

# LLAMAS (MOSTLY) THINK IN ENGLISH: ON CAUSAL INTERVENTIONS IN THE LATENT LANGUAGE OF TRANSFORMERS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Previous research on the Llama-2 family of Large Language Models (LLMs) suggested a correlation indicating the use of English as an intermediary language within these models for tasks in non-English languages. We improve on this by demonstrating a causal relationship. By intervening on the intermediate layers during a forward pass, we show that projecting out the activations onto a subspace corresponding to the correct prediction in English impairs the model’s ability to make correct predictions on non-English translation tasks. Projecting out an unrelated English subspace, or a related subspace in a non-English language, has little effect, demonstrating that this family of models store concepts that have a high similarity to the corresponding concept in English in the residual stream.

## 1 INTRODUCTION

The dramatic abilities of Large Language Models (LLMs) using the transformer architecture Vaswani et al. (2017); Phuong & Hutter (2022) are rather surprising, given the sole goal during training is to predict the next word in a sequence. LLMs generalize to many out-of-distribution tasks, and exhibit abilities typically associated with intelligence, such as solving difficult maths problems, tool use, and demonstrating theory of mind Bubeck et al. (2023). Interestingly, models predominantly trained on English data tend to perform well in other languages, even when other languages constitute a tiny proportion of the training data K et al. (2020); Blevins & Zettlemoyer (2022). Tianyi Tang et al. (2024) show that language-specific neurons in LLMs are responsible for their multilingual capabilities, and that activating or deactivating these neurons can control the output language. Additionally, Julen Etxaniz et al. (2023) show that LLMs trained predominately in English can perform better on tasks in non-English by explicitly prompting the model to translate to English, solving the task in English, and then translating back to the target language. Shi et al. (2022) show that by using chain-of-thought Chu et al. (2023) prompting, models can perform vastly better on even obscure languages.

### 1.1 LLAMA WORKS IN ENGLISH WENDLER ET AL. (2024)

Wendler et al. (2024) claim that the Llama-2 family of multilingual transformers “work in English” by showing that on translation tasks between non-English languages, the model assigns a high probability to the corresponding answer in English midway through a forward pass. More precisely, given the multi-shot translation prompt<sup>1</sup> used by Wendler et al. (2024) from French (*Français*) to Chinese (中文),

Français: " vertu" - 中文: "德"  
 Français: " siège" - 中文: "座"  
 Français: " neige" - 中文: "雪"  
 Français: " montagne" - 中文: "山"  
 Français: " fleur" - 中文: "

<sup>1</sup>We modify the prompt provided in Wendler et al. (2024) slightly and prepend spaces to words using the latin alphabet to aid with tokenization.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

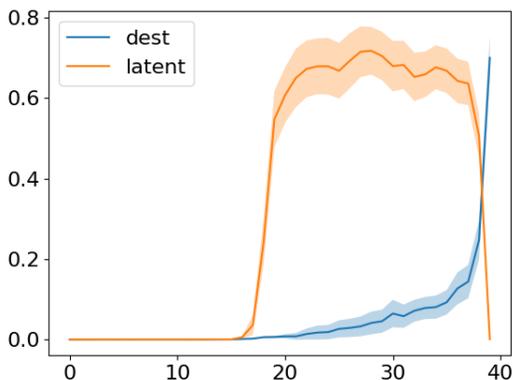


Figure 1: The probability assigned both to the correct answer in Chinese, and in the latent language English when translating from French to Chinese. Probability measured for every layer using *logit lens* for Llama-2-7b. We sum the probability mass over not just the exact correct answer in Chinese, but other semantically similar answers as well.

the model, using in-context-learning Brown et al. (2020) predicts the next token as 花, the correct translation of fleur (flower) from French to Chinese. Surprisingly, if the technique of *logit lens* Nostalgebraist (2020) is used, and the residual stream midway through a forward pass is passed through the unembedding stage of the model, a high probability is assigned to the corresponding token flower in English (Figure 1). We refer to the language being translated from as the *source* (here, French) the language being translated to as the *target* (Chinese) and the language we measure the probability mid-way through the forward pass as the *latent* language.

This effect is observed regardless of which non-English languages are chosen for the source and target, and only if the latent language is English. This is suggestive that the model is using English as an intermediary language for translation tasks. However, this is only a correlation, and may just be an artefact of tokenization, or that the vocabulary of the model is predominantly English.

## 1.2 MAIN CONTRIBUTION

For translation tasks between non-English languages on the Llama-2-7b model Touvron et al. (2023):

- We show a causal relationship between the unembedding vectors of the corresponding answer in English and the downstream prediction of the model in the target language. By computing the vector projection  $\text{proj}_{\mathbb{S}}(h_i)$  of the intermediate activations  $h_i$  onto a subspace  $\mathbb{S}$  corresponding to the correct prediction in English, and subtracting this projection out, this impairs the model’s ability to make correct predictions on non-English translation tasks. This effect is far weaker when projecting out a related subspace in a non-English language, and not present at all when projecting out an unrelated subspace (using words unrelated to that being translated), regardless of the language chosen to intervene with.
- Normally, steering vectors Turner et al. (2023) are obtained by computing a forward pass over inputs of interest, and extracting vectors of interest from the activations. We show a weak effect demonstrating that on translation tasks, we can cause the model to predict a desired counterfactual translation into the target language by steering directly using the unembedding vectors of the corresponding word in English. Note that this effect is not robust, and can in many cases cause the model activations to diverge during a forward pass.

## 2 PRELIMINARIES

### 2.1 MOTIVATION

LLMs are impressive in the capabilities that they have, but it is not at all clear how they can perform the tasks they do. After training, the mechanism by which the model can make good predictions is present, but buried amongst the weights of the network.

This motivates a more transparent approach to understanding the behaviour of models, and what the field of *mechanistic interpretability* aims to achieve: A holistic understanding of how models can perform particular tasks. Meng et al. (2022) show how to extract and modify particular facts known by the model. Wang et al. (2022) demonstrate how particular attention heads in a transformer model can learn how to solve simple in-context-learning tasks. Li et al. (2023) demonstrate that LLMs internally construct models of the world to solve tasks. However, works like this are very costly in terms of human time spent on analysis of models, and it is desirable to find more automatic approaches to understanding models. Zhang & Nanda (2024) introduce *activation patching*, to take the activations from a forward pass on corrupted input, patch them into the activations of a forward pass on clean input, and observe the causal downstream effect. Todd et al. (2024) present *function vectors*, a direction in latent space representing a particular function or operation, and use this to steer the model on other such inputs.

### 2.2 TRANSFORMERS

The transformer architecture Vaswani et al. (2017) is a deep neural network trained on a large corpus of text data, with the goal Radford et al. (2019); Brown et al. (2020) of predicting the next token in a sequence. The training data is composed of strings of text, converted into a sequence of *tokens*  $t_1, \dots, t_{\text{seq}}$  of length  $\text{seq}$ , where each token is a word or subword obtained using *byte pair encoding* Sennrich et al. (2016). Each token is an integer in the range  $t_i \in \{1, \dots, d_{\text{vocab}}\}$ , where  $d_{\text{vocab}}$  is the size of the *vocabulary* for the model.

The model then converts these tokens into embeddings  $e^{t_1}, \dots, e^{t_{\text{seq}}} \in \mathbb{R}^{d_{\text{model}}}$ , by using tokens to index rows in the embedding matrix  $E \in \mathbb{R}^{d_{\text{vocab}} \times d_{\text{model}}}$ . Stacking these embeddings gives the first hidden state  $H^1 \in \mathbb{R}^{\text{seq} \times d_{\text{model}}}$ .

$$e^i = E_{t_i,:}, \quad H_{i,:}^1 = e^{t_i} \quad (1)$$

These embeddings are then passed through a series  $T_1, T_2, \dots, T_n$  of  $n$  layers called *transformer blocks*. For each hidden state  $H^k$ , the subsequent hidden state  $H^{k+1}$  is computed as  $H^{k+1} = T_k(H^k)$ , giving us a sequence of hidden states, or activations,  $H^1, \dots, H^{n+1}$ . Following Elhage et al. (2021), we refer to the sequence of hidden states as the *residual stream*.

Specific to the Llama-2 family of models Touvron et al. (2023), the transformers blocks  $T_i$  are defined as<sup>2</sup>

$$T_k(H^k) = Z^k + \text{GLU}(\text{RMS}(Z^k))$$

where  $Z^k = H^k + \text{MHA}(\text{RMS}(H^k))$

where GLU is a Gated Linear Unit Shazeer (2020) using the Swish activation function, MHA is the multi-head self-attention mechanism Vaswani et al. (2017), and RMS is the Root Mean Square normalization operation Zhang & Sennrich (2019). Each of the internal components of  $T_k$  have a set of learned weights specific to that layer.

The final hidden state  $H^{n+1}$  is then RMS normalized, and then multiplied by the unembedding matrix  $U \in \mathbb{R}^{d_{\text{vocab}} \times d_{\text{model}}}$  to obtain the *logits*  $L \in \mathbb{R}^{\text{seq} \times d_{\text{vocab}}}$ :

$$L = U(\text{RMS}(H^{n+1})) \quad (2)$$

which are finally passed through a *softmax* Goodfellow et al. (2016) operation to obtain a set of probability distributions  $P \in \mathbb{R}^{\text{seq} \times d_{\text{vocab}}}$ , representing (for a given prefix  $t_1, \dots, t_i$  of the input

<sup>2</sup>We omit the details of positional embeddings. Llama-2 uses Rotary Positional Embedding (RoPE) Su et al. (2023) which are performed inside the self-attention mechanism.

sequence) the model’s prediction for the next token  $t_{i+1}$ .

$$P_{i,j} = \frac{\exp(\mathbf{L}_{i,j})}{\sum_{j'} \exp(\mathbf{L}_{i,j'})} \approx \Pr(t_{i+1} = j | t_1, \dots, t_i) \quad (3)$$

During inference we are only interested in the probability of the next unknown token  $t_{\text{seq}+1}$ , given by the last row  $\mathbf{p} := \mathbf{P}_{\text{seq},:} \in \mathbb{R}^{d_{\text{vocab}}}$ , derived from the last logit vector  $\mathbf{l} := \mathbf{L}_{\text{seq},:} \in \mathbb{R}^{d_{\text{vocab}}}$ .

We do not concern ourselves with the internal details of the self-attention mechanism, but focus only on interventions that modify the hidden states  $\mathbf{H}^1, \dots, \mathbf{H}^{n+1}$  between transformer blocks.

