

Monet: Reasoning in Latent Visual Space Beyond Image and Language

Qixun Wang¹ Yang Shi^{1,2} Yifei Wang³
 Yuanxing Zhang² Pengfei Wan² Kun Gai² Xianghua Ying^{1*} Yisen Wang^{1*}
¹Peking University ²Kling Team ³MIT

Abstract

*“Thinking with images” has emerged as an effective paradigm for advancing visual reasoning, extending beyond text-only chains of thought by injecting visual evidence into intermediate reasoning steps. However, existing methods fall short of human-like abstract visual thinking, as their flexibility is fundamentally limited by external tools. In this work, we introduce **Monet**, a training framework that enables multimodal large language models (MLLMs) to reason directly within the latent visual space by generating continuous embeddings that function as intermediate visual thoughts. We identify two core challenges in training MLLMs for latent visual reasoning—high computational cost in latent–vision alignment and insufficient supervision over latent embeddings, and address them with a three-stage distillation-based supervised fine-tuning (SFT) pipeline. We further reveal a limitation of applying GRPO to latent reasoning: it primarily enhances text-based reasoning rather than latent reasoning. To overcome this, we propose **VLPO** (**V**isual-**L**atent **P**olicy **O**ptimization), a reinforcement learning method that explicitly incorporates latent embeddings into policy gradient updates. To support SFT, we construct Monet-SFT-125K, a high-quality text–image interleaved CoT dataset containing 125K real-world, chart, OCR, and geometry CoTs. Our model, Monet-7B, shows consistent gains across real-world perception and reasoning benchmarks and exhibits strong out-of-distribution generalization on challenging abstract visual reasoning tasks. We also empirically analyze the role of each training component and discuss our early unsuccessful attempts, providing insights for future developments in visual latent reasoning. Our model, data, and code are available at <https://github.com/NOVAGlow646/Monet>.*

1. Introduction

Recent work [10, 15, 25, 29, 49] has demonstrated that the incorporation of auxiliary images in the intermediate steps of chain-of-thought (CoT) can improve the visual reasoning of multimodal large language models (MLLMs) [1, 5, 8, 17, 30, 38]. A recent trend is to acquire auxiliary images by training MLLMs to predict key region coordinates [10, 12, 32, 49], invoking visual tools such as grounding or depth estimation models [15, 25, 33], or generating executable codes [46, 47] to modify the input image. While promising, these methods are constrained by a limited set of external tools, lacking the flexible, human-like visual reasoning that arises within an internal perceptual space.

To emulate the flexible visual reasoning of humans, we investigate training MLLMs to reason directly in the continuous latent space. Specifically, we enable MLLMs to generate latent embeddings that serve as intermediate visual thoughts beyond textual descriptions and image embeddings, thus eliminating the need for explicit auxiliary images and overcoming the rigidity of external tools. Recent studies have begun exploring latent visual reasoning strategies [20, 23, 42]. Most existing methods simply align the generated embeddings with those of auxiliary images while applying a next-token prediction (NTP) loss on text tokens during supervised fine-tuning (SFT), and directly apply GRPO [27] for reinforcement learning (RL), which leads to two key limitations: (1) **poor scalability**, as alignment over thousands of image tokens incurs high computational and memory costs, and using mean pooling to compress the image tokens [42] will distort detailed visual features; (2) **insufficient optimization of latent embeddings**, since the NTP objective in SFT can be easily overfit and the GRPO loss can only be computed for text tokens, the optimization of latent embeddings are ignored; Consequently, their improvements remain limited and task-specific.

Motivated by these observations, we propose **Monet**¹, a novel training framework for multimodal latent reasoning that trains a text-output MLLM (Qwen2.5-VL-7B [1])

*Corresponding Author

¹The name **Monet** is derived from the great Impressionist artist, as our method enables MLLMs to think abstractly in the latent visual space.

to perform latent reasoning through SFT and RL. The SFT stage aims at equipping the model with the fundamental ability to generate and reason with latent embeddings. To tackle limitation (1), instead of directly aligning latent embeddings with those of auxiliary images, we introduce dual supervision signals. First, since latent embeddings are intended to facilitate reasoning, we align the hidden representations of text tokens corresponding to crucial visual features observed from the auxiliary images, when conditioned on either auxiliary images or generated latent embeddings. Second, to preserve visual information, we use a controlled attention mask so that latent embeddings can directly attend to auxiliary image embeddings. To overcome limitation (2), we optimize the alignment loss solely through latent embeddings by stopping gradients on non-latent representations. Furthermore, we propose **Vlusal-latent Policy Optimization (VLPO)**, which computes policy gradient directly for latent embeddings by estimating their output probability.

Our contributions are summarized as follows:

1. We propose **Monet-SFT**, a three-stage supervised fine-tuning framework that trains MLLMs to generate and reason with latent embeddings.
2. We propose **VLPO**, a novel RL algorithm tailored for latent reasoning. Unlike GRPO, which targets text reasoning, VLPO incorporates latent embeddings into the total loss by computing an approximate probability for the latent embeddings collected during rollout.
3. We identify the limitations of existing image-text interleaved datasets: the unnecessary usage and inaccuracy of auxiliary images. To address these limitations, we further propose a multi-stage data curation pipeline to construct **Monet-SFT-152K**, a high-quality dataset with image-text interleaved chain-of-thoughts (CoT) for SFT.
4. Extensive experiments on real-world perception and reasoning benchmarks show that Monet-SFT and VLPO outperform conventional SFT + GRPO, cropping-based methods, and prior latent visual reasoning approaches. VLPO further enhances out-of-distribution (OOD) generalization on unseen tasks.

2. Related Work

2.1. Think with Images

These studies can be broadly divided into two categories based on how they obtain intermediate-step auxiliary images: (1) directly emphasizing visual content from the original image, such as grounding, cropping, or re-inputting selected image tokens [3, 7, 25, 34, 44, 49]; and (2) creating new visual content beyond the original image, such as invoking external tools or code interpreters to edit the image (e.g., drawing lines, adding bounding boxes, computing depth maps) [4, 12, 15, 24, 33, 46, 47, 50], or generating new images via text-to-image models [6, 21, 41].

While these approaches significantly enhance visual perception and reasoning, they also introduce key limitations. First, models trained for specific visual tools, such as bounding box prediction, struggle to generalize to tasks requiring more complex visual operations (e.g., visual math, spatial, or graphic reasoning). Second, tool-dependent reasoning increases the cognitive burden for training. Models often fail to generate valid tool calls or executable code without extensive supervision. Finally, reliance on external tools or interpreters necessitates asynchronous, multi-turn inference, complicating deployment and increasing latency.

2.2. Reasoning in Latent Space

Recent studies have explored training LLMs to reason with continuous latent embeddings—either by replacing discrete text tokens with self-generated continuous embeddings [2, 13, 14, 28, 36, 39] or by using reweighted combinations of text embeddings [23, 42]. This line of research aims to enhance reasoning flexibility by enabling exploration in a continuous latent space rather than a discrete linguistic space, while also shortening reasoning chains.

Building on this direction, several works have extended latent reasoning to MLLMs. Li et al. [20], Yang et al. [42] propose to align generated latent embeddings with those of auxiliary images. Yang et al. [42] further compresses image embeddings with mean pooling before alignment, which may potentially distort visual information, while Li et al. [20] focuses only on cropped image regions, which cannot encode visual operations over the entire image. Pham and Ngo [23] removes auxiliary images entirely, optimizing latent embeddings via next-token prediction, but with limited gains. Moreover, these methods neglect latent-space optimization during RL.

3. Method

3.1. Overview

Inference. As illustrated in Figure 1 (left), given a question and its corresponding image, Monet generates text-latent interleaved reasoning chains. Specifically, it can automatically decide when to generate a special token, “<latent>”, to initiate latent reasoning during inference. We modify the decoding process so that the representation from the last layer of the MLLM decoder is fed back as the next input embedding. After generating a predefined number K of latent embeddings, we insert a special stop token, “</latent>”, prompting the model to switch back to language-based reasoning. Consistent with the findings of Li et al. [20], this fixed-length decoding strategy proves simple yet effective. We provide a comprehensive analysis of the effect of K in Section 4.4.

Training. The training of Monet includes a three-stage SFT and an RL stage. We briefly introduce the process and

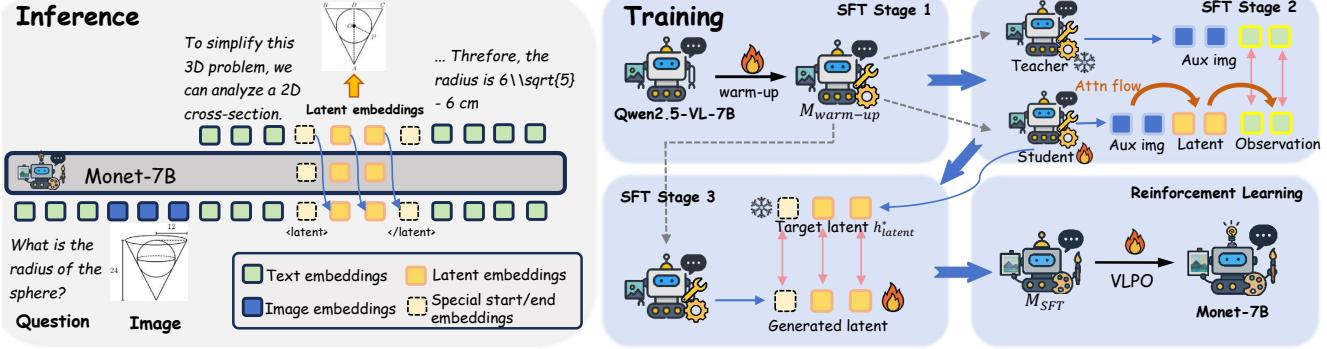


Figure 1. Method overview. *Left:* During inference, Monet can automatically decide when to start latent reasoning by outputting a special start embedding. We fix the output length of the latent embeddings. *Right:* We propose a three-stage SFT (Section 3.3) and RL (Section 3.4) framework. The SFT stages progressively warm up the model, generate high-quality latent embeddings, and distill latent reasoning ability. The RL stage further refines the model using our VLPO algorithm, specifically designed for latent reasoning.

