UNDERSTANDING THE DESIGN SPACE AND CROSS-MODALITY TRANSFER FOR VISION-LANGUAGE MODELS

Anonymous authorsPaper under double-blind review

ABSTRACT

The training of multimodal models involves many design choices, such as the underlying modality-specific tokenizers, fusion mechanisms, and strategies for freezing model layers during different training stages. However, the individual impact of these decisions on downstream multimodal performance remains poorly understood due to the diversity of current practices. In this paper, we systematically investigate how choices in image tokenization, architectural design, and layer-freezing strategies affect the training and cross-modal generalization of vision-language models (VLMs). We train and evaluate over 50 VLM variants across a controlled suite of tokenizers, model architectures, and training recipes. Our experiments reveal several key trends: (1) image tokenizers designed with text alignment in mind, together with training recipes that further enhance image-text alignment, yield the best performance; (2) unfreezing the language model boosts in-domain results but can degrade out-of-domain generalization; and (3) fusion architectures based on the mixture-of-transformers architecture are effective, especially when language parameters are frozen. To further probe cross-modality transfer, we introduce three new synthetic datasets, which we use to evaluate our pretrained models.

1 Introduction

Vision-language models (VLMs) are frequently built by combining a pretrained large language model (LLM) with an image tokenizer via a fusion architecture that integrates image and text representations. Here, we use the term *tokenizer* to refer to any module that processes raw modality inputs (such as images or text) into a sequence of discrete or continuous tokens for downstream modeling. This encompasses both discretization approaches such as VQVAE (van den Oord et al., 2017), and conventional encoders such as CLIP (Radford et al., 2021) and SigLIP (Zhai et al., 2023), which generate token-like embeddings.

Despite rapid progress, the effects of architectural and training choices, such as how to align and fuse modality-specific representations, remain unclear, due to the proliferation of fusion architectures (e.g., joint autoregressive decoders (Liu et al., 2023; Deitke et al., 2024; Bai et al., 2025; Du et al., 2025; Zhu et al., 2025), cross-attention models (Alayrac et al., 2022; Grattafiori et al., 2024; Dai et al., 2024) and mixtures-of-transformers (Liang et al., 2025; Shi et al., 2025b; Deng et al., 2025)), tokenization schemes (Tschannen et al., 2025; Fini et al., 2024; Oquab et al., 2023; Yu et al., 2024; Tian et al., 2024; Bachmann et al., 2025; Miwa et al., 2025), and multi-stage training recipes with varied layer freezing. This diversity of approaches makes it challenging to disentangle how each design choice impacts the performance and generalization behavior of VLMs.

Understanding how design and training strategies of multimodal models enable cross-modality transfer is particularly important. Effective transfer allows models to develop reasoning and understanding that may be more naturally expressed in one modality than another (for example, physical reasoning may be more apparent in vision than in text). Furthermore, robust cross-modal transfer enables models to leverage alternative data sources, which is increasingly valuable as high-quality text data becomes scarce due to the growing scale of LLM training.

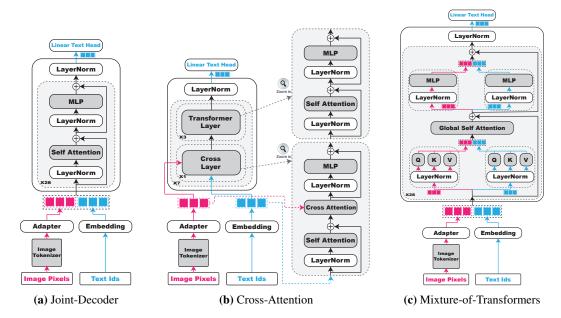


Figure 1: Different multimodal fusion architectures: (a) Joint-Decoder, (b) Cross-Attention, (c) Mixture-of-Transformers.

In this work, we systematically investigate the influence of fusion architectures, image tokenizers, and training recipes on VLM performance. We train and evaluate over 50 VLM variants, all built on a Qwen3-0.6B (Yang et al., 2025) LLM backbone and trained on a mix of image-caption and vision question answering (VQA) data. We assess in-domain and out-of-domain performance using a suite of VQA benchmarks to discover trade-offs between architectural and training choices. To further explore cross-modality transfer, we also introduce three new synthetic datasets designed for evaluating cross-modality generalization.

Our main findings are: (i) Image tokenizers trained to text-alignment objectives lead to better performance trained for image reconstruction, regardless of fusion architecture or whether layers are frozen. (ii) While unfreezing the LLM backbone provides the largest gains for in-domain performance, it often degrades out-of-domain generalization. Unfreezing only the image tokenizer offers a more balanced improvement. (iii) The Mixture-of-Transformer fusion architecture allows freezing the LLM backbone to achieve strong in-domain and out-of-domain performance while preserving text-only performance. (iv) The success of cross-modality transfer depends on both the model's internal design (i.e., architecture and tokenizer) and the representational alignment between the data modalities themselves.

2 EXPERIMENTAL SETTINGS

2.1 ARCHITECTURES

We study three fusion architectures for vision-language modeling, illustrated in Figure 1: Joint-Decoder, Cross-Attention, and Mixture-of-Transformers (MoT).

Joint-Decoder: In the Joint-Decoder architecture (Deitke et al., 2024; Bai et al., 2025; Du et al., 2025; Zhu et al., 2025), image and text tokens are concatenated and then fed into a shared multimodal transformer decoder (Figure 1a).

Cross-Attention: In the Cross-Attention architecture, cross-attention layers allow text tokens to attend to image token representations (Figure 1b), following architectures like Flamingo (Alayrac et al., 2022) and Llama 3-V (Grattafiori et al., 2024). Design variants exist in how and where visual information is injected. For example, NVLM (Dai et al., 2024) explores different placements and use of special visual tokens for cross-modal reasoning.

Mixture-of-Transformers (MoT): For the Mixture-of-Transformers architecture (Liang et al., 2025; Shi et al., 2025b; Deng et al., 2025), the model processes each modality through its own

stack of transformer layers, referred to as *modality-transformers*, each with separate query-key-value (QKV) matrices and feed-forward networks. At every layer, tokens are first routed through their respective modality-transformer for QKV computation, then mixed globally via multimodal self-attention, and finally passed back through that modality's feed-forward sublayer (Figure 1c). This doubles the parameter count of the fusion layers relative to the analogous Joint-Decoder architecture, but, as shown in (Liang et al., 2025), the overall floating-point operations (FLOPs) per forward pass remain comparable.

All of our models are built with Qwen3-0.6B, which contains of 28 transformer layers, as the language backbone. While Qwen3-0.6B ties its text embedding and output head, we untie them for our models. In the **Joint-Decoder** architecture, these 28 layers are repurposed as a single, shared multimodal decoder initialized with the pretrained LLM weights. For the **Cross-Attention** models, we follow the design of Llama 3-V (Grattafiori et al., 2024) and interleave a cross-attention layer every four layers within the backbone, adding a total of seven new cross-attention layers. Finally, the **Mixture-of-Transformers** (MoT) architecture, following Shi et al. (2025b); Deng et al. (2025), creates two parallel modality-transformers, each containing a full copy of the 28 transformer layers, both of which are initialized with the original Qwen3-0.6B weights.

2.2 IMAGE TOKENIZERS AND ADAPTERS

We experiment with a range of image tokenizers, each trained with different objectives:

Continuous tokenizers:

- CLIP (Radford et al., 2021): Trained via contrastive learning to align image and text embeddings.
- SigLIP 2 (Tschannen et al., 2025): Trained using a combination of contrastive learning with a sigmoid loss (Zhai et al., 2023), an autoregressive captioning loss (Wan et al., 2024), and self-distillation (Naeem et al., 2025; Maninis et al., 2025).
- AIMv2 (Fini et al., 2024): Trained with next patch prediction (for images) and an autore-gressive captioning loss (for text).

• Discrete tokenizers:

- TiTok (Yu et al., 2024): Trained to encode images into one-dimensional latent token sequences by reconstructing ground-truth two-dimensional latents.
- VAR (Tian et al., 2024): Trained to autoregressively reconstruct images via multi-scale token maps.

To match tokens to fusion architectures, we include lightweight adapter modules. For Joint-Decoder and MoT, continuous tokenizers are projected using a two-layer MLP, while discrete tokenizers use an embedding layer. In the Cross-Attention architecture, adapters follow Llama 3-V, aligning dimensions for image tokens (using an embedding layer for discrete tokenizers and dimensionality matching for continuous ones).

2.3 Training Setup

We train each of our models in three stages: 1) a pretraining stage; 2) a VQA fine-tuning stage; 3) and a reasoning-transfer stage. The hyperparameters for all three training stages can be found in Appendix D.

Stage 1 (Pretraining): We pretrain models for caption generation on COYO-700M (Byeon et al., 2022) to align the image tokenizer and any uninitialized model weights, improving the representations available to the LLM layers. The language model layers are kept frozen in this stage. For each model with a continuous image tokenizer, we run both frozen and unfrozen variants to enable downstream comparison. Discrete image tokenizers remain frozen throughout, avoiding the need for a straight-through estimator.

During pretraining, we only train the Joint-Decoder and Cross-Attention models. For MoT models, we follow Shi et al. (2025b); Deng et al. (2025) and initialize weights from a trained Joint-Decoder checkpoint by transferring the adapter and image tokenizer weights, and copying the original Qwen3 weights into both modality-transformers.

Stage 2 (Fine-tuning): We fine-tune each pretrained checkpoint on a combination of COCO-Captions (Lin et al., 2014), VQAv2 (Goyal et al., 2017), DocVQA (Mathew et al., 2021), TextVQA (Singh et al., 2019), and ChartQA (Masry et al., 2022). We systematically ablate whether the image tokenizer or LLM layers are frozen at this stage, such that unfrozen layers during pretraining remain unfrozen for fine-tuning. Models are evaluated after each epoch on the validation splits of VQAv2, A-OKVQA (Schwenk et al., 2022), DocVQA, TextVQA, and ChartQA. The checkpoint with the highest mean validation accuracy across VQA datasets is selected for further evaluation.

Stage 3 (Reasoning-Transfer): We further train our models on three synthetic datasets, detailed in Section 4 and Appendix C. Each dataset pairs an image with an equivalent text description, allowing for both image-based VQA and comparable text-only QA training runs. We evaluate on in-distribution and out-of-distribution tasks in both modalities to quantify cross-modality transfer.

2.4 EVALUATION PROTOCOL

To assess model capabilities, we evaluate on a combination of standard academic vision-language benchmarks and our own synthetic datasets. These benchmarks measure both general visual understanding and the ability to transfer knowledge across domains and modalities. We group academic benchmarks into in-domain and out-of-domain: out-of-domain benchmarks not only cover novel topics but also feature visual multiple-choice questions, a format our models were not exposed to during training.

For academic evaluations, our in-domain suite tests a range of capabilities, including general VQA with VQAv2 (Goyal et al., 2017), knowledge-based reasoning with A-OKVQA (Schwenk et al., 2022), and specialized understanding of documents (DocVQA (Mathew et al., 2021)), text in images (TextVQA (Singh et al., 2019)), and charts (ChartQA (Masry et al., 2022)). To assess out-of-domain generalization, we use MathVista (Lu et al., 2024b) for mathematical reasoning, Real-WorldQA (xAI, 2024) for robustness to novel image distributions, and the multi-task benchmark MMTBench (Ying et al., 2024).

