

# Appendix

## A. Mask2Former+ Implementation Details

Our enhancement of the original Mask2Former model focuses on broadening its segmentation capabilities beyond the 134 common object categories it currently handles, which include 80 things and 55 stuffs as defined in the COCO dataset. The primary goal is to enable the model to recognize an expanded range of object categories, as well as segmentation masks of various levels of granularities, such as semantic parts and visual text regions.

**Training data.** We have compiled a comprehensive dataset by combining multiple existing segmentation datasets. This ensemble encompasses a wide spectrum of entities (things and stuff), their semantic parts, and visual text, drawn from sources such as COCO [5], LVIS [25], Entity-v2 [60], Pascal [16], PACO [63], MHP-v2 [40], and TextOCR [67]. The resulting dataset comprises over 200K images and 4.5M masks, as summarized in Table 11. Notably, the annotations from COCO, LVIS, and PACO are based on a shared set of COCO images. We merged these annotations to ensure comprehensive mask proposal coverage, thereby providing holistic instance coverage within each image, as can be illustrated in Figure 7.

Dataset Name	Split	Granularity			Dataset Size	
		Entity	Part	Text	#Image	#Masks
LVIS [25] & PACO [63]	part no.part	✓ ✓	✓		15,089 103,178	596,687 2,062,536
Entity-v2 [60]	cls	✓			31,913	579,076
Pascal [18]	train val	✓ ✓	✓ ✓		4,998 5,105	93,322 95,462
MHP [40]	train		✓		15,403	410,113
TextOCR [67]	train			✓	21,749	714,770
Total		✓	✓	✓	197,435	4,551,966

Table 11. Summary of the training datasets for Mask2Former+. Entity includes both thing and stuff categories.

**Model.** Building on the foundation of the original Mask2Former [10], we developed Mask2Former+, a panoptic segmentation model designed for multi-grained segmentation. We initialize our model from the Mask2Former checkpoint with the Swin-L backbone pre-trained on the COCO panoptic segmentation dataset [34]. Besides the 200 entity queries that are trained for thing and stuff proposals, we added 50 additional expert queries for the segmenting parts and the visual text regions, respectively. Given that not all images have annotations for every type of segmentation (for instance, the TextOCR dataset provides annotations only for visual text regions), our model computes the group-wise matching loss exclusively for the annotations available in each dataset. This approach ensures that the model

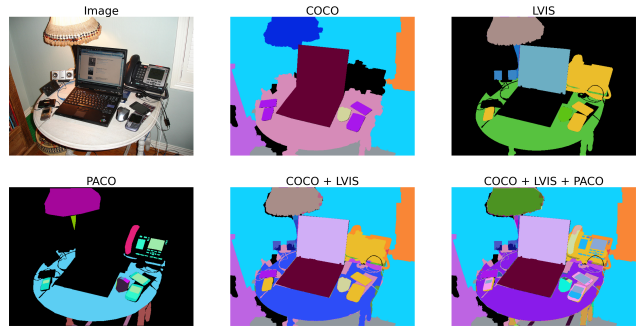


Figure 7. Illustrations of the merged segmentation annotations from COCO Panoptic, LVIS, and PACO datasets.

benefits from partial annotations without compromising its ability to recognize other levels of granularity when certain annotations are unavailable. Although most samples in our dataset also have semantic annotations such as object categories, we do not use them but only train the model for class-agnostic mask proposals. We train the model for 20k iterations on our combined segmentation dataset with a batch size of 16 using the Detectron2 library [77].<sup>1</sup>

## B. The M3G2 Dataset

In this section, we introduce the M3G2 dataset with Multi-Modal Multi-Grained Grounding. M3G2 is a comprehensive dataset consisting of 36 sub-problems, derived and augmented from 27 existing datasets with grounded vision-language annotations. The dataset is categorized into four main types: (1) Grounded Image Captioning (GIC), (2) Grounded Visual Question Answering (GVQA), (3) Referential Expression Segmentation (RES), and (4) Referential Dialog (RD). Details on the dataset sources, image origins, types of grounding annotations, semantic granularity, and data statistics are summarized in Table 12. All datasets are formatted into the conversation format between a human user and a model assistant, where the user provides task objectives as instructions, and model responses are generated automatically based on the annotations.

**Grounded Image Captioning (GIC).** GIC focuses on generating image captions that ground to visual entities presented in the image. We incorporate the Panoptic Narrative Grounding (PNG) [34] and Flickr30K-Entity [58] datasets. PNG, derived from Localize Narrative [59] and COCO Segmentation [5], provides long and detailed narratives with an average of 36.5 words per description, exemplified in Figure 15a. These narratives are rich in detail, offering a high coverage of the visual content including the background. Flickr30K-Entity, offering concise captions with box annotations, complements PNG with its larger vocabulary and finer granularity, as shown in 15b. The example instruction templates used to construct the conversation are listed in Ta-

