Seeing Through Words, Speaking Through Pixels: Deep Representational Alignment Between Vision and Language Models

Zoe Wanying He Sean Trott Meenakshi Khosla

Department of Cognitive Science University of California, San Diego {wah016, sttrott, mkhosla}@ucsd.edu

Abstract

Recent studies show that deep vision-only and language-only models—trained on disjoint modalities—nonetheless project their inputs into a partially aligned representational space. Yet we still lack a clear picture of where in each network this convergence emerges, what visual or linguistic cues support it, whether it captures human preferences in many-to-many imagetext scenarios, and how aggregating exemplars of the same concept affects alignment. Here, we systematically investigate these questions. We find that alignment peaks in mid-to-late layers of both model types, reflecting a shift from modality-specific to conceptually shared representations. This alignment is robust to appearance-only changes but collapses when semantics are altered (e.g., object removal or word-order scrambling), highlighting that the shared code is truly semantic. Moving beyond the one-to-one image-caption paradigm, a forced-choice "Pick-a-Pic" task shows that human preferences for image-caption matches are mirrored in the embedding spaces across all vision-language model pairs. This pattern holds bidirectionally when multiple captions correspond to a single image, demonstrating that models capture fine-grained semantic distinctions akin to human judgments. Surprisingly, averaging embeddings across exemplars amplifies alignment rather than blurring detail. Together, our results demonstrate that unimodal networks converge on a shared semantic code that aligns with human judgments and strengthens with exemplar aggregation.

1 Introduction

The idea of a universal, modality-independent substrate of meaning has long intrigued philosophy, neuroscience and cognitive science—from Plato's ideal forms to Fodor's "Language of Thought" (mentalese) (Fodor, 1975). This motivates a foundamental question: do putatively distinct systems—such as vision and language mod-

els—encode meaning in a shared, abstract space or in modality-specific codes?

Rapid developments in AI—particularly largescale vision and language models—provide novel tools to explore these ideas computationally. Large-scale vision-only and language-only models, trained on massive but disjoint corpora, nonetheless exhibit striking representational convergence. Huh et al. (2024) coined this phenomenon the "Platonic Representation Hypothesis", showing that increasingly capable LLMs align more tightly with larger vision models. Interestingly, this alignment occurs even without explicit cross-modal training. This "Platonic Representation Hypothesis" is further supported by Maniparambil et al. (2024), who demonstrate that this convergence manifests across a range of model architectures and training paradigms.

Critically, cross-modal alignment is not merely correlational. Merullo et al. (2022) show that training just one linear projection is enough to map a frozen vision-transformer's embeddings into the token-embedding space of a frozen language model, letting the stitched system caption images and answer visual questions without any additional multimodal training. Similarly, Koh et al. (2023) show analogous gains for the reverse mapping from text to image, showing that a frozen LLM can be visually grounded with a single learned linear map, achieving strong zero-shot performance on tasks such as contextual image retrieval and multimodal dialogue.

Marjieh et al. (2024) show that even multimodal models like GPT-4 rely predominantly on textual associations rather than direct visual input when predicting human perceptual judgments—highlighting language as a sufficient scaffold for grounding sensory semantics. Bavaresco and Fernández (2025) demonstrate that text alone—when modeled on scale—can implicitly encode rich experiential semantics, echoing Marjieh

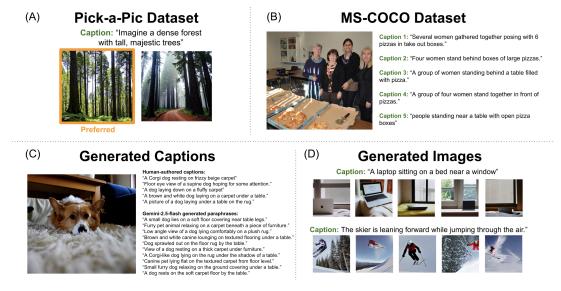


Figure 1: Example data from (**A**) the Pick-A-Pic dataset and (**B**) the MS-COCO dataset. (**C**) Example captions generated by Gemini-2.5-Flash by paraphrasing the human-authored captions in MS-COCO. (**D**) Example MS-COCO captions and synthesized images by the stable-diffusion model.

et al. (2024)'s results on LLMs' ability to recover perceptual hierarchies like the pitch spiral.

Convergent evidence also emerges from neuroscience. Popham et al. (2021) used within-subject fMRI to chart voxel-wise semantic tuning during silent-movie viewing (purely visual) and narrative listening (purely linguistic). They discovered that the two modality-specific maps are topographically contiguous: for every visual category encoded in posterior occipital cortex, a mirror linguistic representation appears immediately anterior to the same cortical border. In other words, visual and linguistic semantics form a single, smoothly joined map that straddles the edge of human visual cortex, implying a tightly aligned cross-modal code rather than two isolated systems. Doerig et al. (2022) asked whether vision already encodes such linguistic semantics. They showed that a vision model trained to translate images directly into sentence embeddings of a language model predicts voxel patterns even better than the embeddings themselves, offering a mechanistic account of how the visual system may recast images into a language-like semantic code by default. Saha et al. (2024) went further, finding that off-the-shelf LLM embeddings sometimes outperform dedicated vision models in explaining activity in high-level visual areas. Together, these findings suggest that the cross-modal alignment observed in artificial networks may reflect, or even recapitulate, the brain's own amodal semantic code.

These findings collectively suggest that modern vision and language models, and possibly even brain systems—like Plato's ideal forms—incrementally discard modality-specific details in favor of a shared, amodal semantic code.

Yet critical gaps remain. First, where along the network hierarchy does this alignment emerge, and is it symmetric across modalities? Second, what visual attributes or linguistic properties drive the effect? Third, all previous demonstrations of crossmodal alignment rely on one-to-one image—text pairs. These analyses inadvertently mask the complexity of real-world semantics where no single description exhausts an image's meaning, and the same sentence can fit many images.

