LANGUAGE MODELS IMPLICITLY LEARN A UNIFIED REPRESENTATION SPACE

Anonymous authors

Paper under double-blind review

ABSTRACT

Modern language and multimodal models can process a wide variety of inputs across different languages and modalities. We hypothesize that models acquire this capability through learning a *unified representation space* across heterogeneous data types. We first show that model representations for semantically equivalent inputs in different languages are similar in the intermediate layers, and that this space can further be interpreted using the model's dominant pretraining language (when it has one) via the logit lens. We also find that models show a similar tendency when processing other kinds of data, including code and visual/audio inputs. Interventions in the unified representation space further affect model outputs in expected ways: for example, replacing the image representations in a vision-language model with language token representations leads to output changes consistent with the language token semantics, suggesting that the unified representations space is not simply a byproduct of large-scale training on broad data, but something that is actively utilized by the model during input processing.

- . 1 Імтр
- 025 026

004

006 007 008

009 010

011

012

013

014

015

016

017

018

019

021

023

1 INTRODUCTION

Modern language and multimodal models (LMs)¹ are capable of processing heterogeneous data 027 types: text in different languages, non-linguistic inputs such as code and math expressions, other modalities such as images and sound, etc. How do LMs process these distinct data types with a 029 single set of parameters? One strategy might be to learn specialized subspaces for each data type that are only recruited when processing it. In many cases, however, data types that are surface-distinct 031 share underlying structures. This is most obvious for sentences in different languages with the same meaning; but such shared structures are present across other data types, e.g., between an image and its caption, or a piece of code and its natural language description. A model could leverage such 033 commonalities by learning to project surface forms of different data types into a *unified* representation 034 space, perform computations in it, and then project back out into surface forms when needed. 035

To what extent is this idealized strategy adopted by actual models? Wendler et al. (2024) find that on simple synthetic tasks, Llama-2 (Touvron et al., 2023b) maps various input languages into a shared "English space" before projecting back out into another language, hinting that it leverages this shared representation scheme to an extent. We show that this is in fact a much more general phenomenon: when a model processes inputs from multiple data types, there is a shared representation space, and this space is scaffolded by the LM's inherently dominant data type (usually English). By scaffolded, we mean that the shared space can be interpreted to an extent in the dominant data type via the logit lens (nostalgebraist, 2020).

We first show that LMs represent semantically similar inputs from different modalities to be close to
one another in the intermediate layers of the LM. Furthermore, we show that we can interpret these
intermediate representations to an extent using the LM's dominant language—e.g., when processing
a Chinese input, an English-dominant LM "thinks" in English before projecting back out to a Chinese
space. This property extends to non-linguistic inputs. When processing code as well as visual and
audio inputs, LMs represent (and perform computations) in a natural-language-adjacent space in its
intermediate layers: for code, for example, the intermediate representations reflect program semantics,
unconstrained by surface syntax; for images, we can probe out properties of an image patch (e.g.,
color, object) from a vision-language model's intermediate representations. Finally, we perform

052

¹Hereafter, we use the term "language model" loosely and also consider multimodal language models that process additional data modalities, since such models are commonly trained on top of a text LM backbone.



Figure 1: Example of the unified representation space across various input data types. For every other layer, we show the closest token in the model vocabulary to the hidden state. Llama-3's (Llama-3-Team, 2024) hidden states are the closest to English tokens when processing Chinese texts and code, in a semantically corresponding way. LLaVA (Liu et al., 2023), a vision-language model, and SALMONN (Tang et al., 2024), an audio-language model, have a similar behavior when processing images/audio. The boldface is for emphasis.

intervention experiments showing that intervening in the shared representation space, using the LM's dominant language, predictably affects model output; that is, the shared representation space (and the processing of these representations through subsequent layers) is not an accidental byproduct of the model's being trained on (say) English-dominant text, but causally impacts model behavior.

080 Our work is complementary and distinct from prior work that 081 found structural similarities between the representation spaces 082 of models trained (usually independently) on different data types, such as those showing that text representations from text-084 only LMs can be aligned, via a transformation, to vision/audio representations of modality-specific models (Ilharco et al., 2021; Merullo et al., 2022; Li et al., 2023; Ngo & Kim, 2024; Huh et al., 2024; i.a.), the literature on cross-lingual word 087 embedding alignment (Mikolov et al., 2013; Artetxe et al., 2017; Conneau et al., 2018; Schuster et al., 2019; i.a.), and work on cross-task transfer (Moschella et al., 2023; Wu et al., 090 2024; *i.a.*). We instead show that an LM trained on multiple 091



Figure 2: Models can be intervened in the dominant data type (here, English, when processing audio) and steered towards corresponding effects.

data types represents and processes them in a shared unified space *without* requiring explicit alignment transformation. We hope our findings shed light on ways to more easily interpret the mechanisms of current models and inspire future work on better model controls using these insights.

094 095 096

2 THE UNIFIED REPRESENTATION SPACE HYPOTHESIS

An LM parameterizes a similar process: it uses M_{LM} to map various input data types into a representation space $S_{LM} \subseteq \mathbb{R}^d$ (early layers), performs computations in the space (middle layers), and verbalizes the output via V_{LM} (end layers and the LM head). However, it is unknown as to how different data types are structured in the representation space. For example, one clearly inefficient possibility is that the LM partitions \mathbb{R}^d into disjoint subspaces for each data type and processes them separately. We instead hypothesize that LMs, through training, learn to represent and process different data types in a *unified* representation space. That is, semantically similar inputs $w_{1:t}^{z_1}$ and $w_{1:t'}^{z_2}$ from distinct data types—for example texts in different languages that are mutual translations are similarly mapped in S_{LM} informally, $M_{\text{LM}}(w_{1:t}^{z_1}) \approx M_{\text{LM}}(w_{1:t'}^{z_2})$. However, absolute similarity measures (i.e., $\sin(M_{\text{LM}}(w_{1:t}^{z_1}), M_{\text{LM}}(w_{1:t'}^{z_1}))$) are generally difficult and unintuitive to interpret in high dimensional spaces.² We thus focus on relative similarity measures, taking a semantically unrelated sequence $u_{1:t'}^{z_2}$, and evaluating whether $w_{1:t}^{z_1}$ is closer to $w_{1:t'}^{z_2}$ than $u_{1:t''}^{z_2}$. Formally, our hypothesis can be formulated as:

116

124

128

129

$$\sin\left(M_{\rm LM}(w_{1:t}^{z_1}), M_{\rm LM}(w_{1:t'}^{z_2})\right) > \sin\left(M_{\rm LM}(w_{1:t}^{z_1}), M_{\rm LM}(u_{1:t''}^{z_2})\right). \tag{1}$$

Moreover, when the LM has a *dominant data type* z^* in training (e.g., English for Llama-2), we hypothesize that this unified representation space is "anchored" by z^* , in the sense that Eq. 1 holds strongly enough that we can probe out to z^* from $M_{\text{LM}}(w_{1:t}^z)$. We further expect this to hold for model representations of the *future*, which autoregressive LMs are trained to model. I.e., the representation of a prefix should better align with a verbalization of the future in z^* than a non-dominant z° (though we now need a different kind of model representation, denoted with repr_{LM}):

$$\sin\left(M_{\mathrm{LM}}(w_{1:t}^{z}), \operatorname{repr}_{\mathrm{LM}}(w_{>t}^{z^{\star}})\right) > \sin\left(M_{\mathrm{LM}}(w_{1:t}^{z}), \operatorname{repr}_{\mathrm{LM}}(w_{>t}^{z^{\circ}})\right).$$
(2)

We hypothesize that this holds even when $z = z^{\circ}$. E.g., with an English-dominant LM, its encoding of the Chinese prefix $w_{1:t}^{\circ}$ ="这篇论文太难" (trans. "This paper is so hard to") should be closer to the representation of the English word $w_{t+1}^{z^{\star}}$ ="write" than its Chinese translation $w_{t+1}^{z^{\circ}}$ ="写".

3 METHOD: TESTING THE UNIFIED REPRESENTATION SPACE HYPOTHESIS

We test the above hypothesis by considering pairs of distinct data types, the dominant one z^* and a non-dominant one z° , which are different for each experiment. Whenever semantically related inputs are available (e.g., an image and its caption), we directly test Eq. 1 by using h_t^ℓ , the LM's hidden state at position t and layer ℓ , as $M_{\text{LM}}(w_{1:t}^z)$, and further using cosine similarity for the similarity function.

