

# UNIFIED-IO: A UNIFIED MODEL FOR VISION, LANGUAGE, AND MULTI-MODAL TASKS



Jiasen Lu



Christopher Clark



Rowan Zellers



Roozbeh Mottaghi



Ani Kembhavi

<https://unified-io.allenai.org/>

# Single-Task Model vs. Unified Model

Single-Task Model

A set of parameters that optimize for single tasks.

Multi-Task Model

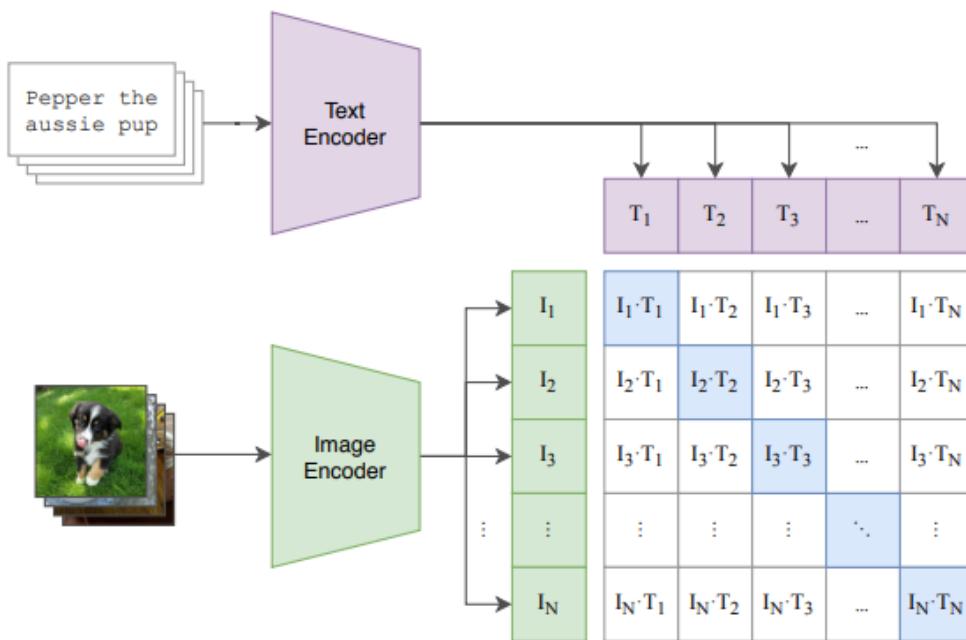
Shared parameters + task specific heads for multiple tasks.

Unified Model

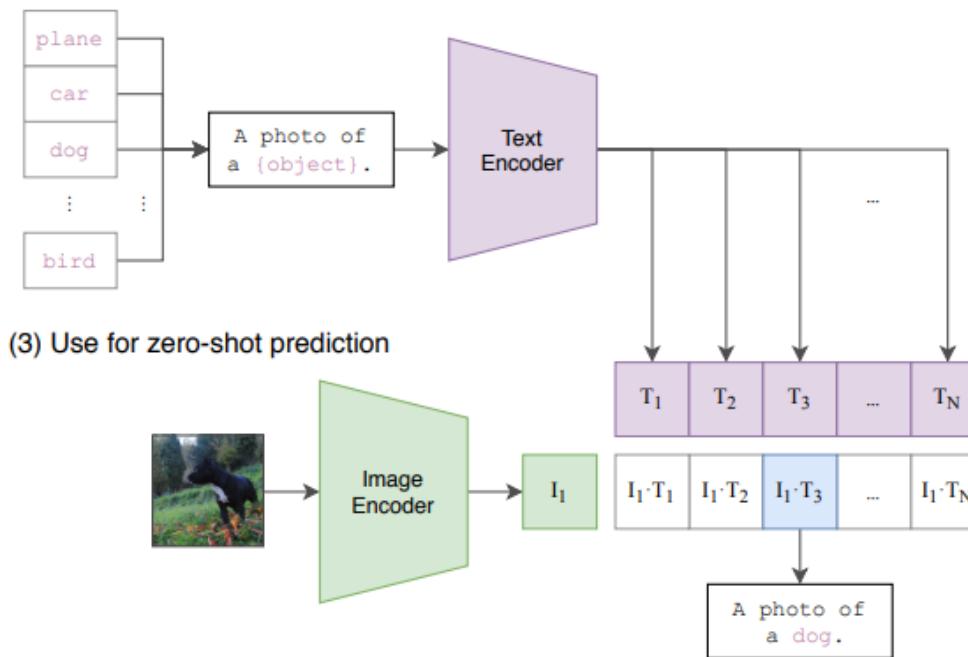
Shared parameters + shared heads for multiple tasks.

# Single-Task Model for Vision

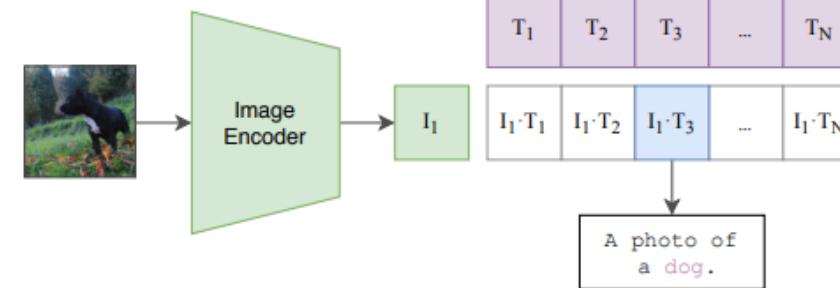
(1) Contrastive pre-training



(2) Create dataset classifier from label text



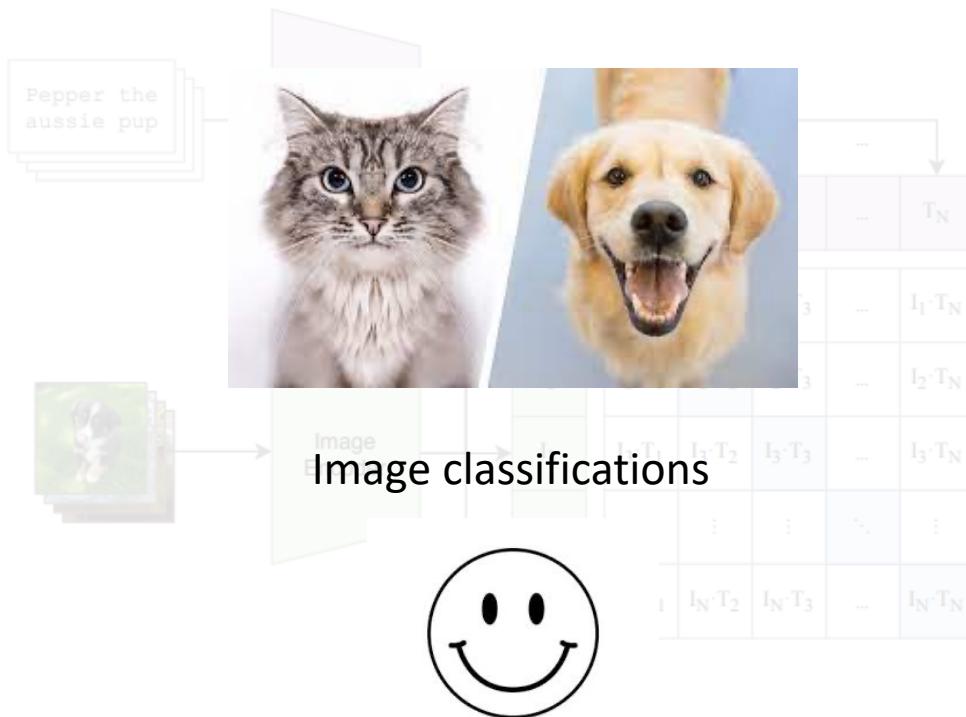
(3) Use for zero-shot prediction



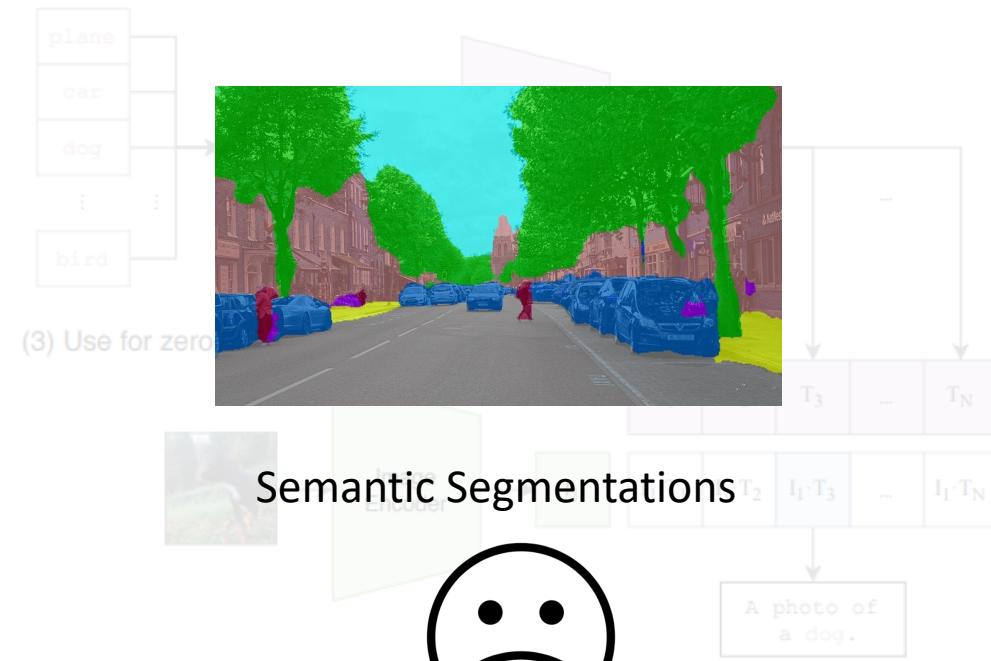
CLIP: Connecting Text and Images [Radford et.al. 2021]

