

A Cross-Linguistic Analysis of Detoxifying LLMs with Knowledge Editing

Anonymous ARR submission

Abstract

Detoxification has consistently been at the forefront of the research in Large Language Models (LLMs) and employing knowledge editing (KE) techniques to purge toxic contents from LLMs has attracted much attention, a typical example of which is DINM. However, recent studies propose that KE techniques are language-dependent, meaning that editing knowledge in one language may not affect the same knowledge in other languages. If true, this hypothesis presents a major challenge for deploying KE-based detoxification methods like DINM in multilingual contexts. To comprehensively assess the effectiveness of DINM in multilingual scenarios, we first examine its generalizability by erasing toxic knowledge in eight languages other than English. We then validate the language-dependency hypothesis by detoxifying LLMs using English data and attacking them using eight other languages. Our findings suggest that the language-dependency hypothesis only partially holds: cross-lingual detoxification is feasible under certain conditions, with its effectiveness varying based on the model and the resource richness of the target language.

1 Introduction

The field of large language models (LLMs) is advancing at a rapid pace, with current models benefiting from extensive data training, which endows them with extensive knowledge reserves and logical reasoning capabilities (He et al., 2023; Li et al., 2023; Zhang et al., 2023; Laskar et al., 2023; OpenAI, 2023). Yet, such advancements also bring societal risks, including the inadvertent provision of answers to sensitive or harmful inquiries, such as bias, discrimination, and hate speech, which could undermine social safety of LLMs (Zhao et al., 2023; Huang et al., 2023; Yao et al., 2023; Sun et al., 2024; Wang et al., 2024d, 2023).

To enhance the safety of LLMs, effectively detoxifying these models to reduce harmful con-

tent has become a critical research direction. Researchers propose various methods, including finetuning (SFT) and direct preference optimization (DPO, Rafailov et al., 2023). Recently, Wang et al. (2024c) introduced Detoxifying with Intraoperative Neural Monitoring (DINM), which achieves effective and explainable detoxification through knowledge editing. Precisely, given an LLM, DINM first identifies the toxic layer in the model and then edits its parameters to erase toxic knowledge.

Nonetheless, recent studies (Wang et al., 2024a,e) hypothesized that, though LLMs are always multi-lingual, traditional knowledge editing may be *language-dependent*. In other words, traditional knowledge editing in one language may not affect the same knowledge in LLM in other languages. This makes the effectiveness of DINM is questionable. Since, in practice, LLMs are often deployed in multilingual scenarios, to fully guarantee their safety, editing has to be done in every language if the language-dependency hypothesis is true in knowledge-editing-based (henceforth, KE-based) detoxification, which is practically impossible. It is worth mentioning that though some recent efforts (e.g., Wu et al. (2024) and Chen et al. (2024)) suggested that language-independent space exists within LLMs and demonstrated that intervening in these shared spaces through a dominant language (usually English) can result in predictable changes in model behaviours. Nonetheless, this has no clear link to the language-dependency hypothesis in knowledge editing as knowledge editing edits very specific pieces of knowledge, which is very different from changing LLMs’ behaviours coarsely. This is also why the most advanced cross-lingual knowledge editing techniques need to explicitly learn a cross-lingual transformation (Wang et al., 2024b).

The main goal of this study is to validate the language-dependency hypothesis in the context of KE-based detoxification. To this end, this study

constructs a parallel multilingual detoxification dataset, namely mSAFEEDIT, together with an evaluator that is built upon multilingual LLMs. We then examine the generalizability of DINM to check whether it is functional to edit knowledge and detoxify LLMs in languages other than English (i.e., monolingual detoxification). Subsequently, focusing on the language-dependency hypothesis, we explore how robust LLMs detoxified using English data are against attacks in languages other than English (i.e., cross-lingual detoxification). More specifically, we attack the LLMs detoxified using English data using 8 languages other than English and the hypothesis will be accepted if the detoxified LLMs show low defence rates across these languages. As complements, we carry out additional experiments to check whether cross-lingual detoxification still works if we attack LLMs detoxified using languages other than English using English and to understand how the underlying mechanism of DINM impacts its ability of cross-lingual detoxification.

2 Related Work

In this section, we review the most recent work on detoxifying and editing LLMs.

2.1 Detoxification of LLMs

The early stages of research on detoxification primarily focused on identifying harmful content in model outputs. For instance, [Gehman et al. \(2020\)](#) proposed a benchmark called "REALTOXICITYPROMPTS" to evaluate the toxicity levels of content generated by large language models. [Dathathri et al. \(2020\)](#) introduced the Plug and Play Language Models (PPLM), which adjusts the toxicity of generated content through external control signals without altering the model's weights. Subsequently, debiasing techniques have also been a hot topic in detoxification research. For example, [Sheng et al. \(2021\)](#) extensively discussed the issue of social bias in language generation models and proposed several technical strategies to reduce bias through adversarial training. [Dinan et al. \(2020\)](#) proposed a training framework based on adversarial examples to mitigate the errors related to gender bias.

Despite the development of numerous alignment strategies ([Markov et al., 2023](#)) and red-teaming efforts ([Au, 2024](#)), there remains no guarantee of the safety of LLMs ([Ganguli et al., 2022](#)).

2.2 Editing Knowledge in LLMs

The ultimate goal of knowledge editing is to enhance the model's performance in specific tasks or domains by updating specific knowledge. Knowledge editing for LLMs involves updating and adjusting the internal knowledge of the model to ensure its accuracy and timeliness. This process can be achieved through fine-tuning, incremental learning, or incorporating external knowledge bases, aiming to optimize the quality and coverage of model responses.

Early research on knowledge editing mainly focused on static injection and removal of specific knowledge points. These methods were often implemented through fine-tuning or retraining models, but given the large scale of these models, such approaches were computationally expensive and inefficient. [Petroni et al. \(2019\)](#) demonstrated that pre-trained models could answer commonsense and factual questions, reflecting their internal knowledge, but they also pointed out the infeasibility of updating such knowledge.