3 Understanding the Design of Multimodal Architectures

3.1 EFFECT OF IMAGE TOKENIZER ON IN-DOMAIN PERFORMANCE

Tokenizer	Lang	VQAv2 test-dev	A-OKVQA val	ChartQA test	TextVQA val	DocVQA test	Average In-domain	Average Out-of-domain
CLIP	F	49.8	22.1	12.2	20.1	12.1	23.2	30.5
	U	66.7	41.6	20.8	31.9	20.7	36.3	32.9
AIMv2	F	62.5	30.3	19.5	33.6	17.9	32.8	32.4
	U	75.3	49.5	30.6	43.2	25.2	44.7	34.5
SigLIP 2	F	55.9	24.0	15.7	30.0	17.0	28.5	31.7
	U	74.8	47.2	28.8	43.6	26.6	44.2	32.8
TiTok	F	3.2	0.2	5.1	1.0	3.0	2.5	26.7
	U	43.1	26.0	13.9	11.8	11.8	21.3	27.1
VAR	F	30.3	2.0	9.2	4.9	6.0	10.5	28.1
	U	46.5	27.7	13.7	11.9	11.9	22.3	22.4

Table 1: Evaluation results (accuracy, %) on the Joint-Decoder architecture with frozen image tokenizer. The table compares results across various image tokenizers while ablating whether the language model (Lang) is frozen (F) or unfrozen (U) during Stage 2.

Table 1 presents the accuracies for different image tokenizers within the Joint-Decoder architecture, with the image tokenizer frozen to normalize for the trainability of the discrete tokenizers.

The results consistently show that image tokenizers trained with text-alignment objectives (AIMv2, SigLIP 2, CLIP) substantially outperform those trained for image reconstruction (TiTok, VAR) for both in-domain and out-of-domain tasks. Among the text-supervised tokenizers, a further hierarchy emerges for in-domain tasks: AIMv2 and SigLIP 2, which incorporate stronger text objectives like autoregressive captioning, outperform the purely contrastively trained CLIP. This general trend holds true across the cross-attention and MoT architectures and is consistent regardless of the layer freezing strategy. Full evaluation results are available in Appendix A.

Takeaway 1. Image tokenizers trained with text-alignment objectives are crucial for strong VLM performance, significantly outperforming those trained solely on image reconstruction on both indomain and out-of-domain tasks.

Takeaway 2. Stronger text-alignment objectives (e.g., autoregressive captioning vs. contrastive loss) provide a clear advantage for in-domain performance.

3.2 VARYING FROZEN LAYERS AND ARCHITECTURAL CHOICES

Stage 1	Stage 2		Joint-D	ecoder	Cross-A	ttention	MoT		
Image	Image	Lang	In	Out	In	Out	In	Out	
F	F	F	28.2 (+0.0)	31.5 (+0.0)	37.7 (+0.0)	28.6 (+0.0)	38.8 (+0.0)	34.0 (+0.0)	
F	F	U	41.8 (+13.6)	33.4 (+1.9)	43.5 (+5.8)	25.6 (-3.0)	42.3 (+3.5)	31.5 (-2.5)	
F	U	F	35.5 (+7.3)	32.9 (+1.4)	41.4 (+3.7)	27.6 (-1.0)	43.2 (+4.4)	35.6 (+1.6)	
F	U	U	45.8 (+17.6)	32.9 (+1.4)	46.8 (+9.1)	25.9 (-2.7)	46.3 (+7.5)	29.8 (-4.2)	
U	U	F	39.9 (+11.7)	34.3 (+2.8)	43.0 (+5.3)	31.3 (+2.7)	45.8 (+7.0)	36.3 (+2.3)	
U	U	U	47.7 (+19.5)	33.1 (+1.6)	47.7 (+10.0)	27.0 (-1.6)	47.8 (+9.0)	31.6 (-2.4)	

Table 2: Average in-domain (In) and out-of-domain (Out) performance (accuracy, %) across different fusion architectures and layer freezing strategies. Scores are averaged over all text-supervised image tokenizers (CLIP, AIMv2, SigLIP 2). The table ablates the freezing status, Frozen (F) or Unfrozen (U), of the image tokenizer (Image) and language model (Lang) during Stage 1 (pretraining) and Stage 2 (fine-tuning). Values in parentheses indicate the change from the fully frozen baseline in the first row.

Table 2 summarizes the effects of unfreezing the image tokenizer and/or LLM during both pretraining (Stage 1) and fine-tuning (Stage 2) across different fusion architectures.

Unfreezing the image tokenizer consistently provides moderate improvements in in-domain accuracy, while unfreezing the LLM backbone yields even larger gains—most notably for Joint-Decoder models, which have less capacity for multimodal integration than Cross-Attention or MoT. When both the image tokenizer and language layers are frozen, MoT outperforms the other architectures, thanks to the extra trainable parameters it has (over 400M more parameters from the image modality-transformer). However, as more layers are unfrozen, all architectures perform more similarly. In contrast, unfreezing the LLM often leads to a reduction in out-of-domain performance, whereas unfreezing only the image tokenizer can provide a modest out-of-domain boost.

Takeaway 3. Unfreezing both image and language layers maximizes in-domain performance, but may hurt out-of-domain generalization.

Our results show that the Mixture-of-Transformers (MoT) architecture is particularly effective. By providing dedicated parameters for multimodal integration without increasing the FLOPs per forward pass (Liang et al., 2025; Shi et al., 2025b), MoT enables a strategy where the language model can be frozen. This approach yields strong performance on both in-domain and out-of-domain tasks, preserves the LLM's original text-only capabilities, and remains computationally efficient. This finding aligns with recent work (Dai et al., 2024; Lin et al., 2024; Shi et al., 2025b) and highlights the value of strategies that limit LLM supervision, especially when high-quality training data is scarce.

Takeaway 4. MoT with an unfrozen image tokenizer and frozen language layers delivers the best overall task performance.

4 Cross-modality Transfer Learning

In this section, we systematically evaluate *cross-modality transfer* across the models we trained. Motivated by similar settings in prior work (Wang et al., 2024; Yamada et al., 2024), we construct

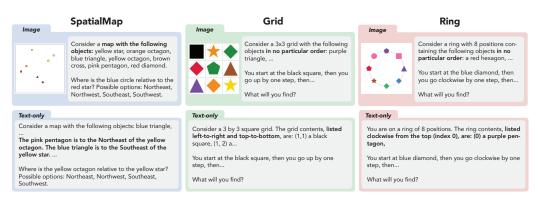


Figure 2: Examples from the three synthetic datasets used to study cross-modality transfer. Each dataset provides paired image and text-only versions with equivalent information but distinct spatial structures. **SpatialMap** (left, blue) places objects on an open canvas and asks about relative positions between two objects. **Grid** (middle, green) arranges objects in a 3×3 lattice and requires tracking movement across cells. **Ring** (right, red) positions objects cyclically and requires tracking clockwise or counterclockwise traversal.

three synthetic datasets designed to isolate reasoning from perception. Each dataset pairs procedurally generated images with equivalent text-only representations and is built around a distinct spatial reasoning task: SpatialMap (objects on an open canvas), Grid (objects in a 2D grid), and Ring (objects arranged cyclically), see Figure 2. Because the image and text modalities contain identical information, we can directly measure a model's ability to transfer learned concepts. Each dataset provides 4,500 training examples, 500 in-distribution (InD) test examples, and 500 out-of-distribution (OOD) test examples designed to be compositionally harder. Further details on dataset generation are in Appendix C.

To measure transfer, we perform further fine-tuning (as described in Section 2.3) on either the image-based (VQA) or text-only (QA) version of a task. We then evaluate each fine-tuned model across two axes of generalization:

- Same-Modality vs. Cross-Modality: We test on tasks using the same modality as the synthetic data (e.g., VQA → VQA) and on tasks using the other modality (e.g., VQA → QA).
- In-Distribution (InD) vs. Out-of-Distribution (OOD): We test on both the standard test set (InD) and the more challenging, compositionally distinct test set (OOD).

In the following subsections, we present results from this further fine-tuning on the image-based tasks (Section 4.1) and the text-only tasks (Section 4.2). An analogous analysis performed on several open-weight VLMs is provided in Appendix B.

4.1 Cross-Modality Transfer from Image to Text

In this section, we evaluate image-to-text cross-modality transfer. To do this, we fine-tune our stage-2 models on the image-based (VQA) version of our synthetic datasets (SpatialMap, Grid, and Ring). We then evaluate performance on both the original image-based task and the unseen, text-only version of the task. This allows us to measure how well knowledge acquired from visual inputs transfers to a purely textual domain. The distinction between these modalities is illustrated in Figure 2.

Table 3 presents these results, averaged across all fusion architectures to isolate the impact of the image tokenizer. A clear trend emerges: models using text-supervised continuous tokenizers (CLIP, AIMv2, SigLIP 2) demonstrate strong performance on the original VQA task and successfully transfer knowledge to the text-only task. In contrast, models with discrete, reconstruction-based tokenizers (TiTok, VAR) struggle, showing substantially less learning on the image task and less transfer. While transfer to out-of-distribution tasks is challenging for all models, the text-supervised tokenizers still show a distinct advantage. This performance gap indicates that text-alignment pretraining is critical for developing representations that support both visual and textual reasoning, a capability not fostered by image reconstruction objectives alone.

Takeaway 5. Text-supervised image tokenizers perform better than the image tokenizers trained for image reconstruction in image-to-text cross-modality transfer.

Dataset	Tokenizer	VQA/InD	VQA/OOD	Text-only/InD	Text-only/OOD
SpatialMap	CLIP	57.6 (+41.3)	54.8 (+38.4)	30.7 (+13.6)	28.5 (+11.5)
	AIMv2	89.6 (+66.8)	86.0 (+62.7)	33.1 (+12.5)	33.2 (+15.6)
	SigLIP 2	94.3 (+75.3)	90.1 (+70.0)	36.5 (+ 7.2)	34.9 (+4.1)
	TiTok	25.6 (+14.6)	25.2 (+14.4)	23.2 (+0.2)	23.6 (+0.7)
	VAR	24.2 (+16.9)	26.7 (+17.9)	26.2 (-3.0)	26.6 (-0.6)
Grid	CLIP	92.0 (+90.0)	17.2 (+16.8)	64.3 (+ 61.0)	19.4 (+18.6)
	AIMv2	95.4 (+91.6)	15.5 (+14.1)	48.4 (+42.6)	18.8 (+16.6)
	SigLIP 2	98.6 (+98.1)	17.4 (+16.8)	55.9 (+53.4)	16.4 (+14.8)
	TiTok	30.4 (+30.2)	12.2 (+11.8)	23.1 (+21.3)	8.0 (+7.2)
	VAR	32.0 (+31.9)	8.6 (+8.6)	31.4 (+31.0)	9.6 (+9.1)
Ring	CLIP	97.0 (+95.3)	10.2 (+9.2)	99.6 (+97.0)	15.2 (+12.4)
	AIMv2	99.4 (+98.3)	8.6 (+8.5)	97.9 (+93.9)	13.5 (+9.9)
	SigLIP 2	99.4 (+99.2)	8.9 (+8.8)	82.2 (+79.3)	10.2 (+8.6)
	TiTok	22.6 (+22.2)	9.3 (+9.3)	38.5 (+35.3)	9.1 (+8.8)
	VAR	25.1 (+25.1)	8.8 (+8.8)	37.4 (+36.2)	10.2 (+9.0)

Table 3: Image-to-text transfer performance (accuracy, %) for each image tokenizer. Models were fine-tuned on the VQA version of the synthetic datasets, and scores are averaged across all fusion architectures. Each score is presented alongside the change in performance due to the additional training.