- goal of each stage to provide an overview (Figure 1 right).
- SFT Stage 1.** Adapts the model to interleaved reasoning patterns. The warmed-up parameters $M_{\text{warm-up}}$ are used to initialize the teacher and student models in Stage 2.
 - SFT Stage 2.** Generates high-quality target latent embeddings $\mathbf{h}^*_{\text{latent}}$ for Stage 3. The student model learns to produce these embeddings by: (2) aligning its key observational token representations with those of the fixed teacher model; and (2) making latent embeddings directly attend to auxiliary image embeddings through controlled attention flow.
 - SFT Stage 3.** Trains the model to generate latent embeddings without access to ground-truth auxiliary images. Initialized with $M_{\text{warm-up}}$, the model aligns the generated latent embeddings with the target embeddings $\mathbf{h}^*_{\text{latent}}$ generated in Stage 2, resulting in the model M_{SFT} .
 - RL Stage.** Further optimizes M_{SFT} using VLPO.

In the following sections, we will elaborate on the design of each component.

3.2. SFT Training Dataset Construction

To train Monet, we require a dataset with image–text interleaved CoTs that provide supervision for intermediate latent thinking steps. However, existing datasets exhibit three major limitations: (1) Many samples can be trivially solved by directly observing the input image, leading the model to bypass learning meaningful feature in the intermediate images. (2) The intermediate images are sometimes inaccurate, introducing noise into training. (3) All text tokens are treated equally, overlooking those that describe crucial visual information, which could serve as useful supervision signals [11]. To overcome these issues, we design a three-stage data curation pipeline that ensures the necessity and correctness of auxiliary images (Figure 2).

Stage 1: We collect raw text-image interleaved CoT data from ReFocus [12], CogCoM [24], Zebra-CoT [18], and Visual-CoT [26], and retain samples that Qwen2.5-VL-7B

[1] **incorrectly** answers using only the question and the original image to ensure the necessity of auxiliary images. *Stage 2:* From these samples, we keep those that Qwen2.5-VL-7B can **correctly** solve using only the auxiliary images, ensuring that the auxiliary images are accurate for reasoning. *Stage 3:* We use Deepseek-V3.1 [22] and Gemini 2.5 Pro [8] to identify text tokens corresponding to crucial visual observations for answering the final question to provide fine-grained supervision for learning latent embeddings. These steps yield **Monet-SFT-125k**, in which the auxiliary images include cropping, grounding, highlighting, and operations that create new visual references for fine-grained perception and complex visual reasoning tasks. Table 1 presents the statistics of Monet-SFT-125k.

3.3. Supervised Fine-tuning

The SFT stage addresses two key challenges in latent reasoning training of MLLMs: (1) **Latent–visual alignment is computationally expensive**, as auxiliary images contain hundreds or thousands of image embeddings. (2) **The next-token-prediction objective provides weak supervision for latent embeddings.** During SFT, current methods [20, 23, 42] employ cross-entropy loss on text tokens following the latent embeddings, allowing gradients to flow back to the latent embeddings. However, the model can simply memorize the following tokens instead of learning effective latent representations.

To overcome these issues, we introduce a three-stage SFT framework (Figure 3) that enables the model to selectively encode useful visual features from auxiliary images while focusing on the downstream reasoning process, without incurring the cost of explicit latent–visual alignment. We now describe the design of each stage in detail.

3.3.1. SFT Stage 1: Warm up.

We first warm up the base model M_{base} (Qwen2.5-VL-7B) by performing vanilla SFT on the image–text interleaved

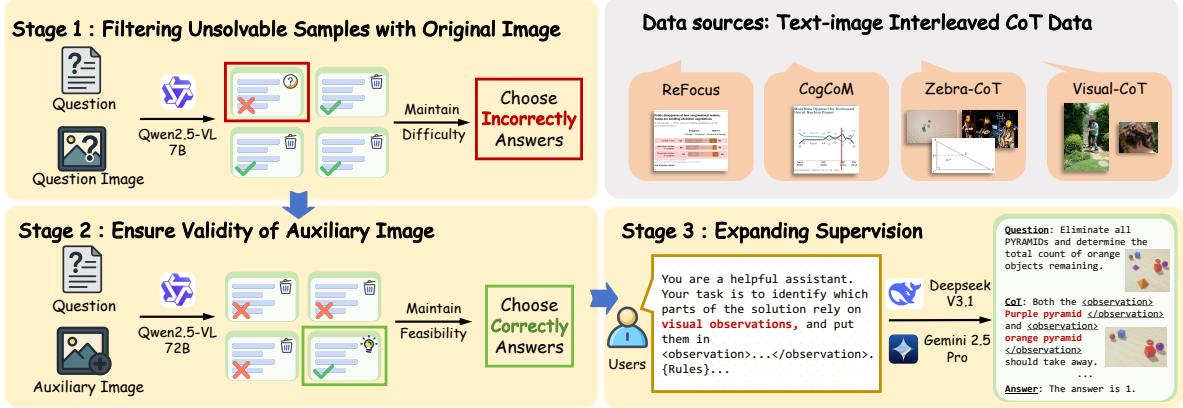


Figure 2. **Construction pipeline of Monet-SFT-125K.** Stage 1 filters hard samples (unsolvable from the original image). Stage 2 keeps those where auxiliary images lead to correct answers, ensuring their necessity and correctness. Stage 3 highlights key visual-observation tokens using advanced LLM judges, providing strong supervision for learning latent embeddings.

Table 1. **Statistics of the Monet-SFT-125K dataset.** It contains 125K real-world, document, chart, and geometry problems with image–text interleaved CoTs, featuring visual operations ranging from simple extractive steps such as cropping and grounding to more complex ones, including drawing auxiliary lines and creating entirely new images that represent intermediate visual states.

Data Source	Problem Domain	Visual Operation Type	Amount
ReFocus [12]	Chart	Drawing bounding boxes, highlighting	0.4K
CogCoM [24]	Real-world, chart	Cropping, drawing auxiliary lines and bounding boxes	0.5K
Visual-CoT [26]	Real-world, documents, chart	Cropping, drawing bounding boxes	118.6K
Zebra-CoT visual search [18]	Real-world, documents, chart	Cropping, drawing bounding boxes	2.7K
Zebra-CoT geometry [18]	Geometry	Drawing auxiliary lines, creating new geometric sketchpads	0.1 K
Zebra-CoT count [18]	3D object counting	Creating new images with 3D objects removed or added	2.9K

CoTs of Monet-SFT-125K, obtaining $M_{\text{warm-up}}$ (see Figure 3 upper left). This stage adapts the model to the image–text interleaved reasoning pattern so it can effectively leverage intermediate-step images when predicting subsequent tokens. Without this adaptation, the model tends to ignore auxiliary images, and the representations of observation tokens may fail to capture sufficient visual information.

Figure 4 illustrates this behavior, in which we compare the token prediction accuracy of the observation tokens with/without the auxiliary images during this warm-up stage: for the unadapted base model, using auxiliary images yields almost no improvement in predicting observation tokens, indicating poor utilization of intermediate visual cues. As the warm-up proceeds, however, the accuracy gain from auxiliary images gradually increases. This demonstrates that SFT on image–text interleaved data drives the model to rely more on intermediate-step visual features rather than memorizing language patterns. Consequently, observation tokens encode meaningful visual information from the auxiliary images and serve as effective supervision signals.

3.3.2. SFT Stage 2: Obtain high-quality target latent embeddings.

In this stage, we train the model to generate latent embeddings that capture useful visual features from auxil-

iary images to support visual reasoning. We initialize both a teacher and a student model from $M_{\text{warm-up}}$. The teacher processes CoTs with ground-truth auxiliary images, while in the student CoT, each auxiliary image segment is followed by autoregressively-generated latent embeddings, and the auxiliary images are made visible only to these latent embeddings via a modified attention mask (Figure 3 right). Now we elaborate on our key designs as follows.

Alignment on Key Observation Tokens. Since latent embeddings are intended to serve the role of auxiliary images in predicting observation tokens, the hidden representations of these tokens should match those obtained when the ground-truth auxiliary images are provided. Inspired by recent text-based latent reasoning work [28, 36], we therefore align the hidden representations of observation tokens under these two conditions.

Specifically, for each training sample, we freeze $M_{\text{warm-up}}$ and extract the observation-token representations from all layers when using auxiliary images, yielding $\mathcal{H}_{\text{obs}}^* = \{\mathbf{h}_{\text{obs}}^{*(i,l)}\}_{i=1}^N$, where i indexes samples and l indexes layers, N is the number of training samples. Let $\hat{\mathbf{h}}_{\text{obs}}^{(i,l)}$ denote the corresponding representations in the student CoT that uses generated latent embeddings. We fix $\mathbf{h}_{\text{obs}}^{*(i,l)}$ and

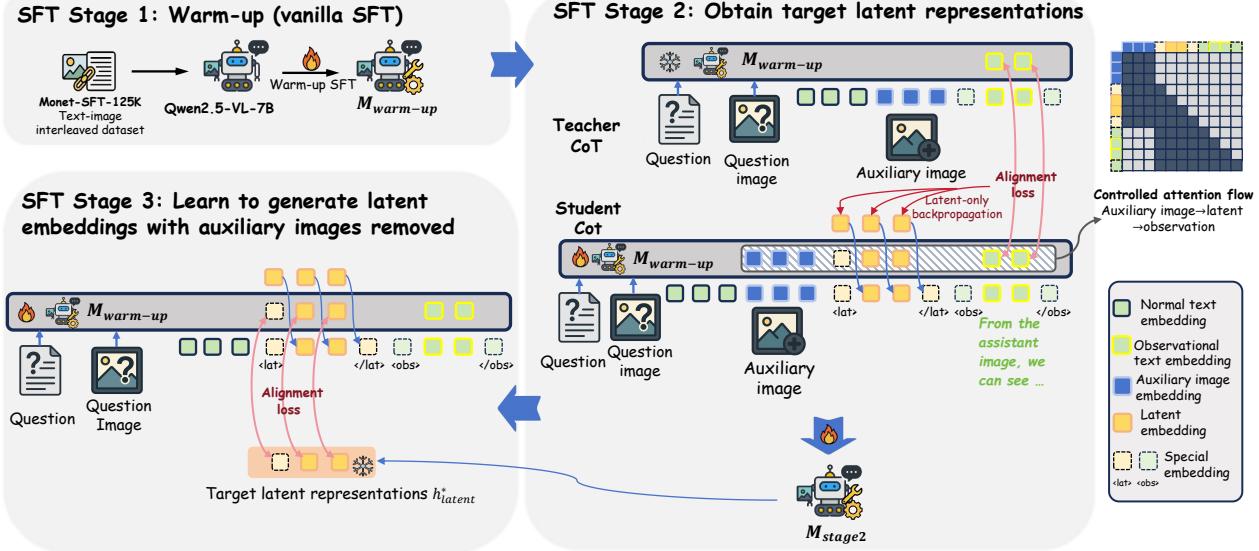


Figure 3. **The proposed three-stage SFT pipeline:** warm-up, supervised latent–observation alignment with controlled attention flow, and latent generation without auxiliary-image access.

