

Appendix for PerceptionGPT: Effectively Fusing Visual Perception into LLM

In this appendix, we first introduce the detailed design choices of our perception encoder and decoders, as well as the training losses for perception tasks; then we elaborate more details of the inference procedure that was skipped in the main paper due to space limitation; afterwards, we illustrate more visualization results of our PerceptionGPT on spot captioning and reasoning-based detection and segmentation.

Algorithm 1 Inference Process of PerceptionGPT

```
1: Input: An image  $I$  and a textual instruction  $T$ 
2: Output: Generated tokens and perception signals
3:  $E_{\text{image}} \leftarrow \text{ImageEncoder}(I)$  {Extract visual tokens}
4:  $E'_{\text{image}} \leftarrow \text{ProjectionLayer}(E_{\text{image}})$  {Map to LLM's embedding space}
5:  $E_{\text{concat}} \leftarrow \text{Concatenate}(V', T)$  {Concatenate image and text features}
6:  $E_{\text{new}} \leftarrow \text{LLM}(E_{\text{concat}})$  {Initial input to the LLM}
7:  $\text{Tokens} \leftarrow []$ 
8:  $\text{Perceptions} \leftarrow []$ 
9:  $\text{token} \leftarrow \text{TokenClassification}(E')$ 
10: while token is not <eos> do
11:    $\text{token} \leftarrow \text{TokenClassification}(E_{\text{new}})$ 
12:    $\text{Tokens} \leftarrow \text{Concatenate}(\text{Tokens}, \text{token})$ 
13:   if token is <vis> then
14:      $\text{Signal} \leftarrow \text{PerceptionDecoder}(E')$ 
15:      $E_{\text{percept}} \leftarrow \text{PerceptionEncoder}(\text{Signal})$ 
16:      $\text{Perceptions} \leftarrow \text{Concat}(\text{Perceptions}, \text{Signal})$ 
17:      $E_{\text{concat}} \leftarrow \text{Concatenate}(E, P)$ 
18:   else
19:      $E_{\text{text}} \leftarrow \text{TextEncode}(\text{token})$ 
20:      $E_{\text{concat}} \leftarrow \text{Concatenate}(E, E_{\text{text}})$ 
21:   end if
22:    $E_{\text{new}} \leftarrow \text{LLM}(E_{\text{concat}})$ 
23: end while
24: return  $\text{Tokens}, \text{PerceptionSignals}$ 
```

1. Detailed Design Choices

1.1. Design of Perception Encoders and Decoders

In this section, we introduce the detailed design choices of our perception encoder and decoders. Notably, although we incorporate detection and segmentation tasks in this paper, other perception tasks can be supported in a similar manner

by leveraging encoder and decoder particularly for the new task.

Bounding Box Encoder The bounding box encoder is implemented with a simple 3-layer MLP, with input dimension equal to 4, corresponding to the coordinates of the top-left and bottom-right bounding box coordinates; the output dimension is set to 4096, which equals to the dimension of the token embedding. The intermediate dimension is set to 512.

Bounding Box Decoder Similar to the bounding box encoder, the decoder is also implemented with a simple 3-layer MLP, with input dimension set to 4096 and output dimension set to 4, the intermediate dimension is set to 512.

Mask Encoder The mask encoder is implemented with a ResNet18-like architecture, the intermediate dimension of all blocks are set to 32. The input to the mask encoder should be mask images with 1 channel. The final linear layer maps the flattened image feature to the LLM's token embedding dimension, which is equal to 4096.