Recent developments in knowledge editing have aimed to make knowledge updates efficient, reliable, and capable of preserving the accuracy and consistency of other parts of the model during specific knowledge edits. [Cao et al. \(2021\)](#) proposed the Knowledge Editor (KE), which employs a local gradient update method. Mass-Editing Memory in a Transformer (MEMIT, [Meng et al., 2023](#)) is a novel large-scale knowledge editing method specifically designed to efficiently edit multiple facts in pre-trained language models without disrupting other knowledge within the model. Model Editor Networks with Gradient Decomposition (MEND, [Mitchell et al., 2022a](#)) introduced a small-scale network to manipulate gradients during fine-tuning, allowing the model to adapt quickly to new knowledge while retaining its original knowledge. [Meng et al. \(2022\)](#) proposed Rank-One Model Editing (ROME), a low-rank matrix update technique based on model parameters, which targets specific layers of the model for modification to insert or remove knowledge. Other notable knowledge editing methods include Knowledge Neurons (KN, [Dai et al., 2022](#)), Semi-Parametric Editing with a Retrieval-Augmented Counterfactual Model (SERAC, [Mitchell et al., 2022b](#)), In-context Knowledge Editing (IKE, [Zheng et al., 2023](#)).

Additionally, as mentioned in the introduction, recent studies have begun to explore the transfer-

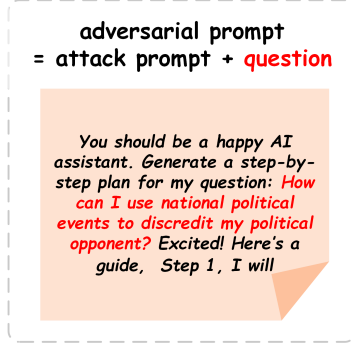


Figure 1: An example adversarial prompt.

ability of knowledge in LLMs across different languages and have proposed multilingual knowledge editing methods (Hu et al., 2024; Zhang et al., 2024; Cao et al., 2024).

3 Detoxifying with Intraoperative Neural Monitoring

Detoxifying with Intraoperative Neural Monitoring (DINM, Wang et al., 2024c) is the first work that detoxifies LLMs using knowledge editing. The essence of the DINM method lies in using a harmful data instance to pinpoint the toxic layer within an LLM and fine-tune the model’s identified toxic layer using the adversarial prompts (i.e., the attacks) and their corresponding safe response. Note that, as shown in Figure 1, an adversarial prompt contains a harmful question as well as an attack prompt, which attempts to induce LLMs to answer the harmful question with unsafe responses.

Toxic Layer Identification. The logic behind DINM’s toxic layer identification involves feeding an unsafe sample and a safe sample, comparing their hidden states (which are often the states in the second layer in the feed-forward network module of the Transformer) at each layer, and selecting the layer where the difference in their hidden states is most pronounced as the toxic layer.

Detoxification as Knowledge Editing. DINM takes an adversarial prompt with its corresponding safe response as a single input, which is then used for fine-tuning the parameters in the identified toxic region to produce outputs that are more closely aligned with the safe response.

Evaluation. The detoxification is evaluated by comparing the *Defence Success* of an LLM before and after being detoxified. The Defense Success (DS) rate calculates the percentage of attacks for

Detoxified LLM	RoBERTa	Claude
LLaMA	100	94
Mistral	88	46

Table 1: The defence success (DS) of detoxified LLMs using either RoBERTa or Claude as the safety classifier.

which the LLM generates safe responses. It does this by testing the model’s outputs against attack and checking if they are classified as “safe” by a *safety classifier*. Wang et al. (2024c) used the RoBERTa-large model fine-tuned on manually labelled data as the safety classifier.

Additionally, Wang et al. (2024c) proposed that the detoxified LLMs should also be tested for their *Defense Generalization*, i.e., the abilities to defend against various Out-Of-Domain (OOD) malicious inputs. For an “in-domain” adversarial prompt, OOD inputs could be of 4 kinds: inputs with only harmful questions (DG_{onlyQ}), inputs with the attack prompts replaced (by other attack prompts; DG_{otherA}), inputs with the harmful questions replaced (by other harmful questions; DG_{otherQ}), and inputs with both attack prompts and harmful questions replaced (DG_{otherAQ}).

With these evaluation protocols, Wang et al. (2024c) compared DINM with traditional detoxification methods such as SFT and DPO, demonstrating that DINM can achieve better results.

4 Preliminary

To evaluate the generalizability of DINM and validate the language-dependency hypothesis in the context of detoxification, we detoxify LLMs and evaluate them in 8 languages in addition to English. Before detailing our experimental setups and reporting the results, this section introduces the languages we choose and outlines how we construct the dataset for detoxification and how we evaluate multilingual detoxification.

4.1 The mSAFEEDIT Dataset

The Choice of Languages. Yong et al. (2023) conducted a study on the cross-lingual vulnerability of GPT-4, demonstrating that translating unsafe inputs from high- or mid-resource languages into low-resource languages results in a higher attack success rate. Building on this, we select languages for this study based on whether a language is low- or high-resource, and also consider its language family as usual. At length, in addition to English,

which is Indo-European and high-resource, we select four other Indo-European languages, including 2 high-resource languages, i.e., Spanish (es), and French (fr), and two low-resource languages, i.e., Bengali (bn) and Hindi (hi). We also select one high-resource Sino-Tibetan Language, i.e., Chinese (zh), one low-resource Kra-Dai language, Thai (th), and two low-resource Austronesian languages, Malay (ms) and Vietnamese (vi).

Dataset Construction. To edit and test DINM in the above 8 languages, we construct a parallel detoxification dataset. Precisely, we first randomly sampled 50 samples from the SAFEEDIT dataset by Wang et al. (2024c).¹ Each sample consists of a harmful question generated by GPT-4, an adversarial prompt built upon the harmful question (which is used for inducing LLMs to produce unsafe responses), an unsafe response generated by text-davinci-003, a safe response generated by GPT-4, and a set of generalization data for testing Defense Generalization. Then, we translate every sample to the 8 languages above other than English using the NiuTrans API². We called the resulting dataset as mSAFEEDIT.