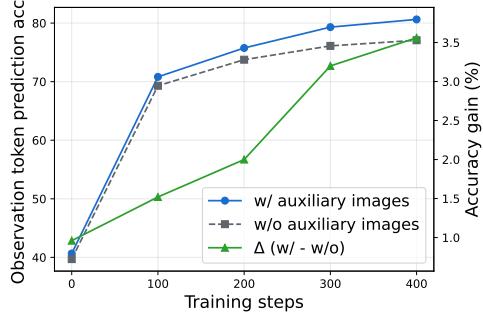


Figure 4. **Prediction accuracy of the observation tokens during warm-up.** Training on image–text interleaved data encourages the model to utilize intermediate visual cues.

maximize their cosine similarity through the alignment loss:

$$\mathcal{L}_{\text{align-obs}} = \frac{1}{N} \sum_i \sum_l \left(1 - \cos(\mathbf{h}_{\text{obs}}^{*(i,l)} \cdot \text{detach}(), \hat{\mathbf{h}}_{\text{obs}}^{(i,l)}) \right). \quad (1)$$

To ensure that the observation alignment loss $\mathcal{L}_{\text{align-obs}}$ can only be minimized by optimizing the latent embeddings instead of through any undesired shortcuts, we restrict gradient flow from $\mathcal{L}_{\text{align-obs}}$ to pass solely through the generated latent embeddings to the model parameters. Implementation details are in the supplementary material. Empirically, removing this latent-only backpropagation leads to significant performance degradation (see Table 2, “w/o latent-only BP”).

“Auxiliary image → latent → observation” Attention Flow.

We observe that using the above alignment loss alone yields suboptimal performance (Table 2, “w/o auxiliary img”), as the observation-token representations may not encode sufficient visual information from the auxil-

iary images. To address this, for the student CoT, we insert the auxiliary image embeddings immediately before each latent-embedding segment and apply a modified attention mask that allows these image embeddings to be attended **only** by the latent embeddings, not by subsequent text tokens. This design offers two benefits: first, latent embeddings can directly access the visual features from auxiliary images without information loss; second, it enforces a structured flow of visual information, i.e., auxiliary images→latent embeddings→observation tokens, encouraging the latent embeddings to selectively encode the relevant visual cues.

Next-Token Prediction Loss. We Additionally apply a standard next-token prediction loss on text tokens to optimize the latent embeddings. Denote the text tokens in student CoT i as $\mathcal{T}^{(i)} = \{\mathbf{y}_t\}_{t=1}^{|\mathcal{T}^{(i)}|}$. For each \mathbf{y}_t , we use $\mathbf{y}_{<t}$ to represent the context, which includes the question text, question image, previous text responses, and latent embeddings (note that the auxiliary images are invisible to the text tokens under our attention design). Denote the output probability distribution as p , the next-token-prediction loss is:

$$\mathcal{L}_{\text{NTP}} = -\frac{1}{N} \sum_i \sum_{\mathbf{y}_t \in \mathcal{T}^{(i)}} \log p(\mathbf{y}_t | \mathbf{y}_{<t}). \quad (2)$$

Total Loss of SFT Stage 2. The total loss is:

$$\mathcal{L}_{\text{stage2}} = \mathcal{L}_{\text{NTP}} + \alpha \mathcal{L}_{\text{align-obs}}, \quad (3)$$

in which we set $\alpha = 2.0$ in all experiments. We denote the student model after training as M_{stage2} . After this stage, we use M_{stage2} to generate the latent embeddings, denoted as $\mathbf{h}_{\text{latent}}^{*(i)}$, which serve as target latent embeddings used for the next stage.

Table 2. **Ablation of the components of the SFT Stage 2.** “Monet-SFT” denotes the full SFT pipeline. “w/o latent-only BP” represents without latent-only backpropagation for the alignment loss. “w/o auxiliary img” disables attention from latent embeddings to auxiliary-image embeddings in the student CoT. Both ablation lead to performance degradation on almost all tasks.

	V*	HR4K	HR8K	MME-RealWorld	VisualPuzzles
Monet-SFT	82.20	68.50	66.00	55.03	30.48
w/o latent-only BP	46.07	40.13	39.00	38.67	33.65
w/o auxiliary img	73.30	63.88	57.63	39.66	28.60

3.3.3. SFT Stage 3: Learn to generate latent embeddings without auxiliary images.

After Stage 2, we obtain target latent embeddings $\mathbf{h}_{\text{latent}}^{*(i)}$ that encode the information contained in auxiliary images. However, the model $M_{\text{stage}2}$ is trained in a setting where latent embeddings can still attend to auxiliary images, which is different from our ultimate goal of generating latent embeddings without ground-truth auxiliary images. To close this gap, we reinitialize the model with $M_{\text{warm-up}}$ and train it to produce latent embeddings $\hat{\mathbf{h}}_{\text{latent}}^{(i)}$ with auxiliary images removed in the CoT (Figure 3 bottom left). We then align them with the fixed target embeddings $\mathbf{h}_{\text{latent}}^{*(i)}$ using the following loss:

$$\mathcal{L}_{\text{align-latent}} = \frac{1}{N} \sum_i \sum_l \left(1 - \cos(\mathbf{h}_{\text{latent}}^{*(i,l)} \cdot \text{detach}(), \hat{\mathbf{h}}_{\text{latent}}^{(i,l)}) \right). \quad (4)$$

Unlike previous latent-visual reasoning work [20, 42], which aligns only the final-layer representations, we align all layers to provide stronger supervision.

We also apply a next-token-prediction loss \mathcal{L}_{NTP} on the text tokens to let the latent embeddings benefit subsequent reasoning, which has the same form as Equation (2).

Total loss of SFT Stage 3. The total loss of Stage 3 is:

$$\mathcal{L}_{\text{stage}3} = \mathcal{L}_{\text{NTP}} + \beta \mathcal{L}_{\text{align-latent}}. \quad (5)$$

in which we set $\beta = 2.0$ in all our experiments.

3.4. VLPO: Visual-latent Policy Optimization

Limitations of GRPO on latent visual reasoning. Prior latent visual reasoning works [20, 42] apply GRPO [27] after SFT. Given a question Q and a question image I , the old policy model π_{old} generates a group of responses $\{\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_G\}$. Then GRPO updates the policy model π_θ

by maximizing:

$$\begin{aligned} \mathcal{J}_{\text{GRPO}}(\theta) &= \mathbb{E}_{Q, I, o \sim \pi_{\text{old}}} \\ &\frac{1}{G} \sum_{i=1}^G \frac{1}{|\mathbf{o}_i|} \sum_{t=1}^{|\mathbf{o}_i|} \min \left[r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip}(r_{i,t}(\theta), 1 - \varepsilon, 1 + \varepsilon) \hat{A}_{i,t} \right] \\ &- \beta \text{KL}(\pi_\theta \| \pi_{\text{ref}}), \end{aligned} \quad (6)$$

where $r_{i,t}(\theta) = \frac{\pi_\theta(\mathbf{o}_{i,t} | Q, I, \mathbf{o}_{i,<t})}{\pi_{\theta_{\text{old}}}(\mathbf{o}_{i,t} | Q, I, \mathbf{o}_{i,<t})}$, $\mathbf{o}_{i,<t}$ are the response tokens before position t , $\hat{A}_{i,t}$ is the advantage computed based on the outcome rewards $\{r_1, r_2, \dots, r_G\}$ of the responses, i.e., $\hat{A}_{i,t} = \frac{r_i - \text{mean}(\{r_1, r_2, \dots, r_G\})}{\text{std}(\{r_1, r_2, \dots, r_G\})}$. π_{ref} is the reference model. A critical limitation is that the GRPO objective **can only be computed on text tokens**. Since latent embeddings have no explicit probability distribution like text tokens, GRPO cannot directly optimize them, leaving the latent reasoning component largely untrained during RL.

To address this limitation, we propose Visual-latent Policy Optimization (VLPO), a novel reinforcement learning objective tailored for latent reasoning. The key idea of VLPO is to estimate the **output probability of the continuous latent embeddings** collected during rollout to compute $r_{i,t}(\theta)$ for the latent embeddings. This enables latent embeddings to be optimized directly with outcome rewards, just like text tokens. Denote the latent embedding generated by π_{old} at position t of the i -th rollout as $\mathbf{h}_{i,t}^{\text{old}}$ with context $[Q, I, \mathbf{o}_{i,<t}]$. To compute $\pi_\theta(\mathbf{h}_{i,t}^{\text{old}} | Q, I, \mathbf{o}_{i,<t})$, we model $\mathbf{h}_{i,t}^{\text{old}}$ as a sample drawn from a latent Gaussian distribution whose mean is the latent embedding generated by π_θ given the same context. We denote this mean by $\mathbf{h}_{i,t}^\theta$. Under this assumption, $\pi_\theta(\mathbf{h}_{i,t}^{\text{old}} | Q, I, \mathbf{o}_{i,<t})$ can be computed as:

$$\pi_\theta(\mathbf{h}_{i,t}^{\text{old}} | Q, I, \mathbf{o}_{i,<t}) = \exp \left(-\frac{1}{2\sigma^2} \|\mathbf{h}_{i,t}^{\text{old}} - \mathbf{h}_{i,t}^\theta\|^2 - \text{const} \right) \quad (7)$$

where σ is a predefined scalar hyperparameter and const is a constant. $\pi_\theta(o_{i,t}^{\text{latent}} | Q, I, \mathbf{o}_{i,<t})$ measures the probability that π_θ chooses $o_{i,t}^{\text{latent}}$ as the “action” at position t of response i , enabling us to compute $r_{i,t}(\theta)$ for a latent step:

$$r_{i,t}(\theta) = \frac{\pi_\theta(\mathbf{h}_{i,t}^{\text{old}} | Q, I, \mathbf{o}_{i,<t})}{\pi_{\theta_{\text{old}}}(\mathbf{h}_{i,t}^\theta | Q, I, \mathbf{o}_{i,<t})} = \exp \left(-\frac{1}{2\sigma^2} \|\mathbf{h}_{i,t}^{\text{old}} - \mathbf{h}_{i,t}^\theta\|^2 \right) \quad (8)$$

Replacing the original $r_{i,t}(\theta)$ for latent steps yields our VLPO objective. During optimization, we fix the rollout rollout $\mathbf{h}_{i,t}^{\text{old}}$ and optimize $\mathbf{h}_{i,t}^\theta$ generated by the policy.

