
ods in Natural Language Processing*, pages 12216–  
12235, Singapore. Association for Computational  
Linguistics. 563  
564  
565  
566  
567  
568  
569

Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine  
Harvey, Dmitrii Troitskii, and Dimitris Bertsimas.  
2023. Finding neurons in a haystack: Case stud-  
ies with sparse probing. *Trans. Mach. Learn. Res.*,  
2023. 570  
571  
572  
573  
574

William Held, Camille Harris, Michael Best, and Diyi  
Yang. 2023. A material lens on coloniality in nlp.  
*arXiv preprint arXiv:2311.08391*. 575  
576  
577

Evan Hernandez, Arnab Sen Sharma, Tal Haklay,  
Kevin Meng, Martin Wattenberg, Jacob Andreas,  
Yonatan Belinkov, and David Bau. 2023. Linearity  
of relation decoding in transformer language models.  
In *The Twelfth International Conference on Learn-  
ing Representations*. 578  
579  
580  
581  
582  
583

Zhengbao Jiang, Antonios Anastasopoulos, Jun Araki,  
Haibo Ding, and Graham Neubig. 2020. X-FACTR:  
Multilingual factual knowledge retrieval from pre-  
trained language models. In *Proceedings of the  
2020 Conference on Empirical Methods in Natural  
Language Processing (EMNLP)*, pages 5943–5959,  
Online. Association for Computational Linguistics. 584  
585  
586  
587  
588  
589  
590

Armand Joulin, Edouard Grave, Piotr Bojanowski, and  
Tomas Mikolov. 2017. Bag of tricks for efficient  
text classification. In *Proceedings of the 15th Con-  
ference of the European Chapter of the Association  
for Computational Linguistics: Volume 2, Short Pa-  
pers*, pages 427–431, Valencia, Spain. Association  
for Computational Linguistics. 591  
592  
593  
594  
595  
596  
597



709	Qinan Yu, Jack Merullo, and Ellie Pavlick. 2023. <a href="#">Characterizing mechanisms for factual recall in language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9924–9959, Singapore. Association for Computational Linguistics.	
710		
711		
712		
713		
714		
715	Xiang Zhang, Senyu Li, Bradley Hauer, Ning Shi, and Grzegorz Kondrak. 2023. <a href="#">Don’t trust ChatGPT when your question is not in English: A study of multilingual abilities and types of LLMs</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7915–7927, Singapore. Association for Computational Linguistics.	
716		
717		
718		
719		
720		
721		
722		

<b>A Appendix</b>		723
<b>A.1 KLAR Dataset Details</b>		724
As discussed in Section 3, BMLAMA17 (Qi et al., 2023) is incompatible with multilingual knowledge probing in auto-regressive models with many objects placed in the middle of sentences, and many relations types with multiple correct answers. To address these limitations, we construct KLAR for reliable multilingual knowledge probing evaluation.		725
		726
		727
		728
		729
		730
		731
		732
		733
BMLAMA17 does not explicitly specify relation types; however, many factual questions share the same templates. We first group sentences with identical templates and use gpt-35-turbo to identify the relation for each template and map them to Wikidata property IDs (Wikidata, 2025). We discard the samples which cannot be mapped to any Wikidata property. This process yields a total of 42 relation types.		734
		735
		736
		737
		738
		739
		740
		741
For each relation, we generate English prompt templates in the format of “<Question> The answer is:” as introduced in Section 3, using gpt-35-turbo. We created five templates per relation and manually verify their clarity. The templates are then translated into 16 additional languages using gpt-35-turbo. Their quality is manually reviewed for Chinese, Spanish, and Japanese. Back-translation is used to verify clarity and consistency in the remaining languages.		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
Finally, we remove relation types with multiple correct answers and those with fewer than 30 samples. The resulting KLAR dataset comprises parallel factual knowledge spanning 17 languages and 20 relation types. For the analysis on LLaMA2 and BLOOM models, we use the intersection of languages supported by these models and included in KLAR, covering 12 languages for LLaMA2 and 7 for BLOOM, see Table 4 for the respective language list. Listing 1 illustrates the example of the KLAR dataset structure for the relation <i>capital</i> in English.		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
<b>A.2 Additional Experimental Results</b>		764
<b>A.2.1 Latent State Similarity</b>		765
Here, we present the complete results for latent state similarity across all language pairs in Figure 9.		766
		767
		768
The plots follow the same trend as in Figure 3b, where similarity across language pairs increases from early to middle layers in both models, indicating that MLMs encode information in a con-		769
		770
		771
		772

