# COCONut-PanCap: Joint Panoptic Segmentation and Grounded Captions for Fine-Grained Understanding and Generation
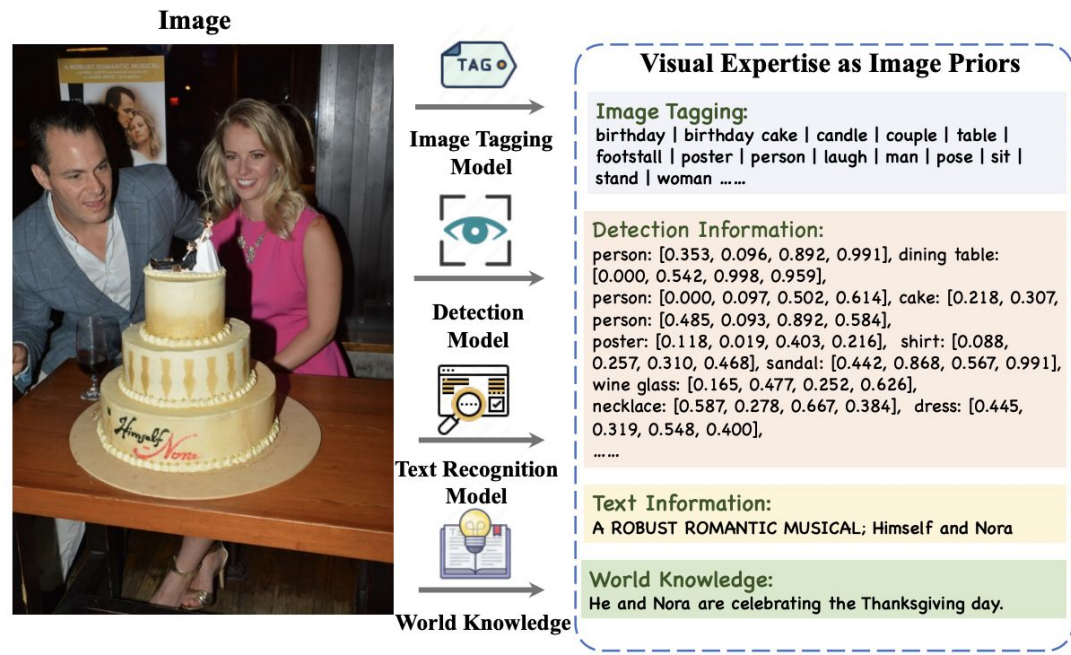
*Xueqing Deng, Linjie Yang, Qihang Yu, Ali Athar, Chenglin Yang, Xiaojie Jin, Xiaohui Shen, Liang-Chieh Chen*

ByteDance | Seed

# Introduction- data curation for VLMs

problems:
1. web crawled data is noisy
2. synthetic data from commercial VLM exists hallucination



**Visual Expertise as Image Priors**

**Image Tagging:**
birthday | birthday cake | candle | couple | table | footstall | poster | person | laugh | man | pose | sit | stand | woman ......

**Detection Information:**
person: [0.353, 0.096, 0.892, 0.991], dining table: [0.000, 0.542, 0.998, 0.959], person: [0.000, 0.097, 0.502, 0.614], cake: [0.218, 0.307, person: [0.485, 0.093, 0.892, 0.584], poster: [0.118, 0.019, 0.403, 0.216], shirt: [0.088, 0.257, 0.310, 0.468], sandal: [0.442, 0.868, 0.567, 0.991], wine glass: [0.165, 0.477, 0.252, 0.626], necklace: [0.587, 0.278, 0.667, 0.384], dress: [0.445, 0.319, 0.548, 0.400],
......

**Text Information:**
A ROBUST ROMANTIC MUSICAL; Himself and Nora

**World Knowledge:**
He and Nora are celebrating the Thanksgiving day.
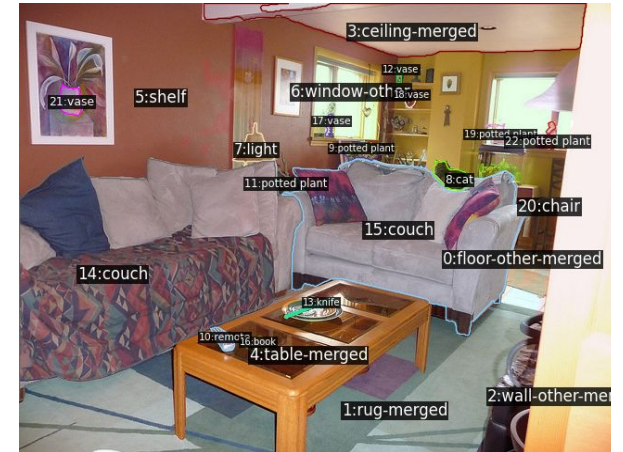
figure credit: DenseFusion-1M

We have a lot of existing vision expert models to help collect more details!

# Our approach - **leveraging panoptic segmentation models to collect comprehensive description**

Contributions:

1) Leverage SoTA panoptic segmentation model to tag image and limited human effort to deliver hallucination free dataset

2) Perception information and image-text data enable novel pretraining on vision and language alignment.

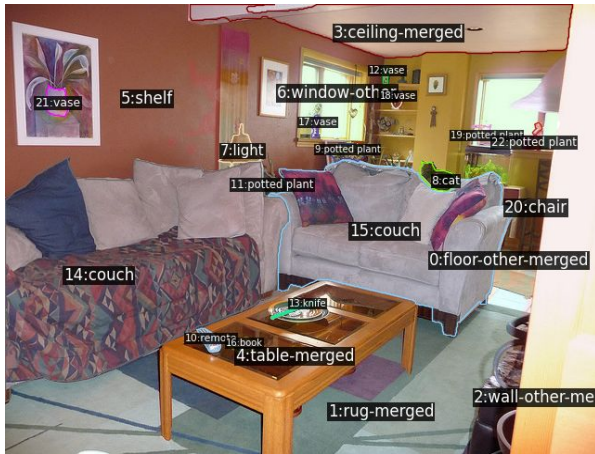3) High-quality instruction dataset enabling novel task and improve existing task



Panoptic segmentation= instance + semantic segmentation
1) instance segmentation: segmenting objects with counting
2) semantic segmentation: segment background stuffs without counting

# Introduction

**COCO:** Cozy living room with couches, coffee table, warm decor, natural light, and potted plants.
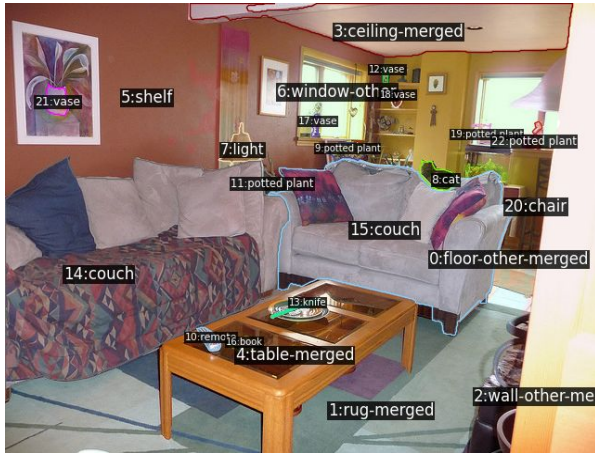
# Introduction