# Single-Task Model for Vision

(1) Contrastive pre-training



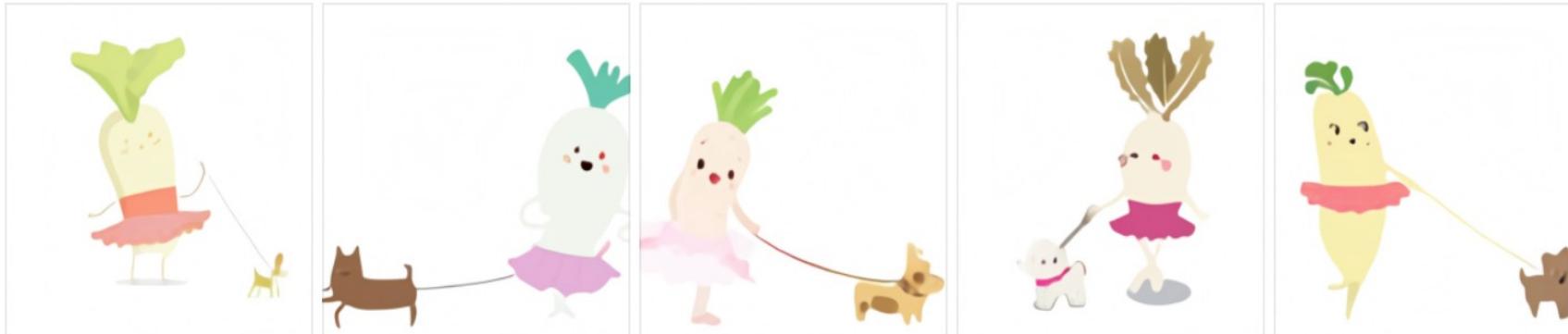
(2) Create dataset classifier from label text



CLIP: Connecting Text and Images [Radford et.al. 2021]

# Single-Task Model for Vision

an illustration of a baby daikon radish in a tutu walking a dog

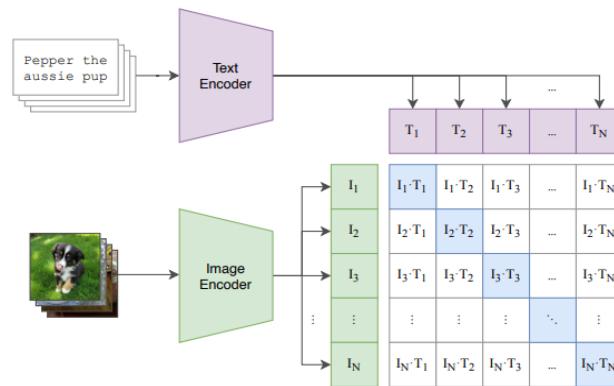


[Edit prompt or view more images↓](#)

Zero-Shot Text-to-Image Generation (DALLE) [Ramesh et.al. 2021]

# Single-Task Model for Vision

(1) Contrastive pre-training



CLIP: Connecting Text and Images [Radford et.al. 2021]



Semantic Segmentations

an illustration of a baby daikon radish in a tutu walking a dog



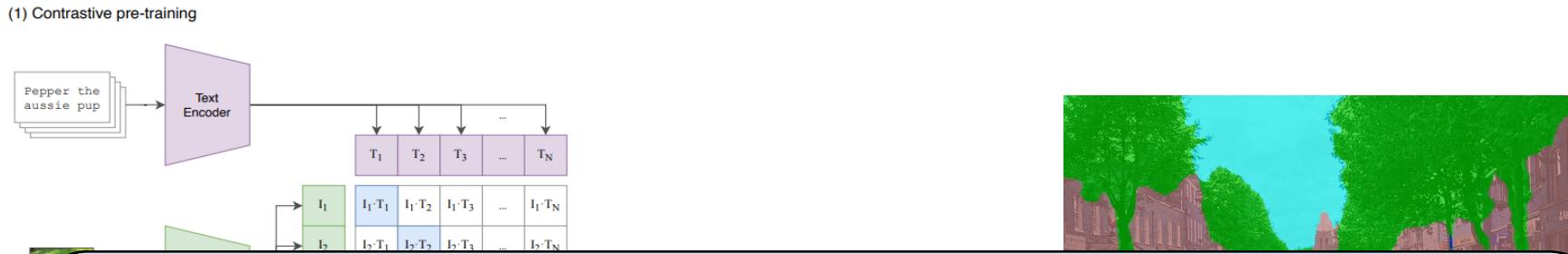
Edit prompt or view more images↓

Zero-Shot Text-to-Image Generation (DALLE) [Ramesh et.al. 2021]



Image classifications

# Single-Task Model for Vision



CLIP:

If the model can learn how to generate images, they should be able to solve any semantic vision tasks.

an illustration of a baby daikon radish in a tutu walking a dog



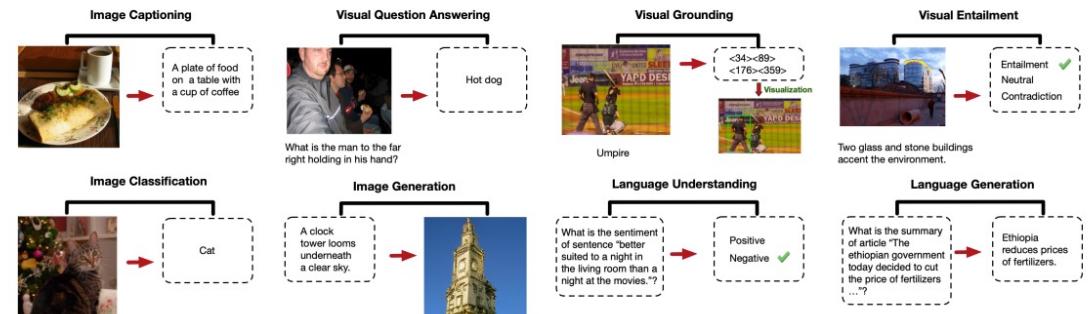
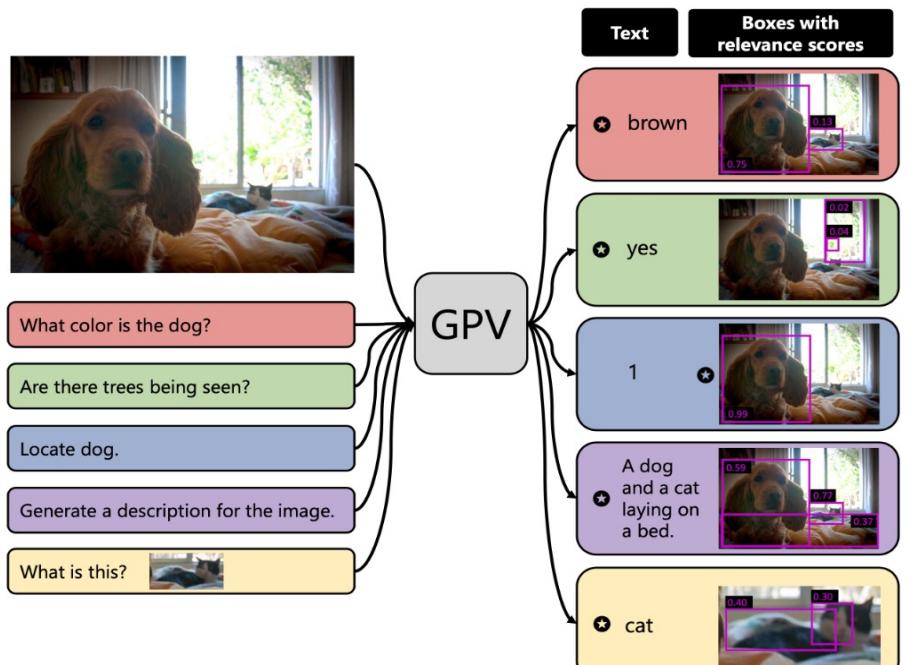
Edit prompt or view more images↓

