# Chain-of-Visual-Thought:
# Teaching VLMs to See and Think Better with Continuous Visual Tokens

**Yiming Qin**[1]   **Bomin Wei**[2]   **Jiaxin Ge**[1]   **Konstantinos Kallidromitis**[3]
**Stephanie Fu**[1]   **Trevor Darrell**[1]   **Xudong Wang**[1†]

[1]UC Berkeley   [2]UCLA   [3]Panasonic AI Research
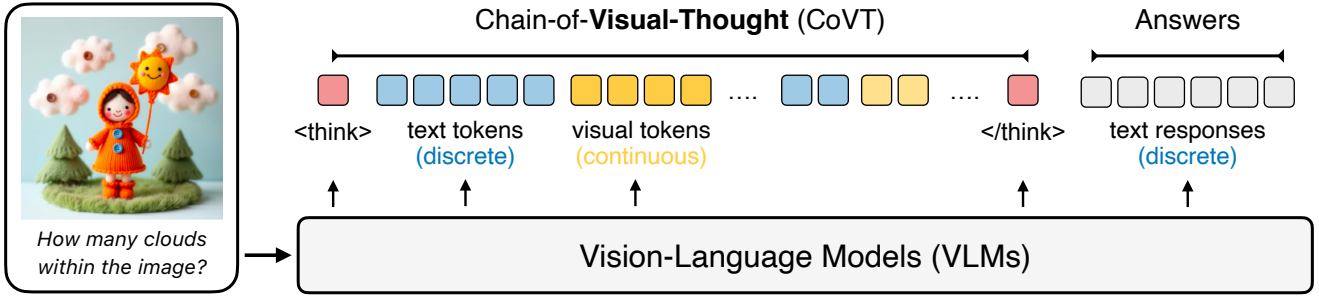
[†]Corresponding authors

Figure 1. **Rather than restricting VLM reasoning to the discrete language space with limited representational capacity, CoVT forms a visual thought chain that enables VLMs to reason in continuous visual space.** By introducing *continuous visual tokens* that encode perceptual cues (*e.g.*, segmentation, depth, instance, and edge structure), CoVT composes *chains of textual and visual thoughts* that link semantic reasoning with perceptual grounding. These visual "thought chains" bridge language and vision, enabling fine-grained understanding, spatial precision, and geometric awareness beyond the reach of text-based reasoning.

## Abstract

*Vision–Language Models (VLMs) excel at reasoning in linguistic space but struggle with perceptual understanding that requires dense visual perception,* e.g., *spatial reasoning and geometric awareness. This limitation stems from the fact that current VLMs have limited mechanisms to capture dense visual information across spatial dimensions. We introduce **Chain-of-Visual-Thought** (CoVT), a framework that enables VLMs to reason not only in words but also through continuous visual tokens—compact latent representations that encode rich perceptual cues. Within a small budget of roughly 20 tokens, CoVT distills knowledge from lightweight vision experts capturing complementary properties such as 2D appearance, 3D geometry, spatial layout, and edge structure. During training, the VLM with CoVT autoregressively predicts these visual tokens to reconstruct dense supervision signals (e.g., depth, segmentation, edges, and DINO features). At inference, the model reasons directly in the continuous visual token space, preserving efficiency while optionally decoding dense predictions for interpretability. Evaluated across more than ten diverse perception benchmarks, including CV-Bench, MMVP, RealWorldQA, MMStar, WorldMedQA, and HRBench, integrating CoVT into strong VLMs such as Qwen2.5-VL and LLaVA consistently improves performance by 3% to 16% and demonstrates that compact continuous visual thinking enables more precise, grounded, and interpretable multimodal intelligence. Our website is available at https://wakalsprojectpage.github.io/comt-website.*

## 1. Introduction

Vision–Language Models (VLMs) [2, 3, 13, 16, 31, 44, 45, 60, 67, 73] have become the cornerstone of modern multimodal intelligence, achieving remarkable progress in understanding and reasoning across text and vision. By projecting visual input into a language-centric token space, VLMs inherit the strong compositional and logical reasoning capabilities of large language models (LLMs), enabling unified multimodal interaction through natural language. Recent advances in text-based Chain-of-Thought (CoT) reasoning [54] further extend this paradigm, showing that structured intermediate reasoning steps can significantly enhance performance on tasks involving logic, mathematics, and knowledge grounding. However, despite these successes,
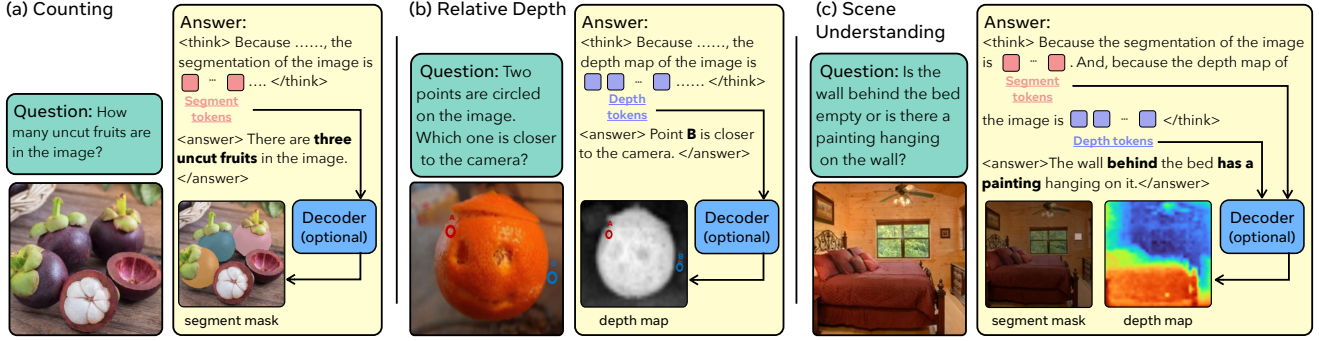
Figure 2. **Continuous visual thinking with CoVT.** CoVT introduces compact, continuous visual tokens that encode fine-grained perceptual cues, such as object localization, spatial structure, and scene semantics, directly into VLM reasoning. These tokens ground multimodal reasoning in visual space, enabling the model to capture fine-grained relationships across vision-centric tasks (*e.g.*, counting, depth ordering, and scene understanding) without relying on external tools. They can also be decoded into dense predictions, offering human-interpretable visualizations of the model's reasoning process.

such reasoning remains fundamentally *language-bound*.

When continuous visual information is projected into discrete text space, *rich perceptual cues*, e.g., *boundaries, layout, depth, and geometry, are lost or poorly represented*. Yet these are precisely the fine-grained signals that humans rely on when reasoning about the visual world. Consequently, current VLMs often struggle with perception-intensive tasks such as counting, spatial correspondence, or relative depth estimation, even when equipped with powerful vision encoders [18, 39, 53], as shown in Fig. 2. Moreover, by forcing vision reasoning through a discrete text bottleneck, the model must verbalize continuous spatial and geometric relations. As a result, text-only CoT can misdirect and even *degrade* visual reasoning performance, as shown by Qwen3-VL-Thinking [3, 61], which performs over 5% worse than Qwen3-VL-Instruct with language CoT on spatial understanding benchmarks such as V [58], HRBench8k [50], and VSI-Bench [62]. This exposes a fundamental limitation: *visual information is inherently continuous and high-dimensional, yet existing models reason over it using symbolic language tokens that lack the fidelity of complex perceptual reasoning.*

A natural solution is to augment VLMs with external vision tools [22, 43], leveraging pre-built specialized models to recover fine-grained perception. While this approach can partially restore spatial and geometric information, it also introduces significant drawbacks: perception is delegated to external tools, and outcomes are bounded by them. It also introduces higher GPU cost. Another solution is generating or cropping images in the thinking process. However, these solutions still project the images into the text space, losing the dense visual information. These limitations motivate a central question: *Can VLMs learn to reason the way humans do, by thinking visually rather than translating everything into words?* More concretely, can we inject fine-grained visual signals directly into a VLM's reasoning process, allowing it to "see" and "think" simultaneously while remaining efficient and self-contained? Yes! We propose **C**hain-**o**f-**V**isual-**T**houghts (**CoVT**).

CoVT enables reasoning over rich perceptual cues by grounding VLMs in continuous visual token space. Each group of visual tokens corresponds to a lightweight perceptual expert (*e.g.*, segmentation, depth, edge detection, or self-supervised representation learning) that encodes specific visual features. During training, the VLM is asked to *predict* these continuous visual tokens within its reasoning chain, compressing rich perceptual information into a compact latent space. These latent tokens are then decoded by task-specific lightweight decoders to reconstruct the corresponding expert targets (*e.g.*, segmentation masks, dense depth maps, edge maps, or DINO features). We backpropagate the reconstruction and distillation losses through the continuous tokens, aligning the model's internal latent representations with expert guidance. This process allows CoVT to internalize fine-grained perceptual knowledge directly into its token space, enabling grounded reasoning without explicit visual maps or external tool calls.

More specifically, we highlight different aspects of fine-grained visual reasoning. We integrate both task-oriented experts (*e.g.*, SAM [27], DepthAnything v2 [63], PIDINet [42]) and representation-based experts (*e.g.*, DINO [7], contrastive encoders), with alignment strategies tailored to each: task-oriented signals are aligned at the prompt level, while representation-based signals are aligned in feature space. Training proceeds through four stages, *including comprehension, generation, reasoning, and efficient reasoning*, gradually teaching the model to reason effectively with visual thoughts.