Dataset	Architecture	VQA/InD	VQA/OOD	Text-only/InD	Text-only/OOD
SpatialMap	Joint-Decoder Cross Attention MoT	79.4 (+54.8) 88.4 (+78.2) 73.6 (+50.3)	75.5 (+50.4) 83.8 (+72.2) 71.6 (+48.6)	34.8 (+2.4) 32.7 (+23.8) 32.8 (+7.1)	34.6 (+2.2) 30.2 (+21.8) 31.9 (+7.2)
Grid	Joint-Decoder Cross Attention MoT	95.4 (+90.4) 94.6 (+94.2) 96.0 (+95.2)	16.8 (+14.4) 15.1 (+15.1) 18.2 (+18.1)	67.3 (+ 62.5) 34.4 (+34.0) 66.9 (+60.4)	22.1 (+20.2) 14.1 (+13.7) 18.4 (+16.2)
Ring	Joint-Decoder Cross Attention MoT	98.8 (+96.6) 99.2 (+99.2) 97.9 (+97.0)	7.6 (+6.5) 9.6 (+9.6) 10.6 (+10.5)	99.2 (+ 94.4) 88.8 (+88.8) 91.8 (+87.1)	12.8 (+9.0) 13.0 (+13.0) 13.1 (+8.9)

Table 4: Image-to-text transfer performance (accuracy, %) by fusion architecture. Scores are averaged over text-supervised image tokenizers (CLIP, AIMv2, SigLIP 2) after models were fine-tuned on the VQA versions of the synthetic datasets. Each score is presented alongside the change in performance due to the additional training.

Table 4 isolates the effect of the fusion architecture by averaging results across the text-supervised tokenizers. Here, we observe that the degree of transfer is highly dependent on the representational alignment between the image and text modalities of a given dataset. For the Grid and Ring datasets, where the textual description is a direct, one-to-one serialization of the visual information, models achieve strong cross-modality transfer. For the SpatialMap dataset, where the text describes relative spatial locations, the textual description of relative coordinates is ambiguous and does not uniquely define the visual layout. Subsequently, transfer is significantly more challenging, resulting in much lower text-only accuracy. This suggests that the representational gap between modalities is a key bottleneck.

Takeaway 6. The success of cross-modality transfer is heavily influenced by the representational alignment between the source and target modalities. Transfer is more effective when the textual representation is a direct, structured description of the visual scene.

4.2 Cross-Modality Transfer from Text to Image

In this section, we evaluate the reverse direction: text-to-image transfer. To do this, we fine-tune our stage-2 models on the text-only (QA) version of our synthetic datasets and then evaluate performance on the corresponding, unseen image-based (VQA) tasks.

Dataset	Tokenizer	Text-only/InD	Text-only/OOD	VQA /InD	VQA/OOD
	CLIP	76.5 (+59.4)	74.9 (+57.9)	17.4 (+1.1)	17.8 (+1.4)
	AIMv2	70.2 (+49.6)	68.8 (+51.2)	16.6 (-6.2)	17.2 (-6.1)
SpatialMap	SigLIP 2	67.4 (+38.1)	69.1 (+38.3)	8.2 (-10.8)	9.9 (-10.2)
	TiTok	61.2 (+38.2)	61.8 (+38.9)	9.0 (-2.0)	9.1 (-1.7)
	VAR	71.2 (+42.0)	74.8 (+47.6)	17.4 (+10.1)	16.4 (+7.6)
	CLIP	99.8 (+96.5)	56.4 (+55.6)	10.2 (+8.2)	5.1 (+4.7)
	AIMv2	99.3 (+93.5)	29.1 (+26.9)	13.1 (+9.3)	7.1 (+5.7)
Grid	SigLIP 2	99.9 (+97.4)	58.6 (+57.0)	12.2 (+11.7)	5.8 (+5.2)
	TiTok	99.9 (+98.1)	35.6 (+34.8)	11.5 (+11.3)	7.3 (+6.9)
	VAR	99.1 (+98.7)	33.3 (+32.8)	11.6 (+11.5)	4.8 (+4.8)
	CLIP	100.0 (+97.4)	13.4 (+10.6)	12.4 (+10.7)	5.2 (+4.2)
	AIMv2	100.0 (+96.0)	13.6 (+10.0)	9.8 (+8.7)	5.3 (+5.2)
Ring	SigLIP 2	100.0 (+97.1)	14.1 (+12.5)	6.8 (+6.6)	4.6 (+4.5)
	TiTok	100.0 (+96.8)	14.3 (+14.0)	8.0 (+7.6)	5.6 (+5.6)
	VAR	100.0 (+98.8)	13.9 (+12.7)	9.0 (+9.0)	6.2 (+6.2)

Table 5: Text-to-image transfer performance (accuracy, %) by image tokenizer. Models were fine-tuned on the text-only (QA) version of the synthetic datasets, and scores are averaged across all fusion architectures. Each score is presented alongside the change in performance due to the additional training.

The results, grouped by image tokenizer in Table 5, reveal a stark asymmetry compared to the previous section. While all models, regardless of tokenizer, successfully learn the source text-only task (often reaching near-perfect in-distribution accuracy), this knowledge largely fails to transfer to the image domain. This is expected, as text-only fine-tuning provides no gradient signal to update the image tokenizer or its alignment with the language model. In many cases, especially for the SpatialMap dataset, performance on the VQA task degrades significantly from the stage-2 training baseline, indicating that adapting the LLM to a new text distribution can disrupt its previously learned alignment with the frozen vision components.

Takeaway 7. Text-to-image transfer is harder than image-to-text transfer due to the fact that the image tokenizer and fusion architecture do not have chance to align with each other while being trained on the text-only version of the datasets.

Dataset	Architecture	Text-only/test	Text-only/test-ood	VQA/test	VQA/test-ood
SpatialMap	Joint-Decoder	76.2 (+43.8)	76.6 (+44.2)	15.6 (-9.0)	18.3 (-6.8)
	Cross Attention	65.4 (+56.5)	66.4 (+58.0)	0.0 (-10.2)	0.0 (-11.6)
	MoT	72.5 (+46.8)	69.8 (+45.1)	26.6 (+3.3)	26.6 (+3.6)
Grid	Joint-Decoder	99.8 (+95.0)	58.2 (+56.3)	11.9 (+6.9)	5.8 (+3.4)
	Cross Attention	99.3 (+98.9)	42.8 (+42.4)	10.5 (+10.1)	5.5 (+5.5)
	MoT	100.0 (+93.5)	43.0 (+40.8)	13.1 (+12.3)	6.6 (+6.5)
Ring	Joint-Decoder	100.0 (+95.2)	15.1 (+11.3)	12.4 (+ 10.2)	6.4 (+ 5.3)
	Cross Attention	100.0 (+100.0)	12.0 (+12.0)	8.2 (+8.2)	3.9 (+3.9)
	MoT	100.0 (+95.3)	14.0 (+9.8)	8.4 (+7.5)	4.8 (+4.7)

Table 6: Text-to-image transfer performance (accuracy, %) by fusion architecture. Scores are averaged over text-supervised image tokenizers (CLIP, AIMv2, SigLIP 2) after models were fine-tuned on the text-only (QA) versions of the synthetic datasets. Each score is presented alongside the change in performance due to the additional training.

When analyzing the results by fusion architecture (Table 6), we see that the representational alignment between modalities remains a key factor. Transfer is marginally better for the Grid and Ring datasets, where the text is a structured serialization of the image, but remains very poor for SpatialMap. Notably, the Cross-Attention architecture appears most vulnerable to this transfer failure, with its performance often collapsing to near-zero on SpatialMap. This suggests its fusion mechanism is less robust to text-only training compared to the Joint-Decoder and MoT architectures.

Takeaway 8. The Cross-Attention architecture is particularly brittle in the text-to-image transfer setting, with its performance sometimes collapsing entirely. This suggests its fusion mechanism may be less robust to text-only fine-tuning compared to Joint-Decoder or MoT, ruining its VQA capability.

5 RELATED WORK

Our work builds on three key areas of VLM research: the design of image tokenizers, the diversity of training recipes, and direct architectural comparisons.

The Importance of the Image Tokenizer. The choice of image tokenizer is a critical factor in VLM performance. Prior work has explored this from several angles. For instance, PaLI (Chen et al., 2023) demonstrated the benefits of scaling up the vision encoder, while Eagle (Shi et al., 2025a) improved performance by combining multiple task-specific image tokenizers. More recent studies, such as Cambrian-1 (Tong et al., 2024), have focused on comparing different language-supervised tokenizers. While these works establish the tokenizer's importance, a systematic comparison of tokenizers trained with different objectives (e.g., contrastive vs. reconstructive) across varied architectural and training setups remains an open area. Our work addresses this by evaluating five distinct tokenizers across three fusion architectures.

Diversity in Training Recipes. Training recipes for VLMs are highly varied, typically involving multi-stage protocols with different layer-freezing strategies. A common approach, seen in the LLaVA family (Liu et al., 2023; 2024; Li et al., 2025), InternVL 2.5 (Chen et al., 2024), and DeepSeek-VL (Lu et al., 2024a), involves freezing both the image tokenizer and LLM during an initial alignment stage where only a small adapter is trained. In contrast, models like Qwen 2.5 VL (Bai et al., 2025) and DeepSeek-VL2 (Wu et al., 2024) unfreeze the image tokenizer from the start to foster better alignment, while keeping the LLM frozen. A less common end-to-end approach, adopted by Molmo (Deitke et al., 2024), trains all parameters simultaneously but requires high-quality data and carefully tuned learning rates. The choices become even more complex in later fine-tuning stages, as shown by Cambrian-1, which ablates freezing the image encoder during its second training phase and finds it can be beneficial. This diversity makes it difficult to attribute performance gains to specific training choices versus other confounding factors. Our work addresses this by systematically ablating freezing strategies across fixed architectures.

Architectural Comparisons. Direct architectural comparisons for VLMs have been conducted, but often with limited scope. For example, NVLM (Dai et al., 2024) provided an early comparison between the Joint-Decoder and Cross-Attention architectures, along with ablations on using special tokens for visual knowledge transfer. More recently, LMFusion (Shi et al., 2025b) compared the Joint-Decoder against using modality-specific MLPs or an MoT architecture, but their study was confined to a setting where the language backbone remained frozen. Our work expands on these studies by providing a unified comparison of all three major architectures (Joint-Decoder, Cross-Attention, and MoT) while also varying layer-freezing strategies for both the vision and language components, thereby offering a more comprehensive understanding of architectural trade-offs.

