# Introducing Visual Perception Token into Multimodal Large Language Model

Runpeng Yu*, Xinyin Ma* and Xinchao Wang†

National University of Singapore

{r.yu,maxinyin}@u.nus.edu and xinchao@nus.edu.sg

(a) Response process of an MLLM equipped with Visual Perception Tokens.
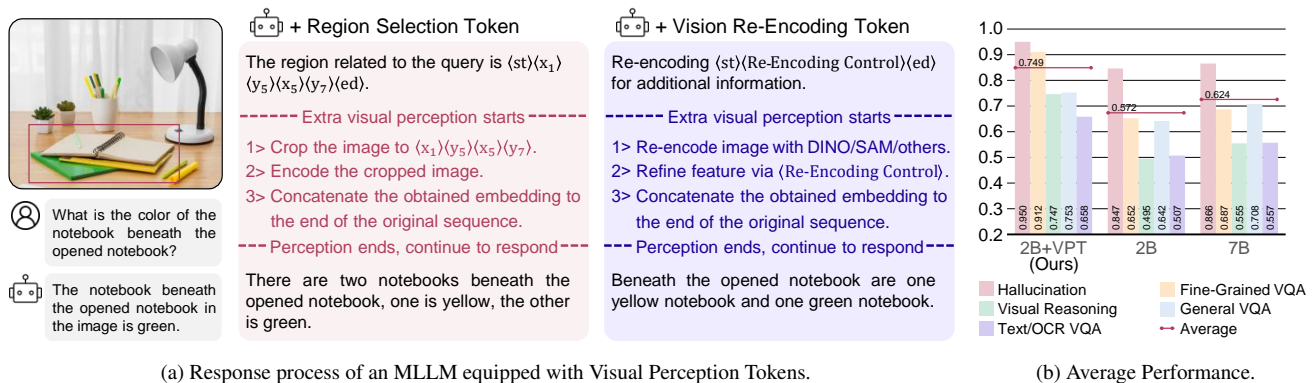
(b) Average Performance.

Figure 1. The Visual Perception Token aids MLLMs by triggering and controlling additional visual perception processes. Fig. 1a illustrates the response process of MLLMs equipped with the Region Selection Token or Vision Re-Encoding Token. The regions marked in the image are selected by the Region Selection Token. Fig. 1b presents the average performance of a 2B model with the Visual Perception Token across various VQA tasks (higher is better, with 1 being the maximum).

## Abstract

*To utilize visual information, Multimodal Large Language Model (MLLM) relies on the perception process of its vision encoder. The completeness and accuracy of visual perception significantly influence the precision of spatial reasoning, fine-grained understanding, and other tasks. However, MLLM still lacks the autonomous capability to control its own visual perception processes, for example, selectively reviewing specific regions of an image or focusing on information related to specific object categories. In this work, we propose the concept of Visual Perception Token, aiming to empower MLLM with a mechanism to control its visual perception processes. We design two types of Visual Perception Tokens, termed the Region Selection Token and the Vision Re-Encoding Token. MLLMs autonomously generate these tokens, just as they generate text, and use them to trigger additional visual perception actions. The Region Selection Token explicitly identifies specific regions in an image that require further perception, while the Vision Re-Encoding Token uses its hidden states as control signals to guide additional visual perception processes. Extensive experiments demonstrate the advantages of these tokens in handling spatial reasoning, improving fine-grained understanding, and other tasks. On average, the introduction of Visual Perception Tokens improves the performance of a 2B model by 30.9%, increasing its score from 0.572 to 0.749, and even outperforms a 7B parameter model by 20.0% (from 0.624). Please check out our repo here.*

## 1. Introduction

Multimodal Large Language Model (MLLM) depend on the perception capabilities of their vision encoder to process and utilize visual information. During this process, MLLM utilizes a vision encoder and a projector to embed visual information into the language space. The quality of Visual Perception determines whether MLLMs can accurately distinguish objects in an image [35], whether MLLMs can rely on visual information to answer questions instead of generating textual hallucinations [14], and whether MLLMs can perform precise reasoning about spatial relationships [3], among other tasks. While current MLLM systems demonstrate strong capabilities in visual information understanding [20, 34, 37, 45], they lack the ability to autonomously

---

* Equal Contribution; † Corresponding Author.

1

control their Visual Perception processes. Instead, these systems depend on manually designed pipelines to perform specific image annotations or visual features enhancement [35, 48].

In this work, we explore the task of enabling MLLMs to autonomously control their Visual Perception processes. Previously, MLLM-based agents and MLLMs equipped with function-calling or tool-use capabilities can be considered as having the ability to control subsequent tasks. They utilize the output of the LLM as arguments for subsequent functions or tool use. However, such control information is confined to the natural language space. The advantage of control signals in the natural language space lies in their interpretability, clear supervision signals, and ease of training data construction. However, these signals are constrained by specific formats. Additionally, natural language inherently contains redundancy, leading to efficiency issues. In this work, we aim to explore control signals beyond the natural language space. However, we also require that these signals remain naturally compatible with the next-token prediction paradigm of LLMs. To address this, we propose the concept of "Visual Perception Tokens". These tokens are integrated into the MLLM vocabulary and can be generated by the MLLM through next-token prediction, similar to natural language generation. These tokens do not correspond to specific words or characters in natural language; instead, their primary function is to trigger additional Visual Perception processes and convey control information for these processes.

We designed two types of Visual Perception Tokens. The first type is Region Selection Token, which instruct the MLLM to crop the input image and encode again the important regions relevant to the query using the vision encoder. The second type is the Vision Re-Encoding Token, which signals the model to input the image into (additional) vision encoder and use the resulting vision features to supplement the original MLLM's vision features. A projector takes both the additional vision features and the hidden state of the Vision Re-Encoding Token as inputs, enabling fine-grained control beyond merely triggering the vision encoder. In this work, we explore using an additional DINO model, a SAM model, or the model's original vision branch as the vision encoder. During the generation process, if the MLLM outputs any Visual Perception Token, the corresponding additional perception process is triggered, and the extra embedding sequence derived from the image is concatenated to the original LLM input. The LLM then continues generating the response in the form of next-token prediction. Fig. 1a illustrates the VQA process incorporating visual perception tokens. Fig. 3 provides a more detailed depiction of how visual perception tokens are generated by the MLLM and how they are utilized to control the visual perception process.

To train the MLLM to use Visual Perception Tokens, we constructed the Visual Perception Token training dataset, which includes 829k samples spanning four task categories: General VQA, Fine-Grained VQA, Spatial Reasoning, and Text/OCR-Related VQA. Experiments demonstrated that Visual Perception Token significantly enhances the MLLM's ability to autonomously control and refine its visual perception.

Our contributions can be summarized as follows:

1. We explored a novel task: enabling MLLMs to autonomously control their visual perception process.
2. We designed two types of Visual Perception Tokens: one enabling the MLLM to select regions of interest, and another allowing the MLLM to incorporate additional vision features and control the final embeddings input to the language model.
3. Experimental results demonstrate the effectiveness of our approach. On tasks such as Spatial Reasoning and Fine-Grained VQA, models equipped with Visual Perception Tokens achieved performance improvements of 34.6% and 32.7% over the 7B baseline model, respectively.

## 2. Related Work

### 2.1. Visual Prompting

From a technical perspective, our approach can also be regarded as a learnable visual prompting method. Visual prompting is a key technique in vision models, especially for segmentation tasks [11, 16, 27]. It uses a prompt encoder to interpret manual annotations, such as points and masks, to control segmentation granularity and assist instance selection. Recent advancements show that LVLMs can interpret visual cues like circles and color masks in a zero-shot manner without an additional encoder [29, 44]. Building on this, [41] and [42] have utilized segmentation model-generated masks as visual prompts, enhancing LVLM performance in segmentation and grounding. However, these methods are query-agnostic, while VQA tasks require adapting visual information based on the query. [48] addresses this by using an auxiliary VLM to generate query-specific visual prompts, overlaying attention maps to guide focus on relevant regions.

