

Localized Cultural Knowledge is Conserved and Controllable in Large Language Models

Veniamin Veselovsky* Berke Argın*² Benedikt Stroebel*¹ Chris Wendler³
Robert West² James Evans^{5,6} Thomas L. Griffiths^{†1,4} Arvind Narayanan^{†1}
¹Princeton University, Department of Computer Science
²EPFL
³Northeastern University
⁴Princeton University, Department of Psychology
⁵University of Chicago, Department of Sociology ⁶Santa Fe Institute

Abstract

Just as humans display language patterns influenced by their native tongue when speaking new languages, LLMs often default to English-centric responses even when generating in other languages. Nevertheless, we observe that local cultural information persists within the models and can be readily activated for cultural customization. We first demonstrate that explicitly providing cultural context in prompts significantly improves the models’ ability to generate culturally localized responses. We term the disparity in model performance with versus without explicit cultural context the *explicit–implicit localization gap*, indicating that while cultural knowledge exists within LLMs, it may not naturally surface in multilingual interactions if cultural context is not explicitly provided. Despite the explicit prompting benefit, however, the answers reduce in diversity and tend toward stereotypes. Second, we identify an explicit cultural customization vector, conserved across all non-English languages we explore, which enables LLMs to be steered from the synthetic English cultural world-model toward each non-English cultural world. Steered responses retain the diversity of implicit prompting and reduce stereotypes to dramatically improve the potential for customization. We discuss the implications of explicit cultural customization for understanding the conservation of alternative cultural world models within LLMs, and their controllable utility for translation, cultural customization, and the possibility of making the explicit implicit through soft control for expanded LLM function and appeal.

1 Introduction

When humans speak a second language, they often display patterns of language use influenced by their

native tongue (Dong et al., 2005; Matsuki et al., 2021). A striking example of this phenomenon can be observed among Japanese speakers learning English. When asked about the color of a pumpkin, they frequently respond “green” (Matsuki et al., 2021), reflecting the concept’s reality in Japan, where pumpkins are predominantly green, rather than the orange variety common in English-speaking countries. Interestingly, these same speakers, when specifically asked about pumpkin colors in America, will correctly identify them as orange.

Large language models (LLMs) exhibit a parallel phenomenon, albeit in reverse. When prompted in Japanese about pumpkin colors, GPT-4o responds with “orange”, revealing what we might call an English pattern of use. When explicitly asked about pumpkin colors in Japan, however, it accurately answers “green” (see Appendix Fig. 4 and Fig. 5). This disparity suggests that cultural knowledge, while present in LLMs, may not naturally surface during multilingual interactions. We dub this phenomenon the *explicit–implicit localization gap*.

The distinction between explicit and implicit prompts reflects two ways models are used. In the explicit setting, we directly measure the capabilities of the model, and in the implicit setting, we measure how real users of, e.g., online chatbots, are likely to prompt the LLM. When an LLM is incapable of adapting its generation depending on the language of the input prompt, then it risks rendering many of its generations “inappropriate” for cultural context (Leibo et al., 2024). Accordingly, understanding how LLMs culturally localize is important, because it will reveal the mechanisms behind how models adapt to different cultures, which can then be used to improve its appropriateness. Quantifying the gap between implicit and explicit localization, on the other hand, contextualizes many existing cultural benchmarks — a lot of which focus on directly measurement capabilities through explicit prompts and thus may not

*Equal contribution. Correspondence to veniamin@princeton.edu¹

[†]Joint senior authors.

transfer to how people actually use them.

Contributions. In this paper, we start by designing a simple cultural localization benchmark and quantify the explicit–implicit localization gap. We find that *implicit prompting*, using the language alone, is insufficient for an LLM to culturally localize its generation and results in a significant degradation in performance across cultural tasks. *Explicit prompting*, also called “cultural prompting” (Tao et al., 2024), results in a significant jump in performance on cultural tasks but comes at the cost of increased homogeneity and stereotypicality in open-ended generation.

Next, using tools from mechanistic interpretability, we explore where cultural localization happens within LLMs and propose a method to overcome the weaknesses of explicit prompting by steering generations to be culturally localized. To this end, we compute linear steering vectors (Rimsky et al., 2024) on the layers that allow us to culturally localize model responses without the need for explicit prompting. Our experiments demonstrate that these cultural steering vectors are effective at culturally localizing generations, as well as being both *language* and *task* agnostic: they remain effective across different languages and can generalize from one localization task (such as identifying a person’s name) to another (such as determining their city of origin). Steered responses also retain the diversity of the implicit prompting setting, reduce stereotypes (oversimplified, fixed beliefs or clichés about a culture) and are more faithful (accurately representative) to the underlying culture. We also find preliminary evidence for a universal cultural localization vector that steers model outputs toward the culture associated with the language of the prompt. This implies that mechanisms underlying cultural localization are also *culture-agnostic*.

As language models become increasingly widely used, understanding their multilingual capabilities and limitations is crucial. Our work introduces the explicit–implicit localization gap as a metric for evaluating these systems, providing both a conceptual framework and practical tools for measuring cultural competence in real-world deployments. By revealing the underlying mechanisms of cultural localization, we not only advance our understanding of these models but also provide concrete paths toward building more culturally aware AI systems that better serve diverse global users. Our findings reveal that cultural localization is not as simple

as explicitly telling a language model the context of the user; instead, we need to think about new approaches for culturally localizing answers.

2 Materials and Methods

First, we define what we mean by language, culture, and task. By language, we refer to the language of the input prompt. Culture represents the beliefs, values, traditions, practices associated with a specific group of people, following the symbolic, discursive approach that has become dominant in modern anthropology (Geertz, 2017). Finally, tasks are various sets of culturally relevant problems.

2.1 Cultural Localization Benchmark

We limit our analysis to English, Turkish, Russian, French, and Bengali — to cover both typographically diverse alongside high- and low-resource languages. To evaluate cultural localization in multilingual language models, we introduce a new benchmark consisting of four datasets: Names, Cities, o1-distilled, and CulturalBench (Chiu et al., 2025). The first three datasets are synthetically generated, while CulturalBench is sourced from prior work. The names dataset consists of paired names, one American and one from the country of the language; similarly, the cities dataset consists of two cities, one American, one from the country of the language we’re evaluating. The o1-distilled dataset comes from few-shot prompting o1-preview (OpenAI, 2024c) to generate culturally relevant questions. For dataset details, see Appendix A.2.

We define cultural localization as the process of tailoring responses so they align with the cultural norms, values, and context implied by the language of the input. For example, if a cultural question is made in Bengali, a culturally localized would answer it in a way that remains faithful to Bengali culture, rather than defaulting to a response rooted in another dominant culture, like American. In the case of asking “what’s a likely name here”, the model should correctly identify that the name “Mohammed” is more likely than “George”. It is worth noting that while language is an important indicator of cultural context, it is an often noisy proxy for a wider phenomena. For this analysis, we will treat language and culture as closely associated by assuming that language embeds within it a world of culturally-specific associations, following recent demonstrative work in large-scale corpus linguistics (Lewis et al., 2023).

2.2 Evaluation

We evaluate cultural localization across the four settings illustrated. For each task, we vary the prompting language between English and Turkish/Russian/French/Bengali as well as whether we include explicit cultural context. Within this setup, we term variants including explicit cultural context in the prompt as *explicitly localized* and model performance on such prompts as their *explicit localization performance*. Variants in which cultural information is only encoded via the language of the prompt we term as *implicitly localized*, and their performance as *implicit localization performance*.

We thus conduct two transformations on each of our data rows. First, we prepend an explicit localization text of the form “I live in [X]”, where [X] is a country that speaks the language. Second, we translate both versions (with and without explicit localization) into the language of study (see Appendix A.2 for further dataset details).

Models. We evaluate seven models: Aya-8b-expanse (Dang et al., 2024) (denoted Aya-8b-it), Gemma2-9b-it, Gemma2-27b-it (Team, 2024), Llama-3.1-70b-it, Llama-3.1-8b-it, Llama-3.1-8b-base (Meta AI, 2025), and GPT-4o (OpenAI, 2024b). With the exception of open-ended generation, we sample at a temperature of 0. For open-ended generation, we use a temperature of 1, top- p of 0.9, and top- k of 50. For the mechanistic interpretability analysis, we study Gemma2-9b-it using nnsight (Fiotto-Kaufman et al., 2025).

Accounting for order bias. The order in which options are presented to a language model has been shown to affect its ability to answer questions (Davidson et al., 2024; Rimsky et al., 2024). Because our benchmark consists of two options for each question, we do two passes over the dataset and average results over possible orderings.

2.3 Analyzing How Localization Occurs

Background on language modeling. When a sentence s is fed into a language model, the text is first tokenized into a sequence of tokens, $x = x_1, \dots, x_t \in V$, where V is the vocabulary of the model. Then when generating the x_{t+1} token, the model, $P : \mathcal{X} \rightarrow \mathcal{Y}$ maps the input sequence to a probability distribution $\mathcal{Y} \in \mathbb{R}^{|V|-1}$. This is done by first embedding each of the input tokens x_i using a learned input embedding matrix E that maps the vocabulary to a set of embeddings denoted

by $h_i^{(0)} = Ex_i$. These embeddings are updated throughout the model as follows:

$$h_i^{(j)} = h_i^{(j-1)} + g^{(j)}(h_1^{(j-1)}, \dots, h_i^{(j-1)})$$

Here $g^{(j)}$ is a function that typically consists of a causal attention layer followed by an MLP layer, alongside normalization layers. Finally, an unembedding matrix W followed by the softmax operation is converts the representation of the last token $h_t^{(L)}$ into a next-token probability distribution

$$P(x) = \frac{e^{Wh_t^{(L)}}}{\sum_{v \in V} e^{(Wh_t^{(L)})_v}} \in \mathbb{R}^{|V|-1}, \quad (1)$$

in which we omitted the dependence of $h_t^{(L)}$ on x .

Activation Steering. We test if we can encourage localization through contrastive activation addition (CAA) (Rimsky et al., 2024). Here, a steering vector is formed by subtracting the mean latent representation of negative prompts from that of positive prompts for the attribute of interest. To build the positive and negative sets, we take the simplest possible contrast. In the explicit setting, a positive example is a hinted prompt where the model picks the culturally localized answer; the negative is the same prompt without the hint, where the model instead defaults to the non-localized choice. In the implicit setting, the positive is a translated prompt (into the target language) that elicits a localized answer; the negative is the English version of that prompt, which yields a non-localized one. We spell out these constructions in the main text. Formally, at layer j , the steering vector shifts the model away from behaviors in D^- and toward those in D^+ , and is defined as:

$$v^{(j)} = \frac{1}{|D^+|} \sum_{x \in D^+} h_t^{(j)}(x) - \frac{1}{|D^-|} \sum_{x \in D^-} h_t^{(j)}(x)$$

where the index t refers to the last token’s position. This steering vector is then added to the generated token representations during inference

$$\tilde{h}_t^{(j)}(x) = h_t^{(j)}(x) + \alpha v^{(j)}.$$

In our case, the target behavior involves aligning the model’s responses with a specific cultural context. Note that the original paper uses a different configuration: prompts for each pair remain identical until the final completion, at which point a token is appended as a hint to desired behavior. By contrast, we pair different versions of the

same question (e.g., translated vs. English, with vs. without cultural context) and omit the final answer. This approach proved more effective and consistent for steering the model toward desired cultural context in our experiments.

2.4 Explicit–Implicit Localization Gap

We define the EI-Gap as the difference in accuracy between the explicit and implicit settings. A response is correct if the model selects the culturally localized option. Formally:

$$\text{EI-Gap} = \mathbb{E}_{(x_q, y)}[\mathbf{1}\{\hat{y}(x_{\text{context}} \circ x_q) = y\}] - \mathbb{E}_{(x_{\text{tr}, q}, y_{\text{tr}})}[\mathbf{1}\{\hat{y}(x_{\text{tr}, q}) = y_{\text{tr}}\}] \quad (2)$$

in which \mathbb{E} denotes an expectation over dataset items, $\mathbf{1}\cdot$ is an indicator for whether the model’s prediction exactly matches the correct culturally localized answer, $\hat{y}(\cdot)$ is the model’s produced output, $x_{\text{context}} \circ x_q$ is the explicitly localized prompt, and $x_{\text{tr}, q}$ is the implicitly localized translated prompt. The first term reflects average correctness under explicit localization, while the second term reflects average correctness under implicit localization.

Open-Ended Generation. To test for the potential downstream effects of explicit localization, we focus on the diversity, stereotypicality, and faithfulness of generations produced. Specifically, we create 24 short story prompts (see Appendix B) and then resample these 30 times for each story. We then calculate the cosine similarity between the embeddings of the generations in a BERTScore-style approach (Zhang* et al., 2020). Specifically, we embed each generation using OpenAI’s text-embeddings-3-small model and measure the cosine similarity across generations from each translated prompt with and without the explicit localization prompt. If the average similarity is higher across answers, we consider those answers to be more semantically similar. We also use LLM-as-a-judge in an arena-style evaluation to rank which generation is more stereotypical and faithful to the underlying culture (see Appendix B.3).