In this study, we fill these gaps through extensive analyses of cross-modal alignment on a broad suite of vision and language models. We map alignment layer-by-layer and probe its dependence on targeted manipulations—semantic (object removal, role shuffling) versus appearance-only. Alignment peaks in mid-to-late layers of both modalities, collapses under semantic changes, and is largely unaffected by superficial appearance edits.

To address the third gap about the many-to-many mapping between images and text, our study employs two complementary analyses that explicitly investigate semantic alignment at a finer granularity using many-to-many mappings. First, using a forced-choice "Pick-a-Pic" task, we show that visual embeddings of human-preferred images

align more closely with the language model embeddings of the caption than non-preferred images. Second, for the same image, we analyze pairs of captions selected based on high and low CLIP-scores—previously validated as proxies for human preferences—and observe analogous alignment patterns. These results indicate that vision and language models converge on a common semantic ground that reflects subtle distinctions aligned with human judgments.

In our second analysis, we investigate the impact of aggregating embeddings across multiple images associated with a single caption and vice versa. Contrary to the intuitive expectation that averaging embeddings would diminish representational specificity, we discover that such aggregation consistently enhances alignment. This suggests that rather than blurring distinctions, averaging distills a more stable, modality-independent semantic core shared across representations. Together, our findings reveal that examining many-to-many correspondences offers richer insights into cross-modal alignment, highlighting a robust convergence toward a shared conceptual space that captures subtle and complex semantic relationships.

We summarize our contributions as follows:

- Layer-wise alignment. Alignment strengthens with depth and converges toward a shared space, asymmetrically across modalities.
- Semantic sensitivity. Alignment drops under semantic edits but is robust to appearance-only changes.
- **Human consistency.** The aligned space reflects human preferences (Pick-a-Pic; CLIP-ranked captions).
- Embedding aggregation. Averaging embeddings across captions/images improves alignment, indicating a modality-independent semantic core.

2 Methods

We compare image representations from large vision models with textual representations of the same images from large language models. For vision models, we employed Vision Transformers (ViTs) trained via DINOv2 (Oquab et al., 2023) on the LVD-142M dataset. DINOv2 learns rich visual representations by solving a self-distillation task

where a student network is trained to match the output distribution of a teacher network (an exponential moving average of the student) while viewing different augmented versions of the same image. For language models, we employed BLOOM (Big-Science et al., 2022), a decoder-only transformer-based architecture trained on a massive multilingual corpus, and OpenLLaMA, an open-source reproduction of the LLaMA model trained on publicly available datasets(Geng and Liu, 2023). Multiple model sizes were selected from repositories including Huggingface (Wolf et al., 2019) and Py-Torch Image Models (TIMM) (Wightman, 2021).

Beyond the primary models, we evaluate three additional LLM families (Qwen, Phi-3, SmolLM) for all core analyses; see Appendix E.

For images, the class token from the penultimate transformer block is used; for language, token activations are averaged from the same layer.

Two datasets are employed:

- Pick-A-Pic: An open dataset of over 500K human preference judgments on text-to-image outputs, collected from 37K real-user prompts; each prompt has two generated images and a binary/tie preference label (Kirstain et al., 2023). Here, we randomly sample 1,000 prompt-image-pair judgments for analysis (Figure 1A).
- MS-COCO: An image captioning dataset of 123K natural images depicting complex every-day scenes, each annotated with five human-authored captions (Lin et al., 2014). Here, we randomly sample 1,000 images (and their associated captions) from the official validation split (Figure 1B).

To assess dataset generalization, we also replicate on Flickr8k for the core analyses; see Appendix D.

Computing Alignment

To quantify alignment between representations from language and vision models, we use *linear predictivity*. For each pair of representations, $\mathbf{X} \in R^{n \times d_X}$ (e.g., from a vision model) and $\mathbf{Y} \in R^{n \times d_Y}$ (e.g., from a language model), we fit a ridge regression from \mathbf{X} to \mathbf{Y} :

$$\hat{\mathbf{W}} = \arg\min_{\mathbf{W}} \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_{2}^{2} + \lambda \|\mathbf{W}\|_{F}^{2}, \quad (1)$$

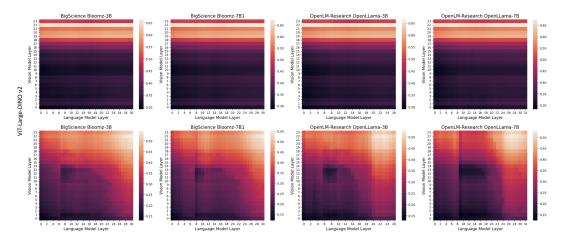


Figure 2: Layer-wise alignment, measured suing linear predictivity score, between one example vision model (ViT-Large-Dinov2) and all language models. Top row: Alignment computed in language-to-vision direction. Bottom row: Alignment computed in vision-to-language direction.

with λ chosen by cross-validation over $10^{-8}...10^{8}$. Alignment is the average Pearson correlation between predicted and actual responses across all units and five cross-validation folds.

We treat this as an *asymmetric* similarity measure and report results for both directions: predicting language representations from vision $(X \to Y)$ and vice versa $(Y \to X)$. This allows us to disentangle directional differences in information content across modalities.

For metric robustness, core analyses are repeated with *CKA* (Appendix A.1) on the same held-out items and layers (Appendix C).

3 Results

We examine (i) layer-wise alignment, (ii) input manipulations, (iii) human preference, and (iv) embedding averaging. Replications hold under CKA and on Flickr8k within scope, and across additional LLM families; see Appendix C, D, E.

3.1 Layer-wise vision-language alignment

To pinpoint where vision-language alignment first appears and how it evolves across the network hierarchy, we performed a layer-by-layer mapping between each pair of vision-transformer and language-model embeddings. As shown in Figure 2, both modalities exhibit low cross-modal predictivity in their earliest layers and increase through the mid and later layers. These patterns hold consistently across different vision-language model pairs (Figure 13). These findings demonstrate that both vision and language models transition from modality-bound encoding toward an abstract,

shared semantic space as depth increases.