Built on these learning-free methods, [28] trains the MLLM to output bounding boxes of important regions. The image is then cropped and re-input for inference, creating a "crop and re-input" CoT process. Compared to [28], our design of Region Selection Tokens does not rely on bounding box information in the natural language space but instead uses specialized tokens to indicate the location of important regions. This design simplifies training and mitigates the issue of MLLMs having difficulty aligning the image coordinate system with coordinates described in natural language.

## 2.2. Visual Perception in MLLM

Currently, MLLMs have developed the ability to extract visual information, enabling them to be used not only for general image-based question answering and chat [1, 20, 21, 23, 34, 37, 38, 43], but also for applications in 3D understanding [8, 9, 51], video analysis [4, 32, 33, 50], domain-specific image question answering [15].

Our objective is to further enhance the visual perception capabilities of MLLMs using Visual Perception Tokens. In addition to the aforementioned visual prompting techniques, there are various other research directions focused on improving MLLM visual perception performance.

[35] identifies the ambiguity within CLIP's features as a limitation that adversely affects MLLMs' visual perception performance. To overcome this, it proposes incorporating additional DINOv2 [26] encoders to boost the visual perception capabilities of MLLMs. [36] observed that even with access to image data, MLLMs can sometimes base their responses on textual information and hallucinations instead of directly leveraging the visual content. To address this, [36] suggests adding image captions to improve visual perception accuracy. In the specific task of spatial reasoning, visual perception focuses on inferring spatial relationships between objects. To enhance a model's spatial reasoning capabilities, [3] proposed using depth maps of scenes as additional input. Moreover, [8] has explored how to leverage multiview scene information to improve the visual perception of MLLMs. To transform multiview information into a suitable input format, [8] employed a neural field representation. Similarly, to address the limitations of 2D images in physical world reasoning, [2, 51] have investigated MLLM visual perception methods based on 3D data. In this case, the input to MLLMs can include additional 3D point cloud features [9] and 3D scene graph features [6].

The above approaches aim to provide MLLMs with better visual inputs to enhance their visual perception abilities. In contrast, our approach features iterative perception, allowing MLLMs to provide feedback on the perception process, conduct multiple rounds of visual perception, and exercise control over the visual perception process.

## 3. Visual Perception Token

We add two types of extra Visual Perception Tokens to the original MLLM vocabulary. Tab. 1 provides a comparison of their key attributes. Previously, MLLM could only generate rationales and answers in natural language. However, with the addition of Visual Perception Tokens, the MLLM can now output valuable information in non-natural language form. These Visual Perception Tokens serve primarily as triggers; when a Visual Perception Token is output, the MLLM initiates additional visual perception processes.



Bounding Box: [1427, 647, 3841, 2823]
Token of the Region: $\langle x_1 \rangle \langle y_0 \rangle \langle x_3 \rangle \langle y_3 \rangle$
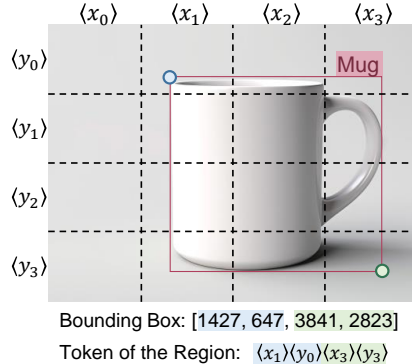
Figure 2. The relationship between the Region Selection Tokens we used and a precise bbox. Region Selection Token uses the cells containing the top-left and bottom-right corners to describe the approximate location of the region. In this example, the image is evenly divided into $4 \times 4$ cells. In our main experiment, we divide images into $8 \times 8$ cells.

| | Region Selection Token | | | | | | Vision Re-Encoding Token | | |
|---|---|---|---|---|---|---|---|---|---|
| Format | st | $x_{\min}$ | $y_{\min}$ | $x_{\max}$ | $y_{\max}$ | ed | st | DINO_Ctrl | ed |
| Function | Crop the Image | | | | | | Input Vision Feature | | |
| Encoder | Original Vision Encoder | | | | | | DINO/SAM/... | | |
| Informative Hidden States | ✗ | | | | | | ✓ | | |
| Supervision Signal | ✓ | | | | | | ✗ | | |
| Clear Semantic Meaning | ✓ | | | | | | ✗ | | |

Table 1. Comparison between the Region Selection Token and Vision Re-Encoding Token.

## 3.1. Region Selection Tokens

Please see Fig. 2 for the relationship between the Region Selection Tokens we used and a precise bbox and see Fig. 1a for an example when MLLM response with Region Selection Tokens.

Each group of Region Selection Tokens represents a bounding box (bbox). After the MLLM outputs a group of Region Selection Tokens, the original image will be cropped according to the bbox, preserving only the regions relevant to the query, and then be re-input into the MLLM. This "crop and re-input" approach enhances visual perception performance by directly increasing resolution. Although simple, it has been shown to be highly effective [20, 28]. For example, in document understanding and OCR-related tasks, the region containing the target text is often small, while the original document or image may be large and might even need to be downsized to fit the model input. In such cases, recognizing the query-relevant text becomes very challenging. However, cropping the image and then upsizing the relevant region before inputting it into the model significantly improves task performance.

Given an $h \times w$ image and a rectangular region $R$, a bounding box $[x_{\min}, y_{\min}, x_{\max}, y_{\max}]$ can precisely describe the position of region $R$, where $(x_{\min}, y_{\min})$ and $(x_{\max}, y_{\max})$

represent the coordinates of the top-left and bottom-right pixels of region $R$, respectively. However, precise pixel coordinates are not necessary for VQA tasks. Previous research also shows that MLLMs struggle to interpret and generate precise bounding boxes [40], especially when the input image resolution is not fixed, as is often the case in practice. In such cases, MLLMs struggle to infer the original image resolution from the patchified image embeddings, making it difficult to generate valid bounding boxes. Therefore, instead of using exact pixel coordinates to describe a region's bounding box, we describe only the approximate location of the region. We divide the $h \times w$ image evenly into a grid of $k \times k$ rectangular cells, with each cell sized $\frac{h}{k} \times \frac{w}{k}$. Each cell can be indexed by its row and column, with the top-left cell indexed as $(0,0)$ and the top-right cell as $(k-1, 0)$. We use the indices of the cells containing the top-left and bottom-right pixels of region $R$ to describe its location. In our implementation, we set $k = 8$.

A group of Region Selection Tokens includes six consecutive tokens: a <Region_Selection_Start> token, two tokens representing the index of the top-left cell, two tokens representing the index of the bottom-right cell, and a <Region_Selection_End> token. To enable the model to better distinguish between horizontal and vertical coordinates, we added $2k$ tokens, $<x_0>, \cdots, <x_{k-1}>$ and $<y_0>, \cdots, <y_{k-1}>$, specifically for indicating cell indices. For example, the top-left cell is represented as $<x_0><y_0>$ and the top-right cell as $<x_{k-1}><y_0>$. Given a group of Region Selection Tokens, the pixel coordinates for the top-left and bottom-right corners of the region to be cropped are calculated as $([x_{\min} \times \frac{w}{k}, y_{\min} \times \frac{h}{k})$ and $((x_{\max}+1) \times \frac{w}{k}, (y_{\max}+1) \times \frac{h}{k}])$, respectively.

Based on this design, Region Selection Tokens have clear and interpretable semantics, indicating specific locations within the image. During training, Region Selection Tokens are learned using the next token prediction loss.