At inference, the model forms *chains of visual thoughts*, reasoning across modalities to produce answers that are

2

both semantically coherent and perceptually grounded. This self-contained, differentiable process enables VLMs to "think" directly in continuous visual space, thereby providing a more faithful bridge between internal reasoning and perceptual understanding. Moreover, this design supports interpretable multimodal intelligence, allowing users to visualize the model's visual thinking process when desired. If visualization is not required, COVT can operate solely on the continuous visual tokens without decoding them into dense predictions, thus maintaining efficiency.

Evaluated across diverse perception benchmark, COVT consistently improves fine-grained visual reasoning, outperforming strong VLM baselines on vision-centric tasks while maintaining competitive performance on general (non-vision-centric) benchmarks. For example, COVT achieves a 5.5% overall gain on CV-Bench [46], delivering a substantial 14.0% improvement on its depth sub-task, and 4.5% overall gain on HRBench [49]. In addition, COVT offers flexible interpretability: the continuous visual tokens can be decoded into human-readable dense predictions, providing a window into the model's underlying visual reasoning process when desired. Together, these results demonstrate that compact continuous visual thinking enables more precise, grounded, and interpretable multimodal intelligence.

***The main elements of our contribution are as follows:***

- We propose Chain-of-Visual-Thought, a framework that equips VLMs with the ability to reason through *continuous visual tokens*, compact perceptual representations that serve as the building blocks for multimodal thinking.
- We develop tailored alignment strategies and a training pipeline (*comprehension, generation, reasoning, and efficient reasoning*) that enable VLMs to learn, interpret, and reason effectively within continuous visual space.
- We demonstrate consistent performance gains across diverse benchmarks, showing that continuous visual tokens enhance both perceptual grounding and interpretability.

## 2. Related Work

**Tool-Augmented Reasoning** Equipping VLMs with external tools enables them to use specialized vision models for targeted visual tasks [22, 34, 43, 56, 65]. While this improves performance, it also introduces computational overhead. Moreover, tool usage is inherently constrained, as the final performance is bounded by the ability of each tool instead of the reasoning process itself. In this work, we consider self-contained visual reasoning, which conduct reasoning flexibly and does not rely on external vision tools.

**Text Space Reasoning** Text space reasoning methods, such as Chain-of-Thought [28], has achieved big success in language reasoning [35, 51, 52, 55], solving problems like math, science, and logical reasoning. The strong performance of LLMs with CoT capabilities has led to its broad adoption and success in models such as DeepSeek-R1 [12].

| Desired Properties | VCoT | MCoT | VChain | Aurora | Ours |
|---|---|---|---|---|---|
| Operates without relying on external tools | ✓ | ✗ | ✓ | ✓ | ✓ |
| Reasons in the continuous visual space | ✗ | ✓ | ✓ | ✗ | ✓ |
| Leverages dense visual information for reasoning | ✗ | ✓ | ✗ | ✓ | ✓ |
| Has 3D-aware perception | ✗ | ✗ | ✗ | ✓ | ✓ |

Table 1. **Comparison of key properties with prior multimodal reasoning methods**. Unlike prior methods such as VCoT [41], MCoT [70], VChain [25], and Aurora[5], COVT uniquely satisfies all desired properties: it reasons in continuous visual space, leverages dense visual cues, maintains 3D awareness, and operates fully without external tools. Desired and undesired properties are shown in green and magenta, respectively.

With the success of text CoT, many works have extended reasoning to the visual modality. A straightforward approach is to generate dense captions and then reason in language space [33], but this process is inherently lossy.

We compare COVT with recent multimodal reasoning paradigms in Tab. 1. Visual CoT [41] relies on textual interpretations of images, limiting reasoning to the discrete text space. MCoT [11] enables continuous visual reasoning by editing or generating supplementary images, but requires substantial compute and lacks flexibility. VChain [26] interleaves images and text in the reasoning chain, yet still loses visual information by projecting images into text space. COVT uniquely combines continuous visual reasoning, dense perceptual cues, and 3D-aware understanding within a single self-contained framework.

**Latent Space Reasoning** Concurrent work shows that reasoning in latent space can strengthen LLMs in complex, multi-step tasks [6, 8]. Coconut [23] finds that continuous latent embeddings are more efficient than explicit CoT, while CCoT [10] compresses CoT into continuous tokens for denser reasoning. Other studies explore specialized reasoning tokens [20] or use hidden states as implicit reasoning paths [14]. Latent reasoning has also been extended to VLMs. Aurora [4] employs VQ-VAE latents of depth and detection signals to enhance depth estimation and counting, whereas Mirage [64] uses latent imagination for visual reasoning tasks. Our work, COVT, builds upon these previous contributions: we introduce a form of tool-use directly embedded in continuous latent space, where the implicit 'tools' are visual thinking tokens tied to specific perceptual experts.

## 3. Chain-of-Visual-Thought (COVT)

We first introduce the preamble in Sec. 3.1. We then show the overall pipeline of COVT in Sec. 3.2. We also discuss how we select the visual token categories and explain how different visual tokens are aligned (Sec. 3.3). Finally, we present the model training pipeline, *e.g.*, the training loss formulation and the data framework design in (Sec. 3.4).
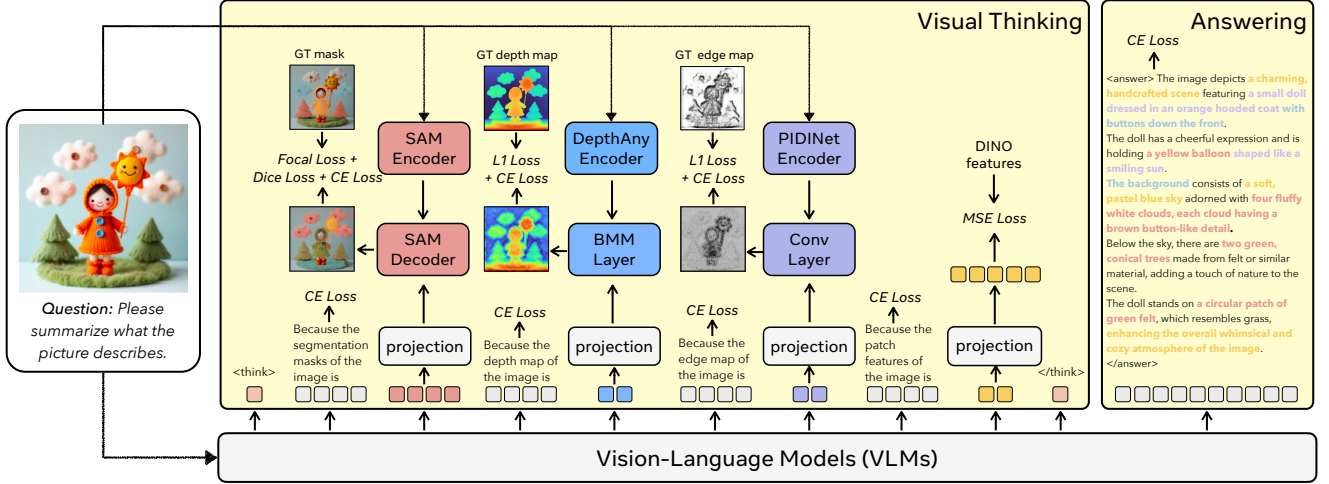
Figure 3. **The training pipeline of CoVT.** CoVT first generates the thinking process, containing visual thinking tokens, and then leverages these visual thoughts to condition next token prediction and reason the final answer. To endow these tokens with perceptual meaning, we align them with lightweight vision experts (*e.g.*, SAM, DepthAnything, PIDINet, DINO) on their respective tasks during training. Specifically: SAM uses 8 visual tokens as mask prompts; DepthAnything uses 4 tokens to reconstruct depth; PIDINet uses 4 tokens to reconstruct edges; and DINO uses 4 tokens to match patch-level features. The VLM is finetuned with LoRA and all the projection layers are trainable. *Note: During inference, dense predictions are decoded only when interpretability is desired; otherwise, reasoning occurs entirely in the latent visual space.*

## 3.1. Preamble

Existing VLMs face two key limitations in fine-grained visual reasoning. **1)** *Text-only CoT accumulates errors.* Text-only CoT executes a long chain of thought, which may generate errors at the early stage. These mistakes will accumulate and ultimately lead to an incorrect final result. Therefore, we need a reasoning that is short and effective. **2)** *Supervision is dominated by text responses,* which provides little incentive for the model to capture *low-level perceptual cues* such as edges, depth, or regions. We need to equip VLMs themselves with the capability of extracting fine-grained visual information from the image, which can be further decoded by vision decoders.

CoVT intends to provides a foundation for the next generation of multimodal reasoning systems, capable of thinking fluidly across both language and vision in a self-contained, interpretable manner.

## 3.2. CoVT **Overall Pipeline**

We propose CoVT, a framework that augments VLMs with *chains of visual thoughts*. Fig. 3 illustrates the overview of CoVT pipeline. Essentially, this framework equips VLMs with the capability of outputting fine-grained visual representations within a continuous visual token space, enabling them to reason directly over rich perceptual information and maintain spatial and geometric coherence throughout the reasoning process.