6 CONCLUDING REMARKS

In this work, we systematically studied how the choice of image tokenizer, fusion architecture, and layer freezing strategies influence the downstream performance of vision-language models. Our results highlight the critical impact of the image tokenizer choice, reveal distinct in-domain and out-of-domain trade-offs associated with layer freezing, quantify the robustness of different fusion architectures to various freezing strategies, and demonstrate the varying degrees of cross-domain transfer enabled by these design choices.

Several open questions arise from our findings. Our experiments were conducted with models in the 1–1.5B parameter range, and it is important to explore how these results scale to larger models. Future work could also extend the comparison to other modalities, such as speech, to investigate whether transfer occurs more readily between specific modality pairs or domains (e.g., speech and poetry vs. vision and physical reasoning). Additionally, our models are text-only generators. Further research is needed to assess the trade-offs of different architectures and training recipes for other generative tasks.

ETHICS STATEMENT

This work does not raise any ethical concerns.

REPRODUCIBILITY STATEMENT

To facilitate reproducibility and future research, we make the following resources available. Our full codebase, all trained model checkpoints, and our newly introduced synthetic datasets will be made publicly available upon publication. Detailed hyperparameters for all training stages are provided in Appendix D, and the generation procedure for our synthetic datasets is described in Appendix C. All academic benchmarks and the Qwen3-0.6B base model used in this work are already publicly available.

REFERENCES

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millicah, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088. 1, 2

Roman Bachmann, Jesse Allardice, David Mizrahi, Enrico Fini, Oğuzhan Fatih Kar, Elmira Amirloo, Alaaeldin El-Nouby, Amir Zamir, and Afshin Dehghan. Flextok: Resampling images into 1d token sequences of flexible length. *arXiv preprint arXiv:2502.13967*, 2025. 1

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv* preprint arXiv:2502.13923, 2025. 1, 2, 9

Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. https://github.com/kakaobrain/coyo-dataset, 2022. 3, 26

Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish V Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme Ruiz, Andreas Peter Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. PaLI: A jointly-scaled multilingual language-image model. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=mWVoBz4W0u.9

Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 9

Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. In *Forty-second International Conference on Machine Learning*, 2025. 23

Wenliang Dai, Nayeon Lee, Boxin Wang, Zhuolin Yang, Zihan Liu, Jon Barker, Tuomas Rintamaki, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. Nvlm: Open frontier-class multimodal llms. *arXiv preprint arXiv:2409.11402*, 2024. 1, 2, 5, 9

Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, Dirk Groeneveld, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. Molmo and PixMo: Open Weights and Open Data for Stateof-the-Art Vision-Language Models. arXiv e-prints, art. arXiv:2409.17146, September 2024. doi: 10.48550/arXiv.2409.17146. 1, 2, 9, 26

Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, Guang Shi, and Haoqi Fan. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025. 1, 2, 3

Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, Congcong Wang, Dehao Zhang, Dikang Du, Dongliang Wang, Enming Yuan, Enzhe Lu, Fang Li, Flood Sung, Guangda Wei, Guokun Lai, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haoning Wu, Haotian Yao, Haoyu Lu, Heng Wang, Hongcheng Gao, Huabin Zheng, Jiaming Li, Jianlin Su, Jianzhou Wang, Jiaqi Deng, Jiezhong Qiu, Jin Xie, Jinhong Wang, Jingyuan Liu, Junjie Yan, Kun Ouyang, Liang Chen, Lin Sui, Longhui Yu, Mengfan Dong, Mengnan Dong, Nuo Xu, Pengyu Cheng, Qizheng Gu, Runjie Zhou, Shaowei Liu, Sihan Cao, Tao Yu, Tianhui Song, Tongtong Bai, Wei Song, Weiran He, Weixiao Huang, Weixin Xu, Xiaokun Yuan, Xingcheng Yao, Xingzhe Wu, Xinxing Zu, Xinyu Zhou, Xinyuan Wang, Y. Charles, Yan Zhong, Yang Li, Yangyang Hu, Yanru Chen, Yejie Wang, Yibo Liu, Yibo Miao, Yidao Qin, Yimin Chen, Yiping Bao, Yiqin Wang, Yongsheng Kang, Yuanxin Liu, Yulun Du, Yuxin Wu, Yuzhi Wang, Yuzi Yan, Zaida Zhou, Zhaowei Li, Zhejun Jiang, Zheng Zhang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Zijia Zhao, and Ziwei Chen. Kimi-VL technical report, 2025. URL https://arxiv.org/abs/2504.07491.1,2

Enrico Fini, Mustafa Shukor, Xiujun Li, Philipp Dufter, Michal Klein, David Haldimann, Sai Aitharaju, Victor Guilherme Turrisi da Costa, Louis Béthune, Zhe Gan, Alexander T Toshev, Marcin Eichner, Moin Nabi, Yinfei Yang, Joshua M. Susskind, and Alaaeldin El-Nouby. Multimodal autoregressive pre-training of large vision encoders, 2024. 1, 3

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 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 4

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Koreney, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin,

595

596

597

598

600

601

602

603

604

605

606

607

608

610

611

612

613

614

615

616

617

618

619

620

621

622

623

625

626

627

629

630

631

632

633

634

635

636

637

638

639

640

641

642

644

645

646

Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon,

Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.1,2,3

- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. LLaVA-onevision: Easy visual task transfer. *Transactions on Machine Learning Research*, 2025. ISSN 2835-8856. URL https://openreview.net/forum?id=zKv8qULV6n. 9
- Weixin Liang, LILI YU, Liang Luo, Srini Iyer, Ning Dong, Chunting Zhou, Gargi Ghosh, Mike Lewis, Wen tau Yih, Luke Zettlemoyer, and Xi Victoria Lin. Mixture-of-transformers: A sparse and scalable architecture for multi-modal foundation models. *Transactions on Machine Learning Research*, 2025. ISSN 2835-8856. URL https://openreview.net/forum?id=Nu6N69i8SB. 1, 2, 3, 5
- Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. Vila: On pretraining for visual language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 26689–26699, June 2024. 5
- 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 *ECCV 2014*, pp. 740–755, 2014. 4
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 1, 9
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL https://llava-vl.github.io/blog/2024-01-30-llava-next/. 9
- Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models. *arXiv preprint arXiv:2505.24864*, 2025. 23
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. Deepseek-vl: Towards real-world vision-language understanding, 2024a. URL https://arxiv.org/abs/2403.05525.9
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *International Conference on Learning Representations* (ICLR), 2024b. 4
- Kevis-Kokitsi Maninis, Kaifeng Chen, Soham Ghosh, Arjun Karpur, Koert Chen, Ye Xia, Bingyi Cao, Daniel Salz, Guangxing Han, Jan Dlabal, Dan Gnanapragasam, Mojtaba Seyedhosseini, Howard Zhou, and André Araujo. TIPS: Text-Image Pretraining with Spatial Awareness. In *ICLR*, 2025. 3
- Ahmed Masry, Do Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 2263–2279, 2022. 4

- Minesh Mathew, Dimosthenis Karatzas, and C.V. Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 2200–2209, 2021. 4
 - Keita Miwa, Kento Sasaki, Hidehisa Arai, Tsubasa Takahashi, and Yu Yamaguchi. One-d-piece: Image tokenizer meets quality-controllable compression. In *Tokenization Workshop*, 2025. URL https://openreview.net/forum?id=lC4xkcLrdv. 1
 - Muhammad Ferjad Naeem, Yongqin Xian, Xiaohua Zhai, Lukas Hoyer, Luc Van Gool, and Federico Tombari. Silc: Improving vision language pretraining with self-distillation. In Aleš Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol (eds.), *Computer Vision ECCV 2024*, pp. 38–55, Cham, 2025. Springer Nature Switzerland. ISBN 978-3-031-72664-4. 3
 - Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023. 1
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language. *arXiv preprint arXiv:2103.00020*, 2021. 1, 3
 - Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge. In *ECCV* 2022, pp. 146–162, 2022. 4
 - Amrith Setlur, Nived Rajaraman, Sergey Levine, and Aviral Kumar. Scaling test-time compute without verification or rl is suboptimal. *arXiv preprint arXiv:2502.12118*, 2025. 23
 - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 27
 - Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv:* 2409.19256, 2024. 27
 - Min Shi, Fuxiao Liu, Shihao Wang, Shijia Liao, Subhashree Radhakrishnan, Yilin Zhao, De-An Huang, Hongxu Yin, Karan Sapra, Yaser Yacoob, Humphrey Shi, Bryan Catanzaro, Andrew Tao, Jan Kautz, Zhiding Yu, and Guilin Liu. Eagle: Exploring the design space for multimodal LLMs with mixture of encoders. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL https://openreview.net/forum?id=Y2RW9EVwhT. 9
 - Weijia Shi, Xiaochuang Han, Chunting Zhou, Weixin Liang, Xi Victoria Lin, Luke Zettlemoyer, and Lili Yu. Lmfusion: Adapting pretrained language models for multimodal generation, 2025b. URL https://arxiv.org/abs/2412.15188.1,2,3,5,9
 - Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8317–8326, 2019. 4
 - Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *Advances in neural information processing systems*, 37:84839–84865, 2024. 1, 3
 - Shengbang Tong, Ellis L Brown II, Penghao Wu, Sanghyun Woo, ADITHYA JAIRAM IYER, Sai Charitha Akula, Shusheng Yang, Jihan Yang, Manoj Middepogu, Ziteng Wang, Xichen Pan, Rob Fergus, Yann LeCun, and Saining Xie. Cambrian-1: A fully open, vision-centric exploration of multimodal LLMs. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=Vi8AepAXGy. 9

Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, Olivier Hénaff, Jeremiah Harmsen, Andreas Steiner, and Xiaohua Zhai. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. *arXiv preprint arXiv:2502.14786*, 2025. 1, 3

Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. Neural discrete representation learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/7a98af17e63a0ac09ce2e96d03992fbc-Paper.pdf. 1

Bo Wan, Michael Tschannen, Yongqin Xian, Filip Pavetic, Ibrahim Alabdulmohsin, Xiao Wang, André Susano Pinto, Andreas Steiner, Lucas Beyer, and Xiaohua Zhai. Locca: Visual pretraining with location-aware captioners. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 116355–116387. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/d303b4flef8d8274ae6b152df70f5406-Paper-Conference.pdf. 3

Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, Sharon Li, and Neel Joshi. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. *Advances in Neural Information Processing Systems*, 37:75392–75421, 2024. 5, 24

Fang Wu, Weihao Xuan, Ximing Lu, Zaid Harchaoui, and Yejin Choi. The invisible leash: Why rlvr may not escape its origin. *arXiv preprint arXiv:2507.14843*, 2025. 23

Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, Zhenda Xie, Yu Wu, Kai Hu, Jiawei Wang, Yaofeng Sun, Yukun Li, Yishi Piao, Kang Guan, Aixin Liu, Xin Xie, Yuxiang You, Kai Dong, Xingkai Yu, Haowei Zhang, Liang Zhao, Yisong Wang, and Chong Ruan. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding, 2024. URL https://arxiv.org/abs/2412.10302.9

xAI. Realworldqa, 2024. URL https://huggingface.co/datasets/xai-org/ RealworldQA. 4

Yutaro Yamada, Yihan Bao, Andrew Kyle Lampinen, Jungo Kasai, and Ilker Yildirim. Evaluating spatial understanding of large language models. *Transactions on Machine Learning Research*, 2024. 5, 24

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. arXiv preprint arXiv:2505.09388, 2025. 2

Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi Lin, Shuo Liu, Jiayi Lei, Quanfeng Lu, Runjian Chen, Peng Xu, Renrui Zhang, Haozhe Zhang, Peng Gao, Yali Wang, Yu Qiao, Ping Luo, Kaipeng Zhang, and Wenqi Shao. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi, 2024. 4

Qihang Yu, Mark Weber, Xueqing Deng, Xiaohui Shen, Daniel Cremers, and Liang-Chieh Chen. An image is worth 32 tokens for reconstruction and generation. *Advances in neural information processing systems*, 37:128940–128966, 2024. 1, 3

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 11975–11986, October 2023. 1, 3

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024. Association for Computational Linguistics. URL http://arxiv.org/abs/2403.13372. 27

Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 1, 2

A FULL EVALUATIONS

We include in this appendix the full list of evaluations for all of the models that we trained for stage 2 and stage 3.