<sup>1</sup><https://github.com/facebookresearch/detectron2>

Metadata			Grounding Annotations			Semantic Granularity					Data Size	
Task Type	Dataset Name	Image Source	Mask	Box	Pointer	Thing	Stuff	Part	Multi.	Text	Train	Val / Test
Grd. Captioning (GCAP)	PNG	COCO	✓	✓		✓	✓		✓		132,045	8,435
	Flickr30K-Entity	Flickr30K	✓	✓		✓	✓	✓	✓		148,915	1,000 / 1,000
Referential Expression Segmentation (RES)	RefCOCO	COCO	✓	✓		✓					113,311	-
	RefCOCO+	COCO	✓	✓		✓					112,441	-
	RefCOCog	COCO	✓	✓		✓					80,322	-
	RefCLEF	ImageCLEF	✓	✓		✓					104,531	-
	gRefCOCO	COCO	✓	✓		✓					194,233	-
	PhraseCut	VG	✓	✓		✓	✓	✓	✓		84,688	-
	D-Cube	GRD	✓	✓		✓			✓		9,499	-
	ReasonSeg	OpenImages & ScanNetV2	✓	✓		✓	✓	✓	✓		1,315	344
	RIO	COCO	✓	✓		✓			✓		27,696	34,170
	SK-VG	VCR	✓	✓		✓					23,404	-
Grounded Visual Question Answering (GVQA)	VizWiz-Grounding	VizWiz	✓	✓		✓	✓			✓	6,494	1,131 / 2,373
	TextVQA-X	OpenImages	✓	✓						✓	14,476	3,620
	GQA	VG		✓		✓	✓	✓	✓		301,623	-
	VQS	COCO		✓		✓					20,380	8,203
	Shikra-BinaryQA	Flickr30K		✓		✓	✓	✓	✓		4,044	1,159
	EntityCount	Entity-v2	✓	✓		✓	✓	✓	✓		11,088	453
	FoodSeg-QA	Recipe1M		✓		✓			✓		7,114	-
	LVIS-QA	COCO	✓	✓		✓	✓	✓			94,860	3,611
Referential Dialog (RD)	RefCOCO-REG	COCO	✓	✓	✓	✓					17,395	-
	RefCOCO+-REG	COCO	✓	✓	✓	✓					17,383	-
	RefCOCog-REG	COCO	✓	✓	✓	✓					22,057	-
	gRefCOCO-REG	COCO	✓	✓	✓	✓					20,282	-
	VG-SpotCap	VG		✓	✓	✓	✓	✓	✓		247,381	232,935
	V7W	COCO		✓	✓	✓					22,805	10,193 / 57,265
	PointQA-Local	VG			✓	✓					27,426	4,855 / 4,880
	PointQA-Twice	VG			✓	✓					36,762	14,668 / 5,710
	VCR-Open	VCR		✓	✓	✓					58,340	-
	VCR-Multichoice	VCR		✓	✓	✓					97,648	26,534 / 25,263
	ShikraRD	Flickr30K		✓	✓	✓	✓	✓	✓		1,878	-
	SVIT-RD	VG		✓	✓	✓	✓	✓	✓		32,571	-
	Guesswhat-Guesser	COCO	✓	✓	✓	✓					92,136	19,665
	Guesswhat-Oracle	COCO	✓	✓	✓	✓					101,256	21,643
	VG-RefMatch	VG		✓	✓	✓	✓	✓	✓		247,381	-
HierText	OpenImages	✓	✓	✓					✓	6,058	3,885	

Table 12. The full list of datasets used in M3G2.

Instruction Templates for Brief Captioning	Instruction Templates for Detailed Captioning
Describe the image briefly.	Describe the image in detail.
Describe the image in a few words.	Describe the picture's every detail.
Describe the image in a short sentence.	Describe the given picture in very detail.
Describe the image in a clear and concise manner.	Make a fine description of the image.
Generate a short caption for the picture.	Generate a long caption for the given image.
Caption the image in a few words.	Give me a detailed caption of this image.

Table 13. Instruction templates for the GIC task.

Instruction Templates For RES	
Highlight "{}" in the image.	Segment "{}" in the image.
Segment: {}.	Help me segment out {}.
Localize "{}" in the image.	Help me localize {}.
Help me highlight the region of {}.	Demonstrate where "{}" is located in this image.
Show me where to find in this photo.	Identify and mark the region of {} for me.
Can you highlight "{}"?	Can you extract the segment: {} for me?
Can you localize "{}" in this image?	Could you please segment out {} in the image?

Table 14. Templates used for the RES task.

ble 13, where we use key words such as "short/briefly" and "in detail" to distinguish between short and long captioning.

**Referential Expression Segmentation (RES).** RES is a task combining language understanding with precise visual segmentation. Our dataset includes 10 diverse sources. To improve the learning efficiency and enhance contextual understanding, we format queries from the same image into a simulated multi-turn dialog, as illustrated in Figures 16 and 17. We employ the widely used RefCOCO+/-g datasets [32, 53] and RefCLEF [65] for single-object RES. gRefCOCO [43] is employed for multi-object and negative queries. To enhance the visual diversity, we also incorporate PhraseCut [76] and D-Cube [79] that use an image source different than COCO. Additionally, ReasonSeg [37], RIO [61], and SK-VG [78] are included, where a textual context is given and the models need to not only understand

that context, but also equips with a certain degree of commonsense knowledge to successfully solve the query, such as shown in Figure 17b, 17c and 17d. The dialogue templates are listed in Table 14.

**Grounded Visual Question Answering (GVQA).** The GVQA task extends the visual question answering by additionally requiring visual grounding of the answer. We include 8 datasets for the grounded VQA task in M3G2. First, we collect and organize some existing datasets that can directly fit into our grounded vision-language task framework, including VizWiz-Grounding [6], TextVQA-X [64], GQA [29], VQA [19] and Shikra-BinaryQA [8] (Figure 18). To further improve the data scale and visual concept coverage, we enlarge the GVQA collection by re-purposing existing panoptic segmentation datasets with templated instruc-

<p>Instruction Templates For Short Response VQA.</p> <p>{ } Answer with a single word or a short phrase.</p> <p>Given the image, answer the question "{ }" with a single word or a short phrase.</p> <p>Give a short answer to the question "{ }" based on the image.</p>
<p>Instruction Templates For Chain-of-Thought Response VQA.</p> <p>{ } Let's think step by step.</p> <p>{ } Please include the reasoning process.</p> <p>{ } Before giving the answer, please explain your reasoning.</p> <p>{ } Explain your logic before giving the answer.</p> <p>Please answer the following question "{ }", and describe your thought process.</p>
<p>Instruction Templates For Grounding Answer to Masks.</p> <p>Show where in the image you found your answer.</p> <p>Mark the part of the image that supports your answer.</p> <p>Please highlight your evidence in the image.</p> <p>Point out the evidence from the image.</p> <p>Indicate the area in the image that justifies your response.</p> <p>Highlight the section of the image that backs up your answer.</p> <p>Shade the section of the image that confirms your reply.</p> <p>Emphasize the part of the image that relates to your answer.</p>
<p>Instruction Templates For Object Presence QA.</p> <p>Is { } present in the image?</p> <p>Is there any { } in this image?</p>
<p>Instruction Templates For Object Counting QA.</p> <p>How many { } can you see in this image?</p> <p>Count the number of { }.</p>
<p>Instruction Templates For Object Segmentation Request.</p> <p>Segment { }.</p> <p>Highlight all the { } in this image.</p> <p>Show me all the { } presented in the picture.</p>