**COCO:** Cozy living room with couches, coffee table, warm decor, natural light, and potted plants.

**GPT-4V:** The image showcases a warm and inviting living room with cozy furniture and thoughtful decor. Two couches, one with a patterned throw and the other with vibrant pillows, are positioned around a wooden coffee table holding decorative items. A black cat is sleeping on the right couch, adding life to the scene. A neutral-toned geometric rug ties the furniture together, while warm red and orange walls adorned with framed artwork enhance the cozy atmosphere. Large windows let in natural light, with potted plants on the windowsills adding a refreshing touch. Shelves in the background display books and decor, completing the harmonious space.

**hallucination**

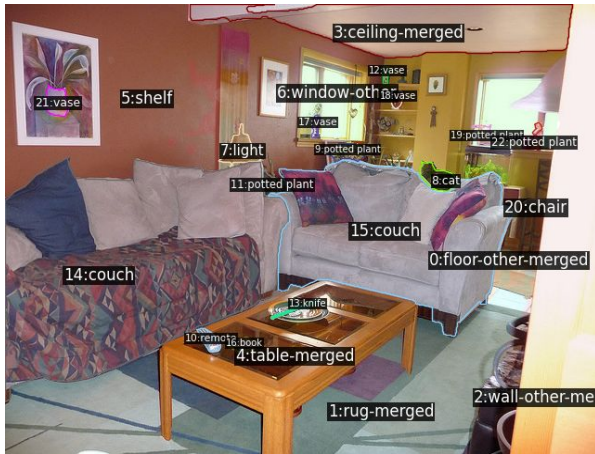**missing details for salient object**

# Introduction



**COCO:** Cozy living room with couches, coffee table, warm decor, natural light, and potted plants.

**GPT-4V:** The image showcases a warm and inviting living room with cozy furniture and thoughtful decor. Two couches, one with a patterned throw and the other with vibrant pillows, are positioned around a wooden coffee table holding decorative items. A black cat is sleeping on the right couch, adding life to the scene. A neutral-toned geometric rug ties the furniture together, while warm red and orange walls adorned with framed artwork enhance the cozy atmosphere. Large windows let in natural light, with potted plants on the windowsills adding a refreshing touch. Shelves in the background display books and decor, completing the harmonious space.

**COCONut-PanCap:** The image features a cozy and well-decorated living room. At the center of the room, <4:a wooden coffee table equipped with glasses> holds various items, including <10:a remote control>, <13: a knife> on the plates, and <16:a square small book>. On the left, The seating arrangement includes <14:a patterned couch with colorful cushions and blanket> and <15:another neutral-toned couch with vibrant throw pillows> , providing balance to the layout.  The rug with colorful patters brings more warm atmosphere to the sitting area. Behind the couch, <20:A chair in the back> complements the seating options. Adding warmth to the room, <8:a black cat> rests comfortably on the couch. Behind the sitting area, there is <5:a 4-layer wall-mounted wooden shelf> with additional decorative items, including <11,12: vases> and other decorative items, enhancing the cozy and inviting atmosphere. Closed to the shelf, there are several <9,19,22:potted plants with green leaves > are placed throughout the room, adding a touch of greenery.  <2:The wall painted in warm tones>, create a cozy atmosphere and are adorned with framed artwork and decorations.  <0: The floor is neutral-toned>, supporting the entire setup. The <3:ceiling painted white>, contrasts subtly with the walls and reflects the natural light entering the room through <6:the large windows>.

# Introduction

| Dataset Name | Samples | Avg. Words | Caption | T2I | Grd. Seg. |
|---|---|---|---|---|---|
| COCO-30K [5] | 30,000 | 11 | ✓ | ✓ | ✗ |
| DOCCI-test [37] | 5,000 | 136 | ✓ | ✓ | ✗ |
| IIW-test [14] | 445 | 217 | ✓ | ✓ | ✗ |

Table 2. **Comparison of Evaluation Datasets.** Our COCONut-PanCap validation set provides detailed captions and supports multiple multi-modal tasks, including image captioning, text-to-image generation (T2I), and grounded segmentation (Grd. Seg.).

# Introduction

| Dataset Name | Samples | Avg. Words | Caption | T2I | Grd. Seg. |
|---|---|---|---|---|---|
| COCO-30K [5] | 30,000 | 11 | ✓ | ✓ | ✗ |
| DOCCI-test [37] | 5,000 | 136 | ✓ | ✓ | ✗ |
| IIW-test [14] | 445 | 217 | ✓ | ✓ | ✗ |
| GenEval [15] | 553 | 8 | ✗ | ✓ | ✗ |
| T2I-CompBench val [19] | 2400 | 9 | ✗ | ✓ | ✗ |

Table 2. **Comparison of Evaluation Datasets.** Our COCONut-PanCap validation set provides detailed captions and supports multiple multi-modal tasks, including image captioning, text-to-image generation (T2I), and grounded segmentation (Grd. Seg.).

# Introduction

| Dataset Name | Samples | Avg. Words | Caption | T2I | Grd. Seg. |
|---|---|---|---|---|---|
| COCO-30K [5] | 30,000 | 11 | ✓ | ✓ | ✗ |
| DOCCI-test [37] | 5,000 | 136 | ✓ | ✓ | ✗ |
| IIW-test [14] | 445 | 217 | ✓ | ✓ | ✗ |
| GenEval [15] | 553 | 8 | ✗ | ✓ | ✗ |
| T2I-CompBench val [19] | 2400 | 9 | ✗ | ✓ | ✗ |
| GLaMM-GCG val-test [42] | 2,000 | 128 | ✓ | ✗ | ✓ |
| COCONut-PanCap val (ours) | 25,000 | 233 | ✓ | ✓ | ✓ |

Table 2. **Comparison of Evaluation Datasets.** Our COCONut-PanCap validation set provides detailed captions and supports multiple multi-modal tasks, including image captioning, text-to-image generation (T2I), and grounded segmentation (Grd. Seg.).