### 2.3 LOGIT LENS

The same argument can be made for transformers, which are also deep models with skip connections, so it is reasonable to expect that the intermediary activations have an interpretable relationship to the final logits (see Appendix A.1). Based on this idea, Nostalgebraist (2020) introduced *logit lens*: During a forward pass, the hidden states  $\mathbf{h}^1 := \mathbf{H}_{\text{seq},:}^1, \dots, \mathbf{h}^{n+1} := \mathbf{H}_{\text{seq},:}^{n+1}$  associated with the last token position/prediction of the next token are cached, and then fed through the final unembedding stage of the model (comprised of an RMS normalization layer, followed by multiplying by the unembedding matrix  $\mathbf{U}$ ) to get increasingly better estimates<sup>3</sup>  $\hat{\mathbf{l}}^1, \dots, \hat{\mathbf{l}}^{n+1}$  of the logits  $\mathbf{l}$ , from which we can recover estimates of  $p$  using the softmax operation.

$$\mathbf{l} = \mathbf{URMS}(\mathbf{h}_{n+1}) \approx \hat{\mathbf{l}}^k = \mathbf{URMS}(\mathbf{h}^k) \quad (4)$$

$$p_i = \frac{\exp(\mathbf{l}_i)}{\sum_{i'} \exp(\mathbf{l}_{i'})} \approx \frac{\exp(\hat{\mathbf{l}}^k[i, j])}{\sum_{j'} \exp(\hat{\mathbf{l}}^k[i, j'])} \quad (5)$$

## 3 EXPERIMENTS

### 3.1 SUBSPACE REJECTION

To demonstrate Llama-2’s reliance on English as an intermediary language, we perform a series of interventions on the residual stream of the model during a forward pass on translation tasks. Information for prediction of the next token must be stored somewhere in the residual stream (it being the only causal path from earlier layers to later layers). We hypothesise that for translation, the information of the concept to translate is stored in a low-dimensional subspace  $\mathbb{S} \subseteq \mathbb{R}^{d_{\text{model}}}$  of the residual stream, and this subspace is similar to a set  $S = \{\mathbf{U}_{i_1,:}, \dots, \mathbf{U}_{i_m,:}\}$  of unembedding vector(s)  $\mathbf{U}_{:,i} \in \mathbb{R}^{d_{\text{model}}}$  for the token(s) describing that concept in English. If so, replacement of  $\mathbf{h}^i$  with the orthogonal projection  $\text{proj}_{\mathbb{S}^\perp}(\mathbf{h}^i)$  of  $\mathbf{h}^i$  onto the orthogonal complement  $\mathbb{S}^\perp$  of  $\mathbb{S}$  should affect the model’s ability to predict the correct answer in the target language.

$$\text{proj}_{\mathbb{S}^\perp}(\mathbf{h}^i) = \mathbf{h}^i - \text{proj}_S(\mathbf{h}^i) = \mathbf{h}^i - S(S^T S)^{-1} S^T \mathbf{h}^i \quad (6)$$

As an abuse of notation, we write  $\mathbb{S}$  as simply  $S$ , and define the subspace  $S$  as  $\text{span}(S)$ . We call  $\mathbf{h}_{\perp S}^i := \text{proj}_{\mathbb{S}^\perp}(\mathbf{h}^i)$  the *rejection* of  $\mathbf{h}^i$  from  $S$ . We intervene on the hidden layers  $\mathbf{h}^i$  between transformer blocks during a forward pass and replace each  $\mathbf{h}^i \leftarrow \mathbf{h}_{\perp S}^i$  during a forward pass over layers of interest. We then observe the downstream effect on the model’s prediction on the target language. We construct two kinds of subspaces to reject,  $S_{\text{lang}}$  and  $S'_{\text{lang}}$  as follows:

- $S_{\text{lang}}(t_{\text{source}})$ : We translate the correct prediction token  $t_{\text{source}}$  from the source language to the corresponding token  $t_{\text{lang}}$  in language *lang*. We consider many such valid translations for a given source word. For example, when translating *livre* (book) from French to English, we consider any of *book*, *manuscript*, *volume* or *tome* as valid translations. We also considered valid translations with and without leading spaces. For clarity, we write spaces as “\_”. If the token  $t_{\text{target}}$  is not present in the vocabulary (in that the tokenizer splits the token into subwords), we use only the first such token (under the assumption that if the model correctly predicts the first subtoken, it will correctly predict the rest Pal et al. (2023)).

<sup>3</sup>Note by definition that  $\hat{\mathbf{l}}^{n+1} = \mathbf{l}$  is exact.

The subspace is then constructed as

$$S_{lang} = \{U[t_{lang}] : t \simeq t_{lang}\} \cup \{U[_t_{lang}] : t \simeq t_{lang}\}$$

where  $t \simeq t_{lang}$  denotes that  $t$  and  $t_{lang}$  are words with semantically identical or similar meanings, but in different languages. For example, given the source word for translation was *livre* (book), and potential translations of *livre* as *book*, *manuscript*, *volume* or *tome*, the corresponding subspace would be<sup>4</sup>

$$S_{en}^{livre} = \text{Span}\{U["book"], U["_book"], U["volume"], U["_volume"], U["man"], U["_manuscript"], U["t"], U["_t"]]\}$$

Here,  $U[x]$  denotes the unembedding vector  $U[i_x]$ , where  $i_x$  is the index of the token  $x$  in the vocabulary.

- $\bar{S}_{lang}(t_{source})$ : Same as above, but we choose an incorrect translation  $t_{lang} \neq t_{source}$  in *lang* of the word  $t_{source}$ , together with other words with similar meanings to  $t_{lang}$ , also in language *lang*.

Following the above example, we might choose the incorrect English translations *dog*, *canine* and *puppy* for *livre*, giving the corresponding subspace

$$\bar{S}_{en}^{livre} = \{U["dog"], U["_dog"], U["can"], U["_can"], U["pu"], U["_pu"]]\}$$

The incorrect translations were chosen by taking a dearrangement of the dataset, ensuring that the correct answer was not chosen, and that the word was not present in the translation prompt.

The subspaces are constructed in this way (with and without spaces) to deal with an artefact of the tokenization process: Often a word may appear in the vocabulary with or without a leading space (or sometimes both are present) and correspondingly, the model may choose to predict the word with or without a leading space.

Wendler et al. (2024) observed the phenomena where the model has three distinct phases during a forward pass: an initial phase where the model is not “thinking” in English, nor the target language (layers 1 to 19), a middle phase where the model is “thinking” in English (layers 20 to 29), and a final phase where the model is “thinking” in the target language (layers 30 to 32), as measured using the logit lens probability (Figure 1).

We perform the intervention over all layers inside the interval  $[a, b]$ , for some  $1 \leq a < b \leq 32$ . A sweep was performed over all possible intervals to perform the rejection intervention on with English as the latent language (see Figure 5a), but we found that the choice of layer was not critical, so long as it was sufficiently deep into the network to cover the region where the model “thinks in English”.

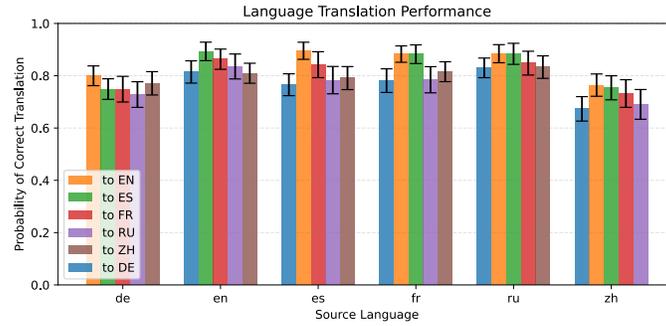
The critical point for where the intervention starts working is around the transition from the initial to middle phase. Note that the rejection intervention does not completely destroy the model’s ability to predict the correct answer, indicating that while the English subspace is important, it does not fully contain all the information required for prediction. So, we fix the rejection intervention to be performed on all layers, and compare to a baseline where no intervention is performed (Figure 2).

An alternative hypothesis is that the residual stream is fragile, and that any similar rejection intervention would have a similar effect. To test this, we perform the same intervention for all choices of source, latent and target languages (where the three languages are distinct), as well as performing the rejection for the unrelated subspace  $S'_{lang}$ .

We find that the rejection intervention has broadly no effect for latent languages that are not English (Figure 3), and no effect for unrelated subspaces in any language (Figure 4), which supports the hypothesis that we are not just damaging the residual stream in general, but performing a targeted intervention on the English subspace.

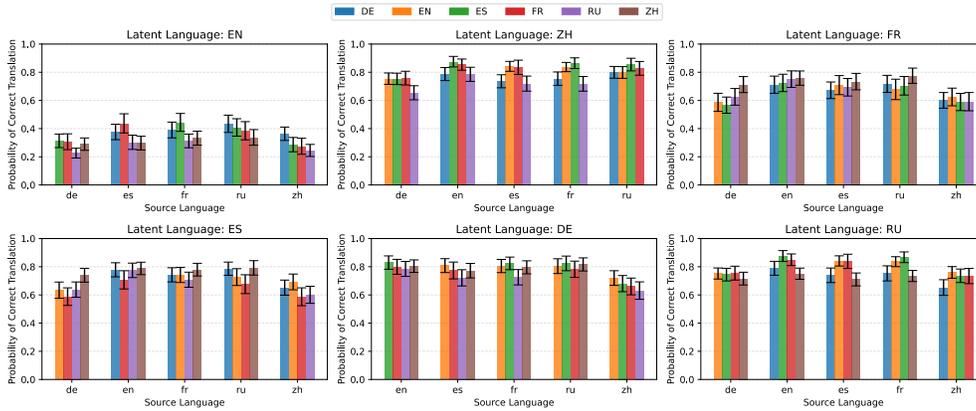
<sup>4</sup>Note that *tome* is not present in the vocabulary for Llama-2, so we only consider the first subtoken  $t$ , and similar for *manuscript*.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281



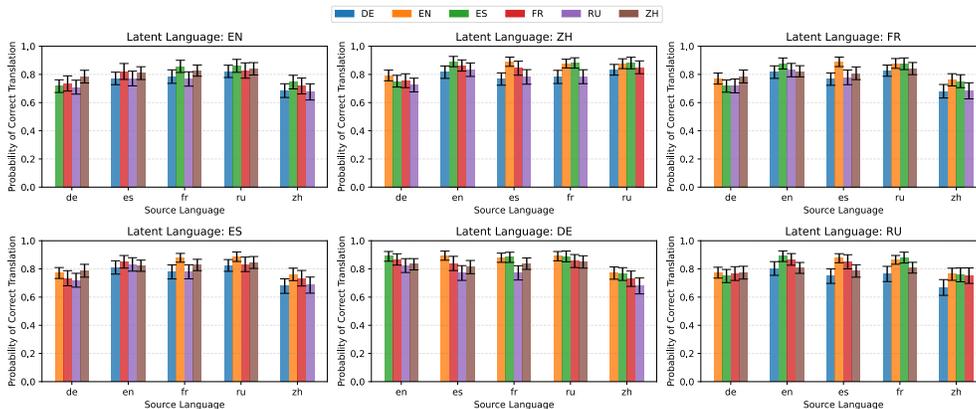
282 Figure 2: A plot of the average probability assigned to the correct answer in the target language when  
283 translating from the source language, for all pairwise translations between French (FR), German  
284 (DE), Chinese (ZH), English (EN), Spanish (EN) and Russian (RU). We see similar performance for  
285 all languages, with slightly worse performance when translating from Chinese. Error bars are 95%  
286 confidence intervals.  
287

288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301



302 Figure 3: A plot of the average probability assigned to the correct answer in the target language  
303 when translating from the source language, and rejecting the latent language, for all possible triplets  
304 of (source, latent, target) languages. We see a stark drop in performance when rejecting the English  
305 subspace, but little effect for other languages. Error bars are 95% confidence intervals.  
306

307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321



322 Figure 4: Same figure as Figure 3, but for the unrelated subspace. We see little effect on the model's  
323 ability to predict the correct answer. Error bars are 95% confidence intervals.

324  
 325  
 326  
 327  
 328  
 329  
 330  
 331  
 332  
 333  
 334  
 335  
 336  
 337  
 338  
 339  
 340  
 341  
 342  
 343  
 344  
 345  
 346  
 347  
 348  
 349  
 350  
 351  
 352  
 353  
 354  
 355  
 356  
 357  
 358  
 359  
 360  
 361  
 362  
 363  
 364  
 365  
 366  
 367  
 368  
 369  
 370  
 371  
 372  
 373  
 374  
 375  
 376  
 377