### 3.2. Vision Re-Encoding Tokens

Please see Fig. 1a for an example when MLLM response with Vision Re-Encoding Tokens.

The Vision Re-Encoding Tokens trigger an additional vision encoder, such as DINO, to re-encode the original image, with the resulting vision features processed by a projector before being input into the MLLM.

Each set of Vision Re-Encoding Tokens consists of three tokens: <Re-Encode_Start> token, <Re-Encoding_Control>, and a <Re-Encode_End> token. The hidden state of the <Re-Encoding_Control> token token together with the vision features is input into the projector, allowing further control over the final embedding sequence fed into the LLM. To enable <Re-Encoding_Control> to freely convey any control information, we do not calculate the loss at the <Re-Encoding_Control> token during

training. As a result, the MLLM's output at the <Re-Encoding_Control> token can be arbitrary. We only require its hidden states to be informative. From this perspective, <Re-Encoding_Control> does not convey clear, interpretable semantics and cannot be decoded into a specific word or phrase. This is also the reason that additional <Re-Encode_Start> and <Re-Encode_End> tokens are used to mark the presence of <Re-Encoding_Control>.

**Mask Modeling**. To enhance the performance of <Re-Encoding_Control> token, we adopt a training approach similar to Masked Language Modeling. During training, 50% of the samples containing the Vision Re-Encoding Token undergo extra masking. We modify the attention mask for these samples, ensuring that the tokens corresponding to the answer in the dialogue can only access the <Re-Encoding_Control> token while being restricted from accessing the tokens corresponding to the original question and image embedding.

## 4. MLLM with Visual Perception Token

### 4.1. Architecture

When the Vision Re-Encoding Token triggers re-encoding, we use either an additional DINO or SAM model or the MLLM's original vision encoder. In all cases, an extra projector is added to align vision features with LLM embeddings. This projector is a cross-attention module that takes the hidden states of the <Re-Encoding_Control> token as the keys and values, and the vision features as the query. This design enables the hidden states of the <Re-Encoding_Control> token to control the vision features finally input to the LLM. See Fig. 3 for an illustration of the generation process with different types of Visual Perception Tokens. Below, we specifically describe the forward process in the modified MLLM.

An MLLM typically includes a LLM, a vision encoder $f_v$, and a projector $g_v$ that connects $f_v$ and the LLM. For an input image $\boldsymbol{x}$, the image features encoded by $f_v$ are denoted as $\boldsymbol{z} = f_v(\boldsymbol{x})$. After alignment through $g_v$, the resulting image embeddings can be represented as $\boldsymbol{h} = g_v(\boldsymbol{z})$. These image embeddings are then concatenated with text embeddings to form the input for the LLM. When the LLM outputs Region Selection Tokens, the cropped image $\boldsymbol{x}'$ is reprocessed through the original vision encoder and projector, resulting in $\boldsymbol{h}' = g_v(f_v(\boldsymbol{x}'))$, which is appended to the previous image and text embeddings as input to the LLM.

Let $f_D$ denote the vision model for re-encoding, $g_D$ the projector between $f_D$ and the LLM, and $\boldsymbol{h}_{DC} \in \mathbb{R}^{1 \times d_h}$ denote the hidden state of the <Re-Encoding_Control> token, where $d_h$ is the hidden size of the LLM. When the LLM outputs <Re-Encoding_Control> token, the original image is input into the $f_D$, yielding features $\boldsymbol{z}_D = f_D(\boldsymbol{x}) \in \mathbb{R}^{N \times d_z}$, where $N$ is the sequence length of the re-encoded image
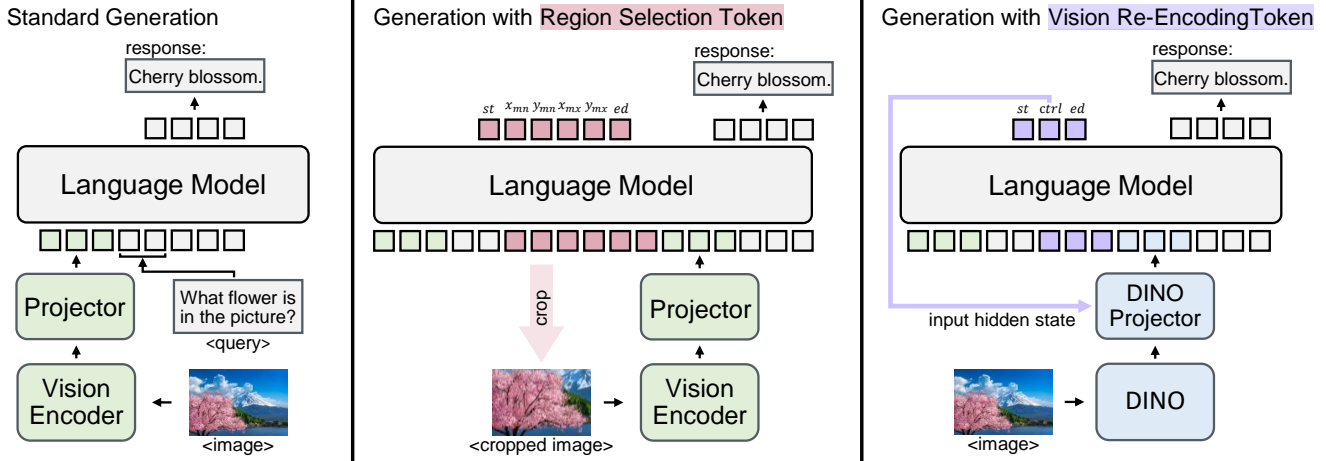
Figure 3. In a standard MLLM generation process, the model directly outputs an response based on the input image and query. However, an MLLM equipped with Visual Perception Tokens can first generate special tokens that trigger additional perception processes before responding. If the MLLM outputs a Region Selection Token, the original image is cropped and reprocessed through the visual encoder. The MLLM then bases its answer on two sets of visual embeddings: the first set contains the global embeddings from the original image, and the second set contains the local embeddings from the cropped image. If the MLLM outputs a DINO Feature Token, the DINO features of the image are used to supplement the original CLIP-based features. Additionally, besides the DINO features, the hidden state of the DINO Feature Token is also input to the projector as a condition to control which features are ultimately passed to the language model.

| Function | Task | Dataset |
|---|---|---|
| Training Region Selection Token (229k) | Text/OCR-Related VQA (84k) | DocVQA (33k), TextVQA (19k), TextCaps (32k) |
| | Spatial Reasoning (145k) | VSR (3k), GQA (10k), OpenImage (43k) |
| Training DINO Feature Token (293k) | General VQA (283k) | LLaVA Instruction Tuning COCO (212k), LLaVA Instruction Tuning GQA (71k) |
| | Fine-Grained VQA (10k) | CUB-200-2011 (10k) |
| Preserving Instruction Following Ability (307k) | General VQA (307k) | Remaining Samples in LLaVA Instruction Tuning (307k) |

Table 2. Composition of the training dataset. Our complete training dataset includes the data for training Visual Perception Tokens discussed in Sec. 4.2 and the remaining portion of the LLaVA-1.5 finetuning dataset. In total, the training dataset comprises 829k samples.

features and $d_z$ is the hidden size of $f_D$. The features $z_D$ and $h_{DC}$ are then input to the projector, resulting in the embeddings $h_D = g_D(z_D, h_{DC}) \in \mathbb{R}^{N \times d_h}$. $h_D$ is then concatenated with the previous image and text embeddings as input to the LLM.

## 4.2. Training Data for Visual Perception Token

We constructed the training dataset for Visual Perception Token based on the datasets from [20] and [28]. Our training data covers four types of tasks: Text/OCR-Related VQA, Spatial Reasoning, General VQA, and Fine-Grained VQA. The Text/OCR-Related VQA and Spatial Reasoning tasks are used to create training samples for Region Selection Token. The General VQA and Fine-Grained VQA tasks are used to construct training samples for Vision Re-Encoding Tokens. Please see Tab. 2 for the composition of training dataset.