Reward Design. We use an accuracy reward (1 for a correct answer; 0 otherwise) and a format reward encouraging the model to place the final answer in “\boxed{}”. Importantly, we do not reward latent-reasoning behavior itself, as doing so would incentivize the model to invoke latent reasoning indiscriminately. We find this simple reward scheme yields strong empirical performance.

How does VLPO work? When $\hat{A}_{i,t} > 0$, maximizing the VLPO objective will increase $r_{i,t}(\theta)$. For text tokens, this means maximizing its output probability; For latent embeddings, Equation (8) shows that maximizing this $r_{i,t}(\theta)$ equals to minimizing $\|\mathbf{h}_{i,t}^{\text{old}} - \mathbf{h}_{i,t}^{\theta}\|^2$, thereby **pulling the policy latent embeddings to a “good-action” latent embedding $\mathbf{h}_{i,t}^{\theta}$ that led to a positive outcome reward**. Thus, VLPO directly optimizes latent embeddings using reward signals—an ability that GRPO fundamentally lacks.

4. Experiment

4.1. Experiment Setup

Training and Evaluation Setup. In SFT Stage 1, we train the model for 4 epochs. For SFT stage 2 and 3, we train 1000 steps (about 1 epoch). For RL training, we adopt a 3.2K subset of Thyme-RL [46] and train for 1 epoch. For Monet-SFT model, we fix the training time latent size to 8; For **Monet-7B** (Monet-SFT + VLPO), we use a training latent size of 10. At test time, we select the best latent size from $\{8, 10, 12, 16\}$. The complete training and experimental setups are provided in the supplementary material.

Evaluated Benchmarks. We adopt the VLMEvalKit [9] framework for fair evaluation. We evaluate Monet on V* [40], HRBench4k [37], HRBench8K [37], MME-RealWorld [45], covering fine-grained perception and reasoning tasks on real-world, chart and OCR data. To evaluate OOD capabilities, we also evaluate on VisualPuzzles [31], which includes visual logical puzzles that demands abstract visual reasoning abilities rather than prior knowledge.

Baselines. We compare Monet against the following baselines: (1) **Qwen2.5-VL-7B**, the base model to train Monet-7B; (2) **Vanilla SFT**, which directly train Qwen2.5-VL-7B with Monet-SFT-125K; (3) **Vanilla SFT + GRPO**, which applies GRPO after Vanilla SFT; (4) **Deepeyes** [49], a representative “think with images” approach that performs cropping to improve perception; (5) **LVR** [20], a recent work on latent visual reasoning that aligns the latent and auxiliary image embeddings; (6) other high-performance MLLMs, including Gemini 2.5 Pro [8], GPT-4o [16], Pangea-7B [43], and LLaVA-OneVision-72B [19].

4.2. Main Results

From Table 3 and Tabel 4, we can see that: (1) **Monet substantially enhances both perception and reasoning on real-world, chart, and OCR tasks, consistently outperforming the baselines.** It improves Qwen2.5-VL-7B by 4.25%–9.75%, surpasses both Vanilla SFT and SFT + GRPO trained on the same data, and outperforms Deepeyes and LVR on most benchmarks. (2) **Monet demonstrates strong OOD generalization.** It achieves the best performance on VisualPuzzles, which mainly contains abstract visual reasoning problems that are unseen during training.

4.3. Analysis of the Role of the Components Design

In this section, we present ablation studies to assess the necessity and effectiveness of each component in our training framework (Table 5). We also discuss several unsuccessful attempts from our early exploration of training objectives and pipeline designs, which we hope can provide insights for future advancements in latent reasoning for MLLMs.

Takeaway 1. Single-signal supervision is less effective. We remove the representation-alignment loss on observation tokens in SFT Stage 2 while retaining auxiliary-image embeddings in the student CoT under the “image→latent→observation” attention flow, which reflects our initial attempt. As shown by “A-SFT w/o obs. token align”, removing the alignment loss causes a substantial drop in performance. This indicates that supervision from observation tokens is crucial. Similarly, removing the auxiliary images visible to latent embeddings while using the alignment loss alone (“w/o auxiliary img”) also significantly degrades performance, showing the importance of visual supervision on the latent embeddings. These two ablations demonstrate the necessity of the dual supervisions we use.

Takeaway 2. GRPO doesn’t consistently improve Monet-SFT, whereas VLPO does. This is evident when comparing “Monet-SFT” with “Monet-SFT + GRPO”, indicating that GRPO is not well-suited for latent reasoning. In contrast, applying VLPO significantly enhances the performance of Monet-SFT.

Takeaway 3. Latent-related losses must backpropagate only through the latent embeddings. “w/o latent-only BP” shows that when the alignment loss is allowed to update non-latent representations, performance drops sharply. This indicates that the model may otherwise exploit shortcut paths by minimizing the loss without actually improving the latent embeddings. Therefore, enforcing latent-only backpropagation is necessary.

4.4. Effect of the Number of the Latent Embeddings on Performance

In this section, we analyze how training and test-time latent size (K_{train} and K_{test}) affect performance. For Monet-SFT, we choose $K_{\text{train}} \in \{8, 10, 12\}$. For RL models, we take Monet-SFT with $K_{\text{train}} = 8$ and further apply VLPO (with RL $K_{\text{train}} = 10$) and GRPO. Results are shown in Figure 5.

Takeaway 1. For in-distribution tasks, improvements indeed come from using abstract visual embeddings; while for OOD tasks, only VLPO incentivizes useful latent reasoning. Across all in-distribution tasks (V*, HRBench, MME-RealWorld), models perform worst when $K_{\text{test}} = 0$. However, on the OOD task (VisualPuzzles), only the VLPO-enhanced model consistently benefits from latent reasoning, i.e., using a $K_{\text{test}} > 0$ induces better performances than $K_{\text{test}} = 0$, indicating that SFT alone cannot

Table 3. **Performance on real-world perception and reasoning benchmarks.** The best-performing open-source model for each dataset is highlighted in **bold**. Results marked with “*” are reported by other papers [20, 46], while the others are from our own evaluations.

Model	V*			HRBench4K			HRBench8K			MME-RealWorld-Lite		
	Overall	Attribute	Spatial	Overall	FSP	FCP	Overall	FSP	FCP	Overall	Reasoning	Perception
<i>Proprietary Model</i>												
GPT-4o [16]	67.5*	72.2*	60.5*	59.0*	70.0*	48.0*	55.5*	62.0*	49.0*	52.0*	48.3*	54.4*
<i>Open-Source Model</i>												
Qwen2.5-VL-7B [1]	76.44	77.39	75.00	68.00	80.25	55.75	63.75	73.75	53.75	45.75	39.73	49.62
+ vanilla SFT	81.68	83.48	78.95	68.38	78.28	58.5	61.63	70.75	52.50	51.28	46.4	54.41
+ vanilla SFT + GRPO	78.53	78.26	78.95	70.00	83.25	56.75	66.75	78.00	55.50	52.42	48.13	55.18
Deepeyes [49]	83.25	84.35	81.58	71.25	83.75	58.75	65.13	77.00	53.25	54.28	50.53	56.63
LVR [20]	80.6*	81.7*	79.0*	-	-	-	-	-	-	-	-	-
<i>Our Model</i>												
Monet-7B (SFT + VLPO)	83.25	83.48	82.89	71.00	85.25	56.75	68.00	79.75	56.25	55.50	51.07	58.34
<i>Relative Improvement</i>	+6.81	+6.09	+7.89	+3.00	+5.00	+1.00	+4.25	+6.00	+2.50	+9.75	+11.34	+8.72

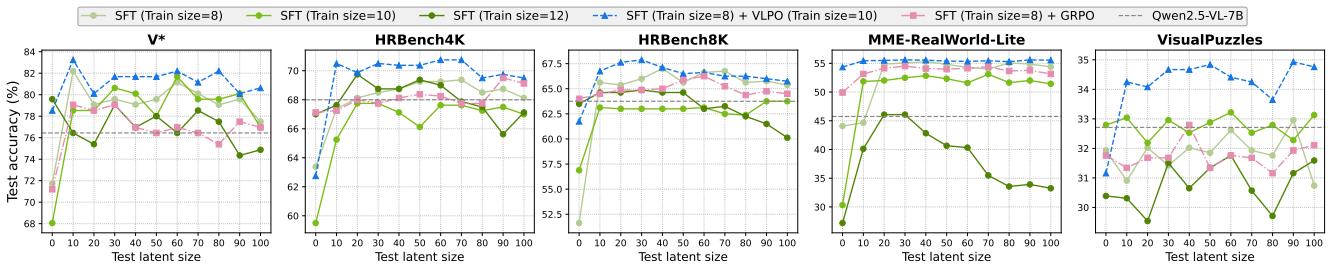


Figure 5. **Effect of the number of abstract visual embeddings used during training and inference on test accuracy.** The dashed line marks the accuracy of Qwen2.5-VL-7B.

Table 4. **Performance on the VisualPuzzles benchmark (out-of-distribution).** The best-performing open-source model for each dataset is highlighted in **bold**. Results marked with “*” are taken from the original VisualPuzzles [31] paper, while the others are obtained from our own evaluations.