A.1 STAGE 2 EVALUATION RESULTS

Tokenizer	$\frac{S1}{Im}$	$\frac{S}{Im}$	2 La	VQAv2 test-dev	A-OKVQA val	ChartQA test	TextVQA val	DocVQA test	Average In-domain	MathVista testmini	MMTBench val	RealWorldQA	Average Out-of-domain	Average All
CLIP	F F F U U	F F U U U	F U F U F U	49.8 66.7 61.8 72.1 66.0 74.1	22.1 41.6 31.4 45.4 35.7 47.2	12.2 20.8 20.4 29.2 23.0 30.6	20.1 31.9 27.9 35.9 34.1 40.1	12.1 20.7 16.5 24.8 20.4 25.2	23.2 36.3 31.6 41.5 35.8 43.4	25.0 27.3 26.7 27.7 22.8 28.4	37.4 34.3 34.7 35.0 38.5 35.4	29.0 37.0 34.5 35.9 31.4 36.9	30.5 32.9 32.0 32.9 30.9 33.6	25.9 35.0 31.7 38.3 34.0 39.7
AIMv2	F F F U U	F U U U U	F U F U F U	62.5 75.3 68.0 77.0 69.1 77.0	30.3 49.5 35.5 50.0 38.1 50.4	19.5 30.6 28.1 36.3 29.6 37.2	33.6 43.2 37.1 43.9 38.8 45.5	17.9 25.2 21.8 27.4 22.6 28.4	32.8 44.7 38.1 46.9 39.6 47.7	26.5 28.2 26.9 26.2 26.4 29.2	43.7 37.0 41.9 36.1 43.5 36.7	26.9 38.3 37.1 41.0 42.7 37.4	32.4 34.5 35.3 34.4 37.6 34.4	32.6 40.9 37.1 42.2 38.9 42.7
SigLIP 2	F F F U U	F F U U U	F U F U F U	55.9 74.8 65.4 77.8 71.1 78.4	24.0 47.2 31.4 50.2 40.9 52.3	15.7 28.8 29.3 43.2 36.5 46.8	30.0 43.6 35.5 44.7 45.1 50.0	17.0 26.6 22.2 29.6 28.1 32.7	28.5 44.2 36.8 49.1 44.3 52.0	25.5 28.0 24.1 25.1 27.4 28.6	38.1 38.1 37.7 35.7 42.1 35.6	31.4 32.4 32.7 33.2 34.1 29.5	31.7 32.8 31.5 31.3 34.5 31.3	29.7 39.9 34.8 42.4 40.7 44.2
TiTok	F F	F F	F U	3.2 43.1	0.2 26.0	5.1 13.9	1.0 11.8	3.0 11.8	2.5 21.3	23.8 26.0	27.7 22.8	28.6 32.5	26.7 27.1	11.6 23.5
VAR	F F	F F	F U	30.3 46.5	2.0 27.7	9.2 13.7	4.9 11.9	6.0 11.9	10.5 22.3	24.4 26.9	27.5 19.3	32.3 21.0	28.1 22.4	17.1 22.4

Table 7: Detailed evaluation results (accuracy, %) for the Joint-Decoder models. The table presents scores on individual benchmarks alongside in-domain, out-of-domain, and overall averages. We ablate the choice of image tokenizer and whether the image tokenizer (Im) or language layers (La) are frozen (F) or unfrozen (U) over the two stages (S1 and S2).

Tokenizer	$\frac{S1}{Im}$	S Im	La	VQAv2 test-dev	A-OKVQA val	ChartQA test	TextVQA val	DocVQA test	Average In-domain	MathVista testmini	MMTBench val	RealWorldQA	Average Out-of-domain	Average All
CLIP	F F F U U	F F U U U	F U F U F U	64.4 72.5 67.3 74.8 69.7 74.5	34.5 46.0 37.2 45.8 40.4 47.5	19.7 24.0 27.4 33.4 30.6 33.4	33.6 37.1 33.9 39.7 39.5 42.7	19.7 23.0 22.0 25.9 24.2 26.9	34.4 40.5 37.6 43.9 40.9 45.0	24.7 24.7 25.0 25.8 26.5 27.5	28.4 19.3 31.1 23.4 28.1 20.6	36.6 30.3 32.9 37.3 33.7 30.7	29.9 24.8 29.7 28.8 29.5 26.3	32.7 34.6 34.6 38.3 36.6 38.0
AIMv2	F F F U U	F U U U U	F U F U F U	69.4 75.4 71.1 76.7 72.0 76.5	36.7 47.0 40.6 49.8 40.5 49.8	24.8 30.7 30.6 35.9 30.8 36.1	40.1 43.6 40.0 43.5 41.8 45.3	21.1 25.3 22.2 25.8 23.6 26.9	38.4 44.4 40.9 46.4 41.7 46.9	24.0 26.9 24.4 27.2 25.0 26.1	26.9 29.7 26.7 28.9 32.5 20.4	28.9 25.0 34.4 28.1 34.8 25.0	26.6 27.2 28.5 28.1 30.7 23.8	34.0 37.9 36.3 39.5 37.6 38.3
SigLIP 2	F F F U U	F F U U U	F U F U F U	70.4 76.5 72.8 78.1 73.5 78.2	37.8 49.1 42.9 50.6 43.2 49.8	23.3 27.2 38.8 43.0 39.2 44.2	44.8 47.4 45.0 48.1 47.4 50.7	25.2 27.7 28.7 30.7 29.4 32.6	40.3 45.6 45.7 50.1 46.5 51.1	24.5 26.2 26.1 25.3 25.9 29.2	33.0 20.8 24.9 17.2 36.1 30.3	30.6 27.3 22.6 19.6 39.1 33.5	29.4 24.8 24.6 20.7 33.7 31.0	36.2 37.8 37.7 39.1 41.7 43.6
TiTok	F F	F F	F U F	40.0 42.5 42.4	13.9 24.6 16.3	10.6 13.6 10.6	7.7 11.1 8.5	8.2 11.8 9.1	16.1 20.7 17.4	24.7 23.8 23.7	25.0 21.3 24.4	37.5 25.8 36.9	29.1 23.6 28.3	20.9 21.8 21.5
VAR	F	F	U	46.6	24.6	13.2	11.6	11.9	21.6	26.3	22.5	24.8	24.6	22.7

Table 8: Detailed evaluation results (accuracy, %) for the Cross-Attention models. The table presents scores on individual benchmarks alongside in-domain, out-of-domain, and overall averages. We ablate the choice of image tokenizer and whether the image tokenizer (Im) or language layers (La) are frozen (F) or unfrozen (U) over the two stages (S1 and S2).

Tokenizer	S1 Im	S Im	2 La	VQAv2 test-dev	A-OKVQA val	ChartQA test	TextVQA val	DocVQA test	Average In-domain	MathVista testmini	MMTBench val	RealWorldQA	Average Out-of-domain	Average All
CLIP	F F F U U	F F U U U	F U F U F U	63.1 67.0 69.5 72.6 72.1 74.2	35.7 42.8 40.1 46.7 46.4 47.4	20.3 22.3 26.6 29.8 29.1 30.4	30.5 33.9 34.0 36.7 37.1 39.5	18.9 21.5 23.5 24.3 25.1 25.7	33.7 37.5 38.7 42.0 42.0 43.4	27.2 28.4 27.5 27.0 26.1 28.2	38.2 31.0 40.6 31.4 41.0 34.6	33.6 22.1 40.8 29.0 43.4 33.9	33.0 27.1 36.3 29.1 36.8 32.2	33.4 33.6 37.8 37.2 40.0 39.2
AIMv2	F F F U U	F F U U U	F U F U F U	72.8 75.4 74.7 77.3 74.8 77.1	42.5 48.0 44.9 50.3 46.6 51.1	28.5 30.7 35.2 37.0 35.1 36.9	40.4 43.4 41.4 44.1 42.5 44.3	23.8 26.4 26.8 28.0 27.6 28.2	41.6 44.8 44.6 47.3 45.3 47.5	27.8 27.3 29.0 27.1 26.4 26.4	40.4 40.8 41.9 31.0 39.0 32.2	36.9 42.2 42.4 28.0 39.0 32.9	35.0 36.8 37.8 28.7 34.8 30.5	39.1 41.8 42.0 40.3 41.4 41.2
SigLIP 2	F F F U U	F F U U U	F U F U F U	71.3 75.0 74.7 78.0 76.4 79.0	42.9 48.1 43.9 49.6 47.9 52.7	26.6 28.8 40.6 43.6 43.9 46.6	41.2 43.7 43.0 45.4 49.1 50.7	23.9 27.4 28.8 30.4 33.7 33.1	41.2 44.6 46.2 49.4 50.2 52.4	25.8 27.5 26.7 28.2 28.4 29.4	38.6 34.3 39.3 34.6 41.6 36.2	37.3 29.5 31.9 32.0 42.0 30.7	33.9 30.4 32.6 31.6 37.3 32.1	38.4 39.3 41.1 42.7 45.4 44.8
TiTok	F F	F F	F U F	39.7 42.3 41.9	21.0 26.9 23.8	10.5 12.9 10.8	7.7 10.5 8.7	7.9 11.2 9.6	17.4 20.8 19.0	23.6 24.9 25.2	27.8 20.4 27.1	30.3 19.3 32.0	27.2 21.5 28.1	21.1 21.1 22.4
VAR	F	F	U	46.9	29.0	13.4	11.7	12.4	22.7	25.0	22.4	24.7	24.0	23.2

Table 9: Detailed evaluation results (accuracy, %) for the Mixture-of-Transformers models. The table presents scores on individual benchmarks alongside in-domain, out-of-domain, and overall averages. We ablate the choice of image tokenizer and whether the image tokenizer (\mathbf{Im}) or language layers (\mathbf{La}) are frozen (\mathbf{F}) or unfrozen (\mathbf{U}) over the two stages ($\mathbf{S1}$ and $\mathbf{S2}$).