# Dataset construction – annotation pipeline



input image

# Dataset construction – annotation pipeline



input image

**human** annotated
panoptic segmentation mask

# Dataset construction – annotation pipeline
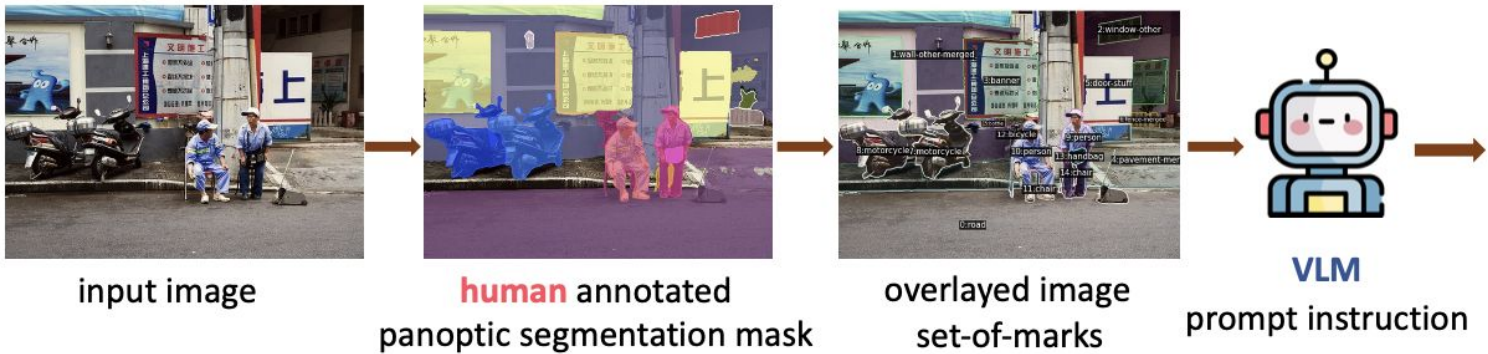


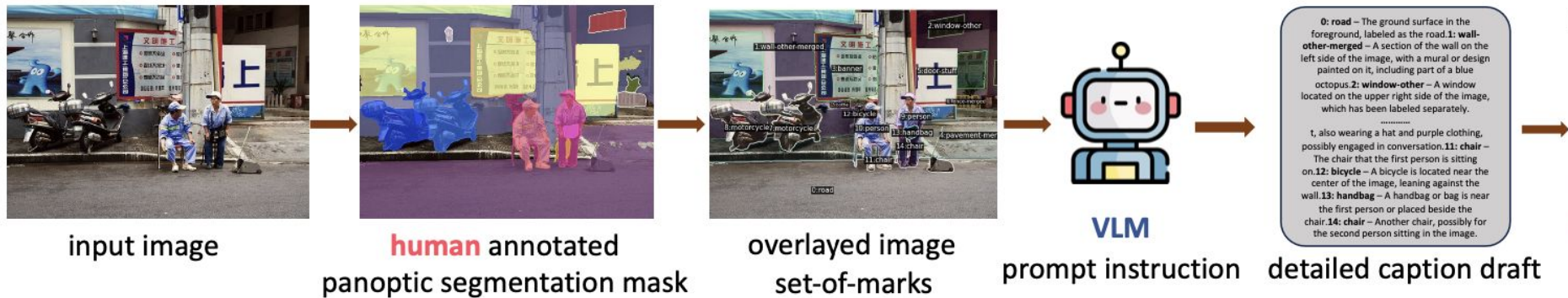input image

**human** annotated
panoptic segmentation mask

overlayed image
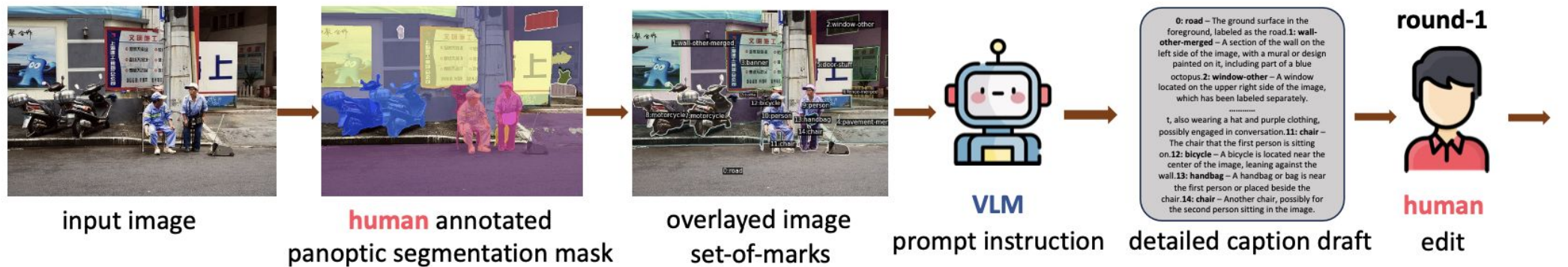set-of-marks

# Dataset construction – annotation pipeline



input image

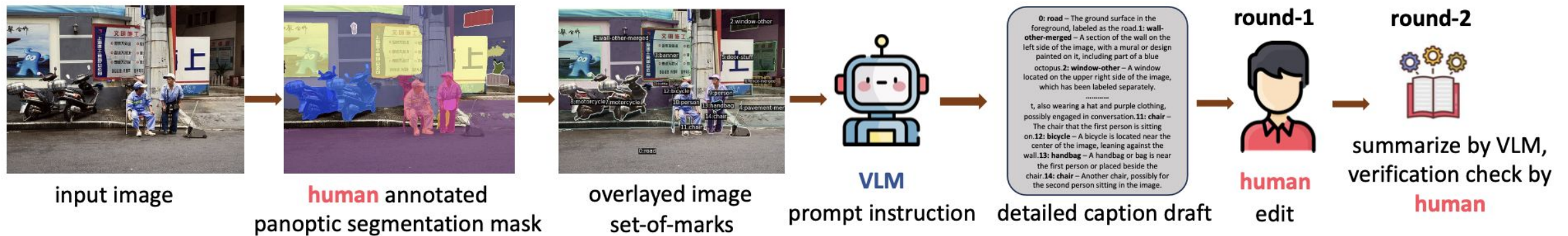**human** annotated
panoptic segmentation mask

overlayed image
set-of-marks

VLM
prompt instruction

# Dataset construction – annotation pipeline



input image

**human** annotated
panoptic segmentation mask

overlayed image
set-of-marks

VLM
prompt instruction

detailed caption draft

# Dataset construction – annotation pipeline



input image

**human** annotated
panoptic segmentation mask

overlayed image
set-of-marks

VLM
prompt instruction

detailed caption draft

round-1

**human**
edit

# Dataset construction – annotation pipeline



input image

**human** annotated
panoptic segmentation mask

overlayed image
set-of-marks

VLM
prompt instruction

detailed caption draft

round-1
**human**
edit

round-2
summarize by VLM,
verification check by
**human**

# Dataset construction – prompt template

# Dataset construction – prompt template



**round-1**

👤 Could you please describe the details of the image for each region?