Model	VisualPuzzles					
	Overall	Algorithmic	Analogical	Deductive	Inductive	Spatial
<i>Proprietary Models</i>						
Gemini 2.5 Pro [8]	49.5*	60.0*	64.0*	60.0*	29.7*	36.4*
GPT-4o [16]	41.3*	49.2*	58.3*	49.0*	27.3*	26.2*
<i>Open-Source Models</i>						
Qwen2.5-VL-7B	32.71	37.02	21.80	47.50	26.32	21.80
+ vanilla SFT	33.99	40.46	30.81	46.00	24.40	29.02
+ vanilla SFT + GRPO	30.99	36.26	25.12	43.50	23.45	27.27
Deepeyes [49]	32.96	37.79	27.01	41.00	26.79	27.01
Pangea-7B [43]	31.3*	32.4*	23.7*	38.5*	28.7*	32.5*
LLaVA-OneVision-72B [19]	30.8*	34.7*	26.5*	37.0*	27.3*	28.7*
<i>Our Model</i>						
Monet-7B	35.02	45.80	30.81	47.50	26.79	25.52
<i>Relative Improvement</i>	+2.31	+6.09	+8.78	+9.01	+0.47	+3.72

Table 5. **Ablation of the components of Monet.**

Model	V*	HRBench4K	HRBench8K	MME-RealWorld-Lite	VisualPuzzles
Monet-7B (full)	83.25	71.00	68.00	55.50	35.02
Monet-SFT (w/o VLPO)	82.20	68.50	66.00	52.68	30.48
+ GRPO	80.10	69.00	64.75	54.19	31.51
w/o latent-only BP	46.07	40.13	39.00	38.67	33.65
w/o auxiliary img	73.30	63.88	57.63	39.66	28.60
w/o obs token align	75.39	67.25	63.50	46.90	27.48

induce strong OOD generalization.

Takeaway 2. For in-distribution tasks, Monet-SFT enables test-time scaling of latent embeddings; VLPO

further extends this trend to OOD scenarios. On V*, HRBench, and MME-RealWorld, performance typically peaks at a test-time latent size K_{test} larger than the training-time one K_{train} . Before the peak, the performance improves as K_{test} increases.

Takeaway 3. VLPO improves robustness to the choice of K_{test} . Compared with other models, Monet-SFT + VLPO exhibits the most stable performance as K_{test} varies.

Takeaway 4. GRPO mainly strengthens non-latent reasoning and provides limited benefits for latent reasoning. When $K_{\text{test}} = 0$, Monet-SFT ($K_{\text{train}} = 8$) + GRPO is comparable to or better than Monet-SFT ($K_{\text{train}} = 8$), showing gains in non-latent reasoning; However, when $K_{\text{test}} > 0$, Monet-SFT + GRPO underperforms the SFT model for most K_{test} , suggesting minimal improvements to latent reasoning. This aligns with our analysis in Section 3.4: GRPO optimizes only text tokens and thus provides no direct supervision for latent embeddings.

5. Conclusion and Limitations

We introduce Monet, a framework that enables MLLMs to perform visual reasoning in the latent space. Our Monet-SFT pipeline enables the model to selectively encode useful visual cues from auxiliary images into compact latent embeddings, while our VLPO algorithm explicitly optimizes these embeddings during RL. We further construct

the Monet-SFT-125K dataset to support SFT. Experiments show that Monet substantially improves real-world perception and reasoning performance and exhibits strong OOD generalization on abstract visual reasoning tasks. Together, these results demonstrate Monet as a promising direction toward more flexible and general multimodal reasoning.

We summarize the limitations of our work as follows. First, Monet relies on a multi-stage SFT pipeline, which may increase the overall training complexity and overhead. Second, we have not yet explored how different reward designs might influence latent visual reasoning in MLLMs, leaving room for exploration and further enhancement.

References

- [1] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. [1](#), [3](#), [8](#)
- [2] Natasha Butt, Ariel Kwiatkowski, Ismail Labiad, Julia Kempe, and Yann Ollivier. Soft tokens, hard truths. *arXiv preprint arXiv:2509.19170*, 2025. [2](#)
- [3] Xinyan Chen, Renrui Zhang, Dongzhi Jiang, Aojun Zhou, Shilin Yan, Weifeng Lin, and Hongsheng Li. Mint-cot: Enabling interleaved visual tokens in mathematical chain-of-thought reasoning. *arXiv preprint arXiv:2506.05331*, 2025. [2](#)
- [4] Yang Chen, Yufan Shen, Wenzuan Huang, Shen Zhou, Qun-shu Lin, Xinyu Cai, Zhi Yu, Botian Shi, and Yu Qiao. Learning only with images: Visual reinforcement learning with reasoning, rendering, and visual feedback. *arXiv preprint arXiv:2507.20766*, 2025. [2](#)
- [5] Zihui Cheng, Qiguang Chen, Xiao Xu, Jiaqi Wang, Weiyun Wang, Hao Fei, Yidong Wang, Alex Jinpeng Wang, Zhi Chen, Wanxiang Che, et al. Visual thoughts: A unified perspective of understanding multimodal chain-of-thought. *arXiv preprint arXiv:2505.15510*, 2025. [1](#)
- [6] Ethan Chern, Zhulin Hu, Steffi Chern, Siqi Kou, Jiadi Su, Yan Ma, Zhijie Deng, and Pengfei Liu. Thinking with generated images. *arXiv preprint arXiv:2505.22525*, 2025. [2](#)
- [7] Jiwan Chung, Junhyeok Kim, Siyeol Kim, Jaeyoung Lee, Min Soo Kim, and Youngjae Yu. Don't look only once: Towards multimodal interactive reasoning with selective visual revisit. *arXiv preprint arXiv:2505.18842*, 2025. [2](#)
- [8] Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blissein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. [1](#), [3](#), [7](#), [8](#)
- [9] Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models. In *Proceedings of the 32nd ACM international conference on multimedia*, 2024. [7](#), [1](#)
- [10] Yue Fan, Xuehai He, Diji Yang, Kaizhi Zheng, Ching-Chen Kuo, Yuting Zheng, Sravana Jyothi Narayananaraju, Xinze Guan, and Xin Eric Wang. Grit: Teaching mllms to think with images. *arXiv preprint arXiv:2505.15879*, 2025. [1](#)
- [11] Lizhe Fang, Yifei Wang, Zhaoyang Liu, Chenheng Zhang, Stefanie Jegelka, Jinyang Gao, Bolin Ding, and Yisen Wang. What is wrong with perplexity for long-context language modeling? In *ICLR*, 2025. [3](#)
- [12] Xingyu Fu, Minqian Liu, Zhengyuan Yang, John Corring, Yijuan Lu, Jianwei Yang, Dan Roth, Dinei Florencio, and Cha Zhang. Refocus: Visual editing as a chain of thought for structured image understanding. In *ICML*, 2025. [1](#), [2](#), [3](#), [4](#)
- [13] Jonas Geiping, Sean McLeish, Neel Jain, John Kirchenbauer, Siddharth Singh, Brian R Bartoldson, Bhavya Kailkhura, Abhinav Bhatele, and Tom Goldstein. Scaling up test-time compute with latent reasoning: A recurrent depth approach. *arXiv preprint arXiv:2502.05171*, 2025. [2](#)
- [14] Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.06769*, 2024. [2](#)
- [15] Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal language models. *Advances in Neural Information Processing Systems*, 2024. [1](#), [2](#)
- [16] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. [7](#), [8](#)
- [17] Pu Jian, Junhong Wu, Wei Sun, Chen Wang, Shuo Ren, and Jiajun Zhang. Look again, think slowly: Enhancing visual reflection in vision-language models. *arXiv preprint arXiv:2509.12132*, 2025. [1](#)
- [18] Ang Li, Charles Wang, Kaiyu Yue, Zikui Cai, Ollie Liu, Deqing Fu, Peng Guo, Wang Bill Zhu, Vatsal Sharan, Robin Jia, et al. Zebra-cot: A dataset for interleaved vision language reasoning. *arXiv preprint arXiv:2507.16746*, 2025. [3](#), [4](#)
- [19] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. [7](#), [8](#)
- [20] Bangzheng Li, Ximeng Sun, Jiang Liu, Ze Wang, Jialian Wu, Xiaodong Yu, Hao Chen, Emad Barsoum, Muham Chen, and Zicheng Liu. Latent visual reasoning. *arXiv preprint arXiv:2509.24251*, 2025. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#)
- [21] Chengzu Li, Wenshan Wu, Huanyu Zhang, Yan Xia, Shaoguang Mao, Li Dong, Ivan Vulić, and Furu Wei. Imagine while reasoning in space: Multimodal visualization-of-thought. *arXiv preprint arXiv:2501.07542*, 2025. [2](#)
- [22] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. [3](#)