We also observe a clear directional asymmetry in these mappings. When mapping from language to vision, we find that even early language layers can successfully predict later vision layers. In contrast, mapping from vision to language reveals a more graded effect: deeper vision layers progressively yield higher predictivity for deeper language layers. Early vision features poorly predict any language layer, while later vision representations align best with higher language layers. This asymmetry suggests that textual representations abstract away from surface form more rapidly than visual ones, while vision networks require deeper processing to reach a comparable semantic level.

3.2 Semantic content, not surface form, drives cross-modal alignment

We next explore whether the cross-modal correspondence we observe is mainly driven by surface form or by deeper semantic content.

3.2.1 Image manipulations

To dissociate appearance-level similarity from semantic correspondence, we performed four controlled perturbations on each MS-COCO image. Two manipulations altered only the appearance while preserving the full meaning: (i) conversion to grayscale and (ii) 15 degree image rotation. The other two manipulations altered the semantic content with different degrees by exploiting the segmentation masks (Figure 3A) provided with COCO-Stuff (Caesar et al., 2018):

• Thing-only views that preserve instances of

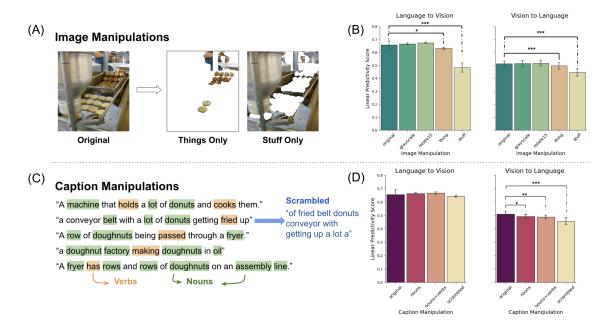


Figure 3: (A) Example thing-only and stuff-only images by manipulating the original image using masks from COCO-Stuff. (B) Alignment by image manipulations. (C) Demonstration of image manipulations: nouns and verbs extraction, and captions scrambling. (D) Alignment by caption manipulations. Paired t-tests (n=8 vision-language model pairs per comparison) were conducted separately for each image manipulation, and p-values were adjusted for four comparisons per mapping direction using the Benjamini–Hochberg procedure (FDR).

the foreground object classes (e.g. person, car) but remove the surrounding context to eliminate spatial and contextual relations;

 Stuff-only views that retain only the background layout and the scene categories (e.g. grass, wall) while removing the foreground objects.

We find that appearance-only manipulations of image inputs have no notable negative effects on alignment (Figure 3B, grayscale: L \rightarrow V: t(7) =-0.8405, p = 0.4284, q = 0.4284; $V \rightarrow L$: t(7) =-1.3386, p = 0.2226, q = 0.2543; rotation: $L \rightarrow V$: t(7) = -1.7569, p = 0.1224, q = 0.1631; $V \rightarrow L$: t(7) = -3.1161, p = 0.0169, q = 0.0271). In contrast, deleting semantic content from images results in substantial alignment degradation (Figure 3B). Isolating only the foreground "thing" pixels and removing contextual relations significantly lowered the alignment scores (L \rightarrow V: t(7) = 3.4304, $p = 0.0110, q = 0.0220; V \rightarrow L: t(7) = 7.2528,$ p = 0.0002, q = 0.0005). Retaining only the "stuff" background further reduced the alignment $(L \rightarrow V: t(7) = 10.1267, p < 0.0001, q = 0.0001;$ $V \rightarrow L$: t(7) = 11.7109, p < 0.0001, q = 0.0001).

Notably, the decline was systematically steeper in the language-to-vision direction, indicating that mapping from textual embeddings to visual layers depends more heavily on intact visual semantics.

3.2.2 Caption manipulations

To explore the linguistic properties driving the alignment, we separately manipulated the captions in the MS-COCO dataset with different levels of semantic disruption by retaining: (i) nouns only, (ii) nouns and verbs, and (iv) all the words but in scrambled order (Figure 3C; Appendix A.2).

Interestingly, only in the vision-to-language mapping direction do caption manipulations negatively affect the alignment (Figure 3D, right). Specifically, nouns-only ($t(7)=3.5956,\ p=0.0088,\ q=0.0176$) and nouns+verbs ($t(7)=5.3561,\ p=0.0011,\ q=0.0032$) show similar moderate decreases, while scrambled captions produce the largest drop ($t(7)=22.8176,\ p<0.0001$, q<0.0001). This suggests that nouns and verbs carry the primary semantic weight in grounding language to visual content, while word order and the full lexical distribution become even more crucial when projecting from vision to language.

The directional asymmetries we observe—greater sensitivity of language→vision mapping to intact visual semantics and of vision→language mapping to linguistic composition—suggest complementary organizational

principles in how each modality abstracts and transmits meaning across the shared representational space.

3.3 Vision-language alignment mirrors human preferences

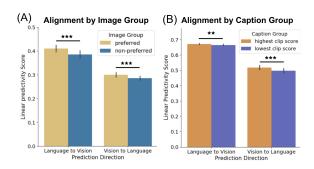


Figure 4: (A) Pick-a-Pic dataset Linear predictivity scores grouped by image variation (preferred vs. non-preferred) based on human judgments. (B) MS-COCO dataset Linear predictivity scores grouped by caption variation based on CLIP Scores. Error bars indicating the standard error across model pairs.

We next evaluated whether cross-modal alignment tracks fine-grained human preferences. Images from the "Pick-a-Pic" dataset, which provides two generated images for the same prompt with human preference judgments, were grouped into high-and low-preference categories. For each group, vision model representations were extracted and linear predictivity scores were computed using the corresponding caption embeddings. This design probes alignment at a finer-grained resolution: can the vision—language mapping replicate the subtle distinctions that lead people to prefer one image over another, even when the linguistic description is identical?