A.2 CROSS-MODALITY TRANSFER LEARNING

			VQA	/InD	VQA	OOD/	Text-o	nly/InD	Text-or	nly/OOD
Dataset	Architecture	Tokenizers	Base	SFT	Base	SFT	Base	SFT	Base	SFT
		CLIP	24.4	52.0	24.0	47.8	22.4	28.2	22.2	30.6
		AIMv2	24.8	88.4	25.6	85.6	36.4	37.2	35.0	37.0
	Joint-Decoder	SigLIP 2	24.8	97.8	25.8	93.2	38.4	39.2	40.0	36.2
		TiTok	1.4	24.6	1.4	26.6	13.0	8.4	13.0	7.8
		VAR	0.0	25.8	0.2	26.4	29.8	26.0	25.4	25.2
		CLIP	2.6	88.6	3.0	83.4	0.2	30.8	0.0	25.6
SpatialMap		AIMv2	21.2	88.8	23.0	84.8	6.0	29.6	4.2	28.6
эрананчар	Cross Attention	SigLIP 2	6.8	88.0	9.0	83.4	20.6	37.8	21.2	36.4
		TiTok	25.2	28.6	25.8	24.4	26.8	35.0	27.8	39.2
		VAR	16.8	22.2	17.2	28.0	31.6	28.0	27.8	28.2
		CLIP	22.0	32.2	22.2	33.2	28.8	33.2	29.0	29.4
		AIMv2	22.4	91.6	21.4	87.8	19.4	32.6	13.8	34.2
	MoT	SigLIP 2	25.6	97.2	25.6	93.8	29.0	32.6	31.4	32.2
		TiTok	6.6	23.6	5.2	24.6	29.2	26.4	28.0	24.0
		VAR	5.2	24.6	9.2	25.8	26.2	24.8	28.4	26.4
		CLIP	5.0	96.8	3.0	9.0	0.4	100.0	3.4	14.6
		AIMv2	0.8	99.8	0.0	7.4	5.6	99.8	3.4	14.8
	Joint-Decoder	SigLIP 2	0.8	99.8	0.4	6.4	8.6	97.8	4.8	9.0
		TiTok	0.6	20.2	0.0	8.6	9.0	13.6	0.0	7.2
		VAR	0.0	24.2	0.0	8.4	3.0	44.0	3.8	11.4
		CLIP	0.0	99.4	0.0	10.4	0.0	99.4	0.0	12.6
Ring		AIMv2	0.0	99.0	0.0	8.4	0.0	98.4	0.0	13.4
King	Cross Attention	SigLIP 2	0.0	99.2	0.0	10.0	0.0	68.6	0.2	13.2
		TiTok	0.0	24.4	0.0	10.0	0.2	14.6	0.0	9.8
		VAR	0.0	24.4	0.0	10.0	0.8	14.2	0.0	8.2
		CLIP	0.2	94.8	0.0	11.2	7.4	99.6	5.2	18.4
		AIMv2	2.6	99.6	0.4	10.2	6.6	95.6	7.4	12.4
	MoT	SigLIP 2	0.0	99.4	0.0	10.4	0.2	80.4	0.0	8.6
		TiTok	0.8	23.2	0.0	9.4	0.6	87.4	1.0	10.4
		VAR	0.0	26.8	0.0	8.0	0.0	54.0	0.0	11.2
		CLIP	4.0	87.8	1.2	18.0	0.0	58.0	0.0	21.6
		AIMv2	9.6	99.2	4.2	13.8	8.0	70.8	2.4	27.4
	Joint-Decoder	SigLIP 2	1.6	99.2	1.8	18.6	6.4	73.2	3.4	17.4
		TiTok	0.4	30.6	1.0	13.8	1.4	5.0	1.0	2.8
		VAR	0.4	33.2	0.0	8.0	1.2	37.0	1.2	10.0
		CLIP	1.4	93.8	0.0	18.4	0.0	65.6	0.0	19.2
Grid		AIMv2	0.0	91.6	0.0	13.8	0.0	15.2	0.0	8.2
Gria	Cross Attention	SigLIP 2	0.0	98.6	0.0	13.2	1.2	22.4	1.4	15.0
		TiTok	0.0	30.0	0.0	11.4	4.0	16.0	1.4	8.8
		VAR	0.0	32.2	0.0	7.8	0.0	17.2	0.0	8.2
		CLIP	0.6	94.6	0.0	15.4	10.0	69.4	2.4	17.6
		AIMv2	2.0	95.4	0.2	19.0	9.6	59.2	4.2	20.8
	MoT	SigLIP 2	0.0	98.2	0.2	20.4	0.0	72.2	0.2	17.0
		TiTok	0.2	30.8	0.2	11.4	0.0	48.4	0.0	12.4
		VAR	0.0	30.8	0.0	10.2	0.0	40.0	0.4	10.8

Table 10: Image-to-text transfer performance (accuracy, %) by different combination of fusion architecture and image tokenizers. Scores are evaluated after models were fine-tuned on the VQA versions of the synthetic datasets. Each score is presented alongside the original score before the additional training.

			Text-o	nly/InD	Text-or	nly/OOD	VQA	/InD	VQA	OOD
Dataset	Architecture	Tokenizers	Base	SFT	Base	SFT	Base	SFT	Base	SFT
		CLIP	22.4	74.8	22.2	74.4	24.4	25.4	24.0	27.8
		AIMv2	36.4	78.8	35.0	78.0	24.8	21.0	25.6	26.6
	Joint-Decoder	SigLIP 2	38.4	75.2	40.0	77.4	24.8	0.4	25.8	0.6
		TiTok	13.0	56.8	13.0	61.0	1.4	0.0	1.4	0.0
		VAR	29.8	71.6	25.4	77.6	0.0	25.4	0.2	22.8
		CLIP	0.2	76.6	0.0	72.8	2.6	0.0	3.0	0.0
SpatialMap		AIMv2	6.0	62.6	4.2	65.8	21.2	0.0	23.0	0.0
эринин тир	Cross Attention	SigLIP 2	20.6	57.2	21.2	60.8	6.8	0.0	9.0	0.0
		TiTok	26.8	64.8	27.8	64.6	25.2	0.0	25.8	0.4
		VAR	31.6	71.8	27.8	75.0	16.8	2.6	17.2	3.6
		CLIP	28.8	78.2	29.0	77.6	22.0	26.8	22.2	25.6
) (T	AIMv2	19.4	69.4	13.8	62.6	22.4	29.0	21.4	25.2
	MoT	SigLIP 2	29.0	70.0	31.4	69.2	25.6	24.2	25.6	29.2
		TiTok	29.2	62.2	28.0	59.8	6.6	27.2	5.2 9.2	27.0
		VAR	26.2	70.4	28.4	72.0	5.2	24.4	9.2	22.8
		CLIP	0.4	100.0	3.4	15.2	5.0	15.4	3.0	5.8
	T 1 . B 1	AIMv2	5.6	100.0	3.4	17.2	0.8	10.8	0.0	5.6
	Joint-Decoder	SigLIP 2	8.6	100.0	4.8	13.0	0.8	11.0	0.4	8.0
		TiTok	9.0	100.0	0.0	14.8	0.6	11.2	0.0	10.4
		VAR	3.0	100.0	3.8	17.2	0.0	9.6	0.0	6.8
		CLIP	0.0	100.0	0.0	9.4	0.0	10.8	0.0	5.2
Ring	a A	AIMv2	0.0	100.0	0.0	12.2	0.0	12.0	0.0	6.4
8	Cross Attention	SigLIP 2	0.0	100.0	0.2	14.4	0.0	1.8	0.0	0.2
		TiTok VAR	0.2	100.0 100.0	$0.0 \\ 0.0$	13.6 13.6	$0.0 \\ 0.0$	5.8 6.6	$0.0 \\ 0.0$	3.6 3.6
	-									
		CLIP AIMv2	7.4 6.6	100.0 100.0	5.2 7.4	15.6 11.6	0.2 2.6	11.2	0.0 0.4	4.8 4.0
	МоТ	SigLIP 2	0.0	100.0	0.0	15.0	0.0	6.6 7.6	0.4	5.8
	IVIOI	TiTok	0.6	100.0	1.0	14.6	0.8	7.0	0.0	3.0
		VAR	0.0	100.0	0.0	11.0	0.0	11.0	0.0	8.2
		CLIP	0.0	99.8	0.0	70.6	4.0	11.6	1.2	4.0
		AIMv2	8.0	99.8 99.8	2.4	70.6 29.4	4.0 9.6	11.0	4.2	7.6
	Joint-Decoder	SigLIP 2	6.4	99.8	3.4	74.8	1.6	13.0	1.8	6.0
	John-Decoder	TiTok	1.4	100.0	1.0	36.2	0.4	14.0	1.0	8.8
		VAR	1.2	97.8	1.2	39.6	0.4	11.8	0.0	4.6
		CLIP	0.0	99.8	0.0	43.0	1.4	10.2	0.0	6.4
		AIMv2	0.0	98.2	0.0	31.6	0.0	9.6	0.0	3.8
Grid	Cross Attention	SigLIP 2	1.2	100.0	1.4	54.0	0.0	11.8	0.0	6.4
		TiTok	4.0	100.0	1.4	45.4	0.0	8.6	0.0	7.2
		VAR	0.0	99.8	0.0	25.2	0.0	8.0	0.0	3.6
		CLIP	10.0	100.0	2.4	55.6	0.6	9.0	0.0	5.0
		AIMv2	9.6	100.0	4.2	26.4	2.0	18.6	0.2	10.0
	MoT	SigLIP 2	0.0	100.0	0.2	47.2	0.0	11.8	0.2	5.0
	•	TiTok	0.0	99.8	0.0	25.4	0.2	12.0	0.2	6.0
		VAR	0.0	99.8	0.4	35.2	0.0	15.2	0.0	6.2

Table 11: Text-to-image transfer performance (accuracy, %) by different combination of fusion architecture and image tokenizers. Scores are evaluated after models were fine-tuned on the text-only (QA) version of the synthetic datasets. Each score is presented alongside the original score before the additional training.

B ADDITIONAL RESULTS: CROSS-MODALITY TRANSFER FOR OPEN-WEIGHT VLMS

To supplement our analysis on cross-modality and out-of-distribution transfer on our trained vision-language models, in this section we evaluate a suite of open-source vision—language models spanning different architectures and scales, including Qwen2.5-VL (3B- and 7B-Instruct), Gemma, InternVL3, and Kimi-VL. Each model is fine-tuned on our synthetic datasets using supervised fine-tuning (SFT) for direct comparison across modalities. For Qwen2.5-VL, we additionally explore

reinforcement learning (RL)-based finetuning, reported in Section B.2, to assess whether optimization beyond standard SFT can further enhance cross-modality transfer. Hyperparameters and other details about training setup are provided in Appendix D.