At its core, CoVT retains the standard next-token pre-

diction paradigm. For standard VLMs, given visual features $\mathcal{V}$ extracted from a frozen vision encoder and text features $\mathcal{T}$ from a language encoder, the VLM estimates the probability of generating a sequence $Y = (y_1, y_2, \ldots, y_n)$ as:

$$P(Y \mid \mathcal{V}, \mathcal{T}; \theta) = \prod_{i=1}^{n} P\left(y_i \mid y_{<i}, \mathcal{V}, \mathcal{T}\right). \quad (1)$$

As shown in Fig. 3, CoVT extends this formulation by introducing *Chain-of-Visual-Thought tokens*, where each token $y_i$ can represent either a visual token or a text token.

To effectively incorporate CoVT tokens into the VLM, we train the model to function as a dense visual encoder capable of generating multiple visual tokens that capture diverse fine-grained perceptual cues. The VLM is trained to generate CoVT tokens that, through task-specific decoders, reconstruct visual outputs under reconstruction supervision. Through this process, CoVT evolves to generate rich, fine-grained visual information across multiple perceptual dimensions within a thinking chain.

## 3.3. CoVT **Tokens**

**Token selection based on core perception ability.** As proposed in [72], the vision-centric perceptual ability of VLMs can be summarized as *(i) instance recognition*, *(ii) 2D and 3D spatial relationships*, *(iii) structure detection*, and *(iv) deep mining of semantic information*. Based on this categorization, we use four vision models to supervise CoVT tokens to learn each ability: *1) Segmentation tokens pro-*
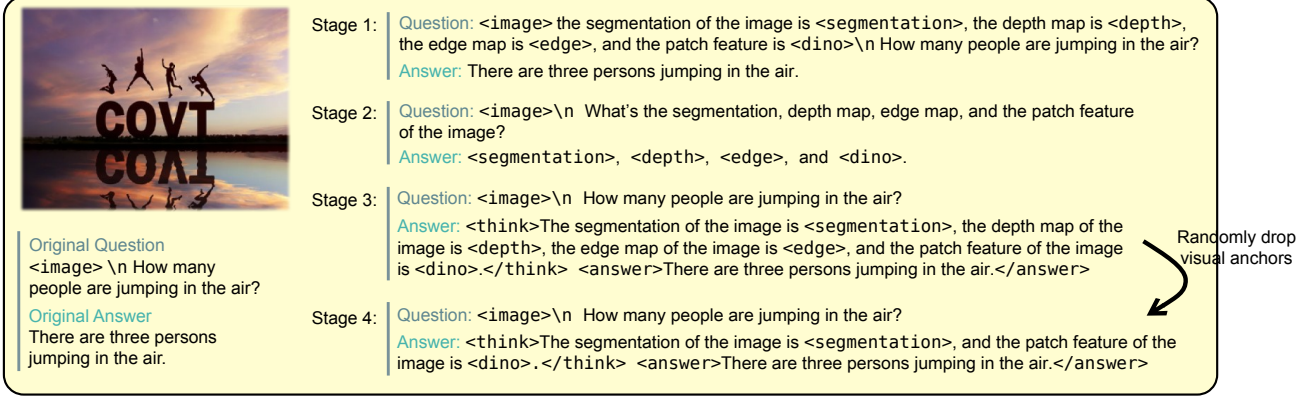
**Figure 4. Four-stage data formatting for CoVT.** The first stage helps the model comprehend the visual tokens, and the second stage guides it to generate them. The third stage enables the VLM to integrate visual tokens into its reasoning process, while the final stage allows the model to efficiently select and utilize visual thinking tokens within visual thought chains.

*vide instance-level position and shape information*, which endow VLMs with the instance recognition signals and 2D spatial perception. *2) Depth tokens provide pixel-level depth information*, equipping VLMs with the capability of figuring out 3D spatial relationships. *3) Edge tokens provide geometry-level details*, which assist models to detect structural cues and partially provide 2D spatial information. *4) DINO tokens provide the patch-level representation of the images*, delivering rich semantic information.

**Tokens alignment based on granularity of visual models.** Task-oriented models and representative models produce outputs at different levels of granularity. In general, task-oriented models tend to be more fine-grained, whereas representative models are usually less fine-grained. We adopt different strategies to align each type of token with the visual models based on different granularities. Essentially, we adopt two main alignment methods: For fine-grained task-oriented models, visual tokens are projected to the prompt space and then aligned at the prompt level with the decoders, while for representative models, alignment with the encoders is applied at the feature level after the projection. The projection layer consists of one multi-head attention layer and two full connected layers.

*1) Segmentation tokens are supervised by SAM* [27], which is a *task-oriented model* that contains dense visual features. Therefore, following LISA [29], we align Segmentation tokens with the SAM decoder. The 8 Segmentation tokens are aligned at the prompt level, and each token prompts one mask, formulated as:

$$\hat{M}_i = \text{Decoder}(T_i^{\text{sam}}, f), \quad \hat{M}_i \in [0,1]^{H \times W}, \quad (2)$$

where $\hat{M}_i$ is the $i$th decoded mask, $T_i^{\text{sam}}$ denotes the $i$th segmentation token predicted in CoVT, serving as the prompt fed into the SAM decoder, and $f$ means the dense embedding from the SAM encoder. During the training process,

the Hungarian matching algorithm is employed to match the predicted masks with the ground truths, while dice loss and focal loss are applied.

*2) Depth tokens are supervised by DepthAnything v2* [63]–*a task-oriented model*. Since these tokens contain dense information, they are also aligned with the decoder at the prompt level. We use 4 Depth tokens to serve as 4 prompts to interact with the dense features extracted by DepthAnythingv2 through batch matrix multiplication (BMM) to reconstruct the depth map, formulated as:

$$\hat{D}_i = \text{softmax}\left(T_i^{\text{depth}} \cdot F_i^{\text{depth}\top}\right), \quad (3)$$

where $\hat{D}_i$ denotes the $i$th reconstructed depth map, $T_i^{\text{depth}}$ represents the $i$th depth visual token, and $F_i^{\text{depth}}$ is the $i$th middle layer feature from DepthAnything encoder. The final depth map is $\hat{D} = \frac{\sum_{i=0}^{3} \hat{D}_i}{4}$. The L1 reconstruction loss is employed for aligning Depth tokens.

*3) Edge tokens are aligned with PIDINet* [42]. Each of 4 Edge tokens functions as an $1 \times 1$ convolutional kernel applied to the dense features from PIDINet encoder to reconstruct the $i$th edge map $\hat{E}_i$. The final edge map is $\hat{E} = \frac{\sum_{i=0}^{3} \hat{E}_i}{4}$, and aligned via L1 loss function.

*4) DINO tokens are supervised by DINOv2* [37], which is trained as the *representative model*, extracting patch-level features. Therefore, the 4 DINO tokens are mapped into the same shape with DINO feature using the projection layer, and aligned under an MSE objective.

### 3.4. CoVT Training

**Training Loss.** During training, the *joint loss function* is defined as:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{ce}} + \gamma\big(\lambda_{\text{seg}} \cdot \mathcal{L}_{\text{visual}}^{\text{seg}} + \lambda_{\text{depth}} \cdot \mathcal{L}_{\text{visual}}^{\text{depth}} \\ + \lambda_{\text{edge}} \cdot \mathcal{L}_{\text{visual}}^{\text{edge}} + \lambda_{\text{dino}} \cdot \mathcal{L}_{\text{visual}}^{\text{dino}}\big), \quad (4)$$

5

| Visual tokens | | | | CV-Bench | | | | Other vision-centric benchmarks | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Seg | Depth | DINO | Edge | CVBench | Count | Depth | Dist. | BLINK | RW-QA | MMT | MMStar-P | MMVP | MME-RW | V* | HR$_{4K}$ | HR$_{8K}$ |
| **Closed-source Models** | | | | | | | | | | | | | | | | |
| Claude-4-Sonnet | | | | 76.3 | 62.2 | 77.7 | 80.5 | 39.6 | 63.7 | - | 58.8 | 48.7 | - | 15.2 | 32.3 | 22.7 |
| GPT-4o | | | | 79.2 | 65.6 | 86.7 | 81.0 | 63.0 | 69.7 | - | 65.2 | 72.0 | - | 42.9 | 50.6 | 46.7 |
| Qwen2.5-VL-7B | | | | 74.5 | 65.0 | 72.8 | 75.5 | 55.7 | 68.6 | 61.7 | 67.1 | 56.0 | 60.0 | 76.4 | 68.6 | 64.9 |
| **CoVT (1 Visual Token)** | | | | | | | | | | | | | | | | |
| ✓ | | | | 77.9 | 66.0 | 80.8 | 80.5 | **57.4** | 71.1 | 62.1 | 68.5 | 58.7 | 62.1 | **79.1** | 71.9 | 69.0 |
| | ✓ | | | 78.7 | 65.4 | 83.2 | 78.2 | 56.4 | 71.5 | **62.7** | **69.9** | 58.7 | 62.0 | **79.1** | 71.9 | 69.4 |
| | | ✓ | | 71.3 | 64.7 | 72.3 | 66.7 | 55.8 | 71.5 | 62.5 | 67.9 | 57.3 | 61.1 | 77.5 | 71.0 | 68.6 |
| **CoVT (3 Visual Tokens)** | | | | | | | | | | | | | | | | |
| ✓ | ✓ | ✓ | | **80.0** | **66.2** | 86.8 | **82.5** | 56.0 | 71.6 | 62.1 | 69.2 | **58.7** | **63.7** | 78.0 | **72.9** | 69.4 |
| Δ *(vs Baseline)* | | | | +5.5 | +1.2 | +14.0 | +7.0 | +0.3 | +3.0 | +0.4 | +2.1 | +2.7 | +3.7 | +1.6 | +4.3 | +4.5 |
| **CoVT (4 Visual Tokens)** | | | | | | | | | | | | | | | | |
| ✓ | ✓ | ✓ | ✓ | 79.8 | 66.1 | **89.2** | 80.5 | 56.2 | **71.8** | 61.9 | 68.4 | 56.7 | 63.3 | 78.5 | 72.5 | **69.9** |
| Δ *(vs Baseline)* | | | | +5.3 | +1.1 | +16.4 | +5.0 | +0.5 | +3.2 | +0.2 | +1.3 | +0.7 | +3.3 | +2.1 | +3.9 | +5.0 |

Table 2. **Comparison of COVT with the baseline and closed-source models.** COVT delivers consistent improvements across all vision-centric benchmarks and further reveals that each visual token type contributes most effectively to the tasks related to its rich information.

where $\mathcal{L}_{ce}$ is the typical cross-entropy loss for VLMs, $\gamma$ is the coefficient of visual loss, and all of the $\lambda$ coefficients are the weighting factors for the losses of the corresponding visual tasks. During the inference process, the visual thinking tokens are not decoded.