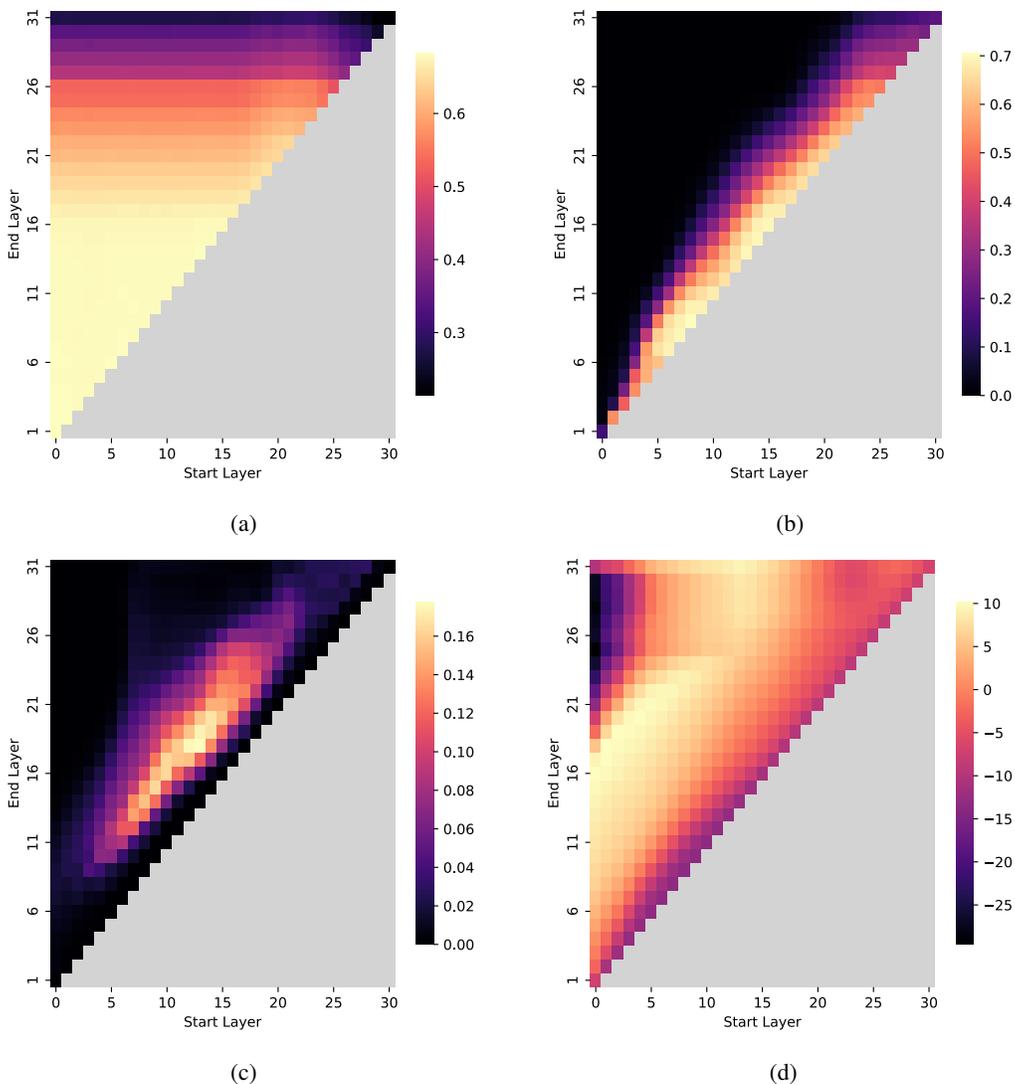


Figure 5: Translation tasks from source=French to target=Chinese, interventions in English.  
 (a) Probability assigned to correct answer with rejection intervention on related subspace (Section 3.1).  
 (b) Probability of correct answer with steering intervention from related subspace, to unrelated subspace (Section 3.2).  
 (c) Same as (b), but probability of desired *counterfactual* answer corresponding to the unrelated subspace we are steering towards (Section 3.2).  
 (d) *Log-probability* of the counterfactual answer, minus the log probability of the correct answer, using steering intervention. Positive values indicate the model is more likely to predict the counterfactual token over token corresponding to the correct translation.  
 Interventions performed between layers  $[a, b]$ , sweeping over all options for  $0 \leq a < b < 32$  (zero-indexing layers), for model Llama-2-7b. Steering interventions use  $c = 8$ .

### 3.2 STEERING TOWARDS A COUNTER-FACTUAL TRANSLATION

The next experiment we perform is to try and steer the model to give a desired counterfactual translation in the target language only by modifying the activations using the unembedding vectors in the latent language. Given a source word  $t_{\text{source}}$  to translate, and subspaces  $S_{\text{lang}}(t_{\text{source}})$  and  $\bar{S}_{\text{lang}}(t_{\text{source}})$  as described in Section 3.1, we perform the intervention

$$\mathbf{h}^i \leftarrow \mathbf{h}^i - \text{proj}_S(\mathbf{h}^i) + c \frac{1}{|\bar{S}|} \sum_{s \in \bar{S}} s \quad (7)$$

where  $c > 0$  is a constant controlling the strength of the steering vector. We then observe the effect on the model’s prediction by sweeping over layers to intervene on, as well as the size of the constant  $c$ .

We found that the strength of the effect grows monotonically with  $c$ , up to the point where the residual stream would diverge for  $c \approx 10$ . As expected, the effect was strongest when using  $U_{\text{en}}$ , but the choice of layer to intervene on is critical (and depends on the choice of  $c$ ). We found the strongest effect (the one that boosts the probability of the counterfactual word in the target language the most) was to intervene on layers 13 to 18, and to choose  $c = 8$ . This is quite a large constant, as the norm of the vector that we project out is  $\approx 4-5$ , so essentially we are adding a larger vector back in to the one we projected out. For more reasonable values of  $c$  (around 1-2), the effect was barely noticeable. This indicates that our steering intervention perhaps is not a suitable one, and that other methods should be investigated. We plot both the probability of the correct answer in the target language, and the probability of the counterfactual answer in the target language to demonstrate the effect is both a suppression of the correct answer, and a boosting of the counterfactual answer (Figure 5b). We also plot the difference in log-probabilities between the counterfactual and correct answer (Figure 5d). See Appendix A.1 for plots of sweeps over choices of  $c$ , and for various choices of languages.

### 3.3 CONCLUSION

Our results are somewhat mixed: We do confirm that to a degree, the LLama-2-7b model is using English as an intermediary language, and that the unembedding vectors in English do explain some of the model’s behaviour in translation tasks. Projecting out in English clearly has a much stronger effect than projecting out in other languages, and the lack of effect for unrelated subspaces does indicate this intervention is not just causing general damage to the model.

However, the effect is not as strong as we would have hoped. There is still a lot of probability mass that the model assigns to the correct prediction after the rejection intervention. We would have liked to have seen the model’s ability to predict the correct answer drop to near zero, but this was not the case. This indicates that the model is storing concepts elsewhere in a subspace that is not easily projected out, or that the model is able to recover from the intervention by using other information stored in the residual stream.

### 3.4 FUTURE WORK

We also performed the same experiment for Llama-2-13b, Gemma-2-2b and Gemma-2-9b, see Appendix A. While the same effect was observed for Llama-2-13b as was for Llama-2-7b, which was expected, the effect was much more destructive for the Gemma models. Any intervention, even those in an unrelated subspace or a non-English language, quite badly affected the model’s ability to predict the correct answer. However, the Gemma models were still more vulnerable to rejections in English than any other language, causing the probability of prediction of the correct token to drop to near zero. We are unsure what to make of this behaviour, perhaps Gemma tends to store concepts in a less English-biased manner, or that the unembedding vectors are similar between languages, or that the Gemma models are more fragile in general. We leave this to future work.

The steering effect that we used was also rather brute-force, and would quite drastically change the activations midway through the network to something well outside the normal distribution. We would like to explore other methods of steering the model using the unembedding vectors only that are less destructive, while still being effective.

432 IMPACT STATEMENT  
433

434 This paper presents work whose goal is to better understand the internals of Large Language Models,  
435 and how we can perform interventions on models in an interpretable fashion. We do not anticipate  
436 any ethical concerns arising from this work, as we are exploring the already existing capabilities of  
437 a pre-existing model.  
438

439 REFERENCES  
440

- 441 AlignmentResearch. Tools for understanding how transformer predictions are built layer-by-layer.  
442 <https://github.com/AlignmentResearch/tuned-lens>, 2023.  
443
- 444 Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella  
445 Biderman, and Jacob Steinhardt. Eliciting Latent Predictions from Transformers with the  
446 Tuned Lens, 2023. URL [https://huggingface.co/spaces/AlignmentResearch/  
447 tuned-lens/tree/main/lens](https://huggingface.co/spaces/AlignmentResearch/tuned-lens/tree/main/lens).
- 448 Terra Blevins and Luke Zettlemoyer. Language Contamination Helps Explain the Cross-lingual  
449 Capabilities of English Pretrained Models, 2022.
- 450 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-  
451 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,  
452 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.  
453 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,  
454 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford,  
455 Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.  
456
- 457 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece  
458 Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi,  
459 Marco Tulio Ribeiro, and Yi Zhang. Sparks of Artificial General Intelligence: Early experiments  
460 with GPT-4. pp. 54–60, 2023.
- 461 Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng,  
462 Ming Liu, Bing Qin, and Ting Liu. A Survey of Chain of Thought Reasoning: Advances, Frontiers  
463 and Future, 2023.
- 464 Clément Dumas, Veniamin Veselovsky, Giovanni Monea, Robert West, and Chris Wendler. How do  
465 llamas process multilingual text? a latent exploration through activation patching. In *ICML 2024*  
466 *Workshop on Mechanistic Interpretability*.  
467
- 468 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,  
469 Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep  
470 Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt,  
471 Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and  
472 Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*,  
473 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- 474 Rahim Entezari, Hanie Sedghi, Olga Saukh, and Behnam Neyshabur. The role of permutation  
475 invariance in linear mode connectivity of neural networks, 2022.
- 476 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. [http:  
477 //www.deeplearningbook.org](http://www.deeplearningbook.org).  
478
- 479 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
480 nition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*  
481 *(CVPR)*, June 2016.
- 482 Stanisław Jastrzębski, Devansh Arpit, Nicolas Ballas, Vikas Verma, Tong Che, and Yoshua Bengio.  
483 Residual connections encourage iterative inference, 2018.  
484
- 485 Julen Etxaniz, Gorka Azkune, Aitor Soroa, Oier Lopez de Lacalle, and Mikel Artetxe. Do Multilin-  
gual Language Models Think Better in English?, 2023.

- 486 Karthikeyan K, Zihan Wang, Stephen Mayhew, and Dan Roth. Cross-Lingual Ability of Multilingual  
487 BERT: An Empirical Study. pp. 1–2, 2020.  
488
- 489 Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Watten-  
490 berg. Emergent world representations: Exploring a sequence model trained on a synthetic task,  
491 2023.
- 492 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual  
493 associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.  
494
- 495 Nostalgebraist. Interpreting gpt: The logit lens, 2020. URL [https://www.lesswrong.com/  
496 posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens](https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens). Accessed:  
497 2024-05-22.
- 498  
499 Koyena Pal, Jiuding Sun, Andrew Yuan, Byron Wallace, and David Bau. Future lens: Anticipating  
500 subsequent tokens from a single hidden state. In *Proceedings of the 27th Conference on Compu-  
501 tational Natural Language Learning (CoNLL)*, pp. 548–560. Association for Computational Lin-  
502 guistics, 2023. doi: 10.18653/v1/2023.conll-1.37. URL [http://dx.doi.org/10.18653/  
503 v1/2023.conll-1.37](http://dx.doi.org/10.18653/v1/2023.conll-1.37).
- 504 Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural  
505 networks, 2013.  
506
- 507 Mary Phuong and Marcus Hutter. Formal algorithms for transformers. Technical report, DeepMind,  
508 London, UK, 2022. LaTeX source available at <http://arxiv.org/abs/2207.09238>.
- 509 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
510 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.  
511
- 512 Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with  
513 subword units, 2016.  
514
- 515 Noam Shazeer. Glu variants improve transformer, 2020.
- 516 Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi,  
517 Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Lan-  
518 guage Models are Multilingual Chain-of-Thought Reasoners, 2022.  
519
- 520 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: En-  
521 hanced transformer with rotary position embedding, 2023.
- 522 Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu  
523 Wei, and Ji-Rong Wen. Language-Specific Neurons: The Key to Multilingual Capabilities in  
524 Large Language Models, 2024.  
525
- 526 Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau.  
527 Function vectors in large language models, 2024.  
528
- 529 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
530 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher,  
531 Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy  
532 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,  
533 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel  
534 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,  
535 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,  
536 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,  
537 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh  
538 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen  
539 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,  
Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models,  
2023.

- 540 Alexander Matt Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDi-  
 541 armid. Activation addition: Steering language models without optimization. (arXiv:2308.10248),  
 542 November 2023. doi: 10.48550/arXiv.2308.10248. URL [http://arxiv.org/abs/2308.](http://arxiv.org/abs/2308.10248)  
 543 10248. arXiv:2308.10248 [cs] Citation Key: turner23steering.
- 544 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
 545 Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. pp. 3–7, 2017.
- 546 Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Inter-  
 547 pretability in the wild: a circuit for indirect object identification in gpt-2 small, 2022.
- 548 Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. Do Llamas Work in  
 549 English? On the Latent Language of Multilingual Transformers, 2024.
- 550 Biao Zhang and Rico Sennrich. Root mean square layer normalization, 2019.
- 551 Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models:  
 552 Metrics and methods, 2024.

