000 IS KNOWLEDGE IN MULTILINGUAL LANGUAGE MOD-001 ELS CROSS-LINGUALLY CONSISTENT? 002 003

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

025 026

027 028 029

037

041

Paper under double-blind review

ABSTRACT

Few works study the variation and cross-lingual consistency of factual knowledge embedded in multilingual models. However, cross-lingual consistency should be considered to assess cross-lingual transferability, maintain the factuality of the model's knowledge across languages, and preserve the parity of language model performance. We are thus interested in analyzing, evaluating, and interpreting cross-lingual consistency for factual knowledge. We apply interpretability approaches to analyze a model's behavior in cross-lingual contexts, discovering that multilingual models show different levels of consistency, subject to either language families or linguistic factors. Further, we identify a cross-lingual consistency bottleneck manifested in middle layers. To mitigate this problem, we try vocabulary expansion, additional cross-lingual objectives, adding biases from monolingual inputs, multi-task fine-tuning, and code-switching training. We find that all these methods, except for multi-task fine-tuning, boost cross-lingual consistency to some extent, with cross-lingual supervision and code-switching training offering the best improvement.

1 INTRODUCTION



Figure 1: Illustration of Cross-lingual knowledge consistency. Frege's theory of reference defines the reference of a sub-sentential expression as the object singled out by the name.

Frege's theory of reference (Frege, 1892) indicates that the knowledge conveyed by a sentence depends on the references of the expressions that compose the sentence. A salient aspect of humanity 040 is that while people may speak different languages, they can share common knowledge. Thus, references and knowledge should be consistent across languages, and a multilingual model serving as 042 a knowledge base (Gupta & Srikumar, 2021; Kassner et al., 2021; Hu et al., 2024) should provide 043 consistent knowledge when queried in different languages. Not only does this theory contribute to 044 the cross-lingual performance and maintain the knowledge across languages, but it ensures parity 045 and self-consistency of model performance (Hupkes et al., 2023; Wang et al., 2023). This motivates us to evaluate the knowledge consistency of multilingual language models across languages. 046

047 Although recent advances have shown that multilingual models are effective for cross-lingual trans-048 fer and generalization (Conneau et al., 2020; Xue et al., 2020; Hu et al., 2020; Muennighoff et al., 2023), Kassner et al. (2021); Fierro & Søgaard (2022); Qi et al. (2023) reported that the prediction varies from one language to others when recalling knowledge in different languages. For given par-051 allel statements, a language model may output predictions for a particular query that differs from one obtained from the query's translation. This examination implicitly instantiates Frege's theory of 052 reference to check knowledge consistency across languages that parallel statements with the same references for sub-sentential expressions such as entities should have the same knowledge.

054 Inspired by that theory, we hypothesize that multilingual language models recall consistent factual 055 knowledge for coreferential statements in cross-lingual settings. To evaluate this hypothesis (Fig-056 ure 1), we create code-mixed coreferential statements from monolingual statements by substituting 057 a subject entity with an equivalent one in another language that shares the same reference and at-058 tempt to answer two questions: 1) do multilingual language models recall factual knowledge for the coreferential statements in a similar manner on different languages? and 2) how does the mechanism of multilingual language models work on the incorporation between entities or references in 060 cross-lingual settings? Our study is related to a broader linguistic phenomenon of entity-level code-061 switching: an entity code-switches between two languages without changing the reference. More 062 recently, we share a similar goal with knowledge incorporation and editing (Beniwal et al., 2024; 063 Li et al., 2024), as we incorporate a coreferential entity from other languages for factual knowledge 064 recall in cross-lingual settings. Our main findings are: 065

- Multilingual language models can leverage coreferential entities across similar languages, e.g., similar writing scripts, to maintain cross-lingual knowledge consistency when recalling factual knowledge but worse across dissimilar languages.
 - Scaling models is not a promising strategy to improve cross-lingual consistency as we observed a bottleneck starting from middle layers across different model families and sizes.
 - Consistency patterns in feed-forward neurons and subject-object attention scores are indistinguishable for similar languages and distinguishable for dissimilar languages.

Based on our main findings, we further evaluate several mitigations that could resolve inconsistency issues with the results described below, particularly those observed in dissimilar languages.

- There is a partial causality between adding monolingual biases and improving cross-lingual knowledge consistency. Thus, adding bias could be a potential method to calibrate consistency across languages while we do need to think a better way to incorporate such bias.
 - Expanding the multilingual vocabulary, adding word alignment objective, and codeswitching training can improve the cross-lingual consistency as such method helps in aligning coreferential entities across languages to alleviate the consistency bottleneck.
 - Multi-task fine-tuning is not promising to benefit the cross-lingual consistency as it potentially refines specific attention heads for in-context information instead of knowledge.

Our contribution is to offer an understanding of multilingual language models' limitations under cross-lingual settings and highlight potential research directions to address such issues.

2 METHODOLOGY

066

067

068

069

071

073 074

075

076

077

078

079 080

081

082

084

085

087 088 089

090 091

106

107

091 2.1 TASK DEFINITION

093 We focus on a code-mixed context-independent cloze task. This setting forces the multilingual 094 model to rely on its internal knowledge base and recall the common knowledge shared by coreferential entities across languages because of cross-lingual generalization ¹. In the following introduction, 095 we will define the evaluation task mathematically. Let $I = \{S^{l1}, \dots, O, \dots\} \in l1^2$ be a statement, 096 where l1 stands for matrix language (the predominant language), $S^{l1} = \{s_1, \dots, s_k\} \in l1$ are subject sub-tokens, and $O = \{o_1, o_2, \dots, o_j\} \in l1$ denote object sub-tokens. This statement is used to create a masked input $I_{mono} = \{S^{l1}, \dots, M, \dots\}$, where $M = \{mask_1, \dots, mask_j\}$ are the 098 mask (or the sentinel token $M = \langle extra_i d_0 \rangle$) used to substitute O in I. We define n-gram 100 prediction for the mask O, denoted as $Cand(O_{\in V}|I_{mono})$, as the top-k n-gram candidates with 101 token lengths ranging from 1 to j obtained from beam search decoding over the model's vocabu-102 lary V. Considering the trade-off between the computational time and the holistic of the language 103 model's prediction, we set the top-k threshold and beam search width to 5. Similarly, we can create 104 a code-mixed coreferential statement I_{cm} by replacing S^{l1} with a coreferential subject S^{l2} in the 105 embedded language l2 (the subsidiary language) in order to obtain $Cand(O_{\in V}|I_{cm})$. Therefore,

¹See limitation in §6.

²The surface structure is not restricted. We use the common subject–object structure as an example.

108 I_{cm} and I_{mono} are coreferential and expected to recall the same knowledge. Finally, we define 109 cross-lingual knowledge consistency as $0 \le f_{metric}(Cand(O_{\in V}|I_{mono}), Cand(O_{\in V}|I_{cm})) \le 1$, 110 where f_{metric} is a consistency metric defined in the next subsection. If $f_{metric} = 1$, it implies that 111 multilingual language models recall factual knowledge for the coreferential statements I_{mono} and 112 I_{cm} in an identical manner. The coreferential statements disagree if $f_{metric} = 0$. Note that we do 113 not consider whether the prediction is correct. Instead, f_{metric} evaluates the parity and consistency 114 across the languages that the model is expected to output similar candidates for I_{mono} and I_{cm} .