Our results indicate that images preferred by human raters exhibit significantly stronger alignment with their associated captions than non-preferred images across all vision-language model pairs (paired t-test, L \rightarrow V: t(7)=19.8225, p<0.001; V \rightarrow L: t(7)=10.2338, p<0.001; Figure 4A). In other words, even when two pictures illustrate the same text, the uni-modal vision and language models collectively "agree" with human raters about which picture is the better semantic fit. This fine-grained sensitivity shows that the cross-modal alignment we measure is not a coarse correlation but captures subtle, human-relevant distinctions within a shared semantic space.

A complementary analysis from the text side re-

inforces this conclusion. We computed the CLIP Score (Hessel et al., 2021)—a reference-free metric based on the cosine similarity of image–caption embeddings—for all MS-COCO captions, as a reasonable proxy for human preferences (Hessel et al., 2021). Our analysis reveals that captions with higher CLIP scores are significantly more aligned with their images than those with lower scores (paired t-test, language-to-vision: t(7) = 3.9231, p = 0.0057; vision-to-language: t(7) = 17.8350, p < 0.001; Figure 4B).

Thus, across seven vision—language model pairs evaluated on MS-COCO and Pick-a-Pic, the model embeddings capture fine-grained semantic distinctions that mirror human evaluative patterns.

3.4 Averaging embeddings across multiple captions and images enhances alignment

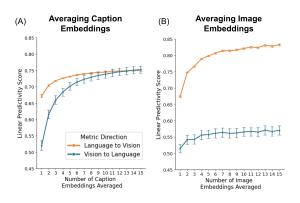


Figure 5: Effect of aggregation on alignment. Cross-modal aggregation: Averaging (A) multiple caption embeddings for the same image or (B) multiple image embeddings for the same caption steadily increases language—vision and vision—language predictivity. Error bars denote standard error across all model pairs.

To quantify the impact of aggregating caption representations, we progressively averaged embeddings from an increasing number of MS-COCO captions per image and computed cross-modal alignment scores. As shown in Figure 5A, alignment improved monotonically with each additional caption. To locate the point of diminishing returns, we expanded the caption pool by paraphrasing each of the five human-authored captions with Gemini-2.5-Flash (Figure 1C, see Appendix A.3 Table 1 for prompt), creating up to 15 captions per image. In the vision-to-language mapping, alignment continued to rise until roughly ten captions were included, after which the curve plateaued.

We performed the complementary analysis in

the opposite direction by synthesizing up to 15 naturalistic images per caption with Stable Diffusion (Figure 1D). Similar to caption aggregation, increasing the number of aggregated image embeddings further improved the alignment (Figure 5B). The alignment gain is larger when predicting vision from language, and plateaued around seven images.

To confirm that these improvements reflect enhanced semantic information rather than a generic averaging artifact, we repeated both analyses after randomly shuffling the image—caption correspondences (Appendix B.1, Figure 8). Under this mismatch baseline, embedding aggregation showed no benefit, demonstrating that the effect depends on semantically matched pairs.

We also observe a clear directional asymmetry both analyses: averaging captions benefited vision-to-language predictions, whereas averaging images benefited language-to-vision predictions. This pattern suggests that aggregation may suppress modality-specific noise within the averaged domain, exposing a cleaner semantic signal that is more easily mapped by the other modality.

3.5 Effect of vision models on vision-language alignment

To assess the generalizability of our findings, we repeated the analyses on seven ViT backbones that differ in objective (strong AugReg, DINO, large-scale DINOv2, supervised distillation DeiT), data scale (ImageNet-1k vs. ImageNet-21k vs. LVD-142 M), and model size (ViT-B/14, ViT-B/16, ViT-L/14, ViT-L/16).

We observe that the improvement of averaging caption embeddings is generalized across different vision model backbones (Figure 6). Notably, when mapping language features into visual space, the alignment differences scores across ViTs were noticeably larger than in the reverse direction. Furthermore, both training methods and data size appear to affect the alignment. When the model size and data were held constant (ViT-B/16, ImageNet-1k), AugReg produced higher alignment than either DINO or DeiT. Keeping the objective similar but increasing the dataset (DINO-ImageNet1k to DINOv2-LVD142m) improved alignment further. However, a larger dataset did not help the AugReg model: its ImageNet-21k checkpoint aligned worse than its ImageNet-1k counterpart. Our current experiment cannot cleanly disentangle the interaction between objective and data distribution. A systematic experiment would be needed to clarify such interactive effects.

4 Discussion

Our results provide new evidence that purely unimodal vision and language models gravitate toward a common semantic manifold. Alignment (i) peaks in their mid-to-late layers where abstract semantic processing occurs, (ii) reduces when we remove or scramble semantic content but survives appearance-only changes, and (iii) exhibits striking correspondence in fine-grained evaluation scenarios with human judgements (e.g., when comparing alignment scores for multiple candidate images corresponding to the same linguistic expression, the model aligns most strongly with the image humans rate as most semantically congruent with the text, and reciprocally for multiple linguistic descriptions of the same image), and (iv) is markedly enhanced when averaging representations corresponding to the same concept in each modality. Together, these findings refine the emerging "Platonic" view of cross-modal representation: the two modalities do not merely share coarse alignment but capture fine semantic gradients that track human judgments. Our work bridges cognitive science and machine learning by suggesting that a shared code for meaning can emerge implicitly in unimodal systems, even without cross-modal training.

Our work opens several promising avenues for future research. Future studies should investigate how alignment strength varies across different types of visual and linguistic content. Are concrete concepts (e.g., "dog", "chair") more strongly aligned than abstract concepts (e.g., "freedom", "justice")? Understanding these variations could reveal fundamental constraints on cross-modal convergence. Different image types—photographs, illustrations, diagrams, artistic renderings—may exhibit varying degrees of alignment with language. Examining these differences could illuminate how visual style and abstraction influence semantic encoding and cross-modal correspondence.