We operationalize Eq. 2 via the *logit lens* (nostalgebraist, 2020), a simple training-free approach for interpreting the hidden states of a model. Transformer-based LMs produce the next-token distribution using softmax (Oh_t^L) (omitting the bias term) where O is the output token embeddings (or "unembeddings") and h_t^L is the final layer hidden state. Logit lens applies the same operation to the intermediate layers to obtain $p^{\text{logitlens}}(\cdot \mid h_t^\ell) := \text{softmax} (Oh_t^\ell)$. Logit lens has been found to produce meaningful distributions that shed light on an LM's internal representations and computations.

141 Under the logit lens, repr_{LM} in Eq. 2 considers the output embedding of single tokens, $\tau^{z^*} \in \mathcal{X}_{z^*}$ 142 and $\tau^{z^\circ} \in \mathcal{X}_{z^\circ}$. We use the first token of $w_{>t}^{z^*}$ and $w_{>t}^{z^\circ}$. Using the dot product for sim(·), Eq. 2 is 143 equivalent to comparing the logit lens probabilities,

$$p^{\text{logitlens}}\left(\tau^{z^{\star}} \mid h_t^{\ell}\right) > p^{\text{logitlens}}\left(\tau^{z^{\circ}} \mid h_t^{\ell}\right),\tag{3}$$

i.e., testing whether the probability of the continuation in the dominant language is more likely
than the continuation in the original input data type. Since the logit lens is tailored for probing out
a single token, we usually consider short-enough verbalizations such that a single BPE token can
reliability identify it. This often means that the two verbalizations are two single words that are
semantic equivalents. Nevertheless, we also consider longer future verbalizations when its first token
unambiguously suggests one interpretation in that context, which allows more flexibility.

¹⁵¹ While the above test is simple, $\tau^{z^{\circ}}$ is unavailable in many multimodal models without vocabulary tokens for z° . We thus only focus on testing Eq. 1, though logit lens enables an additional test:

144

$$p^{\text{logitlens}}\left(\tau^{z^{\star}} \mid h_{t}^{\ell}\right) > p^{\text{logitlens}}\left(\upsilon^{z^{\star}} \mid h_{t}^{\ell}\right),\tag{4}$$

where v is an unrelated token. Also, using the *next* token as τ^{z^*} is unsuitable for multimodal models never been trained to output at non-dominant data type positions t, but the non-language data types are only encoded as prefixes. For such models, we hypothesize that they still fully represent the *current* input in a unified space. We thus take τ^{z^*} to be the last token of $w_{1:t'}^{z^*}$ (which is the interpretation of $w_{1:t'}^{z^*}$ under z^* ; e.g., if $w_{1:t}^{z^*}$ is an image, $w_{1:t'}^{z^*}$ can be the objects in the image described in language).

¹⁶⁰ 161

²See for example Beyer et al. (1999). Most prior work in the probing literature also implicitly uses relative similarity measures since the similarity scores are normalized over a finite label set.



Figure 3: Left: The cosine similarity of intermediate representa- parallel English vs. Chinese tokens when tions of English and Chinese parallel texts and the 95% CI. Right: processing Chinese text and the 95% CI. The same quantity minus a baseline over non-parallel texts. Parallel texts and the 95% CI. The latent representation is closer to the texts have similar representations, particularly in middle layers.

4 EVIDENCE OF THE UNIFIED REPRESENTATION SPACE

We apply our tests across a diverse data types and find evidence of a unified representation space in all cases. We show additional experiments on arithmetic data in §C with similar trends.

181 4.1 MULTILINGUAL

176 177

178

179

Wendler et al. (2024) find that when processing specific in-context learning (ICL) templates for highly synthetic lexical-level tasks (word repetition, word translation, etc.) in non-English languages, the intermediate hidden states of Llama-2 are closer to the unembeddings of English tokens than the output language. This provides some evidence for our unified representation space hypothesis, albeit constrained to a simple synthetic task and one LM. We show that this shared representation space is a general property of LMs when they face naturally occurring text.

188 Experiment 1: Representation similarity of mutual translations. Translation datasets enable 189 a direct test for Eq. 1, with semantically equivalent cross-lingual sentences as $w_{1:t}^{z_1}$ and $w_{1:t'}^{z_2}$ and 190 a randomly chosen non-matching sentence as $u_{1:t''}^{z_2}$. We use the professionally-translated English-Chinese parallel sentences from Chen et al. (2016) (N = 5260). For each sentence pair, we use 191 192 a template to transform each sentence and compute the representation cosine similarity for each layer, using the last token position as the sentence representation following Wu et al. (2023), which 193 preserves sentence information (Morris et al., 2023). We consider two English-dominant LMs, Llama-194 2 and Llama-3, one Chinese-dominant LM, Baichuan-2 (Yang et al., 2023), and one multilingual LM, 195 BLOOM (BigScience, 2023), specifically the 7B/8B variants. §A.1 contains more details. 196

In Figure 3, the raw cosine similarity is high, up to >80% (left), and it is also significantly higher than
the non-matching pairs' similarities (right), but only in the middle layers. These trends corroborate
the unified representation space hypothesis. Notably, this trend also exists for BLOOM, which means
that such a unified space still exists even without a dominating language.

201 **Experiment 2: Probing out continuations in the dominant language.** We next test Eq. 3: whether 202 continuations in the dominant language have higher probability than those in the input language. We 203 use 1,000 Chinese and English sentences from Wikipedia (Wikimedia-Foundation, 2023). For the 204 English-dominant LLama-3, we use a Chinese prefix $w_{1:t}^{z^{\circ}}$ as input and take $\tau^{z^{\circ}}$ to be the next Chinese 205 token (i.e., $w_{t+1}^{z^{\circ}}$) and $\tau^{z^{\star}}$ to be the (first token of the) English translation of $w_{t+1}^{z^{\circ}}$. §A.1 has further 206 details. Figure 4 plots the logit lens probability for the two tokens as well as the uniform distribution 207 probability. In early layers, we cannot read out either token better than random chance. After layer 17, 208 the model representations are substantially closer to the English token than the Chinese token until 209 layer 31, showing that the model hidden space is indeed better scaffolded by English than Chinese.

Next, we extend this analysis to consider global language-level trends across languages. We first compute $p(w \mid z)$, the token distribution under a language z, by running the LM tokenizer on the language-specific split of the mC4 dataset (Xue et al., 2021). We then use Bayes' rule to estimate $p(z \mid w) \propto p(w \mid z)p(z)$ with a uniform prior p(z).³ We finally compute the probability of h_t^ℓ

214 215

³This prior obviously does not reflect the training language distribution, but in fact makes our trends even more salient, since using a real (or estimated) p(z) would make $p(z \mid w)$ even larger for the dominant language.



Figure 5: Latent language probabilities for various models, visualizing the top 3 languages per layer. **Regardless** of the input language, the dominant model language is more salient in the early-middle layers, and the input language is more salient in the final layers. Bloom does not have a clear intermediate latent language.

belonging to each language as $p(z \mid h_t^{\ell}) \propto \sum_{w \in \mathcal{V}} p(z \mid w) p^{\text{logitlens}}(w \mid h_t^{\ell})$. If our hypothesis that the universal representation space is better scaffolded by the dominant language is true, we expect the dominant language z^* to have the highest probability across input languages in the middle layers.