🤖 From your image, here is the analysis:

0: Sky - Clear blue sky with some scattered clouds.

1: Road - Asphalt road, smooth and unobstructed, ideal for vehicles.

2: Building-other-merged - Building structure with ~~large~~ glass windows and beige walls, ~~likely a~~ commercial ~~or public building.~~
  _(x above "commercial")_

3: Tree-merged - Lush green trees with dense foliage, ~~possibly ornamental in a landscaped area.~~

4: Pavement-merged - Grey concrete pavement, typically used for pedestrian pathways.

**is upside down** ✓

5: Stop sign - Red octagonal traffic stop sign with white text, indicating a vehicle stop required.

**is approaching**

6: Car - A car ~~in motion, likely in the process of slowing down as it approaches~~ the stop sign.

7: Truck - A white truck, stationary, parked in a designated area, ~~possibly for~~ loading or unloading goods.
  _(x above "loading")_

---

**human edit legend:**

∨ : add     —— : remove     Ⓧ found hallucination and remove

# Dataset construction – prompt template



**round-1**

👤 Could you please describe the details of the image for each region?

🔵 From your image, here is the analysis:

0: Sky - Clear blue sky with some scattered clouds.

1: Road - Asphalt road, smooth and unobstructed, ideal for vehicles.

2: Building-other-merged - Building structure with ~~large~~ glass windows and beige walls, ~~likely a~~ (commercial) ~~or public building.~~ ✗

3: Tree-merged - Lush green trees with dense foliage, ~~possibly ornamental in a landscaped area.~~

4: Pavement-merged - Grey concrete pavement, typically used for pedestrian pathways.

**is upside down** ✓

5: Stop sign - Red octagonal traffic stop sign with white text, indicating a vehicle stop required.

**is approaching**

6: Car - A car ~~in motion, likely in the process of slowing down as it approaches~~ the stop sign.

7: Truck - A white truck, stationary, parked in a designated area, ~~possibly for~~ (loading or unloading goods.) ✗

**human edit legend:**

∨ : add      — : remove      ⊗ found hallucination and remove

**round-2**

👤 With the given region-level description, can you merge them into one paragraph in the format of grounded caption?

🔵 In this scene, the <u>clear blue sky</u> with some scattered clouds serves as a pleasant backdrop. The <u>smooth asphalt road</u> provides an unobstructed pathway, suitable for vehicles. Nearby, there is <u>a building structure with glass windows</u> and beige walls. <u>Lush green trees</u> with dense foliage add a touch of nature. The <u>grey concrete pavement</u> offers a pedestrian-friendly path alongside the road. <u>A upside down red octagonal traffic stop sign</u> with white text signals vehicles to halt. Approaching the stop sign is <u>a car</u>, likely slowing down. In the background, <u>a stationary white truck</u> is parked in a designated area.
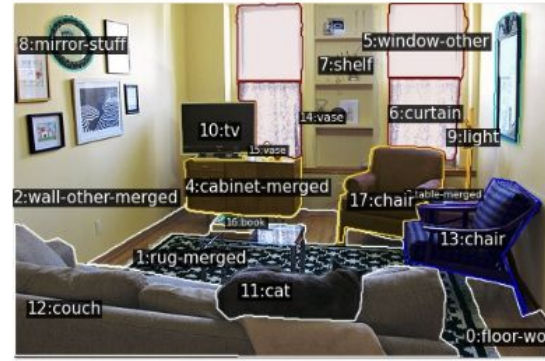
# # masks/image



**#mask <15**

#instance masks = 5    #semantic masks = 2

The image shows a meal plate with various labeled regions. In the center and right side of the plate, the region represents <0:a serving of French fries>. The upper left side shows <1: a slice of onion>. Positioned on the lower left side, it is <2: a small blue bowl> containing a red sauce, likely ketchup, while <3:another bowl> near the top center hold additional dipping sauce. At the bottom right of the image, <4: a silver knife> is partially visible under some of the food items. The main item on the plate, <5:a sandwich with a sesame seed bun> occupies the center and there is <6: another sandwich> on the right. Finally, all these items are put on the <7: dining table surface>. Together, the image captures the different elements of the meal setup, showing a typical plate with a sandwich, fries, and dipping sauce.

**15 < #mask <25**

#instance masks = 13    #semantic masks = 5

The image depicts a cozy living room with various items of furniture and decor. The <0:wooden floor> represents the wooden flooring that spans the room, partially covered by <1: a black and white patterned rug> in the center. The <2:light yellow walls> surrounding the other elements in the room. <3:A coffee table> sits on the rug where there is <16: a book> on top of it, in front of the seating area. Against the back wall, there is <4: wooden cabinet>, which holds the <10: tv with screen off> , <15: a vase> and some decorative items. The <5: window> is located on the back wall, partially covered by <6: white window blind>, which lets light into the room. Above the cabinet, <7:a shelf> holds various decor pieces, including <14:a vase>. <8: A mirror> is mounted on the left wall along with multiple frames, while <9:a light fixture> is visible in the upper right corner. In the foreground, <12: a light-upholstered couch> where <11:a black cat> lies, adds a cozy touch. On the right side, there is <13:a blue-upholstered chair> near the coffee table, and <17:a reddish chair> positioned near the back wall.

**#mask >25**

#instance masks = 54    #semantic masks = 9

The image shows a lively family gathering around <34:along dining table>, with various labeled individuals, dining items, and furniture. In terms of people, the table is surrounded by multiple family members of varying ages, all seated and enjoying the meal together. Key individuals include <50:a lady with blonde short curly hair wearing tank>, <52: a lady with black medium hair in blue shirt>, <53: a women with black medium hair but wearing in tank top>, and <55: a boy in red t shirt > sitting right next to the lady in the left side of the image, as well as on the right there are <32: a women in white t shirt with black dots>, <41: a boy in white t shirt>. < 62,63,48, 47, 13, 56, 53, 54, 55, 52, 50, 33, 42, 43, 41, 32, 45, 44, 47, 46:Additional people> are seated around or standing near the table, filling the room with a sense of warmth and gathering. The dining items are plentiful and include numerous <14, 26, 25, 40: wine glass>, <15, 17, 36, 16, 20, 39, 38: cup>, and <12, 60:other bowl> and <30, 10:spoon> objects scattered across the table. Utensils, such as <5: fork> and <9: knife>, are visible among the plates of food, which include main dishes and side dishes. In the background, the <0: light colored wall> and <1: wooden cabinet> display family photos, framed art, and <4: pink flowers> in <21:a glass vase>, adding a personal touch to the room decor. The combination of food, drink, and a lively gathering suggests a celebratory or holiday meal.