Relation	# Facts	Prompt Example
applies_to_jurisdiction	80	Which country has <subject> as a legal term? The answer is:
capital	336	What is the capital of <subject>? The answer is:
capital_of	213	Where is <subject> the capital of? The answer is:
continent	212	Which continent is <subject> located in? The answer is:
country_of_citizenship	60	Which country is <subject> a citizen of? The answer is:
developer	76	Which company is the developer of <subject>? The answer is:
field_of_work	167	What field does <subject> work in? The answer is:
headquarters_location	51	In which city is <subject>'s headquarter located? The answer is:
instrument	46	Which musical instrument is played by <subject>? The answer is:
language_of_work_or_name	108	What is the original language of <subject>? The answer is:
languages_spoken	104	What language did <subject> use to communicate? The answer is:
location_of_formation	66	Where did the formation of <subject> take place? The answer is:
manufacturer	35	Which company manufactures <subject>? The answer is:
native_language	130	What is the native language of <subject>? The answer is:
occupation	46	What is <subject>'s profession? The answer is:
official_language	602	What is the official language of <subject>? The answer is:
owned_by	50	Who is the current owner of <subject>? The answer is:
place_of_birth	35	In which city was <subject> born? The answer is:
place_of_death	79	In which city did <subject> pass away? The answer is:
religion	125	What is the religious belief of <subject>? The answer is:

Table 3: Relations in the KLAR dataset with fact counts and prompt examples used for knowledge probing.

<b>KLAR languages (17)</b>	Arabic (ar), Catalan(ca), Greek (el), English (en), Spanish (es), Persian (fa), French (fr), Hebrew (he), Hungarian (hu), Japanese (ja), Korean (ko), Dutch (nl), Russian (ru), Turkish (tr), Ukrainian (uk), Vietnamese (vi), Chinese (zh)
<b>LLaMA2 overlap (12)</b>	Catalan(ca), English (en), Spanish (es), French (fr), Hungarian (hu), Japanese (ja), Korean (ko), Dutch (nl), Russian (ru), Ukrainian (uk), Vietnamese (vi), Chinese (zh)
<b>BLOOM overlap (7)</b>	Arabic (ar), Catalan(ca), English (en), Spanish (es), French (fr), Vietnamese (vi), Chinese (zh)

Table 4: KLAR dataset languages and their overlap with LLaMA2 and BLOOM.

cept space independent of the input language. In the final layers, similarity declines as representations transition to a language-specific form. This pattern holds even for linguistically diverse pairs, highlighting that MLMs initially process factual knowledge in a shared latent space before adapting it to the target language.

## A.2.2 Latent Space Language Composition

We examine the language composition of the latent states in LLaMA2 and BLOOM to understand how these MLMs encode information in the concept space. As described in Section 5.4, we apply Logit Lens to project latent states to the vocabulary, and use fasttext to identify the language of the top-10 predicted tokens at each layer.

Figure 10 presents results for languages shared between LLaMA2 and BLOOM, while Figure 11 shows results for languages unique to each model.

LLaMA2’s middle-to-upper layers are dominated by English, aligning with prior findings that “LLaMA2 models think in English” (Wendler et al., 2024). In contrast, BLOOM displays a more

diverse linguistic composition in these layers.

Across different input languages, both models exhibit similar language distributions in the middle-to-upper layers, indicating that MLMs encode knowledge in a concept space largely independent of the input language.

## A.2.3 Rank Plots of Wrong Predictions

Figure 12 presents additional examples, one per language, where the correct English answer ranks highest in the middle-to-upper layers but is later surpassed by an incorrect target-language answer during the language transition phase.

## A.3 Shortcut Experimental Details

### A.3.1 Method

The idea of using linear approximation as a shortcut is inspired by Hernandez et al. (2023), who derive a linear transformation to approximate the mapping from subject to object representations in factual knowledge, showing that relational decoding in transformer models can be effectively modeled with linear functions.

```

{
  "relation_name": "capital",
  "relation_id": "P36",
  "prompt_templates": [
    "Where is <subject>'s capital located?
    ↪ The answer is:",
    "What is the capital of <subject>? The
    ↪ answer is:",
    "Which city serves as the capital of
    ↪ <subject>? The answer is:",
    "Name the capital city of <subject>. The
    ↪ answer is:",
    "Where does <subject> have its capital?
    ↪ The answer is:"
  ],
  "samples": [
    {
      "subject": "Azerbaijan",
      "object": "Baku",
      "index": 6152
    },
    {
      "subject": "Germany",
      "object": "Berlin",
      "index": 6165
    }
  ]
}

```

Listing 1: Example of KLAR for relation *capital* in English.