B.1 SUPERVISED FINE-TUNING TRANSFER RESULTS

Similar to Section 4.1 and Section 4.2, we perform a single-epoch supervised fine-tuning (SFT) on either the image-based VQA task or the equivalent text-only task across our three synthetic datasets to evaluate cross-modality transfer in open-weight vision—language models. Training is conducted on the respective train split, and models are evaluated on four held-out settings: (i) test split in the same modality, (ii) test split in the opposite modality, (iii) out-of-distribution (OOD) test split in the same modality, and (iv) OOD test split in the opposite modality. This setup parallels the fine-tuning procedure used for our in-house multimodal models, enabling direct comparison of how open-source architectures generalize across modalities and dataset variants.

We present evaluation accuracy after fine-tuning on the image version of the datasets in Table 12, We observe consistent image-to-text transfer when models are fine-tuned on the image version of the datasets. Qwen2.5-VL-7B achieves the strongest overall performance, with InternVL3 also showing competitive transfer, particularly on Grid. As in our earlier experiments, SpatialMap remains the most challenging for cross-modal transfer: several models exhibit drops in text-only accuracy after SFT on the image task, underscoring the modality mismatch in how the task is represented.

Conversely, in the text-to-image transfer setting, results are more mixed in Table 13. Qwen2.5-VL-7B shows strong transfer on SpatialMap but weaker performance on Grid and Ring. InternVL3 demonstrates more balanced transfer across datasets, suggesting that pretraining with substantial text-only data may aid robustness when moving from textual to visual reasoning. Overall, these results reinforce the asymmetric nature of cross-modality transfer and highlight model-specific differences in how supervision in one modality propagates to another.

Dataset	Model	VQA/InD	VQA/OOD	Text-only/InD	Text-only/OOD
	Qwen2.5-VL-7B	98.6 (+18.4)	98.2 (+18.0)	63.0 (+0.8)	57.6 (-5.8)
	Qwen2.5-VL-3B	84.4 (+38.8)	81.0 (+38.6)	42.2 (-1.4)	43.0 (-1.0)
SpatialMap	InternVL3-8B	97.2 (+19.4)	95.2 (+19.4)	62.0 (-6.2)	59.2 (-9.4)
	Gemma-3-4B	94.0 (+53.4)	91.4 (+50.0)	52.8 (-1.6)	53.8 (+1.8)
	Kimi-VL-A3B	28.8 (+4.0)	24.8 (-3.6)	34.2 (-24.2)	28.6 (-33.8)
	Qwen2.5-VL-7B	99.4 (+83.2)	84.8 (+77.0)	99.4 (+47.0)	95.6 (+68.2)
	Qwen2.5-VL-3B	91.0 (+79.0)	34.4 (+28.6)	94.4 (+82.0)	60.0 (+53.2)
Grid	InternVL3-8B	99.4 (+80.2)	68.4 (+53.6)	97.8 (+11.6)	87.4 (+18.4)
	Gemma-3-4B	99.0 (+83.6)	59.8 (+56.0)	98.2 (+66.2)	84.4 (+71.6)
	Kimi-VL-A3B	7.0 (-4.4)	2.6 (-4.8)	60.6 (-2.4)	44.8 (-5.6)
	Qwen2.5-VL-7B	99.4 (+83.8)	15.0 (+4.6)	99.6 (+69.4)	15.8 (-4.0)
	Qwen2.5-VL-3B	85.0 (+74.4)	13.8 (+5.4)	82.0 (+69.8)	18.2 (+10.8)
Ring	InternVL3-8B	99.4 (+84.4)	17.8 (+9.4)	95.4 (+41.2)	28.8 (-40.6)
	Gemma-3-4B	99.8 (+89.6)	15.2 (+12.8)	99.4 (+34.6)	15.2 (-15.8)
	Kimi-VL-A3B	7.0 (-3.0)	3.0 (-6.6)	63.4 (+33.0)	34.0 (-0.4)

Table 12: Image-to-text transfer performance (accuracy, %) by base model. Models were fine-tuned on the VQA version of the synthetic datasets. Each score is presented alongside its performance delta from the base model's performance.

B.2 Comparison to Reinforcement Learning Fine-tuning

We compare supervised fine-tuning (SFT) and reinforcement learning (RL) fine-tuning for Qwen-2.5-VL models at two scales (3B and 7B parameters). In both cases, we train exclusively on the image version of each synthetic dataset (SpatialMap, Ring, Grid), using a binary reward signal of 1 for producing the correct final answer and 0 otherwise. We then evaluate on both the indistribution test and OOD test splits, across both image and text modalities. This setup probes not only in-distribution performance, but also cross-modality transfer (image \rightarrow text) and robustness

Dataset	Model	Text-only/InD	Text-only/OOD	VQA/InD	VQA/OOD
	Qwen2.5-VL-7B	86.2 (+24.0)	90.0 (+26.6)	89.0 (+8.8)	86.2 (+6.0)
	Qwen2.5-VL-3B	61.4 (+17.8)	55.0 (+11.0)	27.0 (-18.6)	27.2 (-15.2)
SpatialMap	InternVL3-8B	91.4 (+23.2)	93.6 (+25.0)	78.6 (+0.8)	79.2 (+3.4)
	Gemma-3-4B	90.2 (+35.8)	90.4 (+88.2)	52.6 (+12.0)	48.6 (+7.2)
	Kimi-VL-A3B	70.2 (+11.8)	73.2 (+10.8)	22.8 (-2.0)	25.4 (-3.0)
	Qwen2.5-VL-7B	99.8 (+47.4)	95.4 (+68.0)	19.2 (+3.0)	8.2 (+0.4)
	Qwen2.5-VL-3B	100.0 (+87.6)	49.2 (+42.4)	16.8 (+4.8)	9.8 (+4.0)
Grid	InternVL3-8B	100.0 (+13.8)	98.4 (+29.4)	16.0 (-3.2)	4.8 (-10.0)
	Gemma-3-4B	99.8 (+67.8)	54.4 (+41.6)	16.4 (+1.0)	11.2 (+7.4)
	Kimi-VL-A3B	90.8 (+27.8)	41.4 (-9.0)	60.6 (+49.2)	4.0 (-3.4)
	Qwen2.5-VL-7B	100.0 (+69.8)	16.4 (-3.4)	18.6 (+3.0)	7.0 (-3.4)
	Qwen2.5-VL-3B	93.2 (+81.0)	23.0 (+15.6)	21.4 (+10.8)	8.0 (-0.4)
Ring	InternVL3-8B	99.8 (+45.6)	23.6 (-45.8)	19.8 (+4.8)	9.8 (+1.4)
	Gemma-3-4B	99.2 (+34.4)	16.2 (-14.8)	14.6 (+4.4)	9.2 (+6.8)
	Kimi-VL-A3B	94.0 (+63.6)	41.8 (+7.4)	10.4 (+0.4)	9.0 (-0.6)

Table 13: Text-to-image transfer performance (accuracy, %) by base model. Models were fine-tuned on the text-only (QA) version of the synthetic datasets. Each score is presented alongside its performance delta from the base model's performance.

to distribution shift. We report results at base, after SFT, and at the best RL checkpoint for each configuration.

At the larger 7B scale (see Table 14), RL generally provides stronger generalization than SFT, aligning with similar sentiments from prior work (Chu et al., 2025; Liu et al., 2025; Setlur et al., 2025). For example, RL improves cross-modality transfer on SpatialMap, where SFT hurts image to text performance, and yields large gains on OOD splits for Ring (0.87–0.88 vs 0.14–0.20 for SFT). RL also preserves or slightly enhances in-distribution accuracy, often reaching near-perfect performance. Although there were instances where SFT frequently saturates or even degrades transfer, it is more sample-efficient; we see for certain tasks (eg. Grid), the gains from RL are more modest particularly with respect to cross-modality transfer. However, we observed RL training continued to yield steady improvements for the entire duration of training (15 epochs)— in contrast, SFT offers better sample efficiency.

For the smaller 3B scale (see Table 15), transfer patterns are more varied. RL improves cross-modality and OOD transfer on SpatialMap, but offers less consistent gains on Grid, where SFT remains competitive—likely reflecting that the Grid and Ring tasks have more closely aligned image and text representations. A notable pitfall emerges in Ring: the 3B RL model fails to improve much beyond random chance, collapsing into short outputs with only the final answer token. This illustrates the importance of controlling model output distribution during RL, and the inherent limitations of RLVR to the support of the base model Wu et al. (2025). As mentioned above, when such collapse does not occur, we observed RL training yields steady improvements over epochs, whereas SFT tends to plateau earlier.

Overall, these results indicate that RL at larger scale consistently enhances generalization across modality and distribution shift, while at smaller scale it can either unlock improved transfer (SpatialMap) or suffer from instability (Ring). SFT remains a useful baseline for more aligned tasks but appears less reliable for tasks requiring substantial abstraction across modalities.

	SpatialMap		Grid			Ring			
	Base	SFT	RL	Base	SFT	RL	Base	SFT	RL
VQA/InD	80.2	99.6 (+19.4)	99.8 (+19.6)	16.2	99.8 (+83.6)	98.6 (+82.4)	15.6	100.0 (+84.4)	99.8 (+84.2)
VQA/OOD	80.2	99.6 (+19.4)	99.8 (+19.6)	7.8	89.0 (+81.2)	59.6 (+51.8)	10.4	15.0 (+4.6)	86.6 (+76.2)
Text-only/InD	62.2	65.2 (+3.0)	70.8 (+8.6)	52.4	99.8 (+47.4)	61.8 (+9.4)	30.2	100.0 (+69.8)	98.6 (+68.4)
Text-only/OOD	63.4	63.8 (+0.4)	71.6 (+8.2)	27.4	97.4 (+70.0)	30.4 (+3.0)	19.8	19.8 (+0.0)	87.8 (+68.0)

Table 14: Transfer performance (accuracy, %) for Qwen-2.5-VL-7B-Instruct. Synthetic SpatialMap, Grid, and Ring results at base, after SFT, and best RL checkpoint for Qwen-2.5-VL-7B-Instruct when training with the image version of the respective task.

	SpatialMap		Grid			Ring			
	Base	SFT	RL	Base	SFT	RL	Base	SFT	RL
VQA/InD	45.6	99.2 (+53.6)	100.0 (+54.4)	12.0	99.2 (+87.2)	89.6 (+77.6)	10.6	100.0 (+89.4)	13.2 (+2.6)
VQA/OOD	42.4	99.4 (+57.0)	99.4 (+57.0)	5.8	45.4 (+39.6)	44.8 (+39.0)	8.4	17.4 (+9.0)	12.4 (+4.0)
Text-only/InD	43.6	46.4 (+2.8)	64.4 (+20.8)	12.4	99.2 (+86.8)	42.0 (+29.6)	12.2	99.6 (+87.4)	12.8 (+0.6)
Text-only/OOD	44.0	44.8 (+0.8)	66.0 (+22.0)	6.8	60.0 (+53.2)	20.4 (+13.6)	7.4	21.4 (+14.0)	9.6 (+2.2)

Table 15: Transfer performance (accuracy, %) for Qwen-2.5-VL-3B-Instruct. Synthetic SpatialMap, Grid, and Ring results at base, after SFT, and best RL checkpoint for Qwen-2.5-VL-3B-Instruct when training with the image version of the respective task.