Zero-Shot Text-to-Image Generation (DALLE) [Ramesh et.al. 2021]

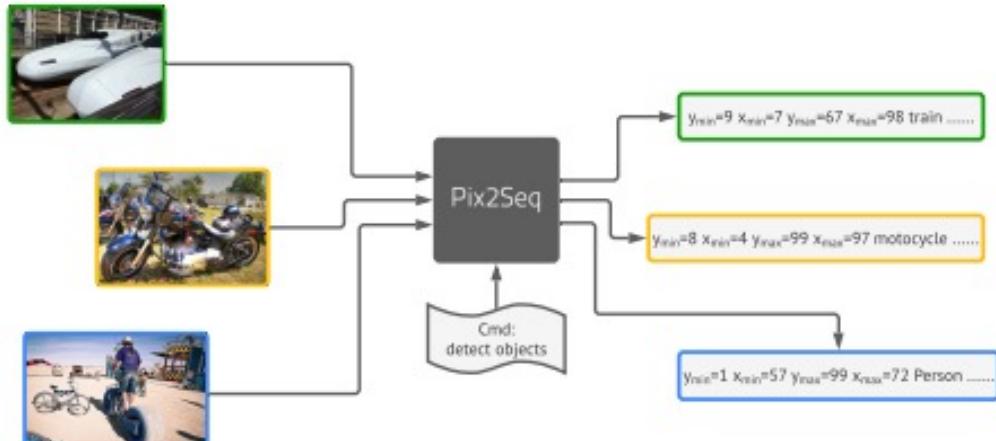


Image classifications

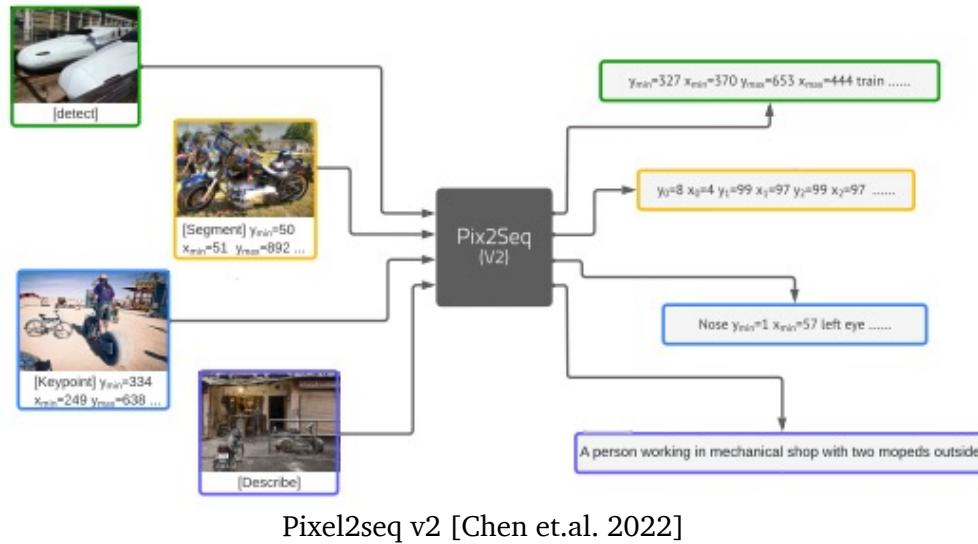
# Prior Work



# Prior Work



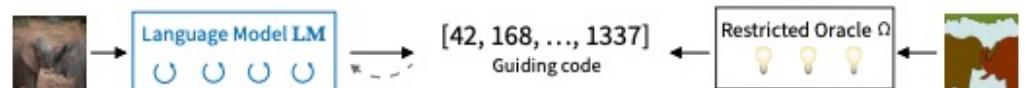
Pixel2seq [Chen et.al. 2021]



Pixel2seq v2 [Chen et.al. 2022]



(a) **Stage I** training: we train the base model  $f$ , which is guided by the code produced by the *restricted oracle* model  $\Omega$ . The oracle has access to the ground-truth label, but is only allowed to communicate with  $f$  by passing a short discrete sequence, which we call a *guiding code*.



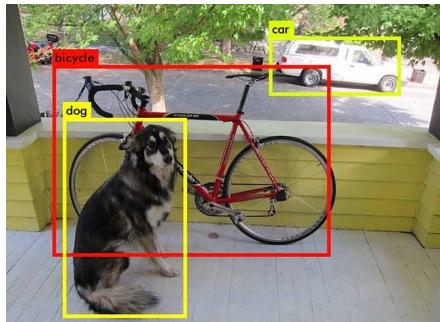
(b) **Stage II** training: we train a *language model* (LM) to output a *guiding code* by learning to mimic the oracle, but using only the image input.

UVIM [Kolesnikov et.al. 2022]

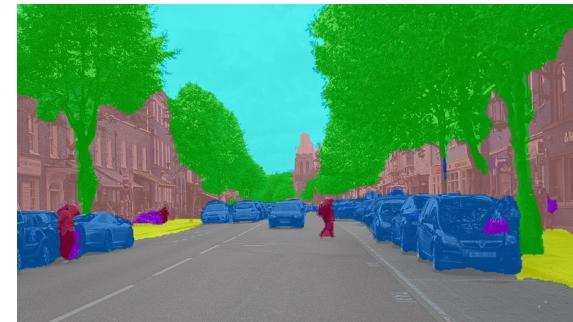
# Vision has diverse output format



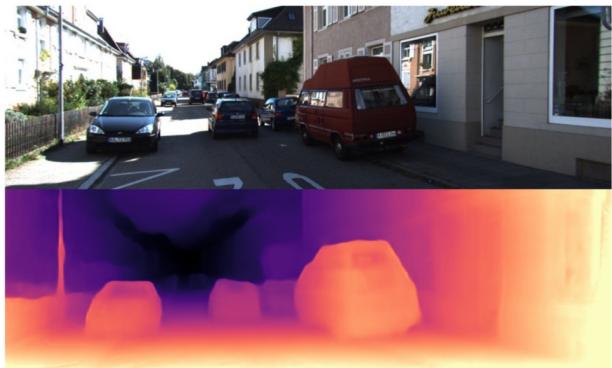
Image classifications



Object detections



Semantic Segmentations



Depth Estimation



Pose Estimation

& MORE

# Vision has diverse output format

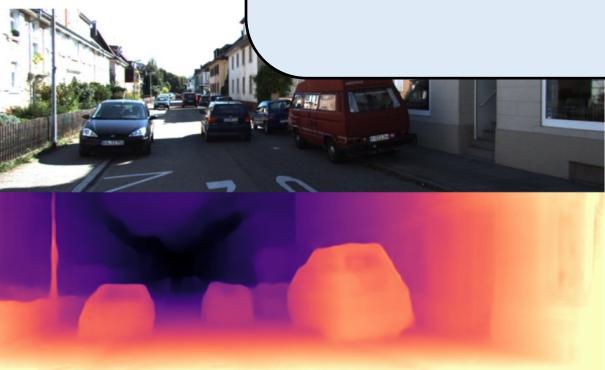
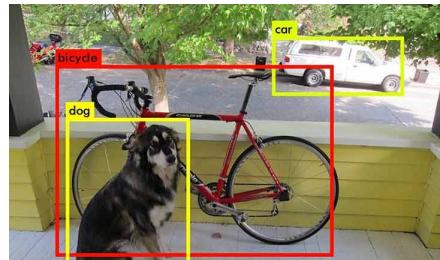
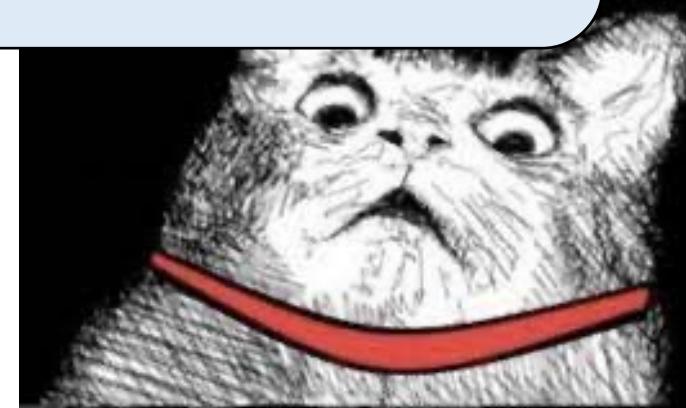


Image classification

Convert diverse vision inputs/outputs into **sequences**.