Building on this idea, we apply linear approximation to address cross-lingual inconsistency by bypassing the language transition process in MLMs. We hypothesize that the mapping from the model’s latent state at layer  $n$  to that at the final layer  $N$ , i.e.,  $h_n \rightarrow h_N$  can be well-approximated by a linear function  $f(h_n) = Wh_n + b \approx h_N$ . Following Hernandez et al. (2023), we use first-order approximation to estimate  $W_r$  and  $b_r$  as the mean Jacobian and bias across  $m$  correctly predicted factual samples  $\{h_{n_i}, h_{N_i}\}_{i=1, \dots, m}$ . That is, we define:

$$\begin{aligned}
 W_r &= \mathbb{E}_{h_{n_i}, h_{N_i}} \left[ \frac{\partial F}{\partial h_n} \Big|_{(h_{n_i}, h_{N_i})} \right], \\
 b_r &= \mathbb{E}_{h_{n_i}, h_{N_i}} \left[ h_{N_i} - \frac{\partial F}{\partial h_n} \Big|_{(h_{n_i}, h_{N_i})} h_{n_i} \right]
 \end{aligned} \quad (1)$$

As noted in Hernandez et al. (2023), the first-order derivative  $W_r$  tends to underestimate the magnitude of changes from  $h_n$  to  $h_N$  in practice.

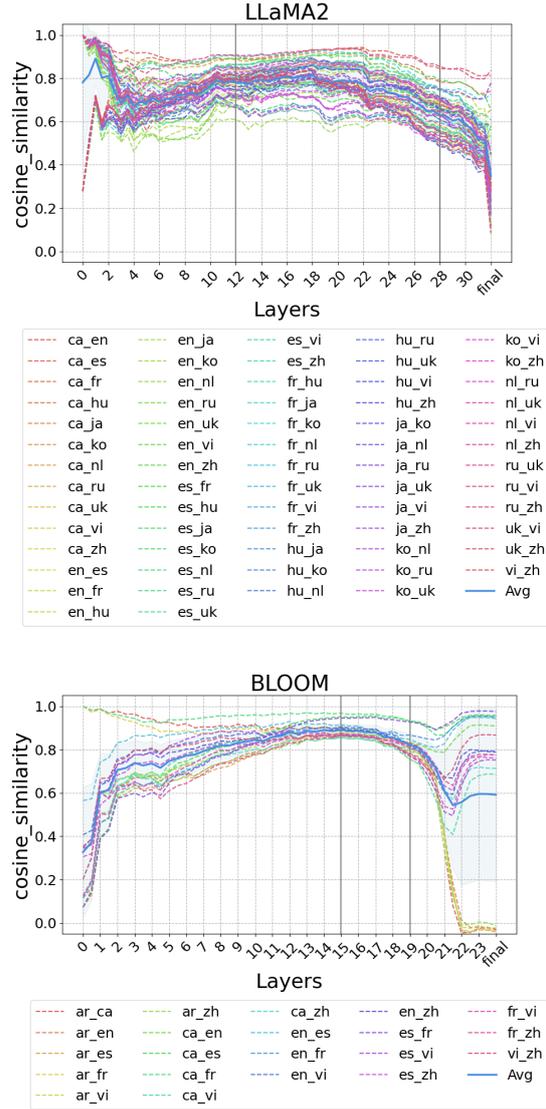


Figure 9: Cosine similarity of latent states between all language pairs averaged across all relation.

They attribute this to the use of layer normalization (Lei Ba et al., 2016) in transformers: which does not transmit changes in scale of inputs to changes in scale of output. Specifically, the input  $h_n$  at layer  $n$  is normalized before being propagated to subsequent layers. To address this underestimation, a scalar constant  $\beta$  is introduced as a hyperparameter and multiplied by  $W_r$  as a corrective factor:

$$f(h_n) = \beta W_r h_n + b_r = Wh_n + b \quad (2)$$

### A.3.2 Hyperparameters

Several hyperparameters are introduced when determining the linear shortcut  $f(\cdot)$ : the layer  $n$  from which the latent state is extracted for linear approximation, the scalar constant  $\beta$  used to adjust the

slope of  $W_r$  to account for the underestimation in the first-order approximation of  $h_n \rightarrow h_N$ , and the number of correct samples used to compute  $f(\cdot)$ . We perform a grid search to select these hyperparameters for each language, aiming to maximize prediction accuracy. For the layer  $n$ , we search within the range of  $[20, 32]$  for LLaMA2 and  $[12, 24]$  for BLOOM. The scalar constant  $\beta$  is searched over the range  $[0, 5.0]$  in increments of 0.25, following Hernandez et al. (2023). The number of samples  $m$  is selected from  $[10, 25, 40, 50]$ . The hyperparameter search is conducted for each language individually. We find that the optimal  $\beta$  value varies across languages, while the other two hyperparameters — the extraction layer  $n$  and the number of samples  $m$  — remain consistent across languages. The selected hyperparameters for both models are summarized in Table 5 and 6, respectively.