235 Figure 5 shows the top 3 languages for each layer on 10,000 English/Chinese Wikipedia sentences. When English-dominant models process Chinese text, Wendler et al.'s finding generalizes, where 236 English dominates in the intermediate layers and Chinese only dominates in the final layers. On the 237 Chinese-dominant LM, this trend flips: when processing English text, its intermediate layers are 238 closer to Chinese space and the final layers are closer to English space. For BLOOM, a multilingual 239 model with a relatively balanced training language mixture, we do not see a clear dominating language 240 in the intermediate layers; when we manually inspect the closest token, in most cases we observe 241 symbols with no clear semantics (though this does not mean it does not have a unified representation 242 space: see experiment 1). 70B model trends in §A.1 are highly similar to the 7B/8B ones. 243

244 4.2 CODE

229

230 231

245 Many recent LMs are trained on code 246 corpora (Touvron et al., 2023a;b; 247 Llama-3-Team, 2024; Gemini-Team, 248 2024). We find that they similarly 249 process code by projecting it into a 250 unified representation space shared with regular language tokens. Figure 6 251 shows examples, where LMs in the intermediate layers tend to verbalize 253 the future in free-form English, un-254 constrained by program syntax. E.g., in the first program, given the Python 256 "... for idx, elem in prefix 257 enumerate(numbers): for idx2, 258 elem2 in enumerate(numbers", in-259 stead of the groundtruth continuation 260 in Python "): if idx != idx2: ...", the most salient intermediate 261 token is "except", likely attempting 262



Figure 6: Logit lens analysis on Llama-2 processing example Python programs. For every other layer, we show the closest token in the model vocabulary (which is sometimes the blank token), where we look at the hidden states *before* the grayed-out texts. The model tends to verbalize the future prediction in English that correspond to the code continuations (in gray).

to predict in English "(for each element in numbers) except if it is equal to idx".
Similarly, in a list expression "[1.0, 2.0," instead of continuing in Python " 3.0", it predicts "and", which is a natural way to continue in English. In these cases, it is difficult to obtain semantically equivalent English-Python pairs, so we only test Eq. 3 across targeted cases in Python below.

Experiment 1: Simple Python list literals. We systematically test the list case, where h_t^l is the hidden state after processing ",", τ^{z^*} ="and", and τ^{z° is the actual next token. Figure 7 shows that this trend holds systematically on all such commas in the MBPP dataset (Austin et al., 2021) (N = 6923, including unit tests): as expected, in the final layers, the representation is closer to the



Figure 7: Llama-2 logit lens log prob. at commas in Python list literals, of the English "and" token (and baseline tokens) vs. the next token in the program. The representation is closer to "and" in early-middle layers, while the code token dominates in later layers.

308

309

310

311

312

313



Figure 8: Llama-2 logit lens probabilities of a function call argument's name (its semantic role) vs. the actual argument expression token, in MBPP. **The latent representation is closer to the semantic argument name in the early-middle layers, and reverses for the final layers.**

ground truth next token's unembedding, and closer to "and" in the middle layers. We also show the probability with two other tokens, "or" and "not", as baselines, both of which are lower than "and".

287 Experiment 2: Python function call arguments. Function arguments have names in the defini-288 tion, such as "range(start, end, step)"; but when invoked, they are filled with actual context-289 appropriate expressions. We call the argument names "semantic roles", and the context-specific 290 expressions the "surface forms", inspired by thematic relations in linguistics (Fillmore, 1968). In the second example in Figure 1, we show that LMs predict the arguments by first "thinking" about their 291 semantic role (τ^{z^*}) and then instantiating with surface-constrained expressions (τ^{z^*}) . We extract all 292 function calls and arguments from MBPP with simple filtering, resulting in 540 arguments (see §A.2 293 for details). For each argument, we use the logit lens to inspect the hidden states at the preceding token ("(" or ","). For each argument, Figure 8 visualizes if the semantic role or the surface form is 295 closer to each layer's hidden state of Llama-2. The semantic role (τ^{z^*}) dominates for the early to 296 middle layers, even though the role token usually does not appear in the context at all, and only in the 297 final layers do the representations converge towards the surface form argument $(\tau^{z^{\circ}})$. 298





Figure 9: The cosine similarity difference between intermediate representations of matching images and captions, over non-matching ones. Left: LLaVa, Right: Chameleon. Semantically matching images and captions have more similar representations than non-matching ones.

Figure 10: The frequency of the closest token to LLaVa's hidden states describing the image color, against a baseline using "white". In many cases, the correct color word is the closest to the image representation.

Past work has investigated the representation of *separately trained* vision and text models, often 314 finding that their representation spaces are similarly structured and alignable (Merullo et al., 2022; 315 Li et al., 2023; Huh et al., 2024; *i.a.*). We show that when trained together, vision-language models 316 learn to project both modalities into a joint representation space. Current vision-language models 317 typically represent images by segmenting them into patches, embedding them into "image tokens", 318 and then feeding them into the transformer model along with other text tokens (Lu et al., 2023; 2024; 319 Liu et al., 2023; *i.a.*). We hypothesize that the intermediate representations of the image patches are 320 close to the corresponding language tokens that describe the scene. Experimental details are in §A.3. 321

Experiment 1: Representation similarity between an image and its caption. Though not constituting exact semantic equivalence, an image paired with its caption provides one possible test for Eq. 1. We take 1000 images and corresponding captions in the MSCOCO dataset (Lin



Figure 11: When processing an image patch, model logit lens probabilities of either the nouns in the corresponding caption or the patch segmentation label, as well as a baseline for each with no correspondence between the patch and the label. The image representations better match the semantically corresponding English words.

et al., 2014) and measure their hidden states cosine similarity in LLaVA-7B (Liu et al., 2023) and
Chameleon-7B (Chameleon-Team, 2024). As in Eq. 1, we subtract the average cosine similarity
between non-matching image-caption pairs as a baseline, separately for each layer. Figure 9 shows
that semantically matching inputs, even though in different modalities, are more similarly mapped in
the models' hidden space, though the similarities are lower than for mutual translations (§4.1).

Experiment 2: Patch-level analysis using logit lens. We now test the image-description similarity using the logit lens, in Eq. 4. First, as a toy setting for illustration, we inspect LLaVA's representations of pure color images, specifically those in red, green, blue, and black. Figure 10 shows that, in up to more than 20% of the time in the intermediate layers (averaged across the patches and the four colors, N = 2304), the closest token is the corresponding color word (out of all vocabulary tokens).

344 We next consider the image cap-345 tions, the same 1000. For each 346 image patch, we compute a patch-347 caption alignment score by sum-348 ming over the logit lens probabilities for all the nouns in the 349 sentence as a proxy for objects in 350 the image. We average this align-351 ment score over all patches in all 352 images, separately for each trans-353 former layer. For the irrelevant to-354 ken baseline, for each image, we 355 compute the alignment score with 356 an unrelated caption where we 357 normalize the number of nouns 358 so that the score is comparable. 359 Figures 11a and 11b show that



Figure 12: The cosine similarity difference between intermediate representations of matching audios and labels, over nonmatching ones. Semantically matching audios and labels have more similar representations than non-matching ones.



Figure 13: When SALMONN processes an audio clip, the logit lens probabilities of the English words in the audio label vs. another random label. The audio representations better match the semantically corresponding label in English.

the matching caption better aligns with the image patch representations than an unmatched caption,reliable across all layers for LLaVa and consistently in the middle-upper layers for Chameleon.

362 Finally, we perform a finer-grained study using not caption information but segmentation information, 363 with object labels in specific image locations. The setup is similar to captions, but the alignment 364 score is not computed using the correspondence between each patch and each noun in the caption, but each patch with the corresponding object label. For the irrelevant token baseline, we compute 366 the alignment by aligning each patch with a different randomly chosen object category from all categories. Figures 11c and 11d show that, for LLaVa, the patches are much better aligned to the 367 corresponding labels than randomly assigned labels (which have near-0 logit lens probability). For 368 Chameleon, this is the case for only one middle layer, and not in a statistically significant way, though 369 we will show in §5 that Chameleon's latent space can be reliably steered using English tokens. 370

371 4.4 AUDIO

Audio is another modality that is often modeled jointly with text (Lu et al., 2024; Gong et al., 2024; 2023; *i.a.*), and we perform similar experiments using SALMONN, an audio-text model. We use the VGGSound dataset (Chen et al., 2020) which contains 10-second audio clips with labels, e.g., "duck quacking" or "playing cello". We use the same two multimodal tests as in the vision case.