3 Results

3.1 Explicit–Implicit Localization Gap

We begin by quantifying the explicit–implicit localization gap within LLMs by comparing the difference in their performance across implicit and explicit settings. We illustrate the results in Fig. 1

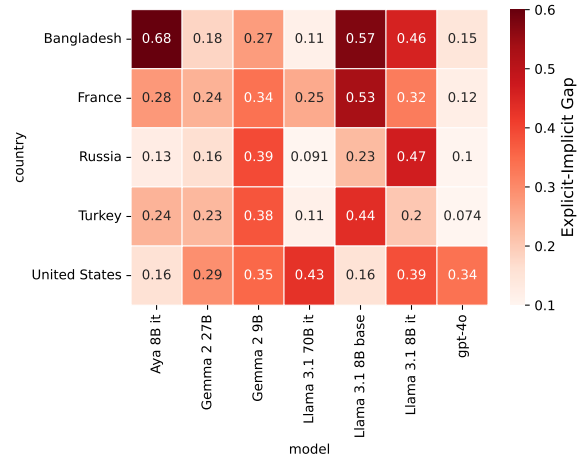


Figure 1: Heatmap showing the explicit–implicit localization gap across models and languages.

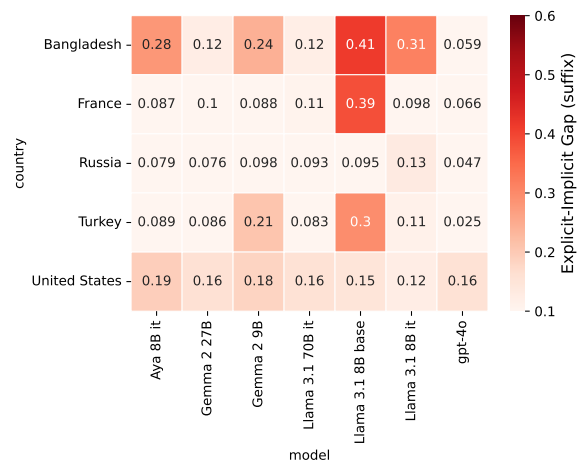


Figure 2: Heatmap showing the explicit–implicit localization gap across models and languages with a culturally relevant prefix prepended.

across language and model. Each cell value consists of the difference between explicit localization (English with cultural context) and implicit localization (local language with no context). We see that across most models, the gap in performance is usually above 10% and sometimes as high as 68% in the case of Bengali and Aya-8b (an explicitly multilingual model). Additionally, it appears that smaller models suffer more from the gap than larger models. The gap in Gemma2-9b-it is 35% whereas in Gemma2-27b-it it is 22%. Similarly, in Llama-8b-it it is 37% whereas in Llama-70b-it it’s 20% (see Appendix C for a subtask breakdown). Another curious result is that all models seem to perform poorly when prompted just in English, as evidenced by the “United States” row — English with USA context vs. just English. We hypothe-

size this is a consequence of the universality of English. In this way, explicitly localizing a prompt to the USA might narrow the output distribution and lead to different responses. This merits further investigation.

Qualitative exploration of failure modes. Our analysis includes cases where models failed to produce valid answers. While most models had failure rates below 1%, GPT-4o exhibited a notably higher rate. This was primarily due to its tendency to refuse answering questions without cultural context, particularly in name and city-related tasks. Specifically, GPT-4o declined to answer 64% of city-related and 32% of name-related questions without explicit cultural context. This cautious behavior suggests the model’s built-in safeguards against potential cultural biases, though these refusals decreased somewhat when prompting in the target language.

Characterizing explicit localization. One natural question that arises is what text is sufficient to explicitly localize the model. Including the background of the user or language they speak seems natural. But what about more subtle approaches like including a concept unique to that culture, e.g., “omakase” for a Japanese speaker. In this section we prepend words from specific cultures and see if these words lead to comparable gains to explicit localization. We test this by including four possible words for each culture: the name of a local dish, currency, city, and cultural object (see Appendix C). In Fig. 2, we show the same difference between explicit and implicit localization as above, only now with the implicit being the correctness with a random prepended cultural word. We observe that across all models and languages, the gap falls dramatically. This implies that having a culturally specific word in the prompt is enough to localize the model to that culture. We also ran experiments with explicit prompting where we modify the explicit prompt to reference language rather than geography (e.g., “I speak Japanese,...”). These results also lead to similar performance as in the explicit setting (see Appendix D).

Effect on Open Generation. A problem with prompting approaches is that they can reduce entropy in the corresponding generation by giving the model context on what to generate (Chu et al., 2025) and unfaithfully simulating ethnic groups (Wang et al., 2025). In this question, we ex-

	Implicit	Explicit
Gemma 2 27B	0.333 \pm 0.008	0.359 \pm 0.006
Gemma 2 9B	0.337 \pm 0.007	0.368 \pm 0.006
Llama 3.1 70B it	0.298 \pm 0.006	0.323 \pm 0.006
Llama 3.1 8B it	0.301 \pm 0.006	0.332 \pm 0.006
GPT-4o	0.317 \pm 0.006	0.324 \pm 0.006

Table 1: Global cosine similarity across generations. Bootstrap standard errors on the means are shown.

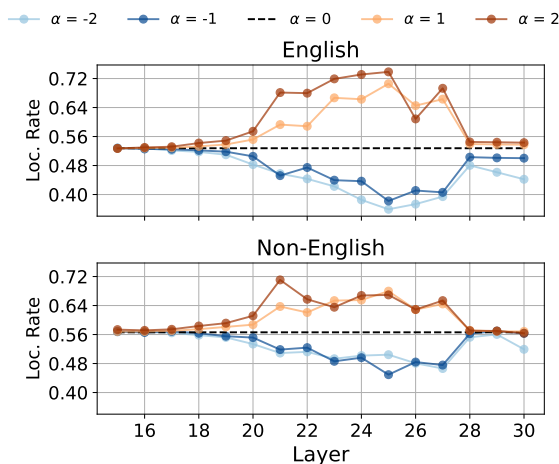


Figure 3: Steering results for per-culture vectors calculated using English pairs with $\alpha \in [-2, 2]$ across layers [15-30], where the horizontal axis represents the layer at which the steering vector is applied, and the vertical axis indicates the ratio of localized responses. Titles denote prompt language.

plore potential downsides with explicit localization: homogeneity and stereotypicality. Specifically, we take open-ended generation prompts that are each translated into the languages, and generate thirty generations for each using a temperature of 1 across explicit and implicit settings. Note that we drop the base models for this task.

In Tab. 1, we show the homogeneity in the implicit and explicit settings. We find that across all models, homogeneity goes up (except for Aya, which was not evaluated). This is most notable for Gemma-2 models, where it increases on average by more than 3%. Additionally, in Fig. 13 we show that explicit prompting results in more stereotypical answers than implicit prompting.

3.2 Steering Localization

We next design a steering experiment by extracting embeddings from two types of prompts: ones with explicit localization (e.g., “I live in Turkey”) and control prompts without cultural context. We conduct all mechanistic analyses on Gemma-2-9b-

	v_{en}	v_{tr}	v_{names}	Implicit	Explicit
cities	0.873 ± 0.021	0.740 ± 0.027	0.896 ± 0.019	0.466 ± 0.032	1.000 ± 0.000
culturebench	0.793 ± 0.042	0.780 ± 0.042	0.798 ± 0.042	0.742 ± 0.044	0.882 ± 0.034
culturedistil	0.614 ± 0.029	0.575 ± 0.031	0.624 ± 0.029	0.484 ± 0.031	0.844 ± 0.023
names	0.734 ± 0.026	0.680 ± 0.028	0.771 ± 0.024	0.552 ± 0.028	0.914 ± 0.016

Table 2: Steering performance across subtasks using v_{names} , v_{en} , and v_{tr} ; details on the α and layers in Appendix D. Each value is the fraction of answers that are culturally localized. Implicit shows the cultural localization rate in the implicit context. Explicit shows the rate when context is directly provided.

it. To extract steering vectors, we use a different prompt structure, detailed in Appendix D. We subtract the control embeddings from their localized counterparts, performing this process in both English and translated versions. When we use all tasks (names, cities, o1-distilled, CulturalBench) we denote these two settings as v_{en} and v_{tr} , respectively. Alternatively, when we steer only using the translated names dataset, we use v_{names} . Using these steering vectors, we run three analyses: (1) How much of the explicit prompting performance can be recovered by adding a single vector? (2) Is the steering vector task agnostic? (3) Are steering vectors culture-specific, or can a steering vector from one culture be used to steer another culture?

Steering performance. In Fig. 3 we show the steering performance across layers and alpha values. The top plot includes results for the English language setting, and in the bottom plot, we show what happens when we apply the English steering vector on other languages. We observe that in both settings, steering around layer 19 begins to increase localization and continues until layer 28. Interestingly, we find that the English steering vector generalizes to other languages. For example, adding the English steering vector to a different language leads to an impressively improved rate of localization. Subtracting the vector leads to a drop in localization with answers more aligned with the United States. We also provide a breakdown of the performance across tasks in Tab. 2.

In this section, we only show results for steering on the explicit prompts. An alternative way to design steering experiments is through translated prompts and their English variants. For example, take the Russian prompt (no explicit localization) and subtract its English counterpart. This approach, however, is less effective (see Appendix D).

Steering across tasks. Given a cultural localization vector from one task, we next show it can be generalized to other tasks. For example, if we cal-

culate the localization for Russian using the names subtask, it can be used to localize in the CulturalBench dataset. To study this question, we extract task-specific cultural localization vectors and apply them to different tasks. In Tab. 2, we show that steering using the names dataset gives comparable results to steering on a specific task. In general, the difference in performance is low, with the exception of cities where we find that the name task leads to a better performance than task-specific steering. These results suggest that cultural localization may generalize across tasks and reflect a broader phenomenon present in large language models.

Open-ended generation. In the previous section, we showed that model steering can be an effective approach for culturally localizing a language model. We now test the effect of steering on open-ended generation. In Tab. 3 we show the homogeneity of generations for Gemma-2-9b-it in implicit, explicit, and steered settings (v_{tr}) across languages. We observe that steering results in more diverse outputs across generations than adding explicit context, but only slightly less diverse than implicit prompting. In Appendix B we show examples of our generations and note that steering causes the model to generate in a format similar to the implicit prompt, but with culturally-relevant information. We also evaluate stereotypicality and faithfulness and find that the steered model is both less stereotypical than the explicit setting, and more faithful to the appropriate cultural context (see Appendix B.3).

Universal culture vector. In our final analysis, we test where cultural localization is universal: does there exist a vector one can add to prompts in any language that automatically localizes it to the language of that prompt? To study this question, we calculate the “universal” steering vector by taking the average of all culture-specific steering vectors, except for the language we are currently evaluating, and then apply the vector to the held-out language and measure its effect on cultural localization. We

	Implicit	Explicit	$v_{tr, \alpha=2}$
Bengali	0.367 \pm 0.006	0.400 \pm 0.005	0.359 \pm 0.005
English	0.285 \pm 0.009	0.339 \pm 0.006	0.286 \pm 0.007
French	0.341 \pm 0.007	0.374 \pm 0.006	0.352 \pm 0.006
Russia	0.265 \pm 0.007	0.312 \pm 0.006	0.283 \pm 0.008
Turkey	0.396 \pm 0.007	0.434 \pm 0.005	0.411 \pm 0.007

Table 3: Homogeneity results for Gemma2 9b it in the implicit, explicit, and steered (v_{tr} , $\alpha = 2$, layer 25). Bootstrap standard errors are shown on the means.

	v_{en}	v_{tr}	v_{names}	$v_{universal (tr)}$	Implicit	Explicit
Bangladesh	0.739 \pm 0.036	0.789 \pm 0.035	0.742 \pm 0.036	0.692 \pm 0.040	0.628 \pm 0.040	0.894 \pm 0.026
France	0.612 \pm 0.041	0.648 \pm 0.040	0.662 \pm 0.038	0.586 \pm 0.040	0.547 \pm 0.045	0.884 \pm 0.028
Russia	0.793 \pm 0.032	0.755 \pm 0.036	0.802 \pm 0.034	0.690 \pm 0.040	0.532 \pm 0.041	0.877 \pm 0.029
Turkey	0.763 \pm 0.030	0.704 \pm 0.033	0.783 \pm 0.030	0.680 \pm 0.036	0.559 \pm 0.037	0.919 \pm 0.020
United States	0.695 \pm 0.028	0.695 \pm 0.028	0.720 \pm 0.026	0.610 \pm 0.028	0.524 \pm 0.029	0.883 \pm 0.019

Table 4: Results for various steering vectors across cultures. Values show the fraction of questions that are correctly culturally localized for each culture (row). $v_{universal (tr)}$ refers to held-out universal steering vector, where we average over all culture-specific translated steering vectors with the exception of the culture.

do this only in the translated context where we take explicitly localized prompts in all the languages (except for the one we’re currently studying) and subtract the implicit prompt.