Table 15. Templates used for the GVQA task.

tions and model responses. Specifically, based on the annotations from LVIS [25] and EntityV2 [60], we design questions about object presence, object counting, and segment query with a possibly negative request (i.e. the target object does not exist in the image), for the model to learn to recognize a diverse set of concepts more faithfully. See Figure 19 for examples of such multi-turn QA, and example question templates used in Table 15.

**Referential Dialog (RD)** . RD features multi-modal conversations where the user can refer to objects or regions in the image by a spatial prompt (e.g. a bounding box). We include various types of RD in our dataset and the templates used are listed in Table 16. First, we add several existing RD datasets such as V7W [94], PointQA [51], VCR [88], ShikraRD [8] and SVIT [90] without much modifications. We then revisit the RefCOCO series [43, 52, 86] for referential expression generation, where the referred object is given and the goal is to generate a unique description that leads to that object. We use the region caption annotations from the VG dataset [36] for region captioning and a region-matching game. We select a set of region pointers and several descriptions to provide to the model, and the goal is to match the pointed regions with the descriptions (Figure 21b). We repurpose the GuessWhat dataset [17] to make it fit into our RD formulation, as shown in Figure 21a. We also construct a referred text reading task based on the HierText [47] dataset and enhance the model’s capability of text recognition, as shown in Figure 21c.

<p>Instruction Templates For REG.</p> <p>Provide a distinct description for that &lt;PTR&gt;</p> <p>Describe the selected area in a unique way. &lt;PTR&gt;</p> <p>Share a unique description of the region &lt;PTR&gt;</p> <p>Offer a one-of-a-kind descriptor &lt;PTR&gt;</p> <p>Describe the selected area &lt;PTR&gt; uniquely.</p> <p>Point out &lt;PTR&gt; in the picture with a unique description.</p> <p>Tell me how &lt;PTR&gt; stands out in the photo.</p> <p>Use your words to highlight just &lt;PTR&gt; in the image.</p> <p>Please describe &lt;PTR&gt; in the image in a way that it can be uniquely identified.</p> <p>If you had to describe just &lt;PTR&gt; to someone, how would you do it?</p> <p>What makes &lt;PTR&gt; different from everything else in the picture?</p> <p>How can you describe &lt;PTR&gt; in the image in a way that it can be uniquely identified?</p> <p>Can you provide a referring expression for &lt;PTR&gt; such that it sets it apart from others?</p> <p>Let's play a game! Describe &lt;PTR&gt; in the photo so I can find it.</p>
<p>Instruction Templates For Region Captioning.</p> <p>Describe it &lt;PTR&gt;.</p> <p>Describe the region &lt;PTR&gt; in a few words.</p> <p>Describe the region &lt;PTR&gt; in a short phrase.</p> <p>Describe the selected area &lt;PTR&gt;.</p> <p>Provide a brief description of this part &lt;PTR&gt;.</p> <p>Give a short caption for this &lt;PTR&gt;.</p> <p>Provide a brief description of the area marked &lt;PTR&gt;.</p> <p>Tell me about the contents in the selected zone &lt;PTR&gt;.</p> <p>Provide a concise description for this spot &lt;PTR&gt;.</p> <p>Narrate what you see in the indicated area &lt;PTR&gt;.</p> <p>What is in the region &lt;PTR&gt;? Describe in a phrase.</p> <p>What can you see in this area &lt;PTR&gt;?</p> <p>How would you describe the content at &lt;PTR&gt;?</p> <p>How would you caption this particular region &lt;PTR&gt;?</p> <p>What's depicted in the marked area &lt;PTR&gt;?</p>

Table 16. Templates used for the RD task.

## C. GROUNDHOG Implementation Details

**Data Balancing.** In constructing the M3G2 dataset, we recognized the need to address the varying scales of the multiple constituent datasets to ensure a balanced data distribution during training. To achieve this, we have implemented dataset-specific sampling strategies, adjusting the volume of data from each source dataset through either up-sampling or down-sampling. The ratios we applied are as follows:

- PNG: up-sampled by a factor of 2.
- Flickr30k-Entities: up-sampled by 1.5 times.
- RefCOCO+: up-sampled by 1.5 times.
- RefCOCOg: up-sampled by 1.5 times.
- SK-VG: up-sampled by a factor of 2.
- Dcube (multiturn): up-sampled by a factor of 10.
- ReasonSeg: up-sampled by a factor of 10.
- Shikra-Binary: up-sampled by a factor of 10.
- VCR-Open (multiturn): down-sampled by half.
- VCR-Multiturn: down-sampled to 10%.
- VizWiz: up-sampled by a factor of 3.
- LVIS-QA: down-sampled by half.
- TextVQAX: up-sampled by a factor of 2.
- EntityCount: up-sampled by a factor of 2.
- VG-SpotCap: down-sampled by half.
- Shikra-RD: up-sampled by a factor of 10.
- HierText: up-sampled by a factor of 5.
- GuessWhat-Oracle: down-sampled to 20%.
- GuessWhat-Guesser: down-sampled to 20%.
- SVIT: up-sampled by a factor of 3.

The balanced sampled dataset contains 1.8 million samples in total.