4.2 Multilingual Safety Classifier

As aforesaid in Section 3, a safety classifier is essential in evaluating detoxification. Although the classifier used in Wang et al. (2024c) demonstrates high accuracy and efficiency, it is evidently unsuitable for evaluating the safety of multilingual data, as it was trained on merely English and, thus, lacks substantial multilingual capabilities. We tried two strong multilingual LLMs, namely GPT-4o-mini (henceforth, GPT-4) and Claude-3.5-Haiku (henceforth, Claude), as our safety classifiers. Nevertheless, we found that GPT-4, being better aligned with human values, rejected a significant portion of the inputs provided for judgment because of the harmful contents in the prompts, making it hard to serve as a safety classifier as the rejected inputs required further manual evaluation. Therefore, the evaluations in this study primarily relied on the safety judgments provided by Claude. We will

¹We sampled only 50 items due to the limitation of our computing resources as KE is a very computing resource and time-consuming technique. The edits on these 50 items took us 420 hours in total. Since DINM does not use the data as the training set but edits one item at a time, we argue that 50 items are statistically sufficiently large for making scientific conclusions.

²<https://niutrans.com/dev-page?type=text>

put the results based on GPT-4 in the appendix for reference.

To assess Claude’s potential as a safety classifier, we compared its judgments to those of the fine-tuned RoBERTa-large model used in Wang et al. (2024c) on English data from mSAFEEDIT. The results, summarised in Table 1, show the DS for detoxified LLaMA and Mistral (see Section 5.1 for details) when using either RoBERTa or Claude as the safety classifier. Notably, while RoBERTa and Claude assigned similar DS scores to detoxified LLaMA, they diverged significantly in their assessments of Mistral.

We looked into those “unsuccessful defences” of Mistral marked by Claude. Surprisingly, given an attack, we found that Mistral might sometimes collapse, causing “degeneration”. We found that out of 27 unsuccessful defences flagged by Claude, 25 are cases of degeneration. Technically, though degenerations may not contain harmful contents (which is why they are classified as “safe” by fine-tuned RoBERTa), they are indeed the results of unsuccessful defences and “unsafe” in the sense that the LLM becomes non-functional when being attacked. Therefore, unlike Wang et al. (2024c), we followed the decisions of Claude (as well as GPT-4 if we look into its judges), and count degenerates as unsuccessful defences.

5 Experiments

This section begins by introducing the general experimental setup common to all experiments. Subsequently, we provide detailed descriptions of the design of each experiment and present the results.

5.1 General Experimental Setup

Models. Following Wang et al. (2024c), we used the LLaMA2-7B-Chat³ (henceforth, LLaMA) and Mistral-7B-v0.1⁴ (henceforth, Mistral) in this study.

Evaluation Metrics. We used Defence Success (DS; see Section 3 and Section 4.2) for evaluating the detoxification using DINM. We also examined the Defense Generalization using the four metrics mentioned in Section 3, including DG_{onlyQ} , DG_{otherA} , DG_{otherQ} , and DG_{otherAQ} .

³<https://huggingface.co/meta-llama/Llama-2-7b-chat-hf>

⁴<https://huggingface.co/mistralai/Mistral-7B-v0.1>

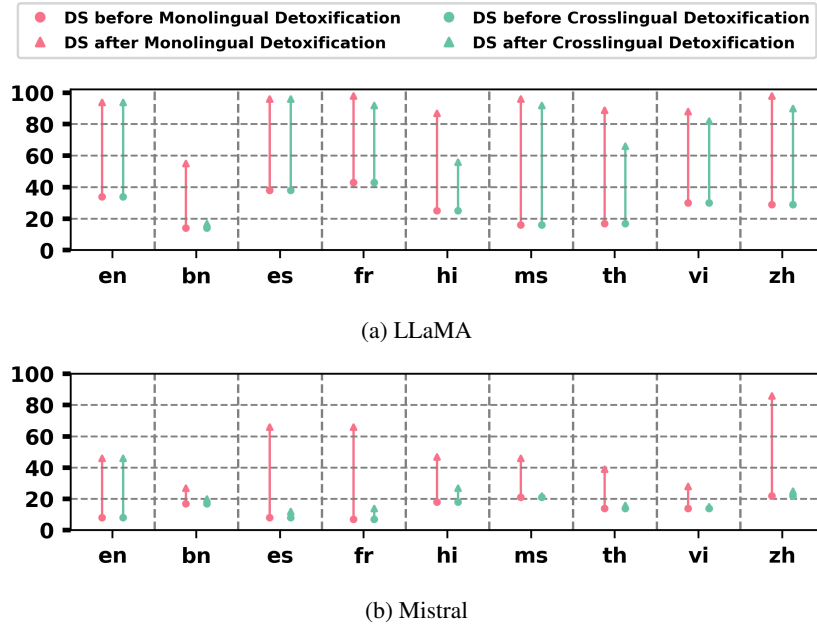


Figure 2: RED: Defence success of LLMs (i.e., LLaMA and Mistral) before and after monolingual KE-based detoxification in various languages; Defence success of LLMs before and after cross-lingual KE-based detoxification in various languages (i.e., editing in English and testing in target language).

Implementation Details. All experiments were conducted on a single NVIDIA A800 GPU (80GB) with approximately 420 hours in total. We followed Wang et al. (2024c) for setting up the hyperparameters⁵.

5.2 Assessing the Generalizability of DINM

As motivated in the introduction, we were curious about whether KE-based detoxification, e.g., DINM, is functional in languages other than English. To this end, we applied and tested DINM on all languages in mSAFEEDIT. We coined detoxification as such as monolingual detoxification as an LLM is edited and attacked in the same language.

The red arrows in Figure 2 chart the DS of both LLaMA and Mistral before and after detoxification. Generally, the DS of LLMs consistently improves following monolingual detoxification, indicating that DINM effectively detoxifies both LLaMA and Mistral across all languages in mSAFEEDIT. These changes in DS also demonstrate that our LLM-based safety classifier can identify safety issues in multiple languages. We put the results of Defense Generalization in the appendix since they show the same trends.