Additionally, our framework supports the flexible integration of new visual token types. Because the pipeline follows a clean next-token prediction paradigm, additional tokens can be incorporated with minimal modification.

**Training Data.** To enable VLMs to progressively learn the visual tokens while not losing ability in the text space, COVT introduces four data formatting stages as shown in Fig. 4. This guides the VLMs to learn progressively through the sequence from understanding visual tokens *(comprehension stage)*, to generating visual tokens *(generation stage)*, to reasoning with chain of visual thoughts *(reasoning stage)*, and finally to dynamically using visual tokens in the thinking chain *(efficient reasoning stage)*.

In *1) comprehension stage*, we insert visual tokens after `<image>` to teach the VLMs to learn the basic semantics of the visual tokens. In *2) generation stage*, we modify the question and answer, as shown in Fig. 4, to guide the VLMs to generate the visual tokens precisely. *3) reasoning stage* introduces the chain-of-visual-thought format, where the visual tokens are used within the thinking process. This teaches the model to leverage the visual tokens to derive the final answers. *4) efficient reasoning stage* randomly drops out some sets (ranging from 0 to $k$, where $k$ is the number of token types) of visual tokens. With a portion of visual token types, COVT learns to utilize all features effectively rather than being constrained by a fixed output pattern.

The dataset used for training include: (1) a vision-centric (and also real-world) subset of the LLaVA-OneVision dataset [30]. (2) spatial perception data, including TallyQA [1] and ADE20K-Depth [4, 71].

| | CV-Bench | | | BLINK | | | | |
|---|---|---|---|---|---|---|---|---|
| | Count | Depth | Dist. | Count | Obj. Loc. | Rel. Depth | Vis. Corr. | Vis. Sim. |
| LLaVA | 59.3 | 61.8 | 50.2 | 56.7 | 54.9 | 52.4 | 29.7 | 51.1 |
| Aurora[†] *(depth)* | 54.9 | 67.7 | **52.3** | 53.3 | 55.7 | 62.9 | 26.2 | 47.4 |
| CoVT *(w/ Depth)* | 60.7 | **71.0** | **52.3** | 56.7 | 59.8 | 75.8 | 31.4 | 53.3 |
| Δ *(vs Aurora)* | +5.8 | +3.3 | +0.0 | +3.4 | +4.1 | +12.9 | +5.2 | +5.9 |
| Aurora[†] *(count)* | 56.0 | 62.2 | 47.8 | 31.7 | 26.2 | 24.2 | 26.7 | 21.5 |
| CoVT *(w/ Seg)* | 61.9 | 60.7 | 51.3 | 58.3 | 56.6 | 69.4 | 29.7 | 52.6 |
| Δ *(vs Aurora)* | +5.9 | -1.5 | +3.5 | +26.6 | +30.4 | +45.2 | +3.0 | +31.1 |

Table 3. **Comparison between COVT and Aurora based on LLaVA-v1.5-13B.** [†] indicates our reproduced results based on the provided checkpoints.

## 4. Experiments

In the experiment section, we first describe the experimental settings of Chain-of-Visual-Thought (COVT) in Sec. 4.1. Second, we introduce the benchmarks results on both vision-centric and non-vision-centric datasets in Sec. 4.2. Third, we present the quantitative results demonstrating the advantages of COVT in Sec. 4.3. Finally, we "visualize" the continuous visual tokens in Sec. 4.4 and conduct ablation studies in Sec. 4.5.

### 4.1. Experiment Details

In our experiments, Qwen2.5-VL-7B [3] is selected as the main baseline, COVT uses LoRA [24] tuning method, while the rank of LoRA is 16 and the LoRA alpha is 32. The learning rate of LoRA is set as $5 \times 10^{-5}$ and the projection layer learning rate is set to $1 \times 10^{-5}$. The first phase steps are 4000, second and third phase steps are 3000, and the fourth phase steps are 5000. Batch size is set to 4. The experiments are carried out on $1 \times$A100 or $4 \times$A6000 GPUs. $\gamma$ and all of the $\lambda$ in Eq. 4 are set as 1.
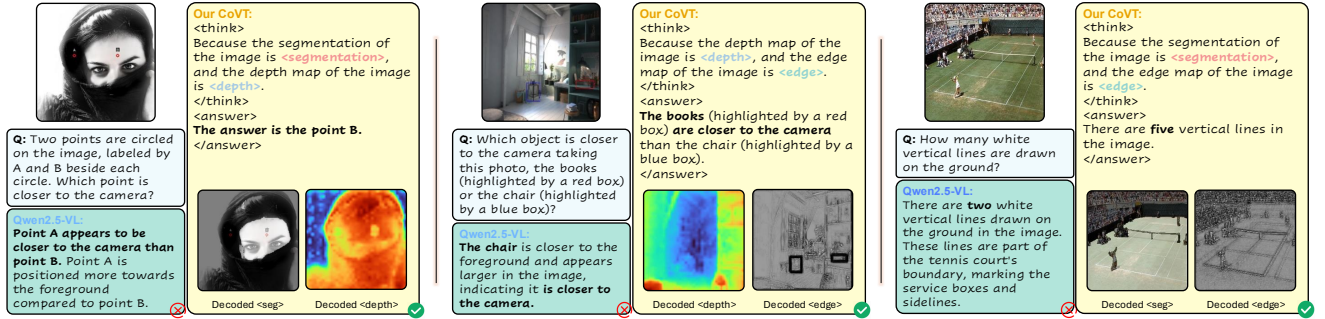
Figure 5. **Visualization of CoVT tokens.** Different visual tokens contribute complementary cues that enable the model to solve complex perceptual reasoning tasks. **Left:** Segmentation tokens localize point B on the face, while the depth tokens capture the relative depth relationships. **Mid:** Depth visual tokens provide depth map information, and the edge tokens help highlight the positions of two boxes. **Right:** The Segmentation tokens identify the attended region, and the edge tokens delineate the fine-grained line structures.

## 4.2. Model Evaluation

All evaluations are performed using VLMEvalKit [15].

**Vision-centric benchmarks.** Our main focus is on CV-Bench. In particular, from CV-Bench we highlight the sub-tasks *Count*, *Depth*, and *Distance*. These sub-tasks act as precise indicators to validate the effectiveness of our method. We further evaluate on other vision-centric benchmarks, including BLINK [19], RealWorldQA (RW-QA) [59], MMT-Bench (MMT) [66], MMStar [9], MMVP [47], MME-RealWorld (MME-RW) [69], V* Bench (V*) [57], and HRBench (HR$_{4K}$ and HR$_{8K}$) [49]. Among them, for MMStar we specifically choose the *Coarse Perception*, *Fine-grained Perception*, and *Instance Reasoning* subsets (MMStar-P), as they are more aligned with real-world reasoning.

**Non-vision-centric benchmarks.** Besides, we also evaluate CoVT on some non-vision-centric visual benchmarks such as OCRBench [32], MME [17], MUIRBench [48], HallusionBench [21], A-OKVQA [40], TaskMeAnything [68], WeMATH [38], and WorldMedQA-V [36]. For MME, we select the text-centric sub-task *text translation* as the evaluation of the text-centric performance.

## 4.3. Quantitative Results

**COVT outperforms the baseline across the vision-centric benchmarks.** As shown in Tab. 2, COVT is capable of incorporating various kinds of visual tokens. We employ three visual tokens, Segmentation, Depth, and DINO, as our main results. Compared to the baseline, COVT consistently achieve large gains across the main vision-centric benchmarks. COVT improves by 5.5% on CV-Bench, 14.0% on the subtask *depth* in CV-Bench, 3.7% on MME-RealWorld, and 4.5% on HRBench8K. These results indicate COVT with visual thinking chain improves across visual-centric and fine-grained perceptual tasks.

**COVT generalizes to the other baseline.** In addition to the experiment based on Qwen, COVT is also implemented based on LLaVA-v1.5-13B, in order to compare COVT with Aurora. As shown in Fig. 3, for COVT with depth tokens, it excels Aurora-*depth* by 12.9% on *relative-depth* in BLINK. For COVT using segmentation COVT tokens, outperforms Aurora-*count* by 26.6% on BLINK-*count* benchmark. These results indicate that COVT generalizes on various baselines across vision-centric tasks.