LLaMA2	$n$	$\beta$	$m$
ca	30	4.75	25
en		1.50	
en		1.50	
es		3.00	
fr		4.25	
hu		2.50	
ja		2.25	
ko		4.50	
nl		3.50	
ru		4.25	
uk		2.25	
vi		1.00	
zh		1.50	

Table 5: Hyperparameters per language for LLaMA2.

BLOOM	$n$	$\beta$	$m$
ar	21	1.25	25
ca		1.00	
en		1.25	
es		1.00	
fr		0.75	
vi		1.25	
zh		1.50	

Table 6: Hyperparameters per language for BLOOM.

### A.3.3 Shortcut Translation Baselines.

As mentioned in Section 7.2, we compare our shortcut method with two translation-based baselines: (1) translation-en (trans-en): We translate all input queries from each language to English using Google Translate, obtain model predictions in English, and then translate them back to the target language to measure accuracy. (2) translation-early-exit (trans-exit): We use Logit

Lens to project the latent states at the same extraction layers as in the shortcut method, i.e., layer 30 for LLaMA2 and layer 20 for BLOOM, and extract the top-predicted tokens. These tokens are then translated into the target language using Google Translate, and their accuracy is calculated against the correct object.

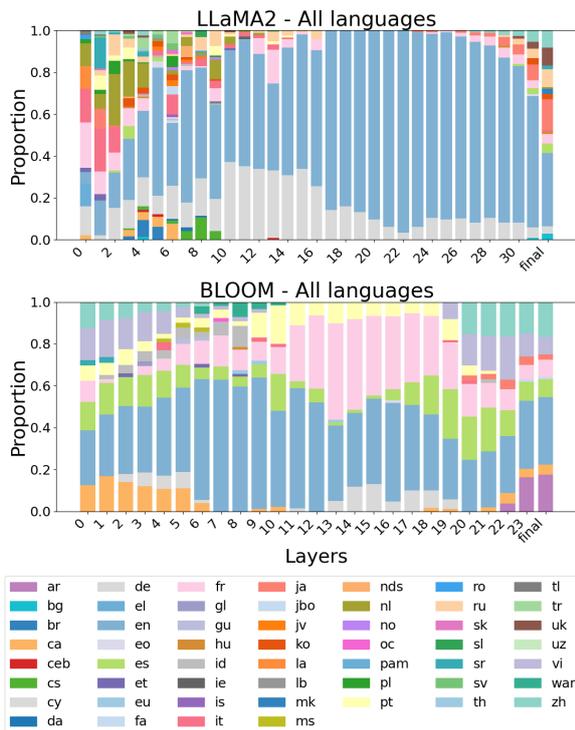
As shown in Table 7 and 8, both translation-based methods perform poorly. The low accuracy of *translation-en* suggests that existing translators struggle with entity translation, especially for languages that are highly dissimilar to English. The poor performance of *translation-early-exit* stems from the inherent unreliability of token-level translations. Overall, these results indicate that translation-based approaches are not a viable solution for cross-lingual factual prediction. In contrast, by directly adapting latent representations from earlier layers, the shortcut method operates at the representation level, capturing richer contextual information. This enables significantly higher prediction accuracy and offers a promising solution for mitigating cross-lingual factual inconsistency.

LLaMA2	original	shortcut	trans-en	trans-exit
ca	76.96	<b>80.54</b>	44.95	24.52
en	81.41	<b>85.06</b>	81.41	43.05
es	78.44	<b>81.16</b>	47.77	28.42
fr	78.14	<b>82.46</b>	53.27	24.85
hu	75.69	<b>79.04</b>	64.60	6.91
ja	63.05	<b>70.45</b>	59.59	0.13
ko	62.14	<b>66.98</b>	49.30	0.28
nl	77.22	<b>80.77</b>	62.07	15.24
ru	67.02	<b>72.71</b>	47.58	2.72
uk	70.46	<b>74.78</b>	46.59	5.62
vi	73.26	<b>77.56</b>	39.07	12.70
zh	53.88	<b>61.40</b>	60.38	1.67

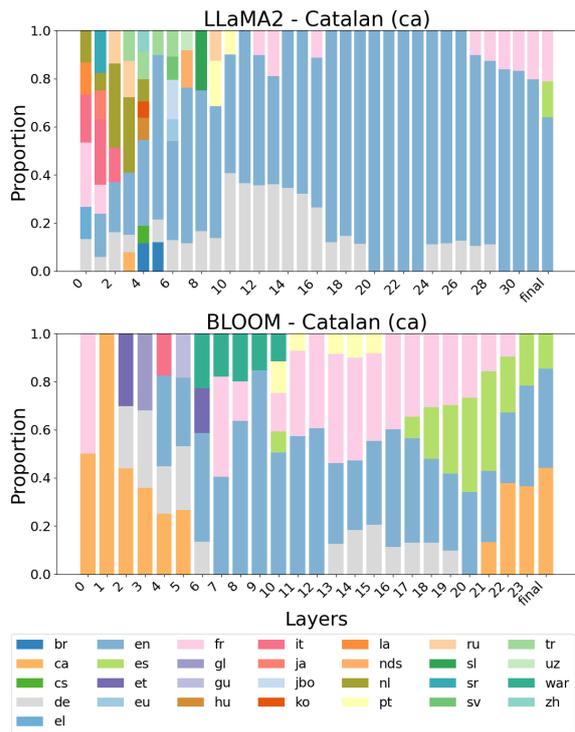
Table 7: Comparison of the prediction accuracy (%) for LLaMA2 across different languages using the original model, the proposed shortcut method, and the translation-based baselines.