## 556 A APPENDIX

### 557 REPLICATION DETAILS

558 All code and datasets used for this paper can be found at [REDACTED]. Scripts are included to  
 559 replicate the main results, as well as to generate the plots.

### 560 DATASET DETAILS

561 In Wendler et al. (2024), the dataset was constructed by taking a corpus of words in English, and  
 562 translating them externally into other languages using DeepL. A similarly constructed dataset was  
 563 also used in the subsequent work by (Dumas et al.), which includes many more languages, as well  
 564 as a list of other possible translations. For each non-English language *lang*, we took the English  
 565 column and the *lang* column from their dataset, and prompted Claude 3.5 Sonnet as follows to  
 566 create a high-quality set of additional possible translation from English to language *lang*. Here,  $\langle x \rangle$   
 567 represents the translation of the English word  $x$  into language *lang*.

```
572 en <lang>
573 book <book>
574 cloud <cloud>
575 ...
576 power <power>
577 summer <summer>
```

578 Generate for me a CSV file, first column is the word in english,  
 579 second column is the  $\langle \text{lang} \rangle$  translation provided, and third is a  
 580 list of synonyms of this word-pair in  $\langle \text{lang} \rangle$ . The list should be  
 581 formatted as a list of strings would be in python. Do not force  
 582 the synonyms, they should all be a suitable replacement for the  
 583 word, or a word strongly associated. For example, big and large  
 584 would be suitable synonyms. Apple and pear would not. Small and  
 585 tiny would be suitable. Bag and handbag would be okay, as well  
 586 as fire and flame. third, three and triple are all suitable.  
 587 Machine and car would not be. Generate as many as you can for  
 588 each entry, but no more than 10. If you can't find any, that's  
 589 okay, don't force words just to make the list long. Some lists  
 590 can be short or even empty if the word is very unique and no  
 591 reasonable synonyms exist.

592 The original work included Japanese. We excluded as often the symbol for a word in Chinese and  
 593 Japanese were identical (for example, *water* is expressed as 水 in both Chinese and Japanese), which  
 would make translation artificially easy.

This gave us a dataset of words from English to Chinese, French, German, Russian and Spanish, and vice-versa. We try all 120 permutations of distinct source, latent and target languages, and perform the rejection intervention, both for the related subspace and the unrelated subspace.

We also slightly modify the prompt used by Wendler et al. (2024) (see Section 1.1) to deal with tokenization: we add a space character before the word if the language is non-Chinese (as the vocabulary of Llama-2-7b has many more tokens corresponding to a space character (e.g. "\_hello"), followed by a word, than just the word itself), but we do not add the space for Chinese (Chinese tokens in the vocabulary are just the symbol itself, "好", or the sequence of bytes in the UTF-8 representation of the symbol, "<0xE5><0xA5><0xBD>"). We ignored Chinese characters that could only be represented in the vocabulary as the raw UTF-8 sequence.

The prompt for each translation is fixed once and then used throughout the experiment. We use 4 words for the translation examples in the prompt, leaving the rest for translation. This was required as to allow for both kv-caching of the shared prefix, as well as to perform inference on a large number of words in parallel, to ensure the experiment was computationally feasible in a reasonable amount of time.

### A.1 SKIP CONNECTIONS

If models can make good predictions over a large class of domains, we would expect that somewhere in the residual stream of the transformer, the model stores useful information for prediction of the next token. In a standard feed-forward neural network, each pair of adjacent weight matrices have permutation invariances Entezari et al. (2022), so we should not expect the intermediary activations to have any clear resemblance to the final logits on the output.

It has been shown that deep vision neural networks pose a problem to train, due to gradients vanishing or exploding Pascanu et al. (2013). This was solved by adding skip connections He et al. (2016), which allowed very deep models to be trained effectively. As a side effect, transformers have been hypothesised to encourage the model to perform iterative inference: the intermediary activations represent the models' best current guess at the output midway through the forward pass, which gradually move towards the model's final prediction Jastrzebski et al. (2018). There are two factors at play here:

- The presence of skip connections in the transformer breaks permutation invariances, so the choice of basis for the intermediary activations is no longer arbitrary.
- Given some neural network  $f_{\Theta}$  with a skip connection, the operation performed is  $x + f_{\Theta}(x)$ . The addition of L2 weight decay Goodfellow et al. (2016) during training encourages the weights of  $f$  to be small, so all else equal, the model will prefer weights  $\Theta$  such that  $x + f_{\Theta}(x)$  is close to  $x$ .

### REVERSE LENS

This technique had no improvement over adding the unembedding vectors directly into the residual stream (Section 3.2), but we include it (and the code) for posterity. *Tuned Lens* (TL) is a technique introduced by Belrose et al. (2023), which is identical to Logit Lens (LL), except for an extra learned linear mapping is applied to the hidden layers  $h_i$  prior to normalization and unembedding. This gives a layer specific transformation  $TL_i$  of the hidden state  $h_i$ . We can write the operation of Tuned Lens (and Logit Lens to compare) as

$$\begin{aligned}
 LL(\mathbf{h}^i) &= URMS(\mathbf{h}^i) \\
 TL_i(\mathbf{h}^i) &= URMS(\mathbf{h}^i + A_i \mathbf{h}^i + \mathbf{b}^i) \\
 RMS(x) &= \frac{x}{\mu_x} \odot \gamma \\
 \mu_x &= \sqrt{\sum_{i=1}^{d_{\text{model}}} x_i^2 + \epsilon}
 \end{aligned}$$

where  $A^i \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$  and  $\mathbf{b}^i \in \mathbb{R}^{d_{\text{model}}}$  are the learned parameters of the tuned lens for layer  $i$ ,  $U \in \mathbb{R}^{d_{\text{vocab}} \times d_{\text{model}}}$  is the unembedding matrix,  $\mu_x \in \mathbb{R}$  is the scale factor for RMS norm, and

648  $\gamma \in \mathbb{R}^{d_{\text{model}}}$  is a learned parameter for RMS norm, taken from the final normalization layer of the  
 649 transformer. Here,  $\odot$  denotes elementwise multiplication. We do not concern ourselves with how  
 650 the weights  $\mathbf{A}^i, \mathbf{b}^i$  are learned, but make use of the pretrained weights for Tuned Lens available at  
 651 AlignmentResearch (2023).

652 The idea behind tuned lens is that rather than unembedding  $\mathbf{h}^i$  directly, an approximation  $\mathbf{x} \mapsto$   
 653  $\mathbf{x} + \mathbf{A}^i \mathbf{x} + \mathbf{b}^i$  of the composition of the subsequent layers  $\mathbf{x} \mapsto (T_n \circ T_{n-1} \dots \circ T_i)(\mathbf{x})$  is learned  
 654 and applied prior to normalization and unembedding, with the goal of aligning the semantics of the  
 655 hidden state  $\mathbf{h}^i$  with the input that  $U$  expects to operate on.

656 With this in mind, *Reverse Lens* is trying to run Tuned Lens backwards: Given a row of the unem-  
 657 bedding matrix  $U_{i,:}$ , what are the activations that Tuned Lens would map to this row? If Tuned Lens  
 658 can recover the logits from the residual stream, then by running it backwards, we should be able to  
 659 compute an approximation of the model’s internal representation of  $U_{i,:}$  at any particular layer.

660 Unfortunately, RMS norm is not invertible, and  $U$  is not square, so we cannot directly invert Tuned  
 661 Lens. To solve the first problem, note that for a fixed scale factor  $\mu_x$ , RMS is a linear function, so  
 662 we can factor  $\mu_x$  out, and the end result is the output logits will be scaled by  $\mu_x$ . Since the argmax  
 663 is invariant to scaling, we end up absorbing  $\mu_x$  elsewhere, and so we just set  $\mu_x = 1$  in the reverse  
 664 operation. For the second, we do not require the ability to invert an arbitrary logit vector, but only  
 665 vectors that correspond to a particular token that we wish to find the model’s internal representation  
 666 of. The logit  $L_{i,:}$  in the  $i^{\text{th}}$  sequence position in the output is computed as  $L_{i,:} = (URMS(\mathbf{H}))_i \propto$   
 667  $(U_{i,:}) \cdot (\mathbf{H}_{i,:} \odot \gamma)$ , the dot product of the  $i^{\text{th}}$  row of  $U$  with the activations  $\mathbf{H}$  (weighed by the  
 668 RMS norm parameter  $\gamma$ ). Since the dot product of a vector with itself is large, and the unembedding  
 669 matrix  $U$  needs to be able to discriminate between many different tokens, we would expect the rows  
 670 of  $U$  are approximately orthogonal. So, the vector that  $U$  would map to the standard basis vector<sup>5</sup>  
 671  $\mathbf{e}_i$  is approximately  $U_{:,i}$ , the  $i^{\text{th}}$  column of  $U$ . We can then use this as the target for the reverse  
 672 tuned lens operation.

673 We define the *reverse tuned lens* (RTL) for layer  $i$  as follows: Taking as input an index  $j \in$   
 674  $\{1, \dots, d_{\text{vocab}}\}$  and a cached scale factor  $\mu \in \mathbb{R}$ , we define the operation

$$675 \text{RTL}_i(j, \mu) = (\mathbf{I} - \mathbf{A}_i)^{-1}(\text{RMS}^{-1}(U_{:,j}, \mu) - \mathbf{b}_i)\text{RMS}^{-1}(x, \mu) = \mu(x \cdot \gamma^{-1})$$

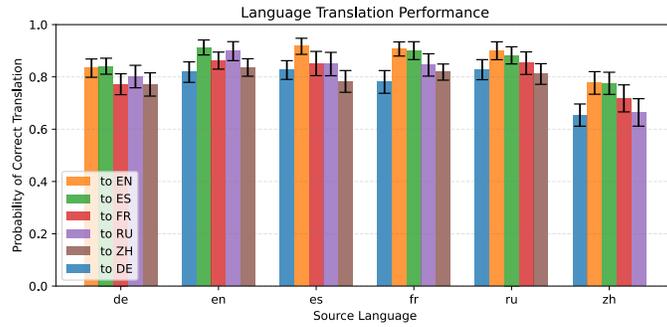
676  
 677  
 678 One can verify that  $\text{argmax}_k(\text{TL}_i(\text{RTL}_i(j, \mu)))_k = j$  for all layers  $i$ , scaling factors  $\mu$  and indices  
 684  $j$ , assuming that for every  $i$ ,  $U_{i,:} \cdot U_{j,:}$  is maximised when  $j = i$ .

685 However, when this was used in practice (repeating the experiment in Section 3.2 with RTL instead  
 686 of the unembedding vectors directly), no appreciable difference in behaviour was observed. The  
 687 rejection effect was equally effective as before. We include the code for the Reverse Tuned Lens  
 688 operation, in the hope that perhaps it may be useful elsewhere.

690  
 691  
 692  
 693 FULL EXPERIMENTAL RESULTS FOR REJECTION

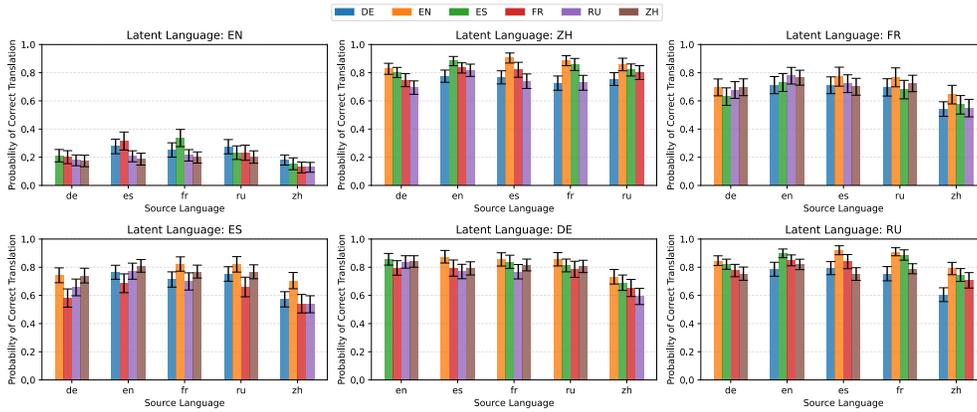
694  
 695  
 696  
 697  
 698  
 699  
 700  
 701 <sup>5</sup>Defined as  $\mathbf{e}_i^i = 1$  and  $\mathbf{e}_j^i = 0$  for  $j \neq i$ .