Experiment 1: Representation similarity between audio and its label. We study the representation cosine similarity between an audio and its label description, and subtract from it a baseline

↑

None

 \downarrow

↑

Chinese

well as the standard deviation across 10 seeds. Cross-lingual steering is consistently successful, someti ren more than monolingual steering, without substantial damage in text fluency and relevance.					
Text Lang.	Steering Dir.	Steering Lang.	Sentiment	Disfluency (\downarrow)	Relevance (†)
	None	None	$ 0.143_{\pm 0.022}$	$7.35_{\pm 1.19}$	$0.861_{\pm 0.002}$
Spanish		Spanish English	$ \begin{array}{c} 0.125_{\pm 0.034} \\ 0.139_{\pm 0.026} \end{array} $	$\frac{10.54_{\pm 2.39}}{8.75_{\pm 2.20}}$	$0.842_{\pm 0.004}$ $0.857_{\pm 0.002}$

 $0.175_{\pm 0.035}$

 $0.159_{\pm 0.026}$

 0.178 ± 0.030

 $0.152 _{\pm 0.040}$

 $0.161_{\pm 0.029}$

 $0.153_{\pm 0.034}$

 $0.179_{\pm 0.032}$

Spanish

English

None

Chinese

English

Chinese

English

 $7.98_{\pm 2.04}$

 $7.35_{\pm 1.01}$

 $11.06_{\pm 3.12}$

 10.78 ± 2.66

 11.36 ± 1.13

 $11.12_{\pm 3.12}$

 $10.90_{\pm 3.25}$

 $0.856_{\pm 0.002}$

 $0.859_{\pm 0.003}$

 $0.869_{\pm 0.004}$

 0.866 ± 0.005

 $0.864_{\pm 0.004}$

 $0.870_{\pm 0.004}$

 $0.869_{\pm 0.003}$

Table 1: Steering Llama-3's output sentiments using trigger words in English vs. the input language (either Spanish or Chinese). We report the mean sentiment, disfluency (perplexity), and relevance of the continuation,

which is the average cosine similarity between non-matching pairs, separately for each layer. On 1000 samples from VGGSound, we see in Figure 12 that semantically matching audios and labels have more similar representations in the intermediate layers.

400 **Experiment 2: Token-level analysis using logit lens.** Unlike for vision-language models where 401 we can map individual image patches to model token positions, such correspondence does not exist in SALMONN. This limits us to position-agnostic evaluations like the captioning study, preventing fine-402 grained analysis such as segmentation. Similar to the captioning experimental design, we measure the 403 average logit lens probabilities of the words in the label, and consider a random label in the dataset 404 with no word overlap as the baseline. On the same 1000 samples, Figure 13 shows a familiar trend, 405 where the audio hidden states are closer to semantically corresponding label words. We note that 406 this is a lower bound—many words in some labels, such as the prepositions in the label "writing on 407 blackboard with chalk", are unlikely to be represented in the audio hidden states. 408

409 410

378

384

386 387

389

390

391

392

393

397

398

399

5 INTERVENING IN THE UNIFIED REPRESENTATION SPACE

411 Prior work has argued that interpretability results should be tested under a causal framework, to 412 ensure that the observation is not an incidental byproduct of model training that has no actual effect on model behavior (Vig et al., 2020; Ravichander et al., 2021; Elazar et al., 2021; Chan et al., 2022; 413 *i.a.*). In this section, we show that the unified representation space does causally affect model output. 414 Specifically, semantically transforming τ^{z^*} in English interpretably leads to corresponding behavior 415 changes in the non-dominant data type. We report relevant hyperparameters in §B. 416

417 **Multilingual.** Past work has shown that (monolingual) interventions in the middle layers can steer 418 the output of LMs in predictable fashions (Subramani et al., 2022; Turner et al., 2024; Rimsky et al., 419 2024; *i.a.*). If the English-dominant LMs have a unified representation space scaffolded by English 420 tokens, we should be able to intervene on this space in English even when processing other languages. 421 We use a popular hidden space intervention technique, Activation Addition (ActAdd; Turner et al., 422 2024), which operates in two stages: (1) taking a steering word (and optionally also a contrastive 423 one) that semantically represents the steering effect, passing it through the same model, and taking its hidden states at an intermediate layer; (2) adding the steering hidden states to the hidden states 424 for the original forward pass of the regular generation process, at the same layer. We take their 425 sentiment-steering experiment but generalize it cross-lingually. See Turner et al. (2024) for details. 426

427 We consider two non-dominant languages, Spanish and Chinese, and take 1000 texts in the InterTASS 428 dataset (Spanish; Díaz-Galiano et al., 2018) and the multilingual Amazon reviews corpus (Chinese; Keung et al., 2020), and generate continuations either without modifications or intervened using 429 ActAdd. As the steering vector, we use the difference between a positive sentiment trigger word and 430 a negative one, in the appropriate direction for negative or positive steering. Specifically, we use 431 "Good" and "Bad" for English, "Bueno" and "Malo" for Spanish, and "好" and "坏" for Chinese. In

441

442

443

444

445

446

447



(a) Steering a single-argument
 (b) Replacing image representations
 (c) Steering mammal sounds to be "range (end)" call to be predicted as of a color with language tokens of predicted as non-mammal sounds, another color, and expecting the and vice versa. model to predict the new color.

Figure 14: For (a) code, (b) images, and (c) audio, steering model output using English words, for various intervention strengths ((a) and (c)) and layers ((b)). (a) and (b) measure successfulness with the proportion of instances steered to the correct output, and (c) measures the probability of predicting mammals. **Overall, intervening in the unified representation space in English reliably leads to desired model output changes.**

addition to sentiment evaluation, we also measure the generation fluency and compute the relevance
of the generation with the prefix using trained models, following Turner et al. (2024). Ideally, the
intervention should achieve the desired sentiment without hurting text fluency/relevance (see §B.1).

451 Because we take some intermediate layer representation of the steering words (step (1)), if the unified 452 representation space is language-agnostic, we expect similar representations and similar steering 453 effectiveness across steering languages. Table 1 shows that this is indeed true on Llama-3: ActAdd in 454 the text language is in most cases effective, achieving the intended effect on sentiment, with usually 455 only a modest decrease in fluency and relevance, often statistically insignificant. This aligns with 456 the English-only findings in Turner et al. (2024). And intervening in English is similarly effective as 457 using the text language. Table 2 (appendix) shows the results for Llama-2, with very similar trends.

458
459
460
460
460
461
461
462
462
462
463
464
464
464
465
465
465
466
466
466
467
468
468
468
469
469
469
460
460
460
461
461
462
462
462
463
464
464
465
465
465
466
466
466
467
468
468
468
469
469
469
460
460
461
461
462
462
461
462
462
462
463
464
464
465
465
465
466
466
466
467
468
468
468
469
469
469
469
469
460
461
461
461
462
462
462
462
462
462
464
465
465
465
466
466
466
467
468
468
468
469
469
469
469
469
469
469
460
461
461
462
461
462
462
462
462
462
462
462
462
461
462
462
462
462
462
461
462
462
462
462
462
462
462
462
463
464
464
464
465
464
465
465
465
465
466
466
466
466
467
468
468
468
469
469
469
469

We take all single-argument "range(end)" calls in the MBPP dataset (N = 159) and attempt to expand it into "range(0, end)". As the intervention, we use an even simpler method than ActAdd: because the unembedding vectors of the semantic roles are close to the intermediate hidden states, we simply compute the difference between the unembeddings of two contrasting trigger tokens ("start" - "end"), scale it by a constant coefficient, and add it to the hidden representation corresponding to the open parenthesis "(" at an intermediate layer (layer 17). For all these "range" call in the dataset, we let Llama-2⁴ generate without and with intervention. Figure 14a shows that, with increasing intervention strength, more instances are successfully steered to "range(0, end)", up to 67%.

470 **Visual inputs.** We show that we can steer the output of vision-language models by intervening on 471 the image patches using language tokens. As the models we examined in §4.3 can only generate text, 472 we analyze how this affects the textual output, specifically focusing on Chameleon which showed a 473 weaker trend in §4.3. Focusing on the color setup, if the representation for a color is similar between 474 visual and language inputs, we hypothesize that we can *replace* the image hidden states corresponding 475 to one pure color image patch with the unembedding of the language token for another color, and 476 mislead the model to "perceive" the new color when asked about it. Note that replacing the hidden 477 state is a more invasive intervention than addition. But there is one confounder: the intervened word 478 may lexically bias the model to generate the same word, without any reasoning that incorporates the new color. To control for this, we show two colors in one image and only intervene at the positions 479 corresponding to one color: if the intervention unconditionally and lexically biases the generation to 480 the new color, this effect would (incorrectly) affect both colors.⁵ 481

 ⁴⁸² ⁴We do not consider Llama-3 in this case because its default behavior usually generates "range(0, end)" in the first place, and it is unclear how to steer from "range(0, end)" to "range(end)".

⁵We tested settings that require more sophisticated reasoning such as asking for a country flag with the two colors, or asking about spatial relationships of the colors. They seem to be beyond the capability of Chameleon-7B—even without interventions, the model cannot answer the questions correctly.