(a) Grounded short caption generation on Flickr30K-Entity. While only box supervisions are available for this dataset, GROUNDHOG generalize to pixel-level grounding after joint training on M3G2.

(b) Grounded detailed narrative generation on PNG. GROUNDHOG successfully generalize to grounding a novel category *watch* in the generated caption, which is not included in the 80 categories of PNG annotation.

Figure 8. Examples of GROUNDHOG’s performance in grounded image captioning.

**Learning from Both Box and Mask Supervision.** In the M3G2 dataset, not all sub-datasets include mask supervision, necessitating a hybrid loss approach to effectively benefit from grounded supervision from both mask and box annotations. We address this by employing different loss functions based on the type of annotation available. When mask annotations are available, we apply the dice loss  $\mathcal{L}_{\text{dice}}$  and binary cross-entropy loss  $\mathcal{L}_{\text{bce}}$  between the predicted grounding masks and the ground truth masks of each phrase, following Cheng et al. [10]. In cases where only box annotations are present, we apply the projection loss  $\mathcal{L}_{\text{proj}}$  as introduced by Tian et al. [69], which selects the mask whose projection on the axis matches the best with the annotated box. Essentially, this can be seen as a 1D dice loss calculated between the projected masks and the edges of the ground truth boxes along both the  $x$  and  $y$  axes. Given that the primary objective of grounding is to accurately select the correct mask, we assign different weights to these loss components. The mask dice loss and box projection loss are both weighted at 1, while the mask bce loss is given a lower weight of 0.1. The final loss calculation is a summation of the language modeling loss  $\mathcal{L}_{\text{lm}}$  and these mask-related losses.

**LLM Configuration** We adopt the Vicuna-7B model [11] as our base LLM, and use the OpenAI CLIP@336 [62] model and DINOv2-L/14-reg[15] pretrained checkpoints. We use the original conversation template from Vicuna, where all the interactions are formatted as `<system_message> <s> USER: <utterance> ASSISTANT: <utterance> </s>`.

**Parameter-Efficient Training.** We adopt the LLaMA2-7B model [70] as our base LLM, and initialized the weight from LLaVA-1.5 [44]. For the vision encoders, we use the OpenAI CLIP@336 [62] model and DINOv2-L/14-reg [15] pretrained checkpoints. We freeze all the parameters of the

Mask2Former+, the CLIP, and the DINOv2 model during training. We use Low-Rank Adaptation (LoRA) [28] with  $r = 16$  and  $\alpha = 16$  to tune the LLM, including all the linear layers, input embeddings, and the LM head. We train all the new components introduced for connecting these models, including the MLP projection layer of CLIP and DINOv2, and the mask retrieval head. As a result, less than 2% of the total parameters are trainable in the whole model. We use the AdamW optimizer [48] with an initial learning rate of  $2e-4$  and a cosine annealing rate. We train our model on the balanced sampled M3G2 dataset for 2 epochs, which takes around 2 days using 8 40G A100 GPUs.

## D. Additional Qualitative Results

We provide additional qualitative examples generated by GROUNDHOG. See Figure 9 & 10 for GIC, Figure 11 for RES, Figure 12 for GVQA, and Figure 13 & 14 for RD, respectively.







Demonstrate where "yellow sauce" is located in this image.

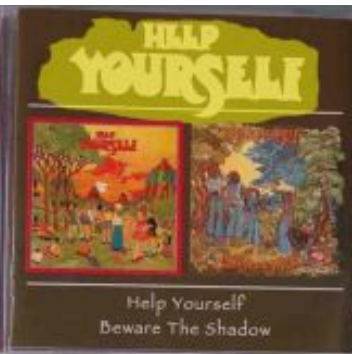
Yellow sauce.



Highlight "not the most prominent umbrella but the one next to it" in the image.

Not the most prominent umbrella but the one next to it.

Figure 11. Additional qualitative examples for referring expression segmentation.



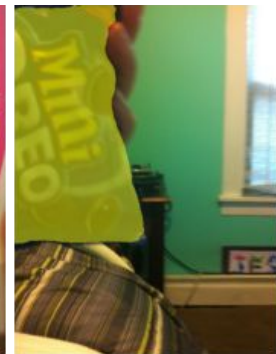
What does the first page of the right say? Mark the part of the image that supports your answer.

Help yourself is printed on the cover.



What is this product?

Toffee cheesecake.



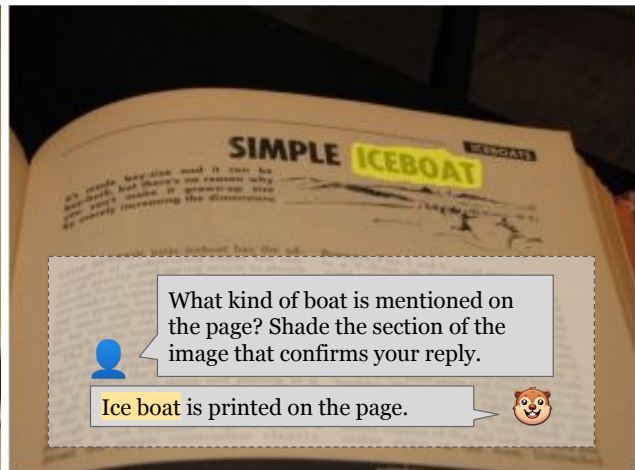
what type of cookies are these?

Oreo.



What was the picture about? Indicate the area in the image that justifies your response.

Vitamix 3600 is printed on the image.

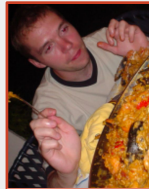


What kind of boat is mentioned on the page? Shade the section of the image that confirms your reply.

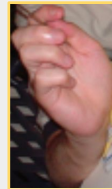
Ice boat is printed on the page.

Figure 12. Additional qualitative examples for grounded VQA of visual text.





<PTR\_0>



<PTR\_1>



<PTR\_2>

Is the man <PTR\_0> touching the food <PTR\_1> directly with his hand <PTR\_2>? Include your thinking process as well.