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713



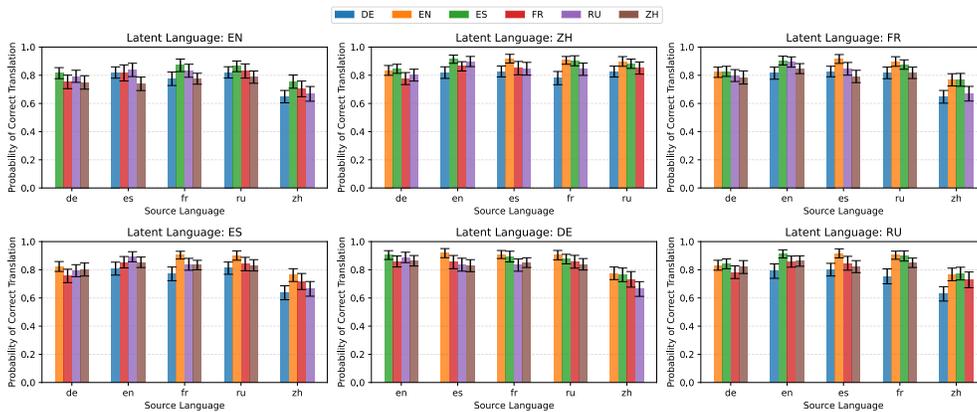
714 Figure 6: A plot of the average probability assigned to the correct answer in the target language when  
715 translating from the source language, for all pairwise translations between French (FR), German  
716 (DE), Chinese (ZH), English (EN), Spanish (EN) and Russian (RU). We see similar performance for  
717 all languages, with slightly worse performance when translating from Chinese. Error bars are 95%  
718 confidence intervals, model is Llama-2-13b.

719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733



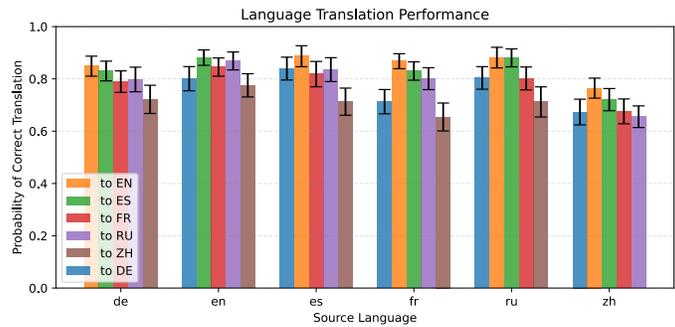
734 Figure 7: A plot of the average probability assigned to the correct answer in the target language  
735 when translating from the source language, and rejecting the latent language, for all possible triplets  
736 of (source, latent, target) languages. We see a stark drop in performance when rejecting the English  
737 subspace, but little effect for other languages. Error bars are 95% confidence intervals, model is  
738 Llama-2-13b.

739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753



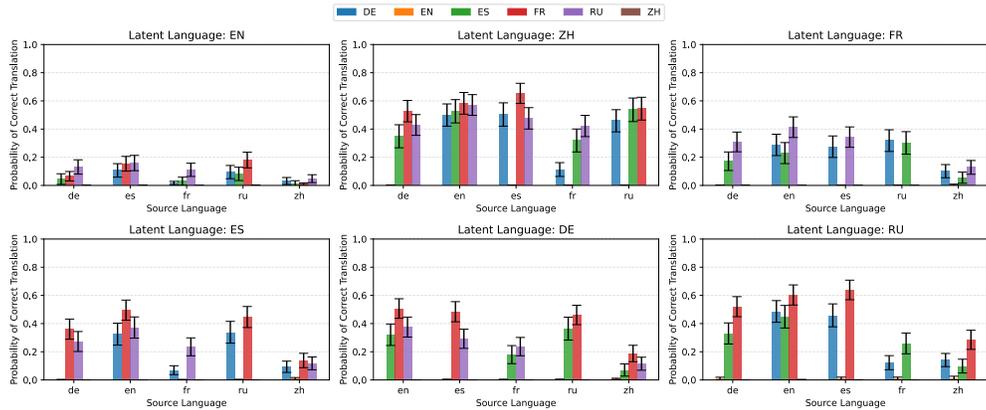
754 Figure 8: Same figure as Figure 3, but for the unrelated subspace. We see little effect on the model's  
755 ability to predict the correct answer. Error bars are 95% confidence intervals, model is Llama-2-13b.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767



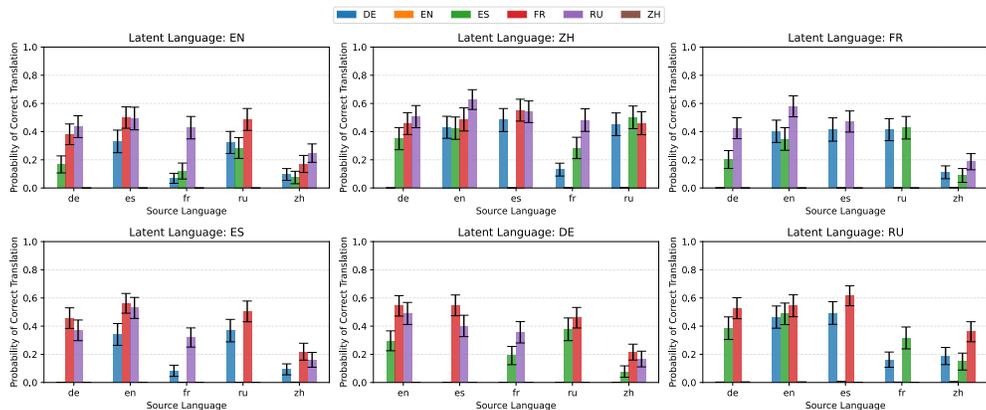
768 Figure 9: A plot of the average probability assigned to the correct answer in the target language when  
769 translating from the source language, for all pairwise translations between French (FR), German  
770 (DE), Chinese (ZH), English (EN), Spanish (EN) and Russian (RU). We see similar performance for  
771 all languages, with slightly worse performance when translating from Chinese. Error bars are 95%  
772 confidence intervals, model is Gemma-2-2b.

773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787



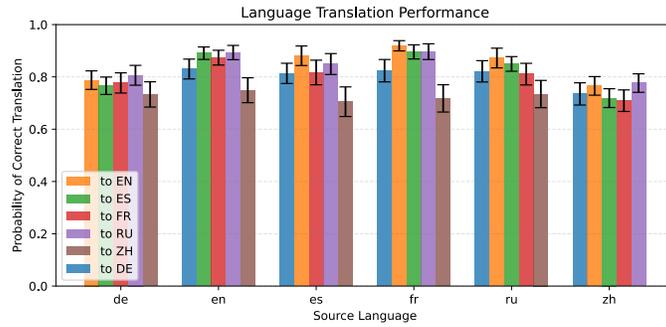
788 Figure 10: A plot of the average probability assigned to the correct answer in the target language  
789 when translating from the source language, and rejecting the latent language, for all possible triplets  
790 of (source, latent, target) languages. We see a stark drop in performance when rejecting the English  
791 subspace, but little effect for other languages. Error bars are 95% confidence intervals, model is  
792 Gemma-2-2b.

793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807



808 Figure 11: Same figure as Figure 3, but for the unrelated subspace. We see little effect on the model's  
809 ability to predict the correct answer. Error bars are 95% confidence intervals, model is Gemma-2-2b.

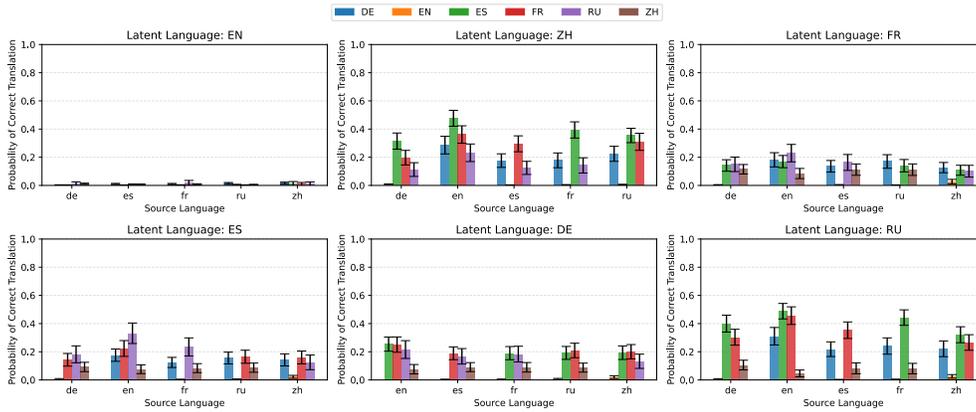
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821



822  
823  
824  
825  
826  
827

Figure 12: A plot of the average probability assigned to the correct answer in the target language when translating from the source language, for all pairwise translations between French (FR), German (DE), Chinese (ZH), English (EN), Spanish (EN) and Russian (RU). We see similar performance for all languages, with slightly worse performance when translating from Chinese. Error bars are 95% confidence intervals, model is Gemma-2-9b.

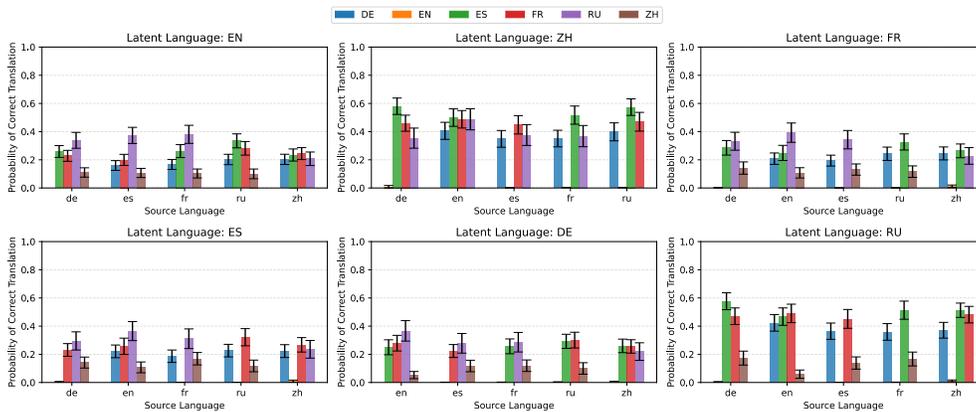
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841



842  
843  
844  
845  
846

Figure 13: A plot of the average probability assigned to the correct answer in the target language when translating from the source language, and rejecting the latent language, for all possible triplets of (source, latent, target) languages. We see a stark drop in performance when rejecting the English subspace, but little effect for other languages. Error bars are 95% confidence intervals, model is Gemma-2-9b.

847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861



862  
863

Figure 14: Same figure as Figure 3, but for the unrelated subspace. We see little effect on the model's ability to predict the correct answer. Error bars are 95% confidence intervals, model is Gemma-2-9b.