## 4.4. Qualitative Results

To better understand why COVT is effective, we select several examples and decode the COVT tokens from the model outputs to visualize whether these tokens provide useful information for reasoning.

Fig. 5 illustrates that different COVT tokens carry different rich fine-grained information, and the cues they provide are highly complementary. To be interpretable, COVT tokens are decoded into fine-grained output (*e.g.* masks, depth maps, edge maps). ***For the left example***, the Segmentation token provides 2D perceptual cues by localizing "point B" on the face, while the Depth token supplies 3D information, indicating that the face region is closer to the camera than the surrounding areas. ***For the middle example***, rhe Depth token encodes depth perception, whereas the Edge token supplies fine-grained boundary cues for the two target objects. This example is from *Depth* sub-task in CV-Bench, thus the synergy explains the 2.4% improvement observed when COVT uses four visual tokens instead of three on the CVBench-Depth task, as shown in Tab. 2. ***For the right example***, the Segmentation token localizes the target region and the Edge token emphasize fine-grained boundaries, which is difficult for the Segmentation token, deriving the correct answer through chains of visual thoughts.

## 4.5. Ablation Studies

**Text-only Chain-of-Thought *vs*. Chain-of-Visual-Thought.** To isolate the contribution of continuous visual tokens, we conduct an ablation comparing text-only CoT with our full COVT framework. For the text-only setting,
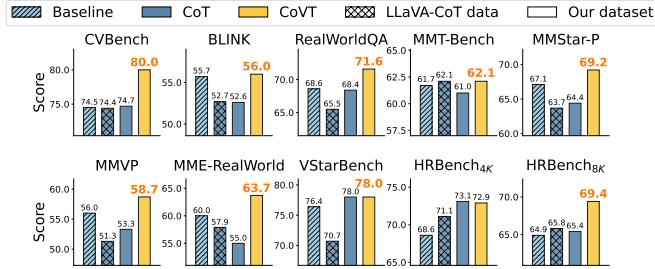
Figure 6. **Text-only CoT *vs* CoVT.** CoVT substantially enhances VLMs' capabilities on vision-centric tasks, whereas text-only CoT can even degrade performance.
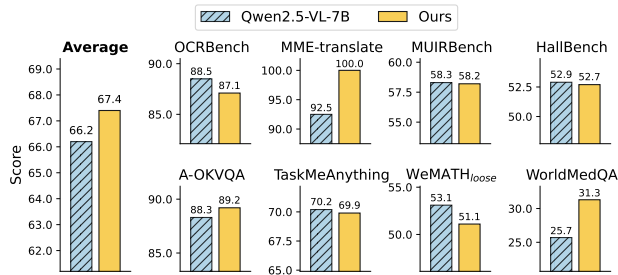


Figure 7. Beyond the gains on vision-centric benchmarks, CoVT also achieves slight improvements on **non–vision-centric tasks**

| Quantity | CVBench | BLINK | RW-QA | MM*-P | MMVP | V* | HR$_{4K}$ |
|---|---|---|---|---|---|---|---|
| 0 token | 76.6 | 55.5 | 70.7 | 68.0 | 55.3 | 77.5 | 68.6 |
| 16 empty | 75.7 | **56.0** | 70.3 | 67.9 | 56.7 | 77.5 | 68.1 |
| 1 token | 78.9 | 55.6 | 70.8 | 68.8 | 56.7 | **78.5** | **73.0** |
| 8 tokens | **80.0** | **56.0** | **71.6** | **69.2** | 58.7 | 78.0 | 72.9 |
| 32 tokens | 73.9 | 54.4 | 68.4 | 62.1 | 55.3 | 77.2 | 70.8 |

Table 4. **CoVT Ablation on segmentation token numbers.** The appropriate token number is essential for CoVT. **0 token** is the control group (direct fine-tuning), while **16 empty tokens** serve to isolate the role of the token embeddings themselves. 8 segmentation tokens perform the best.

| Type | Align | CVBench | BLINK | RW-QA | MM*-P | MMVP | V* | HR$_{4K}$ |
|---|---|---|---|---|---|---|---|---|
| Seg | Feature | 76.8 | 55.2 | 70.6 | 67.7 | 56.0 | 78.0 | 69.8 |
| | Ours | **77.9** | **57.4** | **71.1** | **68.5** | **58.7** | **79.1** | **71.9** |
| Depth | Feature | 77.0 | 54.2 | 70.5 | 67.6 | 55.3 | **78.0** | 71.3 |
| | Ours | **78.7** | **56.4** | **71.5** | **69.9** | **58.7** | 77.5 | **71.9** |

Table 5. Our tailored **alignment strategy** plays a crucial role in further enhancing the performance of CoVT.

we follow the CoT formatting paradigm used in the LLaVA-CoT 100k dataset and apply the same formatting to our own training data, ensuring full consistency with CoVT except for the absence of visual tokens. Fig. 6 shows that text-only CoT not only fails to improve performance on vision-centric reasoning tasks, but often degrades it. In contrast, CoVT consistently enhances performance across vision-centric benchmarks, highlighting the necessity of continuous visual tokens for effective visual reasoning.

**Token Numbers.** We ablate various numbers of segmentation visual tokens, as shown in Tab. 4. The "0 token" setting corresponds to directly fine-tuning the base model on our dataset. The "16 empty" setting replaces our 16 visual thinking tokens with 16 ordinary tokens without any visual alignment, serving as a pure latent-reasoning baseline. Settings with 1, 8, and 32 Segmentation tokens vary the token count while keeping the Depth and DINO tokens fixed at 4 each; the 8-token setting corresponds to our full model.

We observe that using too few Segmentation tokens leads to performance degradation, though still better than the "0 token" baseline. However, increasing the token count to 32 harms performance, maybe due to the difficulty of aligning a large number of segmentation tokens. The poor performance of the "16 empty" variant further highlights the importance of visually aligned tokens. Overall, the results demonstrate that 8 Segmentation tokens, together with 4 Depth and 4 DINO tokens, form a balanced and effective configuration, and that visual alignment is essential for enhancing vision-centric perception in VLMs.

**Segmentation and Depth Alignment Strategies.** We ablate two alignment approaches. Our primary method aligns CoVT tokens through task decoders, enabling the tokens to capture richer and more fine-grained perceptual cues. In contrast, direct feature alignment applies an MSE loss between the visual tokens and the encoder features of the corresponding visual model, which inevitably loses important perceptual details from the image.

As shown in Tab. 5, direct feature alignment consistently underperforms CoVT. These results highlight the importance of our tailored alignment strategies and demonstrate that aligning visual tokens with decoders yields more effective and perceptually grounded representations.

**CoVT remains competitive across various non-vision-centric benchmarks.** Fig. 7 shows our method remains comparable performance, with 1.2% improvement over eight non-vision-centric benchmarks, demonstrating that CoVT does not lead to a notable degradation in the generalization, and even yields an improvement for overall.

## 5. Conclusions

In this paper, we introduced **CoVT**, the chain of continuous visual thoughts that enables vision–language models to reason beyond discrete linguistic space by leveraging compact, dense visual representations. CoVT consistently improves visual-centric reasoning across diverse perception benchmarks and reveals that different types of visual tokens contribute to complementary aspects of multimodal understanding. These findings suggest that CoVT can serve as a general framework for integrating fine-grained perceptual reasoning into broader multimodal systems.

## 6. Acknowledgment

## References

[1] Manoj Acharya, Kushal Kafle, and Christopher Kanan. Tallyqa: Answering complex counting questions. In *AAAI Conference on Artificial Intelligence*, 2018. 6

[2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1

[3] 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, 6

[4] Mahtab Bigverdi, Zelun Luo, Cheng-Yu Hsieh, Ethan Shen, Dongping Chen, Linda Shapiro, and Ranjay Krishna. Perception tokens enhance visual reasoning in multimodal language models. In *Computer Vision and Pattern Recognition*, 2024. 3, 6

[5] Mahtab Bigverdi, Zelun Luo, Cheng-Yu Hsieh, Ethan Shen, Dongping Chen, Linda G Shapiro, and Ranjay Krishna. Perception tokens enhance visual reasoning in multimodal language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3836–3845, 2025. 3

[6] Eden Biran, Daniela Gottesman, Sohee Yang, Mor Geva, and Amir Globerson. Hopping too late: Exploring the limitations of large language models on multi-hop queries. In *Conference on Empirical Methods in Natural Language Processing*, 2024. 3

[7] Mathilde Caron, Hugo Touvron, Ishan Misra, Herv'e J'egou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9630–9640, 2021. 2

[8] Haolin Chen, Yihao Feng, Zuxin Liu, Weiran Yao, Akshara Prabhakar, Shelby Heinecke, Ricky Ho, Phí ThiMùi, Silvio Savarese, Caiming Xiong, and Huan Wang. Language models are hidden reasoners: Unlocking latent reasoning capabilities via self-rewarding. *ArXiv*, abs/2411.04282, 2024. 3

[9] Lin Chen, Jinsong Li, Xiao wen Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and Feng Zhao. Are we on the right way for evaluating large vision-language models? *ArXiv*, abs/2403.20330, 2024. 7