115 From a probability view, we can define our task as measuring the difference between two dis-116 tributions, $Cand(O_{\in V}|I_{cm}) = P(O_{\in V}|K_{\theta})P(K_{\theta}|S^{l_2}, I_{\backslash (S^{l_1}\cap O)})$ and $Cand(O_{\in V}|I_{mono}) =$ $P(O_{\in V}|K_{\theta}^*)P(K_{\theta}^*|S^{l_1}, I_{\setminus (S^{l_1}\cap O)})$, where K_{θ} is the knowledge recalled from the model given the preceding context, and $I_{\setminus (S^{l_1}\cap O)}$ stands for I without both the subject and the object. Then, 117 118 119 cross-lingual knowledge consistency between K_{θ}^* and K_{θ} reflects on the measured difference. The 120 high-level idea of this evaluation task is illustrated in Figure 1 where en entry "Paris is the capital of ____" is evaluated with its possible code-mixed statements (ar entry & ta entry). In this ex-121 ample, S^{l1} , $I_{\backslash (S^{l1} \cap O)}$, and S^{l2} are "Paris", "is the capital of", and the ar or ta entry for "Paris", respectively. If coreferential subject entries are trained to generalize across languages, we could 122 123 observe the cross-lingual consistency. In addition, we are aware of a baseline from this probabil-124 ity view. Specifically, we define the baseline as the difference between $Cand(O_{\in V}|I_{mono})$ and 125 $Cand(O_{\in V}|I_{\setminus (S^{l_1}\cap O)}) = P(O_{\in V}|K_{\theta}^{\alpha})P(K_{\theta}^{\alpha}|I_{\setminus (S^{l_1}\cap O)})$, measuring agnostic consistency without 126 the coreferential subjects S^{l1} and S^{l2} in cross-lingual settings. In implementation, we mask the both 127 subject and object entities to create the "code-mixed" counterpart as the baseline. Readers can refer 128 to Appendix §A.1 for our implementation. 129

- 1302.2METRIC FUNCTION AND INTERPRETABILITY APPROACH
- Readers can refer to Appendix §A.2 for more details, e.g., equations.

Consistency Metrics. For f_{metric} , **Top@1 Accuracy** and **RankC** (i.e., weighted Precision@5) (Qi et al., 2023) are used to evaluate the cross-lingual knowledge consistency between $Cand(O_{\in V}|I_{mono})$ and $Cand(O_{\in V}|I_{cm})$. Since we observe similar experimental results on Top@1 and RankC, Top@1 results are moved to Appendix §A.3.

Consistency Evolution. We analyze the "evolution" of consistency scores as the layer goes deeper to trace the consistency bottleneck and understand the models' behavior. Since the encoder part is crucial for both encoder and encoder-decoder language models to understand the input, we apply LogitLens (nostalgebraist, 2019) (for the xlm-r family) and DecoderLens (Langedijk et al., 2023) (for the mT0 family) to each encoder layer to obtain the layer-wise distribution $O_{\in V}$.

Subject–Object Attention. Inspired by the attention weight analysis method (Clark et al., 2019), we calculate the sum of all subject tokens' attention scores across all masked tokens and average those scores over all possible I_{mono} and I_{cm} statements. Then we collect the difference between average attention scores of I_{cm} and those obtained from I_{mono} . Note that masked tokens are used to prompt the corresponding object tokens of the masked statements, as defined in our task definition.

150 IG^2 Score We do minor modification on IG^2 (Liu et al., 2024) to measure the impact of each 151 feed-forward neuron on the logits of the mask tokens where the higher the value is, the more critical 152 the neuron is to predict the ground truth object on the mask tokens.

153 154

138

2.3 DATASET AND MODEL

Dataset. We use mLAMA dataset (Kassner et al., 2021) that provides parallel triples (object, predicate, subject) in 53 languages written in cloze task format (e.g., "Paris is the capital of [MASK].") to query knowledge in zero-shot settings. In our experiments, *l*1 is set to English. Meanwhile, we set *l*2 to all other 52 languages to report an overall result and rendered deep analysis for 2 similar *l*2 languages (De, Id) and 2 dissimilar *l*2 languages (Ar, Ta)³. Moreover, for the overall result, we also

¹⁶¹

³While Id does not belong to the same language family as En, it has many similarities with En (Krause, n.d.). Ar and Ta are not considered as Indo-European languages and also do not use latin scripts.

categorize these *l*2 languages into two separate categories for each of the three factors (geographics, writing scripts, and language family) using ISO-639 language code information from "localizely"⁴.

Models. We examine distinct language model families: xlm-r (0.3B to 10B) (Conneau et al., 2020)
and mT0 (0.6B to 3.7B) (Muennighoff et al., 2023). Decoder-based language models are excluded
in our study to limit inherent hallucination problems affecting the analysis (Xu et al., 2024; Ji et al., 2023; Fu et al., 2023). In our experiments, we obtain similar findings from both families. Therefore, we only show mT0 results in the main text and move the rest to the Appendix §A.3.

3 OBSERVING CONSISTENCY

3.1 MAIN FINDINGS

170 171

172 173

174 175 176

177 178 179

181

185



Figure 2: Overall cross-lingual consistency in mt0 (red: mt0-large, blue: mt0-base) grouped by 3 factors (left: geographics, mid: language family, right: writing scripts). Note: The dashed line here is the average corresponding consistency scores of mt0-base across languages(*cf.* §A.3.2.)

186 From Figure 2, across all the factors, l^2 that are dissimilar with l_1 , tend to have lower consistency 187 than those are similar to l_1 . The difference in writing scripts plays the most important role among the other two factors. Thus, the number of shared tokens between two languages could affect the 188 cross-lingual consistency, and that is orthogonal to the common view of cross-lingual transfer that 189 shared tokens are not necessary (K et al., 2020; Artetxe et al., 2020). Another intriguing finding 190 is that geographic factor also affects consistency and this could be attributed to common culture 191 and vocabulary (Zhao et al., 2024a). On the other hand, we suppose that other linguistic factors 192 contributing to the cross-lingual performance (de Vries et al., 2022; Kann et al., 2017; Chronopoulou 193 et al., 2023) such as the similarity in linguistic features (Chronopoulou et al., 2023), or borrowing 194 (Tsvetkov & Dyer, 2016), could affect cross-lingual knowledge consistency as well. However, for 195 this study, it is hard to quantify such factors and leave such analyses for future work. Furthermore, 196 we also observe similar results on other models, and this aligns with empirical studies in the literature 197 (Qi et al., 2023). Note that language families and writing scripts have an impact on vocabulary, and 198 we will discuss this vocabulary problem in a later section.

To better understand the cross-lingual consistency bottleneck, we examine the layer-wise consistency patterns across different model sizes, as presented in Figure 3. The noticeable difference lies in the initial consistency, whereby dissimilar language pairs have low consistency scores. The consistency gap between dissimilar and similar languages starts to close at some specific layer while widening again later. This observation provides evidence for empirical studies that scaling benefits the downstream task performance (Conneau et al., 2020), e.g., XNLI, but not be substantially helpful







205

Figure 3: Layer-wise consistency scores in models with different sizes. Scaling models is not a promising strategy to mitigate consistency bottlenecks. (*cf.* §A.3.1)

in refining cross-lingual consistency due to cross-lingual consistency bottlenecks. Moreover on Figure 3, we can see the consistency scores of dissimilar languages are quite similar with baseline thus, the consistency of such code-mixed languages are terrible and needed to be improved. Nonetheless, these dissimilar languages are more consistent than our baseline on xlm-r models (*cf.* 18).