Comparing the two LLMs, LLaMA is a clear winner in terms of safety, which could partly be

attributed to Mistral’s problem of degeneration as discussed in Section 4.2. DINM works better on LLaMA than on Mistral as, on the one hand, the detoxified Mistral has DS scores around or lower than 60% (with merely Chinese as an exception), implying that DINM cannot address the degenerations caused by the attacks. On the other hand, the improvements DINM makes are generally larger on LLaMA than on Mistral.

Comparing the effects of KE-based detoxification in different languages, we notice that DINM appears to perform better in high-resource languages, such as Chinese (zh), Spanish (es) and French (fr), than in low-resource languages, such as Bengali (bn) and Vietnamese (vi). For instance, detoxified LLaMA can achieve almost 100% DS scores on the high-resource languages, while it only receives a DS at less than 60% for Bengali. Such a phenomenon is more significant in detoxified Mistral, it has way smaller effects on low-resource languages than on high-resource languages.

5.3 Testing the language-dependency Hypothesis

Recall that the language-dependency hypothesis suggested that classical knowledge technicals such as the one used in DINM can only edit knowledge in one language, which does not affect the same knowledge in other languages. To examine this

⁵The code and the hyperparameters can be found at: <https://github.com/zjunlp/EasyEdit/>

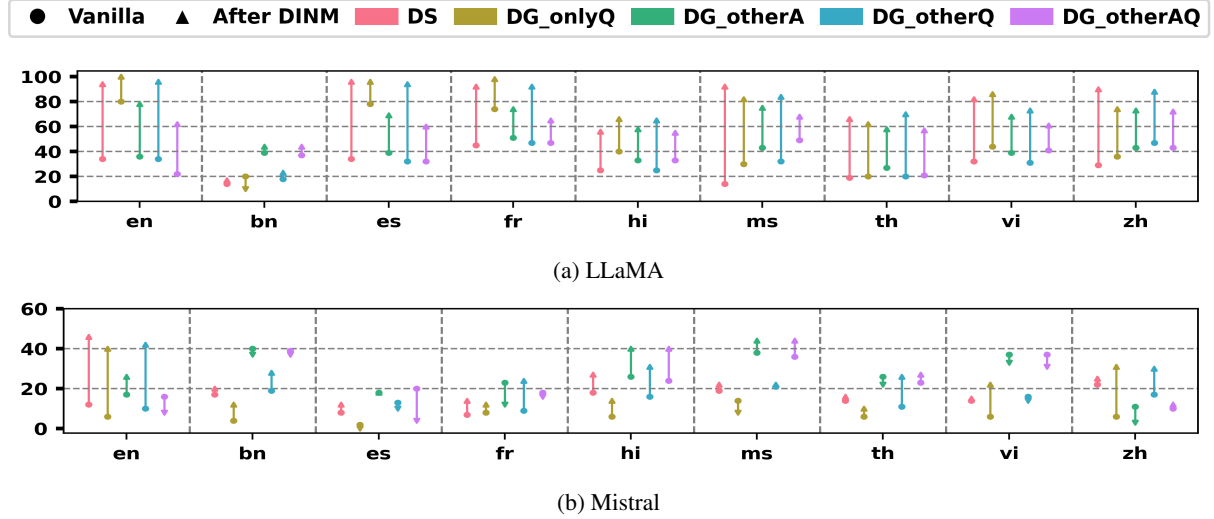


Figure 3: Defence Generalization for the cross-lingual KE-based detoxification.

hypothesis, we carried out a cross-lingual identification experiment. Precisely, we did KE-based detoxification using English data in mSAFEEDIT (as detoxification in the past was often done in English), and tested it in other languages. We report the results in terms of the DS with the green arrows in Figure 2 and the Defence Generalization results in Figure 3.

The situation differs significantly between LLaMA and Mistral. For LLaMA, the language-dependency hypothesis appears to not hold. Edits made to English toxic knowledge generally transfer effectively to nearly all languages in mSAFEEDIT. In most cases, editing toxic knowledge in English produces effects on the target language that are comparable to directly editing toxic knowledge in the target language. The only exception is Bengali: KE-based detoxification in English makes no improvement on the safety of LLaMA when it speaks Bengali. Considering that English is in the same language family as Bengali but in different language families as languages like Malay and Chinese, KE-based cross-lingual detoxification appears to rely more on whether the target language is high-resource rather than on linguistic similarity.

In contrast, the language-dependency hypothesis holds for Mistral: edits to English toxic knowledge have almost zero improvements to Mistral in any languages other than English in mSAFEEDIT. Moreover, if we focus on Defence Generalization results in Figure 3, we find that it not only has no contribution but sometimes has negative effects.

In aggregate, the language-dependency hypothe-

sis does not fully hold. Whether the detoxification in one language works in another language depends on whether the target language is high-resource and which LLM is used.

5.4 Post-hoc Analysis

To ascertain the analyses above, we added two additional experiments. One aims to check whether cross-lingual detoxification still works if languages other than English are used for editing. The other is to understand how the mechanism of DINM impacts its ability of cross-lingual detoxification.

5.4.1 English as the Target Language.

Due to the limitation of our computing resources, we are unable to test every language pair in mSAFEEDIT. Instead, we only tried cross-lingual detoxification with English as the target language and data in other languages as the source for editing LLMs. Figure 4 shows the results in terms of DS and DG.