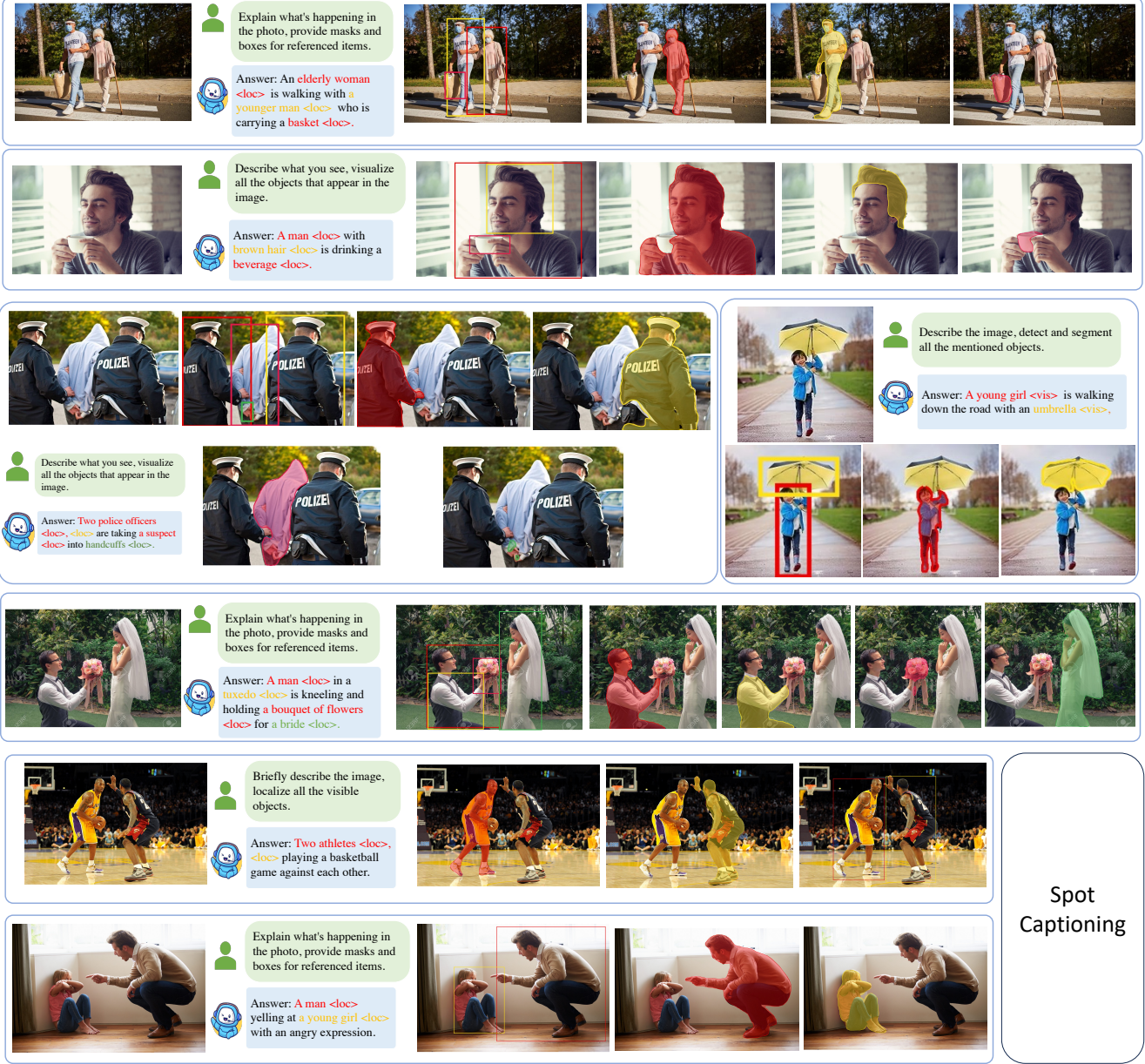
Mask Decoder The design of the mask decoder is inspired by Segment Anything model [1], which consists of two layers of two-way-transformer layers. The decoder block conducts self-attention and cross-attention in two directions (LLM token embedding-to-image embedding and image embedding-to-LLM token embedding). After conducting fusion through two blocks, the updated image embedding is up-sampled and generates the segmentation mask.