BLOOM	original	shortcut	trans-en	trans-exit
ar	31.58	<b>37.93</b>	21.87	0.97
ca	41.50	<b>48.40</b>	22.58	15.88
en	46.81	<b>58.24</b>	46.81	26.85
es	43.56	<b>54.84</b>	25.53	11.26
fr	46.88	<b>56.03</b>	26.15	17.97
vi	56.82	<b>62.38</b>	21.98	25.85
zh	35.54	<b>43.89</b>	31.26	10.96

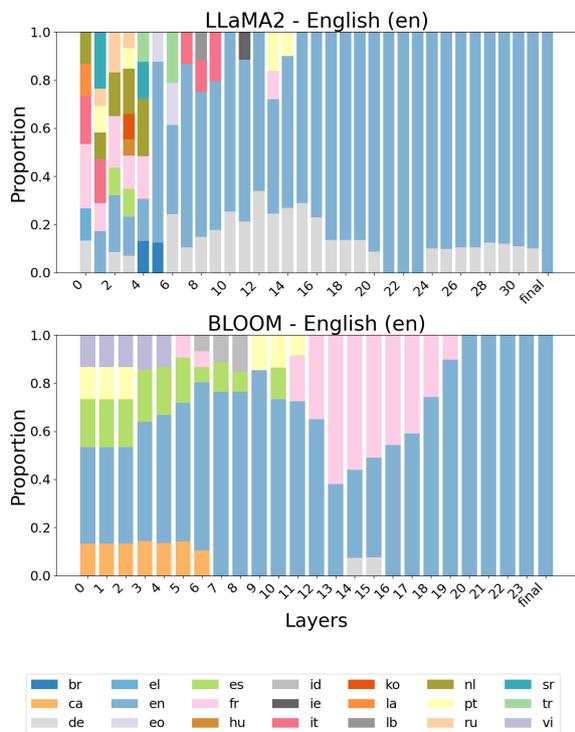
Table 8: Comparison of the prediction accuracy (%) for BLOOM across different languages using the original model, the proposed shortcut method, and the translation-based baselines.



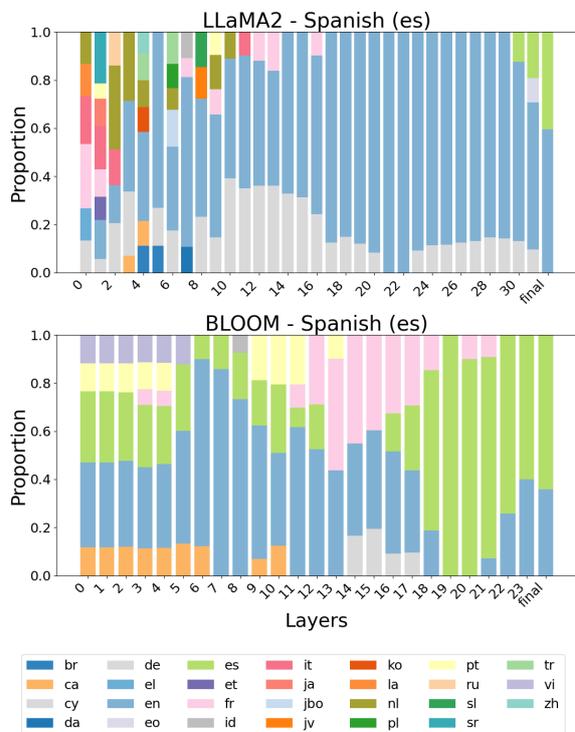
(a) Language composition aggregated across all languages