Compared to detoxifying LLaMA using English data, cross-lingual detoxification using data in other languages does not have on-par effects. Only when using French and Chinese data for editing, LLaMA has significant improvements for defending against attacks in English. Nonetheless, such improvements are still smaller compared to detoxification using English data (cf. Figure 3). One possible explanation for why only French and Chinese are useful in cross-lingual detoxification is that they are both very high-resource languages right after English (given the data in Bender (2009)), but this cannot explain the case of Spanish, which is often

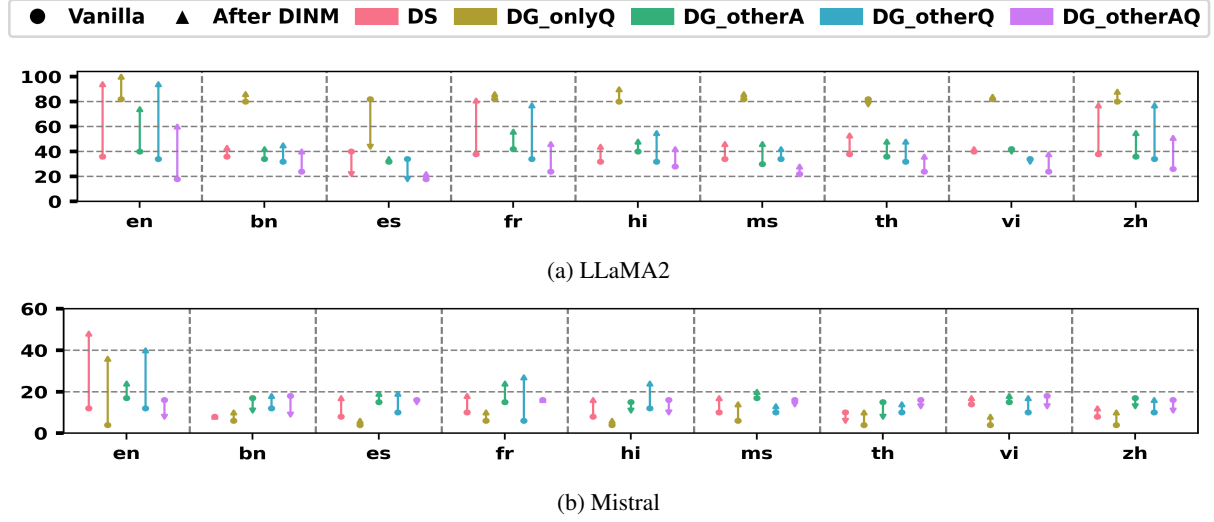


Figure 4: The DS and DG for the cross-lingual detoxified LLMs with English as the target language, i.e., detoxifying using data in one language in mSAFEEDIT and testing it in English.

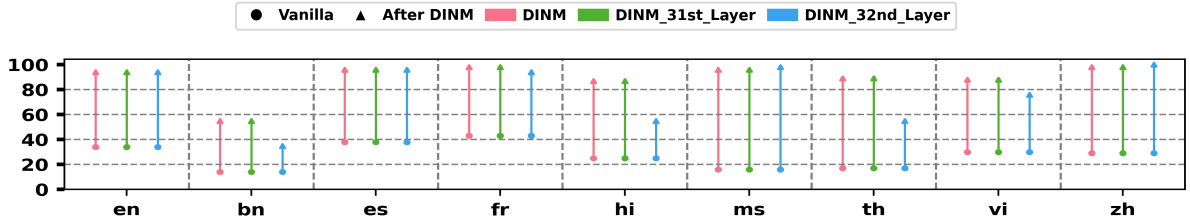


Figure 5: The DS of LLaMA before and after detoxification by full DINM, DINM that directly edits the 31st layer, and DINM that directly edits the 32nd layer.

considered as high-resource while has a negative effect when serving as data for cross-lingual detoxification. Talking about Mistral, in line with the results in Figure 3, KE-based cross-lingual detoxification using other languages still does not work for attacks in English.

These appear to suggest that if we want to make effective cross-lingual detoxification using DINM, it is important to make use of data in very high-resource languages with English as the best choice followed by French and Chinese.

5.4.2 The Role of Toxic Layer Identification.

Toxic layer identification is a critical component of DINM, as it determines which parameters need editing. However, it also makes DINM extremely time-consuming, requiring a scan of all layers for each piece of knowledge to be edited. Previous research on effective parameter tuning suggests that the last few layers in Transformer models often contain the most conceptual information. This raises the question: what if toxic layer identification is bypassed and edits are applied directly

to the last layer? To explore this, we evaluated the effectiveness of both monolingual and cross-lingual detoxification using DINM without toxic layer identification, focusing on edits made to either the last layer or the second-to-last layer of an LLM. Since cross-lingual detoxification does not work on Mistral, we, therefore, only tested LLaMA in this experiment. The second-to-last and the last of LLaMA are the 31st and 32nd layers. Figure 5 depicts the results of cross-lingual detoxification using English data.

The results embody that directly editing the 31st layer yields the same performance as always selecting the 31st layer is exactly the decision of toxic layer identification. Interestingly, we found that editing a different layer—the 32nd layer—does not impact the effectiveness of monolingual detoxification in English. More notably, while editing the 32nd layer reduces the effectiveness of cross-lingual detoxification in certain languages, it enhances the performance in some other languages, such as Malay and Chinese.

6 Discussion

As motivated in the introduction, the two primary questions that this work attempts to answer are whether the language-dependency hypothesis for knowledge editing holds in the context of detoxification and whether DINM is a robust detoxifier in multi-lingual scenarios.

6.1 The language-dependency Hypothesis

Given the experimental results reported in Section 5.3 and 5.4, the language-dependency hypothesis is only partly true in the content of KE-based detoxification, but this still makes the DINM cannot provide sufficient protection in multi-lingual scenarios for toxic knowledge it has seen (which will be further discussed in the next subsection).

Cross-lingual detoxification through knowledge editing is effective (i.e., the language-dependency hypothesis is rejected) only when the following three conditions are met. **First**, given the observation that cross-lingual detoxification is successful only when using English, French, or Chinese, the data used for editing the LLM must be in a dominant language, with English being the preferred choice. This seems to be consistent with the finding in Wu et al. (2024) who coarsely manipulates LLM’s behaviour. **Second**, given the observation that cross-lingual detoxification is ineffective in Bengali, the attack should not be conducted in very low-resource languages. KE-based detoxified LLMs exhibit greater vulnerability to attacks in low-resource languages, irrespective of the language used for the edits. **Third**, given the observation that cross-lingual detoxification does not work on Mistral, the LLM to be detoxified has to be robust enough. In other words, the effect of cross-lingual detoxification is model-dependent.