3.2 ATTENTION WEIGHT ANALYSIS



Figure 4: Subject–Object attention difference with I_{mono} to I_{cm} in mt0-large. Attention scores are inversely proportional to the similarity of l1 and l2 (En–Ta vs En–Id). (*cf.* §A.3.3)

235 The layer-wise analyses help us understand the model behaviors. However, the question remains as to how model 236 components handle statements. The correlation analy-237 sis conducted on each layer in Table 1 shows that there 238 is a moderate correlation between the average negative 239 subject-object attention scores and the consistency met-240 rics. In particular, although I_{mono} and I_{cm} are coreferen-241 tial, I_{cm} has to retrieve the reference via the cross-lingual 242 entry. To identify the difference between the cross-lingual 243 and monolingual entries, we observe the attention scores 244 for subject-object pairs. Figure 4 demonstrates that the 245 attention scores across layers and heads are barely distin-246 guishable in the similar language pair, en-id, but more 247 discernable for the dissimilar language pair, en-ta. Sur-

Table 1: Statistical spearman ρ correlation ($\alpha = 0.05$) between average scores of layers with the patterns on each language model's subject-object attention and IG^2 absolute difference.

	attention		IG^2	
Model	RankC	Acc	RankC	Acc
mT0-base mT0-large xlm-r-base xlm-r-large	0.414^{*} 0.666^{*} 0.433^{*} 0.671^{*}	0.426^{*} 0.661^{*} 0.424^{*} 0.666^{*}	0.528^{*} 0.705^{*} 0.400^{*} 0.508^{*}	0.519^{*} 0.699^{*} 0.397^{*} 0.481^{*}

prisingly, I_{cm} results in higher attention scores than I_{mono} for the dissimilar language pair, which means that the model pays more attention to the subject entity over the predicate across layers and heads in that case. Despite so, we observe from the right side of Figure 4 that having subject-object attention that is too small/big might cause the model to be inconsistent. These insights might be generalizable to different models and other language pairs, as evidenced by Appendix §A.3.3.

3.3 IG^2 Score Analysis



Figure 5: IG^2 scores in mt0-large for en–de and en–ta. We see that the distribution is more contrastive on dissimilar languages (En–Ta) than the similar languages (En–ID). (*cf.* §A.3.4.)

In addition, we inspect the IG^2 scores of all feed-forward neurons across all encoder layers. Our correlation analysis for this factor could show a moderate correlation with the cross-lingual consistency, as shown in Table 1. In Figure 5, the IG^2 scores for similar language pairs are almost the same, while there is a subtle difference for the dissimilar language pairs.

5

251 252 253

> 254 255

> 256

257

261 262 263

264

265 266

220 221

222

232

²⁷⁰ 4 IMPROVING CONSISTENCY

272

273

284

293

295

296

297

298

299

300

301

302 303 304

305

306

307

308

309 310

311

312

313 314 315

4.1 CAN BIAS CALIBRATE CONSISTENCY?

274 From previous findings, we think of one question: can we add biases from Imono to attention layers 275 and feed-forward layers for consistency calibration? Thus, based on two different patterns dis-276 covered from our experiments and having that both are moderately correlated with the consistency score, we do three different causal interventions to align the output of I_{cm} closer to the output of 278 I_{mono} . This experiment measures whether each pattern has a causal relationship with cross-lingual consistency. The experimental setup can be seen in Appendix §A.4.1. We consider: Attention 279 suppression: suppressing I_{cm} 's attention scores to make them closer to I_{mono} 's, Feed-forward 280 neuron activation patching (Vig et al., 2020; Geiger et al., 2021): patching I_{mono}'s activations 281 of all tokens to I_{cm} in selected feed-forward neurons based on IG^2 , and **Hybrid**: using the above 282 methods simultaneously. 283



Figure 6: Intervention scores for En-Ta. (cf. §A.4.1).

Based on Figure 6, there is a causal relationship between the two studied factors to some certain extent. For mT0, intervention approaches can increase the consistency scores in the middle-later layers only. While for xlm-r, none of the interventions manages to improve the consistency for the base model despite there is a rising trend for the hybrid intervention and feed-forward activation patching. However, when we observe the larger model, all intervention methods increase the consistency score on the last layers. Overall, our monolingual bias incorporation offers quite limited improvement subject to the architecture and model size thus a better monolingual bias should be considered. These findings are consistent with those of other languages.

4.2 THE EFFECT OF VOCABULARY EXPANSION TO THE CROSS-LINGUAL CONSISTENCY

We hypothesize that the size of vocabulary plays a crucial role in improving consistency as this enables a language model to align the semantics better due to lesser ambiguity⁵. To test this hypothesis, we consider two similar language models, xlm-r-base and xlm-v-base (Liang et al., 2023), where xlm-v-base has larger vocabularies (901,629 tokens) than xlm-r-base (250,002 tokens).

⁵e.g., if the tokenizer of a language model tokenizes the word "Tokyo" to ["To," "Kyo"], the token "To" is polysemous making thus the alignment of this word would be one-to-many, on the other hand, if a tokenizer keeps the word as it is, the tokenized form of the word is monosemous making it less ambiguous.



Figure 7: Effects of vocabulary expansion to cross-lingual consistency (red: xlm-v, blue: xlm-rbase). (*cf.* §A.4.2)

v



Figure 8: Effects of vocabulary expansion to subject–object attention scores shift to xlm-r-base where we can see there is a shift from earlier layers to middle & last layers. We do not see any notable shift on the baseline input and this is possibly the reason why this method does not offer improvement on the baseline input. (*cf.* §A.4.2)

338 Vocabulary expansion offers slight consistency improvement for 339 similar and dissimilar language on any categorization, as pre-340 sented in Figure 7. A large vocabulary limits sub-tokens to pre-341 vent the model from latching onto shallow local signals or restor-342 ing words from sub-tokens (Levine et al., 2021), which benefits 343 deep semantic learning. However, more samples are required 344 to generalize training. Therefore, vocabulary expansion alone 345 cannot improve the consistency significantly, especially for low-346 resource languages, but it sitll benefits dissimilar languages with lower consistency in the last layers to alleviate the consistency 347 bottleneck to some extent. This can be observed in Figure 9 348 that the layer-wise consistency drops significantly in the base 349 model's last layers but increases in the expanded model's last 350 layers. Meanwhile, it can be further confirmed by a similar shift 351



Figure 9: Effects of vocabulary expansion to cross-lingual consistency. (*cf.* §A.4.2)

phenomenon in attention analysis, as shown in Figure 8. The attention scores on the expanded model's last layer or deep semantic layers are more potent than the base model. Since dissimilar language pairs have minimum language features shared across languages, better-aligned semantics alone is not enough to completely resolve the consistency bottleneck. Nonetheless, it is still a crucial aspect to consider for the cross-lingual consistency improvement as demonstrated on the minor improvement shown on Figure 7. It is also in line with Zhao et al. (2024a) where they found that the one-token P@1 of Afrikaans is higher than the Japanese due to segmentation and tokenziation. ⁶

