



Implicit Multimodal Alignment: On the Generalization of Frozen LLMs to Multimodal Inputs

Mustafa Shukor¹ *

Matthieu Cord^{1,2}

¹Sorbonne University, ²Valeo.ai

Abstract

Large Language Models (LLMs) have demonstrated impressive performance on multimodal tasks, without any multimodal finetuning. They are the *de facto* building block for Large Multimodal Models (LMMs), yet, we still lack a proper understanding of their success. In this work, we expose frozen LLMs to image, video, audio and text inputs and analyse their internal representation aiming to understand their generalization beyond textual inputs. Our work provides the following **findings**. Perceptual tokens (1) are easily distinguishable from textual ones inside LLMs, with significantly different representations (*e.g.* live in different narrow cones), and complete translation to textual tokens does not exist. Yet, (2) both perceptual and textual tokens activate similar LLM weights. Despite being different, (3) perceptual tokens are implicitly aligned to textual tokens inside LLMs, we call this the implicit multimodal alignment effect (IMA), and argue that this is linked to architectural design, helping LLMs to generalize. This provide more evidence to believe that the generalization of LLMs to multimodal inputs is mainly due to their architecture. These findings lead to several **implications**. (1) We find a positive correlation between the implicit alignment score and the task performance, suggesting that this could act as a proxy metric for model evaluation and selection. (2) A negative correlation exists regarding hallucinations (*e.g.* describing non-existing objects in images), revealing that this problem is mainly due to misalignment between the internal perceptual and textual representations. (3) Perceptual tokens change slightly throughout the model, thus, we propose different approaches to skip computations (*e.g.* in FFN layers), and significantly reduce the inference cost. (4) Due to the slowly changing embeddings across layers, and the high overlap between textual and multimodal activated weights, we compress LLMs by keeping only 1 subnetwork (called α -SubNet) that works well across a wide range of multimodal tasks. The code is available here: <https://github.com/mshukor/ima-lmms>.

1 Introduction

Large Language Models (LLMs) [1, 2, 3, 4, 5] represent a noteworthy advancement in recent AI developments. Building upon the success of LLMs, the next stride in this field involves extending beyond the textual modality, giving rise to Large Multimodal Models (LMMs) [6, 7, 8, 4]. A notable line of research involves connecting LLMs to visual encoders, while keeping them frozen and only training a connector with modest number of parameters [9, 10, 11, 12, 13, 14, 15, 16, 17]. These methods yield comparable performance [10] to large-scale multimodal models with significantly reduced computational and data budget.

Keeping all pretrained unimodal models frozen and only training couple of millions of parameters [9, 10, 12] is an interesting phenomenon to understand, with limited research trying to decipher it. To

*Contact: {firstname.lastname}@sorbonne-universite.fr

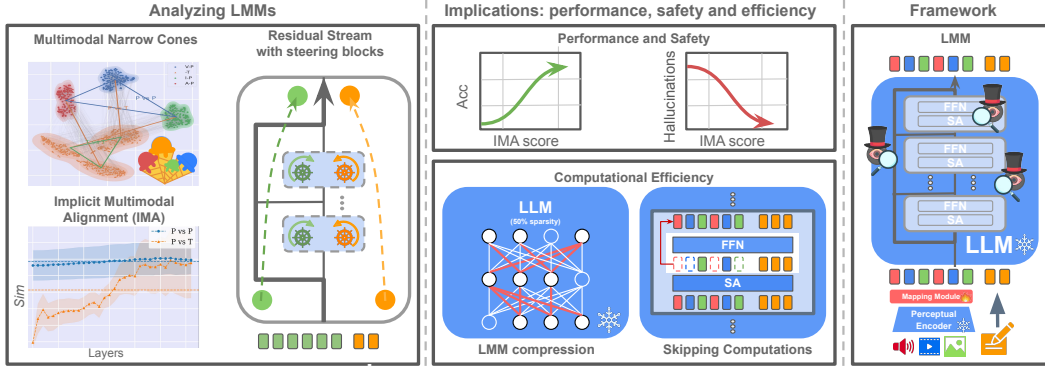


Figure 1: **Summary of the work.** We start by analysing multimodal tokens inside LLMs, and find that they live in different spaces (e.g., multimodal cones). Yet they are implicitly aligned (i.e., IMA), allowing us to see LLMs as residual streams with steering blocks. This lead to implications on performance, safety and efficiency.

explain why frozen LLMs can generalize beyond textual inputs, several hypotheses can be isolated: (1) perceptual tokens are transformed to textual ones, can be simply considered as foreign language [18], and thus the LLM sees text-only tokens. (2) LLMs are able to digest non-textual tokens that are processed by (2a) modality-specific subnetworks or (2b) the same LLM weights that can generalize due to other reasons.

In this study, we expose LLMs to different multimodal inputs, such as image, video, audio and text, and analyse their internal representations. We focus on frozen LLMs and consider two representative setups: single-task (ST) and multitask (MT) finetuning. The former is considered by parameter and data-efficient approaches [11, 12, 9, 10] and consists of training a mapping module for each dataset. The MT setup is considered by recent multimodal assistant models [19, 14, 20, 21], and consists of training the same mapping module on several datasets/tasks.

Our study shows (1) that perceptual and textual tokens still live in significantly different representation spaces inside LLMs (Sec. 3.1): live in different narrow cones, and have different norms, rate of change and vocabulary distributions. (2) We notice high similarity between the weights activated by textual and perceptual tokens (Sec. 3.2), allowing to swap these activated subnetworks between different tasks and modalities. Despite their differences, (3) perceptual tokens are implicitly aligned to textual ones across different stages (Sec. 4.1): during training of the mapping module, and during inference across LLM layers, especially inside each LLM block (e.g. after the self-attention layer). As there is no explicit objective to align these representations, we call it the Implicit Multimodal Alignment effect (IMA) (i.e., increasing similarity between textual and perceptual token distributions). We find that this effect is mostly linked to architectural design (e.g. residual stream with refinement blocks acting as steering blocks Sec. 4.2). This provides more evidence to believe the architecture of LLMs is one of the main factors to generalize to multimodal representations.

We shed light on several practical implications (Sec. 5). (1) We find a positive correlation between the implicit multimodal alignment score and the task performance, suggesting that this score could act as a proxy metric. On the other hand, (2) we find a negative correlation with hallucinations, revealing that the main factor leading to this problem is the lack of alignment between the internal representation of textual and perceptual inputs. (3) The perceptual tokens slightly change inside LLMs, thus, we propose to skip their computations (e.g. inside FFN layers). (4) Due to the slowly changing embedding across layers, and the high overlap between weights activated by different modalities, we compress the LLM by keeping one task-agnostic subnetwork that works well across all modalities. To summarize (Fig. 1), we analyse the internal representations of LLMs when exposed to multimodal inputs, leading to the following **findings**:

- Perceptual and textual tokens live in different representation spaces inside LLMs.
- They activate similar LLM weights.
- They are implicitly aligned (IMA) inside LLMs, during training and during inference.

- The architectural design of LLMs can be perceived as a residual stream with steering blocks. We argue that this is one of the main factors allowing LLMs to: digest very different tokens, drive the implicit multimodal alignment effect, and thus generalize to different modalities.

These findings have several practical **implications** such as:

- The IMA score as a proxy metric candidate for task performance and hallucinations.
- Hallucinations as a result of lack of sufficient multimodal alignment.
- Skipping computations for visual tokens, leading to efficient inference.
- LLMs compression by keeping only 1 subnetwork that generalizes to all multimodal tasks.

2 Framework for analysing perceptually augmented LLMs

General framework We focus on a general family of models that consists of: a frozen language model LLM with L layers, a trainable mapping module C , and a frozen perceptual encoder E_M for different modalities M (e.g. image (I), video (V) and audio (A)). The LLM input X consists of the concatenation of $P = [p_1, \dots, p_{N_p}]$ multimodal/perceptual tokens (referred to as prompt) with $T = [t_1, \dots, t_{N_t}]$ textual tokens. The prompt P is obtained after encoding the modality-specific input XM with the corresponding E_M and using C to project it to the LLM input space. T is obtained from the embedding layer E_T applied to the tokenized input text XT . This can be expressed as follows:

$$P = C(E_M(XM)), \quad T = E_T(XT), \quad O = LLM([P; T]). \quad (1)$$

The k ($k = N_p + N_t$) output tokens $O = [o_1, \dots, o_k]$ are obtained after a normalization, followed by the unembedding layer W_{out} (or LLM head, i.e. $o_i = W_{out}LN_{out}(t_i^L)$). Our focus is on the internal representation of LLMs (i.e. tokens) at different stages, in particular across the L LLM blocks/layers (referred to as B). The mechanism inside the $l + 1$ LLM transformer block can be expressed as follows:

$$X^{l+1} = X_{SA} + FC2(g(FC1(LN2(X_{SA}))), \quad X_{SA} = X^l + SA(LN1(X^l)), \quad (2)$$

where $FC1$, $FC2$, g are the up and down projections and activation inside the FFN , $LN1/2$ are the layer norms and SA the self-attention.

Perceptually augmented LLM baselines. For the single-task (ST) setup, we train many models across different datasets that span image, video and audio-text modalities. Each mapping module is trained on a specific dataset, similar to previous works [10, 12, 9]. Inspired by previous studies [10, 12], we use light-weight transformer consisting of a self-attention to attend to perceptual tokens. In this setup, P refer to perceptual tokens from image, video and audio modalities. For the multitask (MT) setup, we devise different variants of the LLaVA-1.5 [19] model that differ from the original model as follows: LLaVA-1.5-2 (LLM kept frozen), LLaVA-1.5-3 (LLM kept frozen, without pretraining) and LLaVA-1.5-4 (LLM kept frozen, without pretraining and with transformer mapping module similar to the ST setup instead of MLP). In this setup, P refers to image tokens from different datasets. In the paper, we focus on LLaVA-1.5-4 as it is most similar to the ST setup, and analyse other variants in App. E. For analysis (i.e. Sec. 3), we focus on Vicuna-v1.5-7B [22] as it is shared by both setups. For the ST, we use unimodal encoders, such as ViT [23], TimeSformer [24] and AST [25] that are not aligned with text. More implementation details, and experiments with other backbones can be found in App. D and App. E. We report the similarity after averaging the tokens SimAvg (Eq. (3)) and details other measures in App. E.

Analysis tools. We are interested in cross-modal or multimodal alignment, and define the alignment in terms of the cosine similarity; the higher the score, the more the vector representations are pointing in similar directions. This could also indicates how much the two token distributions or vectors are close, in terms of L2 distance (assuming the vectors are normalized and in a narrow cones). In other words, alignment and similarity terms can be used interchangeably in the paper. In addition to cosine similarity we also study their norm, decoded vocabulary distributions and which LLM weights they activate. In the paper, we focus on the global representation per example, by analysing their average across the sequence. More finegrained analysis on the token level with different similarity and norm

measures gives similar observations and are detailed in App. E. For instance, we compute the cosine similarity between perceptual (P) (e.g., tokens corresponding to image patches) and textual (T) tokens (e.g., tokens corresponding to the image caption), after the block l as follows:

$$\text{Sim}(P^l, T^l) = \frac{\hat{P}^l \cdot \hat{T}^l}{\|\hat{P}^l\| \|\hat{T}^l\|}, \quad \hat{P}^l = \frac{\sum_i^{N_p} p_i^l}{N_p}, \quad \hat{T}^l = \frac{\sum_i^{N_t} t_i^l}{N_t}, \quad (3)$$

3 LLMs indeed generalize to non-textual tokens

We investigate the generalization of LLMs to multimodal inputs, by studying the perceptual and textual tokens inside LLMs. We investigate if all tokens are projected to textual ones, or rather they are still different and how so (results with other models and similarity measures in App. E).

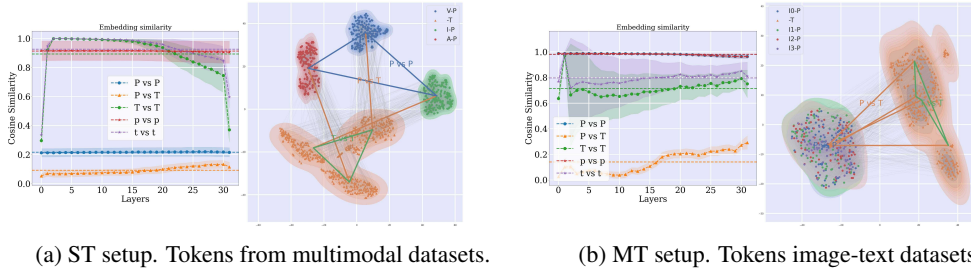


Figure 2: **Multimodal narrow cones.** The cosine similarity after LLM blocks (B) between: perceptual tokens (P vs P), textual tokens (T vs T), perceptual and textual tokens (P vs T). p vs p and t vs t refer to the intra similarity within the same dataset. We also visualize the t-SNE of tokens (at layer 24) showing they stay separated inside the model. V (Video), I (Image), A (Audio).

3.1 How perceptual tokens differ from textual ones?

Multimodal cones: different narrow cones for different modalities (Fig. 2). Previous works [26, 27, 28, 29, 30, 31] have found the representation of contextualized embeddings in language models to be anisotropic: embeddings of different inputs exhibit high cosine similarity, shaping a narrow cone, where all embeddings point in the same narrow direction. In the multimodal domain, the cone effect is also observed [32] in contrastive models (CLIP [33]). In this section, we investigate if textual and multimodal tokens live in narrow cones inside LLMs, and if these cones are distinct. We compute the tokens cosine similarity at different layers. In particular, the unimodal similarity: text-only (T vs T) and perceptual-only (P vs P), and the cross-modal similarity (P vs T) between perceptual and textual tokens. Note that for the ST setup, P vs P covers the similarity between image, video and audio tokens, while for the MT ones cover image tokens from different datasets. Fig. 2 shows a clear narrow cone effect for textual and perceptual tokens. Different perceptual modalities seem to live in different narrow cones, as shown by the low P vs P score for the ST setup. Interestingly, the cross-modal similarity between textual and perceptual tokens (P vs T) is significantly lower, suggesting that textual and perceptual tokens also live in different narrow cones. We also visualize the t-SNE of the tokens embeddings showing they stay separated inside the LLM.