- [23] Tan-Hanh Pham and Chris Ngo. Multimodal chain of continuous thought for latent-space reasoning in vision-language models. *arXiv preprint arXiv:2508.12587*, 2025. 1, 2, 3
- [24] Ji Qi, Ming Ding, Weihan Wang, Yushi Bai, Qingsong Lv, Wenyi Hong, Bin Xu, Lei Hou, Juanzi Li, Yuxiao Dong, et al. Cogcom: A visual language model with chain-of-manipulations reasoning. In *ICLR*, 2025. 2, 3, 4
- [25] Gabriel Sarch, Snigdha Saha, Naitik Khandelwal, Ayush Jain, Michael J Tarr, Aviral Kumar, and Katerina Fragkiadaki. Grounded reinforcement learning for visual reasoning. *arXiv preprint arXiv:2505.23678*, 2025. 1, 2
- [26] Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. Visual cot: Unleashing chain-of-thought reasoning in multi-modal language models. *CoRR*, 2024. 3, 4
- [27] Zhihong Shao, Peiyi Wang, Qihao Zhu, R. Xu, Jun-Mei Song, Mingchuan Zhang, Y. K. Li, Yu Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *ArXiv*, abs/2402.03300, 2024. 1, 6
- [28] Zhenyi Shen, Hanqi Yan, Linhai Zhang, Zhanghao Hu, Yali Du, and Yulan He. Codi: Compressing chain-of-thought into continuous space via self-distillation. *arXiv preprint arXiv:2502.21074*, 2025. 2, 4
- [29] Yang Shi, Yuhao Dong, Yue Ding, Yuran Wang, Xuanyu Zhu, Sheng Zhou, Wenting Liu, Haochen Tian, Rundong Wang, Huanqian Wang, et al. Realunify: Do unified models truly benefit from unification? a comprehensive benchmark. *arXiv preprint arXiv:2509.24897*, 2025. 1
- [30] Yang Shi, Jiaheng Liu, Yushuo Guan, Zhenhua Wu, Yuanxing Zhang, Zihao Wang, Weihong Lin, Jingyun Hua, Zekun Wang, Xinlong Chen, et al. Mavors: Multi-granularity video representation for multimodal large language model. *arXiv preprint arXiv:2504.10068*, 2025. 1
- [31] Yueqi Song, Tianyue Ou, Yibo Kong, Zecheng Li, Graham Neubig, and Xiang Yue. Visualpuzzles: Decoupling multi-modal reasoning evaluation from domain knowledge. *arXiv preprint arXiv:2504.10342*, 2025. 7, 8
- [32] Alex Su, Haozhe Wang, Weiming Ren, Fangzhen Lin, and Wenhui Chen. Pixel reasoner: Incentivizing pixel-space reasoning with curiosity-driven reinforcement learning. *arXiv preprint arXiv:2505.15966*, 2025. 1
- [33] Zhaochen Su, Linjie Li, Mingyang Song, Yunzhuo Hao, Zhengyuan Yang, Jun Zhang, Guanjie Chen, Jiawei Gu, Juntao Li, Xiaoye Qu, et al. Openthinkimg: Learning to think with images via visual tool reinforcement learning. *arXiv preprint arXiv:2505.08617*, 2025. 1, 2
- [34] Guangyan Sun, Mingyu Jin, Zhenting Wang, Cheng-Long Wang, Siqi Ma, Qifan Wang, Tong Geng, Ying Nian Wu, Yongfeng Zhang, and Dongfang Liu. Visual agents as fast and slow thinkers. In *ICLR*, 2025. 2
- [35] Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, Shengyi Huang, Kashif Rasul, and Quentin Galloüdec. Trl: Transformer reinforcement learning, 2020. 1
- [36] Jianwei Wang, Ziming Wu, Fuming Lai, Shaobing Lian, and Ziqian Zeng. Synadapt: Learning adaptive reasoning in large language models via synthetic continuous chain-of-thought. *arXiv preprint arXiv:2508.00574*, 2025. 2, 4
- [37] Wenbin Wang, Liang Ding, Minyan Zeng, Xiabin Zhou, Li Shen, Yong Luo, Wei Yu, and Dacheng Tao. Divide, conquer and combine: A training-free framework for high-resolution image perception in multimodal large language models. In *AAAI*, 2025. 7
- [38] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xinguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3. 5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025. 1
- [39] Xilin Wei, Xiaoran Liu, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Jiaqi Wang, Xipeng Qiu, and Dahua Lin. Simcot: Supervised implicit chain-of-thought. *arXiv preprint arXiv:2509.20317*, 2025. 2
- [40] Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms. In *CVPR*, 2024. 7
- [41] Yi Xu, Chengzu Li, Han Zhou, Xingchen Wan, Caiqi Zhang, Anna Korhonen, and Ivan Vulić. Visual planning: Let's think only with images. *arXiv preprint arXiv:2505.11409*, 2025. 2
- [42] Zeyuan Yang, Xueyang Yu, Delin Chen, Maohao Shen, and Chuang Gan. Machine mental imagery: Empower multimodal reasoning with latent visual tokens. *arXiv preprint arXiv:2506.17218*, 2025. 1, 2, 3, 6
- [43] Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim, Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kantharuban, Lintang Sutawika, Sathyanarayanan Ramamoorthy, and Graham Neubig. Pangea: A fully open multilingual multimodal llm for 39 languages. In *ICLR*, 2024. 7, 8
- [44] Xintong Zhang, Zhi Gao, Bofei Zhang, Pengxiang Li, Xiaowen Zhang, Yang Liu, Tao Yuan, Yuwei Wu, Yunde Jia, Song-Chun Zhu, et al. Chain-of-focus: Adaptive visual search and zooming for multimodal reasoning via rl. *arXiv preprint arXiv:2505.15436*, 2025. 2
- [45] Yi-Fan Zhang, Huanyu Zhang, Haochen Tian, Chaoyou Fu, Shuangqing Zhang, Junfei Wu, Feng Li, Kun Wang, Qingsong Wen, Zhang Zhang, et al. Mme-realworld: Could your multimodal llm challenge high-resolution real-world scenarios that are difficult for humans? *arXiv preprint arXiv:2408.13257*, 2024. 7
- [46] Yi-Fan Zhang, Xingyu Lu, Shukang Yin, Chaoyou Fu, Wei Chen, Xiao Hu, Bin Wen, Kaiyu Jiang, Changyi Liu, Tianke Zhang, et al. Thyme: Think beyond images. *arXiv preprint arXiv:2508.11630*, 2025. 1, 2, 7, 8
- [47] Shitian Zhao, Haoquan Zhang, Shaoheng Lin, Ming Li, Qilong Wu, Kaipeng Zhang, and Chen Wei. Pyvision: Agentic vision with dynamic tooling. *arXiv preprint arXiv:2507.07998*, 2025. 1, 2
- [48] Yaowei Zheng, Junting Lu, Shenzhi Wang, Zhangchi Feng, Dongdong Kuang, and Yuwen Xiong. Easyl1: An efficient, scalable, multi-modality rl training framework, 2025. 1
- [49] Ziwei Zheng, Michael Yang, Jack Hong, Chenxiao Zhao, Guohai Xu, Le Yang, Chao Shen, and Xing Yu. Deep-eyes: Incentivizing “thinking with images” via reinforcement learning. *arXiv preprint arXiv:2505.14362*, 2025. 1, 2, 7, 8

- [50] Zetong Zhou, Dongping Chen, Zixian Ma, Zhihan Hu, Mingyang Fu, Sinan Wang, Yao Wan, Zhou Zhao, and Ranjay Krishna. Reinforced visual perception with tools. *arXiv preprint arXiv:2509.01656*, 2025. [2](#)

Monet: Reasoning in Latent Visual Space Beyond Image and Language

Supplementary Material

A. Additional Experimental Results

We present the complete results on MME-RealWorld-Lite in Table 6. This benchmark spans a wide range of task domains, such as monitoring, autonomous driving, and complex diagram data. Compared with baseline approaches, Monet achieves consistently strong performance across all data sources, demonstrating its generality in both visual reasoning and perception tasks.

B. Implementation Details

B.1. SFT Training

We use TRL [35] to implement our SFT training. In a training CoT, we allow the model to predict the special latent-start token `<latent>` while masking the latent-end token `</latent>`, the observation-start token `<observation>`, and the observation-end token `</observation>` in the SFT labels.

Hyperparameters. We list the training hyperparameters of SFT in Table 7. To avoid OOM, we use a batch size of 1 and set gradient accumulation steps to 16. Additionally, we restrict the total number of pixels (denoted by “max total pixels”) in a training sequence to prevent the OOM issue caused by too many auxiliary images.

Latent only backpropagation in SFT Stage 2. To ensure that gradients flow only through the latent representations to the model parameters, we replace $\mathcal{L}_{\text{align-obs}}$ with the following surrogate loss:

$$\mathcal{L}'_{\text{align-obs}} = \frac{1}{N} \sum_i \text{stop_grad}\left(\frac{\partial \mathcal{L}_{\text{align-obs}}}{\partial \hat{\mathbf{h}}_{\text{latent}}^{(i,L)}}\right)^\top \hat{\mathbf{h}}_{\text{latent}}^{(i,L)}, \quad (9)$$

where `stop_grad` denotes the stop-gradient operation, $\hat{\mathbf{h}}_{\text{latent}}^{(i,L)}$ are the generated latent embeddings in the student CoT, and L is the number of layers in the MLLM’s language model. By differentiating $\mathcal{L}'_{\text{align-obs}}$ with respect to the model parameters θ :

$$\frac{\partial \mathcal{L}'_{\text{align-obs}}}{\partial \theta} = \frac{\partial \mathcal{L}_{\text{align-obs}}}{\partial \hat{\mathbf{h}}_{\text{latent}}^{(i,L)}} \frac{\partial \hat{\mathbf{h}}_{\text{latent}}^{(i,L)}}{\partial \theta}, \quad (10)$$

we can see that minimizing $\mathcal{L}'_{\text{align-obs}}$ is equivalent to optimizing only the generated latent embeddings.

B.2. RL Training

We implement our RL training on EasyR1 [48], an open-source RL training framework for multimodal LLMs.

Hyperparameters. We summarize the RL training hyperparameters in Table 8. Note that we apply an “accuracy threshold” to retain only samples whose group accuracy falls below this threshold but remains non-zero, preventing the advantage from vanishing when the responses in a group are all correct or incorrect.

B.3. Training Data Construction

We provide the prompt for the API models used in the *stage 3* of our training data construction pipeline as below:

System Prompt for Training Data Construction

You are a helpful assistant. Your task is to identify which parts of the solution rely on visual observations (marked by `<latent></latent>`). Wrap these essential observations in `<observation>...</observation>` while leaving all other text unchanged. The input contains only a question (with conditions) and a solution chain-of-thought text — no images. Identifying observations does not require access to images. Segments are delimited by `<STEP_i> ... <END_STEP_i>`.

Rules:

1. Do NOT remove, rename, or reorder any markers. Always output all markers exactly as given.
2. Do NOT add new content to the reasoning.
3. Inside `<observation>...</observation>`, include as little as possible — only the key tokens that truly require image observation.
4. Do NOT put information that can be inferred from the question text or earlier reasoning into `<observation>...</observation>`.
5. If you think there’s no observation, just return the input text unchanged.

C. Detailed Experimental Setup

We use the VLMEvalKit [9] framework for our evaluation. We set the inference-time maximum visual tokens to be $8192 \times 28 \times 28$. Below is the system prompt we used for evaluation.

System Prompt for Evaluation

You are an expert multimodal large language model designed to reason with latent visual embeddings.

Table 6. Full result on MME-RealWorld-Lite.