486 We consider all color pairs using the same colors as in §4.3: red, green, blue, and black, and picking 487 one color in the pair and intervene to a new third color (N = 48). As the intervention, we start 488 from a layer ℓ and replace all hidden states at and after ℓ to be the unembedding of the new color 489 minus the old color. We ask the model what the two colors in the image are, and only consider the 490 intervention successful if the model answers both the new color and the other unintervened color correctly. Figure 14b shows the success rate across all ℓ : it gets as high as above 80%.⁶ We highlight 491 that, for both this experiment and the earlier ones in this section, the interventions are not even 492 necessarily guaranteed to lead sensible outputs, let alone correct ones. 493

494 Audio. We perform a similar intervention with SALMONN, with the same desideratum that the 495 QA process should require reasoning rather than outputting the intervened token as-is. We consider 496 1000 animal sounds in the VGGSound dataset, specifically only single-word animals, and ask "Is 497 this animal a mammal?" We intervene both on mammal sounds with a random non-mammal word 498 and vice versa, in case the intervention only biases the model in a certain direction. We perform the invention similarly to the code case, adding the unembedding difference between the new trigger word 499 and the original animal name, scaled by a constant, at layer 13. We measure the probability of the 500 "Yes" token and the "No" token and compute the normalized "Yes" probability. Figure 14c visualizes 501 the two cases across intervention strengths. As the strength increases, the model is more likely to 502 predict in the steered direction, again demonstrating cross-data-type intervention effectiveness.

503 504

6 RELATED WORK

505 506

507 **Representation alignment between separately trained models.** A long line of work has investi-508 gated the representations of separately trained mono-data-type models, and showed that they can be 509 aligned through a transformation. In the multilingual case, it has been found that separately trained word embeddings for different languages can be aligned (Mikolov et al., 2013; Smith et al., 2017; Cao 510 et al., 2020; i.a.). Similarly, prior work has shown that visual representations and text representations 511 from different models can be mapped together (Merullo et al., 2022; Koh et al., 2023; Maniparambil 512 et al., 2024; i.a.). Huh et al. (2024) argued that these are possible because the different data modalities 513 are projections of the same underlying reality. Our work, in contrast, looks at a *single* model that 514 processes multiple input data types and finds that the resulting representations align, without needing 515 a transformation. Xia et al. (2023) considered an objective that explicitly trains the model to increase 516 such alignment, while we analyze how it organically emerges through autoregressive training. 517

Representation evolution throughout layers. Past work has analyzed the representation evolution throughout transformer layers, inspecting how it affects reasoning (Yang et al., 2024), factual-ity (Chuang et al., 2024), knowledge (Jin et al., 2024), etc. From another angle, work on layer pruning and early exiting also speaks to representation dynamics across layers (Gromov et al., 2024; Sanyal et al., 2024; *i.a.*). Mechanistically, Elhage et al. (2021), Merullo et al. (2024), Todd et al. (2024), Hendel et al. (2023), *i.a.*, more precisely characterized the representation changes algorithmically.

Inspecting model hidden states. We adopted the logit lens for its simplicity which brings few confounders. However, alternatives exist, usually requiring some training (Belrose et al., 2023;
Ghandeharioun et al., 2024; Templeton et al., 2024; *i.a.*). They allow for more expressive explanations, though at the risk of overfitting. Similar methods have been developed for other modalities, such as Toker et al. (2024). Testing our hypothesis using these methods would be valuable future work.

529 530

531

7 CONCLUSION

Throughout extensive experiments across multiple data types, we have shown that language and multimodal models encode inputs from distinct languages and modalities into a joint representation space. Plus, intervening on this space leads to interpretation model behavior changes. We hope our findings future work to build more transparent and efficient models exploiting such properties.

⁵³⁶ ⁶One may argue this is conceptually similar to a half-language half-image input. There are many distinctions: ⁵³⁷ most importantly, a half-image is not processable by Chameleon and severely goes out of its training distribution, ⁵³⁸ since it only ever processes images of size exactly 512×512 . Other distinctions include: the presence of a ⁵³⁹ special token marking the beginning of the image; our intervention repeats the new color token, once for each ⁵⁴⁰ patch, rather than just one; and the token representation is held constant across layers rather than evolving; etc.

540 REFERENCES 541

542 543 544	Mikel Artetxe, Gorka Labaka, and Eneko Agirre. Learning bilingual word embeddings with (almost) no bilingual data. In <i>Proceedings of Annual Meeting of the Association for Computational Linguistics (ACL)</i> , 2017. URL https://aclanthology.org/P17-1042.
545 546 547	Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. <i>arXiv preprint</i> , 2021. URL https://arxiv.org/abs/2108.07732.
548 549 550 551	Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. Eliciting latent predictions from transformers with the tuned lens. <i>arXiv preprint</i> , 2023. URL https://arxiv.org/abs/2303.08112.
552 553 554	Kevin S. Beyer, Jonathan Goldstein, Raghu Ramakrishnan, and Uri Shaft. When is "nearest neighbor" meaningful? In <i>Proceedings of the 7th International Conference on Database Theory</i> , ICDT '99, pp. 217–235, Berlin, Heidelberg, 1999. Springer-Verlag. ISBN 3540654526.
555 556 557	BigScience. Bloom: A 176b-parameter open-access multilingual language model. <i>arXiv Preprint</i> , 2023. URL https://arxiv.org/abs/2211.05100.
558 559 560	Steven Cao, Nikita Kitaev, and Dan Klein. Multilingual alignment of contextual word representations. In <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> , 2020. URL https://openreview.net/forum?id=r1xCMyBtPS.
561 562 563	Chameleon-Team. Chameleon: Mixed-modal early-fusion foundation models. <i>arXiv preprint</i> , 2024. URL https://arXiv.org/abs/2405.09818.
564 565 566 567 568	Lawrence Chan, Adrià Garriga-Alonso, Nicholas Goldwosky-Dill, Ryan Greenblatt, Jenny Nitishinskaya, Ansh Radhakrishnan, Buck Shlegeris, and Nate Thomas. Causal scrubbing, a method for rigorously testing interpretability hypotheses. <i>AI Alignment Forum</i> , 2022. URL https://www.alignmentforum.org/posts/JvZhhzycHu2Yd57RN/ causal-scrubbing-a-method-for-rigorously-testing.
569 570 571	Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio- visual dataset. In <i>Proceedings of the International Conference on Acoustics, Speech, and Signal</i> <i>Processing (ICASSP)</i> , 2020. URL https://ieeexplore.ieee.org/document/9053174.
572 573 574	Song Chen, Gary Krug, and Stephanie Strassel. Gale phase 3 and 4 chinese newswire parallel text. <i>Linguistic Data Consortium</i> , 2016. URL https://catalog.ldc.upenn.edu/LDC2016T25.
575 576 577 578	Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R. Glass, and Pengcheng He. Dola: Decoding by contrasting layers improves factuality in large language models. In <i>Proceedings of the</i> <i>International Conference on Learning Representations (ICLR)</i> , 2024. URL https://openreview. net/forum?id=Th6NyL07na.
579 580 581	Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. Word translation without parallel data. In <i>Proceedings of International Conference on Learning</i> <i>Representations (ICLR)</i> , 2018. URL https://arxiv.org/abs/1710.04087.
582 583 584 585 586	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsuper- vised cross-lingual representation learning at scale. In <i>Proceedings of the Annual Meeting of the</i> <i>Association for Computational Linguistics (ACL)</i> , 2020. URL https://aclanthology.org/2020. acl-main.747.
587 588 589 590 591	Manuel Carlos Díaz-Galiano, Eugenio Martínez-Cámara, Miguel Ángel García Cumbreras, Manuel García Vega, and Julio Villena-Román. The democratization of deep learning in tass 2017. <i>Proces. del Leng. Natural</i> , 2018. URL https://api.semanticscholar.org/CorpusID: 13667878.
592 593	Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. Amnesic Probing: Behavioral Explanation with Amnesic Counterfactuals. <i>Transactions of the Association for Computational Linguistics (TACL)</i> , 2021. URL https://doi.org/10.1162/tacl_a_00359.