# # words/image

# Introduction

| Dataset Name | Image Source | Samples | Captioned by | Avg. Words | Avg. Masks |
|---|---|---|---|---|---|
| COCO-caption [5] | COCO [32] | 118K | **Human** | 11 | ✗ |
| DCI [50] | SA-1B [22] | 8K | **Human** | 144 | ✗ |
| DOCCI [37] | DOCCI [37] | 9.6K | **Human** | 136 | ✗ |
| IIW [14] | WebLI [14] | 8.5K | **Human** | 217 | ✗ |
| BLIP-LCS | LAION [45], CC [3], SBU [38] | 558K | BLIP [27] | 54 | ✗ |
| DenseFusion1M [29] | LAION [45] | 1,059K | Vision Specialist Models | 191 | ✗ |
| LLaVA-Recap118K [35] | COCO [32] | 118K | LLaVA-NEXT [35] | 186 | ✗ |
| LLaVA-Details-23K [34] | COCO [32] | 23K | GPT4 | 105 | ✗ |
| ShareGPT4V [4] | LAION [45], CC [3], SBU [38], COCO [32] etc. | 100K | GPT4-Vision | 162 | ✗ |
| ShareGPT4V-PT [4] | LAION [45], CC [3], SBU [38], COCO [32] etc. | 1,246K | Share-Captioner [4] | 144 | ✗ |
| PixelLM-MUSE [43] | LVIS [16] | 246K | GPT4-Vision | - | 3.7‡ |
| Osprey [58] | COCO [32] | 724K | GPT4-Vision | - | - |
| GLaMM-GCG [42] | Flickr-30K [40], PSG [54], RefCOCOg [20] | 214K | Vision Specialist Models | 128 | 3.6 |
| COCONut-PanCap (ours) | COCO [32] | 118K | GPT-assisted **Human** | 203 | 13.2 |

Table 1. **Dataset Comparison**. Our proposed COCONut-PanCap dataset stands out for its **detailed** (2nd highest in Average Words), **high-quality** (human interactive annotated) *captions* and **high-density** *segmentation masks* (1st in Average Masks). ‡ denotes the mask number for referring segmentation which only counts the targets in QA format. Note that "Samples" means the number of collected annotations, where there may exist one image with multiple different annotation, *i.e.*, in region-level datasets like Osprey.
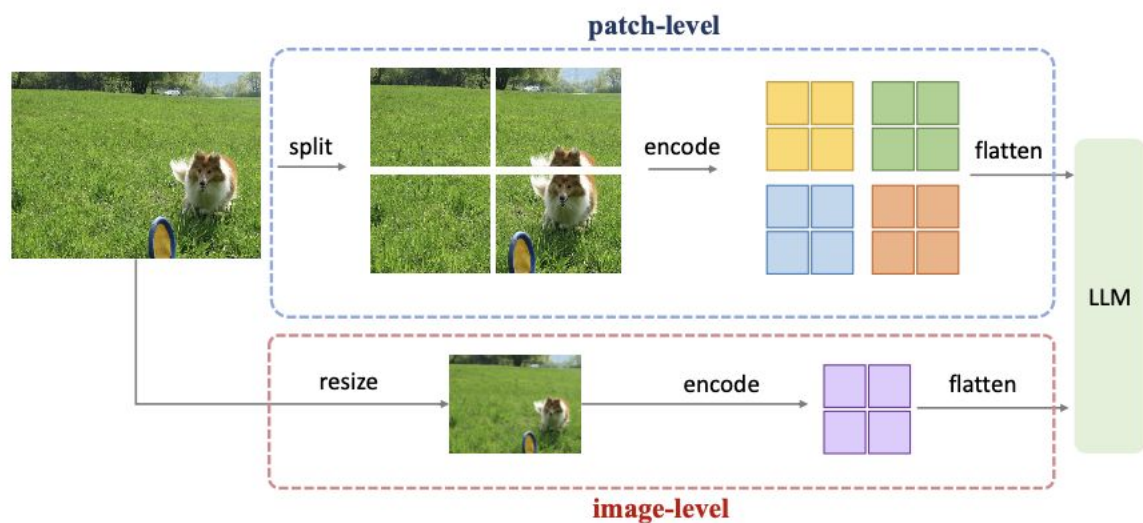
# Question

1. How is the data useful in the VLM training? pre-training? instruction tuning?
2. How can the data motivate the design of architecture modeling?

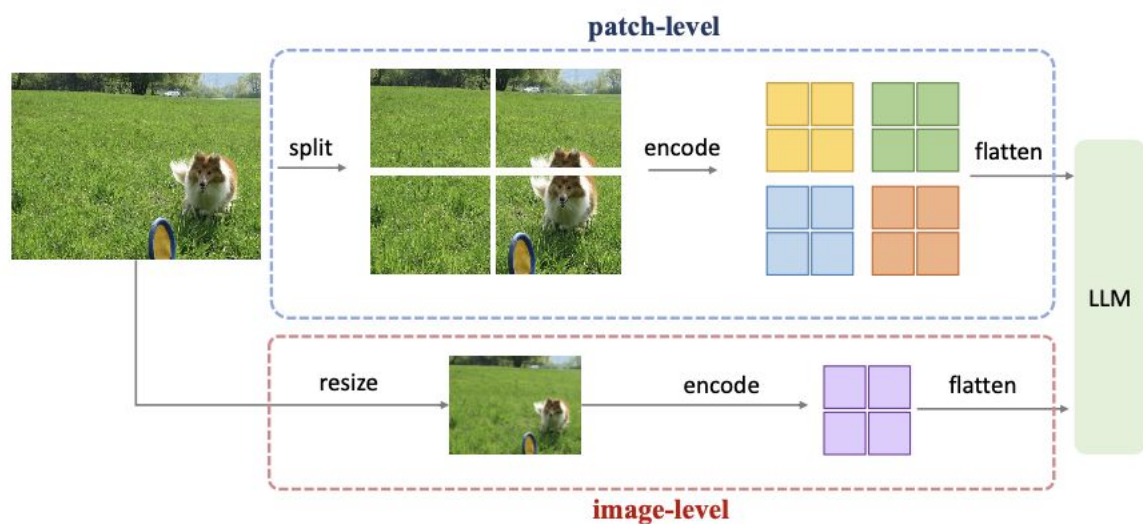# Mask-pooled LLaVA pretraining
## novel pretraining