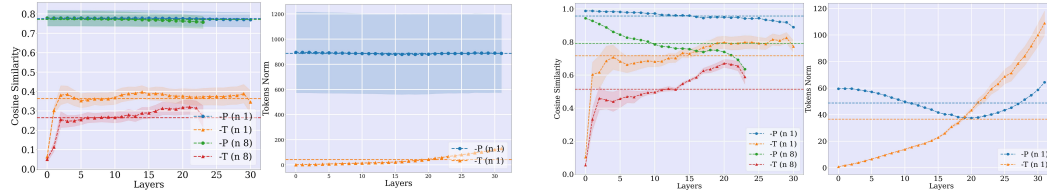


Figure 3: **Tokens norm and evolution across LLM layers.** The tokenwise cosine similarity between consecutive blocks (e.g. X^{l+n} and X^l), and the median token L2 norm after each block (X^l) for the ST (left) and MT (right) setups. Textual and visual tokens evolve differently inside LLMs.

Different token norms and evolution across layers (Fig. 3). We compute the median of the token L2 norms after each LLM block. This shows that textual and perceptual tokens have different norms across layers. Perceptual tokens have significantly higher norm (at the beginning for MT and across all layers for ST), and they change significantly less. When looking at other norm measures, we found perceptual tokens with massive norms, similarly for textual ones [34], especially for the ST setup. We discuss massive tokens more in App. E.2. In addition, we compute the cosine similarity between tokens at block l and block $l + n$, showing that textual and perceptual tokens have different change rates. Textual ones change drastically at the beginning of the LLM, while perceptual ones change significantly less across all layers.

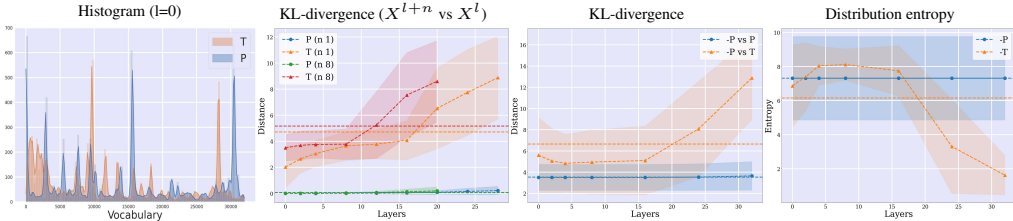


Figure 4: **Tokens vocabulary distribution inside LLMs.** The LLM (Vicuna-v1.5) unembedding layer is used to map each token at different LLM layer, to a probability distribution over the vocabulary. Multimodal tokens exhibit different vocabulary distributions across layers

Different token vocabulary distributions across layers (Fig. 4). For each token, we use the LLM unembedding (*i.e.* LLM head) to decode the latent representation to a probability distribution over the vocabulary. This approach have shown to work well for LLMs at different layers, not just the last one [35, 36, 37, 38]. We show the histogram of this distribution at the first LLM layer for both textual and perceptual prompts. The histograms show clear differences with some overlap between textual and perceptual prompts. In addition, we compute the KL-distance showing that the distributions diverge from each other across LLMs layers. We also notice that the distributions of textual tokens evolve significantly, compared to multimodal ones. This is shown by computing the KL-distance between consecutive blocks and the distribution entropy.

Finding 1. Textual and perceptual tokens live in significantly different representation spaces inside LLMs.

3.2 Do perceptual tokens traverse different paths inside LLMs?

For each trained model, we extract the LLM (frozen) subnetwork activated by each dataset/modality. We study these subnetworks (we refer to as pruning masks) by computing their similarity. We leverage the recent SoTA pruning approach (Wanda [39]), that prune models based on both the weights and the activation norms. Specifically, we use a handful (*e.g.* 256) of calibrated examples coming from different modalities, and keep only $p\%$ (1 - sparsity) of weights with the highest Wanda score, at different sparsity levels (30 % and 50 %). Note that after removing more than 50% of weights we observe a severe degradation of performance.

Similar activated weights across modalities, in the first and deeper layers (Fig. 5). Each subnetwork is represented as a binary mask to indicate which weights are activated. To compute the similarity between these networks, we consider the intersection over union (IoU). Results show an interesting high similarity between subnetworks activated by to different modalities. This high overlap is more seen for the ST setup, for example, the IoU between GQA and VQAv2 is 0.69, similar to GQA vs Audiocaps (0.67) or COCO-Text (0.65). When looking at the IoU across layers, we notice an interesting high score at first layers. It seems that first layers encode general features that are common for all modalities. This similarity increases as we go deeper in the LLM, moving to more abstract and less modality-specific representations, closer to the textual output.

Transfer of multimodal subnetworks across tasks and modalities (Fig. 6). To further validate the previous section, we study if we can simply interchange pruning masks between different tasks.

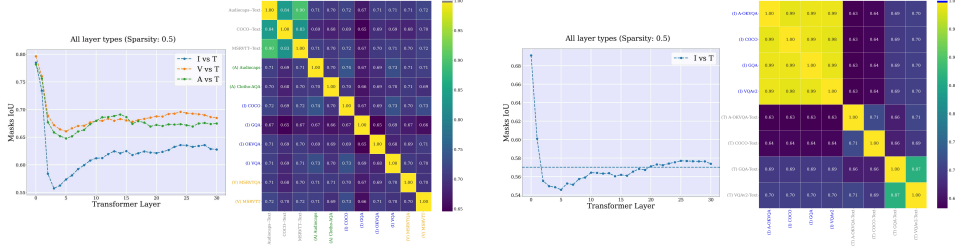


Figure 5: **IoUs of multimodal subnetworks.** IoU of the subnetworks activated by different tasks and modalities, for the ST (left) and MT (right) setups. We show the evolution of IoU across LLM layers and across different multimodal tasks. Different modalities activate similar LLM weights (Fig. 22 for clearer version of the figure).

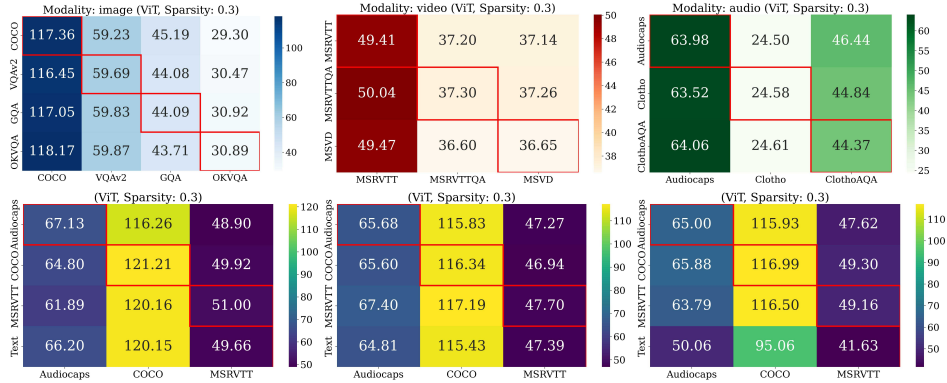


Figure 6: **Transfer of multimodal subnetworks across tasks and modalities.** The subnetwork activated by a given task is used for other tasks for Vicuna-v1.5. From left to right, transfer across: image, video and audio tasks. In each figure, the row corresponds to the subnetwork source dataset and the column to the target dataset. bottom: transfer across modalities for (from left to right): OPT, Llama 2, Vicuna-v1.5.

Specifically, for a given model trained to solve a particular task, we find the pruning mask using calibration data corresponding to different tasks/datasets. The sparsity is set to 30%, which is often used to maintain reasonable performance. Fig. 6 shows that the pruning masks transfer very well across tasks within the same modality (e.g. slight degradation by ~ 1 point CIDEr for captioning with a mask coming from OKVQA). Similarly, we interchange masks across modalities. We fix the task (captioning) and also consider the text modalities (captions without images). In general, we observe similar transfer with a slight performance degradation, especially for OPT and Llama 2. We show similar observations with higher sparsity and with other encoders (e.g., CLIP and MAE) in App. E.4.

Modality-specific subnetworks? The experiments suggest a high overlap between weights activated by different modalities. However, this does not exclude the possibility of finding weights that are generally activated when seeing a particular modality, even if there are small amount of them. More discussion about this can be found in App. E.4.

Finding 2. LLM weights activated by perceptual and textual tokens overlap significantly.

4 What helps LLMs to generalize to multimodal tokens?

Textual and perceptual tokens have very different representations inside LLMs, yet, LLMs are still able to process and generalize to these non-textual tokens. In this section, we try to investigate why this is possible, in particular, we identify which factors facilitate this generalization.

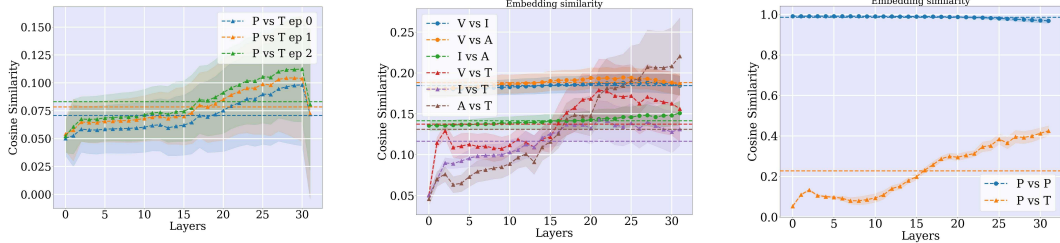


Figure 7: **Multimodal tokens similarity across LLM layers.** The cosine similarity between the textual and multimodal tokens across: training epochs i.e., 0, 1, 2 for Vicuna-v1.5 (first), and across LLMs layers: Vicuna-v1.5 (second) and LLaVA-1.5-4 (last). Textual and multimodal tokens are implicitly aligned during training, and during inference across LLM blocks.

4.1 Observation: the Implicit Multimodal Alignment Effect (IMA)

Implicit alignment during training of the mapping module (Fig. 7). We compute the cosine similarity between the perceptual tokens at the output of the mapping module, and textual tokens at different LLM blocks. Results show that this similarity increases at all layers. This reveals that the mapping module role, is not just to adapt the dimension of the visual tokens, but also to project the visual tokens to be semantically, as similar as possible to the textual ones.

Implicit alignment during inference, across LLM blocks (Fig. 7). We compute the cosine similarity between perceptual and textual tokens after each LLM block. Here we compute the max of tokenwise similarity (MaxSim App. E.1): for each pair of token sequences coming from one example (e.g. image prompt + caption), we take the maximum similarity, then we average across all examples. The tokenwise similarity between perceptual and textual tokens significantly increases, especially in the middle blocks, where the alignment is the highest in deep layers.

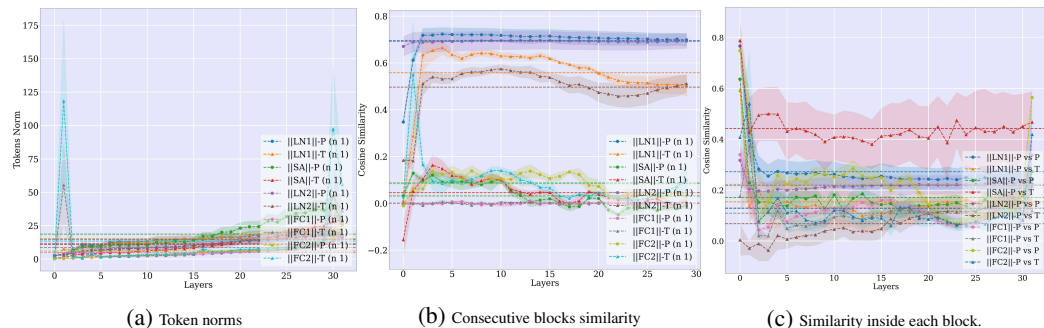


Figure 8: **Multimodal tokens norms and similarity inside LLM blocks.** Token norms (left), tokens cosine similarity between consecutive blocks (middle) and between perceptual and textual tokens (last). The tokens are inside Vicuna-v1.5 blocks (and outside the residual stream): after the self-attention (SA), and FFNs (FC1/2) and layer norms (LN1/2). Multimodal tokens are Implicit alignment inside LLM blocks.

Implicit alignment during inference, inside LLM blocks (outside the residual stream) (Fig. 8). We look deeper to investigate the source of this alignment, and focus on tokens inside the LLM block, which consists mainly of a self-attention (SA), FFN (FC1/2) and layer norms (LN1/2) layers. Interestingly, the similarity between textual and multimodal tokens is the highest after the SA layers.

Finding 3. An implicit multimodal alignment emerges to pull the textual and perceptual tokens closer inside LLMs, during training and inference.

4.2 Explanation: the architectural inductive bias hypothesis

Residual stream with refinement blocks. We notice different observations between the tokens inside and outside the residual stream. In the residual stream, the perceptual and textual tokens exhibit significant representation differences (Sec. 3.1), while outside the residual stream, they are more aligned. Each block contributes slightly to the residual stream (small token norms inside the blocks Fig. 8a), with significantly different contributions (cosine similarity between consecutive blocks close to zero, e.g. after the FC1/2 Fig. 8b). This allows us to view the model as a series of refinement blocks that try to gradually refine the input signals. As the original signals are significantly different, they stay different in the residual stream throughout the model. We argue that this provides a flexibility to handle too different inputs. Moreover, previous works [40] have shown that transformers contain both elements with high and low complexity biases, which helps to build general-purpose architectures [41] that are able to generalize. These works support further our findings.

Refinement blocks as steering blocks. Inside the transformer block, we notice that the layer normalization plays an important role in having comparable norms for both textual and perceptual tokens (Fig. 8a). Perceptual token norms become smaller and closer to textual ones as we traverse several layers in the block. In terms of cross-modal similarity, we notice the highest similarity after the SA, then after the FC2 and LN1. Note that this similarity is higher inside the block, than in the residual stream (e.g. 0.45 vs 0.1 for Vicuna-v1.5 and 0.58 vs 0.15 for LLaVA-1.5-4 in the residual stream Fig. 2). After each block the cross-modal alignment increases, and hence the narrow cones are steered to each other. This suggests that all layers play an important role in steering the textual and perceptual narrow cones to be aligned, with the most contributions coming from the SA.

Finding 4. An LLM can be seen as a residual stream with refinement blocks acting as steering blocks. This architecture design plays an important role in generalizing to very different tokens, and hence other modalities.

5 Implications: performance, safety and efficiency

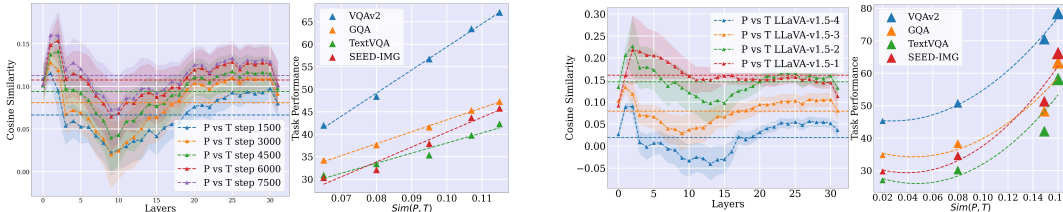


Figure 9: **Implicit alignment as a proxy metric for task performance.** Left: different checkpoints of LLaVA-1.5-4. Right: different variants of the LLaVA-1.5 model. We show the cross-modal token cosine similarity across layers, and the task performance across different benchmarks.