In Tab. 4, we observe that universal steering does lead to an improvement over implicit prompting alone, but still falls behind explicit prompting. It is also worse than steering on the culture of the prompt alone, as evidenced by the v_{tr} vector. Despite this, the fact that universal steering directionally leads to improved cultural localization suggests the existence of a way to universally steer the model to generate culturally relevant answers.

In the current construction of the task, we limit our analysis to binary questions where one of the answers is always from the United States. In this setting, the cultural steering vector may simply learn to provide the “non-American” answer. To account for this, we extend to a multiple choice setting where each of the options comes from a different culture. We notice a similar improvement from universal steering in this context (see Appendix E). We leave it to future research to determine how best to define a universal steering vector or enable model soft control over language-specific vectors, but here provide an initial support for the existence of automatic cultural customization mechanisms.

Effect on downstream tasks. We also experimented with the effects on downstream tasks like MMLU (Hendrycks et al., 2021) and TruthfulQA (Lin et al., 2022), which indicates that there is a tradeoff: steering reduces MMLU accuracy by

2–4 points and TruthfulQA by 8–11 points (Appendix G). A per-category breakdown of TruthfulQA (Tab. 17) suggests the drop may not be uniform: categories tied to cultural and commonsense knowledge — such as Fiction, Stereotypes, Superstitions, Paranormal, and Misinformation — appear more affected, while broad factual domains like Health, Economics, and Science show very small changes. One possible interpretation is that steering perturbs the model’s English, US-default prior, which could affect culturally specific items more than items grounded in domain facts, though this pattern would need to be studied in more depth to be confirmed.

4 Related Work

Multilingual evaluations in LLMs. Many early multilingual benchmarks like MMLU (Hendrycks et al., 2021), math (Shi et al., 2022), commonsense tasks (Lin et al., 2021; Ponti et al., 2020), factual knowledge (Kassner et al., 2021; Zhou et al., 2022), and representations (Conneau et al., 2018), are based on translations of English benchmarks and thus by construction measure mostly culture-agnostic capabilities. Recently, there have been several attempts to address this shortcoming by creating benchmarks aimed at explicitly measuring cultural capabilities (Chiu et al., 2025; Yin et al., 2022; Naous et al., 2024; Singh et al., 2025). The evaluation methodology in these works corresponds to our explicit setting. Implicit prompting has also been studied (Vayani et al., 2024; Zhang

et al., 2023; Arora et al., 2023). One notable work is by Romanou et al. (2025), where they collect real exam questions from 44 written languages and evaluate capabilities in a local context using these questions. Other work has extended multilingual analysis into reward modeling (Gureja et al., 2025) and multimodal models (Kannen et al., 2024; Khanuja et al., 2025). Tao et al. (2024) showed that LLMs are culturally biased in values and propose using “cultural prompting” to reduce that bias; similarly (AlKhamissi et al., 2024) argued in favor of “anthropological prompting”. There is also active research focusing on knowledge transfer across languages (Goldman et al., 2025; Rajae and Monz, 2024) and answer consistency (Qi et al., 2023). As far as we know, our work is the first attempt to rigorously measure the performance gap between the explicit and implicit setting, and to propose techniques beyond prompting to close it.

Multilingual language models. The growth in LM applications has led to a corresponding increase in research into designing multilingual models (Conneau and Lample, 2019; He et al., 2024; Muenighoff et al., 2023; Üstün et al., 2024). Work has created instruction tuning datasets (Singh et al., 2024) to align multilingual models and studied the impacts of multilinguality on overall model performance (Chang et al., 2023; Schäfer et al., 2024; Gurgurov et al., 2024; Held and Yang, 2022). There has been a large country-level push to train multilingual models (Piir, 2023; Martins et al., 2025; OpenAI, 2024a; Hornyak, 2023; Yue et al., 2025).

Mechanistic interpretability for multilinguality. Techniques like activation patching (Meng et al., 2023; Ghandeharioun et al., 2024; Vig et al., 2020), sparse autoencoders (Huben et al., 2024), and steering (Rimsky et al., 2024) have been used to study linguistic representations. These techniques have shown that linguistic bias seems to originate from a deeper representational problem with the underlying concepts biased towards English (Wendler et al., 2024; Alabi et al., 2024; Wu et al., 2024; Schut et al., 2025). Further, it has been shown that in addition to having concept representations biased towards a specific culture, the circuits¹ the model uses to complete multilingual tasks are often language-agnostic (Lindsey et al., 2025; Wang et al., 2024; Held and Yang, 2022). Other work

¹A circuit is a subnetwork within a neural network explaining most of its performance on a specific task. Tasks can be defined, e.g., via an input-output dataset.

has argued that instead of concepts being biased towards English, they really occupy a universal language-agnostic representation (Brinkmann et al., 2025; Dumas et al., 2025; Variengien and Winsor, 2023). Finally, prior work found that language models have specific neurons responsible for generating in a specific language (Tang et al., 2024).

5 Discussion

We motivated cultural localization through the example of pumpkin color, but this phenomenon stretches far beyond simple concept features. Culture forms the way people think about the world from freedom to religion, from happiness to money (Haerpfer et al., 2020; Hofstede, 1984). Many of these cultural associations lie embedded within language (Lewis et al., 2023; Lakoff, 2008). While our benchmark focused on simple examples (like names or cultural facts) it demonstrates strong evidence that language alone is insufficient to culturally localize language models. As LLMs become used in many different cultural contexts around the globe, the lack of adaptation can lead to unfaithful answers inappropriate for context. Addressing this limitation is key to creating truly global models.

Similar to past work (Tao et al., 2024), we find that providing explicit context helps confront some of these issues. We build on that and show that explicit prompting is not a panacea, however, it comes at the cost of increased homogeneity and stereotypicality. Nevertheless, we discover that LLMs contain the capability for cultural specificity within their multi-layered representations which we locate to mechanically unlock that capacity. We show that an approach focusing on model internals can reconcile this problem, demonstrating how even simple model steering via the addition of a single vector at a single layer is able to culturally localize model answers. While steering in this way did not recover the full cultural localization performance of adding explicit context to the prompt, its cultural localization is more culturally faithful, less stereotypical, and more diverse (see Appendix B). We think that this distinct behavior is a result of the steering intervention being more surgical, mechanistically explicit, and having a more localized effect on the forward pass. We add the steering vector only for newly generated tokens and only at a later layer, leaving the remaining forward pass unchanged. In contrast, explicit prompting modifies the forward pass in its entirety.

This distinction between explicit approaches (providing context directly) and implicit approaches (relying on language alone) highlights a broader conceptual divide in how we evaluate cultural localization. We believe that both explicit and implicit localization are important constructs to measure. But it is critical to know what they involve and represent. Whereas the former predominantly measures capabilities, the latter is ecologically valid (Brunswik, 1940). This difference between evaluations for model developers and those informative for downstream users is often underappreciated (Kapoor et al., 2025), and as we argue in this paper, it is also true for multilingual benchmarks. Explicit localization measures whether we can steer models towards outputting faithful content, but it tells us little about how models behave in real-world settings. Others study this problem purely in an implicit setting, and thus come to conclusions like “Don’t trust ChatGPT when your question is not in English” (Zhang et al., 2023) due to biased representations. While this does indeed appear to be true (Wendler et al., 2024; Dumas et al., 2025), it erroneously implies that this knowledge does not exist within language models.

Limitations

There are several limitations to our study. First, we do not measure the downstream effects of prompting and steering outside of simple model-based approaches (embedding similarity and LLM-as-a-judge). While we include a few anecdotal examples in the appendix, a more thorough exploration should be undertaken, ideally run through a user study with people from these cultures. Second, we do not focus on exploring the many possible strategies for guiding a language model’s generations to relate to a specific culture. Future work should explore: (a) different methods for computing steering vectors like distributed alignment search (Geiger et al., 2024; Minder et al., 2025), affine steering (Marshall et al., 2025), or dictionary learning (Huben et al., 2024) and (b) methods for adapting the model through parameter efficient finetuning methods (Li and Liang, 2021; Dettmers et al., 2023; Hu et al., 2022; Wu et al., 2025; Houlsby et al., 2019). Third, we do not rigorously evaluate the universal steering vector. While we find evidence that it exists, we leave it for future work to rigorously study it and its downstream implications. Fourth, we limit our analysis

to language models and do not consider the current paradigm of reasoning models. Future work should extend the results to both and explicitly trained reasoning models and multimodal models, which have been shown to exhibit poor multilingual performance (Bansal et al., 2022; Kannen et al., 2024). Next, our findings also rely on specific prompt designs, and results may vary with different formulations or personalization mechanisms (Bideman et al., 2024). We only evaluate five languages, leaving open questions about how cultural localization behaves across a broader range of linguistic and cultural settings (Khanuja et al., 2024). Finally, we note that we assume a direct mapping between culture and language. This may not be true and future work should better disentangle their connection in language models.

Reproducibility. We make code and data available: <https://github.com/vminvsky/localization-gap>

Acknowledgments

Chris Wendler is supported by a grant from Open Philanthropy. We would also like to thank Julian Minder for providing feedback on the draft, alongside Christiane Fellbaum and Matthew J. Salganik for providing early guidance that shaped the current direction of this paper. Finally, we thank Lisa Schut, Aleskandra Korolova, Hongchuan Zeng, and Rena Nakashima for helpful discussions about multilinguality in language models.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2024. *GPT-4 technical report*. Preprint, arXiv:2303.08774.
- Jesujoba O Alabi, Marius Mosbach, Matan Eyal, Dietrich Klakow, and Mor Geva. 2024. *The hidden space of transformer language adapters*. Preprint, arXiv:2402.13137.
- Badr AlKhamissi, Muhammad ElNokrashy, Mai Alkhamissi, and Mona Diab. 2024. *Investigating cultural alignment of large language models*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12404–12422, Bangkok, Thailand. Association for Computational Linguistics.
- Arnav Arora, Lucie-aimée Kaffee, and Isabelle Augenstein. 2023. *Probing pre-trained language models for cross-cultural differences in values*. In *Proceedings*

- of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP). Association for Computational Linguistics.
- Hritik Bansal, Da Yin, Masoud Monajatipoor, and Kai-Wei Chang. 2022. [How well can text-to-image generative models understand ethical natural language interventions?](#) In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1358–1370, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Behind the Name. 2025. [Behind the name: The etymology and history of first names](#). Accessed: 2025-01-10.
- Stella Biderman, Hailey Schoelkopf, Lintang Sutawika, Leo Gao, Jonathan Tow, Baber Abbasi, Alham Fikri Aji, Pawan Sasanka Ammanamanchi, Sidney Black, Jordan Clive, Anthony DiPofi, Julen Etxaniz, Benjamin Fattori, Jessica Zosa Forde, Charles Foster, Jeffrey Hsu, Mimansa Jaiswal, Wilson Y. Lee, Haonan Li, Charles Lovering, Niklas Muennighoff, Ellie Pavlick, Jason Phang, Aviya Skowron, Samson Tan, Xiangru Tang, Kevin A. Wang, Genta Indra Winata, François Yvon, and Andy Zou. 2024. [Lessons from the trenches on reproducible evaluation of language models](#). *Preprint*, arXiv:2405.14782.
- Jannik Brinkmann, Chris Wendler, Christian Bartelt, and Aaron Mueller. 2025. [Large language models share representations of latent grammatical concepts across typologically diverse languages](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6131–6150, Albuquerque, New Mexico. Association for Computational Linguistics.
- Egon Brunswik. 1940. Thing constancy as measured by correlation coefficients. *Psychological Review*, 47(1):69.
- Tyler A. Chang, Catherine Arnett, Zhuowen Tu, and Benjamin Bergen. 2023. When is multilinguality a curse? language modeling for 250 high- and low-resource languages. In *Conference on Empirical Methods in Natural Language Processing*.
- Yu Ying Chiu, Liwei Jiang, Bill Yuchen Lin, Chan Young Park, Shuyue Stella Li, Sahithya Ravi, Mehar Bhatia, Maria Antoniak, Yulia Tsvetkov, Vered Shwartz, and Yejin Choi. 2025. [CulturalBench: A robust, diverse and challenging benchmark for measuring LMs’ cultural knowledge through human-AI red-teaming](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 25663–25701, Vienna, Austria. Association for Computational Linguistics.
- KuanChao Chu, Yi-Pei Chen, and Hideki Nakayama. 2025. [Exploring and controlling diversity in LLM-agent conversation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 25626–25644, Suzhou, China. Association for Computational Linguistics.
- Alexis Conneau and Guillaume Lample. 2019. *Cross-lingual language model pretraining*. Proceedings of the 33rd International Conference on Neural Information Processing Systems.
- Alexis Conneau, Guillaume Lample, Ruty Rinott, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. [Xnli: Evaluating cross-lingual sentence representations](#). *Preprint*, arXiv:1809.05053.
- John Dang, Shivalika Singh, Daniel D’souza, Arash Ahmadian, Alejandro Salamanca, Madeline Smith, Aidan Peppin, Sungjin Hong, Manoj Govindassamy, Terrence Zhao, et al. 2024. [Aya expanse: Combining research breakthroughs for a new multilingual frontier](#). *arXiv preprint arXiv:2412.04261*.
- Tim R. Davidson, Viacheslav Surkov, Veniamin Veselovsky, Giuseppe Russo, Robert West, and Caglar Gulcehre. 2024. [Self-recognition in language models](#). *Preprint*, arXiv:2407.06946.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *Advances in neural information processing systems*, 36:10088–10115.
- Yanping Dong, Shichun Gui, and Brian MacWhinney. 2005. Shared and separate meanings in the bilingual mental lexicon. *Bilingualism: Language and Cognition*, 8(3):221–238.
- Clément Dumas, Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2025. [Separating tongue from thought: Activation patching reveals language-agnostic concept representations in transformers](#). *Preprint*, arXiv:2411.08745.
- Jaden Fried Fiotto-Kaufman, Alexander Russell Lof-tus, Eric Todd, Jannik Brinkmann, Koyena Pal, Dmitrii Troitskii, Michael Ripa, Adam Belfki, Can Rager, Caden Juang, Aaron Mueller, Samuel Marks, Arnab Sen Sharma, Francesca Lucchetti, Nikhil Prakash, Carla E. Brodley, Arjun Guha, Jonathan Bell, Byron C Wallace, and David Bau. 2025. [NNsight and NDIF: Democratizing access to open-weight foundation model internals](#). In *The Thirteenth International Conference on Learning Representations*.
- Clifford Geertz. 2017. *The interpretation of cultures*. Basic books.
- Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah D. Goodman. 2024. [Finding alignments between interpretable causal variables and distributed neural representations](#). *Preprint*, arXiv:2303.02536.

- Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. 2024. Patchscopes: a unifying framework for inspecting hidden representations of language models. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org.
- Omer Goldman, Uri Shaham, Dan Malkin, Sivan Eiger, Avinatan Hassidim, Yossi Matias, Joshua Maynez, Adi Mayrav Gilady, Jason Riesa, Shruti Rijhwani, Laura Rimell, Idan Szpektor, Reut Tsarfaty, and Matan Eyal. 2025. [Eclektic: a novel challenge set for evaluation of cross-lingual knowledge transfer](#). Preprint, arXiv:2502.21228.
- Srishti Gureja, Lester James Validad Miranda, Shayekh Bin Islam, Rishabh Maheshwary, Drishti Sharma, Gusti Triandi Winata, Nathan Lambert, Sebastian Ruder, Sara Hooker, and Marzieh Fadaee. 2025. [M-RewardBench: Evaluating reward models in multilingual settings](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 43–58, Vienna, Austria. Association for Computational Linguistics.
- Daniil Gurgurov, Tanja Bäuml, and Tatiana Anikina. 2024. [Multilingual large language models and curse of multilinguality](#).
- Christian Haerper, Ronald Inglehart, Alejandro Moreno, Christian Welzel, Kseniya Kizilova, Jaime Diez-Medrano, Marta Lagos, Pippa Norris, Eduard Ponarin, Björn Puranen, et al. 2020. World values survey: Round seven—country-pooled datafile. *Madrid, Spain & Vienna, Austria: JD Systems Institute & WWSA Secretariat*, 7:2021.
- Yifei He, Alon Benhaim, Barun Patra, Praneetha Vadamanu, Sanchit Ahuja, Parul Chopra, Vishrav Chaudhary, Han Zhao, and Xia Song. 2024. [Scaling laws for multilingual language models](#). Preprint, arXiv:2410.12883.
- William Held and Diyi Yang. 2022. [Shapley head pruning: Identifying and removing interference in multilingual transformers](#). Preprint, arXiv:2210.05709.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Geert Hofstede. 1984. *Culture's consequences: International differences in work-related values*, volume 5. sage.
- Tim Hornyak. 2023. Why Japan Is Building Its Own Version of ChatGPT — scientificamerican.com. Accessed 14-10-2024.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*.
- Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. 2024. Sparse autoencoders find highly interpretable features in language models. In *The Twelfth International Conference on Learning Representations*.
- Nithish Kannen, Arif Ahmad, marco Andreetto, Vinodkumar Prabhakaran, Utsav Prabhu, Adji Bousso Dieng, Pushpak Bhattacharyya, and Shachi Dave. 2024. [Beyond aesthetics: Cultural competence in text-to-image models](#). In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Sayash Kapoor, Benedikt Stroebel, Zachary S Siegel, Nitya Nadgir, and Arvind Narayanan. 2025. [AI agents that matter](#). *Transactions on Machine Learning Research*.
- Nora Kassner, Philipp Dufter, and Hinrich Schütze. 2021. [Multilingual LAMA: Investigating knowledge in multilingual pretrained language models](#). Preprint, arXiv:2102.00894.
- Simran Khanuja, Vivek Iyer, Claire He, and Graham Neubig. 2025. [Towards automatic evaluation for image transcreation](#). Preprint, arXiv:2412.13717.
- Simran Khanuja, Sathyanarayanan Ramamoorthy, Yueqi Song, and Graham Neubig. 2024. [An image speaks a thousand words, but can everyone listen? on image transcreation for cultural relevance](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 10258–10279, Miami, Florida, USA. Association for Computational Linguistics.
- George Lakoff. 2008. *Women, fire, and dangerous things: What categories reveal about the mind*. University of Chicago press.
- Joel Z. Leibo, Alexander Sasha Vezhnevets, Manfred Diaz, John P. Agapiou, William A. Cunningham, Peter Sunehag, Julia Haas, Raphael Koster, Edgar A. Duéñez-Guzmán, William S. Isaac, Georgios Pilioras, Stanley M. Bileschi, Iyad Rahwan, and Simon Osindero. 2024. [A theory of appropriateness with applications to generative artificial intelligence](#). Preprint, arXiv:2412.19010.
- Molly Lewis, Aoife Cahill, Nitin Madnani, and James Evans. 2023. Local similarity and global variability characterize the semantic space of human languages. *Proceedings of the National Academy of Sciences*, 120(51):e2300986120.

- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics.
- Bill Yuchen Lin, Seyeon Lee, Xiaoyang Qiao, and Xiang Ren. 2021. Common sense beyond English: Evaluating and improving multilingual language models for commonsense reasoning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL-IJCNLP 2021)*. To appear.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. [TruthfulQA: Measuring how models mimic human falsehoods](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Jack Lindsey, Wes Gurnee, Emmanuel Ameisen, Brian Chen, Adam Pearce, Nicholas L. Turner, Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar, Adly Templeton, Trenton Bricken, Callum McDougall, Hoagy Cunningham, Thomas Henighan, Adam Jermyn, Andy Jones, Andrew Persic, Zhenyi Qi, T. Ben Thompson, Sam Zimmerman, Kelley Rivoire, Thomas Conerly, Chris Olah, and Joshua Batson. 2025. [On the biology of a large language model](#). *Transformer Circuits Thread*.
- Thomas Marshall, Adam Scherlis, and Nora Belrose. 2025. [Refusal in llms is an affine function](#). *Preprint*, arXiv:2411.09003.
- Pedro Henrique Martins, Patrick Fernandes, João Alves, Nuno M. Guerreiro, Ricardo Rei, Duarte M. Alves, José Pombal, Amin Farajian, Manuel Faysse, Mateusz Klimaszewski, Pierre Colombo, Barry Haddow, José G.C. de Souza, Alexandra Birch, and André F.T. Martins. 2025. [Eurollm: Multilingual language models for europe](#). *Procedia Computer Science*, 255:53–62. Proceedings of the Second EuroHPC user day.
- Eriko Matsuki, Yasushi Hino, and Debra Jared. 2021. Understanding semantic accents in Japanese–English bilinguals: A feature-based approach. *Bilingualism: Language and cognition*, 24(1):137–153.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2023. [Locating and editing factual associations in GPT](#). *Preprint*, arXiv:2202.05262.
- Meta AI. 2025. [Llama 3.1 8b instruct](#). Accessed: 2025-01-10.
- Julian Minder, Kevin Du, Niklas Stoehr, Giovanni Monea, Chris Wendler, Robert West, and Ryan Cotterell. 2025. [Controllable context sensitivity and the knob behind it](#). In *The Thirteenth International Conference on Learning Representations*.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. [Crosslingual generalization through multitask finetuning](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.
- Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2024. [Having beer after prayer? measuring cultural bias in large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16366–16393, Bangkok, Thailand. Association for Computational Linguistics.
- OpenAI. 2024a. Government of iceland. <https://openai.com/index/government-of-iceland/>. [Accessed 13-10-2024].
- OpenAI. 2024b. [GPT-4o system card](#). *Preprint*, arXiv:2410.21276.
- OpenAI. 2024c. [OpenAI o1 system card](#). *Preprint*, arXiv:2412.16720.
- OpenAI. 2025. [Introducing OpenAI o1 preview](#). Accessed: 2025-01-10.
- Rait Piir. 2023. [Finland’s chatgpt equivalent begins to think in estonian as well](#). ERR News.
- Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. [XCOPA: A multilingual dataset for causal commonsense reasoning](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2362–2376.
- Jirui Qi, Raquel Fernández, and Arianna Bisazza. 2023. [Cross-lingual consistency of factual knowledge in multilingual language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10650–10666, Singapore. Association for Computational Linguistics.
- Sara Rajae and Christof Monz. 2024. [Analyzing the evaluation of cross-lingual knowledge transfer in multilingual language models](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2895–2914, St. Julian’s, Malta. Association for Computational Linguistics.
- Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. 2024. [Steering llama 2 via contrastive activation addition](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15504–15522, Bangkok, Thailand. Association for Computational Linguistics.