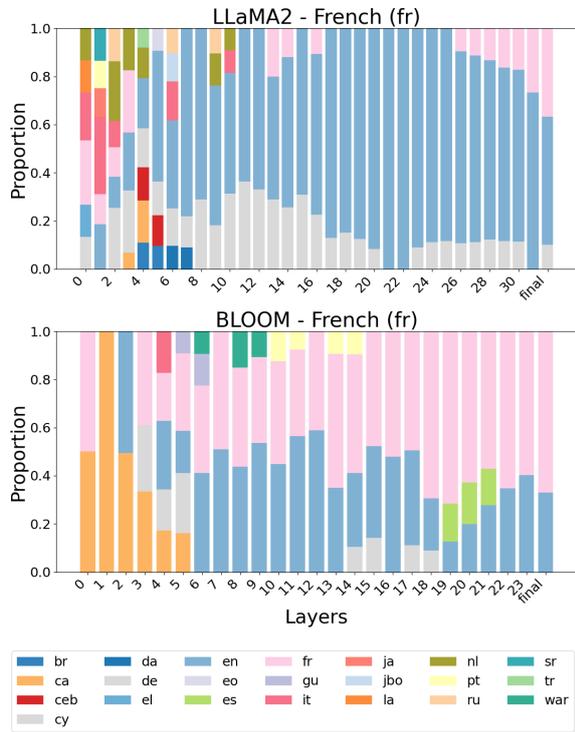
(b) Language composition with Catalan as the input language.



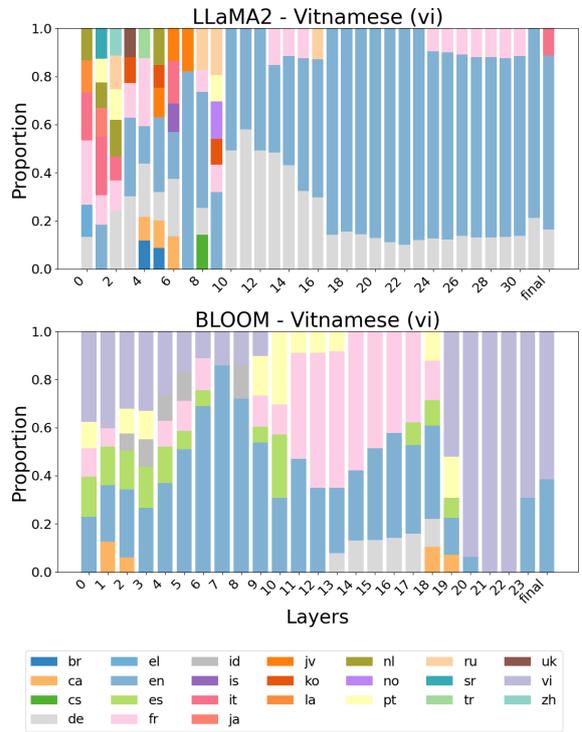
(c) Language composition with English as the input language.



(d) Language composition with Spanish as the input language.

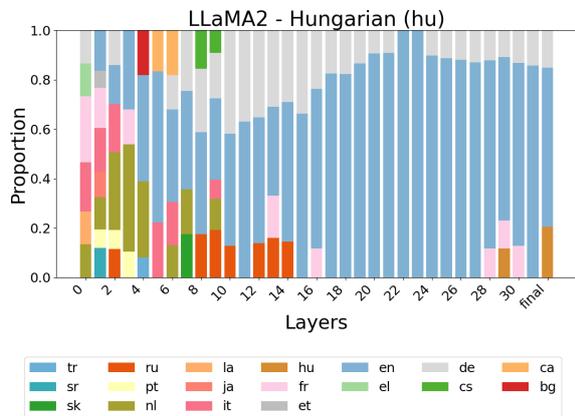


(e) Language composition with French as the input language.

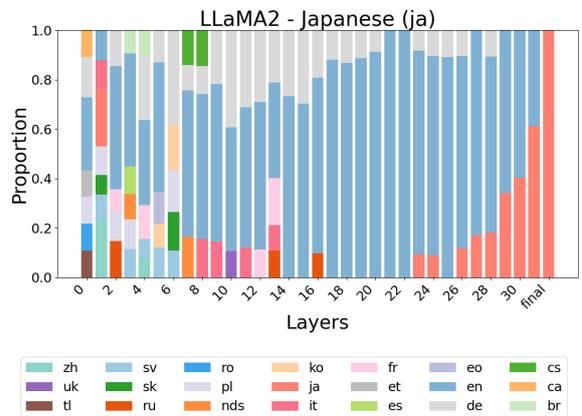


(f) Language composition with Vietnamese as the input language.

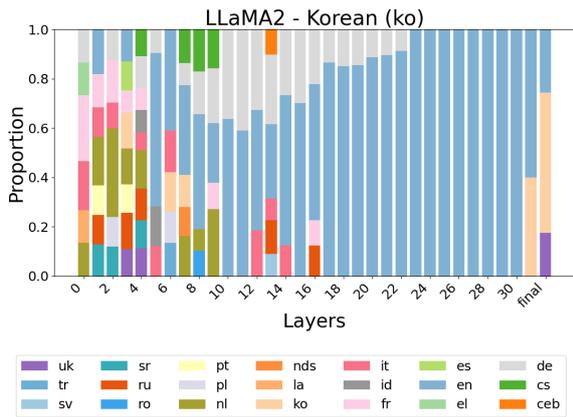
Figure 10: Language composition for languages shared between LLaMA2 and BLOOM.