358 359

360

333

334

335

336 337

4.3 THE EFFECT OF CROSS-LINGUAL SUPERVISION TO THE CROSS-LINGUAL CONSISTENCY

Another possible hypothesis is that there might be an entangle-361 ment of features between linguistic and knowledge features. El-362 hage et al. (2022) discovered that the language model (in partic-363 ular GPT-2 (Radford et al., 2019)) could fit multiple features into 364 one dimension at the price of more entangled features, and this entanglement might cause tokens not cross-lingually aligned as 366 there may be an entanglement between syntactic and semantic 367 features within one dimension. Inspired by that, we suspect this 368 might hinder the consistency of language models. To test this 369 assumption, we evaluate two similar language models in which 370 one model is trained solely on MLM objective (xlm-r), and another similar model is trained on one additional objective to align 371 word translations (xlm-align (Chi et al., 2021)), where this word 372 alignment might be helpful to align references across languages. 373



Figure 10: Effects of crosslingual supervision on the layerwise consistency. (cf. 4.3)

- Word alignment increases cross-lingual consistency monotonically to alleviate the cross-lingual bottleneck. Similar to the vocabulary expansion, this strategy does not improve the consistency for the
- 376 377

⁶We define this as a token parity issue. See more discussion in Figure 34).



Figure 11: Effects of cross-lingual supervision on subject-object attention in xlm-r-base. We observe a slight layer shift. A similar pattern is observed as the attention difference caused by the vocabulary expansion on the baseline input. left: En–Ar, mid: En–Ta, right: baseline. (*cf.* 4.3)



Figure 12: Effects of cross-lingual supervision to xlm-r-base consistency red: xlm-align and blue: xlm-r-base. Note: The dashed line here is the average corresponding consistency scores of xlm-r-base across languages. (*cf.* §A.4.3).

403 baseline as we would expect. The aligned model outperforms the baseline starting from the middle 404 layers in Figure 10. Multiple pre-training objectives that could approximately disentangle different 405 features can help preserve the model's knowledge of different languages. We could also confirm this finding by observing the overall cross-lingual consistency result in Figure 12, as consistency 406 scores jump over the baseline model's cross-lingual consistency. When we look closer at the object-407 subject attention scores of the aligned model in Figure 11, there is a slight shift of subject-object 408 relation features extraction from earlier layers (in the baseline model) into middle-last layers (in the 409 aligned model). In line with the vocabulary expansion, the responsibility shift on feature extraction 410 from the earlier layer into later layers might justify the effectiveness of both approaches on dissim-411 ilar languages. Interestingly, the attention shift is not as strong as the one caused by vocabulary 412 expansion, which is quite counterintuitive as cross-lingual supervision outperforms vocabulary ex-413 pansion in improving cross-lingual consistency. Hence, we leave this interesting finding analysis 414 for future work. In addition, word alignments improve consistency for transliterations or similar 415 orthographical forms, contributing to model's robustness against orthographic variations and nonstandard spellings, but vocabulary expansion can not offer such gains.(cf. §A.4.4) 416

417 418

419

387

388

389 390 391

392 393 394

396 397

399

400

401 402

4.4 THE EFFECT OF CODE-SWITCHING TRAINING TO THE CROSS-LINGUAL CONSISTENCY

Inspired by the experiment on cross-lingual supervision, we 420 further hypothesize that code-switching training, which sub-421 stitutes an entity with alternatives from other languages for 422 intra-sentential alignments in cross-lingual settings, can help 423 the model understand common knowledge across languages for 424 cross-lingual consistency to some extent. To evaluate this hy-425 pothesis, we study xlm-r and xlm-r-cs (Whitehouse et al., 2022), 426 where xlm-r-cs is continuously trained on code-switching cor-427 pus from xlm-r-base and shows high performance in multilin-428 gual fact-checking. From Figure 13, we observe a shift in the 429 consistency bottleneck from the middle layers to the later layers of xlm-r-cs, where the consistency gap between dissimilar and 430 similar languages narrows in xlm-r-cs compared to xlm-r in the 431 middle layers. When observing the attention in Figure 14, we



Figure 13: Effects of codeswitching training on the layerwise consistency. (*cf.* 4.3)



Figure 14: Effects of code-switching training on subject–object attention in xlm-r-base. A similar pattern is observed as the attention difference caused by the vocabulary expansion and the cross-lingual supervision on the baseline input. left: En–Ar, mid: En–Ta, right: baseline. (*cf.* 4.3)



Figure 16: Effects of multi-task fine-tuning on subject–object attention in mt5. We observe a significant head shift but no cross-lingual consistency gains. This is a distinguishable findings from other experiments, where layer behaviors change to mitigate the consistency bottleneck. left: En–Ar, mid: En–Ta, right: baseline. (*cf.* 4.3)

can see one similar layer shift pattern with experiments on the vocabulary expansion and the crosslingual supervision where the attention weights are suppressed on earlier layers. However, unlike these two approaches, there is no amplification of the attention weights on later layers thus we posit that the key of the improvement probably lies on reducing the responsibility on earlier layers. Therefore, code-switching can offer significant gains to the cross-lingual consistency, even without additional objectives. Overall, this finding is consistent with previous experiments.

441

442

443 444

452 453

454

455

456

457 458 459

460

461

462

463

4.5 THE EFFECT OF MULTI-TASK FINE-TUNING TO THE CROSS-LINGUAL CONSISTENCY

468 In previous discussions, we discussed the cross-lingual consis-469 tency in the mt0 family, which is multi-task fine-tuned from 470 the mt5 family (Xue et al., 2020). We hypothesize that this 471 fine-tuning can improve the cross-lingual consistency due to improved cross-lingual generalization across similar tasks in dif-472 ferent languages, as opposed to word-level alignments discussed 473 in previous sections. Surprisingly, multi-task fine-tuning can not 474 offer significant gains to the cross-lingual consistency. As pre-475 sented in Figure 15, the consistency patterns are quite similar 476 across mt0 and mt5. Instead of shifting layers for the attention, 477 which is observed in other experiments, multi-task fine-tuning 478 causes a head shift in Figure 16. We suspect that the model ad-479 justs some neurons at each layer to maintain knowledge but such



Figure 15: Effects of multitask fine-tuning on the layerwise consistency. (cf. 4.3)

adjustment has limited contributions to cross-lingual consistency. While this finding is distinguishable from most of other methods, it is consistent with Figure 4, where we found scaling is not
promising to improve cross-lingual consistency. Specifically, both of them encourage some neurons
to preserve the cross-lingual knowledge consistency, showing limited effectiveness. This finding
aligns with Ortu et al. (2024); Jin et al. (2024) who reported that LLM has attention heads with
contrasting roles in which some of them consider retrieving internal knowledge of language models
and other heads prefer to get the in-context information.