594	Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda
595	Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli,
596	Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal
597	Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris
598	Olah. A mathematical framework for transformer circuits. <i>Transformer Circuits Thread</i> , 2021.
599	https://transformer-circuits.pub/2021/framework/index.html.
600	Charles I Fillmore. The asso for ease. In Emmon Rech and Debort T. Horms (eds.), Universals in
601	Linguistic Theory pp. 0.88 Holt Dinebert and Winston New York 1068
602	Linguistic Theory, pp. 0–66. Holt, Kinchart and Winston, New Tork, 1966.
603	Gemini-Team. Gemini: A family of highly capable multimodal models. arXiv preprint, 2024. URL
604	https://arxiv.org/abs/2312.11805.
605	
606	Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. Patchscopes: A
607 608	unifying framework for inspecting hidden representations of language models. <i>arXiv preprint</i> , 2024. URL https://arxiv.org/abs/2401.06102.
609	
610	Yuan Gong, Alexander H. Liu, Hongyin Luo, Leonid Karlinsky, and James Glass. Joint audio
611	and speech understanding. In Proceedings of the IEEE Automatic Speech Recognition and Understanding Workshop (ASPU) 2022 UPL https://org/abs/2200.14405
612	Onderstanding Workshop (ASKO), 2023. UKL https://arxiv.org/abs/2309.14405.
613	Yuan Gong, Hongvin Luo, Alexander H. Liu, Leonid Karlinsky, and James R. Glass, Listen, think,
614	and understand. In Proceedings of the International Conference on Learning Representations
615	(ICLR), 2024. URL https://openreview.net/forum?id=nBZBPXdJlC.
616	
617	Andrey Gromov, Kushal Tirumala, Hassan Shapourian, Paolo Glorioso, and Daniel A. Roberts. The
618	unreasonable ineffectiveness of the deeper layers. arXiv preprint, 2024. URL https://arxiv.
619	org/abs/2403.17887.
620	Rose Hendel Mor Geva and Amir Globerson. In context learning creates task vectors. In Findings
621	of the Association for Computational Linguistics: FMNLP December 2023 URL https://
622	aclanthology.org/2023.findings.emplp.624
623	
624	Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom
625	embeddings, convolutional neural networks and incremental parsing. To appear, 2017. URL
626	https://spacy.io.
627	
628	tion hypothesis. In <i>Brassadings of Machine Learning Pessarch (ICML)</i> 2024. URL https://
629	(ICML), 2024. UKL HTTPS: //proceedings_mlr_press/v235/bub24a_html
630	//proceedings.met.press/v255/hunz4a.htmt.
631	Gabriel Ilharco, Rowan Zellers, Ali Farhadi, and Hannaneh Hajishirzi. Probing contextual language
632	models for common ground with visual representations. In Proceedings of the Conference of the
633	North American Chapter of the Association for Computational Linguistics: Human Language
634	<i>Technologies</i> (<i>NAACL-HLT</i>), 2021. URL https://aclanthology.org/2021.naacl-main.422.
635	Mingan Lin Olakoj Va Linganon Harro Olarakana Zana Zhatia Wasa Wasa Harri H
636	Zhao Kai Mai Yanda Mang Kaiza Ding at al. Explained constants doubt. How have
637	znao, Kai wiei, Tanua wieng, Kaize Ding, et al. Exploring concept depth: How large language models acquire knowledge at different layers? arViv preprint 2024 UPL https://arviv.arg/
638	abs/2404_07066
639	
640	Phillip Keung, Yichao Lu, György Szarvas, and Noah A. Smith. The multilingual Amazon reviews
641	corpus. In Proceedings of the Conference on Empirical Methods in Natural Language Processing
642	(EMNLP), 2020. URL https://aclanthology.org/2020.emnlp-main.369.
6/3	
644	Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. Grounding language models to images for
645	multimodal inputs and outputs. In <i>Proceedings of the International Conference on Machine</i>
646	Learning (ICLK), 2023. UKL $\Pi(IPS://arXiv.org/abs/2301.13823.$
647	Jiaang Li, Yova Kementchedjhieva, and Anders Søgaard. Implications of the convergence of language and vision model geometries. <i>arXiv preprint</i> , 2023. URL https://arxiv.org/abs/2302.06555.

648	Tsung-Yi Lin Michael Maire Serge Belongie James Hays Pietro Perona Deva Ramanan Piotr
649	Dollar and Larry Zitnick, Microsoft age: Common objects in context. In <i>Proceedings of the</i>
CE0	Donar, and Early Zhinek. Wierosoft coco. Common objects in context. In Proceedings of the
000	<i>European Conference on Computer Vision (ECCV)</i> , 2014. URL https://www.microsoft.com/
651	<pre>en-us/research/publication/microsoft-coco-common-objects-in-context/.</pre>

- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In Advances in Neural Information Processing Systems (NeurIPS), 2023. URL https://arxiv.org/abs/2304. 08485.
- The Llama-3-Team. The llama 3 herd of models. arXiv Preprint, 2024. URL https://arxiv.org/
 abs/2407.21783.
- Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi.
 UNIFIED-IO: A unified model for vision, language, and multi-modal tasks. In *Proceed-ings of the International Conference on Learning Representations (ICLR)*, 2023. URL https://openreview.net/forum?id=E01k9048soZ.
- Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. URL https://arxiv.org/abs/2312.17172.
- Mayug Maniparambil, Raiymbek Akshulakov, Yasser Abdelaziz Dahou Djilali, Mohamed El Amine Seddik, Sanath Narayan, Karttikeya Mangalam, and Noel E. O'Connor. Do vision and language encoders represent the world similarly? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. URL https://arxiv.org/abs/2401. 05224.
- MDBG. Mdbg chinese-english dictionary (cc-cedict). MBDG, 2024. URL https://www.mdbg.net/ chinese/dictionary?page=cc-cedict. Downloaded: 2024-09-25.
- Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. Linearly mapping from image to text space. *arXiv preprint*, 2022. URL https://arxiv.org/abs/2209.15162.
- Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. Language models implement simple Word2Vecstyle vector arithmetic. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2024. URL https://aclanthology.org/2024.naacl-long.281.
- Tomas Mikolov, Quoc V Le, and Ilya Sutskever. Exploiting similarities among languages for machine
 translation. *arXiv preprint*, 2013. URL https://arxiv.org/abs/1309.4168.
- John Morris, Volodymyr Kuleshov, Vitaly Shmatikov, and Alexander Rush. Text embeddings reveal (almost) as much as text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023. URL https://aclanthology.org/2023.emnlp-main. 765.

688

- Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and Emanuele Rodolà. Relative representations enable zero-shot latent space communication. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview. net/forum?id=SrC-nwieGJ.
- Jerry Ngo and Yoon Kim. What do language models hear? probing for auditory representations
 in language models. In *Proceedings of Annual Meeting of the Association for Computational Linguistics (ACL)*, 2024. URL https://arxiv.org/abs/2402.16998.
- 696 697 nostalgebraist. Interpreting gpt: the logit lens. LessWrong, 2020. URL https://www.lesswrong. com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens.
- Abhilasha Ravichander, Yonatan Belinkov, and Eduard Hovy. Probing the probing paradigm: Does probing accuracy entail task relevance? In *Proceedings of the Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, 2021. URL https://aclanthology.org/2021.eacl-main.295.