The man is not touching the food directly with his hand. Instead, he is using a spoon to serve himself a portion of the food.



<PTR>

Provide a distinct clue for that <PTR>.

25. Is that a good description?



<PTR>

Is the rock wall <PTR> in a vertical or horizontal orientation? Please include the reasoning process.

The rock wall is in a vertical orientation, as the man is climbing up it.



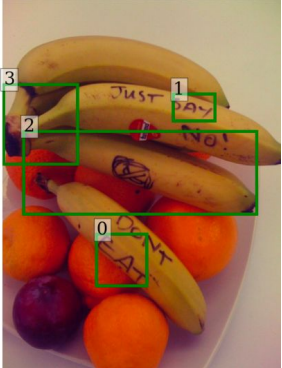
<PTR>

Let's play a game! Describe <PTR> in the photo so I can find it

The boy in the back wearing a white shirt. Did I capture its essence?




Figure 13. Additional qualitative examples for referential dialogue with pointer input.



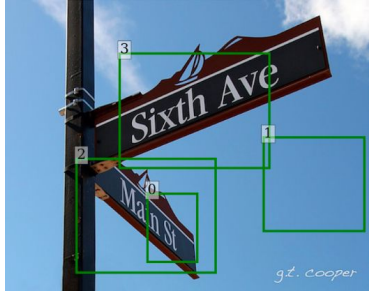
The answer is C. 🐾

Which banana has a drawing of a banana that's crossed out? Identify the best answer from these options:  
A. <PTR\_0> B. <PTR\_1> C. <PTR\_2> D. <PTR\_3>




The answer is C. 🐾

Which part of the restaurant helps customers see what food they offer? Select the most likely answer from:  
A. <PTR\_0> B. <PTR\_1> C. <PTR\_2> D. <PTR\_3>



The answer is A. 🐾

Which banana has a drawing of a banana that's crossed out? Identify the best answer from these options:  
A. <PTR\_0> B. <PTR\_1> C. <PTR\_2> D. <PTR\_3>

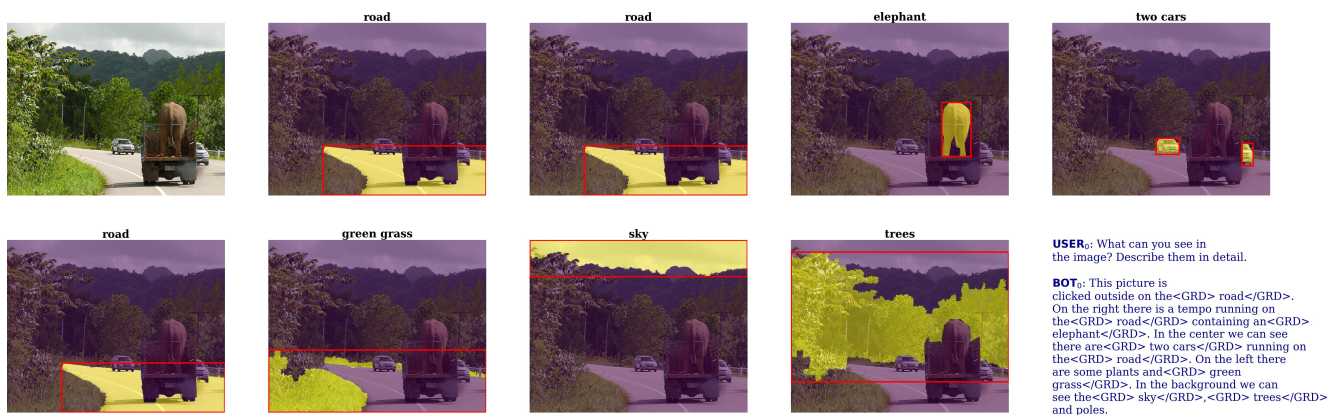


The answer is D. 🐾

Which street name is a number? From the options given, pick the most suitable answer.  
A. <PTR\_0> B. <PTR\_1> C. <PTR\_2> D. <PTR\_3>

Figure 14. Additional qualitative examples for referential dialogue with pointers as multiple choices input.