- Angelika Romanou, Negar Foroutan, Anna Sotnikova, Sree Harsha Nelaturu, Shivalika Singh, Rishabh Maheshwary, Micol Altomare, Zeming Chen, Mohamed A. Haggag, Snegha A, Alfonso Amayuelas, Azril Hafizi Amirudin, Danylo Boiko, Michael Chang, Jenny Chim, Gal Cohen, Aditya Kumar Dalmia, Abraham Diress, Sharad Duwal, Daniil Dzenhaliou, Daniel Fernando Erazo Florez, Fabian Farestam, Joseph Marvin Imperial, Shayekh Bin Islam, Perttu Isotalo, Maral Jabbarishiviari, Börje F. Karlsson, Eldar Khalilov, Christopher Klamm, Fajri Koto, Dominik Krzemiński, Gabriel Adriano de Melo, Syrielle Montariol, Yiyang Nan, Joel Niklaus, Jekaterina Novikova, Johan Samir Obando Ceron, Debjit Paul, Esther Ploeger, Jebish Purbey, Swati Rajwal, Selvan Sunitha Ravi, Sara Rydell, Roshan Santhosh, Drishti Sharma, Marjana Prifti Skenduli, Arshia Soltani Moakhar, Bardia soltani moakhar, Ayush Kumar Tarun, Azmine Tushik Wasi, Thenuka Ovin Weerasinghe, Serhan Yilmaz, Mike Zhang, Imanol Schlag, Marzieh Fadaee, Sara Hooker, and Antoine Bosselut. 2025. [INCLUDE: Evaluating multilingual language understanding with regional knowledge](#). In *The Thirteenth International Conference on Learning Representations*.
- Lisa Schut, Yarin Gal, and Sebastian Farquhar. 2025. [Do multilingual LLMs think in english?](#) In *ICLR 2025 Workshop on Building Trust in Language Models and Applications*.
- Anton Schäfer, Shauli Ravfogel, Thomas Hofmann, Tiago Pimentel, and Imanol Schlag. 2024. [The role of language imbalance in cross-lingual generalisation: Insights from cloned language experiments](#). *Preprint*, arXiv:2404.07982.
- Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2022. [Language models are multilingual chain-of-thought reasoners](#).
- Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David Ifeoluwa Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Raymond Ng, Shayne Longpre, Sebastian Ruder, Wei-Yin Ko, Antoine Bosselut, Alice Oh, Andre Martins, Leshem Choshen, Daphne Ippolito, Enzo Ferrante, Marzieh Fadaee, Beyza Ermis, and Sara Hooker. 2025. [Global MMLU: Understanding and addressing cultural and linguistic biases in multilingual evaluation](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 18761–18799, Vienna, Austria. Association for Computational Linguistics.
- Shivalika Singh, Freddie Vargus, Daniel Dsouza, Börje F. Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividias Matciunas, Laura OMahony, Mike Zhang, Ramith Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Souza Moura, Dominik Krzemiński, Hakimeh Fadaei, Irem Ergun, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Minh Chien Vu, Sebastian Ruder, Surya Guthikonda, Emad A. Alghamdi, Sebastian Gehrmann, Niklas Muennighoff, Max Bartolo, Julia Kreutzer, A. Ustun, Marzieh Fadaee, and Sara Hooker. 2024. [Aya dataset: An open-access collection for multilingual instruction tuning](#). In *Annual Meeting of the Association for Computational Linguistics*.
- Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. 2024. [Language-specific neurons: The key to multilingual capabilities in large language models](#).
- Yan Tao, Olga Viberg, Ryan S Baker, and René F Kizilcec. 2024. [Cultural bias and cultural alignment of large language models](#). *PNAS Nexus*, 3(9):346.
- Gemma Team. 2024. [Gemma](#).
- Alexandre Variengien and Eric Winsor. 2023. [Look before you leap: A universal emergent decomposition of retrieval tasks in language models](#). *Preprint*, arXiv:2312.10091.
- Ashmal Vayani, Dinura Dissanayake, Hasindri Watawana, Noor Ahsan, Nevasini Sasikumar, Omark Thawakar, Henok Biadgign Ademtew, Yahya Hmaiti, Amandeep Kumar, Kartik Kuckreja, et al. 2024. [All languages matter: Evaluating LLMs on culturally diverse 100 languages](#). *Preprint*, arXiv:2411.16508.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. [Investigating gender bias in language models using causal mediation analysis](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 12388–12401.
- Angelina Wang, Jamie Morgenstern, and John P Dickerson. 2025. [Large language models cannot replace human participants because they cannot portray identity groups](#). *Preprint*, arXiv:2402.01908.
- Weixuan Wang, Barry Haddow, Minghao Wu, Wei Peng, and Alexandra Birch. 2024. [Sharing matters: Analysing neurons across languages and tasks in LLMs](#). *Preprint*, arXiv:2406.09265.
- Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2024. [Do Llamas work in English? On the latent language of multilingual transformers](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15366–15394.
- Zhaofeng Wu, Xinyan Velocity Yu, Dani Yogatama, Jiasen Lu, and Yoon Kim. 2024. [The semantic hub hypothesis: Language models share semantic representations across languages and modalities](#). *Preprint*, arXiv:2411.04986.
- Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D Manning, and Christopher Potts. 2025. [Reft: Representation](#)

finetuning for language models. *Advances in Neural Information Processing Systems*, 37:63908–63962.

Da Yin, Hritik Bansal, Masoud Monajatipoor, Lianian Harold Li, and Kai-Wei Chang. 2022. [Geolama: Geo-diverse commonsense probing on multilingual pre-trained language models](#). *Preprint*, arXiv:2205.12247.

Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim, Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kantharuban, Lintang Sutawika, Sathyanarayanan Ramamoorthy, and Graham Neubig. 2025. [Pangea: A fully open multilingual multimodal LLM for 39 languages](#). *Preprint*, arXiv:2410.16153.

Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

Xiang Zhang, Senyu Li, Bradley Hauer, Ning Shi, and Grzegorz Kondrak. 2023. [Don’t trust ChatGPT when your question is not in English: A study of multilingual abilities and types of LLMs](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7915–7927, Singapore. Association for Computational Linguistics.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging LLM-as-a-judge with MT-bench and chatbot arena](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Wenxuan Zhou, Fangyu Liu, Ivan Vulić, Nigel Collier, and Muhao Chen. 2022. [Prix-LM: Pretraining for multilingual knowledge base construction](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, pages 5412–5424.

Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. 2024. [Aya model: An instruction finetuned open-access multilingual language model](#). *Preprint*, arXiv:2402.07827.

A Appendix

A.1 Motivating Example

In the introduction, we share an example of how when you ask GPT-4o what color a pumpkin is in Japanese it states orange; whereas, when explicitly asked it says green. We include those examples here.

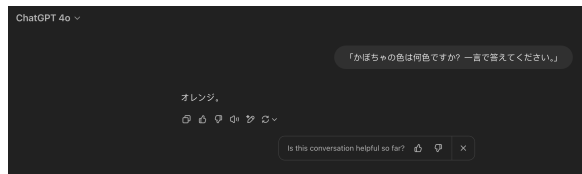


Figure 4: When prompted with “What color is a pumpkin. Please answer in one word.” Model answers orange.

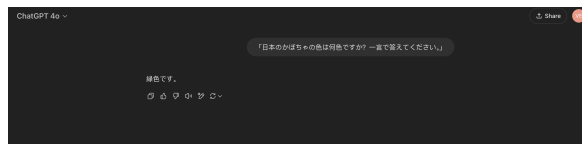


Figure 5: When prompted with “What color is a pumpkin in Japan. Please answer in one word.” Model answers green.

A.2 Cultural Localization Benchmark

The benchmark is constructed from four subtasks, each targeting a different facet of cultural localization: personal identifiers such as names, geographic grounding through cities, open-ended cultural knowledge distilled from a stronger model, and curated cultural reasoning questions drawn from prior work.

- Names:** We selected the 100 most common male and female names for each country and paired names of the same gender with a maximum difference of 2 in their popularity rankings. We then posed the question: “Which of the following names is most likely to be my name?” The names were sourced from the *Behind the Name* website, which reportedly bases its data on official statistics ([Behind the Name, 2025](#)). Note that the website does not contain entries from Bangladesh, so we use GPT-4o to generate these names.
- Cities:** We selected the 20 most populated

Subtask	BD	FR	RU	TR	US	Total
Cities	540	700	700	700	740	3380
CulturalBench	300	156	336	456	120	1368
CultureDistil	624	616	632	808	692	3372
Names	848	848	552	836	770	3854
Total	2312	2320	2220	2800	2322	11974

Table 5: Breakdown of the Cultural Localization Benchmark dataset (transposed). Rows show each subtask, while columns show the countries. The final row and column indicate totals across columns and rows, respectively.

cities (where applicable) in each country and formulated the question: "In which city am I most likely to live?" The list of cities was obtained from the respective country's Wikipedia page.

3. **o1 Distillation:** We provided a few demonstration examples related to food, drinks, traditions, and laws specific to a given country. Using these examples, we prompted the OpenAI o1_preview model to generate 50 additional questions of a similar nature (OpenAI, 2025).
4. **CulturalBench:** We use the dataset from (Chiu et al., 2025) with the Easy setting, where we filter out questions with multiple correct answers and binarize the remaining ones by pairing each correct answer with three incorrect options.

Example prompts from datasets in Fig. 6.

Example prompts from different datasets:

Names Dataset: Which of the following names is most likely to be my name?

Options:

1. Sergey
2. Bill

Cities Dataset: In which city do I most likely live?

Options:

1. Chicago
2. Istanbul

o1-distilled Dataset: Around here, what is the most common alcoholic drink?

Options:

1. Rakı
2. Whiskey

CulturalBench Dataset: Which festival celebration style is most likely appropriate for expressing the collective cultural identity?

Options:

1. Hosting a potluck where each individual brings a dish from their own ethnic background.
2. Celebrating only the most widely recognized national holiday.

Figure 6: Example prompts from the *Cities*, *o1-distilled*, and *CulturalBench* datasets.

A.3 Translation and country removal

We translate all prompts using GPT-4o (Achiam et al., 2024) and the following prompt:

Translate the following text from {src_lang} to {tar_lang}. Respond with just the translation, no additional information.

TRANSLATE: {text}

Figure 7: Prompt used for translation.

For CulturalBench, we remove all mentions of countries using the following prompt also using GPT-4o:

When prompting the base model, we adopt a different few-shot prompt style:

From the following text, remove any explicit mentions of a country.

For example: "In India, what is a common meal to have for lunch?" will become "What is a common meal for lunch?"

Only respond with the rewriting.

TEXT: {text}

Figure 8: Prompt used for removing mentions of countries from CulturalBench.

Question:
What am I?
Options:
1. Human
2. Robot
Answer: 1

Question:
Which planet is closest to the sun?
Options:
1. Mercury
2. Venus
Answer: 2

Question:
I live in Turkey. Which of the following names is most likely to be my name?
Options:
1. Liam
2. Alparslan
Answer:

Figure 9: Prompt used for generation from base model.

B Open-Ended Generation

B.1 Prompts

Here we include the prompts used in the open-ended generation task.

- Write me a recipe for a local dish.
- Write me a short story about a boy and his kite.
- Write me a 100-word bedtime story for kids.
- Tell me about what happened during the Titanic.
- Describe a day in the life of an astronaut on Mars.
- Explain photosynthesis as if I'm five years old.
- Write a letter from a pirate to his long-lost friend.
- Invent a new holiday and describe how people celebrate it.
- Tell me a joke that would make a robot laugh.
- Describe the feeling of standing at the edge of a cliff.
- Write a poem about a lonely lighthouse.
- Explain gravity without using scientific jargon.
- Create a dialogue between a cat and a dog arguing about dinner.
- Write a product review for an imaginary gadget.
- Describe a futuristic city 500 years from now.
- Tell me a legend about a magical forest.

- Explain how to build a sandcastle like a pro.
- Write a diary entry from the perspective of a dragon.
- Imagine you're a time traveler—describe your first day in the past.
- Give me instructions on how to be invisible for a day.
- Write a letter from Earth to an alien civilization.
- Describe a sunset without using the words 'red,' 'orange,' or 'yellow.'
- Tell me about a secret hidden inside an old library.
- Invent a sport that could be played on the moon.

B.2 Additional Details on Homogeneity Results

In this section we include homogeneity results by model and language. Refer to Tab. 6.

B.3 Additional Details on Stereotypicality and Faithfulness Results

To evaluate the stereotypicality and faithfulness of model outputs, we used LLM-as-a-judge (Zheng et al., 2023) with GPT-4o in an arena-style setting. Recall that in the open-ended generation, we had each model answer a prompt 30 times across the different languages in both the implicit and explicit setting. For the arena evaluation we looped over each story and sampled 15 generations in the explicit, implicit, and steered settings. We then took pairwise comparisons between explicit-implicit and explicit-steered and ask GPT-4o to rate which one is more stereotypical and faithful to the culture. Fig. 10 shows the stereotypicality prompt, and Fig. 11 shows the faithfulness prompt.

You will be given two text excerpts from different sources. I would like you to compare them based on how stereotypical they are of culture culture. Stereotypical means that a text relies on oversimplified, fixed beliefs or clichés about a culture. It reduces the complexity of a cultural group to a small number of traits—often leading to misrepresentation or reinforcing biases.

Text A: {text_a}
Text B: {text_b}

Which text relies more on cultural stereotypes about {CULTURE} people or culture? Answer with just 'A' if Text A is more stereotypical, 'B' if Text B is more stereotypical, or 'TIE' if they are equally stereotypical.

Figure 10: LLM-as-a-judge arena-style prompt for rating the stereotypicality of generations.

Tab. 7 shows the results for stereotypicality. Across all languages with the exception of French, the steered generation generates less stereotypical results. We also run a cultural faithfulness evaluation and show the results in Tab. 8. Observe that

Language	Llama-3.1-70B		Llama-3.1-8B		GPT-4o	
	Implicit	Explicit	Implicit	Explicit	Implicit	Explicit
bn	0.303 ± 0.004	0.324 ± 0.004	0.295 ± 0.004	0.328 ± 0.005	0.426 ± 0.005	0.435 ± 0.005
en	0.235 ± 0.007	0.240 ± 0.006	0.243 ± 0.008	0.246 ± 0.007	0.240 ± 0.007	0.236 ± 0.007
fr	0.304 ± 0.006	0.326 ± 0.006	0.308 ± 0.006	0.334 ± 0.006	0.308 ± 0.007	0.312 ± 0.007
ru	0.231 ± 0.006	0.257 ± 0.006	0.247 ± 0.006	0.285 ± 0.007	0.235 ± 0.007	0.251 ± 0.006
tr	0.386 ± 0.008	0.436 ± 0.008	0.379 ± 0.006	0.439 ± 0.006	0.364 ± 0.006	0.377 ± 0.006

Language	Gemma 2 27B		Gemma 2 9B	
	Implicit	Explicit	Implicit	Explicit
bn	0.348 ± 0.009	0.404 ± 0.006	0.368 ± 0.006	0.402 ± 0.005
en	0.286 ± 0.009	0.286 ± 0.006	0.286 ± 0.008	0.297 ± 0.006
fr	0.339 ± 0.007	0.360 ± 0.005	0.341 ± 0.007	0.375 ± 0.006
ru	0.274 ± 0.008	0.296 ± 0.007	0.266 ± 0.007	0.312 ± 0.007
tr	0.386 ± 0.007	0.417 ± 0.006	0.398 ± 0.007	0.435 ± 0.005

Table 6: Homogeneity scores for different models across languages.