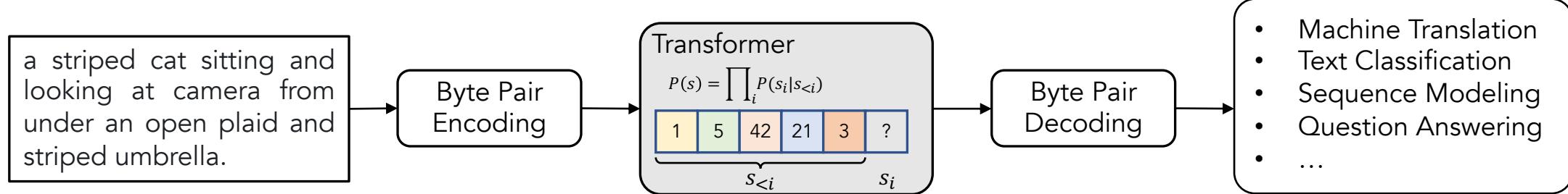


Pose Estimation



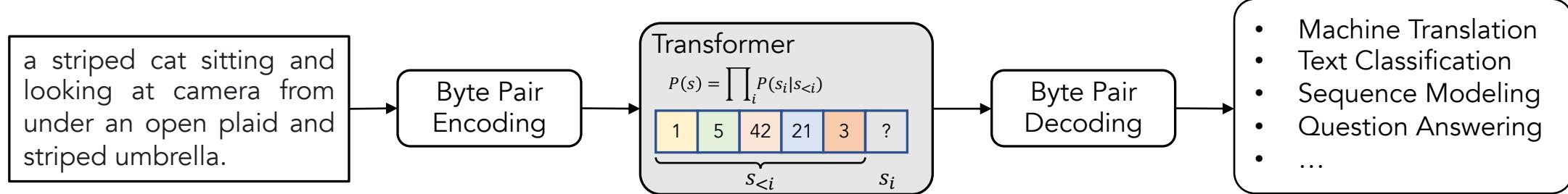
# Unified IO

- NLP Task Pipeline

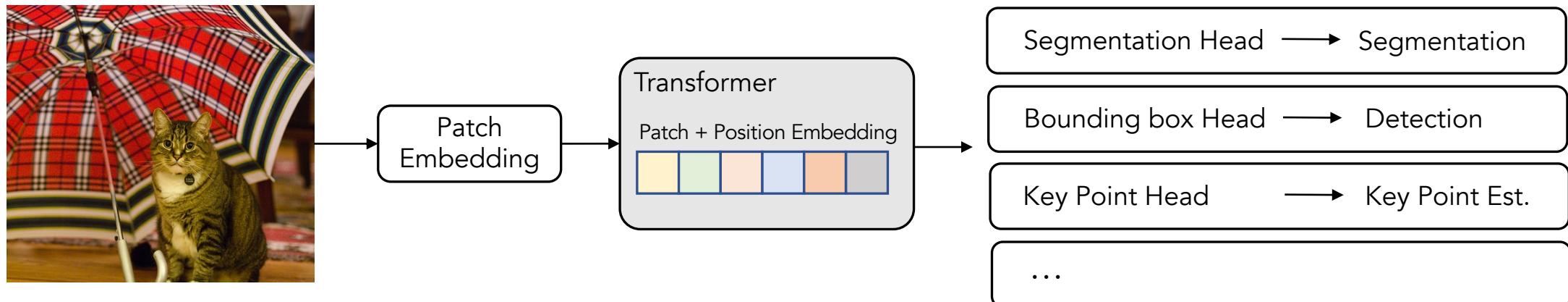


# Unified IO

- NLP Task Pipeline

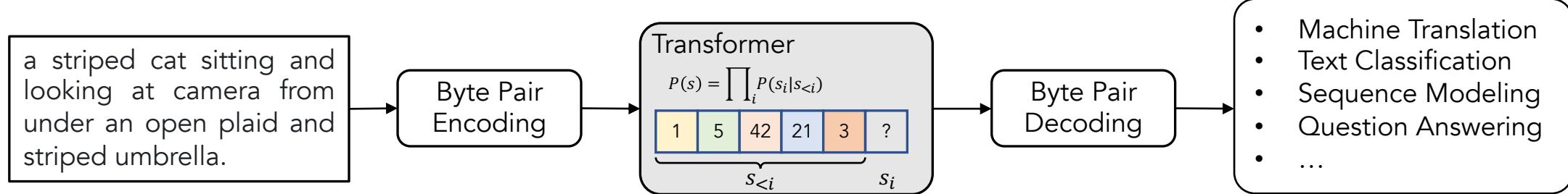


- Vision Task Pipeline

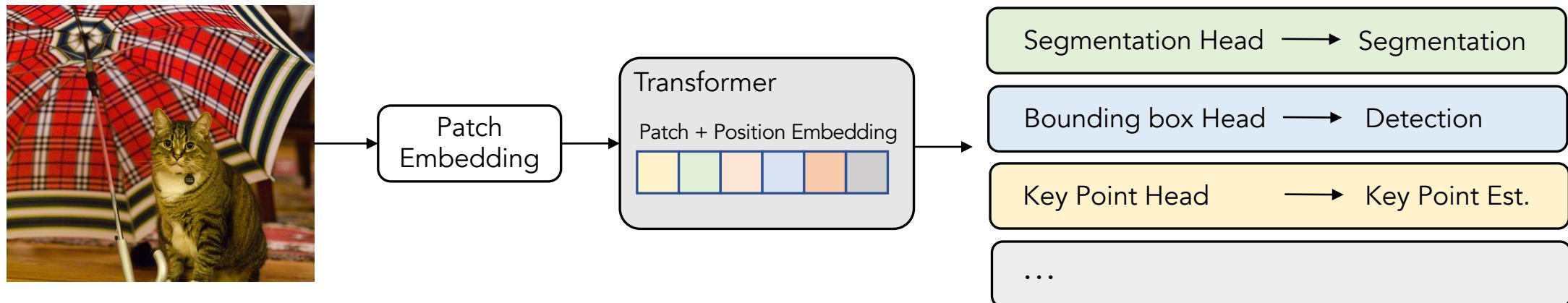


# Unified IO

- NLP Task Pipeline

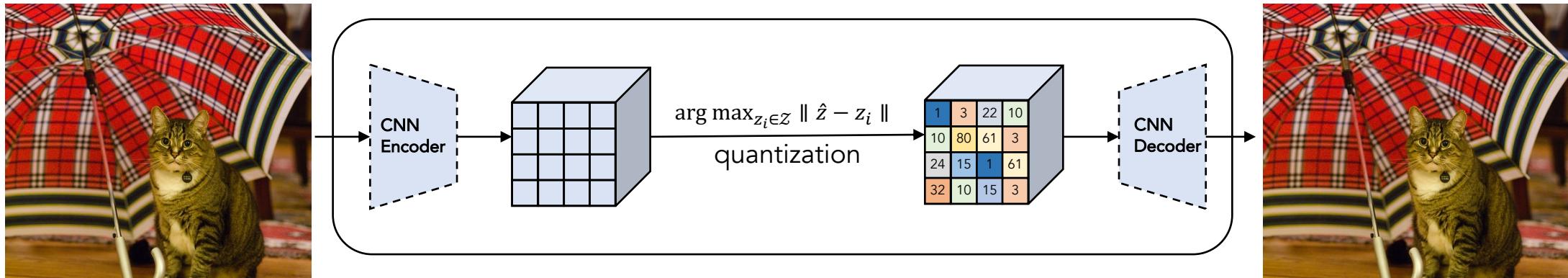


- Vision Task Pipeline



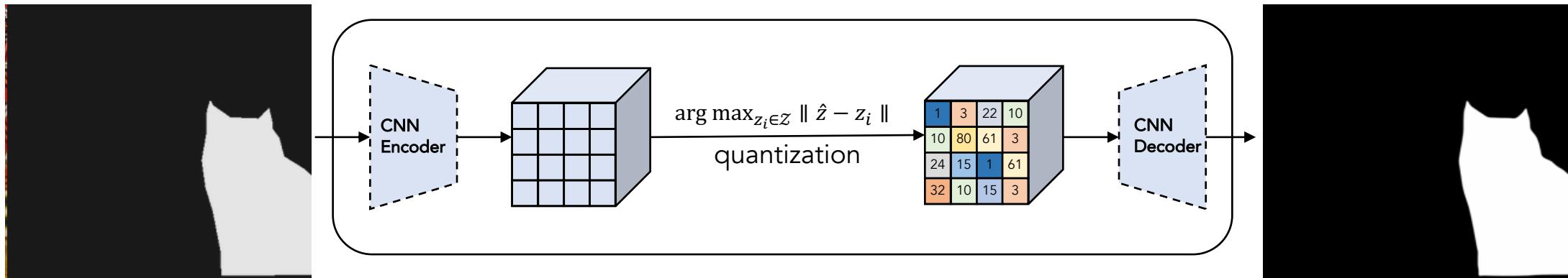
# Unified IO

- Image Quantization Using VQ-VAE



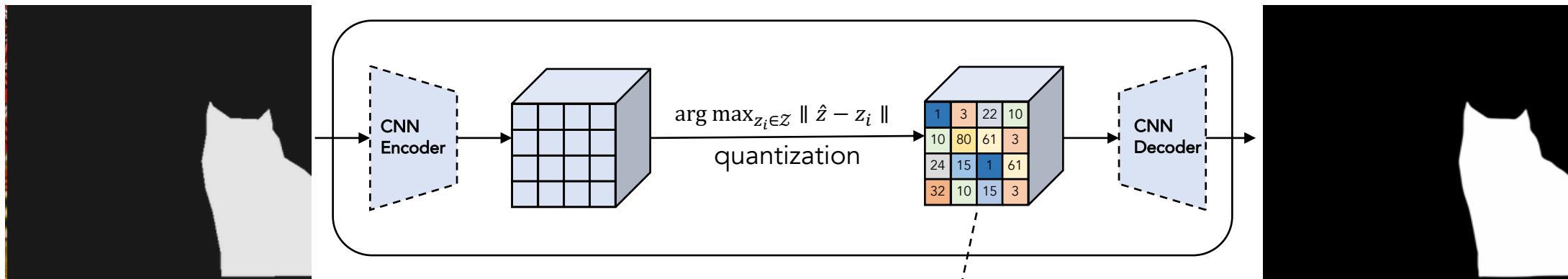