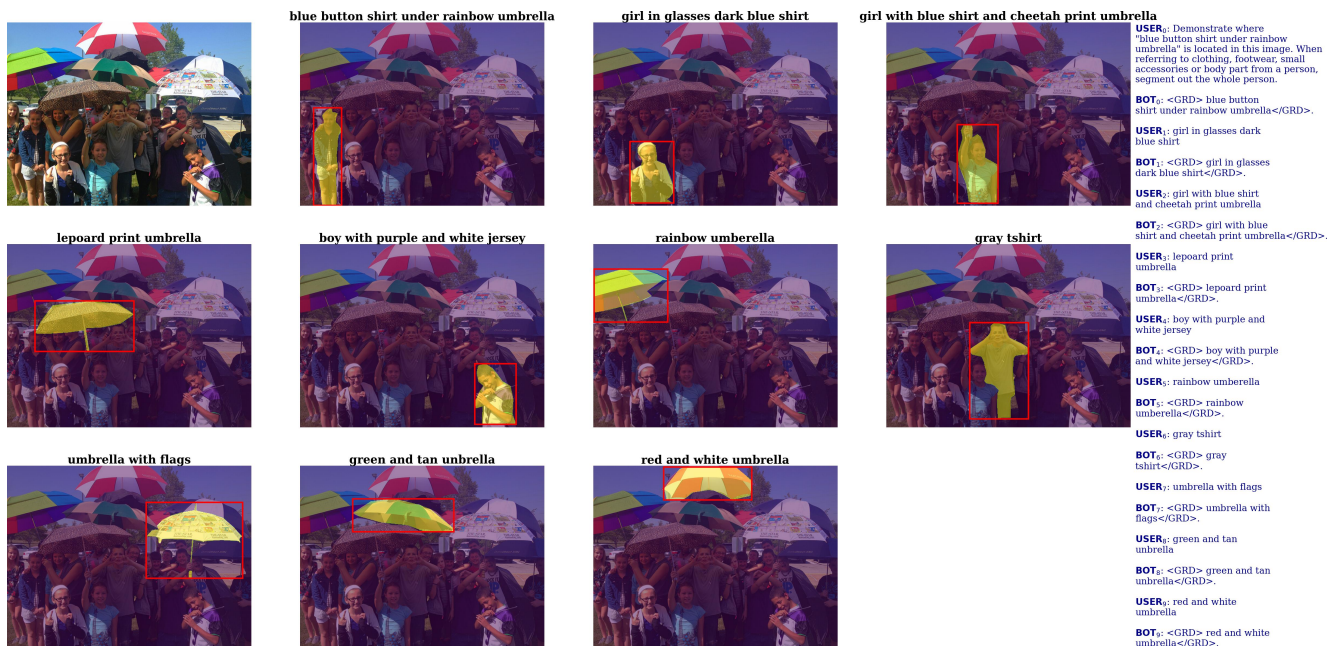


(a) Example of the detailed captioning task originated from PNG.



(b) Example of the short captioning task originated from Flickr30K-Entity.

Figure 15. Data samples for the GCAP task.

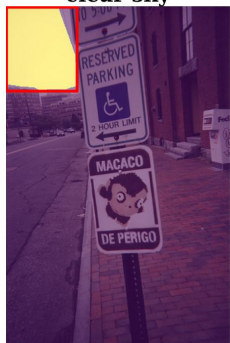
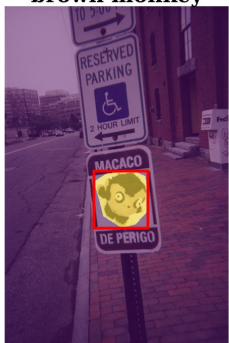
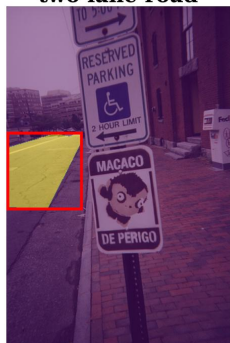


(a) Example of the referential expression segmentation task originated from RefCOCO+.

**two lane road**

**brown monkey**

**clear sky**



**USER<sub>0</sub>**: Can you highlight "two lane road"? Try to segment the region that most accurately matches the description.

**BOT<sub>0</sub>**: Here it is:<GRD> two lane road</GRD>.

**USER<sub>1</sub>**: brown monkey

**BOT<sub>1</sub>**: Here it is:<GRD> brown monkey</GRD>.

**USER<sub>2</sub>**: clear sky

**BOT<sub>2</sub>**: Here it is:<GRD> clear sky</GRD>.

(b) Example of the referential expression segmentation task originated from PhraseCut.

**right lady**

**woman in black shirt and woman holding wine glass on the right**

**hat man**



**USER<sub>0</sub>**: Show me where to find right lady in this photo. When referring to clothing, footwear, small accessories or body part from a person, segment out the whole person.

**BOT<sub>0</sub>**: Here it is:<GRD> right lady</GRD>.

**USER<sub>1</sub>**: woman in black shirt and woman holding wine glass on the right

**BOT<sub>1</sub>**: Here it is:<GRD> woman in black shirt and woman holding wine glass on the right</GRD>.

**USER<sub>2</sub>**: hat man

**BOT<sub>2</sub>**: Here it is:<GRD> hat man</GRD>.

**USER<sub>3</sub>**: woman in black shirt

**BOT<sub>3</sub>**: Here it is:<GRD> woman in black shirt</GRD>.

**USER<sub>4</sub>**: woman on left and hat man

**BOT<sub>4</sub>**: Here it is:<GRD> woman on left and hat man</GRD>.

**USER<sub>5</sub>**: orange or red in middle

**BOT<sub>5</sub>**: No region can be identified as "orange or red in middle".

**USER<sub>6</sub>**: four people in front

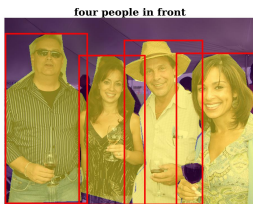
**BOT<sub>6</sub>**: Here it is:<GRD> four people in front</GRD>.



**woman in black shirt**



**woman on left and hat man**

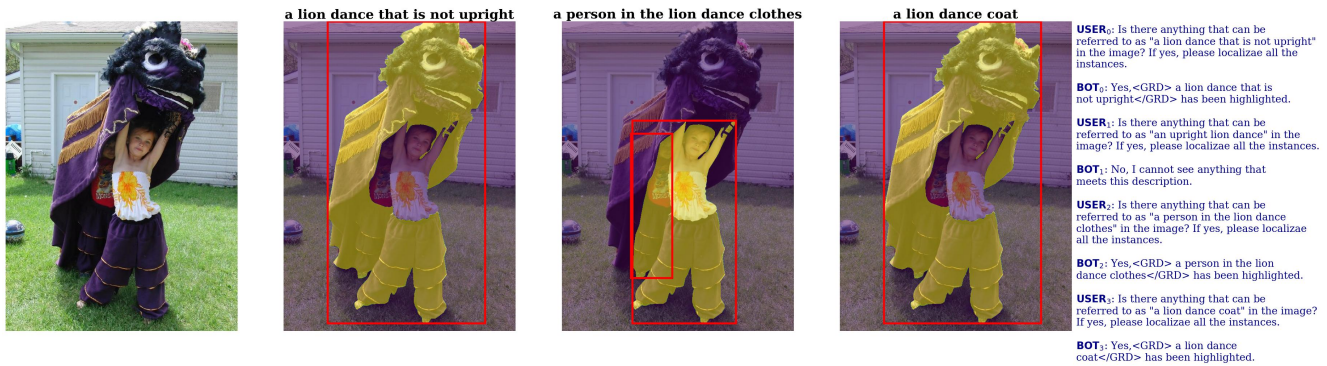


**four people in front**

(c) Example of the generalized referential expression segmentation task originated from gRefCOCO.