486 5 RELATED WORK

487 488

Petroni et al. (2019) found the availability of relational knowledge within the pre-trained language 489 model by evaluating the language models on the cloze task dataset they proposed, namely LAMA. 490 Kassner et al. (2021) introduced the multilingual version of it, mLAMA, and discovered that the 491 language's relational knowledge capability varies in different languages and other works also found 492 similar findings (Schott et al., 2023; Zhao et al., 2024a). Nevertheless, Zhao et al. (2024a) showed 493 that multilingual language models exhibited limited cross-lingual knowledge recall capability on 494 low-resource languages. Following this line, Fierro & Søgaard (2022); Qi et al. (2023) studied the 495 final predictions in different languages and reported inconsistencies across languages. Moreover, Jin et al. (2024) proposed a method to mitigate such conflicting mechanisms by nullifying heads 496 having significant impact in either of both roles. We take a different angle from those works where 497 we evaluate the cross-lingual knowledge consistency against references in different languages by 498 creating coreferential statements in cross-lingual settings. 499

500 Bhattacharya & Bojar (2023); Kojima et al. (2024); Zhao et al. (2024b) discovered the languagesensitive neurons of decoder in the early and last layers while a considerable portion of language-501 agnostic ones in the middle layers encode universal concepts and utilize the latent language (in 502 this case English) (Wendler et al., 2024; Dumas et al., 2024). Tan et al. (2024) observed encoder-503 decoder language models that neurons tend to be more language-agnostic in the later layers of the 504 encoder part while language-specific in the later layers of the decoder part. Zhao et al. (2024b); 505 Wang et al. (2024b); Zhang et al. (2024) further showed the cross-lingual downstream performance 506 is potentially proportional to the amount of language-agnostic neurons. Ferrando & Costa-jussà 507 (2024) discovered a shared circuit or subnetwork that is responsible for subject-verb agreement task 508 for English & Spanish and Stanczak et al. (2022); Wang et al. (2024a) found that morpho-syntax 509 attributes have noticeable neuron overlapping degree over notable amount of language pairs. We 510 push this line further to trace consistent information and knowledge throughout the layers in cross-511 lingual settings, attempting to understand and interpret how commonly used strategies to improve 512 multilingual models for downstream tasks could impact the cross-lingual knowledge consistency.

513 514

6 CONCLUSION

515 516

517 Do multilingual language models demonstrate cross-lingual consistency? Is it worthwhile to op-518 timise for cross-lingual knowledge consistency? We find the answer to both of these questions is 519 'yes', but with the caveat that performance is tied to language characteristics. In our work, we code 520 mix source monolingual sentences containing a coreferential named entity to control and analyze 521 cross-lingual knowledge consistency.

522 Our analysis reveals that knowledge consistency is heavily dependent on language-specific informa-523 tion such as geography, language family, and writing script. Our layer-by-layer analysis of multi-524 lingual models discovers a consistency bottleneck in the middle layers of models. This bottleneck 525 can be alleviated by expanding the vocabulary, injecting cross-lingual supervision and in training, or including code-swithcing corpus. Our work highlights promising directions in post-calibration, vo-526 cabulary formation, pretraining with cross-lingual objectives, and code-switching training to achieve 527 knowledge consistency across languages, which will better preserve parity of language model per-528 formance. As our experiment discovers that pretraining objective and code-switching training cause 529 most significant positive impact on the cross-lingual consistency, we encourage researchers to em-530 phasize more on representation learning approaches to make the language models more consistent 531 across different languages. 532

533 534

535

References

 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of monolingual representations. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4623–4637, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/ 2020.acl-main.421. URL https://aclanthology.org/2020.acl-main.421.

- 540 Himanshu Beniwal, Kowsik D, and Mayank Singh. Cross-lingual editing in multilingual lan-541 guage models. In Yvette Graham and Matthew Purver (eds.), Findings of the Association 542 for Computational Linguistics: EACL 2024, pp. 2078-2128, St. Julian's, Malta, March 2024. 543 Association for Computational Linguistics. URL https://aclanthology.org/2024. 544 findings-eacl.140.
- Sunit Bhattacharya and Ondrej Bojar. Unveiling multilinguality in transformer models: Exploring 546 language specificity in feed-forward networks. arXiv preprint arXiv:2310.15552, 2023. 547
- 548 Zewen Chi, Li Dong, Bo Zheng, Shaohan Huang, Xian-Ling Mao, Heyan Huang, and Furu Wei. Improving pretrained cross-lingual language models via self-labeled word alignment. arXiv preprint 549 arXiv:2106.06381, 2021. 550
- 551 Alexandra Chronopoulou, Dario Stojanovski, and Alexander Fraser. Language-family adapters for 552 low-resource multilingual neural machine translation. In Atul Kr. Ojha, Chao-hong Liu, Ekate-553 rina Vylomova, Flammie Pirinen, Jade Abbott, Jonathan Washington, Nathaniel Oco, Valentin 554 Malykh, Varvara Logacheva, and Xiaobing Zhao (eds.), Proceedings of the Sixth Workshop on 555 Technologies for Machine Translation of Low-Resource Languages (LoResMT 2023), pp. 59–72, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics. doi: 10.18653/v1/ 556 2023.loresmt-1.5. URL https://aclanthology.org/2023.loresmt-1.5.
- 558 Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What does BERT look 559 at? an analysis of BERT's attention. In Tal Linzen, Grzegorz Chrupała, Yonatan Belinkov, and Dieuwke Hupkes (eds.), Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and 561 Interpreting Neural Networks for NLP, pp. 276–286, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-4828. URL https://aclanthology. 563 org/W19-4828.
- 564 Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, 565 Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Un-566 supervised Cross-lingual Representation Learning at Scale. In Proceedings of the 58th Annual 567 Meeting of the Association for Computational Linguistics. Association for Computational Lin-568 guistics, 2020. URL https://www.aclweb.org/anthology/2020.acl-main.747. 569
- Wietse de Vries, Martijn Wieling, and Malvina Nissim. Make the best of cross-lingual transfer: 570 Evidence from POS tagging with over 100 languages. In Smaranda Muresan, Preslav Nakov, 571 and Aline Villavicencio (eds.), Proceedings of the 60th Annual Meeting of the Association for 572 Computational Linguistics (Volume 1: Long Papers), pp. 7676–7685, Dublin, Ireland, May 2022. 573 Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.529. URL https: 574 //aclanthology.org/2022.acl-long.529. 575
- Clément Dumas, Veniamin Veselovsky, Giovanni Monea, Robert West, and Chris Wendler. How do llamas process multilingual text? a latent exploration through activation patching. In ICML 2024 Workshop on Mechanistic Interpretability, 2024. 578
- 579 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, 580 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition. arXiv preprint arXiv:2209.10652, 2022. 581 582
 - Javier Ferrando and Marta R Costa-jussà. On the similarity of circuits across languages: a case study on the subject-verb agreement task. arXiv preprint arXiv:2410.06496, 2024.
- 585 Constanza Fierro and Anders Søgaard. Factual consistency of multilingual pretrained language 586 models. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), Findings of the Association for Computational Linguistics: ACL 2022, pp. 3046–3052, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.240. URL 588 https://aclanthology.org/2022.findings-acl.240. 589
- Gottlob Frege. On sense and reference, 1892. 591

577

583

584

Zihao Fu, Wai Lam, Qian Yu, Anthony Man-Cho So, Shengding Hu, Zhiyuan Liu, and Nigel Collier. 592 Decoder-only or encoder-decoder? interpreting language model as a regularized encoder-decoder. arXiv preprint arXiv:2304.04052, 2023.

- Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. Causal abstractions of neural networks. *Advances in Neural Information Processing Systems*, 34:9574–9586, 2021.
- Ashim Gupta and Vivek Srikumar. X-fact: A new benchmark dataset for multilingual fact checking. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pp. 675–682, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-short.86.
 URL https://aclanthology.org/2021.acl-short.86.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson.
 Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generaliza tion. In *37th International Conference on Machine Learning, ICML 2020*, volume PartF16814,
 pp. 4361–4371, 2020. ISBN 9781713821120. URL https://sites.
- Kuming Hu, Junzhe Chen, Xiaochuan Li, Yufei Guo, Lijie Wen, Philip S. Yu, and Zhijiang Guo. Towards understanding factual knowledge of large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum? id=90evMUdods.
- Dieuwke Hupkes, Mario Giulianelli, Verna Dankers, Mikel Artetxe, Yanai Elazar, Tiago Pimentel, Christos Christodoulopoulos, Karim Lasri, Naomi Saphra, Arabella Sinclair, Dennis Ulmer, Florian Schottmann, Khuyagbaatar Batsuren, Kaiser Sun, Koustuv Sinha, Leila Khalatbari, Maria Ryskina, Rita Frieske, Ryan Cotterell, and Zhijing Jin. A taxonomy and review of generalization research in nlp. *Nature Machine Intelligence*, 5:1161–1174, 10 2023. doi: 10.1038/s42256-023-00729-y.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. ACM
 Computing Surveys, 55(12):1–38, 2023.
- ⁶²¹ Zhuoran Jin, Pengfei Cao, Hongbang Yuan, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and Jun Zhao. Cutting off the head ends the conflict: A mechanism for interpreting and mitigating knowledge conflicts in language models. *arXiv preprint arXiv:2402.18154*, 2024.
- Karthikeyan K, Zihan Wang, Stephen Mayhew, and Dan Roth. Cross-lingual ability of multilingual
 bert: An empirical study. In *International Conference on Learning Representations*, 2020. URL
 https://openreview.net/forum?id=HJeT3yrtDr.
- Katharina Kann, Ryan Cotterell, and Hinrich Schütze. One-shot neural cross-lingual transfer for paradigm completion. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1993–2003, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10. 18653/v1/P17-1182. URL https://aclanthology.org/P17-1182.
- Nora Kassner, Philipp Dufter, and Hinrich Schütze. Multilingual LAMA: Investigating knowledge in multilingual pretrained language models. In Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty (eds.), Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pp. 3250–3258, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.284. URL https: //aclanthology.org/2021.eacl-main.284.
- Takeshi Kojima, Itsuki Okimura, Yusuke Iwasawa, Hitomi Yanaka, and Yutaka Matsuo. On the multilingual ability of decoder-based pre-trained language models: Finding and controlling language-specific neurons. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 6919–6971, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. naacl-long.384. URL https://aclanthology.org/2024.naacl-long.384.
- 647 Wayne B. Krause. English & indonesian similarities & differences, n.d. URL https:// indodic.com/SimilaritiesDiffs.htm. Accessed: 2024-10-01.

667

- Anna Langedijk, Hosein Mohebbi, Gabriele Sarti, Willem Zuidema, and Jaap Jumelet. Decoderlens: Layerwise interpretation of encoder-decoder transformers. arXiv preprint arXiv:2310.03686, 2023.
- Yoav Levine, Barak Lenz, Opher Lieber, Omri Abend, Kevin Leyton-Brown, Moshe Tennenholtz, and Yoav Shoham. {PMI}-masking: Principled masking of correlated spans. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum? id=3Aoft6NWFej.
- Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. Unveiling the pitfalls of knowledge editing for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=fNktD3ib16.
- Davis Liang, Hila Gonen, Yuning Mao, Rui Hou, Naman Goyal, Marjan Ghazvininejad, Luke
 Zettlemoyer, and Madian Khabsa. Xlm-v: Overcoming the vocabulary bottleneck in multilin gual masked language models. *arXiv preprint arXiv:2301.10472*, 2023.
- Yan Liu, Yu Liu, Xiaokang Chen, Pin-Yu Chen, Daoguang Zan, Min-Yen Kan, and Tsung-Yi Ho.
 The devil is in the neurons: Interpreting and mitigating social biases in language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven 668 Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, 669 Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert 670 Webson, Edward Raff, and Colin Raffel. Crosslingual generalization through multitask fine-671 tuning. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), Proceedings of the 672 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Pa-673 pers), pp. 15991–16111, Toronto, Canada, July 2023. Association for Computational Linguis-674 tics. doi: 10.18653/v1/2023.acl-long.891. URL https://aclanthology.org/2023. 675 acl-long.891.
- nostalgebraist. Interpreting gpt: The logit lens, 2019. URL https://www.lesswrong.com/
 posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens. Accessed:
 2024-08-04.
- Francesco Ortu, Zhijing Jin, Diego Doimo, Mrinmaya Sachan, Alberto Cazzaniga, and Bernhard
 Schölkopf. Competition of mechanisms: Tracing how language models handle facts and counter *arXiv preprint arXiv:2402.11655*, 2024.
- Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*, 2019.
- Adithya Pratapa, Gayatri Bhat, Monojit Choudhury, Sunayana Sitaram, Sandipan Dandapat, and Kalika Bali. Language modeling for code-mixing: The role of linguistic theory based synthetic data. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1543–1553, 2018.
- Jirui Qi, Raquel Fernández, and Arianna Bisazza. Cross-lingual consistency of factual knowledge
 in multilingual language models. *arXiv preprint arXiv:2310.10378*, 2023.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Tim Schott, Daniel Furman, and Shreshta Bhat. Polyglot or not? measuring multilingual ency clopedic knowledge in foundation models. In Houda Bouamor, Juan Pino, and Kalika Bali
 (eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Pro cessing, pp. 11238–11253, Singapore, December 2023. Association for Computational Linguis tics. doi: 10.18653/v1/2023.emnlp-main.691. URL https://aclanthology.org/2023.
 emnlp-main.691.