1.2. Choices of Perception Task Losses

In this section, we elaborate on the choices of perception task losses. Incorporating such losses during training helps leveraging the prior knowledge and characteristics of the tasks, which further helps learning.

Bounding Box Loss We adopt a combination of L1-norm and GIoU [3] losses, which can be expressed as follows:

$$\mathcal{L}_{\text{box}}(b_{\text{gt}}, S_{\text{in}}, I) = \mathcal{L}_{\text{L1}}(b_{\text{gt}}, b_{\text{pred}}) + \mathcal{L}_{\text{GIoU}}(b_{\text{gt}}, b_{\text{pred}}(S_{\text{in}}, I)) \quad (1)$$



Spot Captioning

Figure 1. Visualization of results from PerceptionGPT on spot captioning task, where the model is asked to generate a description regarding the image content, while detecting and segmenting the mentioned objects at the same time.

The Generalized IoU (GIoU) loss is then defined as:

$$\mathcal{L}_{\text{GIoU}}(b_{\text{gt}}, b_{\text{pred}}) = 1 - \left(\text{IoU} - \frac{\text{area}(C) - \text{area}(b_{\text{pred}} \cup b_{\text{gt}})}{\text{area}(C)} \right) \quad (2)$$

where C is the smallest enclosing box that contains both b_{pred} and b_{gt} .

Mask Loss We adopt a combination of L1-norm and GIoU losses, which can be expressed as follows:

$$\mathcal{L}_{\text{mask}}(m_{\text{gt}}, S_{\text{in}}, I) = \mathcal{L}_{\text{L1}}(m_{\text{gt}}, m_{\text{pred}}) + \mathcal{L}_{\text{DICE}}(m_{\text{gt}}, m_{\text{pred}}(S_{\text{in}}, I)) \quad (3)$$

The DICE loss can be formulated as follows:

$$\mathcal{L}_{\text{DICE}}(m_{\text{gt}}, m_{\text{pred}}) = 1 - \text{Dice} = 1 - \frac{2 \times |m_{\text{pred}} \cap m_{\text{gt}}|}{|m_{\text{pred}}| + |m_{\text{gt}}|} \quad (4)$$