**Samples with Region Selection Tokens**. We used the DocVQA [24], TextVQA [31], and TextCaps [30]

datasets to build training data for Text/OCR-Related tasks. DocVQA images include photos of documents such as books and invoices, while TextVQA and TextCaps images feature natural scenes such as billboards and store signs. Both datasets involve questions requiring reasoning about text within images. For these datasets, we use the text regions' bounding boxes obtained in [28], and convert the bounding boxes into Region Selection Token sequences as described in Sec. 3.1. The training data for Spatial Reasoning tasks included images from the VSR [19], GQA [10], and OpenImage [12] datasets, which contain real-world scenes and questions about the spatial relationships between objects. We used the filtered bounding boxes and question-answer pairs from [28] and converted the bounding boxes into the corresponding Region Selection Tokens.

**Samples with Vision Re-Encoding Tokens**. For General VQA tasks, we used data from COCO [18] and GQA [10] as provided by LLaVA 1.5 [20]. We inserted an additional MLLM response with Vision Re-Encoding

| Model | Max Resolution | Visual Reasoning | | | General VQA | | Fine-Grained VQA |
|---|---|---|---|---|---|---|---|
| | | GQA | OpenImage | VSR | LLaVA Instruction Tuning | Flickr* | CUB Birds |
| Qwen2-VL-2B | 224 | 0.448 | 0.413 | 0.561 | 0.655 | 0.521 | 0.621 |
| Qwen2-VL-7B | 224 | 0.464 | 0.442 | 0.632 | 0.703 | 0.524 | 0.697 |
| Qwen2-VL-2B-VPT (DINO) | 224 | **0.606** | **0.842** | **0.657** | **0.705** | **0.558** | **0.892** |
| Qwen2-VL-2B | 512 | 0.487 | 0.418 | 0.580 | 0.728 | 0.557 | 0.652 |
| Qwen2-VL-7B | 512 | 0.569 | 0.456 | 0.641 | 0.780 | 0.636 | 0.687 |
| Qwen2-VL-2B-VPT (DINO) | 512 | 0.625 | 0.872 | 0.738 | 0.797 | 0.663 | 0.898 |
| Qwen2-VL-2B-VPT (DINO, Free Choice) | 512 | **0.635** | 0.874 | 0.733 | 0.802 | 0.705 | 0.911 |
| Qwen2-VL-2B-VPT (CLIP) | 512 | 0.621 | 0.872 | 0.746 | 0.791 | 0.660 | 0.913 |
| Qwen2-VL-2B-VPT (SAM) | 512 | 0.617 | 0.876 | 0.746 | 0.796 | 0.657 | 0.905 |
| Qwen2-VL-7B-VPT (CLIP) | 512 | 0.633 | **0.878** | **0.790** | **0.831** | **0.680** | **0.921** |

| Model | Max Resolution | Text/OCR Related VQA | | | | Hallucination | Average |
|---|---|---|---|---|---|---|---|
| | | DocVQA | TextVQA | TextCaps | DUDE* | POPE* | |
| Qwen2-VL-2B | 224 | 0.051 | 0.383 | 0.390 | 0.063 | 0.821 | 0.448 |
| Qwen2-VL-7B | 224 | 0.063 | 0.421 | 0.431 | 0.095 | 0.827 | 0.482 |
| Qwen2-VL-2B-VPT (DINO) | 224 | **0.125** | **0.537** | **0.466** | **0.103** | **0.843** | **0.576** |
| Qwen2-VL-2B | 512 | 0.301 | 0.765 | 0.710 | 0.253 | 0.847 | 0.572 |
| Qwen2-VL-7B | 512 | 0.360 | 0.816 | 0.732 | 0.322 | 0.866 | 0.624 |
| Qwen2-VL-2B-VPT (DINO) | 512 | 0.573 | 0.860 | 0.766 | 0.430 | 0.893 | 0.738 |
| Qwen2-VL-2B-VPT (DINO, Free Choice) | 512 | 0.576 | 0.861 | 0.758 | 0.438 | 0.950 | 0.749 |
| Qwen2-VL-2B-VPT (CLIP) | 512 | 0.567 | 0.856 | 0.770 | 0.433 | 0.887 | 0.738 |
| Qwen2-VL-2B-VPT (SAM) | 512 | 0.558 | 0.858 | 0.750 | 0.431 | 0.894 | 0.735 |
| Qwen2-VL-7B-VPT (CLIP) | 512 | **0.658** | **0.906** | **0.788** | **0.532** | **0.903** | **0.773** |

Table 3. Performance comparison of MLLMs with and without Visual Perception Tokens. Datasets marked with "*" are not used in the training process. The best performance is highlighted in **bold**. A 2B model with Visual Perception Tokens can even outperform the 7B model without Visual Perception Tokens.

Tokens between the original question and response. For Fine-Grained VQA tasks, we used the CUB-200-2011 [7] dataset, which includes images of 200 bird species and annotations about detailed attributes of these birds. We used the question-answer pairs generated in [28] and inserted an MLLM response with Vision Re-Encoding Tokens between the questions and answers.

# 5. Experiments

## 5.1. Main Results

We used Qwen2-VL-2B or Qwen2-VL-7B [37] as the base MLLM model and DINOv2 [26] or SAM [11] as the additional vision feature extractor. We use CLIP to denote the case when the original vision branch is used in the re-encoding process. The 2B models are fully fine-tuned, while the 7B model are tuned using LoRA. The datasets described in Sec. 4.2 were used to finetune the model. We then compared its performance with the original Qwen2-VL-2B and Qwen2-VL-7B models.

We evaluated the performance of Qwen2-VL-2B-VPT on the test sets of the datasets used in Sec. 4.2. When an official test split was available, we used it directly. When no official test split was provided, we randomly split the data into training and testing sets and filter out any images that appeared in both sets to prevent data leakage. Additionally, we also conducted evaluations on Flickr [46], DUDE [13], and POPE [17] datasets, which were not used in training, to assess the generalization capability of our method.

Following established practices [20, 49], we used GPT-4o (2024-08-06) to evaluate the alignment between the model's responses and the ground truth for each question. A higher score indicates a higher degree of matching, with 0 representing no match and 1 indicating a perfect match. We reported the average matching scores for each dataset. The prompts used in the evaluation are included in the supplementary material.