Our discovery that alignment strengthens when averaging concept-specific representations raises intriguing questions about the geometric properties of these embeddings. Future work should explore whether averaging acts as a denoising mechanism that preserves core semantic content while reducing modality-specific variations. Additionally, it would be interesting to investigate whether averag-

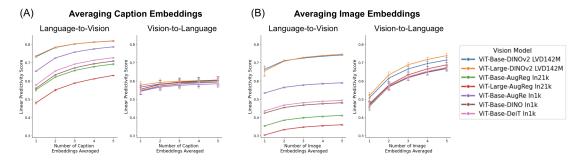


Figure 6: Comparing the alignment of different vision models with language models after averaging (A) caption embeddings and (B) image embeddings.

ing techniques applied to paraphrases of the same linguistic expression could enhance performance on downstream tasks involving natural language inference. Prompt ensembling in CLIP—where averaging multiple engineered text prompts boosts zero-shot accuracy in a multimodal model (Radford et al., 2021)—offers a useful parallel.

While our study demonstrates alignment at the representation level, identifying which specific features or dimensions drive this alignment remains an open question. Future research should develop techniques to isolate the most aligned dimensions between vision and language models and analyze their semantic properties.

Further, investigating how alignment patterns evolve during training could provide insights into the developmental trajectory of cross-modal correspondence. Do alignment patterns appear early in training and strengthen over time, or do they emerge suddenly after sufficient exposure to domain-specific data? This temporal perspective could reveal fundamental insights about how semantic convergence develops in neural networks trained on different modalities.

Limitations

Our analysis primarily relies on linear predictivity, complemented by CKA to verify robustness. While these provide complementary perspectives, they still represent only a subset of possible approaches to assessing representational alignment. Future work could benefit from employing a broader spectrum of alignment metrics to provide a more complete understanding of vision-language relationships. For instance, more constrained mapping approaches—such as orthogonal transformations in Procrustes analysis (Williams et al., 2021) or permutation-based methods like permutation score and soft matching score (Khosla

and Williams, 2024)—might reveal unit-level correspondences between visual and language model representations that linear regression or CKA cannot capture. Other families of metrics, including Representational Similarity Analysis (Kriegeskorte et al., 2008) for population-level relationships and neighborhood-based approaches (e.g., mutual k-NN) for local structure, could further enrich the picture. These complementary metrics would provide a multi-faceted view of the nature of alignment between vision and language models. Our analysis does not fully reveal which specific features drive the observed alignment between vision-only and language-only models, nor does it identify the scenarios where these models systematically diverge in their representations. Investigating these questions would require more extensive probing of representations across diverse stimuli and large-scale datasets.

The synthetic nature of our image dataset introduces another limitation. While diffusion models generate high-quality images corresponding to text prompts, some generated images may not perfectly capture the semantic content or nuances present in the texts. This potential mismatch between text and generated images could influence our alignment measurements and subsequent interpretations.

Furthermore, our work examines models trained at a specific point in time, with particular architectures and training objectives. As model architectures and training paradigms evolve, the nature of cross-modal alignment may change significantly.

Finally, representational similarity is descriptive. It does not prove shared processing mechanisms or functional interchangeability. Causal interventions are needed to determine whether the aligned dimensions are necessary for each model's downstream behavior.

References

- M. Abdin and 1 others. 2024. Phi-3 technical report: A highly capable language model locally on your phone. arXiv preprint arXiv:2404.14219. 3.8B "mini" model, plus larger phi-3 variants.
- Anna Bavaresco and Raquel Fernández. 2025. Experiential semantic information and brain alignment: Are multimodal models better than language models? *arXiv preprint arXiv:2504.00942*.
- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgren, Xuan-Son Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, and 2 others. 2025. Smollm2: When smol goes big data-centric training of a small language model. arXiv preprint arXiv:2502.02737. "Small" (1.7B parameter) LM trained over 11 trillion tokens using multi-stage process.
- BigScience, T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, and 1 others. 2022. Bloom: A 176bparameter open-access multilingual language model. *Preprint*, arXiv:2211.05100.
- Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. 2018. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1209–1218.
- Adrien Doerig, Tim C Kietzmann, Emily Allen, Yihan Wu, Thomas Naselaris, Kendrick Kay, and Ian Charest. 2022. Semantic scene descriptions as an objective of human vision. *arXiv preprint arXiv:2209.11737*, 10.
- Jerry A Fodor. 1975. *The language of thought*, volume 5. Harvard university press.
- X. Geng and H. Liu. 2023. Openllama: An open reproduction of llama. https://github.com/ openlm-research/open-llama.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. Clipscore: A reference-free evaluation metric for image captioning. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP).
- Micah Hodosh, Peter Young, and Julia Hockenmaier. 2013a. Flickr8k image captioning dataset. Downloaded from Kaggle. 8,000 images; five captions per image.
- Micah Hodosh, Peter Young, and Julia Hockenmaier. 2013b. Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 47:853–899.

- Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. 2024. The platonic representation hypothesis. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- Meenakshi Khosla and Alex H Williams. 2024. Soft matching distance: A metric on neural representations that captures single-neuron tuning. In *Proceedings of UniReps: the First Workshop on Unifying Representations in Neural Models*, pages 326–341. PMLR.
- A. Kirstain and 1 others. 2023. Pick-a-pic: A dataset for evaluating the robustness of vision-language models. *Preprint*, arXiv:2305.01569.
- Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. 2023. Grounding language models to images for multimodal inputs and outputs. In *International Conference on Machine Learning*, pages 17283–17300. PMLR.
- Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. 2019. Similarity of neural network representations revisited. In *Proceedings of the 36th International Conference on Machine Learning*, pages 3519–3529. PMLR.
- Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bandettini. 2008. Representational similarity analysis-connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, 2:249.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. *Preprint*, arXiv:1405.0312.
- Mayug Maniparambil, Raiymbek Akshulakov, Yasser Abdelaziz Dahou Djilali, Mohamed El Amine Seddik, Sanath Narayan, Karttikeya Mangalam, and Noel E O'Connor. 2024. Do vision and language encoders represent the world similarly? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14334–14343.
- Raja Marjieh, Ilia Sucholutsky, Pol van Rijn, Nori Jacoby, and Thomas L Griffiths. 2024. Large language models predict human sensory judgments across six modalities. *Scientific Reports*, 14(1):21445.
- Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. 2022. Linearly mapping from image to text space. *arXiv preprint arXiv:2209.15162*.
- M. Oquab, T. Darcet, T. Moutakanni, H. V. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, R. Howes, P.-Y. Huang, H. Xu, V. Sharma, S.-W. Li, W. Galuba, M. Rabbat, M. Assran, N. Ballas, and 7 others. 2023. Dinov2: Learning robust visual features without supervision. Preprint.