	MME-RealWorld-Lite										
	Reasoning					Perception					
	Overall	MO	AD	OCR	DT	overall	MO	AD	OCR	DT	RS
Qwen2.5-VL-7B-Instruct	39.73	30.67	27.00	74.00	70.00	49.62	31.03	31.43	90.80	85.00	39.33
+ SFT	46.40	52.00	33.00	73.00	65.00	54.41	40.44	36.29	88.80	86.00	48.00
+ SFT + GRPO	48.13	59.33	33.75	71.00	66.00	55.18	43.26	37.71	89.60	83.00	45.33
Deepeyes	50.53	46.67	40.25	78.00	70.00	56.63	43.89	38.86	90.00	84.00	51.33
Monet-SFT	50.53	48.67	43.00	65.00	69.00	54.06	41.69	44.57	76.40	71.00	54.00
Monet-7B (Monet-SFT + VLPO)	51.07	46.00	41.50	73.00	75.00	58.34	41.07	48.86	85.60	84.00	54.67

Table 7. Hyperparameters for SFT.

Hyperparameter	Value
learning rate	0.00001
batch size	1
gradient accumulation steps	16
weight decay	0.01
SFT Stage 1 steps	3885 (3 epochs)
SFT Stage 2 steps	1000 (1.03 epoch)
SFT Stage 3 steps	1000 (1.03 epoch)
latent size	{8, 10, 12}
SFT Stage 1 max total pixels	2000×28×28
SFT Stage 2 max total pixels	1500×28×28
SFT Stage 3 max total pixels	2000×28×28
SFT Stage 1,2,3 max pixels per img	2000×28×28
SFT Stage 2 alignment weight	2.0
SFT Stage 3 alignment weight	2.0

Table 8. Hyperparameters for RL.

Hyperparameter	Value
learning rate	0.000001
batch size	64
weight decay	0.01
rollout size	8
temperature	0.5
max response length	4096
VLPO σ	10.0
SFT Stage 1 steps	3885 (3 epochs)
SFT Stage 2 steps	1000 (1.03 epoch)
SFT Stage 3 steps	1000 (1.03 epoch)
latent size	10
max pixels per img	2000×28×28
accuracy threshold	0.6

To ensure reliable evaluation, we first apply a rule-based judge (exact match) and then employ DeepSeek-V3.1 or Gemini-2.5-Pro as secondary judges.

D. Case Studies

D.1. Monet-SFT-125K Examples

We select three examples from Monet-SFT-125K to illustrate the data structure and corresponding visual operations. These examples cover three representative types of visual operations: cropping key regions, generating new visual states, and drawing auxiliary lines or bounding boxes. As shown in Figures 6, 7, and 8, each training instance contains an image–text interleaved CoT with intermediate auxiliary images that facilitate the visual reasoning process. Additionally, tokens that encode key observations from the auxiliary images are wrapped in `<observation>` and `</observation>`, which are used in the alignment loss during SFT Stage 2.

D.2. Inference Examples

We present some inference examples generated by Monet-7B in this section to show the generality of tasks that Monet-7B can solve. For clarity, we don't map the generated latent embeddings to the language space as this will produce meaningless tokens. Instead, we use `<latent><latent_embeddings></latent>` to represent the latent reasoning segments.

The tasks of the selected examples include 3D Spatial Reasoning (Figure 9), 2D transformation (Figure 10), complex diagram reasoning (Figure 11), commonsense QA (Figure 12), fine-grained OCR (Figure 13), and math reasoning (Figure 14). Note that Monet not always uses latent thinking. For example, for the pure-text math problem in Figure 14, and Monet directly rely text-based reasoning to solve the problem.

Monet-SFT-125K Example 1: Cropping the Crucial Region (from Zebra-CoT Visual Search)

Question: According to the survey presented in the infographic, what percentage of parents in the U.S. delay the Varicella vaccine for their children beyond the CDC's recommendations?

Put your final answer within `\boxed{ }.`

Question Image:

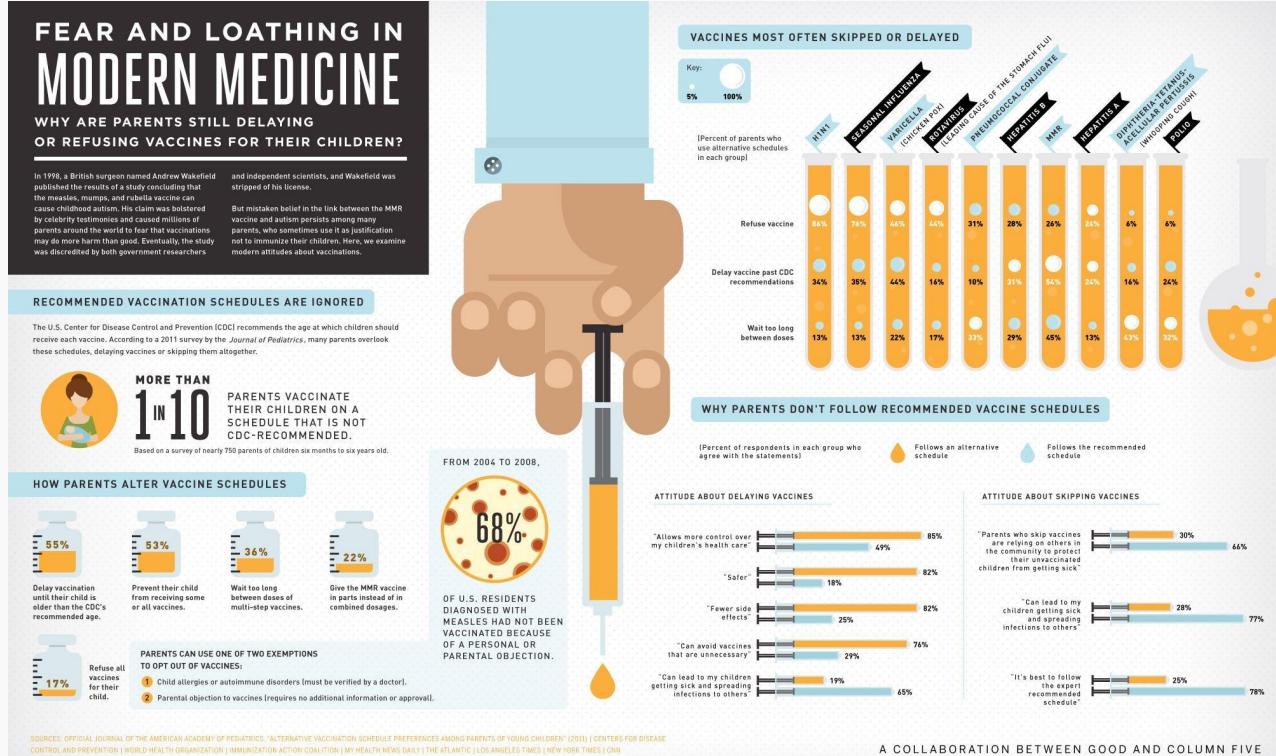


Image-text Interleaved CoT:

The question asks for the percentage of parents who delay the Varicella vaccine according to CDC recommendations, based on the provided infographic. I need to find the section of the infographic that lists vaccines and the percentages for delaying them past CDC recommendations. I will focus on the Varicella vaccine data.

The cropped image shows the specific data point for the `<observation>Varicella vaccine</observation>` under the `<delay vaccine past CDC recommendations>` row. Therefore, the final answer is `\boxed{44\%}`.

Figure 6. Example of Monet-SFT-125K: cropping the crucial region.

Monet-SFT-125K Example 2: Creating New Visual States (from Zebra-CoT 3D Object Counting)

Question: Consider the initial arrangement of objects presented in the image below. You are required to perform a series of removal operations in the specified order:

Step 1: Eliminate all objects identified as PYRAMIDS.

Step 2: Eliminate all objects that are BLUE in color.

Step 3: Eliminate all objects that are RED in color.

After completing these steps, determine the total count of orange objects remaining.

Put your final answer within `\boxed{ }.`"

Question Image:

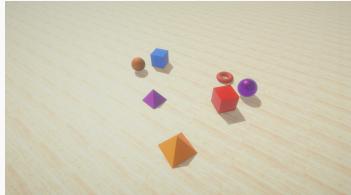
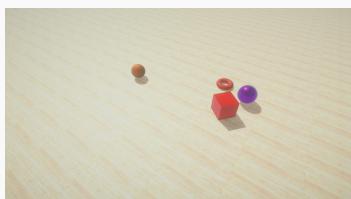


Image-text Interleaved CoT:

I will begin by analyzing the initial set of objects shown in the problem image. There are 7 distinct objects in the initial configuration. My approach will be to systematically apply each removal step and visually track the remaining objects before arriving at the final count. The first operation requires the removal of all pyramid-shaped objects from the scene. This means both the purple pyramid and the orange pyramid will be taken away. To visualize the updated set of objects, I will generate an image reflecting this change.



After removing the pyramids, as seen in the previous visual, we are left with 5 objects: an orange sphere, a blue cube, a red torus, a red cube, and a purple sphere. The next instruction is to remove all objects that are blue. This specifically targets the blue cube. I will generate a new visual to show the state after this removal.



...

Having completed all specified removal operations, the final visual shows 2 objects remaining: an orange sphere and a purple sphere. The question asks for the count of orange objects. From the two remaining objects, only one is orange. The final answer is `1`.

Figure 7. Example of Monet-SFT-125K: creating new visual states.

Monet-SFT-125K Example 3: Drawing Auxiliary Lines and Bounding Boxes (from CogCoM)

Question: What is the sum of real users below 1 million?

Put your final answer within `\boxed{ }.`"

Question Image:

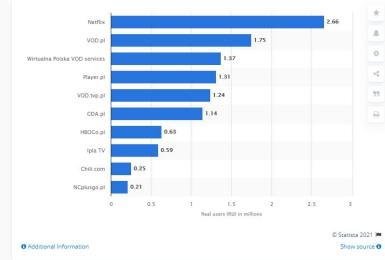
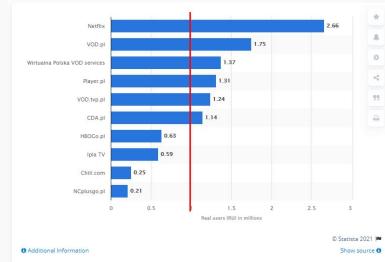
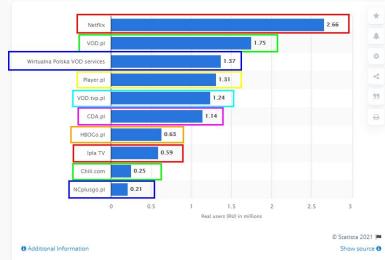


Image-text Interleaved CoT:

Draw a line to obtain the new image after the line is drawn.