Implicit alignment as a proxy metric for task performance? (Fig. 9) We compute the cosine similarity between perceptual tokens at the LLM input and the textual tokens across LLM layers. The similarity increases during training. Interestingly, we notice a clear and positive correlation with the task performance on several multimodal benchmarks. In addition, we find that this correlation exists across different models, as shown for different LLaVA-1.5 variants.

Implicit alignment as a proxy metric for hallucination? (Fig. 10) Previous works have shown that LMMs suffer from severe hallucinations [42, 43, 44], and generally try to tackle this problem by training on better datasets [45], using RLHF or RLAIIF [46, 47] or post-training heuristics [48, 49]. Here we highlight one of the main causes of hallucinations: which is the lack of internal alignment between textual and perceptual representations. We show the cosine similarity between textual and perceptual tokens after each LLM block, and report the hallucinations on POPE [50] and COCO [51] benchmarks. The curves show clear correlation between the implicit alignment and the hallucinations.

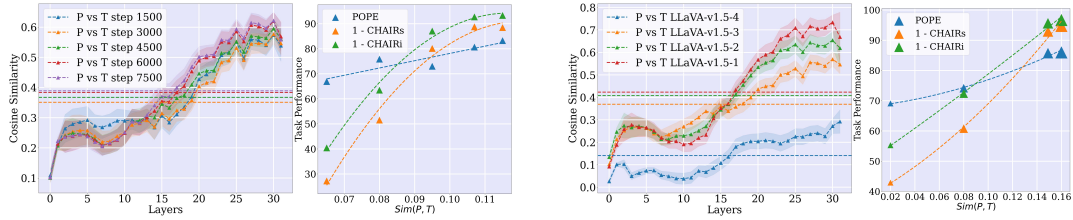


Figure 10: **Implicit alignment as a proxy metric for hallucinations.** Left: different checkpoints of LLaVA-1.5-4. Right: different variants of the LLaVA-1.5 model. We show the cross-modal token cosine similarity across layers, and the hallucinations across different benchmarks.

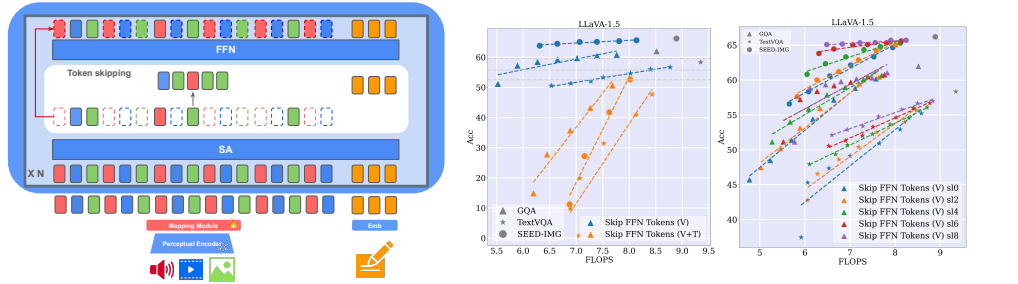


Figure 11: **Skipping computations for visual tokens.** Skipping (Skip ratio)% of the tokens in the FFN layers. sl: skipping start layer. (V): visual tokens. (T): textual tokens. Results on the MT (with LLaVA-1.5) setup.

Skipping computations for visual tokens (Fig. 11). In Sec. 3.1 we show that perceptual tokens change significantly less across layers, compared to textual ones. Sec. 4 highlights the importance of SA layers for cross-modal alignment. In this section, we leverage these observations to reduce the LLM computation overhead by skipping the computations of visual tokens. Specifically, starting from a given start layer (sl), we reduce computations in FFN layers, which accounts for almost 2/3 of model weights, by skipping p% (Skip ratio) of visual tokens. Fig. 11 shows that skipping the visual tokens leads only to slight decrease in performance, while reducing significantly the amount of compute. We provide additional results with the ST setup, and ablation study in App. F.3.

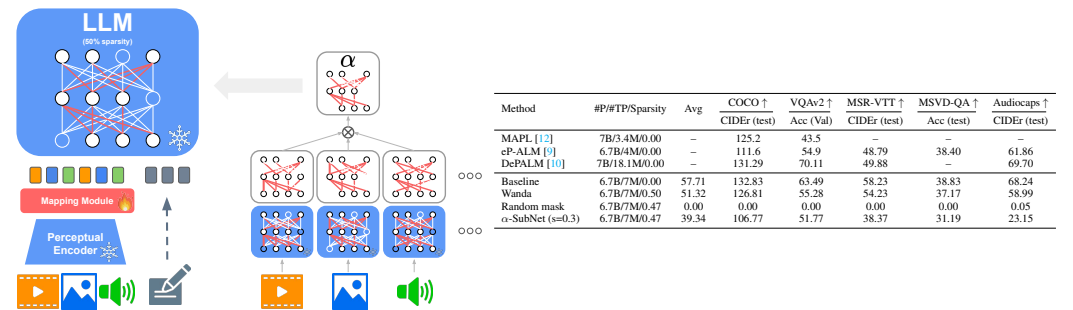


Figure 12: **α -SubNet: a modality-agnostic subnetwork.** Left: illustration of how we obtain the α -SubNet. Right: different methods to compress multimodal LLMs (OPT). Table 2.

α -SubNet: one LLM Subnetwork for all multimodal tasks (Fig. 12) . Despite their differences, multimodal tokens share an important property: slowly changing embeddings across layers (Sec. 3.1). This suggests the possibility of compressing the model while retaining reasonable performance. In addition, textual and multimodal tokens are pulled closer inside the LLM (Sec. 3.2), and processed by almost the same LLM weights (Sec. 4), especially for the ST setup. This suggests the possibility of finding a common subnetwork (α -SubNet) that works well for all multimodal tasks. Thus, we focus on the ST setup with the OPT and CLIP encoders that are currently used by previous works. We consider two representative tasks: COCO image captioning and VQAv2 and provide similar results for other tasks, and modalities in App. F.4. First we use Wanda for task-specific pruning

and show in Fig. 12 that we can obtain scores close to the original ones while removing 50% of the weights. To find the task and modality agnostic α -SubNet, we first extract many pruning masks (e.g. at 30% sparsity) for different modalities, then take the intersection of all these masks (e.g. leading to a global mask at $\sim 50\%$ sparsity). This approach is significantly better than other baselines such as magnitude pruning or a random mask, and leads to comparable performance compared to the task-specific Wanda pruning, especially for VQAv2.

6 Discussion

Limitations. The paper focuses on open-source and frozen LLMs up to 7B parameters, LMMs that concatenate perceptual tokens at the LLM input and are relatively efficient. The generalization of our findings, to larger and more powerful models, with different architectures, including proprietary ones remains to be seen. Detailed discussion in App. B.

Conclusion. We propose the first study of the internal representation of frozen LLMs when exposed to multimodal inputs. We find very different representations for perceptual and textual tokens, yet LLMs are still able to generalize to these non-textual tokens. The implicit multimodal alignment (IMA) effect, linked mostly to architectural design, facilitates this generalization by bringing multimodal tokens closer inside the LLM. Our findings have several implications, such as reducing the computation resources at inference time, understanding better the performance as well as safety-related problems such as hallucinations. We hope that this study will have positive impact, pushing for more works to understand multimodal LLMs, and pave the way to devise more powerful models that are better aligned to human preferences, while targeting safety-related issues.

7 Acknowledgments

The authors would like to thank Arnaud Dapogny and Edouard Yvinec for fruitful discussions, and Damien Teney and Alexandre Ramé for their helpful feedback on the paper. This work was partly supported by ANR grant VISA DEEP (ANR-20-CHIA-0022), and HPC resources of IDRIS under the allocation 2024-[AD011013415R2] made by GENCI.

References

- [1] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022. 1, 24
- [2] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022. 1, 24, 26
- [3] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 1, 24
- [4] OpenAI. Gpt-4 technical report. *arXiv*, 2023. 1, 24, 25
- [5] 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. *arXiv preprint arXiv:2307.09288*, 2023. 1, 24, 26
- [6] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:23716–23736, 2022. 1, 24, 25, 26
- [7] Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu, Carlos Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, et al. Pali-x: On scaling up a multilingual vision and language model. *arXiv preprint arXiv:2305.18565*, 2023. 1, 24, 25

- [8] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 1, 24
- [9] Mustafa Shukor, Corentin Dancette, and Matthieu Cord. ep-alm: Efficient perceptual augmentation of language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22056–22069, October 2023. 1, 2, 3, 9, 24, 37
- [10] Théophane Vallaëys, Mustafa Shukor, Matthieu Cord, and Jakob Verbeek. Improved baselines for data-efficient perceptual augmentation of llms. *arXiv preprint arXiv:2403.13499*, 2024. 1, 2, 3, 9, 24, 26, 37
- [11] Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. Linearly mapping from image to text space. *arXiv preprint arXiv:2209.15162*, 2022. 1, 2, 24
- [12] Oscar Mañas, Pau Rodriguez, Saba Ahmadi, Aida Nematzadeh, Yash Goyal, and Aishwarya Agrawal. Mapl: Parameter-efficient adaptation of unimodal pre-trained models for vision-language few-shot prompting. *arXiv preprint arXiv:2210.07179*, 2022. 1, 2, 3, 9, 24, 26, 37
- [13] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 1
- [14] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*, 2023. 1, 2, 24
- [15] Meghal Dani, Isabel Rio-Torto, Stephan Alaniz, and Zeynep Akata. Devil: Decoding vision features into language. *arXiv preprint arXiv:2309.01617*, 2023. 1
- [16] Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021. 1
- [17] Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Zero-shot video question answering via frozen bidirectional language models. In *NeurIPS 2022-36th Conference on Neural Information Processing Systems*, 2022. 1, 24
- [18] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*, 2022. 2
- [19] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 2, 3, 24, 26
- [20] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. 2
- [21] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023. 2
- [22] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [23] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. 3, 26

- [24] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *International Conference on Machine Learning (ICML)*, volume 2, page 4, 2021. 3, 26
- [25] Yuan Gong, Yu-An Chung, and James Glass. AST: Audio Spectrogram Transformer. In *Proc. Interspeech 2021*, pages 571–575, 2021. 3, 26
- [26] Kawin Ethayarajh. How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 55–65, Hong Kong, China, November 2019. Association for Computational Linguistics. 4, 24
- [27] Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tiejun Liu. Representation degeneration problem in training natural language generation models. In *International Conference on Learning Representations*, 2018. 4, 24
- [28] Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. On the sentence embeddings from pre-trained language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9119–9130, 2020. 4, 24
- [29] Sara Rajaei and Mohammad Taher Pilehvar. How does fine-tuning affect the geometry of embedding space: A case study on isotropy. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3042–3049, 2021. 4, 24
- [30] Sara Rajaei and Mohammad Taher Pilehvar. A cluster-based approach for improving isotropy in contextual embedding space. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 575–584, 2021. 4, 24
- [31] William Rudman and Carsten Eickhoff. Stable anisotropic regularization. In *The Twelfth International Conference on Learning Representations*, 2024. 4
- [32] Victor Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Y Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. *Advances in Neural Information Processing Systems*, 35:17612–17625, 2022. 4, 24
- [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, pages 8748–8763. PMLR, 2021. 4, 26
- [34] Mingjie Sun, Xinlei Chen, J Zico Kolter, and Zhuang Liu. Massive activations in large language models. *arXiv preprint arXiv:2402.17762*, 2024. 5, 24, 30
- [35] nostalgebraist. interpreting gpt: the logit lens. *Less Wrong*, 2020. 5, 31
- [36] J Alammari. Ecco: An open source library for the explainability of transformer language models. In Heng Ji, Jong C. Park, and Rui Xia, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 249–257, Online, August 2021. Association for Computational Linguistics. 5, 31
- [37] Mor Geva, Avi Caciularu, Guy Dar, Paul Roit, Shoval Sadde, Micah Shlain, Bar Tamir, and Yoav Goldberg. Lm-debugger: An interactive tool for inspection and intervention in transformer-based language models. *arXiv preprint arXiv:2204.12130*, 2022. 5, 31
- [38] 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. *arXiv preprint arXiv:2303.08112*, 2023. 5, 31
- [39] Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. A simple and effective pruning approach for large language models. *arXiv preprint arXiv:2306.11695*, 2023. 5, 31

- [40] Damien Teney, Armand Nicolicioiu, Valentin Hartmann, and Ehsan Abbasnejad. Neural redshift: Random networks are not random functions. *arXiv preprint arXiv:2403.02241*, 2024. 8, 24
- [41] Micah Goldblum, Marc Finzi, Keefer Rowan, and Andrew Gordon Wilson. The no free lunch theorem, kolmogorov complexity, and the role of inductive biases in machine learning. *arXiv preprint arXiv:2304.05366*, 2023. 8
- [42] Mustafa Shukor, Alexandre Rame, Corentin Dancette, and Matthieu Cord. Beyond task performance: evaluating and reducing the flaws of large multimodal models with in-context-learning. In *The Twelfth International Conference on Learning Representations*, 2024. 8, 25
- [43] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 292–305, 2023. 8
- [44] Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18135–18143, 2024. 8, 25
- [45] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023. 8
- [46] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023. 8, 25
- [47] Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. Aligning modalities in vision large language models via preference fine-tuning. *arXiv preprint arXiv:2402.11411*, 2024. 8, 25
- [48] Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. Analyzing and mitigating object hallucination in large vision-language models. *arXiv preprint arXiv:2310.00754*, 2023. 8, 25
- [49] Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. Woodpecker: Hallucination correction for multimodal large language models. *arXiv preprint arXiv:2310.16045*, 2023. 8
- [50] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023. 8, 25, 26, 36
- [51] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4035–4045, 2018. 8, 25, 36
- [52] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:1877–1901, 2020. 24
- [53] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022. 24
- [54] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. OpenFlamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023. 24, 26