[10] Jeffrey Cheng and Benjamin Van Durme. Compressed chain of thought: Efficient reasoning through dense representations. *ArXiv*, abs/2412.13171, 2024. 3

[11] Zihui Cheng, Qiguang Chen, Xiao Xu, Jiaqi Wang, Weiyun Wang, Hao Fei, Yidong Wang, Alex Jinpeng Wang, Zhi Chen, Wanxiang Che, and Libo Qin. Visual thoughts: A unified perspective of understanding multimodal chain-of-thought. *ArXiv*, abs/2505.15510, 2025. 3

[12] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. 3

[13] Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025. 1

[14] Yuntian Deng, Kiran Prasad, Roland Fernandez, Paul Smolensky, Vishrav Chaudhary, and Stuart Shieber. Implicit chain of thought reasoning via knowledge distillation. *ArXiv*, abs/2311.01460, 2023. 3

[15] Haodong Duan, Junming Yang, Yu Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiao wen Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models. *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024. 7

[16] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407, 2024. 1

[17] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models. *ArXiv*, abs/2306.13394, 2023. 7

[18] Stephanie Fu, Tyler Bonnen, Devin Guillory, and Trevor Darrell. Hidden in plain sight: Vlms overlook their visual representations, 2025. 2

[19] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A. Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. *ArXiv*, abs/2404.12390, 2024. 7

[20] Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. Think before you speak: Training language models with pause tokens. *ArXiv*, abs/2310.02226, 2023. 3

[21] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. Hallusionbench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14375–14385, 2023. 7

[22] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14953–14962, 2023. 2, 3

[23] Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason E. Weston, and Yuandong Tian. Training large language models to reason in a continuous latent space. *ArXiv*, abs/2412.06769, 2024. 3

[24] J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685, 2021. 6

[25] Ziqi Huang, Ning Yu, Gordon Chen, Haonan Qiu, Paul Debevec, and Ziwei Liu. Vchain: Chain-of-visual-thought for reasoning in video generation. *arXiv preprint arXiv:2510.05094*, 2025. 3

[26] Ziqi Huang, Ning Yu, Gordon Chen, Haonan Qiu, Paul E. Debevec, and Ziwei Liu. Vchain: Chain-of-visual-thought for reasoning in video generation. *ArXiv*, abs/2510.05094, 2025. 3

[27] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloé Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B. Girshick. Segment anything. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3992–4003, 2023. 2, 5

[28] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *ArXiv*, abs/2205.11916, 2022. 3

[29] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9579–9589, 2023. 5

[30] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *ArXiv*, abs/2408.03326, 2024. 6

[31] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26286–26296, 2023. 1

[32] Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xucheng Yin, Cheng lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Sci. China Inf. Sci.*, 67, 2023. 7

[33] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering, 2022. 3

[34] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. *ArXiv*, abs/2304.09842, 2023. 3

[35] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback, 2023. 3

[36] João Matos, Shan Chen, Siena Placino, Yingya Li, Juan Carlos Climent Pardo, Daphna Idan, Takeshi Tohyama, David Restrepo, Luis Filipe Nakayama, Jose M. M. Pascual-Leone, Guergana K Savova, Hugo Aerts, Leo Anthony Celi, An-Kwok Ian Wong, Danielle S. Bitterman, and Jack Gallifant. Worldmedqa-v: a multilingual, multimodal medical examination dataset for multimodal language models evaluation. *ArXiv*, abs/2410.12722, 2024. 7

[37] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Q. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russ Howes, Po-Yao (Bernie) Huang, Shang-Wen Li, Ishan Misra, Michael G. Rabbat, Vasu Sharma, Gabriel Synnaeve, Huijiao Xu, Hervé Jégou, Julien Mairal, Patrick Labatut, Armand

Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. *ArXiv*, abs/2304.07193, 2023. 5

[38] Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu, Chong Sun, Xiaoshuai Song, Zhuoma Gongque, Shanglin Lei, Zhe Wei, Miaoxuan Zhang, Runfeng Qiao, Yifan Zhang, Xiao Zong, Yida Xu, Muxi Diao, Zhimin Bao, Chen Li, and Honggang Zhang. We-math: Does your large multimodal model achieve human-like mathematical reasoning? *ArXiv*, abs/2407.01284, 2024. 7

[39] Pooyan Rahmanzadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. Vision language models are blind: Failing to translate detailed visual features into words, 2025. 2

[40] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge. In *European Conference on Computer Vision*, 2022. 7

[41] Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. Visual cot: Advancing multi-modal language models with a comprehensive dataset and benchmark for chain-of-thought reasoning. *Advances in Neural Information Processing Systems 37*, 2024. 3

[42] Z. Su, Wenzhe Liu, Zitong Yu, Dewen Hu, Qing Liao, Qi Tian, Matti Pietikäinen, and Li Liu. Pixel difference networks for efficient edge detection. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5097–5107, 2021. 2, 5

[43] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11888–11898, 2023. 2, 3

[44] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 1

[45] Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025. 1

[46] Shengbang Tong, Ellis L Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, Austin Wang, Rob Fergus, Yann LeCun, and Saining Xie. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *ArXiv*, abs/2406.16860, 2024. 3

[47] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9568–9578, 2024. 7

[48] Fei Wang, Xingyu Fu, James Y. Huang, Zekun Li, Qin Liu, Xiaogeng Liu, Mingyu Derek Ma, Nan Xu, Wenxuan Zhou, Kai Zhang, Tianyi Yan, Wenjie Jacky Mo, Hsiang-Hui Liu,

Pan Lu, Chunyuan Li, Chaowei Xiao, Kai-Wei Chang, Dan Roth, Sheng Zhang, Hoifung Poon, and Muhao Chen. Muirbench: A comprehensive benchmark for robust multi-image understanding. *ArXiv*, abs/2406.09411, 2024. 7

[49] Wenbin Wang, Liang Ding, Minyan Zeng, Xiabin Zhou, Li Shen, Yong Luo, and Dacheng Tao. Divide, conquer and combine: A training-free framework for high-resolution image perception in multimodal large language models. *ArXiv*, abs/2408.15556, 2024. 3, 7

[50] 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 *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7907–7915, 2025. 2

[51] Xuezhi Wang and Denny Zhou. Chain-of-thought reasoning without prompting, 2024. 3

[52] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models, 2023. 3

[53] XuDong Wang, Xingyi Zhou, Alireza Fathi, Trevor Darrell, and Cordelia Schmid. Visual lexicon: Rich image features in language space. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19736–19747, 2025. 2

[54] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022. 1

[55] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. 3

[56] Chenfei Wu, Sheng-Kai Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *ArXiv*, abs/2303.04671, 2023. 3

[57] Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13084–13094, 2023. 7

[58] Penghao Wu and Saining Xie. V?: Guided visual search as a core mechanism in multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13084–13094, 2024. 2

[59] XAI. Grok-1.5 vision preview, 2024. 7

[60] Ji Xie, Trevor Darrell, Luke Zettlemoyer, and XuDong Wang. Reconstruction alignment improves unified multimodal models. *arXiv preprint arXiv:2509.07295*, 2025. 1

[61] 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, Jingren Zhou, Jingren Zhou,

Junyan Lin, Kai Dang, Keqin Bao, Ke-Pei Yang, Le Yu, Li-Chun Deng, Mei Li, Min Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shi-Qiang Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yi-Chao Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *ArXiv*, abs/2505.09388, 2025. 2

[62] Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in space: How multimodal large language models see, remember, and recall spaces. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 10632–10643, 2025. 2

[63] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *ArXiv*, abs/2406.09414, 2024. 2, 5

[64] Zeyuan Yang, Xueyang Yu, Delin Chen, Maohao Shen, and Chuang Gan. Machine mental imagery: Empower multimodal reasoning with latent visual tokens, 2025. 3

[65] Shaofeng Yin, Ting Lei, and Yang Liu. Toolvqa: A dataset for multi-step reasoning vqa with external tools. 2025. 3

[66] Kaining Ying, Fanqing Meng, Jin Wang, Zhiqiang 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, Yuning Qiao, Ping Luo, Kaipeng Zhang, and Wenqi Shao. Mmtbench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. *ArXiv*, abs/2404.16006, 2024. 7

[67] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11941–11952, 2023. 1

[68] Jieyu Zhang, Weikai Huang, Zixian Ma, Oscar Michel, Dong He, Tanmay Gupta, Wei-Chiu Ma, Ali Farhadi, Aniruddha Kembhavi, and Ranjay Krishna. Task me anything. *ArXiv*, abs/2406.11775, 2024. 7

[69] Yi-Fan Zhang, Huanyu Zhang, Haochen Tian, Chaoyou Fu, Shuangqing Zhang, Jun Wu, Feng Li, Kun Wang, Qingsong Wen, Zhang Zhang, Liang Wang, Rong Jin, and Tien-Ping Tan. Mme-realworld: Could your multimodal llm challenge high-resolution real-world scenarios that are difficult for humans? *ArXiv*, abs/2408.13257, 2024. 7

[70] Zhuosheng Zhang, Aston Zhang, Mu Li, George Karypis, Alex Smola, et al. Multimodal chain-of-thought reasoning in language models. *Transactions on Machine Learning Research*. 3

[71] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127:302 – 321, 2016. 6

[72] Chenyue Zhou, Mingxuan Wang, Yanbiao Ma, Chenxu Wu, Wanyi Chen, Zhe Qian, Xinyu Liu, Yiwei Zhang, Junhao Wang, Hengbo Xu, Fei Luo, Xiaohua Chen, Xiaoshuai Hao, Hehan Li, Andi Zhang, Wenxuan Wang, Lingling Li, Zhiwu Lu, Yang Lu, and Yi wang Guo. From perception to cognition: A survey of vision-language interactive reasoning in

multimodal large language models. *ArXiv*, abs/2509.25373, 2025. 4

[73] 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

# Chain-of-Visual-Thought:
# Teaching VLMs to See and Think Better with Continuous Visual Tokens
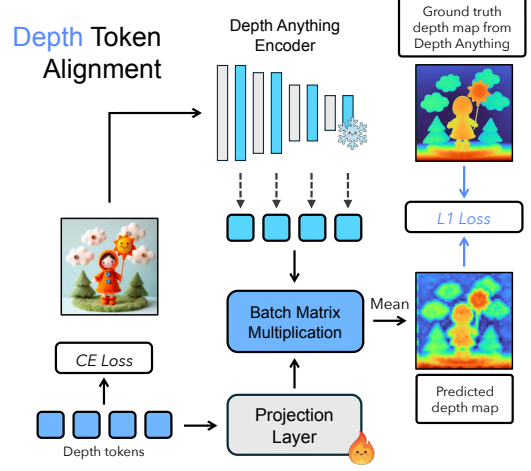
## Supplementary Material

## Table of Contents

Figure 9. Detailed framework for the depth token alignment.