Outline each bar column, which are located at Netflix, VoD.pI, Wirtualna Polska VOD services, Player.pl, VoD.tvp.pI, CDA.pI, HBOGo.pI, lpla TV, Chili.com, NCplusgo.pI.



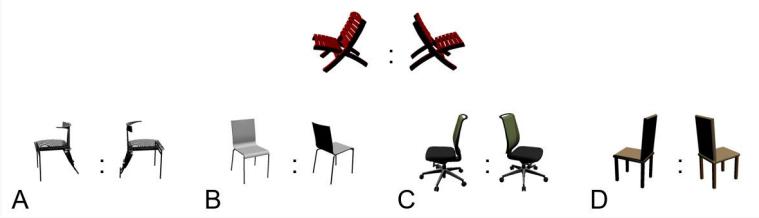
Analyze the value of each bar in relation to $x=1$. The values obtained from HBOGo.pI, Ipla TV, Chili.com, and NCplusgo.pI are all less than one million, with respective values of 0.63, 0.59, 0.25, and 0.21. Calculate their total sum as $0.63 + 0.59 + 0.25 + 0.21 = 1.68$. The final answer is $\boxed{1.68}$.

Figure 8. Example of Monet-SFT-125K: drawing auxiliary lines and bounding boxes.

Monet-7B Inference Example 1: 3D Spatial Reasoning (from VisualPuzzles)

Question: Given a pair of chairs arranged at a specific angle, select the option from the four provided that matches the angle of the chairs in the question. Options: Choose from (A) (B) (C) (D) in the image. ”

Question Image:



Response by Monet-7B:

The problem asks to identify the option that matches the given angle between the chairs. I will analyze the provided image to determine the angle between the chairs and then compare it with the given choices. To accurately assess the angle between the chairs, I will focus on the visual representation of the chairs in the problem image.
The zoomed-in view of the chairs in the problem image clearly shows the angle between them. This visual confirms the specific spatial relationship between the two chairs. Therefore, the final answer is A.

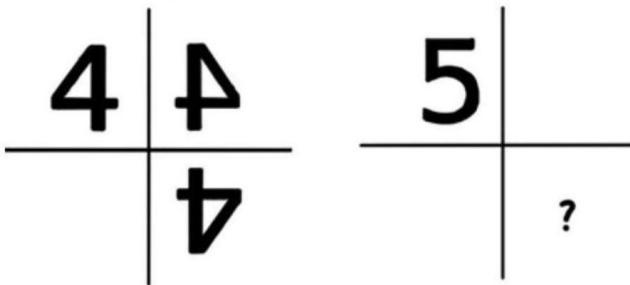
Figure 9. **Inference example: 3D spatial reasoning.** Instead of describing the angles between the chairs in language, Monet directly reasons with latent embeddings before giving the final answer.

Monet-7B Inference Example 2: 2D Transformation (from VisualPuzzles)

Question: The number 4 is reflected twice in <image1>. Which among the options in <image2> should appear in the field with the question mark if we do the same with the number 5? ”

Question Image:

<image1>:



<image2>:

- A) ፩
- B) ፪
- C) ፫
- D) ፬

Response by Monet-7B:

To determine the correct answer, we need to identify the pattern of reflection for the number 5, similar to how the number 4 is reflected in the provided image. To better understand the reflection pattern, I will generate a zoomed-in view of the relevant section of the image. <latent><latent_embeddings></latent>The zoomed-in image clearly shows how number 4 is reflected. The reflection pattern for the number 4 is obtained by flipping it twice. Once horizontally, once vertically. Therefore, the final answer is C.

Figure 10. **Inference example: 2D transformation.** By generating latent embeddings, Monet successfully identifies the flipping rule of the number.

Monet-7B Inference Example 3: Complex Diagram Reasoning (from MME-RealWorld-Lite)

- Question:** Which country is the top sales country? The choices are listed below:
- Canada
 - Germany
 - UK
 - United States
 - The image does not feature the number.

Question Image:



Response by Monet-7B:

To determine the top sales country, I need to locate the section of the image that displays sales by country. *<latent><latent_embeddings></latent>* The highlighted area in the image clearly shows the "Top Sales Countries" section. It lists various countries along with their respective sales figures. The sales figures are in millions of dollars. From the highlighted section, I can see that the United States has the highest sales at 1,241.60M, followed by the UK, Germany, and other countries. This directly answers the question about the top sales country. The final answer is **D**.

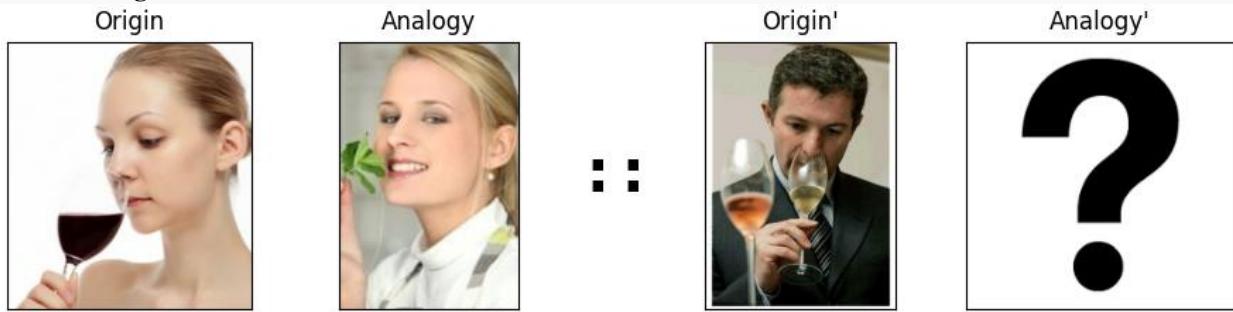
Figure 11. **Inference example: complex diagram reasoning.** Monet-7B exhibit a **hierarchical** reasoning pattern. Firstlt, it focus on the relevant section in the image by generating latent embeddings: "*The highlighted area in the image clearly shows the "Top Sales Countries" section.*" Then, it accurately identifies the contents in the "Top Sales Countries" section and gives the correct answer.

Monet-7B Inference Example 4: Commonsense QA (from VisualPuzzles)

Question: Based on the relationship between Origin and Analogy, select the option from the candidates that best completes this analogy.

Options: Choose from (A) (B) (C) (D) in the image.

Question Image:



Candidates



Response by Monet-7B:

To identify the correct image that completes the analogy, I will analyze the relationship between the "Origin" and "Analogy" images in the provided problem image. The goal is to find a similar relationship between the "Origin'" image and the missing "Analogy'" image. To better understand the relationship, I will focus on the "Origin" and "Analogy" images to identify the specific action being performed by the person in each image.
<latent><latent_embeddings></latent>

The zoomed-in view of the "Origin" image clearly shows a person drinking wine, while the "Analogy" image depicts a person smelling a flower. This indicates that the action involves the sense of smell. Therefore, the final answer is D.

Figure 12. **Inference example: commonsense QA.** Monet-7B correctly identifies the connection between the "Origin" and the "Analogy" image, showing its commonsense reasoning capability.

Monet-7B Inference Example 5: Fine-grained OCR (from MME-RealWorld-Lite)

Question: Where does Joseph work? The choices are listed below:

- (A) Boston.
 - (B) Miami.
 - (C) Chicago.
 - (D) New York.
 - (E) The image does not feature the related information

Question Image:



Response by Monet-7B:

To determine where Joseph works, I need to locate his name or relevant information within the article. I will zoom in on the section of the article that mentions Joseph to clearly identify his location.

The zoomed-in image clearly shows that Joseph Higgins is a portfolio manager in New York. This directly answers the question about where Joseph works. Therefore, the final answer is **D**.

Figure 13. Inference example: fine-grained OCR. Monet-7B accurately identifies the key information, which is located in the middle of the rightmost region of the image.

Monet-7B Inference Example 6: Math Reasoning (from VisualPuzzles)

Question: Every night the wizard Tilim makes the weather forecast for the king. When Tilim gets it right he gets 3 gold coins, but when he makes a mistake, he pays a fine of 2 gold coins. After making the prediction for 5 days, Tilim did the math and discovered that he neither won nor lost coins. How many times did he get the weather forecast right in those 5 days?

Options:

- (A) 1
- (B) 2
- (C) 3
- (D) 4

Response by Monet-7B:

To determine the number of correct weather forecasts, we can set up an equation based on the given conditions. Let's define: - x as the number of correct forecasts (where he earns 3 gold coins per correct forecast).

- y as the number of incorrect forecasts (where he loses 2 gold coins per incorrect forecast)

The wizard made a total of 5 forecasts, so:

$$x + y = 5$$

He neither won nor lost any gold, which means the total gain from correct forecasts equals the total loss from incorrect forecasts:

$$3x = 2y$$

We now have a system of two equations:

$$1. \quad x + y = 5$$

$$2. \quad 3x = 2y$$

We can solve this system by expressing y in terms of x from the first equation:

$$y = 5 - x$$

Next, substitute $y = 5 - x$ into the second equation:

$$3x = 2(5 - x)$$

Solve for x :

$$3x = 10 - 2x$$

$$3x + 2x = 10$$

$$5x = 10$$

$$x = 2$$

So, Tilim made 2 correct forecasts. To verify, we calculate y :

$$y = 5 - x = 5 - 2 = 3$$

Now, check the total gold:

- Correct forecasts: $2 \times 3 = 6$ gold

- Incorrect forecasts: $3 \times (-2) = -6$ gold

- Total gold: $6 + (-6) = 0$ gold. The calculations confirm that Tilim neither won nor lost any gold. Therefore, the number of correct forecasts is 2. The final answer is B.

Figure 14. **Inference example: math reasoning.** Monet-7B can perform pure-text math reasoning. Since the problem is not visual-related, it doesn't activate latent thinking mode.