- [55] Hugo Laurençon, Lucile Saulnier, Leo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. OBELICS: An open web-scale filtered dataset of interleaved image-text documents. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. 24, 25, 26
- [56] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. PaLI: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022. 24
- [57] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR, 2022. 24
- [58] Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi. Unified-io: A unified model for vision, language, and multi-modal tasks. In *The Eleventh International Conference on Learning Representations*, 2022. 24
- [59] Shizhe Diao, Wangchunshu Zhou, Xinsong Zhang, and Jiawei Wang. Write and paint: Generative vision-language models are unified modal learners. In *The Eleventh International Conference on Learning Representations*, 2023. 24
- [60] Mustafa Shukor, Corentin Dancette, Alexandre Rame, and Matthieu Cord. UnIVAL: Unified model for image, video, audio and language tasks. *Transactions on Machine Learning Research*, 2023. 24, 25
- [61] David Mizrahi, Roman Bachmann, Oguzhan Fatih Kar, Teresa Yeo, Mingfei Gao, Afshin Dehghan, and Amir Zamir. 4m: Massively multimodal masked modeling. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 24
- [62] 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. *arXiv preprint arXiv:2312.17172*, 2023. 24
- [63] Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. Grounding language models to images for multimodal generation. *arXiv preprint arXiv:2301.13823*, 2023. 24
- [64] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in Neural Information Processing Systems (NeurIPS)*, 34, 2021. 24
- [65] Mustafa Shukor, Guillaume Couairon, and Matthieu Cord. Efficient vision-language pre-training with visual concepts and hierarchical alignment. In *33rd British Machine Vision Conference (BMVC)*, 2022. 24
- [66] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Zicheng Liu, Michael Zeng, et al. An empirical study of training end-to-end vision-and-language transformers. *arXiv preprint arXiv:2111.02387*, 2021. 24
- [67] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. *arXiv preprint arXiv:2201.12086*, 2022. 24
- [68] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15638–15650, 2022. 24
- [69] Zhanyu Wang, Longyue Wang, Zhen Zhao, Minghao Wu, Chenyang Lyu, Huayang Li, Deng Cai, Luping Zhou, Shuming Shi, and Zhaopeng Tu. Gpt4video: A unified multimodal large language model for instruction-followed understanding and safety-aware generation. *arXiv preprint arXiv:2311.16511*, 2023. 24

- [70] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal llm. *arXiv preprint arXiv:2309.05519*, 2023. 24
- [71] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Jiao Qiao. LLaMA-Adapter: Efficient fine-tuning of language models with zero-init attention. 2303.16199, 2023. 24
- [72] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 24
- [73] Artemis Panagopoulou, Le Xue, Ning Yu, Junnan Li, Dongxu Li, Shafiq Joty, Ran Xu, Silvio Savarese, Caiming Xiong, and Juan Carlos Niebles. X-instructblip: A framework for aligning x-modal instruction-aware representations to llms and emergent cross-modal reasoning. *arXiv preprint arXiv:2311.18799*, 2023. 24
- [74] Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 2023. 24
- [75] Zhengqing Yuan, Zhaoxu Li, and Lichao Sun. Tinygpt-v: Efficient multimodal large language model via small backbones. *arXiv preprint arXiv:2312.16862*, 2023. 24
- [76] Yichen Zhu, Minjie Zhu, Ning Liu, Zhicai Ou, Xiaofeng Mou, and Jian Tang. Llava-phi: Efficient multi-modal assistant with small language model. *arXiv preprint arXiv:2401.02330*, 2024. 24
- [77] Baichuan Zhou, Ying Hu, Xi Weng, Junlong Jia, Jie Luo, Xien Liu, Ji Wu, and Lei Huang. Tinyllava: A framework of small-scale large multimodal models. *arXiv preprint arXiv:2402.14289*, 2024. 24
- [78] Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, En Yu, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Small language model meets with reinforced vision vocabulary. *arXiv preprint arXiv:2401.12503*, 2024. 24
- [79] Jiwon Song, Kyungseok Oh, Taesu Kim, Hyungjun Kim, Yulhwa Kim, and Jae-Joon Kim. Sleb: Streamlining llms through redundancy verification and elimination of transformer blocks. *arXiv preprint arXiv:2402.09025*, 2024. 24
- [80] Andrey Gromov, Kushal Tirumala, Hassan Shapourian, Paolo Glorioso, and Daniel A Roberts. The unreasonable ineffectiveness of the deeper layers. *arXiv preprint arXiv:2403.17887*, 2024. 24
- [81] Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava, Ce Zhang, Yuandong Tian, Christopher Re, et al. Deja vu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*, pages 22137–22176. PMLR, 2023. 24
- [82] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35:30318–30332, 2022. 24
- [83] Ibrahim Alabdulmohsin, Vinh Q Tran, and Mostafa Dehghani. Fractal patterns may unravel the intelligence in next-token prediction. *arXiv preprint arXiv:2402.01825*, 2024. 24
- [84] Sarah Schwettmann, Neil Chowdhury, Samuel Klein, David Bau, and Antonio Torralba. Multimodal neurons in pretrained text-only transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2862–2867, 2023. 24
- [85] Haowen Pan, Yixin Cao, Xiaozhi Wang, and Xun Yang. Finding and editing multi-modal neurons in pre-trained transformer. *arXiv preprint arXiv:2311.07470*, 2023. 24
- [86] Zhiyuan Fang, Jianfeng Wang, Xiaowei Hu, Lijuan Wang, Yezhou Yang, and Zicheng Liu. Compressing visual-linguistic model via knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1428–1438, 2021. 24

- [87] Kan Wu, Houwen Peng, Zhenghong Zhou, Bin Xiao, Mengchen Liu, Lu Yuan, Hong Xuan, Michael Valenzuela, Xi Stephen Chen, Xinggang Wang, et al. Tinyclip: Clip distillation via affinity mimicking and weight inheritance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21970–21980, 2023. 24
- [88] Zhe Gan, Yen-Chun Chen, Linjie Li, Tianlong Chen, Yu Cheng, Shuohang Wang, Jingjing Liu, Lijuan Wang, and Zicheng Liu. Playing lottery tickets with vision and language. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 652–660, 2022. 24
- [89] Jia Huei Tan, Chee Seng Chan, and Joon Huang Chuah. End-to-end supermask pruning: Learning to prune image captioning models. *Pattern Recognition*, 122:108366, 2022. 24
- [90] Dachuan Shi, Chaofan Tao, Ying Jin, Zhendong Yang, Chun Yuan, and Jiaqi Wang. UPop: Unified and progressive pruning for compressing vision-language transformers. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 31292–31311. PMLR, 23–29 Jul 2023. 24
- [91] Ali Furkan Biten, Lluís Gomez, and Dimosthenis Karatzas. Let there be a clock on the beach: Reducing object hallucination in image captioning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1381–1390, 2022. 25
- [92] Junyang Wang, Yiyang Zhou, Guohai Xu, Pengcheng Shi, Chenlin Zhao, Haiyang Xu, Qinghao Ye, Ming Yan, Ji Zhang, Jihua Zhu, et al. Evaluation and analysis of hallucination in large vision-language models. *arXiv preprint arXiv:2308.15126*, 2023. 25
- [93] Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. *arXiv preprint arXiv:2311.16922*, 2023. 25
- [94] Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale Fung. Plausible may not be faithful: Probing object hallucination in vision-language pre-training. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2136–2148, May 2023. 25
- [95] Bohan Zhai, Shijia Yang, Xiangchen Zhao, Chenfeng Xu, Sheng Shen, Dongdi Zhao, Kurt Keutzer, Manling Li, Tan Yan, and Xiangjun Fan. Halle-switch: Rethinking and controlling object existence hallucinations in large vision language models for detailed caption. *arXiv preprint arXiv:2310.01779*, 2023. 25
- [96] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023. 25
- [97] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 25
- [98] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 26
- [99] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024. 26
- [100] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning (ICML)*, pages 10347–10357. PMLR, 2021. 26

- [101] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16000–16009, 2022. 26
- [102] Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang, and Rongrong Ji. X-clip: End-to-end multi-grained contrastive learning for video-text retrieval. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 638–647, 2022. 26
- [103] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *Advances in neural information processing systems*, 35:10078–10093, 2022. 26
- [104] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. In *NeurIPS*, 2022. 26
- [105] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017. 26
- [106] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6700–6709, 2019. 26
- [107] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 26
- [108] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision (ECCV)*, pages 740–755. Springer, 2014. 26
- [109] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM International Conference on Multimedia, MM '17*, page 1645–1653, New York, NY, USA, 2017. Association for Computing Machinery. 26
- [110] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 26
- [111] Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In *NAACL-HLT*, 2019. 26
- [112] Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. Clotho: An audio captioning dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 736–740. IEEE, 2020. 26
- [113] Samuel Lipping, Parthasaarathy Sudarsanam, Konstantinos Drossos, and Tuomas Virtanen. Clotho-aqa: A crowdsourced dataset for audio question answering. In *2022 30th European Signal Processing Conference (EUSIPCO)*, pages 1140–1144. IEEE, 2022. 26
- [114] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023. 26
- [115] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326, 2019. 26

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: For findings Sec. 3 and implications Sec. 5.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Sec. 6 and more detailed in App. B

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Details to reproduce the experiments are in Sec. 2 and more detailed in App. D. Also the code will be made public.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The paper uses solely public datasets.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The experimental details can be found in Sec. 2 and App. D

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The paper includes different statistical measures to report scores such as the average, median and the variance.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The paper mostly about analysis, but the considered models are trained by authors with details included in App. D. In addition we report the FLOPs during inference in the implications section Sec. 5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: No violation for the code of ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: App. C

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The paper cites all used assets.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

Supplementary material

This supplementary material is organized as follows:

- App. **A**: detailed related work.
- App. **B**: detailed discussion about the work, limitations and other implications.
- App. **C**: the broader impact of the work.
- App. **D**: implementation details, including the trained models, datasets and metrics.
- App. **E**: additional experiments analysing LLMs.
- App. **F**: additional experiments for the implications.

A Detailed related work

Large multimodal models. Motivated by the success of large-scale training of LLMs [52, 1, 3, 53, 2, 5, 4], the multimodal community has embarked on a parallel journey, striving to develop larger and more powerful models capable of processing multiple modalities. Typical Large Multimodal Models (LMMs) are constructed either by building upon frozen LLMs [6, 54, 55] or by training them end-to-end after initialization [56, 7, 8]. These models have demonstrated success in numerous general and intricate multimodal tasks, achieving performance levels close to human capability. Another important line of research focuses on unified models, where a single model is designed to handle diverse tokenized modalities, such as image-text [57, 58, 59], or even beyond two modalities [60, 61, 62].

Efficient large multimodal models. Recently, to mitigate the training cost associated with training large multimodal models, efficient adaptation of unimodal models has emerged as a promising direction. Models like [11, 9, 17, 12, 10, 63] maintain LLMs frozen and train only a small subset of adaptation parameters for different multimodal tasks. These approaches achieve competitive performance compared to end-to-end trained models [64, 65, 66, 67, 68] on image-text tasks and also on audio and video-text tasks [9, 69, 70, 10]. Beyond single-task tuning, many approaches do relatively light-weight pretraining and/or instruction tuning [71, 72, 14, 73, 19] and achieve good zero-shot generalization and instruction following abilities. To make these models more efficient, previous works have trained models with smaller LLMs showing competitive performance [74, 75, 76, 77, 78]

Analyzing LLMs. Previous research has highlighted the highly anisotropic nature of embeddings within language models, characterized by high cosine similarity [26, 27, 28, 29, 30]. Studies focusing on efficiency have shown that textual tokens exhibit small changes across layers [79, 80, 81]. Additionally, LLMs contain outlier features [82] and massive activations [34], which significantly influence model performance. Work by [83] suggests that LLMs may generalize due to the fractal structure of language. Moreover, [40] demonstrate that the building blocks of LLMs implicitly bias towards approximating both complex and simple functions. In the multimodal domain, [32] identify a modality gap in CLIP models attributable to the narrow cone effect. Previous studies have also explored neurons in LLMs that encode multimodal representations [84, 85].

Compression and pruning for multimodal models. Few works have targeted multimodal model compression, focusing mostly on image-text models. Notably, some works have concentrated on distilling knowledge from larger models through attention maps [86] or the affinity matrix in CLIP models [87]. Recent efforts have successfully applied the Lottery Ticket Hypothesis (LTH) to these models, optimizing both the model weights and masks jointly [88, 89]. A unified framework for structured pruning based on iterative training and adaptive sparsity allocation was proposed by [90]. It's worth noting that these approaches, initially designed for relatively small image-text models, encounter scalability challenges when applied to very large models.

Hallucinations in multimodal models. Hallucinations in multimodal models involve generating text that refers to objects not present in the input image [51, 50]. This pervasive issue affects a wide range of multimodal models, varying in architecture, training data, and scale [91, 60, 42, 6]. To gain deeper insights into this phenomenon, numerous studies have proposed evaluation benchmarks to quantify hallucinations across different dimensions [51, 50, 92, 44], shedding light on underlying causes. Specifically, co-occurrences and uncertainty [48], as well as visual uncertainty stemming from lower image resolution [93, 94, 95], have been identified as contributing factors. Additionally, it has been demonstrated that multimodal in-context learning exacerbates hallucinations [42]. To address this issue, various techniques have been proposed, including training on improved datasets [96, 44], aligning models with reinforcement learning [46, 47], refining training objectives [94], and employing post-training heuristics [48, 93]. Our study highlights the misalignment between internal representations of textual and perceptual tokens as a key cause of hallucinations.

B Discussion

Study across LLMs and setups. Our investigation primarily centers on Vicuna-v1.5 across single-task and multitask setups. We find that our conclusions remain consistent across various LLMs (e.g., OPT, Llama 2) and different settings (e.g., with and without pretraining, using different mapping modules), as detailed in the appendix. Extending our analysis to encompass other multimodal models, potentially with diverse architectures [6, 55], could offer additional valuable insights.

Study on larger models. While our work primarily focuses on frozen LLMs to provide insights relevant to future multimodal models, we also present results involving trained LLMs such as in LLaVA-1.5. These experiments yield observations akin to those with frozen LLM variants. However, the applicability of our experiments or the generalizability of our findings to larger LLMs (beyond 7B parameters), larger multimodal models [7, 6], or massively-trained multimodal foundation models like Gemini [97] or GPT4-V [4] remains an open question.

Remaining questions to understand LLMs. While we primarily investigate why and how LLMs generalize to multimodal inputs, and offer insights into issues such as object hallucinations, numerous unanswered questions persist. For instance, further exploration is needed to discern the encoded information in tokens and how LLMs extract information from visual tokens. Deeper inquiries are also required to address safety-related issues in large models, including the inability to abstain from answering, compositionality, and the precise adherence to user instructions [42].

Other implications. Our paper discusses several practical implications with potential benefits. Future extensions of our study could focus on specific aspects, such as enhancing model efficiency during training and inference by reducing redundant computations or model size. Additionally, addressing alignment with human preferences, such as faithfulness and safety, remains a significant challenge requiring further investigation. Our study may also inform model architecture design, such as developing mapping modules explicitly aligning multimodal tokens before entering the LLM.