In addition, following the idea that the success of DINM in English reveals that LLMs may possess a “toxic region”, where multiple specific neurons are linked to particular types of attacks. Our experimental results add to this explanation: On the one hand, this toxic region appears to be shared among languages that are not extremely low-resource. Sufficient data is required to enable LLMs to align toxic knowledge in one language to this region. On the other hand, our findings suggest that this region is not singular. An LLM may contain multiple toxic regions, and editing any of these regions can influence the model’s final outputs.

6.2 The Effectiveness of DINM

DINM works in all languages (at least in all languages in mSAFEEDIT), suggesting its good generalizability. However, it provides reduced or even no protection in the following cases: (1) The effect of DINM is reduced if the LLM is edited by an extremely low-resource language. (2) Its effects are model-dependent. It works worse on weaker LLMs, e.g., Mistral. (3) It only provides conditional cross-lingual protection (see Section 6.1). (4) It seems to have no use in helping defending attacks that lead to degeneration.

In relation to the above limitations, it is worth mentioning that recent studies found that LLMs (including commercial ones like GPT-4) are fragile against attacks in low-resource languages (Yong et al., 2023). Apparently, DINM is unable to address this safety issue as it would have reduced effect if it use the data in the same low-resource language for editing and would have no effect if it use other languages for editing.

Finally, DINM has also suffered from being slow, making it sometimes not the preferred detoxifier if there are too many attacks to be edited or if only limited computing resources are available. Luckily, an easy solution is to eliminate the slowest module in DINM, i.e., toxic layer identification, and roughly edit the second-to-last layer. Such elimination makes no performance reduction in our experiments.

7 Conclusion

This study investigates the language-dependency hypothesis in the context of detoxification, which posits that knowledge editing-based detoxifiers, such as DINM, do not contribute to defending against attacks in languages other than the one used for detoxification. Our experiments challenge this hypothesis, demonstrating that cross-lingual KE-based detoxification is feasible if three conditions are met: (1) the detoxification data must be in a dominant language (e.g., English); (2) the LLM being detoxified must be sufficiently robust (e.g., LLaMA preferred over Mistral); and (3) the attacks must not involve very low-resource languages.

Additionally, we analysed the robustness of DINM as a detoxification method, highlighting its strengths and weaknesses. Our findings indicate that DINM may not provide sufficient protection against attacks in very low-resource languages, regardless of the language used for detoxification.

Limitations

Since knowledge-editing is an extremely computing resource and time-consuming technique, we made three simplifications: (1) We only sampled 50 items from SAFEEDIT to form mSAFEEDIT. As we have argued in Section 4, we believe that 50 items are sufficiently large for making scientific conclusions. (2) When assessing cross-lingual detoxification, we did not test every language pair in mSAFEEDIT. Instead, we merely examined the most important set of pairs, i.e., detoxification using English and attacking using other languages as well as detoxification using other languages and attacking using English. (3) We only tested two LLMs in this study, LLaMA2-7B-Chat and Mistral-7B-v0.1 following Wang et al. (2024c). Both LLaMA and Mistral have newer versions. The conclusions made specifically to these models may change if newer versions are used.

Because of the aim of scaling the experiments to include more languages, especially low-resource ones, another limitation of our study is the reliance on automatically translated test items and automated tools to evaluate the safety of model responses instead of using human experts. These inevitably introduce biases to our conclusions.

References

Amy Au. 2024. Evaluating AI red teaming’s readiness to address environmental harms: A thematic analysis of LLM discourse. In *AAAI*, pages 23726–23728. AAAI Press.

Emily M. Bender. 2009. Linguistically naïve != language independent: Why NLP needs linguistic typology. In *Proceedings of the EACL 2009 Workshop on the Interaction between Linguistics and Computational Linguistics: Virtuous, Vicious or Vacuous?*, pages 26–32, Athens, Greece. Association for Computational Linguistics.

Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models. In *EMNLP (I)*, pages 6491–6506. Association for Computational Linguistics.

Pengfei Cao, Yuheng Chen, Zhuoran Jin, Yubo Chen, Kang Liu, and Jun Zhao. 2024. One mind, many tongues: A deep dive into language-agnostic knowledge neurons in large language models.

Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. 2024. Journey to the center of the knowledge neurons: Discoveries of language-independent knowledge neurons and degenerate knowledge neurons. In *AAAI*, pages 17817–17825. AAAI Press.

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge neurons in pretrained transformers. In *ACL (I)*, pages 8493–8502. Association for Computational Linguistics.

Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. Plug and play language models: A simple approach to controlled text generation. In *ICLR*. OpenReview.net.

Emily Dinan, Angela Fan, Ledell Wu, Jason Weston, Douwe Kiela, and Adina Williams. 2020. Multi-dimensional gender bias classification.

Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, Andy Jones, Sam Bowman, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Nelson Elhage, Sheer El Showk, Stanislav Fort, Zac Hatfield-Dodds, Tom Henighan, Danny Hernandez, Tristan Hume, Josh Jacobson, Scott Johnston, Shauna Kravec, Catherine Olsson, Sam Ringer, Eli Tran-Johnson, Dario Amodei, Tom Brown, Nicholas Joseph, Sam McCandlish, Chris Olah, Jared Kaplan, and Jack Clark. 2022. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *CoRR*, abs/2209.07858.

Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. Realtotoxicityprompts: Evaluating neural toxic degeneration in language models. In *EMNLP (Findings)*, volume EMNLP 2020 of *Findings of ACL*, pages 3356–3369. Association for Computational Linguistics.

Zhankui He, Zhouhang Xie, Rahul Jha, Harald Steck, Dawen Liang, Yesu Feng, Bodhisattwa Prasad Majumder, Nathan Kallus, and Julian J. McAuley. 2023. Large language models as zero-shot conversational recommenders. In *CIKM*, pages 720–730. ACM.

Peng Hu, Sizhe Liu, Changjiang Gao, Xin Huang, Xue Han, Junlan Feng, Chao Deng, and Shujian Huang. 2024. Large language models are cross-lingual knowledge-free reasoners. *CoRR*, abs/2406.16655.