# Unified IO

- Image Quantization Using VQ-VAE

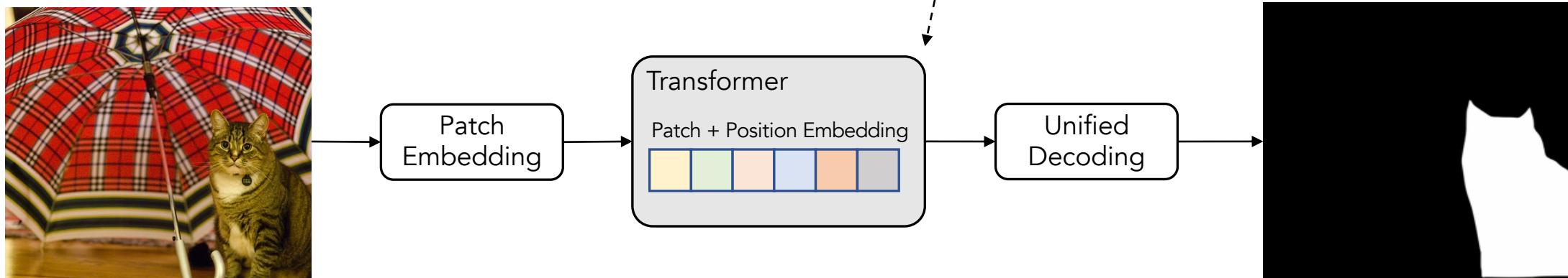


# Unified IO

- Image Quantization Using VQ-VAE

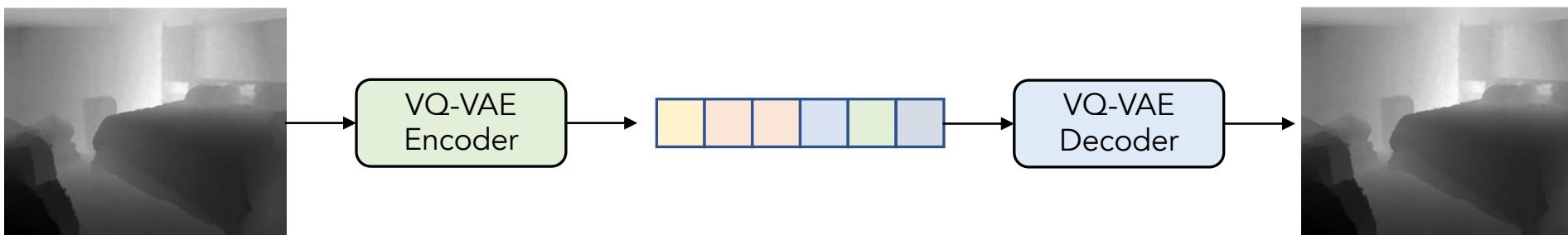
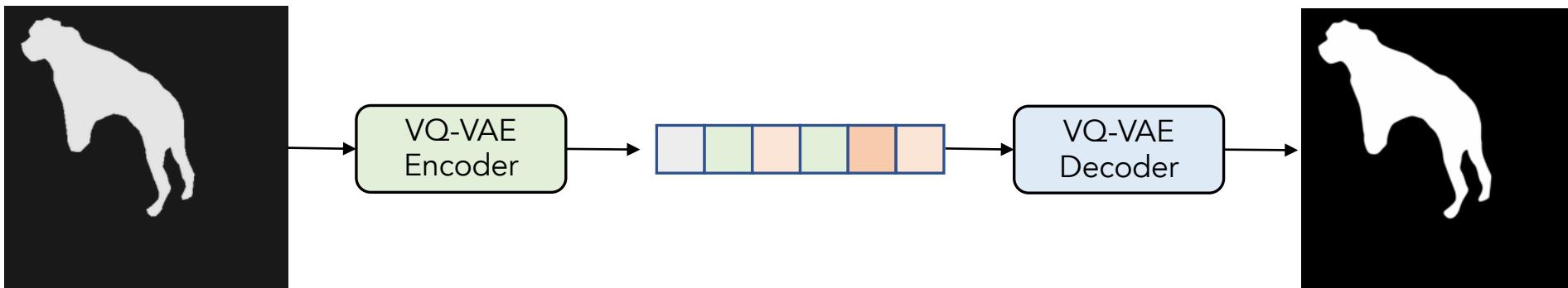
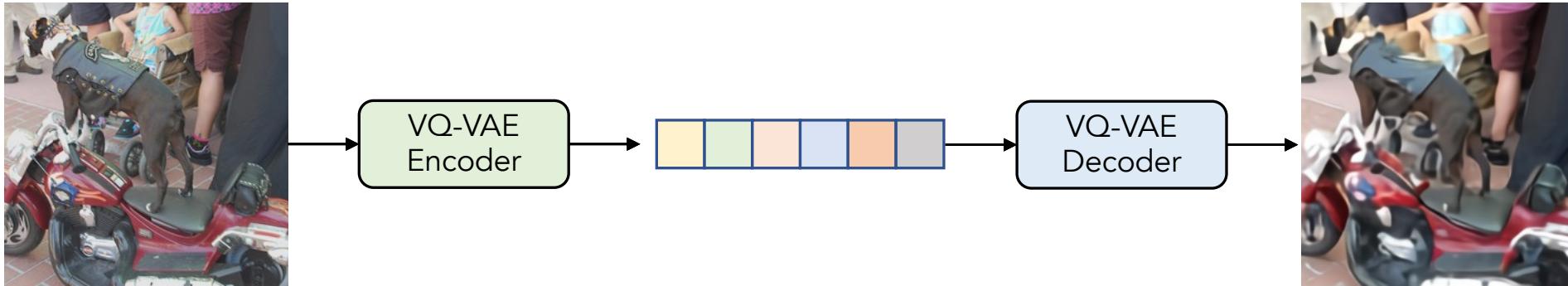


- Image Segmentation Using Unified Decoding



# Image Output Quantization

- Task with 2D structured outputs.

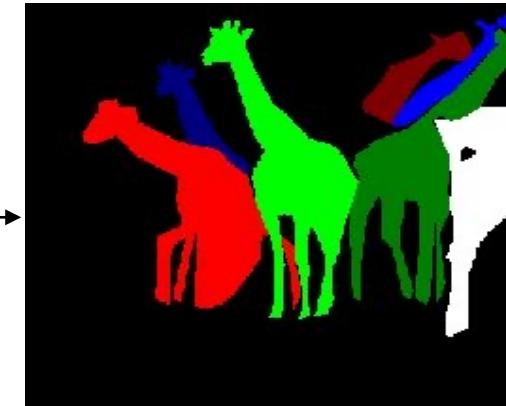
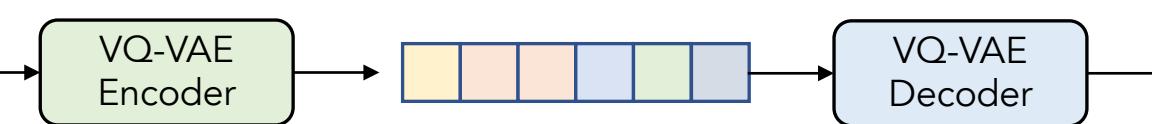
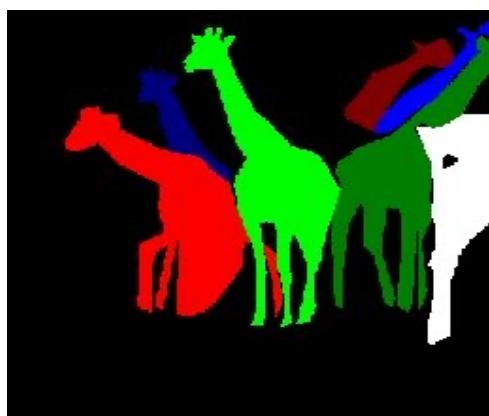


# Image Output Quantization

- Task with 2D structured outputs.