C Broader impacts

The paper aims to enhance our comprehension of LLMs within the realm of multimodal inputs. We contend that a deeper understanding of these models can yield positive societal impacts, which we partially address in this study. For instance, our findings may contribute to mitigating the consumption of large models and their potential societal harms. Moreover, our work may inspire future research endeavors with various impacts, none of which we think must be specifically discussed here.

D Implementation details.

D.1 Perceptually augmented LLM baselines

D.1.1 ST setup

We train many models across different datasets that span image, video and audio-text modalities. We first devise powerful baselines based on 3 tenets: (a) having the smallest number of trainable

parameters, (b) general architecture that span or similar to many existing models and (3) good performance. To this end, and inspired by previous studies showing the effectiveness of using transformer-based mapping module [10, 12, 6, 54, 55, 98], we use light-weight transformer with learnable queries and self-attention to attend to perceptual tokens. This transformer operates in low dimension space (*i.e.*, due to down/up projection layers and the the number of learnable query are limited to 10. We also favor a deeper architecture (5 blocks) compared to a wider one [10]. Our baselines are close to [10], but with significantly less trainable parameters. We train these baselines with different LLMs: OPT-6.7B [2], Llama 2-7B [5] and Vicuna-v1.5-7B [99] and different encoders for: image (ViT [23, 100], CLIP[33], MAE[101]), video (TimesFormer[24], X-CLIP[102], VideoMAE[103]) and audio (AST[25], AudioMAE[104]).

To train these baselines, we use AdamW optimizer with a learning rate of $2e-4$ that decreases with a cosine annealing scheduler to a minimum of $1e-5$. We train with a total batch size of 16 for captioning and 64 for VQA datasets. The number of epochs is set to 20 to ensure that all models converged, though most of these models converge after only couple of epochs. We select the best checkpoint for evaluation. For example the model for image captioning converged after ~ 4 epochs. All models are trained on 8 V100 GPUs and the training time depends on the task, *e.g.*, for the large VQAv2 dataset each epoch takes ~ 30 mins, for other smaller datasets it takes less time, *e.g.*, ~ 10 mins for Audiocaps and MSVD-QA. Unless specified otherwise, we fix the hyperparameters for all baselines to isolate the variations that could results from this.

We refer to all text-aligned models as CLIP, trained for classification as ViT, and self-supervised with MAE objective as MAE.

D.1.2 MT setup

To study the impact of different factors (*e.g.* pretraining, mapping module) on the internal representations (*e.g.* implicit alignment), we devise different variants of the LLaVA-1.5 [19] model that differ from the original model as follows: LLaVA-1.5-2 (LLM kept frozen), LLaVA-1.5-3 (LLM kept frozen, without pretraining) and LLaVA-1.5-4 (LLM kept frozen, without pretraining and with transformer mapping module similar. The latter is very similar to the models used in the ST setup, which ensure comparable observatoins. All these models are based on the Vicuna-v1.5-7B LLM.

We follow the same training setup of LLaVA-1.5 [19], including the training data, steps and hyperparameters.

For analysis in the paper (*e.g.* Sec. 3), we focus on Vicuna-v1.5 as it is shared by both setups. For the ST setup, we use unimodal encoders, such as ViT, TimeSformer and AST that are not aligned with text.

D.2 Datasets and metrics

ST setup. We consider a wide range of public multimodal datasets that cover 2 representative tasks: captioning and question-answering (QA) across image (VQAv2 [105], GQA [106], OKVQA [107], COCO caption [108]), video (MSVD, MSRVTQA [109], MSRVT [110]), audio (Audiocaps [111], Clotho [112], Clotho-AQA [113]) and language tasks. For QA datasets we report the accuracy (in open-ended generation setup with exact match), and for captioning we report the CIDEr metric.

MT Setup. We also evaluate the MT setup on recent datasets such as SEED [114], TextVQA [115] and POPE [50].

E LLMs generalize to multimodal inputs: additional experiments

E.1 Tokens evolution across layers

Different LLaVA-1.5 variants. In Fig. 13, we illustrate the differences between perceptual and textual tokens across different LLaVA-1.5 variants. Comparing the two variants with multimodal pretraining, we observe higher cross-modal similarity in LLaVA-1.5 compared to LLaVA-1.5-2, which freezes the LLM. This suggests that training the LLM enhances alignment between representations. For models without pretraining, we note that using an MLP (LLaVA-1.5-3) to connect both models yields better results than using transformer-based pooling (LLaVA-1.5-4), potentially explaining the

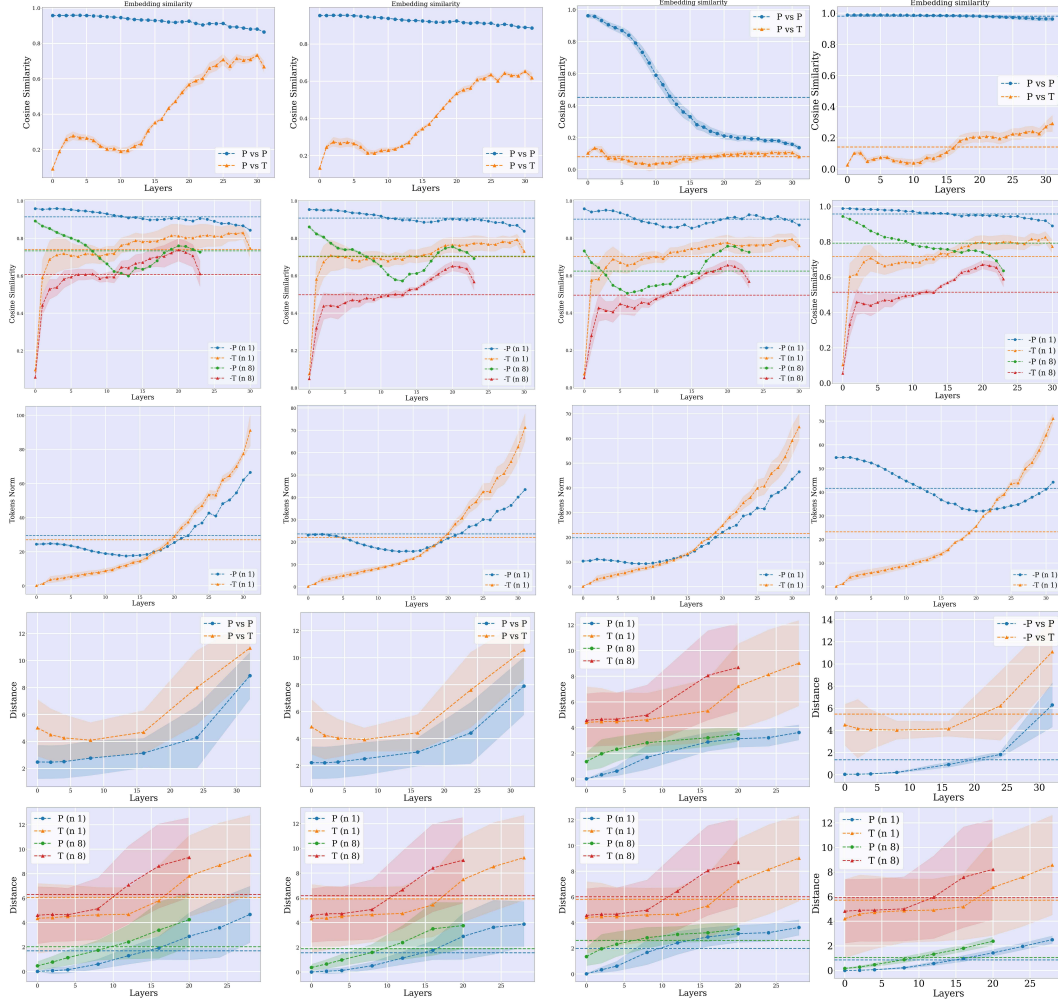


Figure 13: Textual and multimodal tokens for LLaVA-1.5 variants (MT setup). From top to bottom: (1) the cosine similarity between the textual and multimodal tokens across LLM blocks. (2) the cosine similarity between consecutive blocks. (3) token norms, (4) KL-distance between vocabulary distributions decoded from textual and perceptual tokens, (6) cosine similarity between vocabulary distribution at consecutive layers. From left to right: LLaVA-1.5, LLaVA-1.5-2, LLaVA-1.5-3, LLaVA-1.5-4.

superior scores. Utilizing all visual tokens appears to bolster alignment with textual tokens, with pretraining further enhancing this alignment. Notably, vocabulary distributions undergo significant changes in middle layers, particularly for textual tokens. Similar observations hold across different variants, indicating that our findings generalize to broader setups and that training the LLM does not substantially alter token behavior.

Different similarity measures for cross-modal alignment. In this section, we compare the following similarity measures to compute the similarity between perceptual ($P = [p_1, \dots, p_{N_p}]$) and

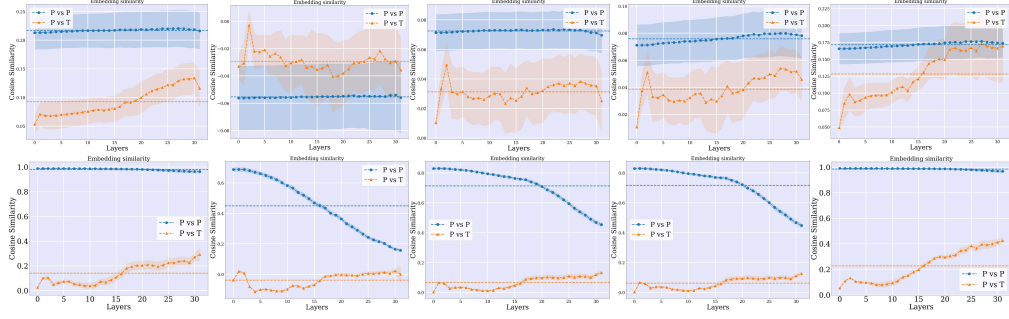


Figure 14: **Different similarity measures.** From left to right: SimAvg, MinSim, AvgSim, MedSim, MaxSim. ST setup (top) and MT setup (bottom).

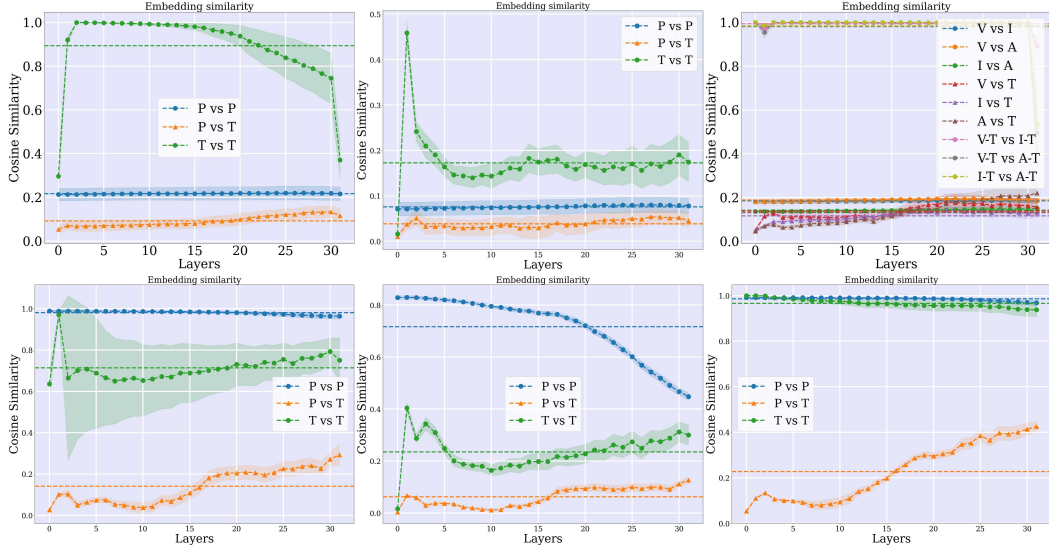


Figure 15: **Different similarity measures and the narrow cone effect.** From left to right: SimAvg, MedSim and MaxSim. Vicuna-v1.5 (top), LLaVA-1.5-4 (bottom).

textual ($T = [t_1, \dots, t_{N_t}]$) tokens:

$$\text{Sim}(X, Y) = \frac{X \cdot Y}{\|X\| \|Y\|}, \quad (4)$$

$$\text{SimAvg}(P, T) = \text{Sim}(\hat{P}, \hat{T}), \quad \hat{P} = \frac{\sum_{i \in [N_p]} p_i}{N_p}, \quad \hat{T} = \frac{\sum_{i \in [N_t]} t_i}{N_t}, \quad (5)$$

$$\text{MaxSim}(P, T) = \max_{i \in [N_p], j \in [N_t]} \text{Sim}(p_i, t_j), \quad (6)$$

$$\text{MinSim}(P, T) = \min_{i \in [N_p], j \in [N_t]} \text{Sim}(p_i, t_j), \quad (7)$$

$$\text{AvgSim}(P, T) = \frac{\sum_{i \in [N_p], j \in [N_t]} \text{Sim}(p_i, t_j)}{N_p + N_t}, \quad (8)$$

$$\text{MedSim}(P, T) = \text{Med}_{i \in [N_p], j \in [N_t]} \text{Sim}(p_i, t_j), \quad (9)$$

Where $[N_p] = \{1, \dots, N_p\}$ and $[N_t] = \{1, \dots, N_t\}$ and Med is the median operation.

Fig. 14 shows the inter (P vs T) and intra (P vs P) similarity. According to all measures, except AvgSim and MinSim, we have similar observations: increasing inter similarity and higher intra similarity that increases in last layers. For MinSim and AvgSim for the ST setup, we do not see such observations, indicating that not all perceptual tokens are, or should be aligned to text.

Fig. 15 and Fig. 16 show the measures comparison for the narrow cone experiments. Interestingly, the narrow cone effect is less seen when looking at the median of the token similarities (MedSim), indicating that this effect is not driven by all tokens, and at the token level the representation is not always anisotropic.

In spite of having similar observations between several measures, we focus on SimAvg, as it is much faster to compute, especially when there is a large number of tokens (as in LLaVA-1.5).

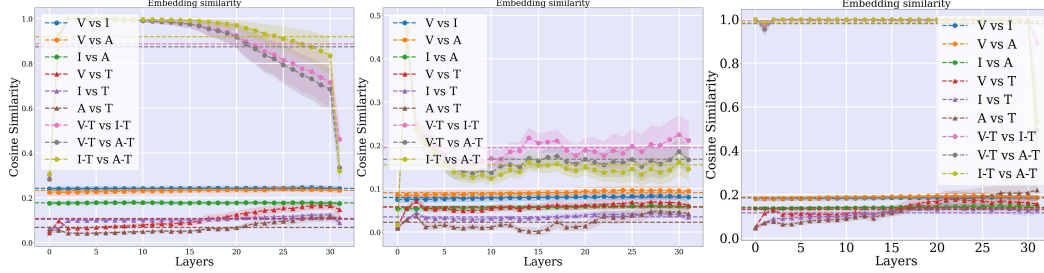


Figure 16: **Narrow cones for image, video, audio and text modalities.** From left to right: SimAvg, MedSim and MaxSim. Vicuna-v1.5 (top), LLaVA-1.5-4 (bottom).