702 Karolina Stanczak, Edoardo Ponti, Lucas Torroba Hennigen, Ryan Cotterell, and Isabelle Augen-703 stein. Same neurons, different languages: Probing morphosyntax in multilingual pre-trained 704 models. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.), 705 Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 1589-1598, Seattle, United States, 706 July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.114. 707 URL https://aclanthology.org/2022.naacl-main.114. 708 709 710 Shaomu Tan, Di Wu, and Christof Monz. Neuron specialization: Leveraging intrinsic task modularity for multilingual machine translation. arXiv preprint arXiv:2404.11201, 2024. 711 712 Yulia Tsvetkov and Chris Dyer. Cross-lingual bridges with models of lexical borrowing. Journal of 713 Artificial Intelligence Research, 55:63-93, 2016. 714 715 Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Simas Sakenis, Jason 716 Huang, Yaron Singer, and Stuart Shieber. Causal mediation analysis for interpreting neural nlp: 717 The case of gender bias. arXiv preprint arXiv:2004.12265, 2020. 718 719 Hetong Wang, Pasquale Minervini, and Edoardo M Ponti. Probing the emergence of cross-lingual 720 alignment during llm training. arXiv preprint arXiv:2406.13229, 2024a. 721 722 Weixuan Wang, Barry Haddow, Minghao Wu, Wei Peng, and Alexandra Birch. Sharing matters: 723 Analysing neurons across languages and tasks in llms. arXiv preprint arXiv:2406.09265, 2024b. 724 725 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha 726 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language 727 models. In The Eleventh International Conference on Learning Representations, 2023. URL 728 https://openreview.net/forum?id=1PL1NIMMrw. 729 730 Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. Do llamas work in en-731 glish? on the latent language of multilingual transformers. arXiv preprint arXiv:2402.10588, 732 2024. 733 734 Chenxi Whitehouse, Fenia Christopoulou, and Ignacio Iacobacci. EntityCS: Improving zeroshot cross-lingual transfer with entity-centric code switching. In Yoav Goldberg, Zornitsa 735 Kozareva, and Yue Zhang (eds.), Findings of the Association for Computational Linguistics: 736 EMNLP 2022, pp. 6698–6714, Abu Dhabi, United Arab Emirates, December 2022. Associa-737 tion for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.499. URL https: 738 //aclanthology.org/2022.findings-emnlp.499. 739 740 Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of 741 large language models. arXiv preprint arXiv:2401.11817, 2024. 742 743 Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya 744 Barua, and Colin Raffel. mt5: A massively multilingual pre-trained text-to-text transformer. arXiv 745 preprint arXiv:2010.11934, 2020. 746 747 Zhihao Zhang, Jun Zhao, Qi Zhang, Tao Gui, and Xuanjing Huang. Unveiling linguistic regions in 748 large language models. arXiv preprint arXiv:2402.14700, 2024. 749 750 Xin Zhao, Naoki Yoshinaga, and Daisuke Oba. Tracing the roots of facts in multilingual language 751 models: Independent, shared, and transferred knowledge. arXiv preprint arXiv:2403.05189, 752 2024a. 753 754 Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji Kawaguchi, and Lidong Bing. How do large 755 language models handle multilingualism? arXiv preprint arXiv:2402.18815, 2024b.

756 LIMITATIONS

We only cover the transformer encoder and encoder-language models for this work. Another promis-ing avenue for this work is evaluating cross-lingual knowledge consistency on other language mod-els. Moreover, we only analyze each crucial component independently due to the time constraint and left scrutinizing the interaction between each component for future work. In the future, we may expand this work by analyzing how the interaction among these components could affect the crosslingual consistency of multilingual models. Another thing is that our causal intervention method needs to be done manually, and we suspect that this method could produce a side effect on the model because the representations encoded by language models are more likely to be polysemous. Additionally, we only evaluate the language models in context-independent settings. Thus, in the future, we plan to evaluate the consistency of the models' knowledge and observe whether language models utilize their parametric knowledge more or emphasize the knowledge from the given context under the cross-lingual setting. Another thing to consider is that we only evaluate our solution using one particular model due to the time constraint. Also, we do not explore various pre-training ob-jectives and evaluate solutions to the encoder-decoder model; hence, we leave such things as future work. One interesting thing to explore in this aspect is to see whether adversarial training could help to enhance cross-lingual consistency. Another thing that we want to consider is that we use an assumption that one reference is represented as a single English object entity to make the evaluation tractable; hence, we do not take into account the real-world setting where one reference can be interpreted in different ways on multiple languages (e.g., "China" is written as "ZhongGuo" in Chinese rather than "China"). Lastly, our research scope assumes that the knowledge we want to evaluate is factual and not dependent on subjective aspects (e.g., cultural context). With that assumption, we assume that references here generally have one-to-one mapping to representation in one language where the representation here is considered common knowledge.

ETHICS STATEMENT

This work aims to evaluate the consistency of the language model across different senses (particularly between a monolingual input and its code-mixed counterparts) and the impact of different factors on that metric. Doing such a study could shed light on the limitations of language models and think of the mitigations of such matters.

Reproducibility Statements

We used open-source pretrained models and also dataset for all of the reported experiments thus no undisclosed assets utilized in our work. Additionally, we also provide necessary experiments' output and codes on https://anonymous.4open.science/r/knowledgeConsistencyAndConflict-4827.

	xlm-r input	mt0 input
I_{mono} I_{cm} $I_{(S^{l1} \cap O)}$	Paris is the capital of < mask > باریس is the capital of < mask > <mask> is the capital of <mask></mask>	Paris is the capital of < extra_id_0 > باریس is the capital of < extra_id_0 > <extra_id_0> is the capital of <extra_id_1></extra_id_0>

Table 2: Input sample for the evaluation task. We only predict the object in bold. $I_{\setminus (S^{l_1} \cap O)}$ is the baseline input.

A APPENDIX

A.1 INPUT FORMAT

In our task definition, we introduce our evaluation task in both intuition and math perspective. Here is the input sample in Table 2. Meanwhile, as presented in the task definition, we do not consider whether predictions are true but focus on the same prediction distributions regardless of languages. Note that we did not perturb the surface structure in order to minimize variables to affect factual knowledge recall because S^{l2} "switches-in" at grammatically correct point as the new subject (Pratapa et al., 2018).

A.2 METRIC FUNCTION AND INTERPRETABILITY APPROACH

RankC RankC (Qi et al., 2023) is used to evaluate the cross-lingual knowledge consistency. Given a set of statements S where each of the statement having each own I_{mono} and I_{cm} , the number of candidates $Cand(O_{\in V}|I_{mono})$ of i-th statement N_i , $mono^j$ stands for the j-th candidate of $Cand(O_{\in V}|I_{mono})$, cm^j stands for the j-th candidate of $Cand(O_{\in V}|I_{cm})$, and the RankC score of $Cand(O_{\in V}|I_{mono})$ concerning $Cand(O_{\in V}|I_{cm})$ can be written as

$$RankC(cm, mono) = \frac{\sum_{i=1}^{|S|} \sum_{j=1}^{N_i} \frac{e^{N_i - j}}{\sum_{k=1}^{N_i} e^{k - j}} * P@j}{|I_{mono}|},$$
(1)

$$P@j = \frac{1}{j} | \{ cm_i^1, cm_i^2, \cdots, cm_i^j \} \cap \{ mono_i^1, mono_i^2, \cdots, mono_i^j \} |.$$
(2)

Top@1 Accuracy The Top@1 accuracy is defined as the average number of exact matches between the top-1 predictions given I_{mono} and I_{cm} .

Subject-Object Attention Let $A_{a,b}^{(k)}$ be the attention score between *a*-th token and *b*-th token in a 881 statement *k*, O_k is the set of indices of the masked tokens in *k*, S_k is the set of indices of subject tokens in *k*, and *K* is a set of statements, the average attention weight of head *l* in *i* layer can be 883 defined as

$$Attn(h^{(l,i)}) = \frac{\sum_{k \in K} \frac{\sum_{o \in O_k} \sum_{s \in S_k} A_{o,s}^{(k)}}{|O_k|}}{|K|}$$
(3)

 IG^2 Score If $w_j^{(l)}$ is the activation value of *j*-th neuron in the *l*-th layer of a particular input (either code-mixed or not), *m* is the approximation step, and *t* as a token of the whole ground truth object entity, the score for a given I_{mono} or I_{cm} is defined as

$$IG^{2}(w_{j}^{(l)}) = \sum_{t \in T} \frac{\frac{w_{j}^{(l)}}{m} \sum_{k=1}^{m} \frac{\partial P(t|\frac{k}{m}w_{j}^{(l)})}{\partial(\frac{k}{m}w_{j}^{(l)})}}{|T|}$$
(4)

918 A.3 FINDINGS IN DETAILS

A.3.1 LAYER-WISE CONSISTENCY





RankC across Encoder Layers in mt0-base





Accuracy across Encoder Layers in mt0-base



Figure 17: mT0 layer-wise cross-lingual consistency scores (left: RankC, right: Top@1)