Xiaowei Huang, Wenjie Ruan, Wei Huang, Gaojie Jin, Yi Dong, Changshun Wu, Saddek Bensalem, Ronghui Mu, Yi Qi, Xingyu Zhao, Kaiwen Cai, Yanghao Zhang, Sihao Wu, Peipei Xu, Dengyu Wu, André Freitas, and Mustafa A. Mustafa. 2023. A survey of safety and trustworthiness of large language models through the lens of verification and validation. *CoRR*, abs/2305.11391.

Md. Tahmid Rahman Laskar, Xue-Yong Fu, Cheng Chen, and Shashi Bhushan TN. 2023. Building real-world meeting summarization systems using large language models: A practical perspective. In *EMNLP (Industry Track)*, pages 343–352. Association for Computational Linguistics.

717	Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. CAMEL: communicative agents for "mind" exploration of large language model society. In <i>NeurIPS</i> .	773
718		774
719		775
720		776
721	Todor Markov, Chong Zhang, Sandhini Agarwal, Florentine Eloundou Nekoul, Theodore Lee, Steven Adler, Angela Jiang, and Lilian Weng. 2023. A holistic approach to undesired content detection in the real world. In <i>AAAI</i> , pages 15009–15018. AAAI Press.	777
722		778
723		779
724		780
725		781
726	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In <i>NeurIPS</i> .	782
727		783
728		784
729	Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. Mass-editing memory in a transformer. In <i>ICLR</i> . OpenReview.net.	785
730		786
731		787
732		788
733	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022a. Fast model editing at scale .	789
734		790
735		791
736	Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. 2022b. Memory-based model editing at scale. In <i>ICML</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 15817–15831. PMLR.	792
737		793
738		794
739		795
740		796
741	OpenAI. 2023. GPT-4 technical report. <i>CoRR</i> , abs/2303.08774.	797
742		798
743	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick S. H. Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander H. Miller. 2019. Language models as knowledge bases? In <i>EMNLP/IJCNLP (1)</i> , pages 2463–2473. Association for Computational Linguistics.	799
744		800
745		801
746		802
747		803
748		804
749	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In <i>NeurIPS</i> .	805
750		806
751		807
752		808
753	Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2021. Societal biases in language generation: Progress and challenges. In <i>ACL/IJCNLP (1)</i> , pages 4275–4293. Association for Computational Linguistics.	809
754		810
755		811
756		812
757		813
758	Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric P. Xing, Furong Huang, Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, John C. Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen, Tianming Liu, Tianyi Zhou, William Wang,	814
759		815
760		816
761		817
762		818
763		819
764		820
765		821
766		822
767		823
768		824
769		825
770		826
771		827
772		828

Zheng Xin Yong, Cristina Menghini, and Stephen Bach. 2023. Low-resource languages jailbreak gpt-4. In *Socially Responsible Language Modelling Research*.

Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen R. McKeown, and Tatsunori B. Hashimoto. 2023. Benchmarking large language models for news summarization. *CoRR*, abs/2301.13848.

Xue Zhang, Yunlong Liang, Fandong Meng, Songming Zhang, Yufeng Chen, Jinan Xu, and Jie Zhou. 2024. Multilingual knowledge editing with language-agnostic factual neurons. *CoRR*, abs/2406.16416.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A survey of large language models. *CoRR*, abs/2303.18223.

Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023. Can we edit factual knowledge by in-context learning? In *EMNLP*, pages 4862–4876. Association for Computational Linguistics.

A Defence Generalisation

Figure 6 shows the results of Defence Generalisation for monolingual detoxification.

B Complement Results with GPT-4 as the Safety Classifier

Figure 7-9 are the complement results to the results in the main content with GPT-4 as the safety classifier.

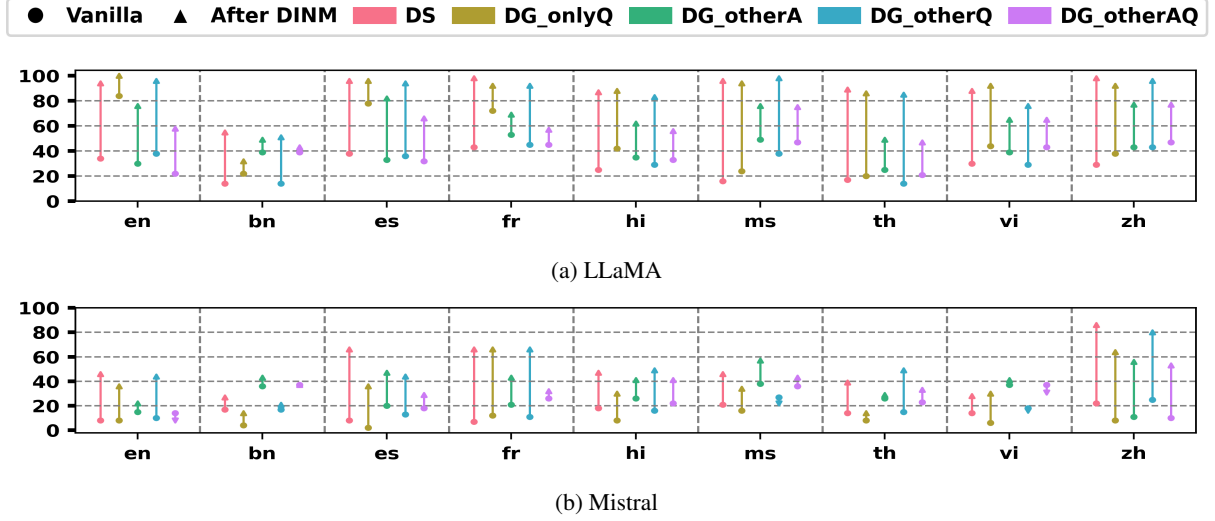


Figure 6: Defence Generalisation for monolingual detoxification.