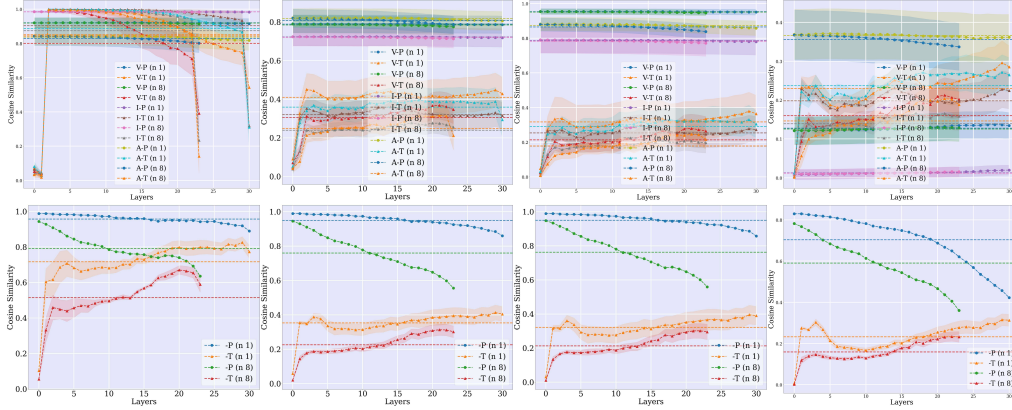


Figure 17: **Different similarity measures between tokens at consecutive layers.** From left to right: SimAvg, AvgDiagSim, MedDiagSim and MedSim. Vicuna-v1.5 (top), LLaVA-1.5-4 (bottom).

Different similarity measures for the similarity across consecutive layers. we compare the following similarity measures to compute the token similarity between consecutive blocks (e.g. between tokens at block $X^a = x_1^a, \dots, x_N^a$ and $b X^b = x_1^b, \dots, x_N^b$):

$$\text{Sim}(X^a, X^b) = \frac{X^a \cdot X^b}{\|X^a\| \|X^b\|}, \quad (10)$$

$$\text{SimAvg}(X^a, X^b) = \text{Sim}(\hat{X}^a, \hat{X}^b), \quad \hat{X} = \frac{\sum_i^N x_i}{N}, \quad (11)$$

$$\text{AvgDiagSim}(X^a, X^b) = \frac{\sum_{i \in [N]} \text{Sim}(p_i^a, p_i^b)}{N}, \quad (12)$$

$$\text{MedDiagSim}(X^a, X^b) = \text{Med}_{i \in [N]} \text{Sim}(p_i^a, p_i^b), \quad \text{MedSim}(X^a, X^b) = \text{Med}_{i, j \in [N]} \text{Sim}(p_i^a, p_j^b), \quad (13)$$

Where $[N] = \{1, \dots, N\}$ and Med is the median operation.

Fig. 17 shows similar observations across all different measures when each token is compared with the token at the same position in different layers. However, when taking the median of similarities (MedSim) across all tokens, this similarity is significantly smaller, especially for the ST setup. This reveals that tokens can be very different within the same modality or example.

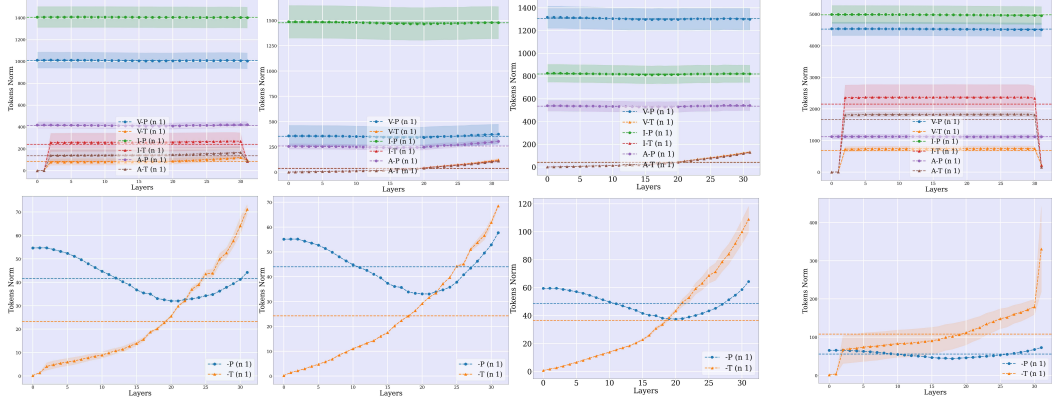


Figure 18: **Different token norm measures.** We compute the token L2 norm at consecutive blocks (e.g. B^{l+n} and B^l) for the ST (top) and MT (bottom) setups. From left to right: NormAvg, MinNorm, MedianNorm and MaxNorm.

Tokens evolution for different modalities. Fig. 19, shows that textual and multimodal tokens evolve differently inside LLMs.

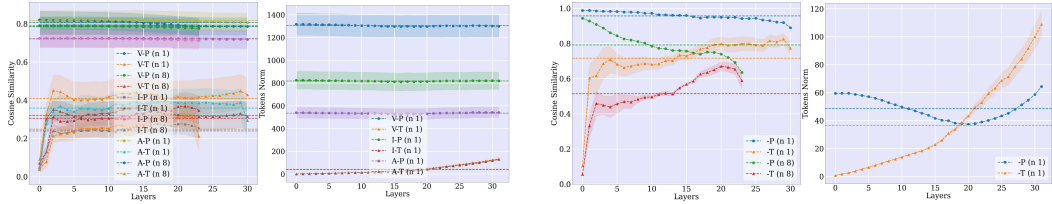


Figure 19: **Textual and multimodal tokens evolve differently inside LLMs.** We compute the tokenwise cosine similarity and the median token L2 norm at consecutive blocks (e.g. X^{l+n} and X^l) for the ST (left) and MT (right) setups.

E.2 Token norms across layers

Massive token norms. In this section we highlight the presence of tokens with massive norms, this becomes clearer when looking at different norm measures. We compare the following measures to compute the token L2 norms across blocks (e.g. $X = x_1, \dots, x_N$):

$$\text{Norm}(X) = \sqrt{\sum_i^M X_i^2}, \quad (14)$$

$$\text{NormAvg}(X) = \text{Norm}(\hat{X}), \quad \hat{X} = \frac{\sum_i^N x_i}{N}, \quad (15)$$

$$\text{MinNorm}(X) = \min_{i \in [N]} \text{Norm}(x_i), \quad (16)$$

$$\text{MedianNorm}(X) = \text{Med}_{i \in [N]} \text{Norm}(x_i), \quad (17)$$

$$\text{MaxNorm}(X) = \max_{i \in [N]} \text{Norm}(x_i), \quad (18)$$

Where $[N] = \{1, \dots, N\}$ and Med is the median operation and M is the total number of elements in the tensor. For the ST setup, Fig. 18 shows a very high token norms when looking at NormAvg and MaxNorm, compared to MinNorm and MedianNorm. These massive norms are present for both textual and perceptual tokens, and they are larger for perceptual ones. When looking closely, we find that these tokens correspond to start or split tokens as seen in [34]. For the MT setup, we notice that these massive tokens presents mainly in the system message, which we remove for our study as it is common for all examples. Interestingly the perceptual tokens for the MT setup do not seem to have massive norms.

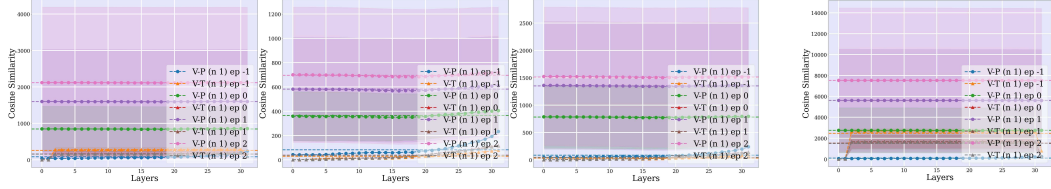


Figure 20: **L2 token norm increases with training.** We compute the token L2 norm during training and across the LLM blocks for the ST setup (Vicuna-v1.5). From left to right: NormAvg, MinNorm, MedianNorm and MaxNorm.

Increasing token norm during training. We try to investigate why we have very high perceptual token norms. To this end, we compute the norm across different epochs. Fig. 20 shows that during training of the mapping module, the norms increase significantly.

E.3 Vocabulary distribution

For each token, we use the LLM unembedding (*i.e.* LLM head) to decode the latent representation to a probability distribution over the vocabulary. This approach have shown to work well for LLMs at different layers, not just the last one [35, 36, 37, 38]. In Fig. 21, we show the histogram of this distribution at the first LLM layer for both textual and perceptual tokens, the KL-distance between the 2 distributions, KL-distance between consecutive layers and the entropy. Here we report additional results for the LLaVA-1.5 baseline showing similar observations to those of ST reported in the main paper.

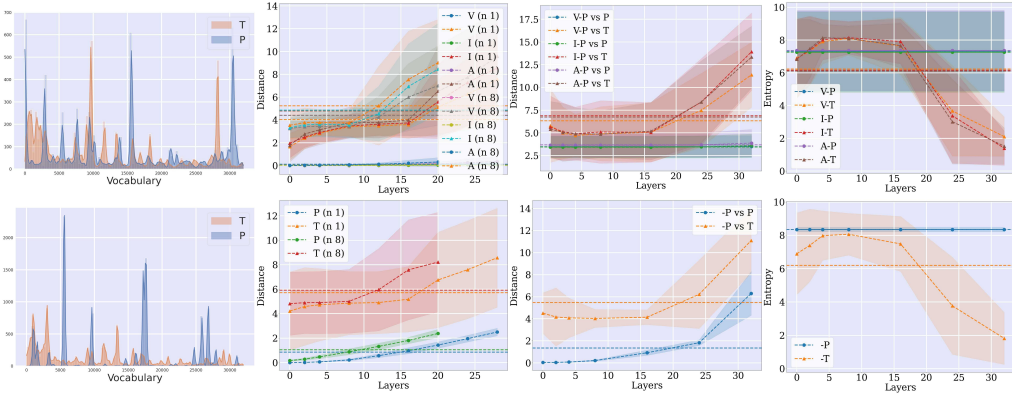


Figure 21: **Textual and visual tokens have different vocabulary distributions inside LLMs.** We use the LLM unembedding layer to map each token to a probability distribution over the vocabulary. We then show (from left to right): the histograms at the input of the LLM, the KL divergence between the distributions at consecutive layers, the KL divergence between textual and perceptual distributions and the distribution entropy. Top: Vicuna-v1.5, Bottom: LLaVA-1.5-4.

E.4 Similar activated weights by different modalities

Experimental setup. In this section, we analyse the subnetworks activated by different multimodal inputs. We use the Wanda score [39] to extracted these subnetworks or pruning masks, then compute the IoU. For multimodal datasets we consider only the perceptual tokens, for example the visual tokens without the questions for VQAv2. We also use the text in these datasets as a source for textual tokens (*e.g.*, COCO-text consider only the captions in the COCO dataset).

Different LLaVA-1.5 variants. In Fig. 23, we show the overlap between the weights activated by different modalities. Across different LLaVA-1.5 variants, we find similar observations: high overlap between perceptual and textual activated weights (*e.g.* 0.6 IoU), which is less than the overlap between weights activated by the same modality (*e.g.* 0.95 for perceptual tokens and 0.87 for textual

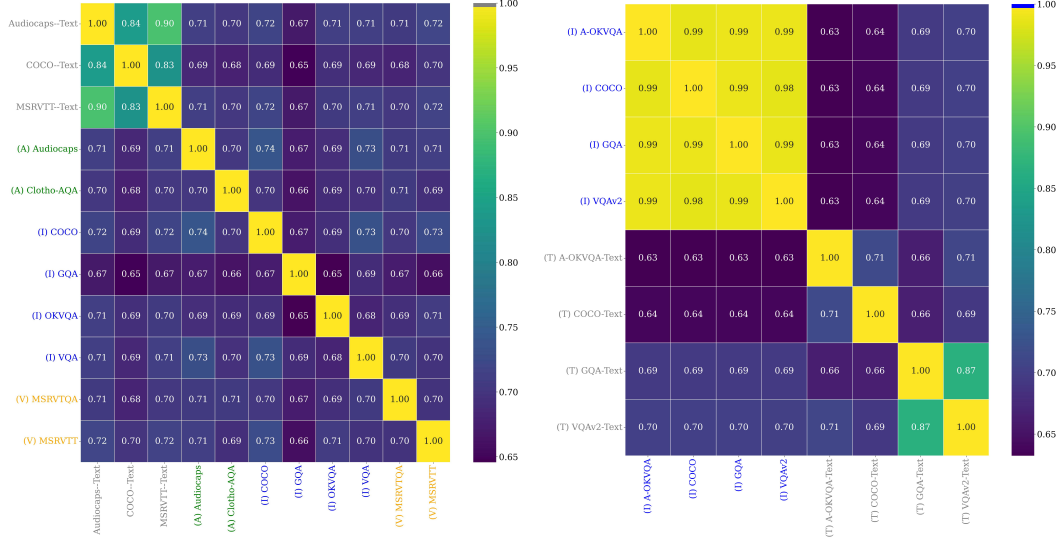


Figure 22: **IoUs of multimodal subnetworks.** IoU of the subnetworks activated by different tasks and modalities, for the ST (left) and MT (right) setups. Different modalities activate similar LLM weights.