The experimental results are presented in Tab. 3. In terms of overall average performance, 2B model with Visual Perception Tokens outperforms the 7B model without Visual Perception Tokens by a significant margin. A more detailed analysis is as follows: (1) The trend of using different vision encoders for re-encoding remains consistent, with the 2B model enhanced by Visual Perception Tokens outperforming the 7B model. (2) We implemented two methods for triggering Visual Perception Tokens. The first method explicitly enforces the use of either Region Selection or DINO Feature Tokens across the entire dataset, aligning with real-world scenarios where users may require explicit control. This corresponds to most results in Tab. 3. The second allows the model to autonomously decide whether to use Perception Tokens and which type to apply, corresponding to the "Free Choice" results in Tab. 3. This "Free Choice" mechanism achieves even better performance than the first approach. (3) The Visual Perception Token shows a clear advantage in Visual Reasoning and Fine-Grained
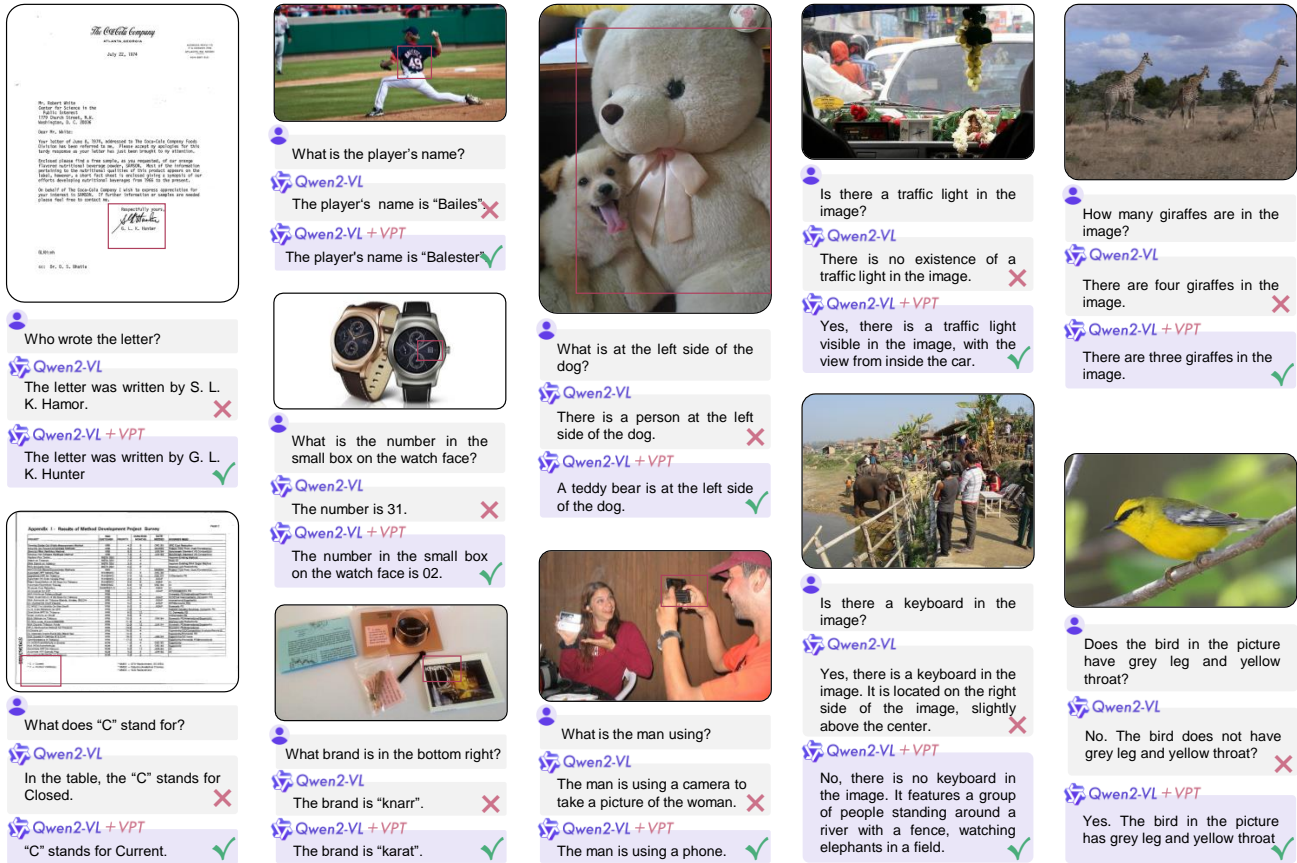
Figure 4. Examples collected from the testing sets. The responses were generated by the 7B model and the 2B+VPT model. During the generation process, if Region Selection Tokens were utilized, the region selected by these tokens are highlighted with red boxes in the images. For additional examples, please refer to the supplementary material.

VQA tasks, with improvements of 0.161 and 0.207 over the 7B model, respectively. However, in some datasets for General VQA and Text/OCR-related VQA tasks, the 2B model with Visual Perception Tokens only performs comparably to the 7B model, without a significant performance gain. (4) The Visual Perception Token remains effective in zero-shot settings. During evaluation, we included three datasets that were not used in training. On these datasets, the 2B-VPT model still outperforms or matches the performance of the 7B model. (5) The Visual Perception Token is effective for both high- and low-resolution images. On average, when using low-resolution images as input, the 2B+VPT model improves by 0.094 over the 7B model, while for high-resolution images, it shows an improvement of 0.084 over the 7B model.

In Fig. 4, we show examples demonstrating the effectiveness of Visual Perception Tokens. Visual Perception Tokens are particularly effective for several types of queries: The first category involves locating small regions within large documents or images. The small regions can be signatures, footnotes, page numbers, product brands, logos, or small

text on clothing, which are often too small for direct OCR by MLLMs. However, Region Selection Tokens can pinpoint these areas. The second category addresses hallucination issues. For example, when asked what is next to a dog, text-based hallucination might generate an answer like "a person". Region Selection Tokens mitigate such errors by grounding responses in image. The third category involves identifying or counting objects in complex scenes. Vision Re-Encoding Tokens excel in these tasks because encoding the image twice enhances segmentation and detection.

## 5.2. Discussion on Region Selection Token

In this subsection, we assess the necessity of introducing Region Selection Tokens and examine whether directly using bounding boxes is an effective approach for indicating selected regions. Additionally, we conduct an ablation study on the granularity parameter $k$.

Due to the cost of experiments, we did not train models on all the datasets used in the main experiment for ablation. Instead, we focused on training using the DocVQA, TextVQA, and TextCaps datasets and evaluated on their test

splits. This choice was made because, for text/OCR-related VQA tasks, the query-related region is typically small, and the images often contain complex objects that distract the MLLM. Therefore, the quality of region selection substantially affects the final VQA performance. We trained five models in total. The first four models employed Region Selection Tokens with $k$ values of 4, 8, 16, and 32, where a larger $k$ indicates finer granularity in the image partition. The fifth model directly used bounding boxes, as described in [28]. The results are presented in Tab. 4.

The findings are twofold: (1) Region Selection Tokens prove to be more effective than bounding boxes. As shown in the results, models with $k = 8$ and $k = 16$ significantly outperformed model that used bounding boxes as region indicators. We also observed that when MLLMs utilized bounding boxes, they generated many invalid bounding boxes, such as those where the height or width exceeded the original image dimensions. (2) There is an optimal granularity $k$. In our experiments, $k = 8$ generally yielded the best performance. When $k$ is too small, the cropped region remains large and can still contain complex scenes or multiple objects, which fails to guide the model's attention effectively. Conversely, when $k$ is too large, more new tokens are required, increasing the learning difficulty for the MLLM. Moreover, overly fine-grained region descriptions do not contribute to VQA performance. These two factors together leads to an overall decline in performance.

### 5.3. Discussion on Vision Re-Encoding Token

In this subsection, we validate the effectiveness of the Vision Re-Encoding Token and conduct an ablation on it.

To verify the effectiveness of the Vision Re-Encoding Token, Tab. 5 presents the performance comparison between two model variants. The model with control information (last row in the table) includes the hidden state of the <Re-Encoding_Control> token as an input to the projector, enabling fine-grained control over the embeddings fed into the LLM. In contrast, the model without control information (middle row in the table) uses directly the re-encoded image feature as input to the LLM, after aligning its dimension with the embedding dimension of the LLM. In this experiment, DINO-v2 is used as the additional vision encoder. The experimental results show that the additional control information carried by the <Re-Encoding_Control> token significantly enhances model performance.