Sara F Popham, Alexander G Huth, Natalia Y Bilenko, Fatma Deniz, James S Gao, Anwar O Nunez-Elizalde, and Jack L Gallant. 2021. Visual and linguistic semantic representations are aligned at the border of human visual cortex. *Nature neuroscience*, 24(11):1628–1636.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR.

Shreya Saha, Ishaan Chadha, and 1 others. 2024. Modeling the human visual system: Comparative insights from response-optimized and task-optimized vision models, language models, and different readout mechanisms. *arXiv* preprint arXiv:2410.14031.

Ross Wightman. 2021. Pytorch image models. https://github.com/rwightman/pytorch-image-models.

Alex H Williams, Erin Kunz, Simon Kornblith, and Scott Linderman. 2021. Generalized shape metrics on neural representations. *Advances in Neural Information Processing Systems*, 34:4738–4750.

T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, and 1 others. 2019. Huggingface's transformers: State-of-the-art natural language processing. *Preprint*, arXiv:1910.03771.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 23 others. 2024. Qwen2.5 technical report. arXiv preprint arXiv:2412.15115. "A comprehensive series of large language models", includes open-weight and instruction tuned models; 0.5B - 72B parameters.

5 Appendix

A Methods

A.1 Alignment metrics: implementation details.

Linear Predictivity. For each pair of vision model and language model, we construct $\mathbf{X} \in R^{N \times d_X}$ and $\mathbf{Y} \in R^{N \times d_Y}$ where N=1,000 image—caption combinations in a dataset (e.g. MS-COCO-val2017 or Pick-a-Pic), d_X is the language-feature dimensionality (e.g. 2560 for BLOOMZ-3B), and d_Y is the vision-feature dimensionality (e.g. 1024 for ViT-Large14-DINOv2).

A ridge map is fit with 5-fold cross-validation (outer KFold with shuffling) to pick the best λ on the training data, where $\lambda \in \{10^{-8}, 10^{-7}, \dots, 10^{8}\}$ (17 logarithmically spaced values). Features are z-scored using training statistics and the transform is applied to the test split.

On the test split, $\hat{\mathbf{Y}} = \mathbf{X}\hat{\mathbf{W}}$ is evaluated by Pearson r between $\hat{\mathbf{Y}}$ and \mathbf{Y} , averaged over target dimensions; the final score is the mean across folds. Both directions are reported ($\mathbf{X} \rightarrow \mathbf{Y}$ and $\mathbf{Y} \rightarrow \mathbf{X}$).

Centered Kernel Alignment (CKA). CKA (Kornblith et al., 2019) measures similarity between representational spaces in a way that is invariant to isotropic scaling and orthogonal transformations, and is symmetric between modalities. We applied CKA to the same vision–language layer pairs as in the linear predictivity analysis.

For each vision-language model pair, item×feature matrices $\mathbf{X} \in R^{N \times d_X}$ and $\mathbf{Y} \in R^{N \times d_Y}$ are constructed over the same heldout items. With kernels κ_X, κ_Y and Gram matrices $K_{ij} = \kappa_X(x_i, x_j)$ and $L_{ij} = \kappa_Y(y_i, y_j)$, CKA is computed via the (biased) HSIC normalization:

$$\mathrm{CKA}(K,L) \; = \; \frac{\mathrm{HSIC}(K,L)}{\sqrt{\mathrm{HSIC}(K,K)\,\mathrm{HSIC}(L,L)}},$$

 $HSIC(K,L)=1N^2tr(\tilde{K}\tilde{L}),$

where $\tilde{K} = HKH$, $\tilde{L} = HLH$, and $H = I - 1N\mathbf{1}\mathbf{1}^{\mathsf{T}}$ is the centering matrix. In the *linear* case used here, $\kappa_X(u,v) = u^{\mathsf{T}}v$ and $\kappa_Y(u,v) = u^{\mathsf{T}}v$, so $K = \mathbf{X}\mathbf{X}^{\mathsf{T}}$ and $L = \mathbf{Y}\mathbf{Y}^{\mathsf{T}}$ (double-centered by H inside HSIC).

A.2 Caption manipulation procedure

Captions are tokenized and *part-of-speech (POS)* tagged using spaCy (en_core_web_sm). We use spaCy's Universal POS (UPOS) labels

(Token.pos_) to create two filtered variants per caption: N keeps only tokens labeled {NOUN}, and NV keeps {NOUN, VERB}. All other tokens are removed, and the remaining tokens are rejoined with single spaces.

A scrambled baseline (random permutation of the word tokens within a caption, fixed seed) is implemented separately.

A.3 MS-COCO caption generation.

Caption paraphrases for MS-COCO were generated using Gemini-2.5-Flash to support the embeddingaveraging analyses (§3.4). For each image, the five human captions are provided as context, and the model is asked to produce 10 new captions that preserve meaning while varying wording and surface form (Table 1).

Gemini-2.5-Flash Prompt

f"""You are an expert image captioner. I'll show you some existing captions for an image, and your task is to generate 10 NEW captions that:

- 1. Are similar in style and detail level to the existing captions
- 2. Capture the same meaning but with different wording
- 3. Are direct, concise descriptions (around 10-15 words each) 4. Are worded differently from each existing caption and from
 - each other

Here are the existing captions: {insert all captions text for the image here}

Generate 10 new captions formatted exactly as

- Ifirst new caption
- [Second new caption]
- [Third new caption] [Fourth new caption]
- [Fifth new caption]
 [Sixth new caption]
- [Seventh new caption]
- [Eighth new caption]
- [Ninth new caption] 10. [Tenth new caption]"""

Table 1: Prompt used to generate MS-COCO caption paraphrases with Gemini-2.5-Flash.