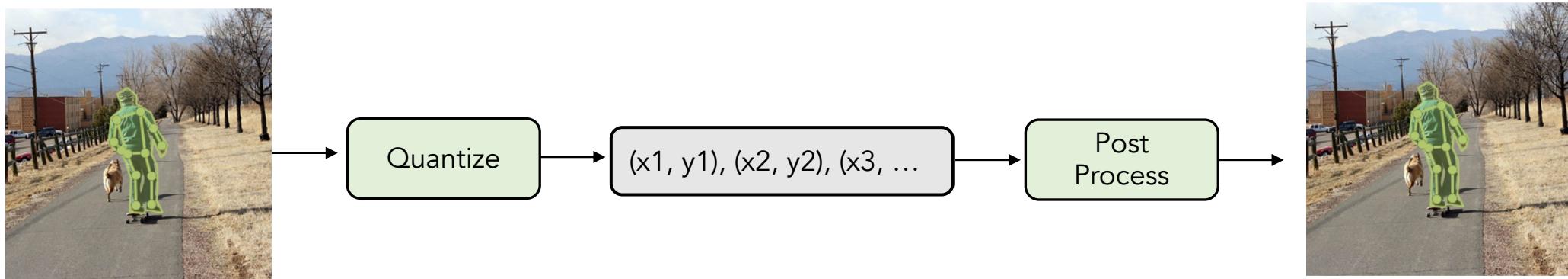
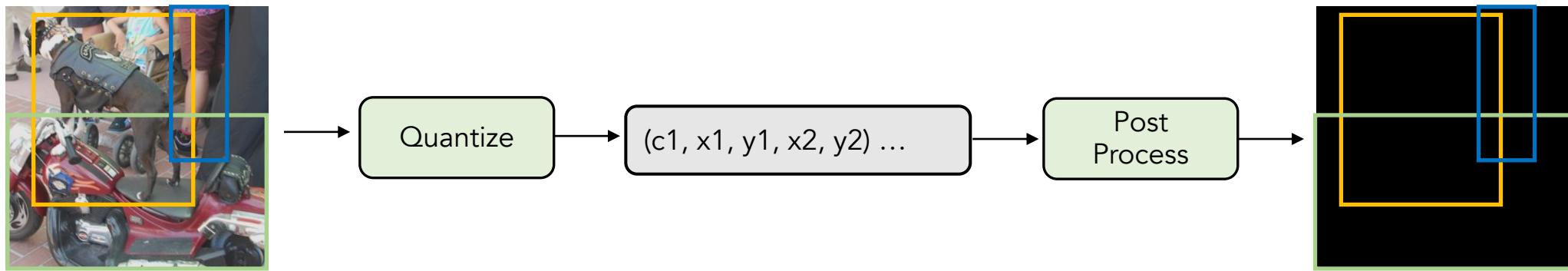


What is the segmentation of " giraffe " ?



# Image Output Quantization

- Task with 1D structured outputs



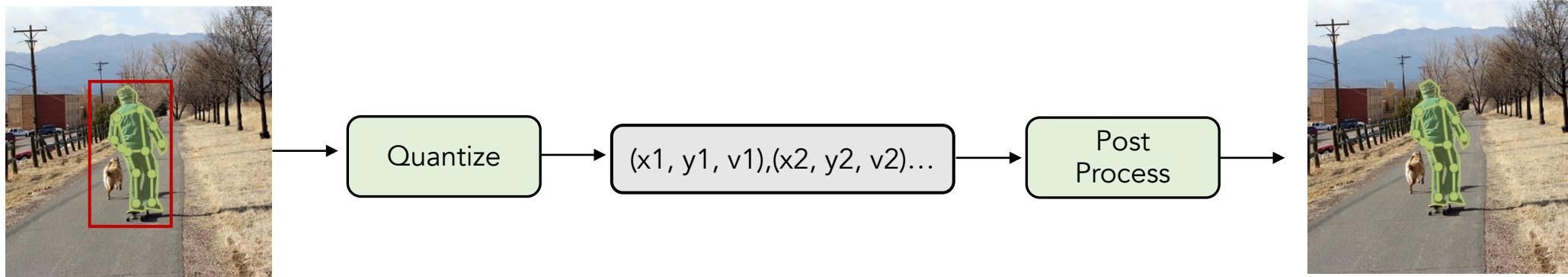
# Image Output Quantization

- Task with 1D structured outputs

Sequence length:  $17 \times 3 = 51$  per person

Take bounding box as input and detect one person each time.