Next, we conduct an ablation study on the number of <Re-Encoding_Control> tokens and mask modeling. Considering the computational cost, in these experiments, we trained models using only the training data from the CUB Bird dataset and LLaVA Instruction Tuning data from COCO and GQA. Results are presented in Tab. 6 and indicate the following conclusions. First, increasing the number of <Re-Encoding_Control> tokens from 1 to 2 provides

|  |  | DocVQA | TextVQA | TextCaps |
|---|---|---|---|---|
| Region Selection Token | $k = 4$ | 0.251 | 0.599 | 0.521 |
|  | $k = 8$ | **0.369** | 0.686 | **0.578** |
|  | $k = 16$ | 0.307 | **0.690** | 0.569 |
|  | $k = 32$ | 0.264 | 0.634 | 0.539 |
| Bounding Box |  | 0.219 | 0.620 | 0.547 |

Table 4. Ablation on the parameter $k$, which controls the granularity of the Region Selection Token and whether Bounding Boxes should be used directly.

| Control Info | CUB Bird | LLaVA COCO | LLaVA GQA |
|---|---|---|---|
| ✗ | 0.7459 | 0.6226 | 0.4354 |
| ✓ | **0.8943** | **0.7804** | **0.7908** |

Table 5. Validation of the designed Vision Re-Encoding Token's ability to convey valuable control information.

| Method | Mask Modeling | No. of Ctrl Tokens | CUB Bird | LLaVA COCO | LLaVA GQA |
|---|---|---|---|---|---|
| Qwen2-VL-2B |  | 0 | 0.882 | 0.718 | 0.734 |
| 2B+VPT (DINO) |  | 1 | 0.892 | 0.730 | 0.758 |
| 2B+VPT (DINO) |  | 2 | 0.896 | 0.731 | 0.761 |
| 2B+VPT (DINO) |  | 4 | 0.876 | 0.715 | 0.759 |
| 2B+VPT (DINO) | ✓ | 1 | 0.901 | 0.739 | 0.763 |
| +Tune Projector | ✓ | 1 | **0.918** | **0.755** | 0.769 |
| 2B+VPT (CLIP) | ✓ | 1 | 0.902 | 0.728 | 0.770 |
| +Tune Projector | ✓ | 1 | 0.909 | 0.734 | **0.775** |

Table 6. Ablation on the number of Vision Re-Encoding Token. On one hand, with Mask Modeling the performance is improved. On the other hand, the results do not suggest a significant improvement when using more <Re-Encoding_Control> tokens.

limited performance gains, while increasing to 4 leads to a decline. This is due to the projector's lightweight design, making it susceptible to over-parameterization and overfitting. Second, introducing masked modeling improves the performance. Further experiments show that fine-tuning the linear layer connecting the hidden states of <Re-Encoding_Control> in the projector for 1 extra epoch further improves performance, confirming that these hidden states encode valuable control information.

## 6. Conclusion

In this work, we propose Visual Perception Tokens to enable MLLM to autonomously controls its visual perception process. MLLM can generate these Visual Perception Tokens in a manner similar to generating natural language and use them to convey information to the visual perception process. Experiments validate the effectiveness of Visual Perception Token over various tasks.

# Introducing Visual Perception Token into Multimodal Large Language Model
## *- Supplementary Material -*

## S1. Implement Details

### S1.1. Training Details

Our training process consists of two phases: alignment and finetuning. The alignment stage aligns the additional vision features with the LLM embeddings. If the original vision encoder is used for re-encoding, the alignment stage is omitted. We use the same image-text pair data for the LLaVA 1.5 alignment, and only use the additional vision branch as the LLM's input. During training, all components except the projector are frozen. In this phase, we train the model for 1 epoch with a learning rate of 2e-3 and a batch size of 128. The second finetuning stage allows the model to learn to output the correct Region Selection Tokens and to transmit information through the Vision Re-Encoding Tokens. We finetune the model using our constructed dataset, as well as remaining samples from the LLaVA 1.5 finetuning dataset that were not included in our dataset. In this stage, all components except the original visual encoder and the additional vision encoder are unfrozen. In this phase, we train the model for 1 epoch with a learning rate of 2e-5 and a batch size of 256. For both the first and the second phase, we use AdamW optimizer. The experiments are deployed on 8 A100 GPU. The total training time is about 20 hours. For the 7B model, the rank of the LoRA is set to 512.

### S1.2. Evaluation Prompt

Following established practices [20, 49], we used GPT-4o (2024-08-06) to evaluate the alignment between the model's responses and the ground truth for each question. We use the evaluation prompt in [28].

> **Evaluation Prompt**
>
> You are responsible for proofreading the answers, you need to give a score to the model's answer by referring to the standard answer, based on the given question. The full score is 1 point and the minimum score is 0 points. Please output the score in the form 'score: <score>'. The evaluation criteria require that the closer the model's answer is to the standard answer, the higher the score.
> Question: <question>
> Ground Truth: <ground truth>
> Answer: <answer>

### S1.3. Template of the Training Data Examples

Here, we show the format of our training examples. The training example for the Region Selection Token is essentially the same as the samples used in [28], except that the method for representing regions has changed from bounding boxes to region tokens. The training example for the Vision Re-Encoding Token is almost identical to the data in the original LLava [21] finetuning dataset, with the only difference being the insertion of an additional round of dialogue between the original question and answer. This added dialogue includes the Vision Re-Encoding Token.

> **Template of Training Example for Region Selection Token**
>
> **User**: <image> <question> Please identify the region that can help you answer the question better, and then answer the question.
> **Assistant**: <Region_Selection_Start> <x_min> <y_min> <x_max> <y_max> <Region_Selection_End>.
> **User**: <image>
> **Assistant**: <ground truth>

|  | MME | | MMB | |
| --- | --- | --- | --- | --- |
|  | Cognition | Perception | en | cn |
| Qwen2-VL-2B | 1434 | 280 | 78.20 | 77.30 |
| Qwen2-VL-7B | 1664 | 335 | 78.70 | **83.30** |
| 2B+VPT (DINO) | 1511 | 274 | 79.11 | 76.64 |
| 2B+VPT (CLIP) | 1510 | 273 | 79.53 | 77.41 |
| 2B+VPT (SAM) | 1475 | 270 | 80.22 | 76.99 |
| 7B+VPT (CLIP) | **1706** | **336** | **83.80** | **83.30** |

Table S1. Performance comparison of MLLMs with and without Visual Perception Tokens on MME and MMBench Benchmarks.

---

**Template of Training Example for Vision Re-Encoding Token**

**User**: <image> <question> Please require additional perception features, and then answer the question.
**Assistant**: <Re-Encoding_Start> <Re-Encoding_Control> <Re-Encoding_End>.
**User**: <image>
**Assistant**: <ground truth>

---

The training for the free-choice experiment differs from other experiments only in the sample template. For the free-choice experiment, we removed the additional prompt from the questions. The training sample template is as follows.

---

**Template of Training Example for Region Selection Token (Free Choice)**

**User**: <image> <question>
**Assistant**: <Region_Selection_Start> <x_min> <y_min> <x_max> <y_max> <Region_Selection_End>.
**User**: <image>
**Assistant**: <ground truth>

---

**Template of Training Example for Vision Re-Encoding Token (Free Choice)**

**User**: <image> <question>
**Assistant**: <Re-Encoding_Start> <Re-Encoding_Control> <Re-Encoding_End>.
**User**: <image>
**Assistant**: <ground truth>

---

## S2. Supplementary Experiments

We conducted experiments on the MME [5] and MM-Bench [22] benchmarks without using the Visual Perception Token, allowing the model to generate answers directly. This assessed the impact of our fine-tuning on general benchmarks. Results in Tab. S1 show that our model does not cause degeneration and even improves performance on these benchmarks.

To verify the advantage of the Region Selection Token over direct BBox prediction, we compared the predicted regions with ground truth using IoU and Intersection over Ground Truth (IoGT), defined as:

$$(\text{IoGT} = \frac{\text{Area of } (GT \cap \text{Pred})}{\text{Area of } GT}).$$

Results in Tab. S2 show that Region Selection Token significantly outperforms direct BBox prediction in accuracy.