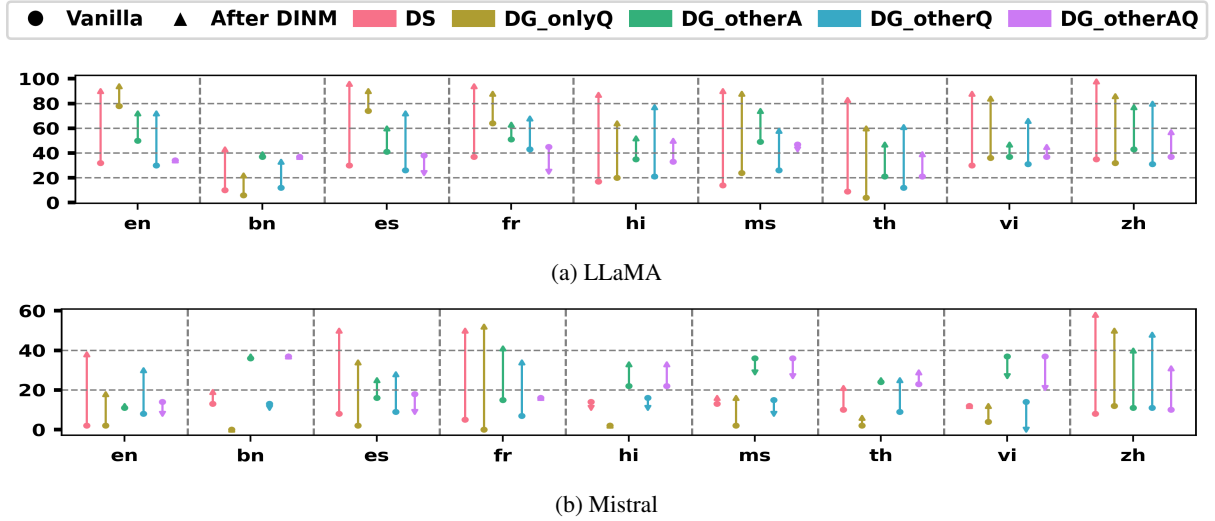


Figure 7: The DS and DG for monolingual detoxification with GPT-4 as the safety classifier.

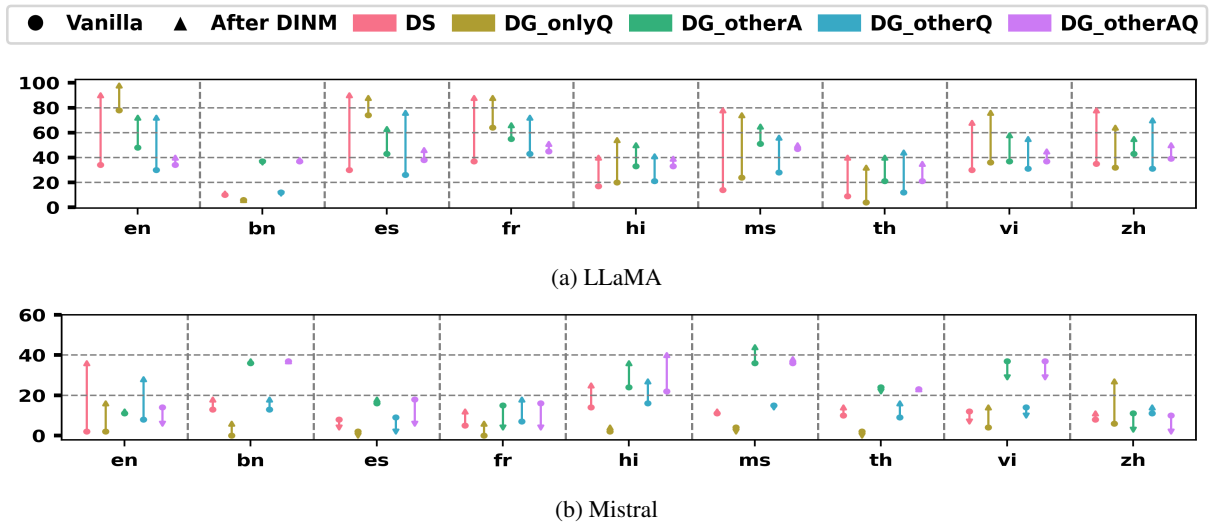


Figure 8: The DS and DG for cross-lingual detoxification with GPT-4 as the safety classifier.

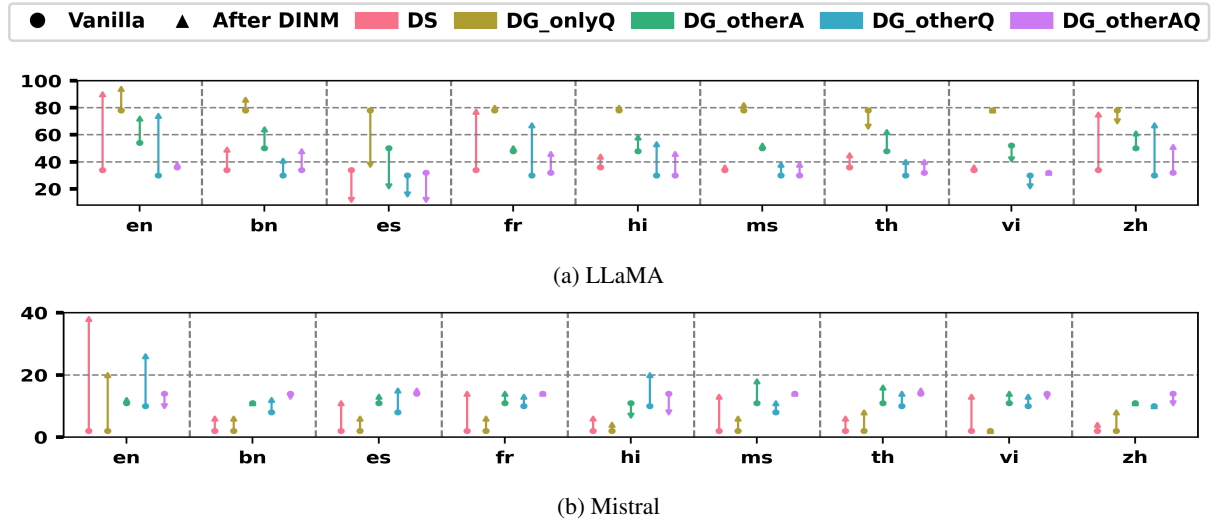


Figure 9: The DS and DG for cross-lingual detoxification with English as the target language and GPT-4 as the safety classifier.