You will be given two text excerpts from different sources. I would like you to compare them based on how faithful they are to the culture. Faithful means that the text represents cultural practices, beliefs, or values in a nuanced, accurate, and respectful manner. It acknowledges internal diversity and context, avoids homogenizing or flattening a group’s identity, and strives for factual correctness.
Text A: {text_a}
Text B: {text_b}
Which text is more faithful to the {CULTURE} culture? Answer with just 'A' if Text A is more faithful, 'B' if Text B is more faithful, or 'TIE' if they are equally faithful.

Figure 11: LLM-as-a-judge arena-style prompt for rating the faithfulness of generations.

You will be given two text excerpts. I would like you to compare them based on their fluency. Fluency means that the text is written in natural, grammatically correct language with coherent sentence structure and smooth flow. A fluent text is easy to read and understand.
Text A: {text_a}
Text B: {text_b}
Which text is more fluent? Answer with:
- 'A' (if Text A is more fluent)
- 'B' (if Text B is more fluent)
- 'TIE' (if they are equally fluent)
Do only answer with the letter, no other text.

Figure 12: LLM-as-a-judge arena-style prompt for rating the fluency of generations.

across all languages with the exception of Turkish, the steering provides more faithful generations. Fig. 13 shows the fraction GPT-4o says a specific setting is more stereotypical.

	Implicit	Explicit	v_{tr}
bn	0.011 ± 0.020	0.125 ± 0.067	0.114 ± 0.064
ru	0.008 ± 0.017	0.158 ± 0.078	0.089 ± 0.058
fr	0.006 ± 0.014	0.192 ± 0.086	0.267 ± 0.089
en	0.008 ± 0.019	0.403 ± 0.103	0.156 ± 0.078
tr	0.003 ± 0.008	0.222 ± 0.083	0.103 ± 0.058

Table 7: Stereotypicality scores for different models across languages for Gemma 2 9B. Steered using layer 25, $\alpha = 2$.

	Implicit	Explicit	v_{tr}
bn	0.117 ± 0.067	0.183 ± 0.081	0.456 ± 0.100
ru	0.083 ± 0.056	0.206 ± 0.081	0.250 ± 0.092
fr	0.031 ± 0.036	0.169 ± 0.078	0.306 ± 0.097
en	0.039 ± 0.036	0.147 ± 0.072	0.164 ± 0.078
tr	0.083 ± 0.056	0.267 ± 0.092	0.200 ± 0.078

Table 8: Faithfulness scores for different models across languages for Gemma 2 9B. Steered using layer 25, $\alpha = 2$.

	Explicit	Steering
bn	0.108 ± 0.016	0.886 ± 0.017
ru	0.269 ± 0.023	0.667 ± 0.026
fr	0.361 ± 0.023	0.492 ± 0.025
en	0.283 ± 0.024	0.508 ± 0.026
tr	0.269 ± 0.023	0.706 ± 0.023

Table 9: Fluency comparison. Win rates for explicit vs. steering for Gemma 2 9B across languages.

B.4 Sample generations

Below we include sample generations. In Fig. 16, we show apply an American steering vector affects generation for a prompt about a recipe. Afterwards, we show what happens when we apply each cultures steering vector onto the same English prompt — illustrating that steering is possible cross lan-

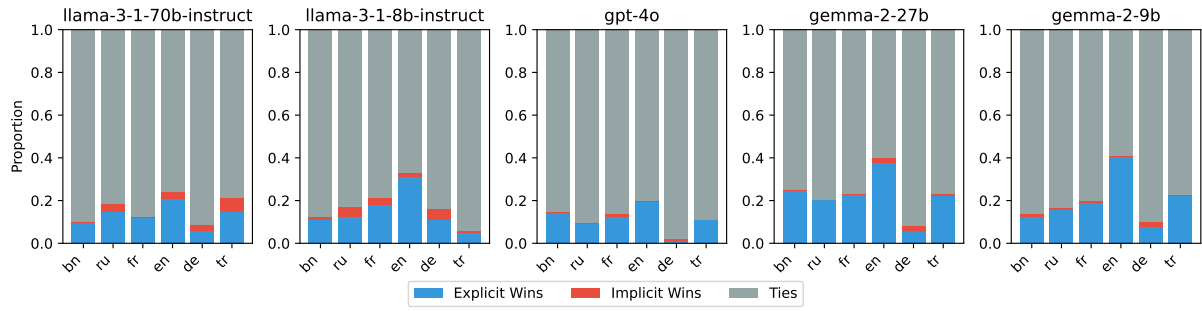


Figure 13: Stereotypicality win rate between the explicit and implicit. Winning here means that the models generation is more stereotypical.

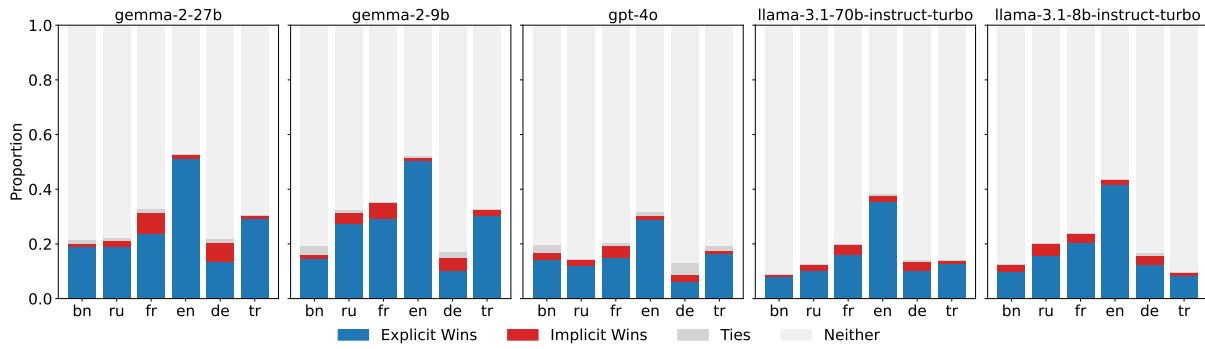


Figure 14: Stereotypicality win rate between the explicit and implicit judged by GPT-OSS-120B with reasoning effort set to medium.

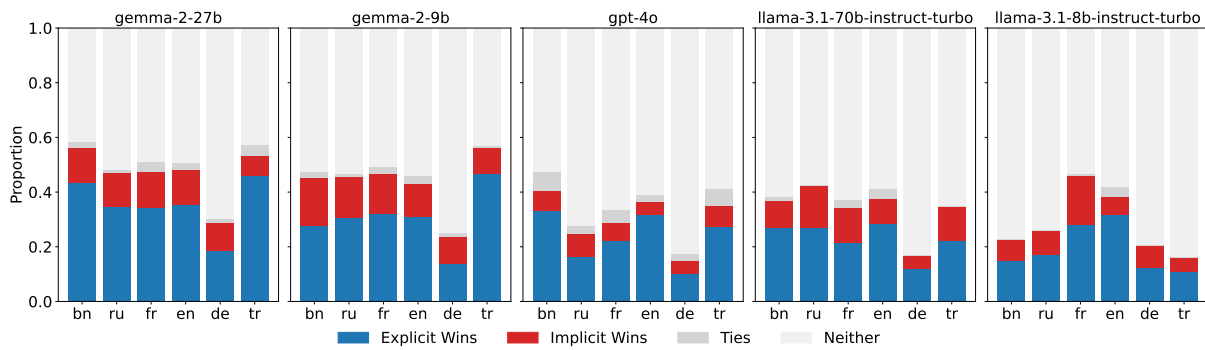


Figure 15: Faithfulness win rate between the explicit and implicit judged by GPT-OSS-120B with reasoning effort set to medium.

guage. A few observations. First of all, the steered generations have the same syntactic structure to the implicit generation prompt, using similar generation structure. Second, the steering vector doesn't force the generation to only generate for a specific context.

Explicit: Alright, let me whip up something delicious for you, something distinctly American, but with a little southern twist! We're gonna make some **Shrimp and Grits**. It's a classic Lowcountry dish hailing from the coastal areas of South Carolina and Georgia. **Here's what you'll need:**

For the Grits:

- * 1 cup stone-ground grits
- * 4 cups water
- * 1 teaspoon salt
- * 4 tablespoons unsalted butter
- * 1/2 cup grated sharp cheddar cheese

For the Shrimp:

- * 1 pound large shrimp, peeled and deveined
- * 1/4 cup butter
- * 1/4 cup chopped onion
- * 2 cloves garlic, minced

Implicit: ## A Taste of Tuscany: Ribollita
Ribollita, a hearty Tuscan bread soup, is perfect for a comforting meal on a chilly evening.

Ingredients:

- * 1 tbsp olive oil
- * 1 large onion, chopped
- * 2 carrots, chopped
- * 2 celery stalks, chopped
- * 2 cloves garlic, minced
- * 8 cups vegetable broth
- * 1 (14.5 oz) can diced tomatoes, undrained

Steered: ## Texas Chili (Classic Texas Style)
This recipe is for a true Texas chili, made with beef, chili peppers, and minimal additions. It's hearty, flavorful, and perfect for a chilly day.

Ingredients:

- * 2 pounds ground beef
- * 1 large onion, chopped
- * 2 bell peppers (red or green), chopped
- * 1 teaspoon chili powder
- * 1 teaspoon cumin
- * 1 teaspoon garlic powder
- * 1/2 teaspoon cayenne pepper (or more, to taste)
- * 1 (28 ounce) can crushed tomatoes
- * 1 (15 ounce) can kidney beans, drained and rinsed
- * Salt and pepper to taste

Figure 16: Example generations in the explicit, implicit, and steered setting for the prompt "Write me a recipe for a local dish." This is done with the American steering vector and explicit prompt.

Turkey: "Let's describe the **Dabke**, a traditional dance from the Levant region of the Middle East, including countries like Lebanon, Syria, Jordan, and Palestine.

Here's a glimpse of the Dabke[...]

France: Let's describe the **Cancan**, a quintessential French dance!

The Cancan:

Origin: Born in the 1830s in Paris, it was a staple of the Parisian cabaret scene at the Moulin Rouge[...]

Russia: Let's describe the **Khorovod**, a traditional Russian dance:

The Dance:

Formation: The Khorovod is a circle dance, with dancers holding hands and moving in a circle around a central point, often[...]

Bangladesh: Let's describe the **Bharatanatyam**, a classical dance form from **Tamil Nadu, India**.

Here's what makes it traditional:

Ancient Roots: Bharatanatyam is rooted in ancient Hindu temple rituals[...]

United States: Let's talk about the **square dance**, a quintessential American tradition!

The Basics:

Formation: Square dancers form a square, with four couples facing each other.

Music: Square dances are usually[...]

Figure 18: Anecdotal steering examples for prompt "Can you describe me a traditional dance?" with per-culture steering vectors applied on layer 25, $\alpha = 2$.

C Explicit–Implicit Localization Gap

We include Fig. 20 to illustrate performance of each model by subtask.

Prefix analysis. Examples of prefix words that are added to test what words lead to localization, relevant to Sec. 3.1.

- **France:** Baguette, Paris, Euro, Guillotine
- **Turkey:** Baklava, Istanbul, Lira, Nazar
- **USA:** Burger, New York, Dollar, Cowboy Hat
- **Bangladesh:** Biriyani, Dhaka, Taka, Rickshaw
- **Russia:** Borscht, Moscow, Ruble, Matryoshka

Turkey: “**Muhammara - Syrian Red Pepper and Walnut Dip**”

Muhammara is a delicious and vibrant dip from Syria, made with roasted red peppers, walnuts, pomegranate molasses, and spices. It’s traditionally served with[...]

France: “Cassoulet Toulousain (Toulouse Cassoulet)

This is a classic dish from Toulouse, France, known for its rich flavors and hearty ingredients.

Ingredients:

* **For the Cassoulet:[...]

Russia: “Borscht (Ukrainian Beetroot Soup) **This recipe is for a traditional Ukrainian borscht, a hearty and flavorful beetroot soup that is a staple in Ukrainian cuisine.**”

Ingredients:

* **For the borscht:[...]

Bangladesh: “Chana Masala (Bengali Chickpea Curry) **This is a flavorful and aromatic chickpea curry, a staple in Bengali cuisine. It’s traditionally made with Bengal gram (chana), potatoes, and a blend of spices[...]

United States: “New England Clam Chowder (Classic) This recipe is a classic New England staple, packed with clams, potatoes, and bacon. It’s a hearty and comforting dish perfect for a chilly day. **Ingredients:[...]