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

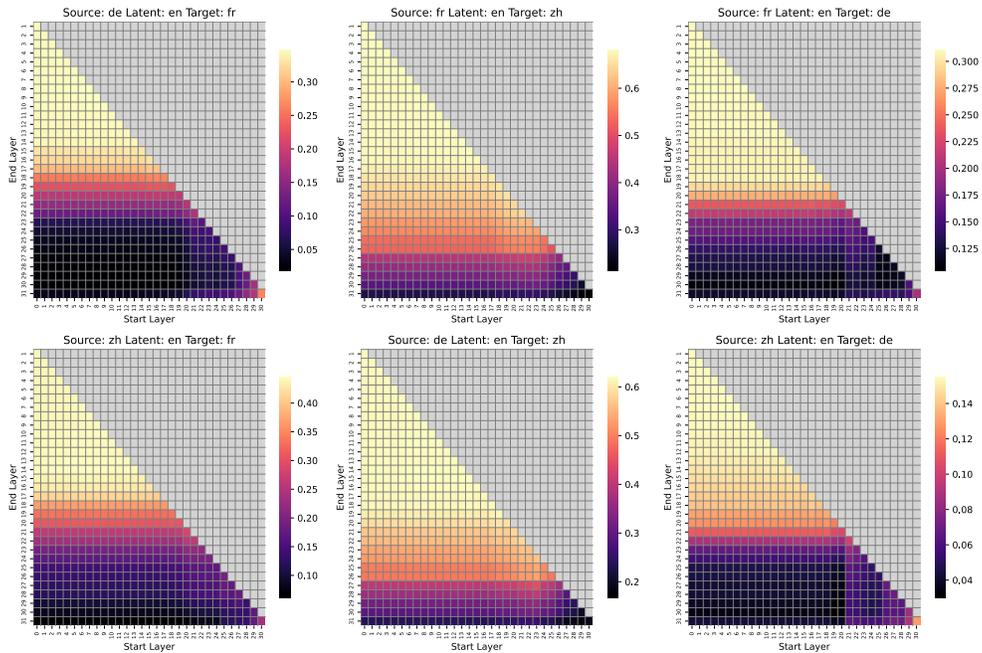


Figure 15: The probability of Llama-2-7b predicting the correct translation, given the rejection intervention on the related English subspace. The intervention was performed across all layers  $\{0, 1, \dots, 31\}$  that lie in the interval (start, end). We sweep over all possible values  $0 \leq \text{start} < \text{end} \leq 31$ , as well as sweep over the choices for source and target languages. Columns sorted by target language.

#### FULL EXPERIMENTAL RESULTS FOR STEERING

Here, we focus only on translating from Chinese to French and vice versa, and perform the steering interventions with either German or English as the latent language. Grey regions either indicate an invalid range for the intervention (the start layer is greater than the end layer), or that the residual stream diverged to infinity, which results in a NaN value when fed into the final RMS norm.

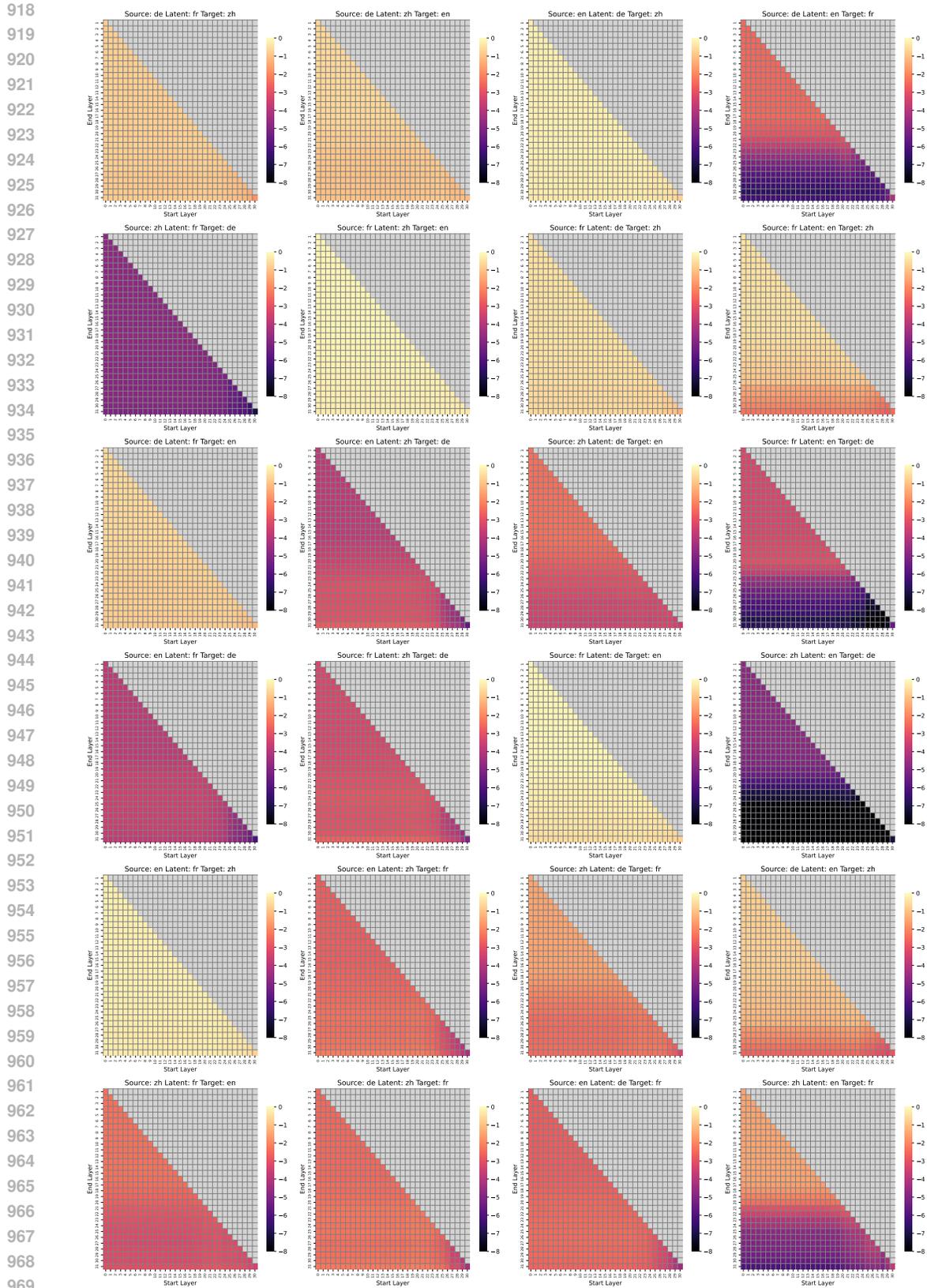


Figure 16: Same plot as Figure 15, but now plotting log-probabilities all over the same range  $-8 \leq \log p \leq 0$ , and with every choice of latent language. Lower values mean a greater rejection effect. Clearly visible is the greater rejection effect when the latent language is English. Columns sorted by latent language.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

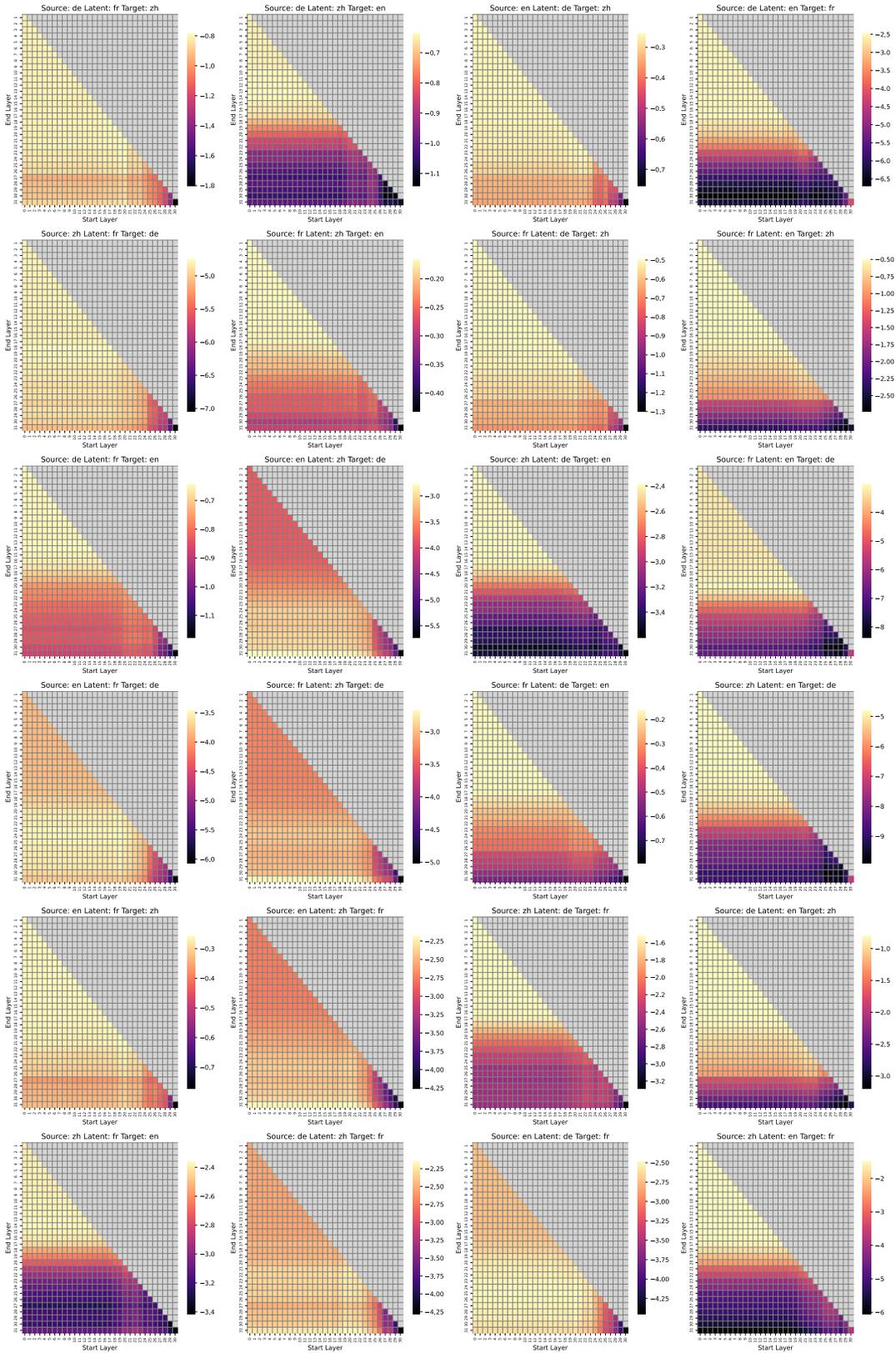


Figure 17: Same plot as Figure 16, but plotting log-probabilities with ranges determined per experiment to make the effect more visible. As before, the range of values is much greater for the English intervention, though interventions in other languages do have a small effect. Columns sorted by latent language.