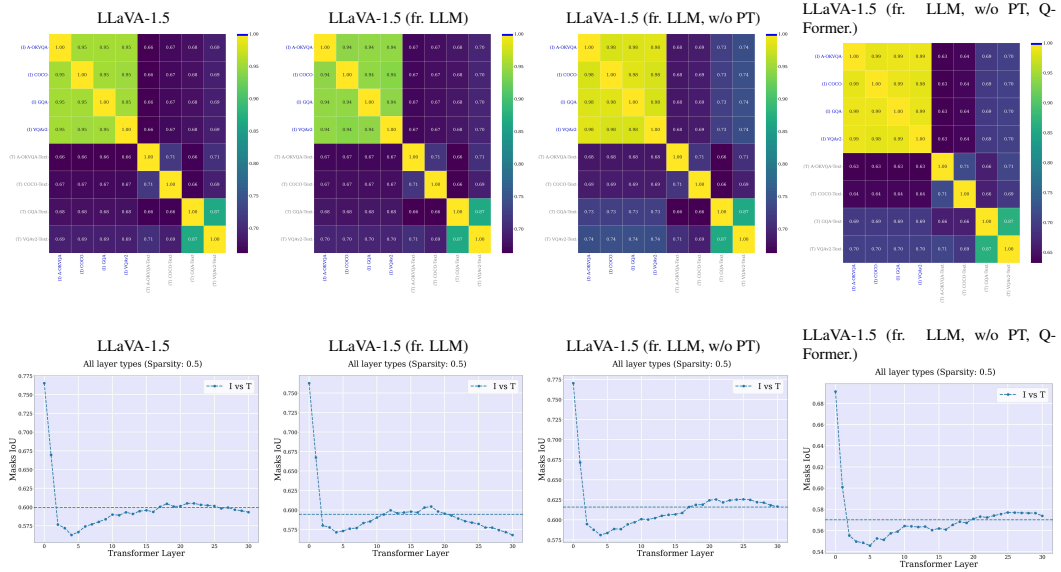


Figure 23: **IoUs of activated subnetworks for LLaVA-1.5 variants.** We compute the IoU of weights activated by different multimodal tokens. From left to right: LLaVA-1.5, LLaVA-1.5-2, LLaVA-1.5-3, LLaVA-1.5-4.

ones). We also notice a significant decrease in IoU at the first layers, which might reveal that the first layers encode general features that are shared across modalities.

Different LLMs for the ST setup. In Fig. 24, we show the IoUs of activated weights for different frozen LLMs (*i.e.*, OPT, Llama 2 and Vicuna-v1.5) for the ST setup. We notice similar observations across LLMs, and relatively higher overlap for OPT. We notice similar observation compared to the MT setup, where we have a significant decrease in the IoU in first layers. For this setup, it is clearer that the overlap increase for deeper layers.

Different sparsity levels. In Fig. 25, we study how the overlap between activated weights changes with the size of the extracted subnetworks. This size depends on the sparsity of the final model. We

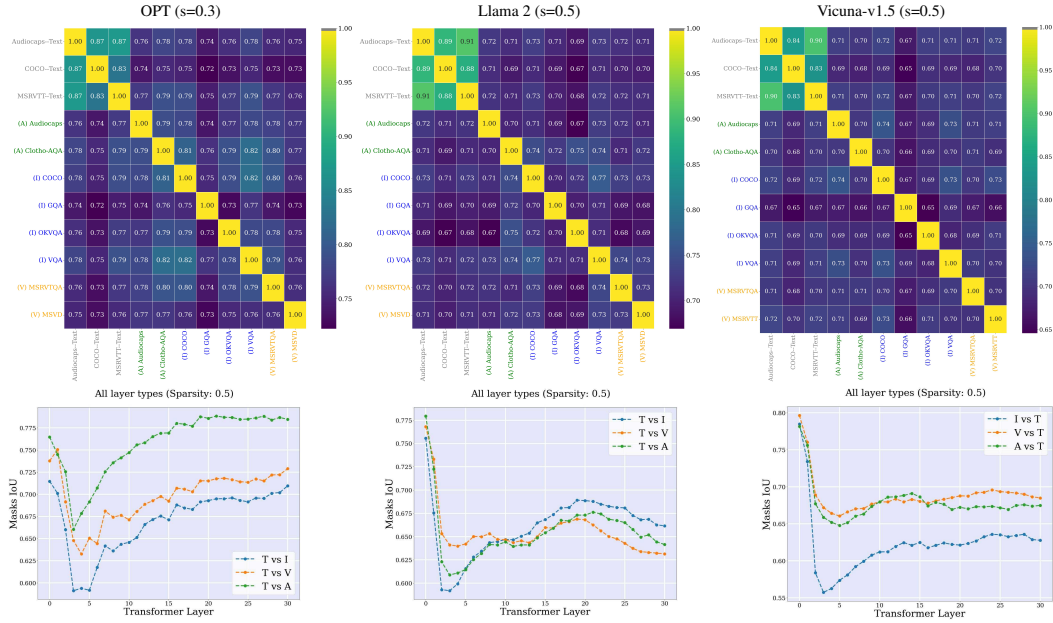


Figure 24: IoUs of activated subnetworks for different LLMs. We compute the IoUs for weights activated by different multimodal tokens. From left to right for the ST setup: OPT, Llama 2, Vicuna-v1.5.

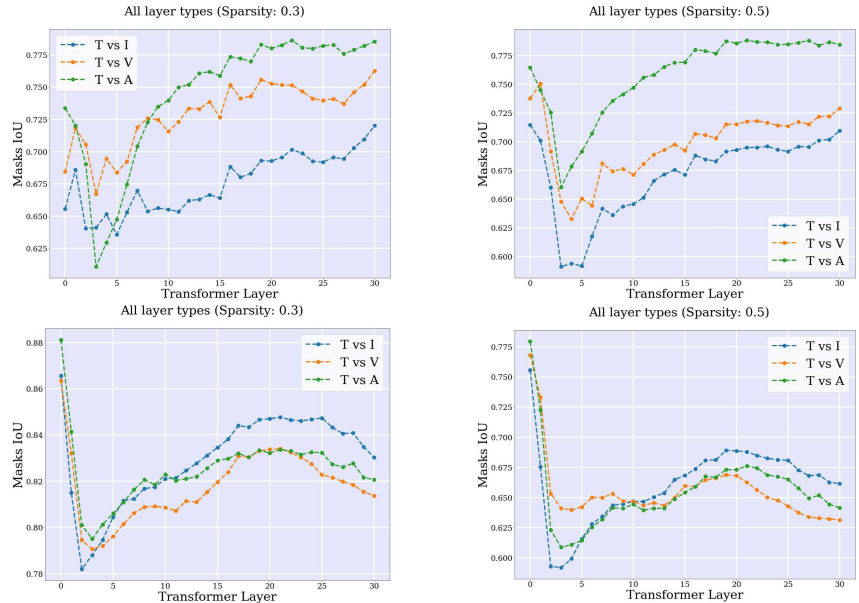


Figure 25: Overlap between multimodal subnetworks at different sparsity levels.. We compute the IoU of activated weights across layers at 0.3 and 0.5 sparsity levels. OPT (top), Llama 2 (bottom).

notice that the lower the sparsity, the higher the overlap, revealing that higher sparsity allows to extract more modality-specific activated weights.

Pruning weights by streaming different modalities. In addition to the IoU, we also compare the differences between the task performance of activated weights. For each dataset, we give to the model either the perceptual prompt, the textual prompt or both, knowing that each example in the dataset consists of a perceptual prompt followed by the textual one. Fig. 26 shows slight differences in overall performance when the LLM is pruned by different modalities, with the best performance is

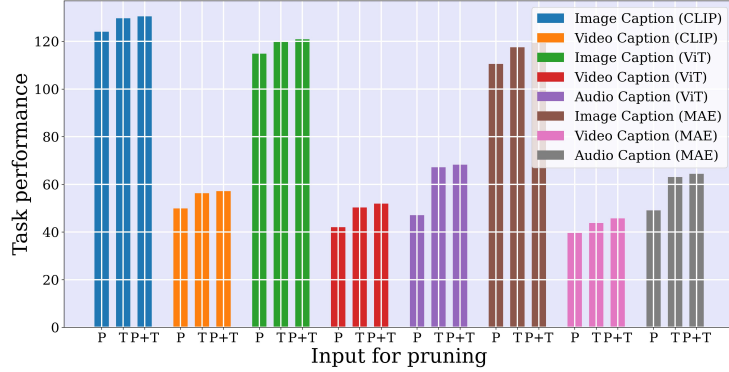


Figure 26: **Similar LLMs weights are activated by different modalities.** We report the task performance when keeping only the subnetwork activated by: multimodal prompt (P), text (T), or both (P+T) for OPT. Sparsity level: 0.3.

by considering both the textual and the perceptual tokens. This also show the high overlap between the weights used to activate textual and perceptual tokens.

Transfer of pruning masks to other tasks and modalities. To further highlight the overlap between weights, we report the performance when the model is pruned given data from other modalities or datasets. We notice similar observations across LLMs, such Llama 2 Fig. 29 and Vicuna-v1.5 Fig. 30. Interestingly, we find the similar overlap with the unsupervised MAE encoders Fig. 27 compared to text aligned ones Fig. 28. We notice a performance degradation when the model is pruned at high sparsity levels (0.5).

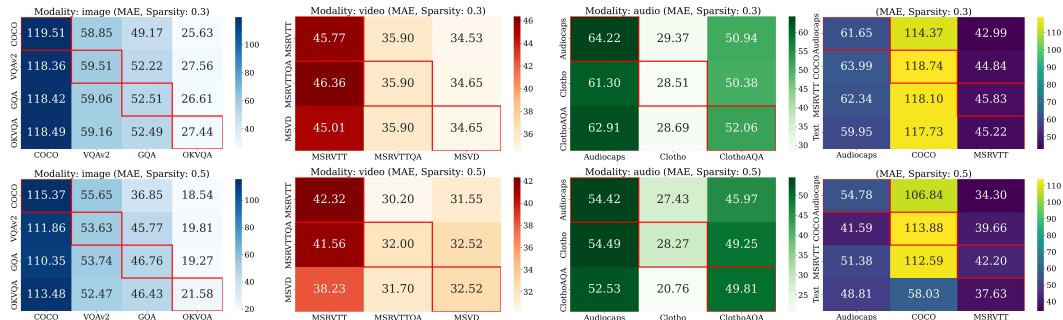


Figure 27: **Transfer of multimodal subnetworks across tasks and modalities with OPT and MAE encoders.** We use the subnetwork activated by a given task/modality to other tasks/modalities and report the task performance. From left to right, transfer across: image tasks, video tasks, audio tasks and across modalities for the captioning task. In each figure, the row corresponds to the source dataset of the subnetwork and the column to the target dataset.

Modality-specific subnetworks? The experiments suggest a high overlap between weights activated by different modalities. For instance, the pruning masks similarity (IoU) between datasets within the same modality is on par with those across modalities for the ST setup. However, this does not exclude the possibility of finding weights that are generally activated when seeing a particular modality, even if there are small amount of them. The overlap is smaller with LLaVA-1.5 variants making this possibility more likely for large scale multitask models.

E.5 Implicit multimodal alignment effect

Alignment inside each LLM Block. Fig. 32 reports the tokens similarity and norms for both LLaVA-1.5 and Vicuna-v1.5.

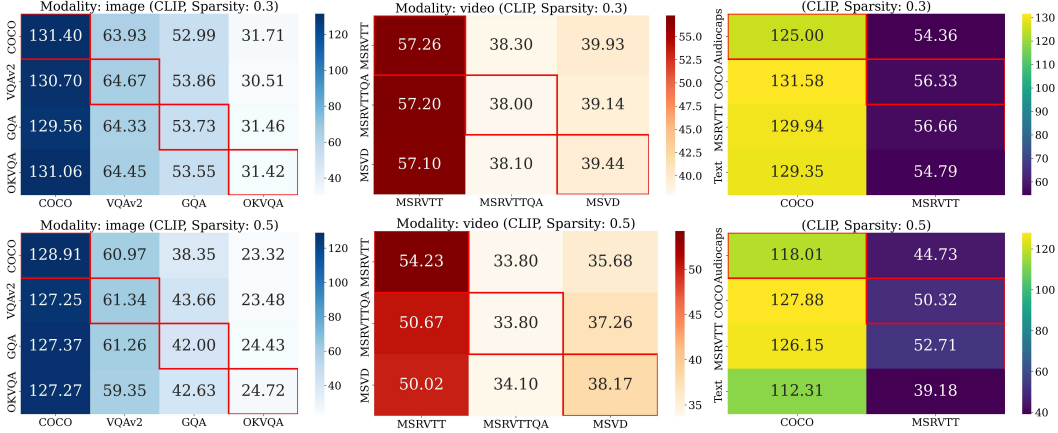


Figure 28: **Transfer of multimodal subnetworks across tasks and modalities with OPT and CLIP encoders.** We use the subnetwork activated by a given task/modality to other tasks/modalities and report the task performance. From left to right, transfer across: image tasks, video tasks, audio tasks and across modalities for the captioning task. In each figure, the row corresponds to the source dataset of the subnetwork and the column to the target dataset.

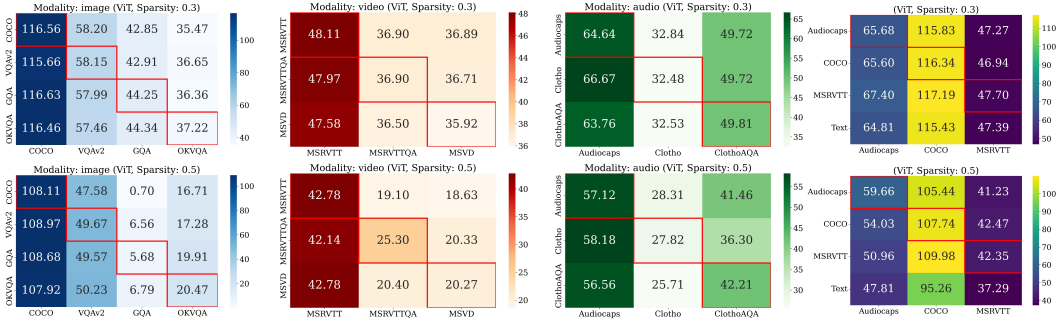


Figure 29: **Transfer of multimodal subnetworks across tasks and modalities with Llama 2 and ViT encoders.** We use the subnetwork activated by a given task/modality to other tasks/modalities and report the task performance. From left to right, transfer across: image tasks, video tasks, audio tasks and across modalities for the captioning task. In each figure, the row corresponds to the source dataset of the subnetwork and the column to the target dataset.

F Implications on performance, safety and efficiency: additional experiments

F.1 Implicit multimodal alignment as proxy metric for task performance?

IMA score across epochs. Fig. 33 shows an increasing similarity between textual and perceptual tokens during training with OPT and Vicuna-v1.5.

Table 1: **IMA score across different encoders.** We report the IMA score and the task performance with the ST setup (OPT). A positive correlation exists between IMA score and the performance; the most aligned encoders (CLIP) have the best accuracy/CIDEr on VQA and captioning tasks.