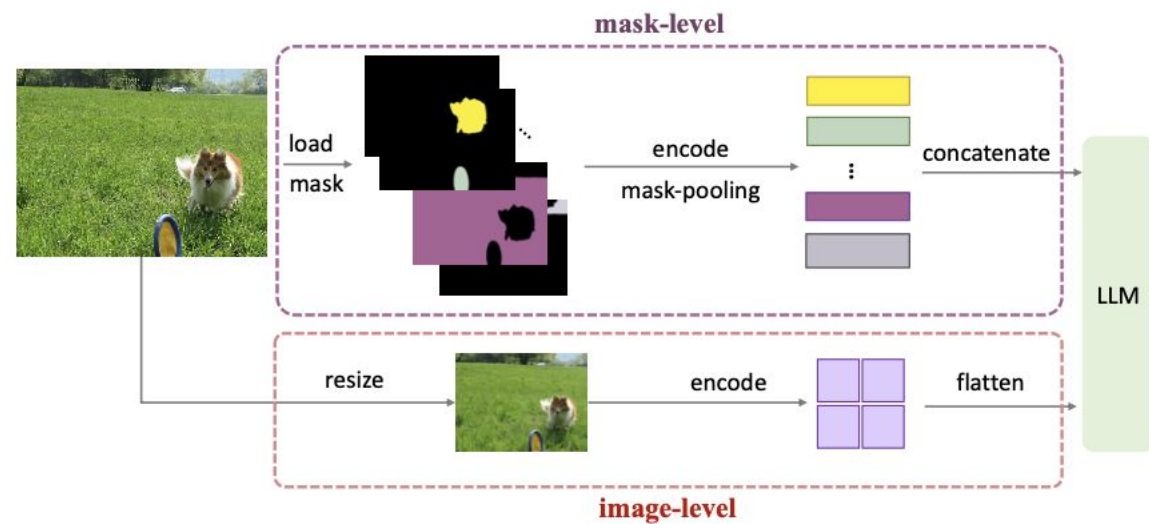

(a) LLaVA-NeXt-AnyRes

# Mask-pooled LLaVA pretraining
## novel pretraining



(a) LLaVA-NeXt-AnyRes

(b) our LLaVA-NeXt-pool

# High-quality instruction data boost performance

same data amount

| Training recipe | Method | Pretrain Dataset | Instruction-tuning dataset | Mask pooled | CAPTURE | CIDEr | BLEU@4 | METEOR | ROUGE-L |
|---|---|---|---|---|---|---|---|---|---|
| finetune | LLaVA-NeXT* | LAION-CC-SBU | LLaVA 665K | ✗ | 55.4 | 10.8 | 4.2 | 13.2 | 23.1 |
| | LLaVA-NeXT | LAION-CC-SBU | LLaVA 665K-COCONut-PanCap | ✗ | 58.7 | 11.2 | 4.8 | 16.2 | 24.6 |
| | LLaVA-NeXT-pool | LAION-CC-SBU | LLaVA 665K-COCONut-PanCap | ✓ | 61.4 | 13.1 | 5.3 | 17.1 | 26.8 |
| | LLaVA-NeXT-I | LAION-CC-SBU | LLaVA 665k-InternVL2-Cap | ✗ | 53.9 | 9.4 | 4.4 | 11.5 | 21.4 |
| | LLaVA-NeXT-Q | LAION-CC-SBU | LLaVA 665k-Qwen2VL-Cap | ✗ | 55.4 | 8.9 | 4.6 | 12.9 | 22.5 |
| | LLaVA-NeXT-G | LAION-CC-SBU | LLaVA 665k-GPT4V-Cap | ✗ | 56.2 | 9.6 | 4.7 | 13.3 | 22.8 |

human refined

synthetic data

Table 3. **Caption Benchmark Results Evaluated on Our COCONut-PanCap Val Set.** Note that the amount of data in the instruction dataset remains the same; only the sources of the detailed captions vary, with a total of 23K images that have detailed captions.

# Various VLM understanding tasks facilitate by instruction data

## Detailed Captioning

**USER:** Could you please describe the image in detail?

**ASSISTANT:** Three teddy bears, each a different shade of brown, are positioned together on a blue bedsheet. On the left, a golden brown teddy ....
slightly with the others, also embracing its adjacent bears. The soft blue bedsheet in the background provides a cozy base for this heartwarming scene.

# Various VLM understanding tasks facilitate by data

## Detailed Captioning

**USER:** <u>Could you please describe the image in detail?</u>

**ASSISTANT:** Three teddy bears, each a different shade of brown, are positioned together on a blue bedsheet. On the left, a golden brown teddy ....
 slightly with the others, also embracing its adjacent bears. The soft blue bedsheet in the background provides a cozy base for this heartwarming scene.

## Panoptic Segmentation Grounded Caption

**USER:** <u>Please provide a detailed description of the image and segment each part.</u>

**ASSISTANT:**



n the image, there are two individuals<seg1> <seg2> sitting on chairs<seg3> on a sidewalk<seg4>. They are wearing blue niforms and appear to be street cleaners, as they have cleaning tools with them. Next to them are two parked cooters<seg5><seg6>, one of which has a helmet hanging from its handlebar. The background<seg7> features various signs nd posters<seg8>, including a large blue and white sign with Chinese characters and a cartoon character and light<seg9> anging on the wall. On the other side, there is another building with windows<seg10> and barrier in front of the uilding<seg11> and the plants<seg12>. The scene is set against a backdrop of buildings and a street, indicating a busy urban nvironment.

# Various VLM understanding tasks facilitate by data

## Detailed Captioning

**USER:** Could you please describe the image in detail?

**ASSISTANT:** Three teddy bears, each a different shade of brown, are positioned together on a blue bedsheet. On the left, a golden brown teddy ....
 slightly with the others, also embracing its adjacent bears. The soft blue bedsheet in the background provides a cozy base for this heartwarming scene.



## Panoptic Segmentation Grounded Caption

**USER:** Please provide a detailed description of the image and segment each part.

**ASSISTANT:**



the image, there are two individuals<seg1> <seg2> sitting on chairs<seg3> on a sidewalk<seg4>. They are wearing blue niforms and appear to be street cleaners, as they have cleaning tools with them. Next to them are two parked cooters<seg5><seg6>, one of which has a helmet hanging from its handlebar. The background<seg7> features various signs nd posters<seg8>. including a large blue and white sign with Chinese characters and a cartoon character and light<seg9> anging on the wall. On the other side, there is another building with windows<seg10> and barrier in front of the uilding<seg11> and the plants<seg12>. The scene is set against a backdrop of buildings and a street, indicating a busy urban nvironment.

## Text-to-Image Generation

**USER:** Generate an image with the following description. This image shows a cozy, well-lit **kitchen** with white cabinetry, giving it a bright, clean look. The dark countertops contrast with the white **cabinets**. ...

**ASSISTANT:**

# Various VLM understanding tasks facilitate by data

## Detailed Captioning

**USER:** Could you please describe the image in detail?

**ASSISTANT:** Three teddy bears, each a different shade of brown, are positioned together on a blue bedsheet. On the left, a golden brown teddy ....
slightly with the others, also embracing its adjacent bears. The soft blue bedsheet in the background provides a cozy base for this heartwarming scene.