Figure 19: Anecdotal steering examples for prompt “Write me a recipe for a local dish.” with per-culture steering vectors applied on layer 25, $\alpha = 2$.

Error analysis. When no context is provided, some models have a tendency to refuse to provide the answer. In Tab. 10 we show the fraction of answers by model that fail to provide an answer.

subtask	cities	culturalbench	o1-distilled	names
aya_8b_it	0.125	0.169	0.113	0.183
gemma2_27b_it	0.000	0.000	0.000	0.000
gemma2_9b_it	0.004	0.009	0.000	0.001
gpt4o	0.638	0.044	0.086	0.319
llama31_70b_it	0.000	0.001	0.000	0.005
llama31_8b_base	0.000	0.000	0.000	0.000
llama31_8b_it	0.000	0.000	0.000	0.000

Table 10: Error rate with no context.

D Steering Prompts and Parameters

First, we report the prompt format we used for defining the steering vectors since it deviated a little from the prompt used for activation patching. We added a “My guess is **” suffix at the end of the prompt to encourage generation, and removed numbering in the prompt options. This encouraged the model to generate the entire output of the question and empirically worked better for steering.

Hyperparameters for steering. In Tab. 11, we show the best alpha and layer for steering the model. The best parameter is determined by the task the steering vector was defined to solve.

- v_{tr} : Best performance on localizing the translated prompts with no localization for each culture.
- v_{en} : Best performance on localizing the English prompts for each culture.
- v_{names} : Best performance on localizing the English prompts for each culture on the names subtask.
- $v_{universal (tr)}$: Best performance on localizing the translated prompts with no localization for each culture in the held-one-out universal vector.
- $v_{universal (en)}$: Best performance on localizing the English prompts for each culture in the held-one-out universal vector English.

	v_{tr}		v_{en}		v_{names}		$v_{universal (tr)}$	
	l	α	l	α	l	α	l	α
Bangladesh	25	2	27	2	21	2	21	2
France	25	2	25	1	23	2	25	1
Russia	27	2	25	2	25	2	25	2
Turkey	21	2	24	2	25	2	27	2
United States	22	2	22	2	22	2	21	-2

Table 11: Best-performing steering configuration per culture across all steering vectors.

Implicit steering results. In the main paper, we limited our analysis to steering on the explicit vector, however, we note that steering is also possible in the implicit context. The way we construct the implicit steering is by taking one translated prompt that correctly localized and subtracting its English variant. In Fig. 22, we illustrate the steering performance across various alpha and layers. We see that steering leads to a minor improvement (on average around 7%), but still not close to the improvement from explicit steering.

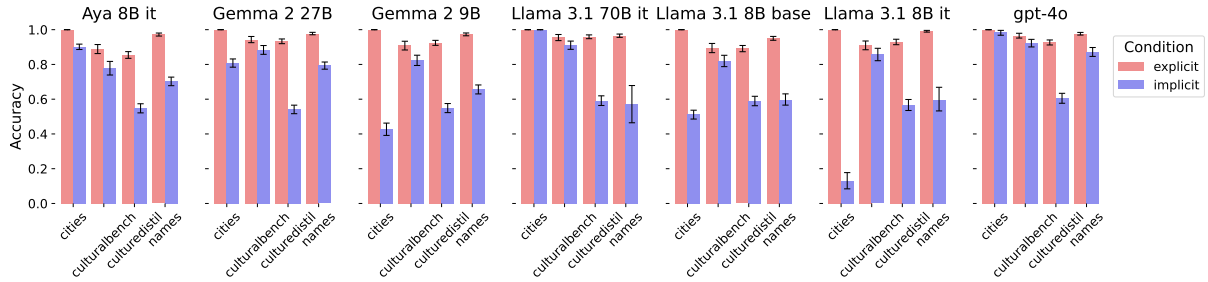


Figure 20: Performance of the different models by subtask in the explicit and implicit localization setting.

α	Layer	Country	Implicit steering	Explicit steering	No context
1	25	Bangladesh	0.696	0.739	0.628
2	26	France	0.595	0.612	0.547
1	25	Russia	0.614	0.793	0.532
2	23	Turkey	0.654	0.763	0.559
1	25	United States	0.621	0.695	0.518

Table 12: Steering results on the implicit questions. The α and layer show the best layer and alpha parameter for the implicit steering. Implicit steering column shows performance on cultural tasks when the implicit vector is added to the translated question. No context is baseline implicit performance. And explicit steering shows performance when explicit steering.

```

<bos><start_of_turn>user

Question:You must select one option and answer. First, state the selected option in full, then explain your guess. What is a common living arrangement for children after they reach adulthood?

Options:

Children go to live with their distant relatives for better education or job opportunities.

Children often continue to live with their parents, or have their parents move into their homes to take care of them<end_of_turn>

<start_of_turn>model
My guess is **

```

Figure 21: Example prompt used for model steering.

Effect of order position on steering. One way of constructing the dataset would be to include two rows for each question with the question order changed. In Fig. 23, we test to see the difference in steering performance with swapping options and not. Overall, we find that swapping the order of the questions has a negligible effect on steering vector performance. For this reason, we limit our analysis to non-swapped context for the steering analysis.

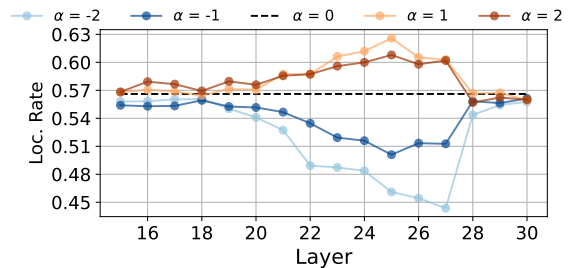


Figure 22: Steering results for per-culture vectors calculated in implicit setting on translated questions with $\alpha \in [-2, 2]$ across layers [15-30], where the x-axis represents the layer at which steering vector is applied, and the y-axis indicates the ratio of localized responses. For United States, localization performance is calculated in the reverse direction (i.e. English minus translated)

Language Hint Experiment. A potential concern with our experimental setup is that it assumes a one-to-one mapping between local knowledge, country, and language. To disentangle the effect of language from explicit cultural context, we conduct an additional experiment where we provide the model only with information about the speaker’s language (e.g., “I am currently speaking English”) without any country-specific information. We refer to this setting as *language hint*. In Table 13, we compare the localization accuracy across the implicit, explicit, and language hint settings. We find that for most models, the language hint alone is sufficient to achieve localization performance

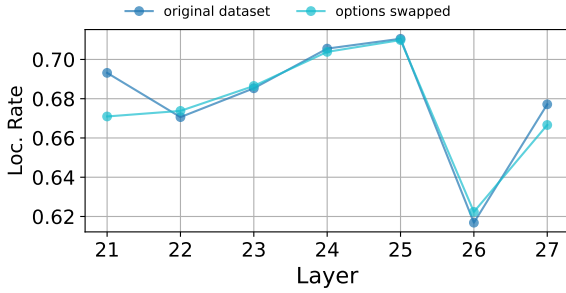


Figure 23: Steering results on the original and option-swapped datasets for per-culture vectors calculated using English pairs with $\alpha = 2$ across layers [21-27]. The x-axis represents the layer where the steering vector is applied, and the y-axis indicates the ratio of localized responses.

Table 13: Localization accuracy across different prompt settings. *Language hint* provides only the speaker’s language (“I am currently speaking [language]”) without explicit country information.

Model	Setting	Mean	95% CI
Aya Expans 8B	explicit	0.763	[0.733, 0.793]
	implicit	0.697	[0.681, 0.713]
	lang. hint	0.834	[0.822, 0.846]
Gemma 2 27B	explicit	0.878	[0.859, 0.896]
	implicit	0.727	[0.714, 0.739]
	lang. hint	0.881	[0.873, 0.891]
Gemma 2 9B	explicit	0.881	[0.862, 0.899]
	implicit	0.604	[0.588, 0.620]
	lang. hint	0.838	[0.827, 0.848]
Llama 3.1 70B-it	explicit	0.880	[0.859, 0.898]
	implicit	0.739	[0.721, 0.758]
	lang. hint	0.872	[0.862, 0.881]
Llama 3.1 8B-base	explicit	0.849	[0.829, 0.871]
	implicit	0.543	[0.528, 0.560]
	lang. hint	0.605	[0.592, 0.619]
Llama 3.1 8B-it	explicit	0.898	[0.879, 0.917]
	implicit	0.555	[0.530, 0.580]
	lang. hint	0.645	[0.633, 0.661]
GPT-4o	explicit	0.892	[0.873, 0.910]
	implicit	0.776	[0.758, 0.793]
	lang. hint	0.885	[0.877, 0.894]

comparable to the explicit setting.

E Multiple Choice Question

E.1 Steering

Throughout our analysis we focused on the binary setting where one option was always from the West for simplicity. This is a problem for the steering analysis, since the steering vector may simply learn to answer the non-American answer. If we imagine a prompt where two of the options are non-American, the model then may not give the answer related to the language of the prompt. We now extend to the five-choice setting where multiple options are non-American. To do this we augment our prior datasets with more options. Since the questions are the same across languages, we add five randomized options from each culture for these datasets. But since the o1-distill data was culture-specific, we used o3-mini-high to regenerate a set of culturally relevant questions with one option for each of our five contexts, in a way similar to the o1-preview approach previously done. In this setting we evaluate both the translated steering vector and the universal translated steering vector.

In Tab. 14 we show the performance of the steering vector on the multiple choice alongside the performance of the culture-specific steering vector (explicit translated subtract implicit translated) and implicit/explicit baselines. We observe that the universal steering vector gets similar results to the culture-specific steering, but still lags behind the explicit setting.

	$v_{\text{universal (tr)}}$	v_{tr}	Implicit	Explicit
Bangladesh	0.537	0.719	0.418	0.863
France	0.292	0.308	0.233	0.860
Russia	0.565	0.542	0.286	0.914
Turkey	0.441	0.526	0.285	0.841
United States	0.371	0.596	0.219	0.848

Table 14: Steering results for the universal (translated) and culture-specific (translated) steering vectors.

	v_{tr}		$v_{\text{universal (tr.)}}$	
	l	α	l	α
Bangladesh	24	2	21	2
France	27	2	22	2
Russia	23	2	23	2
Turkey	27	2	27	2
United States	25	2	22	2

Table 15: Optimal steering parameters for the multiple choice question answering task. The optimal parameters were done from a sweep of $\alpha \in [-2, -1, 0, 1, 2]$ and $l \in [21, \dots, 27]$

F Pinpointing Cultural Localization using Activation Patching

F.1 Activation Patching

We use activation patching (Vig et al., 2020; Ghandeharioun et al., 2024; Dumas et al., 2025) to determine which layers of the LLM are the most important ones for our cultural localization tasks. Due to the causal masking applied in the attention layers, the latent representation of the i th input token after the j th transformer block always depends on all preceding tokens $h_i^{(j)} = h_i^{(j)}(x_1, \dots, x_i)$. For notational convenience, we either omit this dependence when it is obvious from context (see above) or use the following short-hand notation $h_i^{(j)}(x)$.

Now, given a latent representation $h_i^{(j)}(x_{\text{source}})$ from a forward pass² on a source prompt x_{source} , we can patch this latent into another forward pass $h_i^{(j)}(x_{\text{target}})$ on a target prompt, effectively overwriting the embedding at that point, while observing how this changes the prediction $P(x_{\text{target}})$. We use $\hat{P}(x_{\text{target}})$ to denote the target forward pass perturbed via activation patching.

More specifically, for our analysis we calculate the latent representations after each layer on the last token position t_{source} during the source forward pass, i.e., $h_{t_{\text{source}}}^{(j)}(x_{\text{source}})$. Next, we select a target prompt x_{target} and, during its forward pass, replace its corresponding latent for the last token with those from the source prompt $h_{t_{\text{target}}}^{(j)}(x_{\text{target}}) \leftarrow h_{t_{\text{source}}}^{(j)}(x_{\text{source}})$. To understand the mechanism behind implicit and explicit localization, we use source and target pair configurations differing in the culture of the answer.

Following from our activation patching approach, we design a steering experiment by extracting embeddings from two types of prompts: ones with explicit localization (e.g., “I live in Turkey”) and control prompts without cultural context. To extract steering vectors, we use a different prompt structure, detailed in Appendix D. We subtract the control embeddings from their localized counterparts, performing this process in both English and translated versions. When we use all tasks (names, cities, o1-distilled, CulturalBench) we denote these two settings as v_{en} and v_{tr} , respectively. Alternatively, when we steer only using the translated names dataset, we use v_{names} . Using these steering vectors, we run three analyses: (1) How much of the explicit prompting performance can be recov-

²Evaluating (1) is called forward pass.

ered by adding a single vector? (2) Is the steering vector task agnostic? (3) Are steering vectors culture-specific, or can a steering vector from one culture be used to steer another culture?