LLM	Encoder	IMA Score \uparrow	COCO \uparrow	VQAv2 \uparrow	GQA \uparrow
			CIDEr (test)	Acc (Val)	Acc (Val)
Vicuna-v1.5	CLIP-ViT-L	0.130	127.63	63.05	54.34
Vicuna-v1.5	ViT-L (ImageNet)	0.105	116.76	61.27	51.57
Vicuna-v1.5	MAE-L	0.060	76.40	59.57	52.88

IMA score across different encoders. In Table 1 we compare different image encoders and report the the IMA score and the task performance on several VQA and image captioning tasks. Encoders

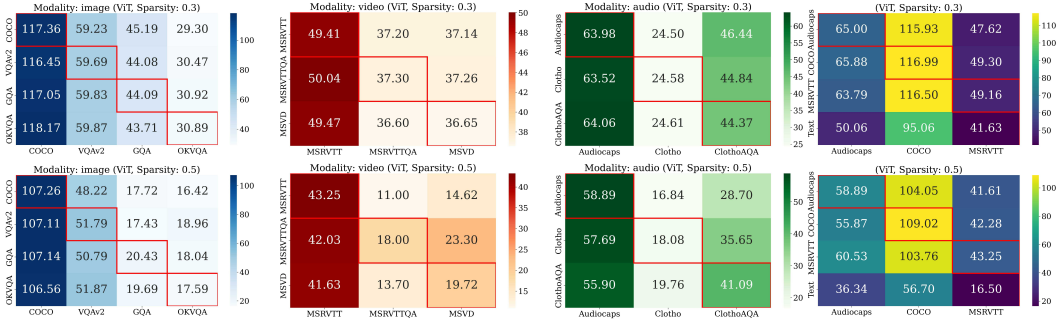


Figure 30: Transfer of multimodal subnetworks across tasks and modalities with Vicuna-v1.5 and ViT encoders. We use the subnetwork activated by a given task/modality to other tasks/modalities and report the task performance. From left to right, transfer across: image tasks, video tasks, audio tasks and across modalities for the captioning task. In each figure, the row corresponds to the source dataset of the subnetwork and the column to the target dataset.

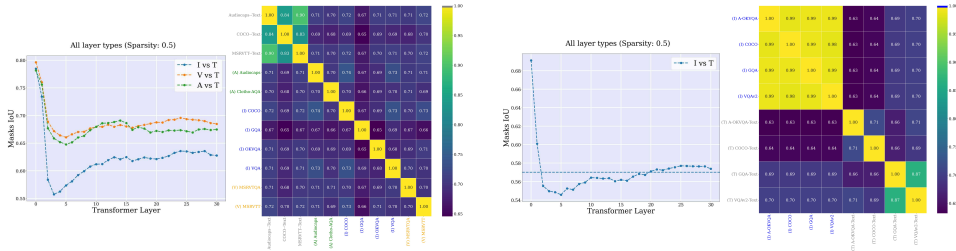


Figure 31: High similarity between LLM weights activated by different modalities. We compute the IoU of the subnetworks activated by different tasks across modalities, for the ST (left) and MT (right) setups.

that are most aligned to textual tokens inside LLMs (highest IMA *e.g.* CLIP) have also the best task performance.

F.2 Implicit multimodal alignment as proxy metric for hallucination?

We provide more details regarding the hallucinations metrics. These metrics are supposed to measure multimodal or object hallucination.

CHAIR on COCO captioning [51]. On the COCO image captioning dataset, the model is asked to describe the images. We compute the CHAIR metrics based on the generated captions and the ground truth annotations of all objects in the image. If a caption contains non-existent objects, we classify it as a hallucinated caption. The CHAIRs score is the ratio of hallucinated captions to the total number of captions. Additionally, we calculate the ratio of hallucinated objects to the total number of objects across all captions, which is referred to as CHAIRi. A CHAIR score of 0 indicates no hallucinations. In the paper, we report $(1 - \text{CHAIR}) \times 100$, thus a higher score indicates fewer hallucinations.

POPE benchmark [50]. This is a question-answering task involving questions about the existence of objects in images. The metric used is accuracy; the fewer the hallucinations, the higher the accuracy.

F.3 Skipping computations for visual tokens.

In this section, we propose to skip computations for the visual tokens.

Skip FFN Tokens. We randomly skip a number of tokens (we refer to this amount as skip ratio), both the textual and the remaining visual tokens are processed in the FFN layers. ??, shows a linear relationship between the skip ratio and task performance, the higher the ratio the lower the scores.

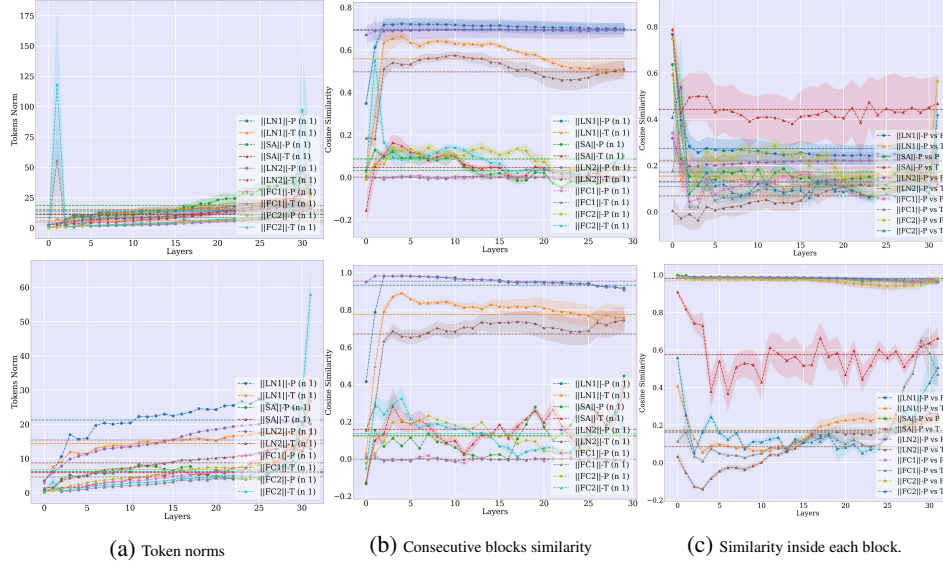


Figure 32: **Implicit alignment inside the LLM blocks.** We compute the token norms (left), tokens cosine similarity between consecutive blocks (middle) and across modalities (last). The tokens are inside the LLM blocks (and outside the residual stream): after the self-attention (SA), and FFNs (FC1/2) and layer norms (LN). From top to down: Vicuna-v1.5, LLaVA-1.5-4.

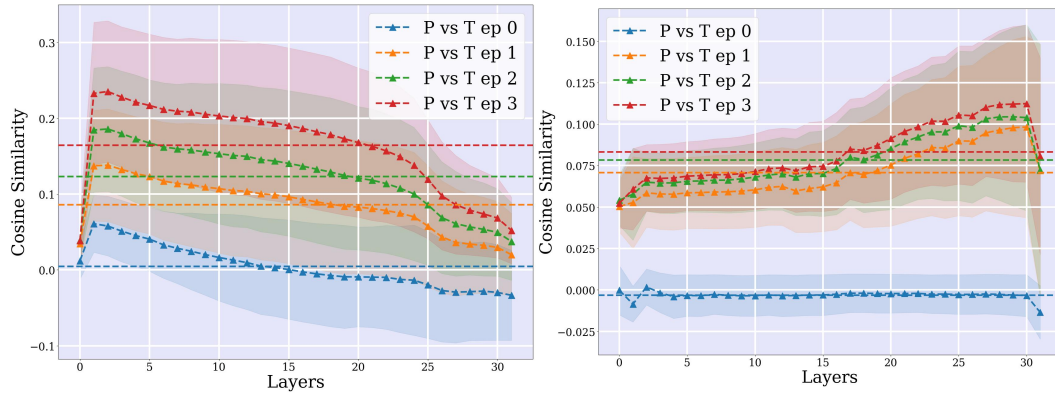


Figure 33: **Implicit alignment score across epochs.** We report the implicit alignment score for OPT (left) and Vicuna-v1.5 (right) during training of the mapping module.

F.4 α -SubNet

In Table 2, we evaluate our α -SubNet on additional multimodal tasks. Compared to other task-agnostic baselines such as magnitude pruning the scores of α -SubNet are significantly higher. This support more that the multimodal tokens activate similar weights inside LLMs.

Table 2: **α -SubNet a task and modality-agnostic subnetwork.** We prune the LLMs using different post-training pruning methods, including our α -SubNet.

Method	#P/#TP/Sparsity	Avg	COCO \uparrow CIDEr (test)	VQAv2 \uparrow Acc (Val)	OKVQA \uparrow Acc (Val)	GQA \uparrow Acc (Val)	MSR-VTT \uparrow CIDEr (test)	MSRVTT-QA \uparrow Acc (test)	MSVD-QA \uparrow Acc (test)	Audiocaps \uparrow CIDEr (test)	Clotho \uparrow CIDEr (test)	Clotho-AQA \uparrow Acc (test)
MAPL [12]	7B/3.4M/0.00	-	125.2	43.5	18.7 / 31.6	-	-	-	-	-	-	-
eP-ALM [9]	6.7B/4M/0.00	-	111.6	54.9	-	42.91	48.79	35.90	38.40	61.86	-	-
DePALM [10]	7B/18.1M/0.00	-	131.29	70.11	37.69	-	49.88	-	-	69.70	-	-
Baseline	6.7B/7M/0.00	57.71	132.83	63.49	33.01	55.29	58.23	38.84	38.83	68.24	35.66	52.72
Wanda	6.7B/7M/0.50	51.32 (88.93%)	126.81	55.28	24.72	42.00	54.23	33.80	37.17	58.99	31.81	48.41
Random mask	6.7B/7M/0.47	0.00 (0%)	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.05	0.00	0.00
α -SubNet ($\alpha=0.3$)	6.7B/7M/0.47	39.34 (68.17%)	106.77	51.77	17.72	38.09	38.37	29.80	31.19	23.15	8.52	48.03

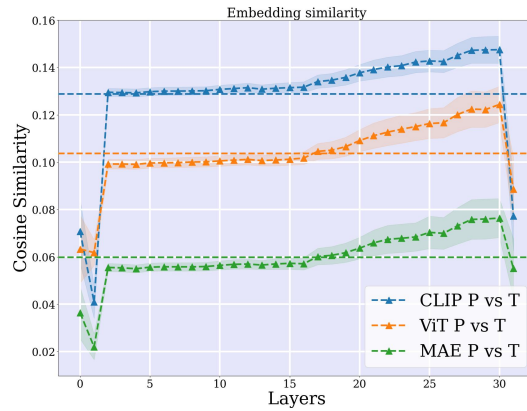


Figure 34: **Comparison of implicit multimodal alignment score across layers for different encoders.** CLIP models produce features that are most aligned to textual tokens across LLM layers. On the other hand, self-supervised encoders (*e.g.* MAE) produce the least text-aligned features. However, the relatively low cosine similarity score (closer to 0), reveals that the modality gap (*e.g.* Narrow cones) still exists in LLMs, even for text-aligned encoders.

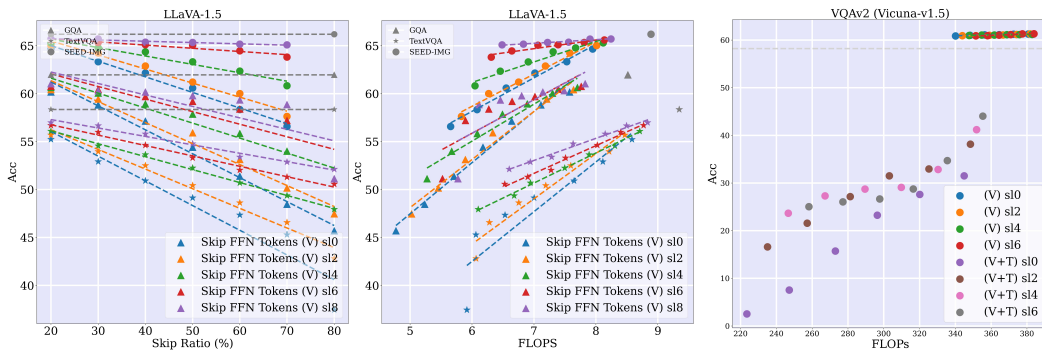


Figure 35: **Skipping computations for visual tokens.** Skip Tokens: we skip (Skip ratio)% of the tokens in the FFN layers. sl: skipping start layer. (V): visual tokens. (T): textual tokens. Results MT (with LLaVA-1.5) and ST (last column) setups.

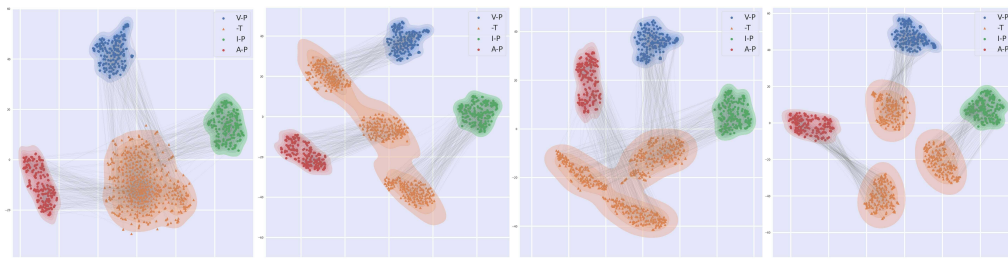


Figure 36: **t-SNE visualization of tokens inside LLMs.** From left to right: layer 0, 1, 24 and 32 for Vicuna-v1.5.

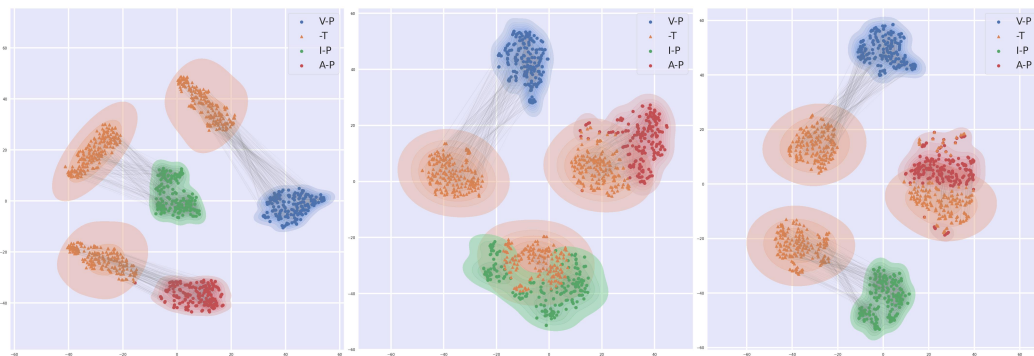


Figure 37: **t-SNE visualization of tokens inside LLM blocks (after SA layers).** From left to right: layer 0, 1, 24 and 32 for Vicuna-v1.5. There is less separation inside the LLM block.