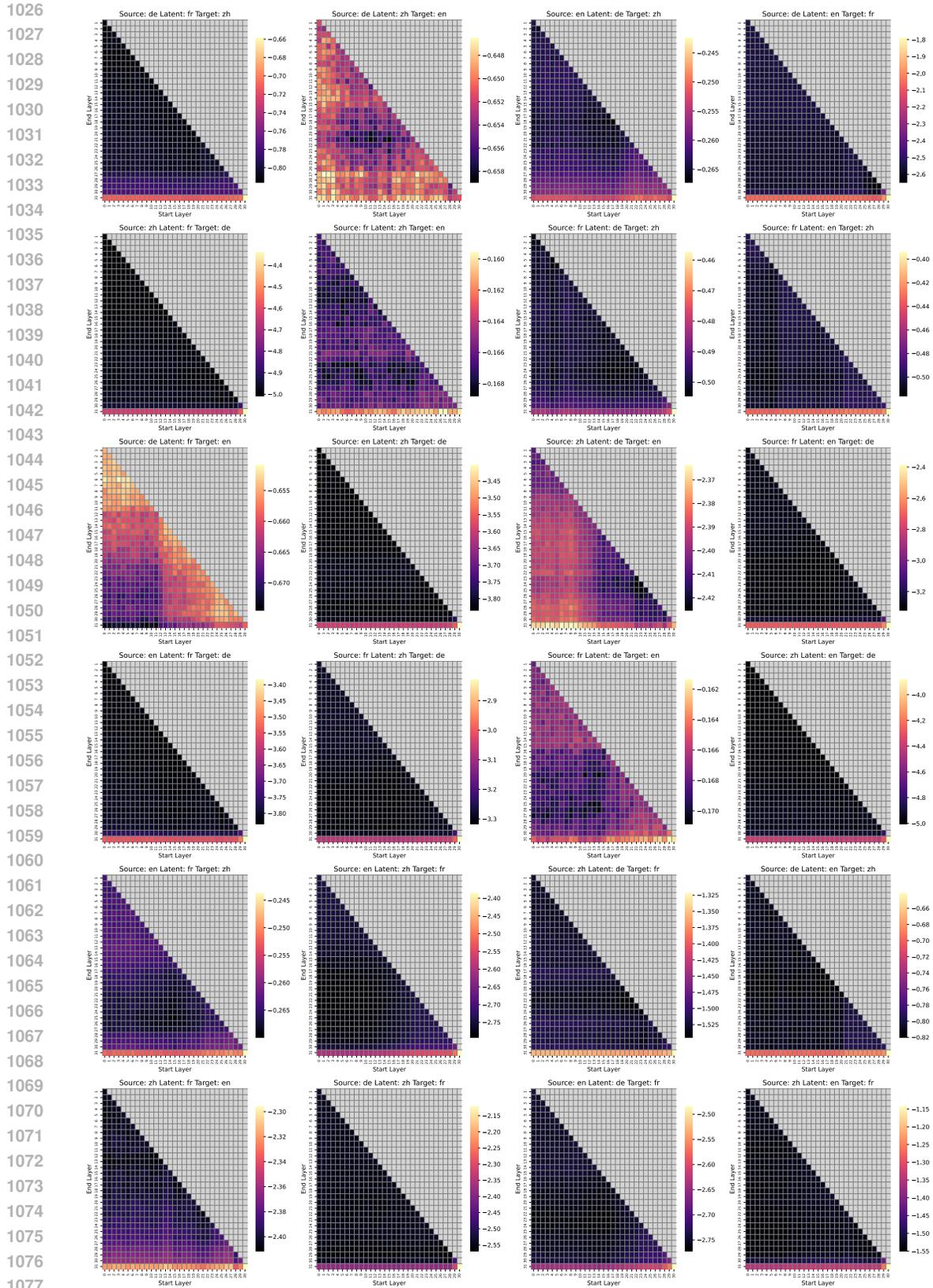


Figure 18: Same plot as Figure 16, but now the intervention is the rejection on the subspace of an unrelated word in the latent language, plotting log-probabilities with ranges determined per experiment to make the effect more visible. This intervention as expected has very little effect on the models prediction.

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

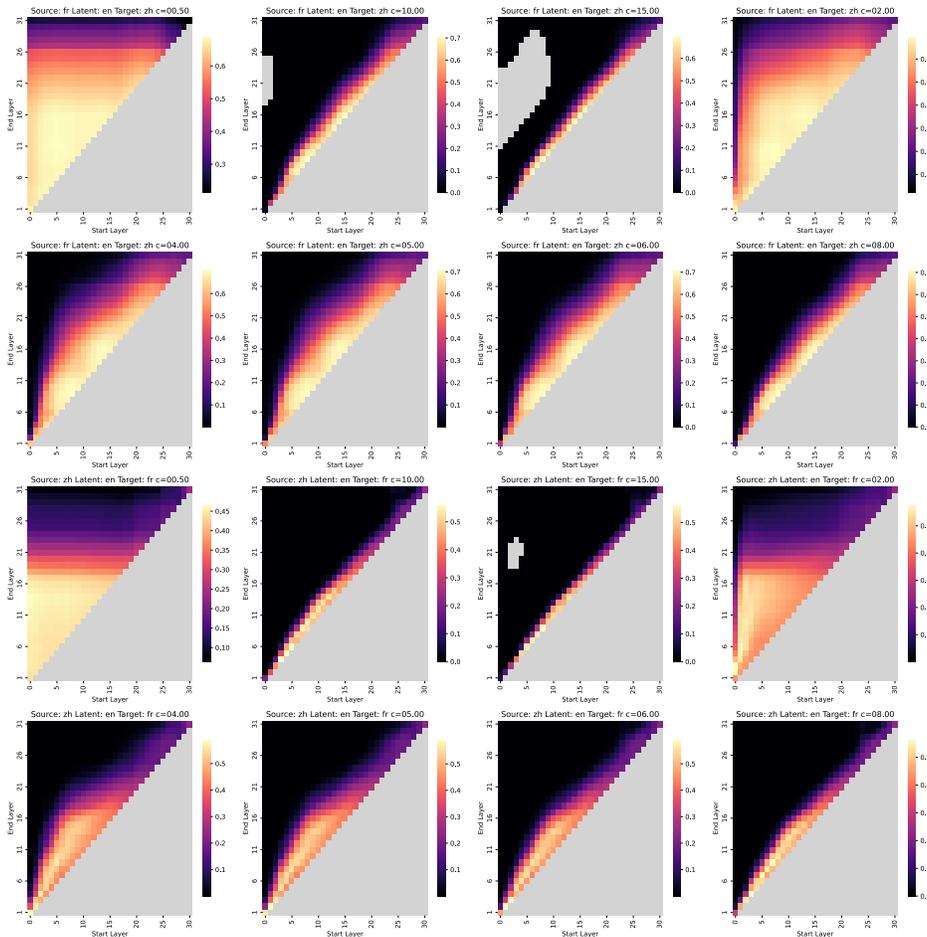


Figure 19: The probability of Llama-2-7b predicting the correct translation, steering towards a different concept in **English**. The intervention was performed across all layers  $\{0, 1, \dots, 31\}$  that lie in the interval (start, end). We sweep over all possible values  $0 \leq \text{start} < \text{end} \leq 31$ , as well as sweep over the choices for source and target and intervention languages. Lower values indicate a greater steering effect, as we want to steer the model away from predicting the correct translation.

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

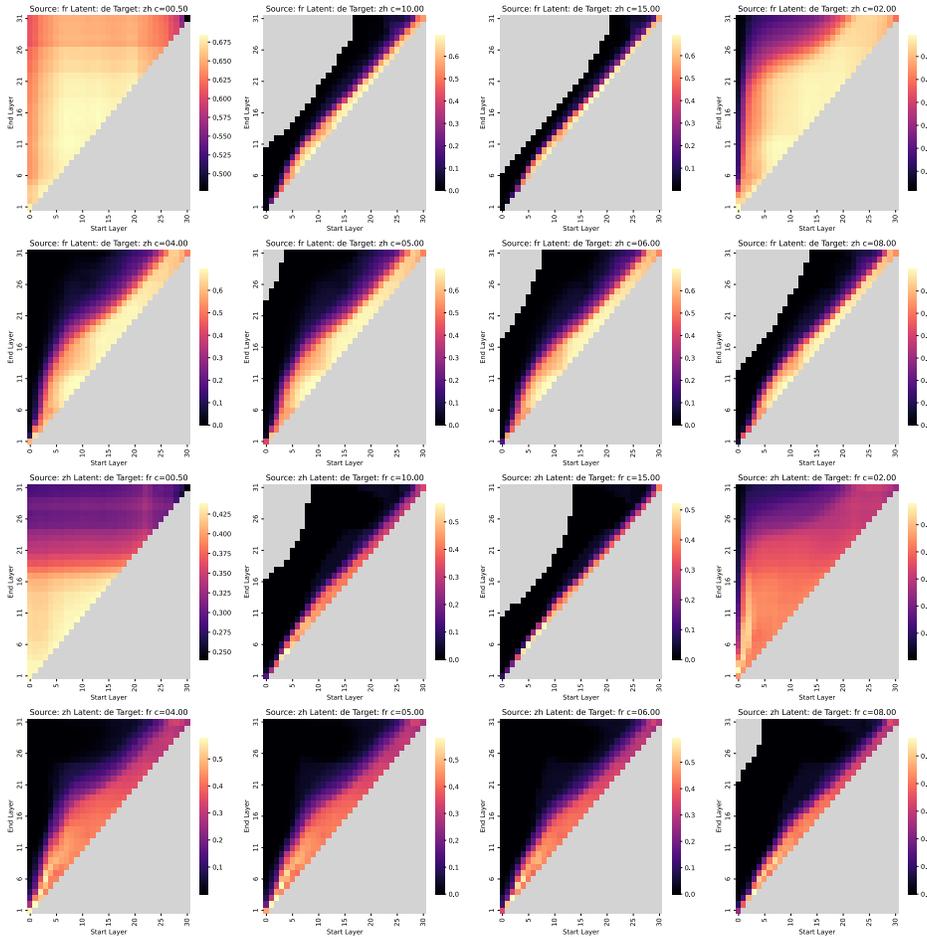


Figure 20: The probability of Llama-2-7b predicting the correct translation, steering towards a different concept in **German**. The intervention was performed across all layers  $\{0, 1, \dots, 31\}$  that lie in the interval (start, end). We sweep over all possible values  $0 \leq \text{start} < \text{end} \leq 31$ , as well as sweep over the choices for source and target and intervention languages. Lower values indicate a greater steering effect, as we want to steer the model away from predicting the correct translation.

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

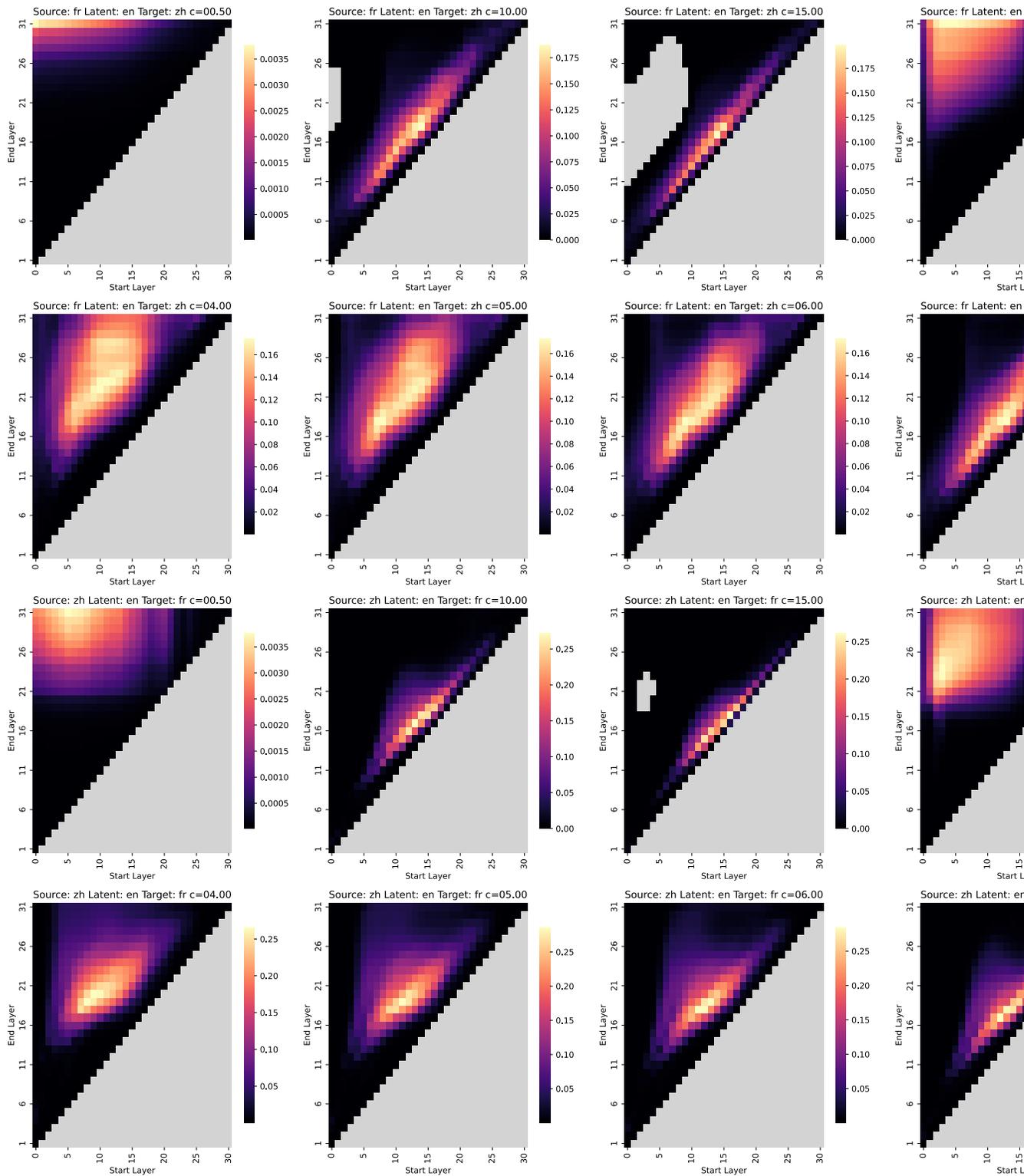


Figure 21: Same plot as Figure 19, but now plotting probabilities of the counterfactual answer we are steering towards, using **English** as the latent language. Higher values mean a greater steering effect.