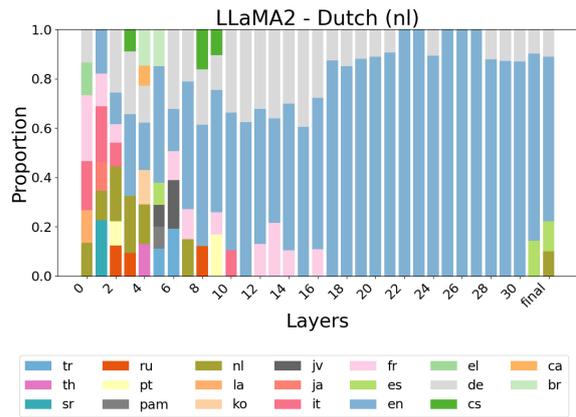
(a) Language composition in LLaMA2 with Hungarian as the input language.



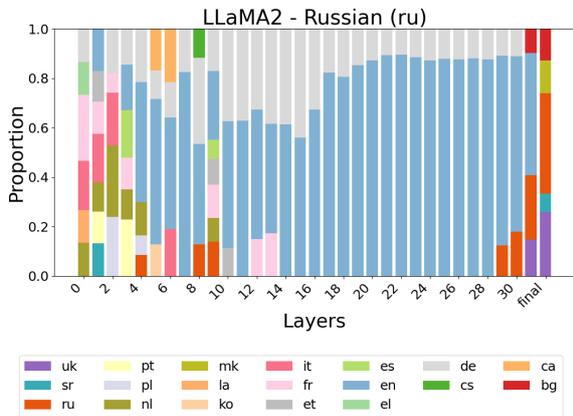
(b) Language composition in LLaMA2 with Japanese as the input language.



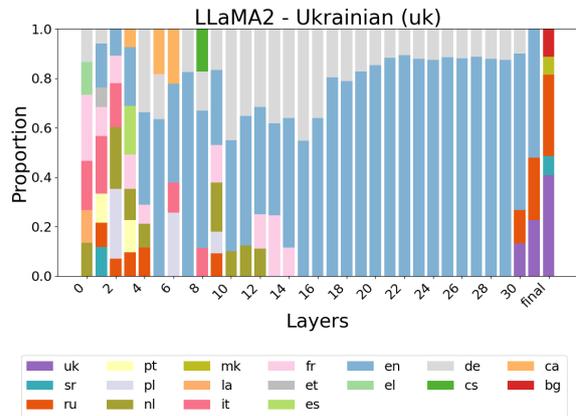
(c) Language composition in LLaMA2 with Korean as the input language.



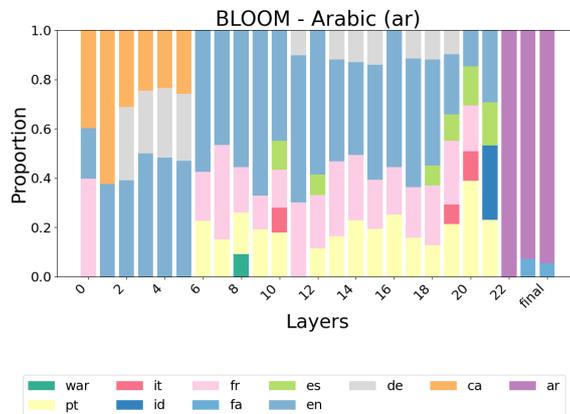
(d) Language composition in LLaMA2 with Dutch as the input language.



(e) Language composition in LLaMA2 with Russian as the input language.

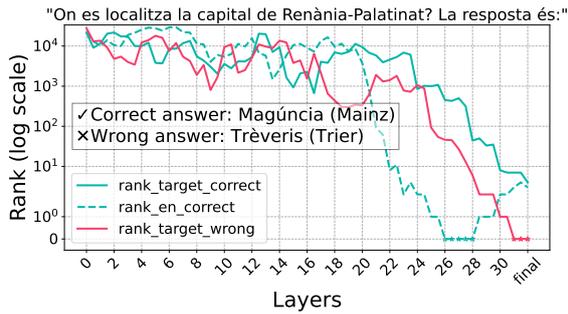


(f) Language composition in LLaMA2 with Ukrainian as the input language.