A.4 MS-COCO image generation.

Synthetic images for MS-COCO are generated with the Diffusers StableDiffusionPipeline initialized from CompVis/stable-diffusion-v1-4. Each caption text of an MS-COCO image is used as the prompt and K=2 variants are sampled per caption with num_inference_steps = 50, yielding 10 synthesized images per MS-COCO image.

Extended analyses on MS-COCO with Linear Predictivity

B.1 Baseline alignment with shuffled image-caption correspondences.

Under the image-caption mismatch baseline, averaging multiple embeddings does not improve vision-language alignment: the alignment score remains around 0 in both mapping directions (Figure 8).

B.2 Embedding aggregation effect on manipulated captions.

Given that averaging caption embeddings enhances vision-language alignment, we also explored whether the embeddings of semantically manipulated captions would also benefit from embedding aggregation (Figure 9). Interestingly, the alignment was enhanced even though the embeddings come from manipulated captions.

Alternative metric: CKA

To assess metric robustness, we replicated our core analyses with linear CKA on the same held-out items, models, and layers. Because CKA is symmetric, it does not encode the $L\rightarrow V$ vs. $V\rightarrow L$ directionality.

CKA reproduces the qualitative patterns (Figure 7A-C). For captions, semantic disruptions (nounsonly, nouns+verbs, and scrambled) reduce alignment, and for images, retaining only stuff (things masked out) yields a significant decrease relative to originals, whereas retaining only things (stuff masked out) shows a smaller, non-significant decrease (Figure 7A). Human-preferred pairs have higher alignment (Figure 7B). Lastly, averaging embeddings increases alignment, although the gain is weaker yet persistent for image-embedding averaging (Figure 7C).

Additional datasets: Flickr8k D

Using the same procedure as the main text (ridgebased linear predictivity), we test dataset generalization on an additional dataset—Flickr8k (Hodosh et al., 2013b)—a captioning dataset of 8,000 images, each annotated with five humanauthored captions (Hodosh et al., 2013a,b). Because Flickr8k lacks instance-level segmentations, we replicate only the core analyses that do not require foreground/background masks.

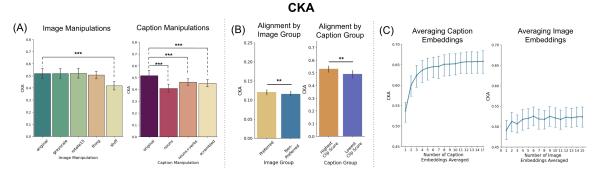


Figure 7: CKA replication across analyses. (A) Manipulations on images (left) and captions (right). (B) Preference. Left: Pick-a-Pic pairs, preferred vs. non-preferred. Right: MS-COCO comparison of the highest- vs. lowest-CLIP-similarity caption group. (C) Embedding averaging.

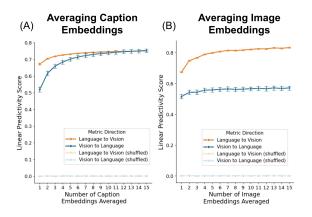


Figure 8: Effect of aggregation on alignment with a mismtach baseline.

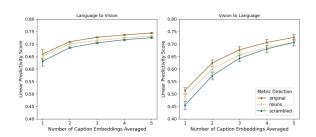


Figure 9: Effect of aggregation on alignment with manipulated captions which either only includes nouns or are scrambled in word order.

We randomly sample 1,000 images (and their associated five captions) for analysis. We also generate caption paraphrases and synthesized image variants using the same protocols as in the main text to evaluate embedding averaging.

The main patterns replicate for Flickr8k: alignment increases from mid to late layers and the L \rightarrow V > V \rightarrow L asymmetry holds (Figure 14). The CLIP-based preference proxy yields higher alignment for higher-ranked captions (paired t-test, L \rightarrow V: t(7)=

5.0520, p = 0.0015; V \rightarrow L: t(7) = 13.8867, p < 0.0001; Figure 10); and averaging multiple caption/image embeddings improves alignment and plateaus as the number of embeddings increases (Figure 11).

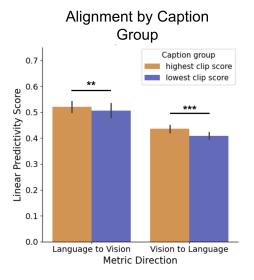


Figure 10: Flickr8k dataset linear predictivity scores grouped by caption variation based on CLIP Scores.

E Expanded LLM families: Qwen, Phi-3, SmolLM

Using the same procedure as the main text, we evaluate three additional LLM families—Qwen (3B, 7B) (Yang et al., 2024); Phi-3 (mini) (Abdin et al., 2024), and SmolLM (1.7B) (Ben Allal et al., 2025)—on the same analyses (Figure 12A–D). The main patterns replicate across LLM families.

Layer-wise alignment. The mid-to-late layer rise holds across families, and the $L\rightarrow V>V\rightarrow L$ asymmetry is preserved (Figure 15).

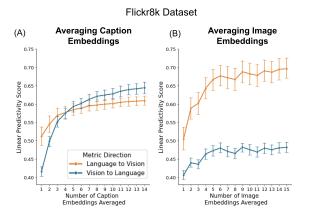


Figure 11: Effect of aggregation on alignment on Flickr8k.

Image manipulations. Retaining only "things" foreground and retaining only "stuff" background both significantly reduce alignment. A small rotation does not yield a reliable decrease, whereas converting images to grayscale produces a small but significant reduction (Figure 12A, grayscale: L→V $t(7) = 8.1393, p = 0.0001, q = 0.0002; V \rightarrow L$ t(7) = 2.7209, p = 0.0297, q = 0.0396; rotate15: $L\rightarrow V t(7) = -1.1300, p = 0.2957, q = 0.3379;$ $V \rightarrow L t(7) = -0.5698, p = 0.5866, q = 0.5866;$ stuff-only: L \rightarrow V t(7) = 12.8386, p < 0.0001, q < 0.0001; V \rightarrow L t(7) = 12.1947, p < 0.0001, q < 0.0001; things-only: L \rightarrow V t(7) = 21.4355, $p < 0.0001, q < 0.0001; V \rightarrow L t(7) = 5.5535,$ p = 0.0009, q = 0.0014). A plausible interpretation is that color carries captioned semantics (e.g., color words in text), and these families are sensitive to that finer-grained correspondence.