## Panoptic Segmentation Grounded Caption

**USER:** Please provide a detailed description of the image and segment each part.

**ASSISTANT:**



n the image, there are two individuals<seg1> <seg2> sitting on chairs<seg3> on a sidewalk<seg4>. They are wearing blue niforms and appear to be street cleaners, as they have cleaning tools with them. Next to them are two parked cooters<seg5><seg6>, one of which has a helmet hanging from its handlebar. The background<seg7> features various signs nd posters<seg8>, including a large blue and white sign with Chinese characters and a cartoon character and light<seg9> anging on the wall. On the other side, there is another building with windows<seg10> and barrier in front of the uilding<seg11> and the plants<seg12>. The scene is set against a backdrop of buildings and a street, indicating a busy urban nvironment.

## Text-to-Image Generation

**USER:** Generate an image with the following description. This image shows a cozy, well-lit **kitchen** with white cabinetry, giving it a bright, clean look. The dark countertops contrast with the white **cabinets**. ...

**ASSISTANT:**



## Downstream Task

**VQA:**

**USER:** How many kinds of vegetables ?

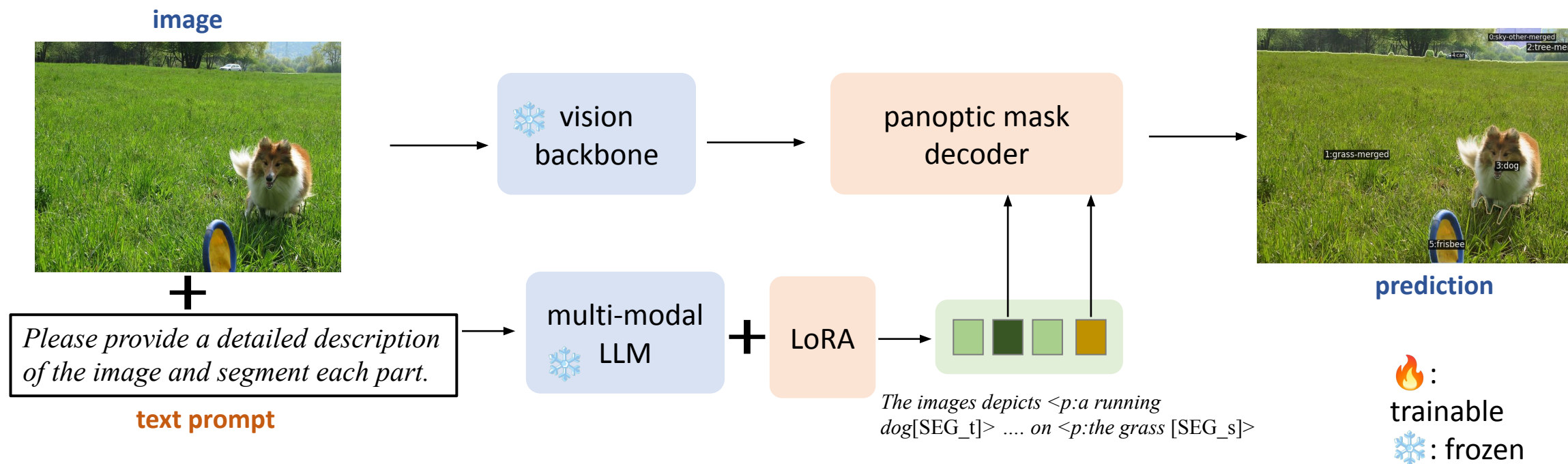**ASSISTANT:** There are broccoli, carrot ....

**Referring Segmentation:**

**USER:** Can you segment the cat which is on the toilet?

**ASSISTANT:** It is **<SEG>**.

# Panoptic segmentation Grounded Caption - PGC (novel task)

# Experimental results – **PGC (novel task)**

| Method | Pretrain dataset | Instruction dataset | Mask pooled | Caption | | | | Grounding segmentation | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | CAPTURE | CIDEr | BLEU@4 | METEOR | PQ | $PQ^{thing}$ | $PQ^{stuff}$ |
| LISA+ * | LAION-CC-SBU | GranDf | ✗ | 46.2 | 6.6 | 3.8 | 9.8 | 0.43 | 0.41 | 0.45 |
| LISA+ | LAION-CC-SBU | COCONut-PanCap (ours) | ✗ | 57.9 | 8.1 | 4.9 | 13.8 | 0.50 | 0.49 | 0.44 |
| GLaMM GCG * | LAION-CC-SBU+GranD | GranDf | ✗ | 43.2 | 6.5 | 3.6 | 10.6 | 0.27 | 0.35 | 0.21 |
| GLaMM GCG | LAION-CC-SBU+GranD | COCONut-PanCap (ours) | ✗ | 56.8 | 7.8 | 5.2 | 14.3 | 0.55 | 0.54 | 0.46 |
| PanCaper (ours) | LAION-CC-SBU | COCONut-PanCap (ours) | ✗ | 62.6 | 12.0 | 5.8 | 15.4 | 0.56 | 0.55 | 0.66 |
| PanCaper-Pro (ours) | LAION-CC-SBU | COCONut-PanCap (ours) | ✓ | 64.3 | 12.5 | 6.4 | 17.9 | 0.61 | 0.58 | 0.68 |

Table 4. **Joint Panoptic Segmentation and Grounded Captioning (PGC) on COCONut-PanCap Val Set.** * denotes reproduced results.

# Thank!

# Data gain: text-to-image generation

| | w/o FT | COCO-caption [5] | DOCCI [37] | COCONut-PanCap (ours) |
|---|---|---|---|---|
| color attribution | 0.37 | 0.34 | 0.38 | 0.40 |
| colors | 0.73 | 0.70 | 0.74 | 0.75 |
| position | 0.33 | 0.30 | 0.36 | 0.36 |
| counting | 0.65 | 0.64 | 0.65 | 0.70 |
| single object | 0.96 | 0.94 | 0.95 | 0.96 |
| two objects | 0.80 | 0.78 | 0.81 | 0.89 |
| overall score | 0.64 | 0.62 | 0.65 | 0.68 |

Table 6. **Effects of Fine-tuning the SD3-medium (T2I model) with Different Datasets on GenEval [15].** w/o FT denotes the model is not finetuned with any datasets (*i.e.*, zero-shot testing).