702 703 704 705	Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. Steering llama 2 via contrastive activation addition. In <i>Proceedings of the 62nd Annual Meeting of</i> <i>the Association for Computational Linguistics (ACL)</i> , 2024. URL https://aclanthology.org/ 2024.acl-long.828.
706 707 708 709	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint</i> , 2020. URL https://arxiv.org/abs/1910.01108.
710 711	Sunny Sanyal, Sujay Sanghavi, and Alexandros G. Dimakis. Pre-training small base lms with fewer tokens. <i>arXiv preprint</i> , 2024. URL https://arxiv.org/abs/2404.08634.
712 713 714 715 716	Tal Schuster, Ori Ram, Regina Barzilay, and Amir Globerson. Cross-lingual alignment of contextual word embeddings, with applications to zero-shot dependency parsing. In <i>Proceedings of the Con-</i> <i>ference of the North American Chapter of the Association for Computational Linguistics: Human</i> <i>Language Technologies (NAACL-HLT)</i> , 2019. URL https://aclanthology.org/N19-1162.
717 718 719 720	Samuel L. Smith, David H. P. Turban, Steven Hamblin, and Nils Y. Hammerla. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. In <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> , 2017. URL https://openreview.net/forum? id=r1Aab85gg.
721 722 723	Nishant Subramani, Nivedita Suresh, and Matthew Peters. Extracting latent steering vectors from pretrained language models. In <i>Findings of the Association for Computational Linguistics: ACL</i> , 2022. URL https://aclanthology.org/2022.findings-acl.48.
724 725 726	Junyi Sun. Jieba: Chinese text segmentation tool. Github, 2024. URL https://github.com/fxsjy/jieba. Accessed: 2024-09-25.
727 728 729 730	Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao Zhang. SALMONN: Towards generic hearing abilities for large language models. In <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> , 2024. URL https://openreview.net/forum?id=14rn7HpKVk.
731 732 733 734 735 736 737	Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. <i>Transformer Circuits Thread</i> , 2024. URL https://transformer-circuits.pub/2024/scaling-monosemanticity/ index.html.
738 739 740 741	Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. Function vectors in large language models. In <i>Proceedings of the International Conference on Learning Representations (ICLR)</i> , 2024. URL https://openreview.net/forum?id=AwyxtyMwaG.
742 743 744 745	Michael Toker, Hadas Orgad, Mor Ventura, Dana Arad, and Yonatan Belinkov. Diffusion lens: Interpreting text encoders in text-to-image pipelines. In <i>Proceedings of the Annual Meeting of the</i> <i>Association for Computational Linguistics (ACL)</i> , 2024. URL https://aclanthology.org/2024. acl-long.524.
746 747 748 749	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. <i>arXiv preprint</i> , 2023a. URL https://arxiv.org/abs/2302.13971.
750 751 752 753	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint</i> , 2023b. URL https://arxiv.org/abs/2307.09288.
754 755	Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. Activation addition: Steering language models without optimization. <i>arXiv</i> preprint, 2024. URL https://arxiv.org/abs/2308.10248.

756 757 758 759 760	Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. Investigating gender bias in language models using causal mediation analysis. In <i>Advances in Neural Information Processing Systems</i> (<i>NeurIPS</i>), 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/ 92650b2e92217715fe312e6fa7b90d82-Paper.pdf.
761 762 763 764	Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Multilingual e5 text embeddings: A technical report. <i>arXiv preprint</i> , 2024. URL https: //arxiv.org/abs/2402.05672.
765 766 767 768	Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. Do llamas work in English? on the latent language of multilingual transformers. In <i>Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)</i> , 2024. URL https://aclanthology.org/2024.acl-long.820.
769 770	Wikimedia-Foundation. Wikimedia downloads, 2023. URL https://dumps.wikimedia.org.
771 772 773	Zhaofeng Wu, William Merrill, Hao Peng, Iz Beltagy, and Noah A. Smith. Transparency Helps Reveal When Language Models Learn Meaning. <i>Transactions of the Association for Computational</i> <i>Linguistics (TACL)</i> , 2023. URL https://doi.org/10.1162/tacl_a_00565.
774 775 776 777	Zijun Wu, Yongkang Wu, and Lili Mou. Zero-shot continuous prompt transfer: Generalizing task semantics across language models. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=26Xphug0cS.
778 779 780	Yan Xia, Hai Huang, Jieming Zhu, and Zhou Zhao. Achieving cross modal generalization with multimodal unified representation. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023. URL https://openreview.net/forum?id=t7ZowrDWVw.
781 782 783 784 785	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mT5: A massively multilingual pre-trained text-to-text transformer. In <i>Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)</i> , 2021. URL https://aclanthology.org/2021.naacl-main.41.
786 787 788 789 790 791 792 793 794	Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, JunTao Dai, Kun Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu, MingAn Lin, Nuolan Nie, Peidong Guo, Ruiyang Sun, Tao Zhang, Tianpeng Li, Tianyu Li, Wei Cheng, Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan, Yanjun Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou, and Zhiying Wu. Baichuan 2: Open large-scale language models. <i>arXiv Preprint</i> , 2023. URL https://arxiv.org/abs/2309.10305.
795 796 797 798 799 800	Sohee Yang, Elena Gribovskaya, Nora Kassner, Mor Geva, and Sebastian Riedel. Do large language models latently perform multi-hop reasoning? In <i>Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)</i> , 2024. URL https://aclanthology.org/2024. acl-long.550.

A EXPERIMENTAL DETAILS FOR §4

A.1 MULTILINGUAL

801

802 803

804

For Experiment 1, for each sentence pair, we use a template to transform each sentence. This is due to the automatic code-switching behavior of LMs. For an English model processing Chinese text, we expect the Chinese tokens to have high probabilities in the final layer because they need to be output; however, we observe these models tend to code-switch back to their dominant language after a full sentence, which confounds our analysis. We therefore put the parallel sentences into a template, "{English Sentence} This represents" (and the corresponding Chinese version), as the model is less



Figure 15: Latent language probabilities for various models, visualizing the top 3 languages per layer. Regardless of the input language, the dominant model language is more salient in the early-middle layers, and the input language is more salient in the final layers. Bloom does not have a clear intermediate latent language.

likely to code-switch mid-sentence after "represents". We experimented with other templates that
led to similar results. Furthermore, for sentence in GALE (Chen et al., 2016), we make sure the
transcript; unicode is not empty for both the source and the translation.

For Experiment 2, due to tokenization, it is challenging to obtain exactly parallel English-Chinese tokens, and hence we perform aggressive filtering. We consider only text positions where the next BPE token (1) is a valid Chinese word (as segmented by Jieba (Sun, 2024)), and (2) has an English translation (using the English-Chinese dictionary CC-CEDICT (MDBG, 2024)). For example, "今 天是开心的一天" (Today is a happy day), Llama-3 tokenizes it as { '今天', '是', '开', '心', '的 -', '天' }, while Jieba segments it to be { '今天', '是', '开心', '的', '一天' }. We only keep { '今天', '是' }. Furthermore, only '今天''s translation is a single token, the only token that survives the cutoff is '今天'.

The full result of Figure 5 is in Figure 15. We observe that the 7B/8B and 70B models of the same model family have highly similar trends, so we only consider the 7B/8B models in other experiments.

846 847 848

832

A.2 CODE

849 We consider all non-zero-argument function calls in the MBPP dataset, excluding unit tests. We 850 automatically identify the argument names (the "semantic roles") by function inspection for built-in 851 functions and by looking at the function definition for those defined in-context, and skip when this 852 is not possible. We also ignore arguments whose semantic roles are generically called "obj" or "object", and instances where the instantiated surface-form argument is the same as the semantic 853 role. We look the hidden state corresponding to the previous token, either "(" or ",", except when 854 tokenization renders this impossible (e.g., when the previous token is merged with a part of the 855 surface argument). This leaves 540 arguments. 856

857

859

858 A.3 VISION-LANGUAGE

To pass the images through the model, we embed them in templates, only for the logit
lens experiments. For the color experiment, we use the template "USER: What is the color
in the image?<image>\n ASSISTANT:". For the caption and segmentation experiments, we
use "USER: What is in the image?\n<image> ASSISTANT:" for LLaVA and "What is in the
image?\n<image>" for Chameleon.

	Ð	5,	U	J.	
Text Lang.	Steering Dir.	Steering Lang.	Sentiment	Disfluency (\downarrow)	Relevance (†)
	None	None	$0.144_{\pm 0.014}$	$8.58_{\pm 0.57}$	$0.850_{\pm 0.006}$
Spanish	\downarrow	Spanish English	$\left \begin{array}{c} 0.143_{\pm 0.012}\\ 0.097_{\pm 0.024}\end{array}\right $	$\begin{array}{c} 8.84_{\pm 0.79} \\ 8.99_{\pm 0.72} \end{array}$	$\begin{array}{c} 0.847_{\pm 0.006} \\ 0.847_{\pm 0.005} \end{array}$
	 ↑	Spanish English	$\left \begin{array}{c} 0.164_{\pm 0.018}\\ 0.149_{\pm 0.015}\end{array}\right.$	$\begin{array}{c} 9.11_{\pm 0.50} \\ 8.35_{\pm 0.30} \end{array}$	$\begin{array}{c} 0.844_{\pm 0.005} \\ 0.849_{\pm 0.006} \end{array}$
	None	None	$0.223_{\pm 0.036}$	$14.63_{\pm 2.65}$	$0.844_{\pm 0.009}$
Chinese	\downarrow	Chinese English	$ \begin{vmatrix} 0.117_{\pm 0.080} \\ 0.156_{\pm 0.076} \end{vmatrix} $	$\begin{array}{c} 15.29_{\pm 2.47} \\ 14.80_{\pm 2.24} \end{array}$	$\begin{array}{c} 0.840_{\pm 0.011} \\ 0.842_{\pm 0.008} \end{array}$
	 ↑	Chinese English	$ \begin{vmatrix} 0.359_{\pm 0.077} \\ 0.227_{\pm 0.038} \end{vmatrix} $	$545.94_{\pm 1544.36}_{14.14_{\pm 2.42}}$	$\begin{array}{c} 0.839_{\pm 0.010} \\ 0.845_{\pm 0.009} \end{array}$

Table 2: Steering Llama-2's output sentiments using trigger words in English vs. the input language (eitherSpanish or Chinese). We report the mean sentiment, disfluency (perplexity), and relevance of the continuation,as well as the standard deviation across 10 seeds. Cross-lingual steering is consistently successful, sometimeseven more than monolingual steering, without substantial damage in text fluency and relevance.