F.2 Pinpointing Cultural Localization

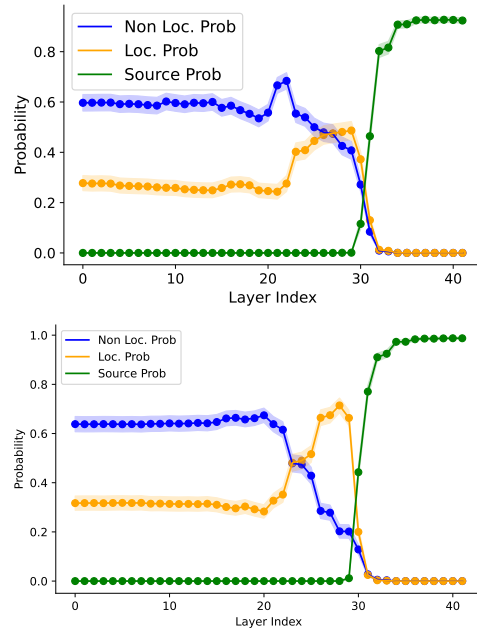


Figure 24: Activation patching results, where target prompt localized token probability (Loc. Prob) is shown in yellow, and non-localized target prompt token probability (Non Loc. Prob) is shown in blue. Finally, green shows the probability of answering the question from the source prompt. Shaded regions around plot lines represent 95% confidence intervals (CI), calculated as $\text{mean} \pm 1.96 \times \text{SEM}$. (Left) Source-translated prompt and target English prompt. (Right) Source-translated prompt with cultural context and target non-context translated prompt.

We now examine a deeper explicit–implicit gap based on how language models culturally localize text in the implicit and explicit setting through activation patching. For the implicit localization analysis, we construct pairs of prompts: a target prompt in English that generates a non-localized response, and a source prompt that is its translation with modified option labels (using letters A, B instead of numbers). The modified option keys help isolate the source prompt’s influence on the behavior. Concretely, the English version may look like “A common name here is (1) George, (2) Sergey. Answer:” where the model would answer (1). The source version would be its Russian translation.

For explicit localization, we use only translated examples, because they show the most consistent

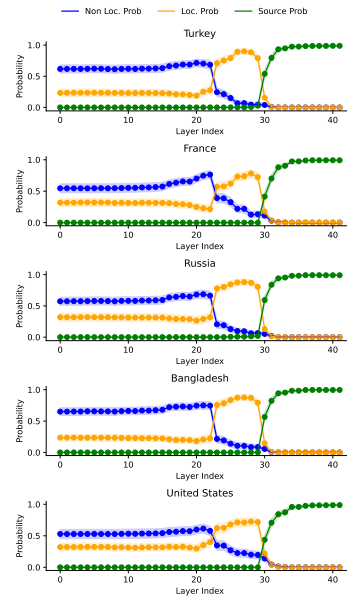
localization when given cultural context. Here, we pair target prompts without context with their counterparts including explicit cultural context as source prompts, again using modified option labels. In this case, the text we pass in would not be the translated version of the prompt above, but instead one with “I live in the United States.” as a prefix.

Fig. 24 presents three probability trajectories across model layers: the likelihood of decoding the non-localized response from the target prompt (Non. Loc. Prob), the probability of culturally localizing the target prompt to the locality of the source prompt (Loc. Prob), and the probability of generating the modified option character (A, B instead of 1, 2) from the source prompt (green line). This modification helps us identify the layers at which the source prompt starts to dominate.

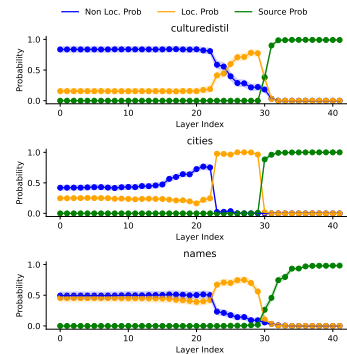
Substituting the layer early leads the context from the target prompt to effectively “overwrite” latent information from the passed-in source. When we pass it across later layers, the generation is almost always simply the source probability because the attention unable to write the relevant context. The interesting points are in the middle, where we see the correct localized answer (from the target forward pass) spike. This means that the source latent is writing in the cultural localization answer *without* overwriting the answer.

These results indicate: (1) Implicit and explicit localization spike and drop at the same layer, namely 23 and 30, suggesting that the mechanisms behind cultural localization may indicate the consolidation of a world model that becomes culturally customized within these specific middle layers. We decompose our analysis across languages and tasks, and find a similar pattern (see Appendix F). (2) Implicit localization is less pronounced than explicit localization, implying that language may not contain sufficient context for the language model to culturally localize.

F.3 Explicit (English)



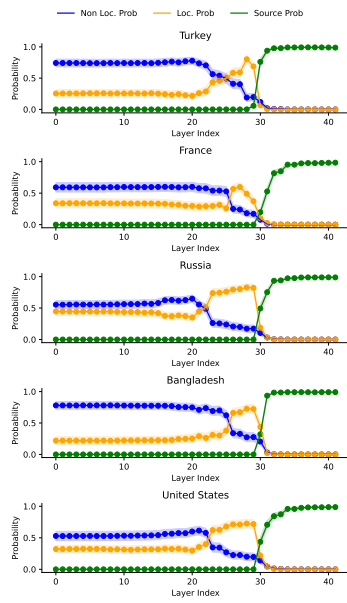
(a) Activation patching results breakdown per culture.



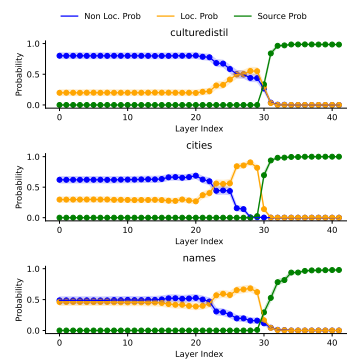
(b) Activation patching results breakdown per task.

Figure 25: Comparison of activation patching results per culture and per task. Target prompt localized token probability (Loc. Prob) is shown in yellow, and non-localized target prompt token probability (Non Loc. Prob) is shown in blue. Green shows the probability of answering the question from the source prompt. Shaded regions represent 95% confidence intervals (CI) as $\text{mean} \pm 1.96 \times \text{SEM}$. Source-English prompt with context and target English prompt.

F.4 Explicit (Translated)



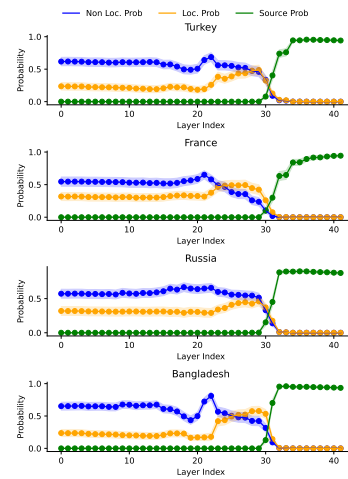
(a) Activation patching results breakdown per culture.



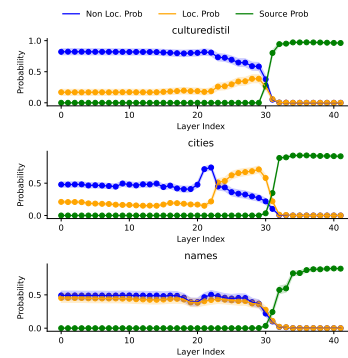
(b) Activation patching results breakdown per task.

Figure 26: Comparison of activation patching results per culture and per task. Target prompt localized token probability (Loc. Prob) is shown in yellow, and non-localized target prompt token probability (Non Loc. Prob) is shown in blue. Green shows the probability of answering the question from the source prompt. Shaded regions represent 95% confidence intervals (CI) as $\text{mean} \pm 1.96 \times \text{SEM}$.

F.5 Implicit



(a) Activation patching results breakdown per culture.



(b) Activation patching results breakdown per task.

Figure 27: Comparison of activation patching results per culture and per task. Target prompt localized token probability (Loc. Prob) is shown in yellow, and non-localized target prompt token probability (Non Loc. Prob) is shown in blue. Green shows the probability of answering the question from the source prompt. Shaded regions represent 95% confidence intervals (CI) as $\text{mean} \pm 1.96 \times \text{SEM}$. Source-translated prompt and target English prompt.

G Effect of steering on downstream capabilities

	Accuracy	Δ
bn	0.537 ± 0.018	-0.114 ± 0.014
ru	0.547 ± 0.018	-0.104 ± 0.013
fr	0.548 ± 0.018	-0.103 ± 0.014
tr	0.556 ± 0.017	-0.095 ± 0.013
us	0.573 ± 0.017	-0.078 ± 0.012
baseline	0.651	—

(a) TruthfulQA MC1 ($n = 797$).

	Accuracy	Δ
bn	0.650 ± 0.008	-0.038 ± 0.006
ru	0.663 ± 0.008	-0.025 ± 0.005
fr	0.646 ± 0.008	-0.042 ± 0.006
tr	0.662 ± 0.008	-0.026 ± 0.005
us	0.652 ± 0.008	-0.037 ± 0.006
baseline	0.688	—

(b) MMLU ($n = 14042$).

Table 16: Steered model performance for Gemma 2 9B, steered with the English vector at layer 25, $\alpha = 2$. Baseline is the unsteered model ($\alpha = 0$). Error bars are half-widths of 95% bootstrap CIs (10,000 paired resamples). All Δ CIs exclude zero.

To assess whether the localization gains from steering come at the cost of general model capabilities, we evaluate the steered model on two standard benchmarks: MMLU (Singh et al., 2025) and TruthfulQA MC1 (Lin et al., 2022). We use the same operating point as in the main text (Gemma 2 9B, layer 25, $\alpha = 2$) and apply the English perculture steering vector derived from each of the five cultures (bn, ru, fr, tr, us). Confidence intervals are computed via paired bootstrap with 10,000 resamples.

Results are shown in Tab. 16. Steering causes a statistically significant drop on both benchmarks across all steering vectors (all Δ CIs exclude zero). The effect is more pronounced on TruthfulQA, where accuracy drops by roughly 8–11 points from a baseline of 0.651, than on MMLU, where the drop is 2–4 points from a baseline of 0.688.

To get a better sense of where the TruthfulQA drop comes from, we break down accuracy by the benchmark’s 38 categories and compare the steered model to the unsteered baseline. Tab. 17 shows a representative subset of these categories. The pattern appears qualitatively consistent across the five culture vectors: categories tied to cultural and commonsense knowledge — Misinformation,

Paranormal, Fiction, Superstitions, Stereotypes, Religion, and Myths and Fairytales — show the largest drops, while broader factual domains such as Health, Economics, Sociology, and Science are affected much less. One category, Misquotations, improves slightly under steering. A possible interpretation is that steering perturbs the model’s English, US-default prior, which could affect items that depend on culturally specific knowledge more than items grounded in domain facts. We note that the us vector also causes a drop, though the smallest among the five on TruthfulQA (-0.078 vs. -0.095 to -0.114). This is consistent with the interpretation above, but suggests that any steering intervention at $\alpha = 2$ adds some noise to the residual stream regardless of the vector’s direction. Disentangling this from the cultural effect, combined with the small sample sizes in many categories, calls for a more in-depth study.

Category	n	base	us	tr	fr	ru	bn	Δ_{avg}
Misinformation	12	0.750	0.250	0.167	0.250	0.167	0.250	-0.533
Paranormal	24	0.750	0.583	0.417	0.417	0.417	0.417	-0.300
Fiction	30	0.567	0.433	0.367	0.400	0.300	0.367	-0.193
Superstitions	21	0.476	0.381	0.286	0.286	0.286	0.286	-0.171
Stereotypes	24	0.833	0.750	0.708	0.667	0.708	0.625	-0.142
Religion	15	0.533	0.467	0.400	0.333	0.400	0.400	-0.133
Myths and Fairytales	21	0.571	0.524	0.429	0.429	0.429	0.476	-0.114
Misconceptions	100	0.800	0.780	0.730	0.740	0.730	0.700	-0.064
Sociology	54	0.593	0.519	0.537	0.574	0.556	0.574	-0.041
Economics	30	0.367	0.267	0.367	0.333	0.333	0.367	-0.033
Science	6	0.667	0.500	0.667	0.667	0.667	0.500	-0.067
Health	55	0.764	0.727	0.764	0.764	0.782	0.745	-0.007
Misquotations	16	0.500	0.562	0.562	0.625	0.562	0.625	+0.087

Table 17: Per-category TruthfulQA accuracy for Gemma 2 9B steered with the English perculture vector (layer 25, $\alpha = 2$). Baseline is the unsteered model. Δ_{avg} is the mean accuracy change across the five culture vectors. A representative subset of the 38 categories is shown; categories with baseline accuracy of 0 or 1 (no headroom) are excluded. Many categories have small n , so the exact ordering should be interpreted with caution.