# Data gain: VQA

adding more human-annotated detailed caption

| Method | LLM | Instruction-tuning Dataset | MM-Vet | Seed-IMG | MMBench-en | TextVQA | POPE | MME |
|--------|-----|----------------------------|--------|----------|------------|---------|------|-----|
| LLaVA-NeXT * | Llama3-8B | orginal LLaVA 665K [38] | 43.5 | 70.1 | 71.4 | 68.9 | 85.4 | 1523 |
| LLaVA-NeXT-20K | Llama3-8B | LLaVA 665K-COCONut-PanCap-20K | 44.1 | 72.5 | 73.6 | 69.8 | 86.1 | 1552 |
| LLaVA-NeXT-50K | Llama3-8B | LLaVA 665K-COCONut-PanCap-50K | 44.6 | 73.1 | 74.2 | 70.0 | 87.1 | 1600 |
| LLaVA-NeXT-Full | Llama3-8B | LLaVA 665K-COCONut-PanCap-118K | 45.5 | 74.3 | 75.1 | 70.7 | 87.9 | 1612 |
| LLaVA-1.5 | Vicuna-7B | LLaVA 665K-ShareGPT4V-100K | 37.8 | 67.4 | 70.5 | 64.6 | 84.7 | 1519 |
| LLaVA-1.5 | Vicuna-7B | LLaVA 665K-COCONut-PanCap-20K | 38.5 | 67.7 | 70.9 | 64.5 | 84.9 | 1521 |

Table 7. **Benchmark Results and Ablation Study on VQA.** By adding extra detailed caption data for instruction tuning, the models show increased improvement. * denotes reproduced results. Using only **20K human labeled data** can still achieve **comparable performance** to 100K synthetic data.

# Better dense understanding for VLM: referring segmetation

| Method | refCOCO | | | refCOCO+ | | | refCOCOg | |
|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val | test |
| GLaMM* [42] | 77.5 | 79.2 | 74.9 | 71.3 | 74.7 | 61.5 | 71.3 | 71.9 |
| PixelLM [43] | 73.0 | 76.5 | 68.2 | 66.3 | 71.7 | 58.3 | 69.3 | 70.5 |
| LISA-7B [25] | 74.1 | 76.5 | 71.1 | 62.4 | 67.4 | 56.5 | 66.4 | 68.5 |
| PanCaper$^+$ | 74.5 | 76.7 | 69.9 | 69.9 | 73.4 | 59.5 | 69.8 | 70.6 |
| PanCaper$^+$ + COCONut-PanCap | 76.2 | 77.1 | 72.3 | 70.5 | 73.9 | 60.1 | 72.1 | 71.6 |

Table 8. **Benchmark Results on Referring Segmentation.** * denotes reproduced results. It is noted that GLaMM uses extra data from the GranD dataset for pretraining. $^+$ denotes our PanCaper model is adapted for referring segmentation task.

# Undergoing project

1. leverage perception data in video generation
2. building multi-modal CoT data for improving reasoning ability to image generation

# 1.Perception data to improve video generation

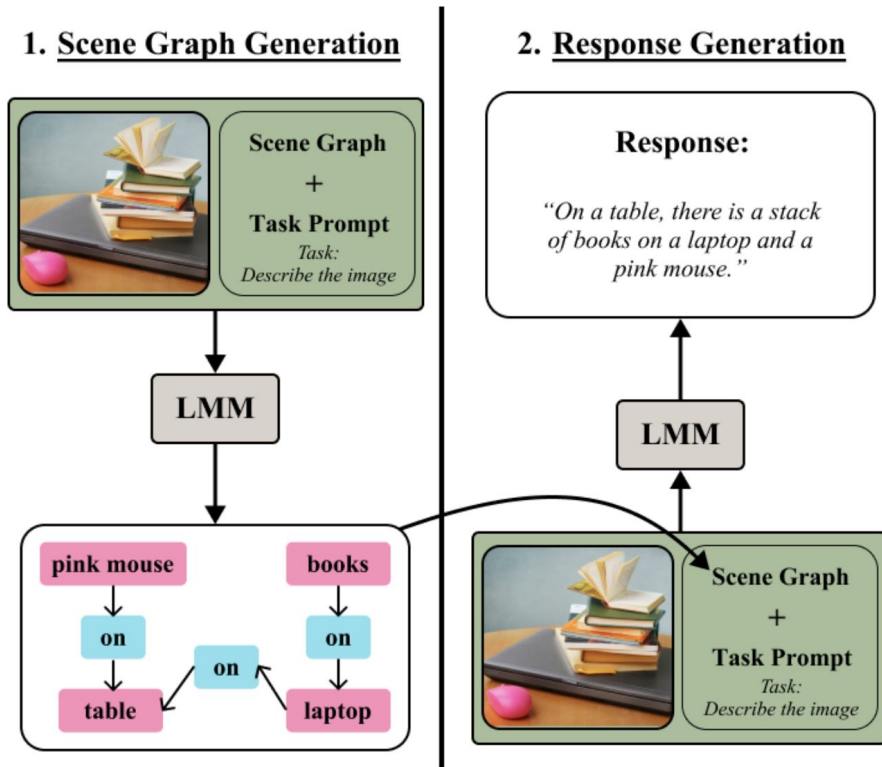| Method | Subject consisitency | Background consisitency | Dynamic degree | Motion smoothness |
|---|---|---|---|---|
| RGB finetuning | 0.958 | **0.959** | **0.39** | 99.47 |
| w/ depth prediction | **0.967** | 0.956 | **0.67** | 99.11 |
| w/ flow prediction | 0.965 | 0.954 | **0.71** | **99.56** |
| w/ depth and flow | 0.953 | 0.954 | **0.85** | 99.38 |

# 1. Perception data to improve image to video generation
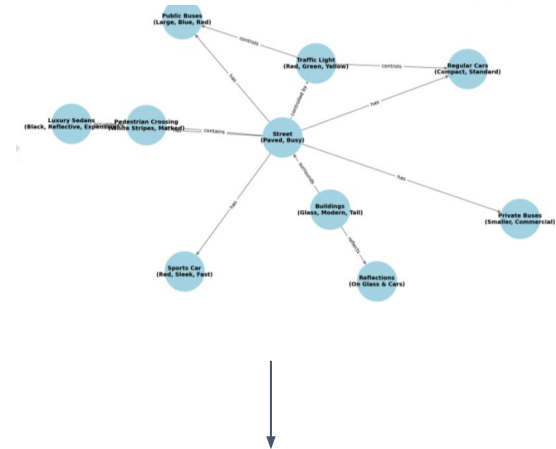
prompt: grab a pen and put it into a bowl.

# 2. Build CoT data to help improve image generation

Leverage panoptic segmentation to create CCoT data



figure credit: compositional CoT (CCoT)

image-to-text understanding

**extend simple prompt with CCoT**

extended prompt

text-to-image generation

# Panoptic Scene Graph Evaluation

Follow the OpenPSG[1] evaluation setting to report the recall rate, when calculate the recall

Evaluation set: 1000 selected long and complex scene text prompts (COCO images)

| Method | R@10 ↑ | R@15 ↑ | R@20 ↑ | FID ↓ |
|---|---|---|---|---|
| SD3: Caption only | 31.4 | 41.7 | 46.8 | 23.4 |
| SD3: Caption->SG->SG text | 35.6 | 47.8 | 51.9 | **19.5** |
| SDXL + Long-CLIP (no finetuning) | 21.7 | 32.5 | 37.8 | 35.7 |

OpenPSG, ECCV2024