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

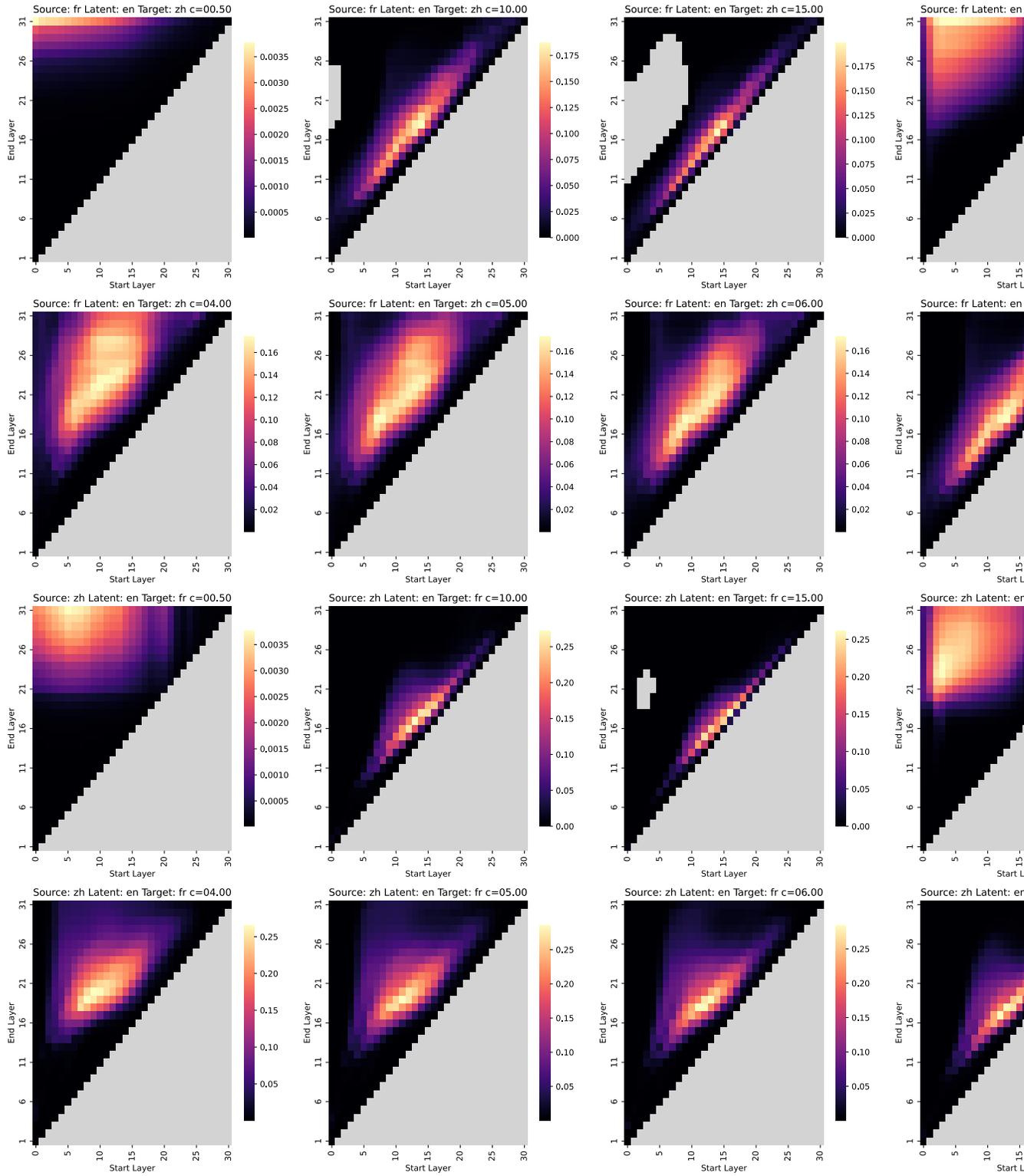


Figure 22: Same plot as Figure 20, but now plotting probabilities of the counterfactual answer we are steering towards, using **German** as the latent language. Higher values mean a greater steering effect.

1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

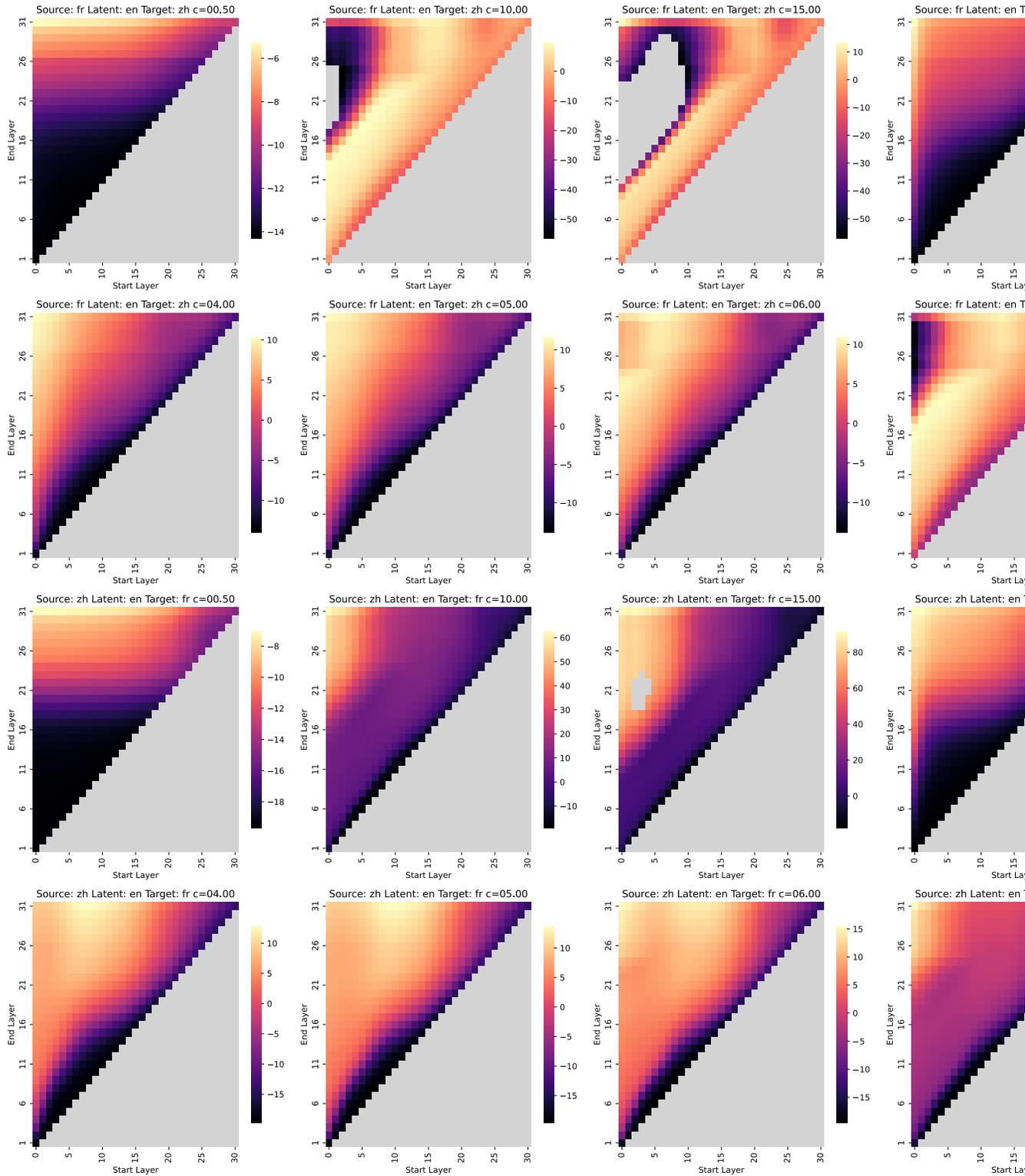


Figure 23: Same plot as Figure 19, but now plotting the log-probabilities of the counterfactual answer we are steering towards, minus the log-probabilities of the correct answer we are steering away from. Intervention is performed with **English** as the latent language. Higher values mean a greater steering effect.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

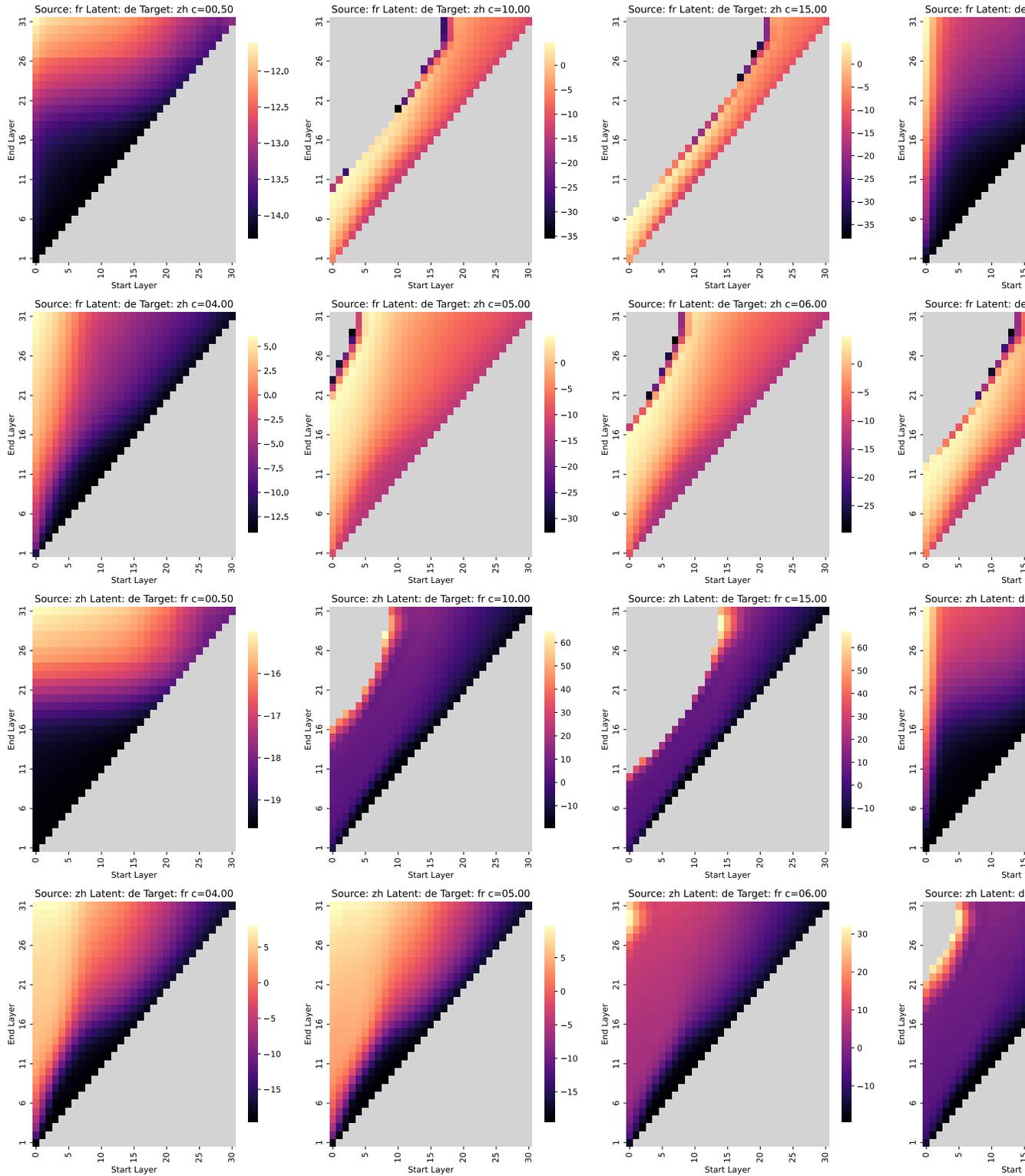


Figure 24: Same plot as Figure 23, but now with **German** as the latent language.