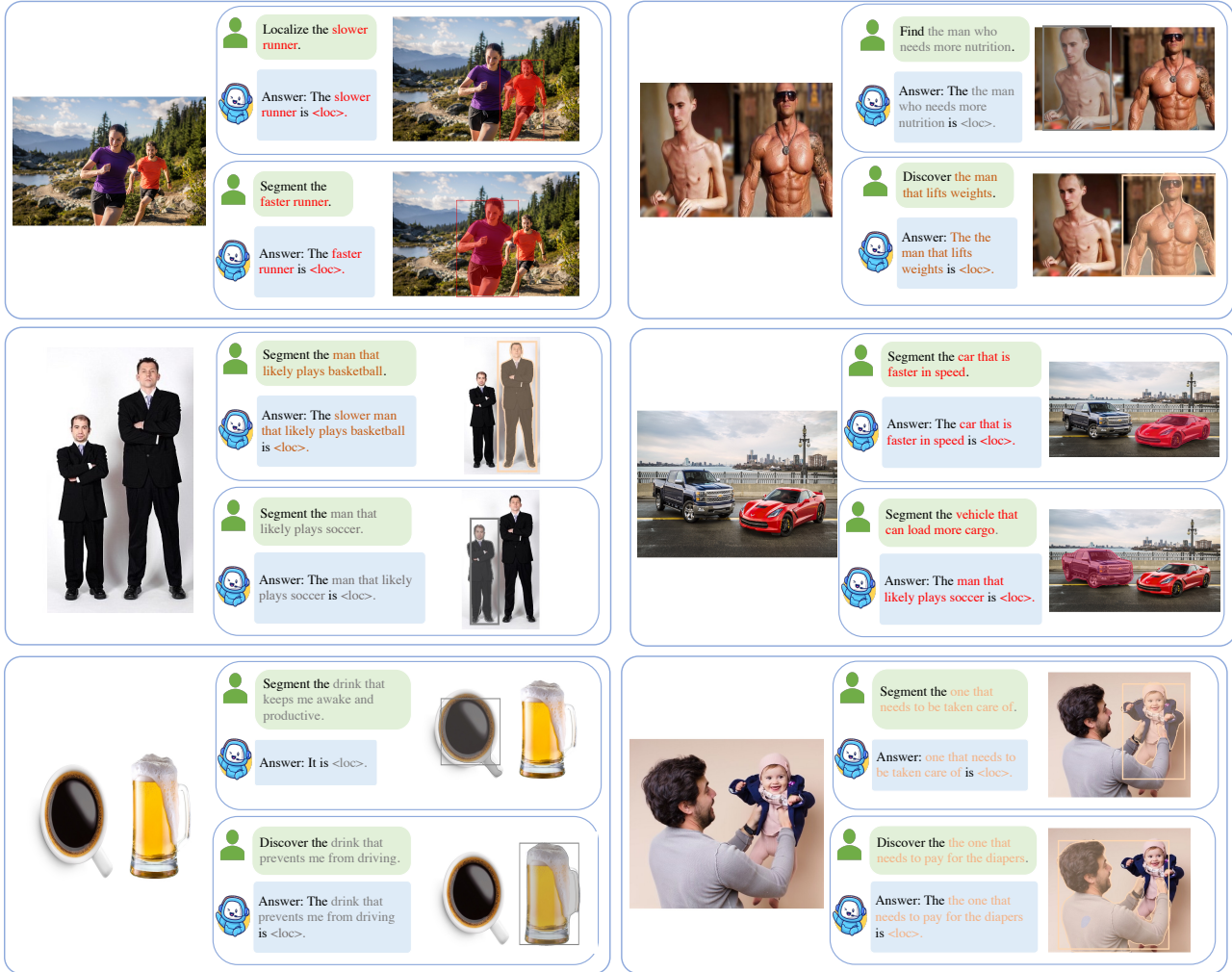


Figure 2. Visualization of results from PerceptionGPT on reasoning-based detection and segmentation task. Our PerceptionGPT is able to detect and segment the object related to the instruction via reasoning, even if the object is not explicitly given by the user.

where

$$\text{Dice} = \frac{2 \times |m_{\text{pred}} \cap m_{\text{gt}}|}{|m_{\text{pred}}| + |m_{\text{gt}}|} \quad (5)$$

2. Detailed Inference Procedure

We demonstrate the detailed inference procedure in Algorithm 1. During inference, given an image and a textual instruction, the image encoder first extracts the visual tokens from the image, which are then mapped to the dimension of LLM’s embedding space via the projection layer. Then, the mapped image features are concatenated with text embeddings to serve as the input to the LLM. Subsequently, the LLM begins to iteratively perform next-token-generation similar to previous VLLMs [2, 4], until the end of sentence token is generated.

Perception Signal as Input. When the input contains a perception signal, our lightweight perception encoder maps it into the embedding space of the large language model (LLM), which is treated as the embedding for $\langle vis \rangle$ token. This embedding is then concatenated with other embeddings before being processed by the LLM.

Perception Signal as Output. During inference, when a token is decoded as $\langle vis \rangle$, its associated embedding is extracted and processed by the perception decoder to reconstruct the original signal.

Notably, during inference, when a $\langle vis \rangle$ is generated, we use the embedding of bounding box in place of $\langle vis \rangle$ to generate the next token, while this choice is arbitrary and the mask embedding can also be used.

References

- [1] 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, 2023. [1](#)
- [2] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. [3](#)
- [3] Hamid Rezaatofghi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression, 2019. [1](#)
- [4] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models, 2023. [3](#)