Figure 16. Data samples for the RES task (part 1).





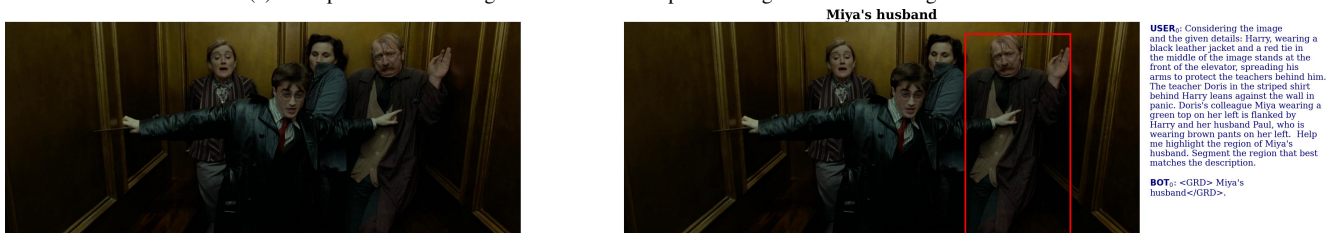
(a) Example of the generalized referential expression segmentation task originated from D-Cube.



(b) Example of the reasoning-based referential expression segmentation task originated from ReasonSeg.



(c) Example of the reasoning-based referential expression segmentation task originated from RIO.



(d) Example of the reasoning-based referential expression segmentation task originated from SK-VG.

Figure 17. Data samples for the RES task (part 2).

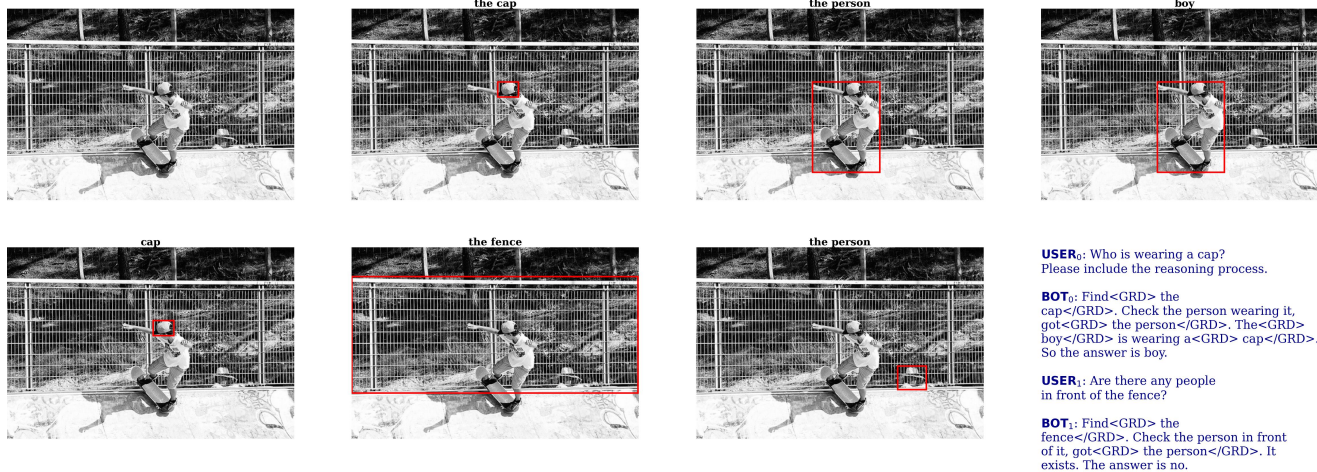




(a) Example of grounded VQA originated from VizWiz-Grounding.



(b) Example of grounded VQA originated from TextVQA-X.



(c) Example of grounded VQA with chain-of-thought reasoning originated from GQA.



(d) Example of grounded VQA task originated from VQS.



(e) Example of the grounded VQA with chain-of-thought reasoning originated from Shikra-BinaryQA.

Figure 18. Data samples of the GVQA task (part 1).



	<b>painting</b> 	<b>towel</b> 	<b>cabinets</b> 	<p><b>USER<sub>0</sub></b>: Have you spotted any sky in this picture?</p> <p><b>BOT<sub>0</sub></b>: After taking a close look at the image, I couldn't find any trace of sky.</p> <p><b>USER<sub>1</sub></b>: Can you give me a count of the women present in the image?</p> <p><b>BOT<sub>1</sub></b>: I've searched the image high and low, but no woman can be found.</p> <p><b>USER<sub>2</sub></b>: Have you spotted any painting in this picture?</p> <p><b>BOT<sub>2</sub></b>: Yes, the only &lt;GRD&gt; painting&lt;/GRD&gt; is highlighted in the image.</p> <p><b>USER<sub>3</sub></b>: Can you see any floor in this image?</p> <p><b>BOT<sub>3</sub></b>: After taking a close look at the image, I couldn't find any trace of floor.</p> <p><b>USER<sub>4</sub></b>: Can you see any towel in this image?</p> <p><b>BOT<sub>4</sub></b>: Certainly, there's a single &lt;GRD&gt; towel&lt;/GRD&gt;.</p> <p><b>USER<sub>5</sub></b>: Does the image contain any cabinets?</p> <p><b>BOT<sub>5</sub></b>: Certainly, there's a single &lt;GRD&gt; cabinets&lt;/GRD&gt;.</p> <p><b>USER<sub>6</sub></b>: How many boxes are there?</p> <p><b>BOT<sub>6</sub></b>: Let me count! There are &lt;GRD&gt; 7 boxes&lt;/GRD&gt; in total.</p> <p><b>USER<sub>7</sub></b>: How many wineglass can you identified in this image?</p> <p><b>BOT<sub>7</sub></b>: The count comes up to &lt;GRD&gt; 9 wineglass&lt;/GRD&gt;.</p> <p><b>USER<sub>8</sub></b>: How many cups are in the image?</p> <p><b>BOT<sub>8</sub></b>: The only &lt;GRD&gt; cup&lt;/GRD&gt; is highlighted in the image.</p> <p><b>USER<sub>9</sub></b>: Can you see any ceiling in this image?</p> <p><b>BOT<sub>9</sub></b>: Yes, here it is: &lt;GRD&gt; ceiling&lt;/GRD&gt;.</p>
<b>7 boxes</b> 	<b>9 wineglass</b> 	<b>cup</b> 	<b>ceiling</b> 	