human pose estimation --> detection + single human pose estimation



# Text Input for Different Tasks

- Prompt driven task specification

Categorization



Localization



VQA



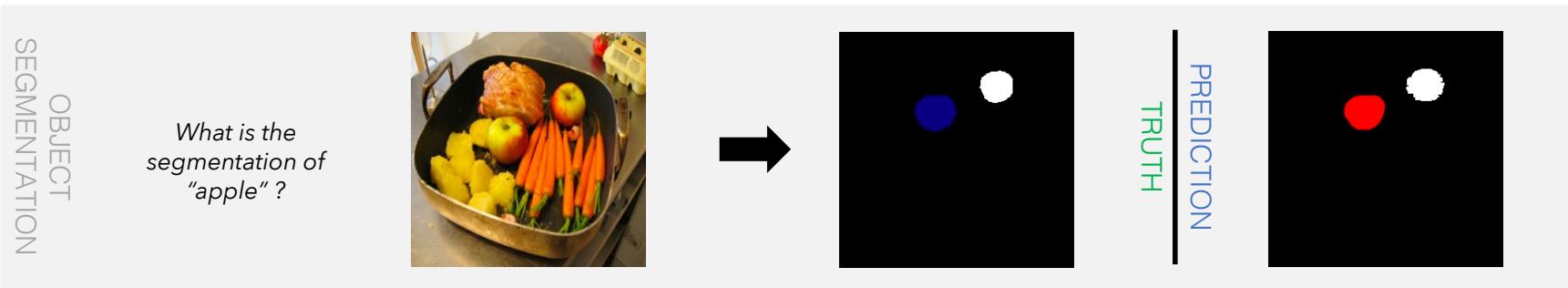
# Text Input for Different Tasks

- Prompt driven task specification

Refer Expression



Segmentation



Keypoints



# Text Input for Different Tasks

- Prompt driven task specification

Surface Normal

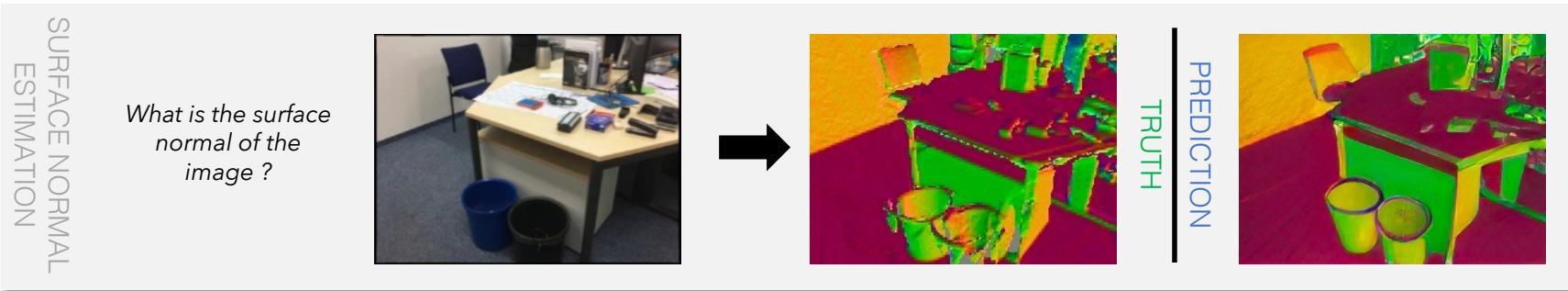


Image Generation



Seg based Image Generation



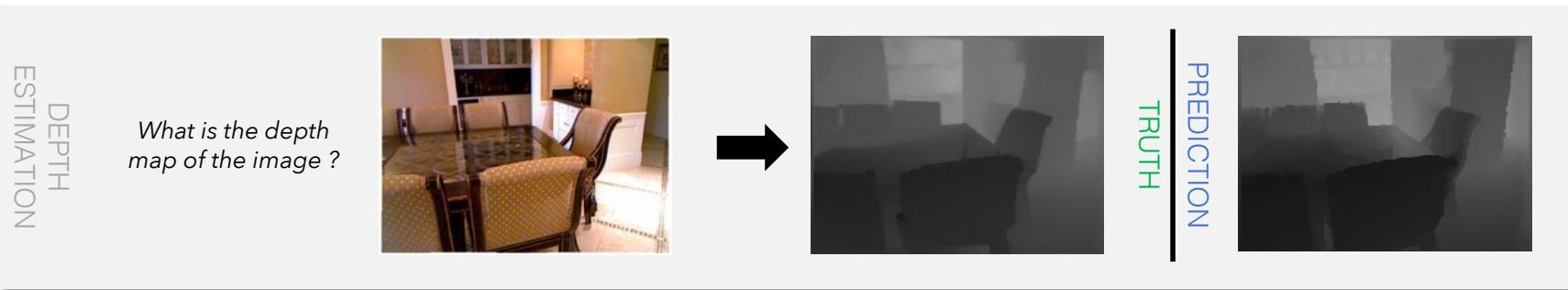
# Text Input for Different Tasks

- Prompt driven task specification

Image Inpainting



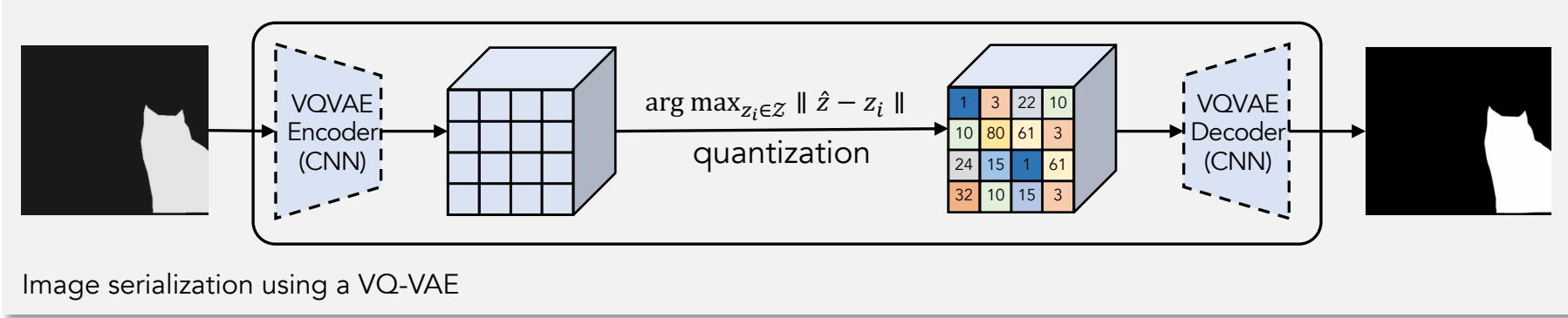
Depth Prediction



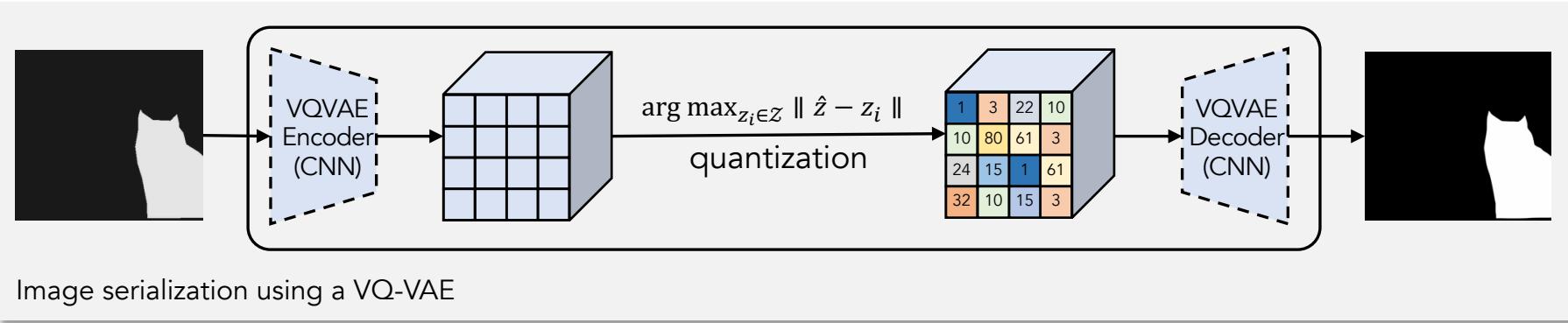
NLP Tasks



# Model



# Model



*Segment the cat*

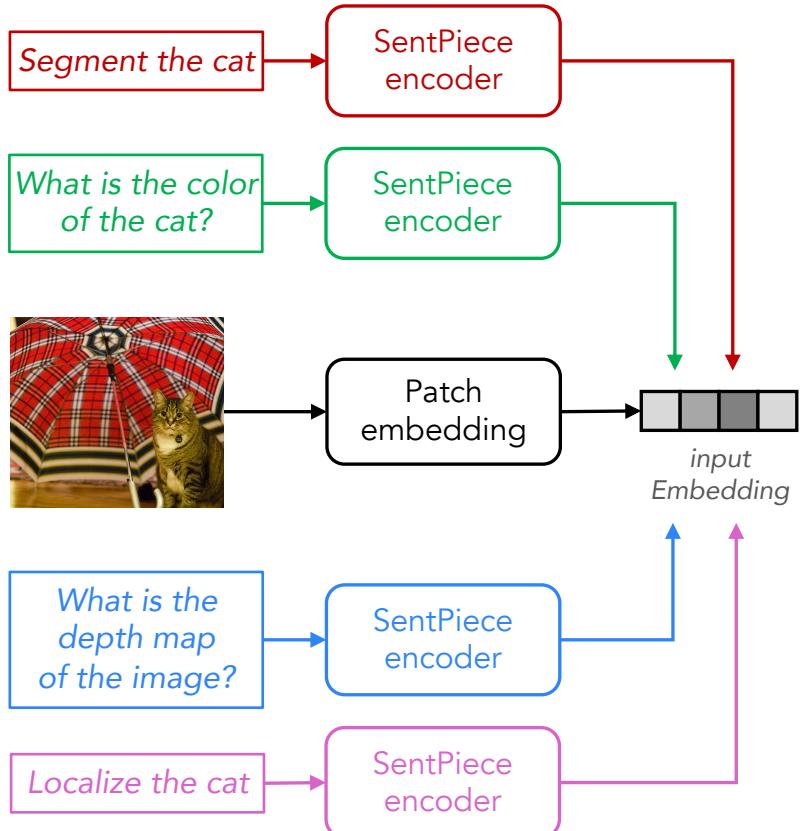
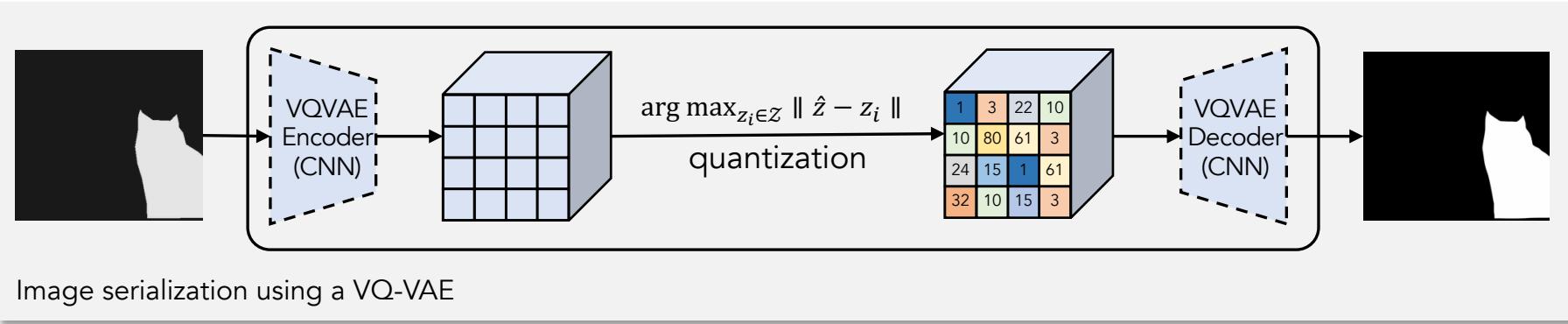
*What is the color  
of the cat?*