For all caption and segmentation experiments, we use the MSCOCO 2017 dataset. For the caption
evaluation, for each image patch (with the associated caption for the entire image), we compute a scalar
patch-caption alignment score (for each layer separately), by averaging the logit lens probabilities of
all nouns in the caption at that image patch position.⁷ The baseline is similarly calculating by taking
an irrelevant caption (see §4.3). We then average the patch-caption alignment score across all patches
in all images to obtain the curves in Figure 11a and 11b.

For the segmentation evaluation, we use the MSCOCO 2017 panoptic segmentation labels. The metric calculation is similar as above. Instead of a scalar patch-caption alignment score, we consider a scalar patch-label alignment score between a patch and its matching segmentation label, computed likewise using the logit lens. We consider a patch and a label as matching if there is an image segment with that label that occupies more than half of the pixels in the patch. Under this definition, a patch cannot have more than one label. When a patch is not matched with any label, we disregard it. In the baseline, we use a randomly chosen incorrect label from all possible labels for the alignment score. Finally, we average this alignment score across all patches in all images to obtain the curves in Figure 11c and 11d.

B EXPERIMENTAL DETAILS FOR §5

For both the code and vision-language intervention experiments, we use argmax decoding.

904 B.1 MULTILINGUAL

For each language, we sample N = 1000 instances from the training set of InterTASS for Spanish and the multilingual Amazon reviews corpus for Chinese. Following Turner et al. (2024), we use trained models for various metrics. We automatically evaluate the sentiment of the generation using a DistillBERT-based (Sanh et al., 2020) model finetuned for multilingual sentiment analysis,⁸ judge the generation fluency by taking the conditional perplexity of the generation given the prefix from Llama-3.1-70B (Llama-3-Team, 2024),⁹ and compute the relevance of the generation with the prefix by computing the cosine similarity between the generation and the prefix using a XLM-R-Large-based (Conneau et al., 2020) model finetuned for sentence representation (Wang et al., 2024). All these models support both Spanish and Chinese.

⁸https://huggingface.co/lxyuan/distilbert-base-multilingual-cased-sentiments-student

⁷Words with NOUN or PROPN tags given by SpaCy's en_core_web_trf model (Honnibal & Montani, 2017).

⁹We also tried using Mistral-Nemo-Base (https://huggingface.co/mistralai/ Mistral-Nemo-Base-2407), and found similar trends.



(a) Cosine similarity between an (b) Same as (a), but only the exact (c) Logit lens log probability when arithmetic expression in Arabic nu- translation similarities subtracted by predicting a number, between either merals vs. English words, broken the others. down into separate categories.

the number itself or its English equivalent.

Figure 16: Results for the arithmetic experiments. The 95% CI is plotted in all. Expressions in Arabic numerals have similar representation as corresponding expressions in English, as well as the unembeddings of corresponding number words in English.

934 We perform ActAdd by passing both the positive and negative steering words through the LM, taking 935 their hidden states at layer 17, computing their difference, scaling it by a constant, and adding it to 936 the normal generation forward pass also at layer 17, exactly following Turner et al. (2024), except we 937 use a scaling coefficient of 5, rather than 2 in their experiments, for which we observed a larger effect. 938 For generation, we use a temperature of 1, top-p=0.3, and a frequency penalty of 1, all following Turner et al. (2024), without tuning. 939

940 We showed the Llama-3 intervention results in Table 1, and here in Table 2 we show the results on 941 Llama2, with similar trends. 942

943 944 945

927

928

929

930

931

932

933

С EXPERIMENTS WITH ARITHMETIC EXPRESSIONS

946 We hypothesize that a similar trend exists when LMs process arithmetic expressions where they route 947 to a shared space anchored by numerical words in English in intermediate layers. We consider simple 948 expressions in the form of "a=b+c" or "a=b*c"; for simplicity, we restrict "a" and "b" to be at most 949 two digits and "c" to be a single positive digit.

950 951

Experiment 1: Representations are similar for translations. Here, we only consider the right-952 hand side, "b+c" and "b*c", as $w_{1:t}^{z_1}$ in Eq. 1. Like in the multilingual case, we translate them 953 into English (e.g., "five plus three") as $w_{1:t'}^{z_2}$, and evaluate the representation cosine similarity 954 between every English expression and every numeric expression, throughout layers. We group 955 the pairwise cosine similarities in three buckets: (1) exact translation (e.g., "5+3" and "five plus 956 three"; N = 1123), (2) non-exact but same value (e.g., "5+3" and "two plus six"; N = 13293), 957 and (3) different value (N = 1247836). Figure 16a shows that exact translations have high cosine 958 similarity, although this is to be expected since embeddings of numbers and their corresponding 959 English words are near one another (thus even a bag-of-word-embeddings should also have high 960 similarity). More interestingly, we that the similarities are still higher when the surface forms are 961 distinct but the "meaning" of the expression (i.e., the value of the expression) is the same. Next, like in §4.1, we subtract the cosine similarities among non-translation pairs as a baseline $(u_{1:t'}^{z_2})$. 962 Figure 16b shows high similarity in the early-middle layers for translations over the baseline, but 963 gradually decreasing to near 0. 964

965 966

Experiment 2: Representations are anchored by semantically-equivalent English words. We 967 hypothesize that, for some prefix such as "a=b+", the intermediate representations h_t^{ℓ} are close to the 968 English word for "c" that would make the equality hold. First, we randomly sample 100 such prefixes 969 and take the representation of the last token at all layers. For each prefix, we plot the representation 970 evolution throughout layers using PCA, as well as the unembeddings of numbers in English τ^{z^*} vs. 971 numerals $\tau^{z^{\circ}}$. Figure 18 shows that the representations indeed go through the space occupied by the



Figure 17: Steering arithmetic expressions' results to a different value.

Figure 18: The Llama-3 hidden representation evolution when predicting a number, projected by PCA where the principal components are learned on the output embeddings of 20 number tokens, 10 in English and 10 numerals.

English words in intermediate layers. Next, we repeat our logit lens experiments, inspecting the log probability of the following numeral token vs. its English version (N = 1123). Figure 16c shows that the two tokens have similar log probability until around layer 25, after which the numeral token dominates.

Experiment 3: Intervention. We perform intervention using our arithmetic expressions, for example "4=1+3". We intervene by attempting to modify the token after "+" to be one smaller, e.g. "2" here, and expect this to not only lead the model to output "2" instead of "3", but also fundamentally affects the model's reasoning process and causes the model to patch this error with an additional suffix "+1", i.e., "4=1+2+1". We use ActAdd except for adding the intervention vector (e.g., "three" -"two") only at the position of "+".¹⁰ For all addition expressions in our data (N = 846), we perform such intervention at an intermediate layer (25 for Llama-3 and 30 for Llama-2) and measure how often this leads to the model correctly outputting the decremented number followed by "+1", versus unchanged, or changed to some other output. Figure 17 shows that, as the intervention coefficient (i.e., the scaling constant of the vector) increases, this procedure leads to the expected output for up to > 90% of the instances.

^{1024 &}lt;sup>10</sup>Another difference is that we do not use the hidden representation after seeing e.g. "three", because that usually represents the *next* token. Instead, we use a prefix that uniquely determines the number, e.g. "Eight equals to five plus", and take the last token hidden representation, which *is* supposed to represent "three".