(a) Example of the grounded VQA task originated from EntityCount.

	<b>onion</b> 	<b>potato</b> 	<b>steak</b> 	
<b>corn</b> 	<b>pork</b> 	<b>ice cream</b> 		<p><b>USER<sub>0</sub></b>: Detail the ingredients for this dish for me.</p> <p><b>BOT<sub>0</sub></b>: Here's what goes into this dish: &lt;GRD&gt; onion&lt;/GRD&gt;, &lt;GRD&gt; potato&lt;/GRD&gt;, &lt;GRD&gt; steak&lt;/GRD&gt;, &lt;GRD&gt; corn&lt;/GRD&gt;, &lt;GRD&gt; pork&lt;/GRD&gt;, and &lt;GRD&gt; ice cream&lt;/GRD&gt;.</p>

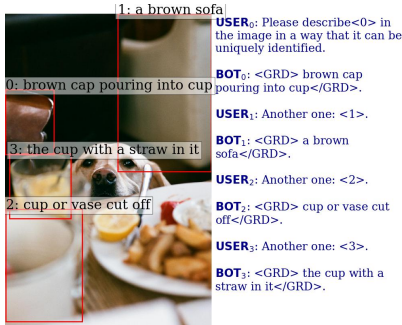
(b) Example of the grounded VQA task originated from FoodSeg-QA.

	<b>camera</b> 	<b>clock tower</b> 	<b>pants</b> 	<p><b>USER<sub>0</sub></b>: Can you show me all the hats in the image?</p> <p><b>BOT<sub>0</sub></b>: Emm... It seems that there is no hat in the image. Do you want to try another object?</p> <p><b>USER<sub>1</sub></b>: How many cameras are there?</p> <p><b>BOT<sub>1</sub></b>: There are &lt;GRD&gt; 3 cameras&lt;/GRD&gt; in the image.</p> <p><b>USER<sub>2</sub></b>: Segment clock tower.</p> <p><b>BOT<sub>2</sub></b>: Sure. The &lt;GRD&gt; clock tower&lt;/GRD&gt; is highlighted.</p> <p><b>USER<sub>3</sub></b>: How many pants can you see?</p> <p><b>BOT<sub>3</sub></b>: I can see &lt;GRD&gt; 9 pants&lt;/GRD&gt; in this photo.</p>
--	-------------------	------------------------	------------------	--

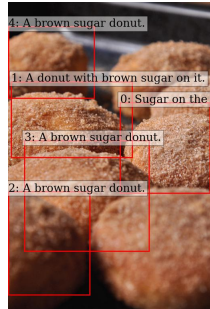
(c) Example of the grounded VQA task originated from LVIS-QA.

Figure 19. Data samples of the GVQA task (part 2).

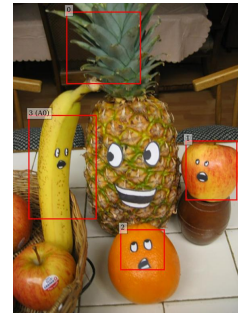




(a) Refer expression generation (RefCOCO+).



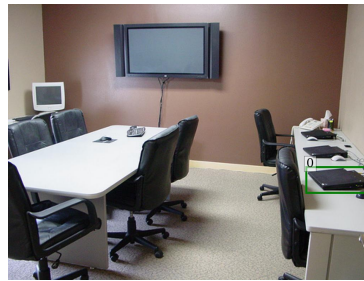
(b) Region captioning (VG).



(c) Example of the referring QA task from V7W.



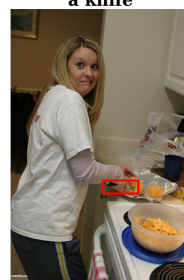
(d) Example of the referring QA task originated from PointQA-Local.



(e) Example of the referring QA task originated from PointQA-Twice.



(f) Example of the referring QA task originated from VCR.



(g) Example of the referential dialog task originated from ShikrARD.



(h) Example of the referential dialog task originated from SVIT.

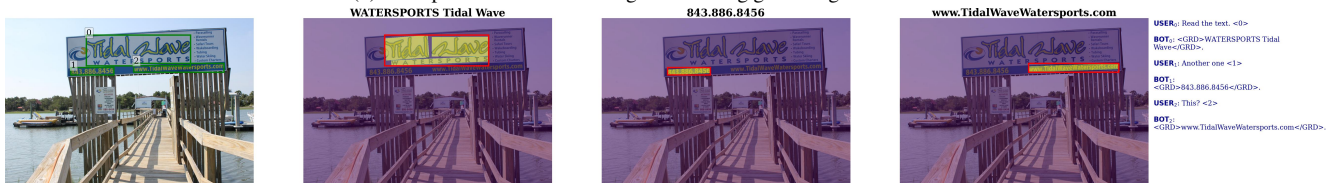
Figure 20. Data samples of the RD task (part 1).



(a) Example of the guesswhat game originated from GuessWhat.



(b) Example of the referential region matching game originated from VG.



(c) Example of the referred text reading task originated from HierText.

Figure 21. Data samples of the RD task (part 2).