## S3. Further Examples

Here we present additional examples obtained using the visual perception token. Figs. S1 and S2 include the responses generated with the Vision Re-Encoding Token. Figs. S3 and S4 present the responses generated with the Region Selection Token, with the regions selected by the Region Selection Token highlighted in the images.

|  | Metric | DocVQA | TextVQA | TextCap |
|---|---|---|---|---|
| Directly Predicting BBox | IoU | 0.15 | 0.26 | 0.25 |
|  | IoGT | 0.20 | 0.28 | 0.27 |
| Using Region Selection Token | IoU | 0.26 | 0.56 | 0.50 |
|  | IoGT | 0.38 | 0.71 | 0.66 |

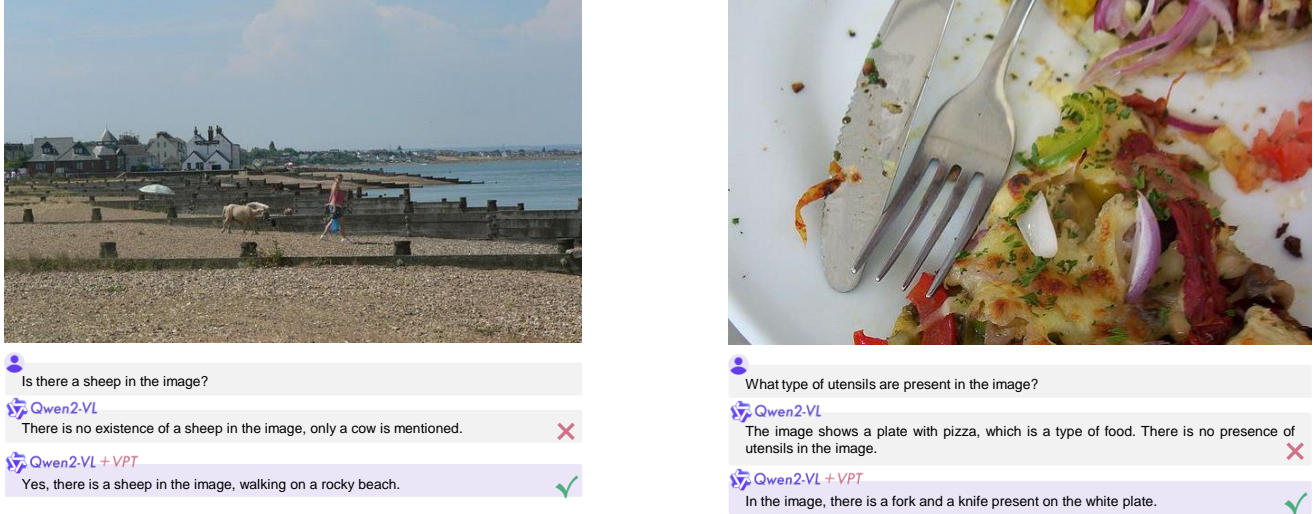Table S2. Performance comparison of MLLMs with and without Visual Perception Tokens on MME and MMBench Benchmarks.



Figure S1. This set of images demonstrates how the DINO Feature Token assists MLLMs in identifying specific objects within images. These objects are often difficult for MLLMs to recognize directly due to their small size or interference from surrounding objects.

## S4. Additional Related Works

### S4.1. Reasoning Token

In Large Language Model (LLM), there are tokens, similar to Visual Perception Token, which are designed to control the generation process of LLM. These token are termed reasoning tokens or planning token and have recently been introduced in OpenAI's o1 model [25] and other LLMs. For example, to enhances models' reasoning capabilities, reasoning tokens were explicitly integrated into OpenAI's o1 models to segment prompts into smaller, manageable parts, exploring multiple response strategies before generating the final output [25]. Similar methods aim to incorporate CoT reasoning into language models through planning tokens or distillation techniques. For example, a hierarchical generation framework using planning tokens has been proposed, embedding high-level plans at each reasoning stage with minimal parameter increase [39]. Moreover, techniques like Rephrase and Respond have been distilled back into models, improving efficiency and accuracy in reasoning, as demonstrated in [47].

Our work focuses on MLLMs, where we design visual perception tokens to enhance the visual perception capabilities of MLLMs, not for LLM. Moreover, our exploration goes beyond LLM reasoning tokens. Unlike these tokens, which merely trigger specific actions and lack the ability to convey detailed instructions or rich information, we focus on designing tokens capable of transmitting nuanced control information for fine-grained visual perception.

## S5. Discussion

**Adaptability of Visual Perception Token.** The design of the visual perception token depends on the specific visual perception method. In this paper, we use Crop and the addition of vision features as examples to introduce two types of visual perception tokens. However, our approach can be extended to other visual prompting techniques or visual encoder models, and even to LLM-agent or LLM-tool systems beyond vision.

How many airplanes are visible in the image?

Qwen2-VL
There are three airplanes visible in the image. ✗

Qwen2-VL + VPT
There are two airplanes visible in the image. ✓



How many elephants are in the image, and can you describe their age difference?

Qwen2-VL
There is only one elephant in the image, and there is no mention of age difference. ✗

Qwen2-VL + VPT
There are two elephants in the image, one being a baby elephant and the other being an older elephant. ✓



How many laptops are there on the desk?

Qwen2-VL
There is one laptop on the desk. ✗

Qwen2-VL + VPT
There are two laptops on the desk. ✓



How many laptops are there on the desk?

Qwen2-VL
There is one laptop on the desk. ✗

Qwen2-VL + VPT
There are two laptops on the desk. ✓

Figure S2. This set of images illustrates how the DINO Feature Token assists MLLMs in counting the number of objects in an image. Counting has long been a significant limitation for MLLMs. By leveraging the DINO Feature, the DINO Feature Token enables precise localization of individual objects within the image, thereby improving the counting capability of MLLMs.

**Document 1 (form):**

DATE: March 22, 1991

COUNTRY - U. S.

GRADE - CG1 1989 Chinese Flue Cured

| Dealer | Lbs. Strips Packed | % Packed | No. Rejects/Reruns - Reason |
|---|---|---|---|
| A. C. Monk | 597.472 | 100.0 | 1 stem |

Total  597.472  100.0  1

Foreign Matter Found In Core Samples

DEALER: A. C. Monk

TYPE: No. of Pieces

| Grass/Straw | 2 |
| Lint/String | 22 |
| Paper | 4 |
| Plastic | |
| Feathers | |
| Foam | |
| Wood | |
| Foil | |
| Other | |
| Total Pieces F. M. | 28 |
| Lbs. Core Sample | 352 |
| No. Pieces F. M. / Lb. | .1 |

51336 0089

👤 What is the name of the Dealer?

Qwen2-VL
The dealer is C. York  ✗

Qwen2-VL + VPT
The dealer is A. C. Monk.  ✓

**Document 2 (letter):**

ARee/Vate/Amstar Crp

Amstar
CORPORATION

1251 AVENUE OF THE AMERICAS, NEW YORK, N. Y. 10020 • (212) 489-9000

Research & Development Division
266 Kent Avenue
Brooklyn, New York  11211

July 23, 1976

Dr. Paul Saltman
University of California, San Diego
Office of the Vice Chancellor-
  Academic Affairs
LaJolla, CA  92093

Dear Paul:

Enclosed is a check for $419.00 representing your expenses in connection with our Food and Nutrition Advisory Council meeting. We were very pleased that you were able to attend.

Sincerely yours,

A. B. Rizzuto

A. B. Rizzuto
Director of Research

ABR/jf
Enc.

👤 What is the amount of the check enclosed?