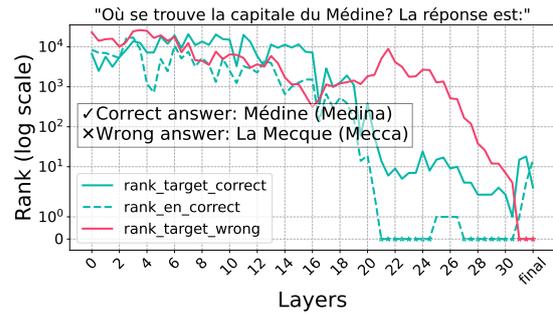


(g) Language composition in BLOOM with Arabic as the input language.

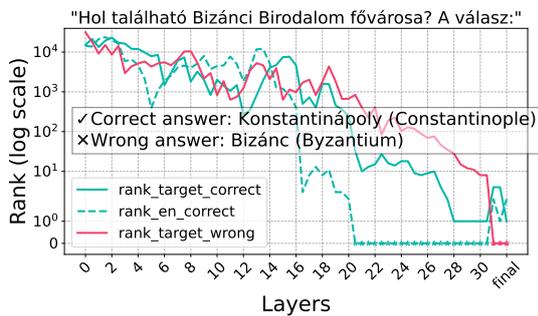
Figure 11: Language composition for unique languages in LLaMA2 and BLOOM, respectively.



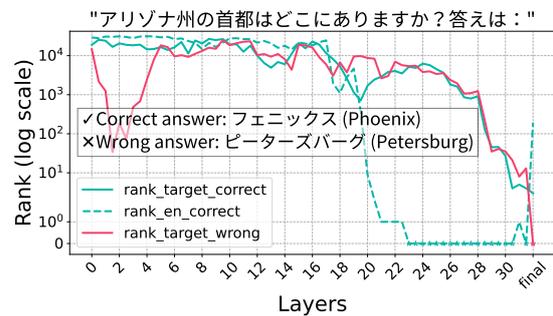
(a) Prompt in Catalan; English translation: "What is the capital of Rhineland-Palatinate? The answer is:".



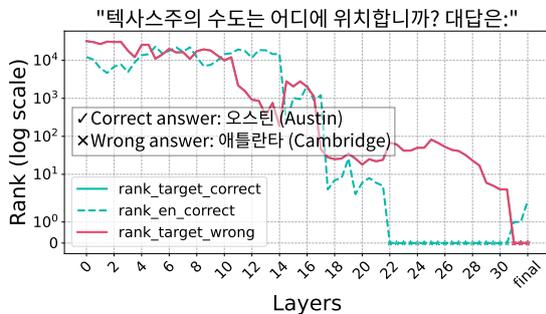
(b) Prompt in French; English translation: "What is the capital of Medina? The answer is:".



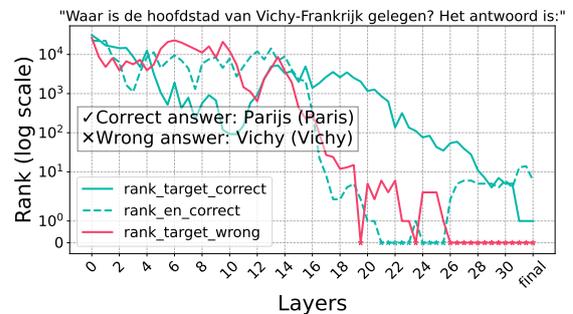
(c) Prompt in Hungarian; English translation: "What is the capital of Byzantine Empire? The answer is:".



(d) Prompt in Japanese; English translation: "What is the capital of Arizona? The answer is:".



(e) Prompt in Korean; English translation: "What is the capital of Texas? The answer is:".



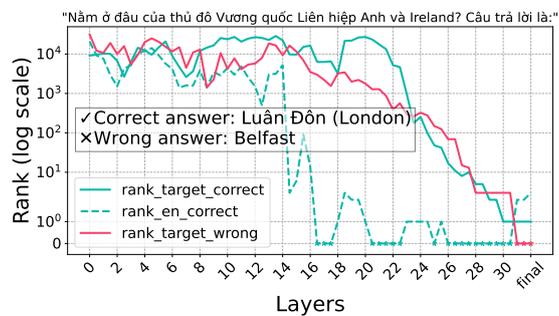
(f) Prompt in Dutch; English translation: "What is the capital of Vichy France? The answer is:".



(g) Prompt in Russian; English translation: "What is the capital of Andalusia? The answer is:".



(h) Prompt in Ukrainian; English translation: "What is the capital of Guyana? The answer is:".



(i) Prompt in Ukrainian; English translation: “What is the capital of United Kingdom of Great Britain and Ireland? The answer is:”.

Figure 12: Rank evolution for prompts in different languages. rank\_target\_wrong represents the rank of the model’s final incorrect prediction across layers, while rank\_target\_correct and rank\_en\_correct denote the ranks of the correct answer in the target language and the English equivalent, respectively. The plots show the impact of errors during language transition, where the rank of the incorrect answer surpasses the correct answer in the final layers.