Caption manipulations. All caption manipulations (nouns-only, nouns+verbs, scrambled) produce significantly lower alignment than the original captions in each family (Figure 12B; nouns: L \rightarrow V t(7) = 6.3039, p = 0.0004, q = 0.0008; V \rightarrow L t(7) = 5.3781, p = 0.0010, q = 0.0015; nouns+verbs: L \rightarrow V t(7) = 3.6133, p = 0.0086, q = 0.0086; V \rightarrow L t(7) = 3.6944, p = 0.0077, q = 0.0086; scrambled: L \rightarrow V t(7) = 8.9154, p < 0.0001, q = 0.0001; V \rightarrow L t(7) = 9.8068, p < 0.0001, q = 0.0001.)

Human preference / CLIP proxy. The preference effect replicates across families: preferred > non-preferred image group (paired t-test, L \rightarrow V: t(7) = 14.0585, p < 0.0001; V \rightarrow L: t(7) = 9.1631, p < 0.0001), and high CLIP Score > low CLIP score (paired t-test, L \rightarrow V: t(7) = 3.5055, p = 0.0099; V \rightarrow L: t(7) = 14.4357, p < 0.0001).

(Figure 12C)

Embedding averaging. Averaging embeddings increases alignment and plateaus with larger number of examplars for both caption- and image-embedding averaging (Figure 12D).

Code Availability

Code and scripts are available at https://github.com/zoewhe/vision-language-alignment.

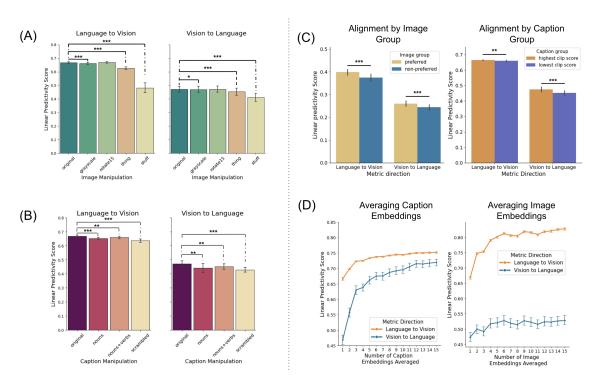


Figure 12: Analysis replication on additional LLM families. (**A**) Image Manipulations. (**B**) Caption Manipulations. (**C**) Preference. Left: Pick-a-Pic pairs, preferred vs. non-preferred. Right: MS-COCO comparison of the highest-vs. lowest-CLIP-similarity caption group. (**D**) Embedding averaging.

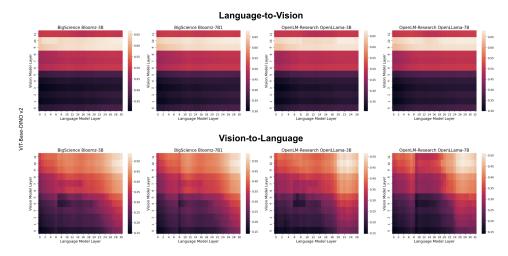


Figure 13: Layer-wise alignment for additional vision-language model pairs (with ViT-Base-DINO v2).

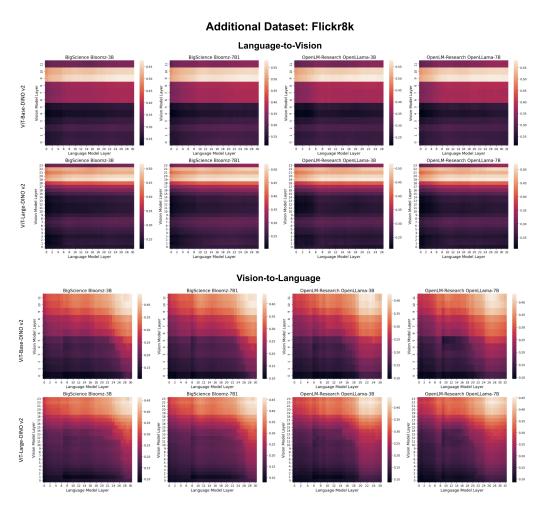


Figure 14: Layer-wise alignment on additional dataset Flickr8k. Top two rows: Alignment computed in language-to-vision direction. Bottom two rows: Alignment computed in vision-to-language direction.

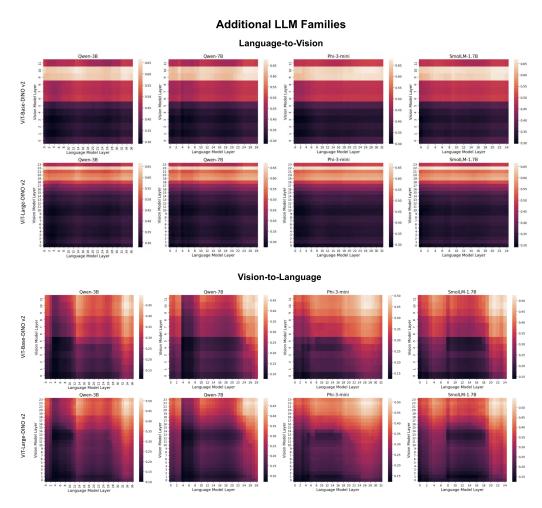


Figure 15: Layer-wise alignment on MS-COCO with additional LLMs: Qwen2.5-3B, Qwen2.5-7B, Phi-3-mini-128k-instruct, and SmolLM2-1.7B. Top two rows: Alignment computed in language-to-vision direction. Bottom two rows: Alignment computed in vision-to-language direction.