Qwen2-VL
The check enclosed is for $1,000.00.  ✗

Qwen2-VL + VPT
The amount of the check enclosed is $419.00.  ✓

**Document 3 (budget table):**

SWANSON CENTER FOR NUTRITION
1979 Budget
Revised, July, 1979

Budget Analysis for September, 1979

👤 What is the name of the institute mentioned in the title?

Qwen2-VL
The name of the institute mentioned in the title is "Sanford Center for Nutrition."  ✗

Qwen2-VL + VPT
The name is "SWANSON CENTER FOR NUTRITION".  ✓

**Document 4:**

Procter & Gamble: Diapers

In an effort to regain branded market share, Procter & Gamble is returning to a brand equity strategy.

1993 Initiatives
• Increase advertising
   -Performance, Value Positionings
   -Major TV Campaign
• Roll-out disposable training pants
• Test 50% thinner diapers (same as KC new product)

16

👤 What is the page no mentioned in this document?

Qwen2-VL
The page number mentioned in the document is 72.  ✗

Qwen2-VL + VPT
16.  ✓

Figure S3. This group of examples shows how the Region Selection Token aids MLLMs in understanding textual information within images by correctly identifying the corresponding regions. The image inputs primarily consist of large but structured documents, such as tables, forms, or letters.

**What is the name of the bar?**

*Qwen2-VL*
The name of the bar is "Sushi Bar." ✗

*Qwen2-VL + VPT*
The name of the bar is "moon bar." ✓

**What is the number on the sign behind the player?**

*Qwen2-VL*
The number on the sign behind the player is 24. ✗

*Qwen2-VL + VPT*
The number on the sign behind the player is 340. ✓





**What brand is the clock?**

*Qwen2-VL*
The clock in the image is a mechanical clock, and it does not appear to be branded. ✗

*Qwen2-VL + VPT*
The brand is Breitling. ✓

**What kind of brick does the box mention?**

*Qwen2-VL*
The box mentions "American Beauty" bricks. ✗

*Qwen2-VL + VPT*
The box mentions "arto" bricks. ✓

Figure S4. This set of images illustrates how the Region Selection Token enables MLLMs to comprehend textual information within real-world scenes by accurately identifying the corresponding regions. The image inputs consist of real-world scenarios, such as signboards and trademarks, where the text occupies only a small portion of the overall scene and is highly susceptible to interference from the surrounding context.

# References

[1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023. 3

[2] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14455–14465, 2024. 3

[3] An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. In *NeurIPS*, 2024. 1, 3

[4] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, and Lidong Bing. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024. 3

[5] 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. *CoRR*, abs/2306.13394, 2023. 10

[6] Qiao Gu, Alihusein Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen, Aditya Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, Chuang Gan, Celso Miguel de Melo, Joshua B. Tenenbaum, Antonio Torralba, Florian Shkurti, and Liam Paull. Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning. *arXiv*, 2023. 3

[7] Xiangteng He and Yuxin Peng. Fine-grained visual-textual representation learning. Technical Report 2, 2020. 6

[8] Yining Hong, Chunru Lin, Yilun Du, Zhenfang Chen, Joshua B. Tenenbaum, and Chuang Gan. 3d concept learning and reasoning from multi-view images. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3

[9] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models, 2023. 3

[10] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019. 5

[11] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. In *arXiv*, 2023. 2, 6

[12] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International Journal of Computer Vision*, 128(7):1956–1981, 2020. 5

[13] Jordy Van Landeghem, Rafał Powalski, Rubèn Tito, Dawid Jurkiewicz, Matthew Blaschko, Łukasz Borchmann, Mickaël Coustaty, Sien Moens, Michał Pietruszka, Bertrand Ackaert, Tomasz Stanisławek, Paweł Józiak, and Ernest Valveny. Document understanding dataset and evaluation (dude). In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 19471–19483, 2023. 6

[14] Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13872–13882, 2024. 1

[15] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *arXiv preprint arXiv:2306.00890*, 2023. 3

[16] Feng Li, Qing Jiang, Hao Zhang, Tianhe Ren, Shilong Liu, Xueyan Zou, Huaizhe Xu, Hongyang Li, Chunyuan Li, Jianwei Yang, Lei Zhang, and Jianfeng Gao. Visual in-context prompting. *CoRR*, 2023. 2

[17] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. 6

[18] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 5

[19] Fangyu Liu, Guy Edward Toh Emerson, and Nigel Collier. Visual spatial reasoning. *Transactions of the Association for Computational Linguistics*, 2023. 5

[20] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 1, 3, 5, 6, 9

[21] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Conference on Neural Information Processing Systems (NeurlPS)*, 2023. 3, 9

[22] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-modal model an all-around player? *CoRR*, abs/2307.06281, 2023. 10

[23] Xinyin Ma, Gongfan Fang, and Xinchao Wang. LLM-pruner: On the structural pruning of large language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 3

[24] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209, 2021. 5

[25] OpenAI. How reasoning works, 2024. Accessed: 2024-11-05. 11

[26] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZ-IZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. 3, 6

[27] Ting Pan, Lulu Tang, Xinlong Wang, and Shiguang Shan. Tokenize anything via prompting. *CoRR*, 2023. 2

[28] 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, 2024. 2, 3, 5, 6, 8, 9

[29] Aleksandar Shtedritski, Christian Rupprecht, and Andrea Vedaldi. What does CLIP know about a red circle? visual prompt engineering for vlms. In *International Conference on Computer Vision (ICCV)*, 2023. 2

[30] Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for image captioningwith reading comprehension. 2020. 5

[31] Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*, pages 8317–8326, 2019. 5

[32] Guangzhi Sun, Wenyi Yu, Changli Tang, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, Yuxuan Wang, and Chao Zhang. video-SALMONN: Speech-enhanced audio-visual large language models. In *Forty-first International Conference on Machine Learning*, 2024. 3

[33] Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao Zhang. SALMONN: Towards generic hearing abilities for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. 3

[34] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 1, 3

[35] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual short-comings of multimodal llms, 2024. 1, 2, 3

[36] Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, and Neel Joshi. Is a picture worth a thousand words? delving into spatial reasoning for vision language models, 2024. 3

[37] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 1, 3, 6

[38] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models, 2023. 3

[39] Xinyi Wang, Lucas Caccia, Oleksiy Ostapenko, Xingdi Yuan, William Yang Wang, and Alessandro Sordoni. Guiding language model reasoning with planning tokens. In *First Conference on Language Modeling*, 2024. 11

[40] Yang Wu, Shilong Wang, Hao Yang, Tian Zheng, Hongbo Zhang, Yanyan Zhao, and Bing Qin. An early evaluation of gpt-4v(ision). *CoRR*, abs/2310.16534, 2023. 4

[41] Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark prompting unleashes extraordinary visual grounding in GPT-4V. *CoRR*, 2023. 2

[42] Lingfeng Yang, Yueze Wang, Xiang Li, Xinlong Wang, and Jian Yang. Fine-grained visual prompting. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2023. 2

[43] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of lmms: Preliminary explorations with gpt-4v(ision), 2023. 3

[44] Yuan Yao, Ao Zhang, Zhengyan Zhang, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. CPT: colorful prompt tuning for pre-trained vision-language models. *CoRR*, 2021. 2

[45] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *CoRR*, 2023. 1

[46] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014. 6

[47] Ping Yu, Jing Xu, Jason Weston, and Ilia Kulikov. Distilling system 2 into system 1, 2024. 11

[48] Runpeng Yu, Weihao Yu, and Xinchao Wang. Api: Attention prompting on image for large vision-language models, 2024. 2

[49] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities, 2023. 6, 9

[50] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023. 3

[51] Chenming Zhu, Tai Wang, Wenwei Zhang, Jiangmiao Pang, and Xihui Liu. Llava-3d: A simple yet effective pathway to empowering lmms with 3d-awareness. *arXiv preprint arXiv:2409.18125*, 2024. 3