## A. Additional Details of CoVT

In this section, we provide a comprehensive description of CoVT architecture. We first detail the design and functionality of the projection layer in (Sec. A.1). Then we provide more specific alignment architecture about segmentation token, depth token, and edge token from (Sec.A.2) to (Sec.A.4). Finally, we provide the detailed composition of the dataset used in CoVT in (Sec. A.5).

### A.1. Projection Layer

As shown in Fig. 8, the details of the projection layer are illustrated in the yellow block. It contains a single linear layer that projects the VLM latent space into the decoder's prompt space (or the encoder's feature space while aligning DINO tokens), formulated as:

$$\mathbf{z}_m = W\mathbf{z} + \mathbf{b}, \qquad (5)$$

where $\mathbf{z}$ denotes the VLM latent feature and $\mathbf{z}_m$ is the mapped prompt-space feature after the linear layer. We then introduce a learnable query $q$, while the mapped feature serves as both the key $k$ and value $v$ in the cross-attention layer, defined as:

$$\mathbf{z}_p = \text{Attn}(q, k, v) = \text{softmax}\left(\frac{qk^\top}{\sqrt{d_k}}\right)v, \qquad (6)$$

where $\mathbf{z}_p$ is the projected tokens, functioning as the prompts for the subsequent visual model decoding.

### A.2. Segmentation CoVT Token Alignment

Our model first predicts eight Segmentation tokens $\{T_i^{\text{seg}}\}_{i=0}^7$. Shown in Fig. 8, each token is then projected into the SAM decoder's prompt space through the projection layer,

$$T_i^{\text{sam}} = \text{proj}(T_i^{\text{sam}}), \qquad (7)$$

and the projected token serves as an individual prompt for mask decoding. Given the projected prompt $T_i^{\text{sam}}$ and the



Figure 8. Detailed frameworks for the projection layer and segmentation token alignment.

Figure 10. Detailed framework for edge token alignment.

| | Qwen2.5-VL-7B | | |
|---|---|---|---|
| | 1 type | 3 types | 4 types |
| **Optimization** | | | |
| Optimizer | | AdamW | |
| Learning rate | | 5e-5 | |
| Projection layer lr | | 1e-5 | |
| lr schedule | | cosine | |
| $\beta$ | | (0.9, 0.999) | |
| Weight decay | | 0.1 | |
| Warmup ratio | | 0.05 | |
| First stage steps | 4K | 6K | 8K |
| Second stage steps | | 3K | |
| Third stage steps | | 3K | |
| Forth stage steps | 4K | 5K | 6K |
| Per-GPU batch size | | 4 | |
| $\gamma$ | | 1.0 | |
| **LoRA settings** | | | |
| LoRA rank | | 16 | |
| LoRA alpha | | 32 | |
| **Visual models** | | | |
| SAM Encoder | | ViT-H | |
| Depth Anything v2 Encoder | | ViT-L | |
| PIDINet Encoder | | Table5-Baseline | |
| DINO v2 | | ViT-L | |

Table 6. Fine-tuning hyperparameter setup.

dense embedding $f$ from the SAM encoder, the SAM decoder produces one mask per token:

$$\hat{M}_i = \text{Decoder}(T_i^{\text{sam}}, f), \qquad \hat{M}_i \in [0, 1]^{H \times W}. \quad (8)$$

To construct reliable supervision, we generate all masks from SAM on the input image and apply a quality filter based on *mask area* and *stability score*. From these, we retain eight high-quality masks,

$$\mathcal{G} = \{M_j\}_{j=0}^7, \quad (9)$$

which serve as ground truths. We employ the Hungarian algorithm to match each predicted mask with one SAM mask. Unlike using similarity-based costs, we define the matching cost directly using the segmentation losses. For each pair $(\hat{M}_i, M_j)$, the Dice loss and Focal loss are

$$\mathcal{L}_{\text{dice}}(\hat{M}_i, M_j) = 1 - \frac{2 \sum \hat{M}_i M_j}{\sum \hat{M}_i + \sum M_j}, \quad (10)$$

$$\mathcal{L}_{\text{focal}}(\hat{M}_i, M_j) = -(1 - \hat{M}_i)^{\gamma_F} M_j \log \hat{M}_i, \quad (11)$$

where $\gamma_F$ is set to 2. Therefore, the matching cost becomes

$$C_{i,j} = \mathcal{L}_{\text{dice}}(\hat{M}_i, M_j) + \alpha \cdot \mathcal{L}_{\text{focal}}(\hat{M}_i, M_j), \quad (12)$$

where the $\alpha$ is set to 1 in our experiments. The optimal assignment is then obtained via

$$\sigma^* = \arg \min_\sigma \sum_{i=0}^7 C_{i, \sigma(i)}. \quad (13)$$

After obtaining the matching, the final mask loss is computed using the same Dice and Focal losses:

$$\mathcal{L}_{\text{mask}} = \sum_{i=0}^7 \left[ \mathcal{L}_{\text{dice}}\left(\hat{M}_i, M_{\sigma^*(i)}\right) + \alpha \cdot \mathcal{L}_{\text{focal}}\left(\hat{M}_i, M_{\sigma^*(i)}\right) \right]. \quad (14)$$

### A.3. Depth COVT Token Alignment

As shown in Fig. 9, our model predicts four Depth tokens $\{T_i^{\text{depth}}\}_{i=0}^3$, each of which is first projected into the DepthAnything decoder's prompt space through a linear projection layer:

$$T_i^{\text{depth-s}} = W_d T_i^{\text{depth}} + b_d. \quad (15)$$

DepthAnything v2 provides four dense intermediate-layer features

$$\{F_i^{\text{depth}}\}_{i=0}^3,$$

where $F_3^{\text{depth}}$ is the final-layer feature. Each projected depth token interacts with its corresponding feature map through batch matrix multiplication (BMM) to produce one depth map. This process is formulated as:

$$\hat{D}_i = \text{softmax}\left(T_i^{\text{depth-s}} \cdot F_i^{\text{depth}\top}\right), \quad i = 0, \dots, 3, \quad (16)$$

where $\hat{D}_i$ denotes the $i$th reconstructed depth map. Then the four reconstructed depth maps are averaged to produce
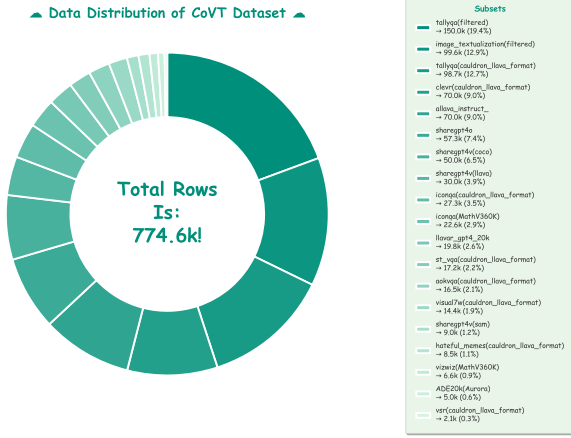
Figure 11. COVT dataset utilizes some subsets of LLaVA-OneVision, and merges the filtered TallyQA dataset and ADE20K-Depth from Aurora.

the final prediction:

$$\hat{D} = \frac{1}{4}\sum_{i=0}^{3}\hat{D}_i. \tag{17}$$

For supervision, we use the depth map predicted by DepthAnything v2 as the ground truth, denoted as $D^{\text{gt}}$. Depth tokens are aligned through an L1 reconstruction loss between the final depth prediction and the ground truth:

$$\mathcal{L}_{\text{depth}} = \left\| \hat{D} - D^{\text{gt}} \right\|_1. \tag{18}$$