Figure 18: xlm-r layer-wise crosslingual consistency scores (left: RankC, right: Top@1)

1026 A.3.2 OVERALL CONSISTENCY



Figure 19: Cross-lingual consistency scores across languages of mt0 (top: RankC, bottom: Top@1
 Accuracy). Note: The dashed line here is the average corresponding consistency scores of mt0-base
 across languages





Figure 21: Overall cross-lingual consistency in mt0 grouped by 3 factors (left: geographics, middle: language family, right: writing scripts.). Metrics legend: top: RankC, bottom: Top@1 Accuracy. Models legend: red: mt0-large, blue: mt0-base



Figure 22: Overall cross-lingual consistency in xlm-r grouped by 3 factors (left: geographics, middle: language family, right: writing scripts). Metrics legend: top: RankC, bottom: Top@1 Accuracy.
Models legend: red: xlm-r-large, blue: xlm-r-base



Figure 23: Subject–Object attention difference with I_{mono} to I_{cm} in mT0 for some code-mixed languages. Models legend: left: mt0-base, right: mt0-large. Languages legend: from top to bottom, en–de, en–ar, en–id, and en–ta.



Figure 24: Subject–Object attention difference with I_{mono} to I_{cm} in xlm-r for some code-mixed languages (From left to right, base model and large model. from top to bottom, en–de, en–ar, en–id, and en–ta).





1350 A.3.4 FEED-FORWARD NEURONS' GRADIENTS SUM

Figure 26: IG^2 scores in mt0 for en–de, en–ta, en–id, and en–ar. Models legend: upper two rows: mt0-base, lower two rows: mt0-large.





1458		M. 1.1	C. L							
1459		Model	Codemixing Language	α	Patched FFN Layers					
1460		mt() hasa	en-ta	0.7	[0,3,10,11]					
1461		Into-base	en–ar	0.7	[0,1,9,10]					
1462		mt() large	en-ta	0.8	[0,1,19,20,21]					
1/63		into-rarge	en–ar	0.8	[0,1,19,20,21]					
1464		vlmr basa	en-ta	0.7	[5,8,9,10]					
1404		XIIIII-Dase	en–ar	0.7	[5,7,8,10]					
1400		vlmr large	en-ta	0.7	[0,2,5,19,20]					
1466		xiiii-iaige	en–ar	0.8	[17,18,19,20,21]					
1467										
1468		Table	e 3: Causal Intervention H	yperpa	arameters Setup					
1469										
1470										
1471	A.4 IMPROV	ING CONSIST	TENCY							
1473	A.4.1 ADDIN	NG MONOLIN	IGUAL BIAS							
1474	This experimen	it aims to me	asure whether each pattern	n has a	a causal relationship w	ith cross-lingual				
1475	consistency.									
1470		•	······			· · ·				
14//	• Attention score suppression: Using the definition from A.2 and define a suppression con-									
1478	stant $\alpha, \alpha \in [0, 1)$, the patched attention weight of every object-subject relation will be									
1479	$A_{a,b} =$	$= \alpha A_{a,b}.$								
1480	• Feed-f	forward neuro	on activation patching (Vi	g et al	., 2020; Geiger et al.,	2021): consider				
1481	$a_i^{(l,p)}$	as the activat	ion of <i>i</i> -th token on I_{mo}	no pro	duced by p-th neuron	in <i>l</i> -th encoder				
1482	layer's	feed-forwar	d network, then patched a	activat	ion value for the <i>i</i> -th	token on I_{cm} is				
1483	$\bar{a}_i^{(l,p)}$	$=a_{i}^{(l,p)}$, in w	hich we apply this for eve	ry mas	sk token.					
1484	• Hybrid	i. We apply a	ttention weight suppressio	n and	feed forward neuron a	ctivations natch				
1485	ing sir	nultaneously	uchtion weight suppressio	in and	iccu-ioi waru iicuroii a	cuvations paten-				
1486	ing sir	indituneousry.								
1487	For the hyperpa	arameters use	d in the causal intervention	on exp	eriment, we set the α	value that is not				
1488	too big so that	did not signi	ficantly diminish the atter	tion w	veight yet making the	attention weight				
1489	distribution for	I_{cm} closer to	that weight distribution for	or I _{mor}	no. While for FFN-laye	ers, we intervene				
1490	4 different enco	der layers for	base models and 5 differe	nt enco	oder layers that have la	nguage-sensitive				
1491	neurons based	on IG^2 (i.e.	layer which has noticeable	$e IG^2$	distribution difference	e between I _{mono}				
1492	and $I_c m$). Read	lers can refer	to table 3 to see the hyper	param	eters used in this expe	riment.				
1493										
1494										
1495										
1496										
1497										
1498										
1499										
1500										
1501										
1502										
1503										
1504										
1505										
1506										
1507										
1507										
1500										
1509										
1510										
1511										



Figure 28: Intervention scores in mt0. Metrics legend: left: RankC, right: Top@1 Accuracy. Model
 legend: upper two rows: mt0-base, lower two rows: mt0-large



1620 A.4.2 IMPACT OF LARGER VOCABULARY



Figure 30: Layer-wise cross-lingual knowledge consistency of xlm-v vs xlm-r-base





Figure 31: Effects of vocabulary expansion to overall cross-lingual consistency (top: RankC, bot-tom: Top@1 Accuracy). Note: The dashed line here is the average corresponding consistency scores of xlm-r-base across languages





Figure 34: Regression analysis between parity ratio and RankC improvement offered by xlm-v to xlm-r. Spearman $\rho = 0.06$. We define parity ratio as the token length ratio between tokenized subjects for xlm-v-base and xlm-r-base. Our analysis discovers that many languages have a token parity ratio average within 0.8-1, which means that many of the subject entities are known on both tokenizers of the models.

1836 A.4.3 THE EFFECT OF PRE-TRAINING OBJECTIVE





1944 A.4.4 CASE STUDY FOR TRANSLITERATION

Instead of using translations, we transliterate bn^7 and ar^8 to understand the impact of writing sys-1946 tems, particularly transliterations. As presented in Figure 38, word alignments (or the similar effect 1947 from CS training) contribute to the model's cross-lingual consistency against writing systems be-1948 cause xlm-align and xlm-rcs show similar performance in both original and transliteration settings. 1949 Meanwhile, we can observe that xlm-align and xlm-r-cs significantly improve the overall perfor-1950 mance for non-Lattin scripts in §A.4.3. This is reasonable as word alignments or CS training help 1951 the model link original words with their translations or transliterations, depending on the training 1952 corpus, thereby enhancing cross-lingual consistency. We suspect that these word alignments might 1953 also improve robustness for handling non-standard spellings and orthographic variations. However, 1954 xlm-v-base and xlm-r-base without word alignment benefit from transliterations, which means that 1955 xlm-v-base and xlm-r-base do not sufficiently align original words with their transliterations to main cross-lingual consistency. It is also confirmed by the overall performance of vocabulary expansions 1957 in §A.4.2, where vocabulary expansions can not offer significant gains for cross-lingual consistency. Overall, the evaluation task does not inadequately boost consistency for languages using Latin script 1958 because word alignments resulting in cross-lingual consistency are the main factor. 1959