C SYNTHETIC DATASET DETAILS

Below we provide details regarding the three synthetic datasets.

SpatialMap: The Synthetic SpatialMap dataset tests models on spatial reasoning over symbolic objects. Each sample consists of a set of n colored shapes placed randomly on a blank canvas. Questions probe pairwise spatial relations, for example: "Where is the blue circle relative to the red star? Possible options: Northeast, Northwest, Southeast, Southwest." The dataset follows the structure of the spatial mapping tasks in prior work (Wang et al., 2024), with object names replaced by colored geometric shapes for greater visual simplicity and compositional control. We provide two parallel modalities: an image version, where the model must infer relations directly from the visual configuration, and a text version, which encodes equivalent information as a series of binary relation statements (e.g., "The purple pentagon is to the Northeast of the blue circle."). The text description is complete and sufficient to solve the task, but requires chaining relational statements. The training split contains 4,500 procedurally generated examples with 7 objects per image, while the OOD split introduces 8 objects per image, slightly increasing task complexity without changing the basic query format.

Grid: The Synthetic Grid dataset evaluates navigation and reasoning in discrete two-dimensional environments. We generate an $n \times n$ grid where each cell contains a distinct colored shape. Questions specify a starting shape and a sequence of relative moves, such as "Begin at the yellow triangle and go down one step, then right two steps. Which shape do you land on?" The image version provides the grid visualization alongside the question, while the text version lists the grid contents in row-major order, thereby encoding object positions without requiring visual perception. This dataset is adapted from the "global grid" task of Yamada et al. (2024), though we substitute ImageNet categories with geometric shapes to simplify object recognition and emphasize spatial reasoning. The training set uses 3×3 grids with navigation sequences of length 8, while the OOD split increases difficulty with 4×4 grids and up to 12-step sequences. Unlike the image version, the text version bypasses object recognition, highlighting how modality differences impact reasoning difficulty.

Ring: The Synthetic Ring dataset parallels the grid task but in a circular layout, also present in its text version in Yamada et al. (2024). Objects are evenly arranged around a ring, and queries specify a starting shape and a number of clockwise or counterclockwise steps (e.g., "Starting from the red square, move four steps clockwise. Which shape do you land on?"). The image version presents the circle of shapes with the question, while the text version linearizes the ring into an ordered list beginning from a designated reference point and continuing clockwise. The training split contains rings of 8 objects with navigation sequences up to 8 steps, while the OOD split expands to 12 objects and 12-step sequences.

 Dataset generation: We describe how we procedurally generate the three datasets below. We plan on releasing the dataset generation code and example datasets.

- For each **SpatialMap** example, we sample a set of distinct color-shape pairs and randomly assign them positions along two independent orderings (vertical and horizontal). These orderings determine the relative row and column of each object, which are rendered on a blank canvas. We then select an unambiguous query pair (q,r) whose vertical and horizontal offsets uniquely determine a diagonal relation (Northeast, Northwest, Southeast, or Southwest). The **text modality** question (question) begins with a full relational description of the scene expressed as binary statements (e.g., "The purple pentagon is to the Northeast of the blue circle."). The query is appended to this description, requiring the model to chain multiple statements to answer correctly. The **image modality** question (question_direct) instead lists only the set of objects without relational information, with the same query appended, such that solving requires interpreting the image directly. Solutions are generated in parallel: the text solution (solution) provides a multi-step chain-of-thought reasoning through vertical and horizontal relations, while the image solution (solution_direct) gives a concise explanation phrased as direct visual inspection. The final label (answer) is the correct diagonal relation.
- For each **Grid** example, we construct an $n \times n$ grid and assign a unique color-shape object to each cell, rendering the grid with light boundaries. A navigation query is created by sampling a start cell and a valid sequence of directional moves. The **text modality** question (question) specifies the grid contents in row-major order, ensuring that the layout can be reconstructed entirely from text, followed by the navigation program (e.g., "Start at the yellow triangle, move down one step, then right two steps. Which object do you land on?"). The **image modality** question (question_direct) instead lists the objects in random order without positional information, requiring the model to resolve the navigation over the visual grid. The solutions mirror these formats: the text solution (solution) details each step of the navigation trace, while the image solution (solution_direct) is identical, describing the sequence of moves until the final object is reached. The final label (answer) is the object found at the destination cell.
- For each **Ring** example, we arrange a sampled set of unique color-shape objects evenly around a circle in clockwise order, selecting a starting position and a sequence of clockwise or counterclockwise steps to form a navigation query. The **text modality** question (question) lists the objects deterministically in clockwise order from a fixed reference point, then provides the navigation program (e.g., "Starting from the red square, move four steps clockwise. Which object do you land on?"). The **image modality** question (question_direct) instead lists the same objects in random order without positional information, such that the model must rely on the ring image to resolve the query. Both the text and image solutions (solution, solution_direct) provide an explicit step-by-step trace of the walk around the ring, concluding with the identified object. The final label (answer) is the object at the destination position.

D TRAINING DETAILS

D.1 STAGE 1 (PRETRAINING) HYPERPARAMETERS

For Stage 1 pretraining, models are trained on the COYO-700M dataset (Byeon et al., 2022) for 100,000 steps with a global batch size of 1536. Input text captions are truncated to a maximum length of 256 tokens. During this stage, only the Joint-Decoder and Cross-Attention models are trained directly. For stage 2, MoT checkpoints are initialized from the resulting Joint-Decoder weights, as mentioned in Section 2.3.

We follow the learning rate strategy of Deitke et al. (2024). The adapter modules for the Joint-Decoder and Cross-Attention models use a learning rate of 2×10^{-4} . When the image tokenizer is unfrozen, it is trained with a learning rate of 6×10^{-6} . All models use a cosine-decay learning rate schedule with a linear warmup, decaying to 10% of the maximum learning rate by the end of training. A comprehensive list of all hyperparameters is provided in Table 16.

	Joint-Decoder	Cross-Attention
Adapter LR	2×10^{-4}	2×10^{-4}
Cross-Attention LR	N/A	2×10^{-4}
Image Tokenizer LR	6×10^{-6}	6×10^{-6}
Optimizer	AdamW	AdamW
Betas	(0.9, 0.999)	(0.9, 0.999)
Weight decay	0.01	0.01
LR Schedule	Cosine Decay	Cosine Decay
Min LR	10% of Max	10% of Max
Linear Warmup	2000 steps	2000 steps
Global Batch size	1536	1536
Num Training Steps	100k steps	100k steps

Table 16: Hyperparameters for stage 1 training.

D.2 STAGE 2 (FINE-TUNING) HYPERPARAMETERS

For the Stage 2 fine-tuning, models are trained on a combined dataset comprising COCO-Captions, VQAv2, ChartQA, TextVQA, and DocVQA. Following the second-stage setup of Deitke et al. (2024), we set the learning rate to 1×10^{-5} for the image tokenizer (when unfrozen), the adapters in the Joint-Decoder and MoT models, and the cross-attention layers. A higher learning rate of 5×10^{-5} is used for the LLM backbone (when unfrozen) and the image modality-transformer in the MoT models. All models use a cosine-decay learning rate schedule with a linear warmup, decaying to 10% of the maximum learning rate. A complete summary of these hyperparameters is available in Table 17.

D.3 STAGE 3 (REASONING-TRANSFER) HYPERPARAMETERS

For the Stage 3 reasoning transfer experiments, models are fine-tuned on either the text-only (QA) or image-based (VQA) training split of one of the three synthetic datasets. A uniform learning rate is applied to all unfrozen layers. As in previous stages, discrete image tokenizers remain frozen to avoid the need for a straight-through estimator.

Although the training was configured for 5 epochs using a cosine-decay schedule, we exclusively use the checkpoint saved after the first epoch for all evaluations. Under this setup, the learning rate only decays to approximately 90.5% of its initial value by the end of the first epoch ($\frac{1+\cos(\pi/5)}{2}\approx 0.905$), effectively creating a near-constant learning rate. While we expect performance to be very similar to using a true constant learning rate, we provide these specifics for full reproducibility. All hyperparameters are detailed in Table 18.

	Joint-Decoder	Cross-Attention	МоТ
Adapter LR	1×10^{-5}	1×10^{-5}	1×10^{-5}
Cross-Attention LR	N/A	1×10^{-5}	N/A
Image Tokenizer LR	1×10^{-5}	1×10^{-5}	1×10^{-5}
Image Transformer LR	N/A	N/A	5×10^{-5}
Language LR	5×10^{-5}	5×10^{-5}	5×10^{-5}
Optimizer	AdamW	AdamW	AdamW
Betas	(0.9, 0.999)	(0.9, 0.999)	(0.9, 0.999)
Weight decay	0.01	0.01	0.01
LR Schedule	Cosine Decay	Cosine Decay	Cosine Decay
Min LR	10% of Max	10% of Max	10% of Max
Linear Warmup	750 steps	750 steps	750 steps
Global Batch size	1536	1536	1536
Num Training Steps	10860 steps (15 epochs)	10860 steps (15 epochs)	10860 steps (15 epochs)

Table 17: Hyperparameters for stage 2 training.

	Joint-Decoder	Cross-Attention	MoT
Learning Rate	1×10^{-5}	1×10^{-5}	1×10^{-5}
Optimizer	AdamW	AdamW	AdamW
Betas	(0.9, 0.999)	(0.9, 0.999)	(0.9, 0.999)
Weight decay	0.01	0.01	0.01
LR Schedule	Cosine Decay	Cosine Decay	Cosine Decay
Linear Warmup	0 steps	0 steps	0 steps
Global Batch size	32	32	32
Num Training Steps	140 steps (1 epoch)	140 steps (1 epoch)	140 steps (1 epoch)

Table 18: Hyperparameters for stage 3 training.

D.4 OPEN-WEIGHT VISION-LANGUAGE MODELS

For our fine-tuning results on open-weight vision-language models presented in Appendix B, we use the LLaMA-Factory framework (Zheng et al., 2024) for SFT and the verl (Sheng et al., 2024) implementation of Group Relative Policy Optimization (GRPO) (Shao et al., 2024) for RL fine-tuning. We specify hyperparameters in Table 19 below.

	SFT	RL
Learning Rate	1×10^{-6}	1×10^{-6}
Optimizer	AdamW	AdamW
Betas	(0.9, 0.999)	(0.9, 0.999)
Warmup Steps	0	0
Scheduler	Cosine Decay	Constant
Num Epochs	1	15
Training Global Batch Size	64	128
Rollout Global Batch Size	N/A	128
N Samples per Prompt	N/A	5
KL Coeff	N/A	1×10^{-3}

Table 19: Hyperparameters for SFT and RL training.

E LARGE LANGUAGE MODEL USAGE

We used large language models to help improve the clarity and style of the manuscript. All drafts, including the main text, citations, and tables, were originally written by hand. We then used ChatGPT-4.1 and Gemini 2.5 Pro to suggest revisions for clarity, conciseness, and academic tone. No content generation or data analysis was performed by language models; all substantive contributions and data interpretation are the authors' own.