### A.4. Edge COVT Token Alignment

Similarly, as shown in Fig. 10, we first project the four predicted Edge tokens into the PIDINet prompt space, obtaining $\{T_i^{\text{edge}}\}_{i=0}^3$. Each projected token is then used as a $1 \times 1$ convolutional kernel and applied to the dense intermediate features extracted from the PIDINet encoder. Let $\{F_i^{\text{edge}}\}_{i=0}^3$ denote the four intermediate feature maps. For each level, the reconstructed edge map is obtained by

$$\hat{E}_i = T_i^{\text{edge}} * F_i^{\text{edge}}, \quad \hat{E}_i \in \mathbb{R}^{H \times W}, \tag{19}$$

where "$*$" denotes a $1 \times 1$ convolution operation. The four reconstructed edge maps are then aggregated by averaging, followed by a sigmoid normalization:

$$\hat{E} = \sigma\left(\frac{1}{4}\sum_{i=0}^{3}\hat{E}_i\right). \tag{20}$$

For supervision, we use the edge map predicted directly by PIDINet, denoted as $E^{\text{gt}}$. The alignment between the predicted and ground-truth edges is enforced using the L1 loss:

$$\mathcal{L}_{\text{edge}} = \left\| \hat{E} - E^{\text{gt}} \right\|_1. \tag{21}$$

### A.5. COVT Dataset Composition

In order to fully leverage the value of COVT, we select the vision-centric subsets from the LLaVA-OneVision dataset, including:

- *IconQA(MathV360K)*
- *VizWiz(MathV360K)*
- *allava_instruct_*
- *aokvqa(cauldron_llava_format)*
- *clevr(cauldron_llava_format)*
- *hateful_memes(cauldron_llava_format)*
- *iconqa(cauldron_llava_format)*
- *image_textualization(filtered)*
- *llavar_gpt4_20k*
- *sharegpt4o*
- *sharegpt4v(coco)*
- *sharegpt4v(llava)*
- *sharegpt4v(sam)*
- *st_vqa(cauldron_llava_format)*
- *tallyqa(cauldron_llava_format)*
- *visual7w(cauldron_llava_format)*
- *vsr(cauldron_llava_format)*

In addition, we re-filtered the TallyQA dataset. Since TallyQA is a counting dataset but contains many samples with the answer 0, we reduce the proportion of zero-count samples and construct a 150k-sample subset. Moreover, following the methodology used in the Aurora paper, we generate 5k samples related to relative depth from the ADE20K dataset using the same procedure. We integrate these three components to form the complete COVT dataset.

## B. Additional Experiments

In this section, we first describe the experimental settings used throughout our study in (Sec. B.1). We then present ablation studies on the first two training stages in (Sec. B.2). Subsequently, we investigate the impact of varying the number of COVT tokens in (Sec. B.3). Finally, we provide additional output examples in (Sec. B.4).

### B.1. More settings

In Tab. 6, we present the complete hyperparameter configurations to ensure full reproducibility of our experiments. In this table, *1 type* denotes that the model is aligned with a single supervision signal, chosen from segmentation, depth, or DINO tokens. *3 types* corresponds to jointly aligning the model with segmentation, depth, and DINO tokens. *4 types* further incorporates edge tokens, thereby enabling simultaneous alignment across segmentation, depth, DINO, and edge tokens. For hyperparameters that remain consistent across all three experimental settings, we consolidate the corresponding columns and report them using a single centered entry for clarity and conciseness.

15

| Quantity | CVBench | BLINK | RW-QA | MM*-P | MMVP | V* | HR$_{4K}$ |
|---|---|---|---|---|---|---|---|
| Stage 3 & 4 | 78.2 | 53.8 | 70.0 | 68.3 | **60.7** | **78.0** | 71.2 |
| Ours | **80.0** | **56.0** | **71.6** | 69.2 | 58.7 | **78.0** | **72.9** |

Table 7. The first two stages in COVT 4-stage training strategy enhance the stability of the improvement.

## B.2. Training Stage Impact Ablation

To elucidate the pivotal contribution of the first two stages in the four-stage training strategy of COVT, we perform an ablation study comparing the full model with a variant trained solely on Stages 3 and 4, as reported in Tab. 7. The model trained across all four stages exhibits consistent and robust improvements over all evaluated benchmarks. In contrast, when restricting training to only the last two stages, COVT experiences notable degradation on the *BLINK* benchmark and yields only marginal gains on *RealWorld-QA* and *MMStar-Perception*. These results underscore the critical role of the early-stage training signals and highlight their importance in shaping the model's downstream performance.

## B.3. Token Numbers Ablation

To determine the optimal number of segmentation tokens, we conduct a detailed ablation while fixing both depth and DINO tokens at 4. As shown in Fig. 12, increasing the number of segmentation tokens in COVT from 1 to 32 yields an initial performance gain followed by a gradual decline, whereas the computational overhead rises steadily. Notably, the overall time cost of COVT is considerably higher than that of the baseline. Only a minor fraction of this overhead originates from the additional COVT tokens; a more substantial portion stems from COVT producing richer and more fine-grained responses for tasks such as image captioning, which naturally leads to longer output sequences (examples provided in Sec. B.4). Under our experimental conditions, eight segmentation tokens—combined with four depth tokens and four DINO tokens—offer the most favorable balance between performance and efficiency.

## B.4. More Results

We provide additional VQA examples in Fig. 13 through Fig. 17, including detailed image captioning, counting, instance identification, depth-aware questions, real-world OCR, and so on. In Fig. 13, we compare COVT with the baseline and show that COVT offers more fine-grained captioning ability. In Fig. 14, COVT improves instance identification (*e.g.*, describing an object as the white car hood) and better counting performance. In the subsequent figures from Fig. 15 to Fig. 17, we demonstrate that our model maintains stable performance across various vision-centric tasks and is also capable of tackling text-centric tasks such as real-world OCR. Specifically, in Fig. 15, COVT demonstrates its ability to correctly identify the NBA teams and



Figure 12. Equipped with 4 depth tokens and 4 DINO tokens, the model achieves its best performance when using 8 segmentation tokens in COVT. Allocating more COVT tokens leads to a slight diminishing performance and increased computational cost.

their scores in the first example, and to recognize visually ambiguous backgrounds in the second example. In Fig. 16, the first example shows that COVT can identify the farthest object and classify it correctly, while the second example illustrates that COVT can handle common-sense VQA tasks (e.g., identifying that the tall buildings are from Times Square). The last figure, Fig. 17, shows that COVT maintains stable performance on OCR tasks. In the first example, the model accurately detects the text "Sales Tax 4.24", and in the second example, it identifies the partially visible text "Mer" in the background, along with several car logos on the side.

## C. Limitations and Future Work

Despite the strong performance of COVT, several limitations remain. First, *we have not exhaustively explored the design space of visual models or token combinations.* Our current setup focuses on representative perception axes such as segmentation, depth, edges, and DINO features, but alternative or hybrid visual experts may yield more expressive or complementary tokens. Systematic exploration of this space could further enhance reasoning performance. Second, *our current framework does not yet realize fully interleaved multimodal reasoning.* COVT generates continuous visual thoughts but does not integrate them with free-form textual reasoning in an interleaved sequence. Developing such a unified chain that seamlessly blends textual and visual thoughts represents a promising direction for future work.

We believe that COVT provides a foundation for the next generation of multimodal reasoning systems—those capable of thinking fluidly across both language and vision in a self-contained, interpretable manner.

**Question**

Describe the scene in the picture in detail, and find out how many clouds are there in the sky.

**Qwen2.5-VL**

The image depicts a charming, handcrafted scene featuring a small, knitted figure dressed in an orange hooded outfit with matching boots. The figure is holding a cheerful, smiling sun-shaped balloon on a stick. The background includes three fluffy, cloud-like shapes in various shades of pink and white, floating above a small, green, grassy area with two coniferous trees. The sky is a soft, light blue color.

There are four clouds visible in the sky: one large pink cloud on the left, one smaller white cloud near the top right, and two more white clouds further to the right.

**CoVT**

<answer> The image depicts a charming, handcrafted scene featuring a small doll-like figure standing on a grassy patch. The figure is dressed in an orange outfit with matching boots, and it has a cheerful expression. The doll is holding a stick with a smiling sun at the top, which appears to be made from felt or similar material.

In the background, there are two green felt trees, each with a simple triangular shape and a brown trunk. Above the trees, there are four fluffy white clouds, each adorned with a brown button-like detail. The sky is painted in a soft blue color, providing a serene backdrop for the scene.

The overall composition is whimsical and playful, likely intended to evoke a sense of joy and innocence. The use of bright colors and simple shapes contributes to the overall aesthetic appeal of the image.

To answer the question about the number of clouds: There are four clouds visible in the sky above the trees. </answer>

Figure 13. Example of CoVT compared to the baseline Qwen2.5-VL-7B.

**Question**

What is the beautiful woman sitting on in the picture?

**Qwen2.5-VL**

The woman in the picture is sitting on the hood of a car.

**CoVT**

<answer> The beautiful woman is sitting on the hood of a white car. </answer>

**Question**

How many men are facing the camera in the picture?

**Qwen2.5-VL**

In the picture, there are three men facing the camera.

**CoVT**

Figure 14. Examples of CoVT compared to the baseline Qwen2.5-VL-7B.



**Question**

In the NBA basketball game shown in the picture, which team is leading?

**CoVT**

<answer> The Golden State Warriors are leading with a score of 107 to the Houston Rockets' 106. </answer>

**Question**

What's the background of the horses in the image?

**CoVT**

<answer> The background of the horses in the image consists of a beautiful mountain range with snow-capped peaks, creating a picturesque scene. </answer>

Figure 15. More examples of CoVT.

Figure 16. More examples of COVT.



Figure 17. More examples of COVT.