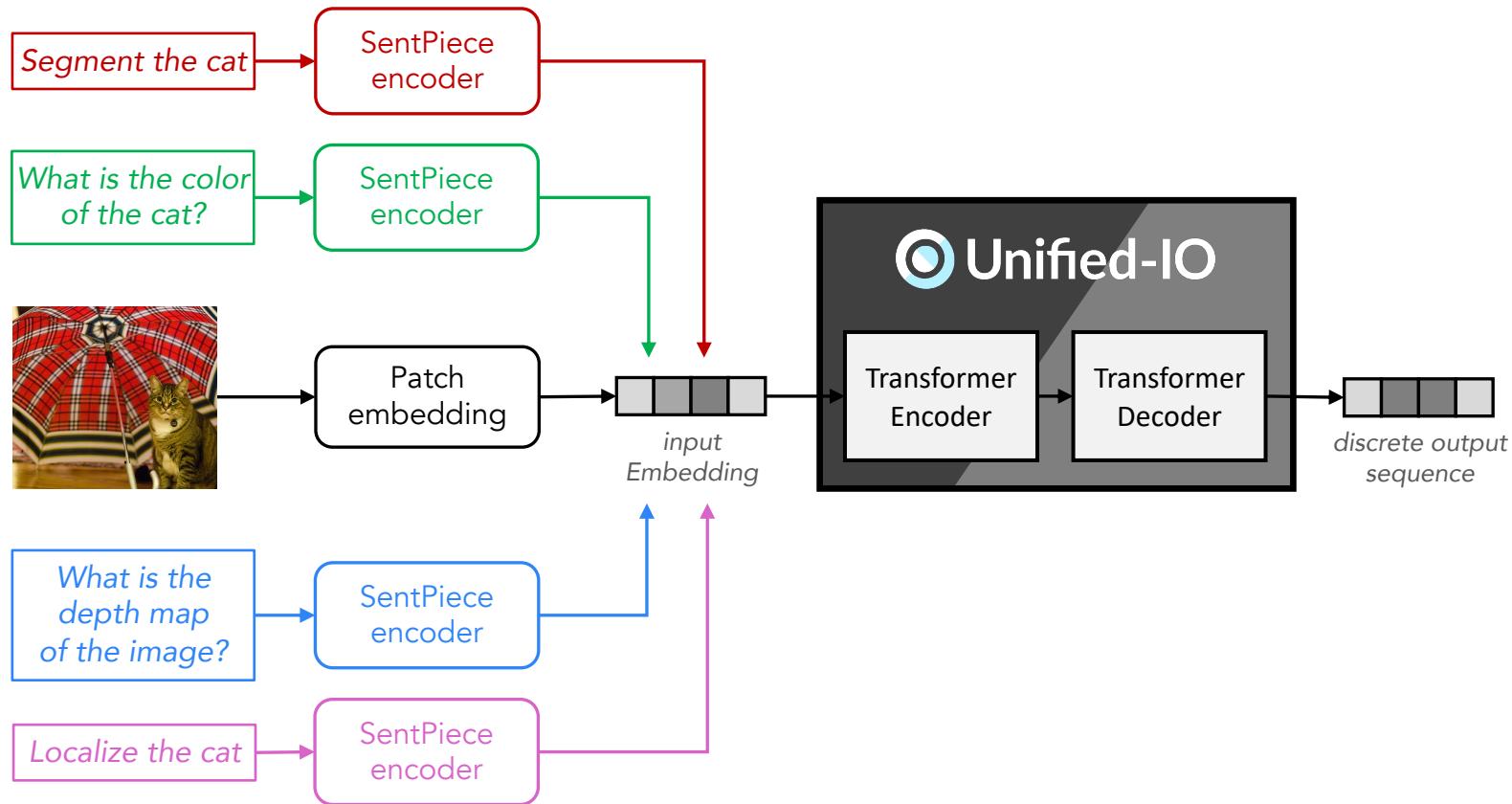
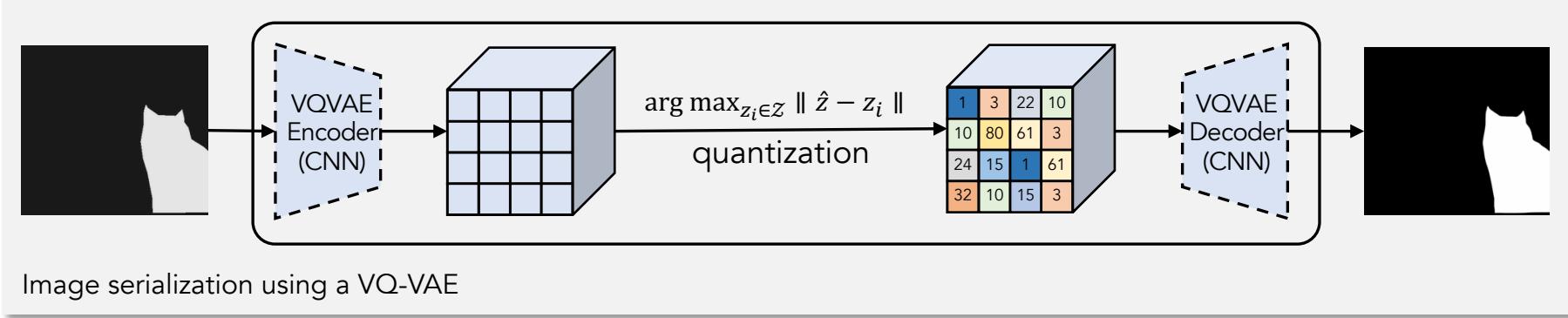
*What is the  
depth map  
of the image?*

*Localize the cat*

# Model



# Model



# Model

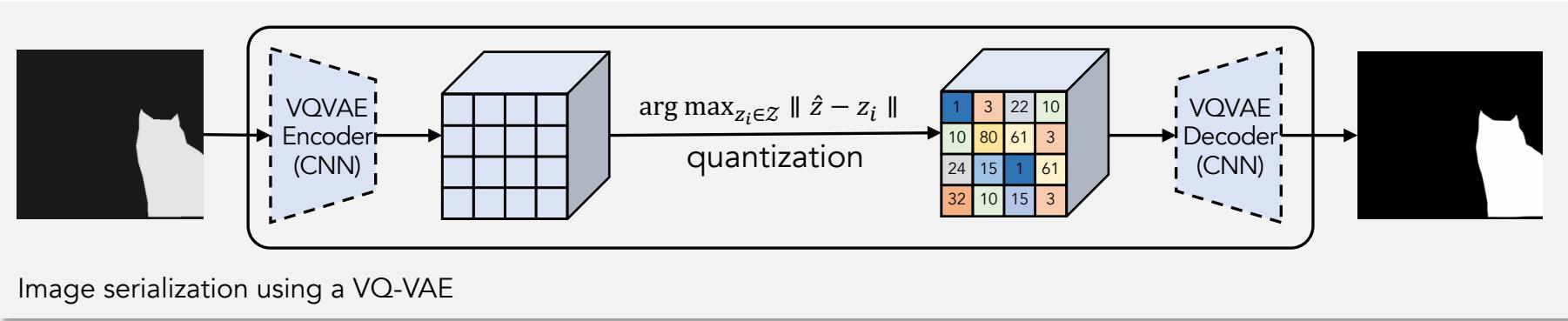
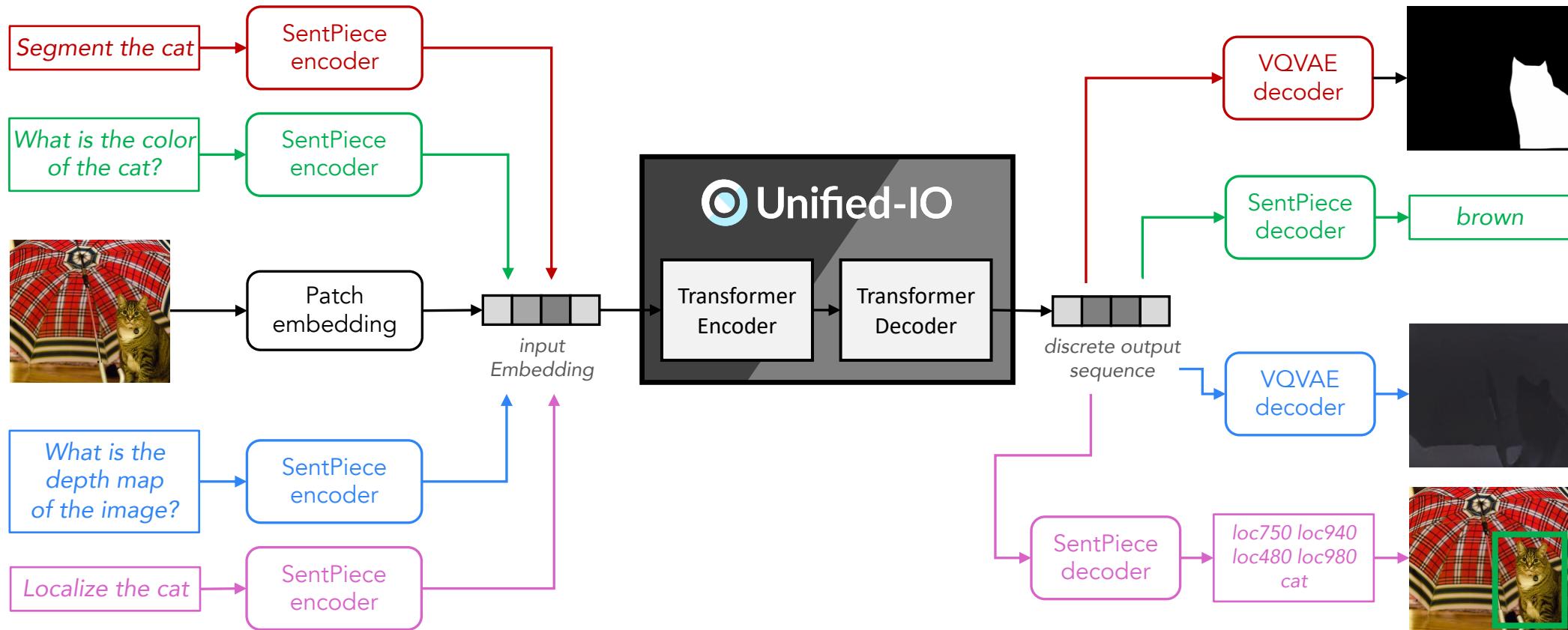
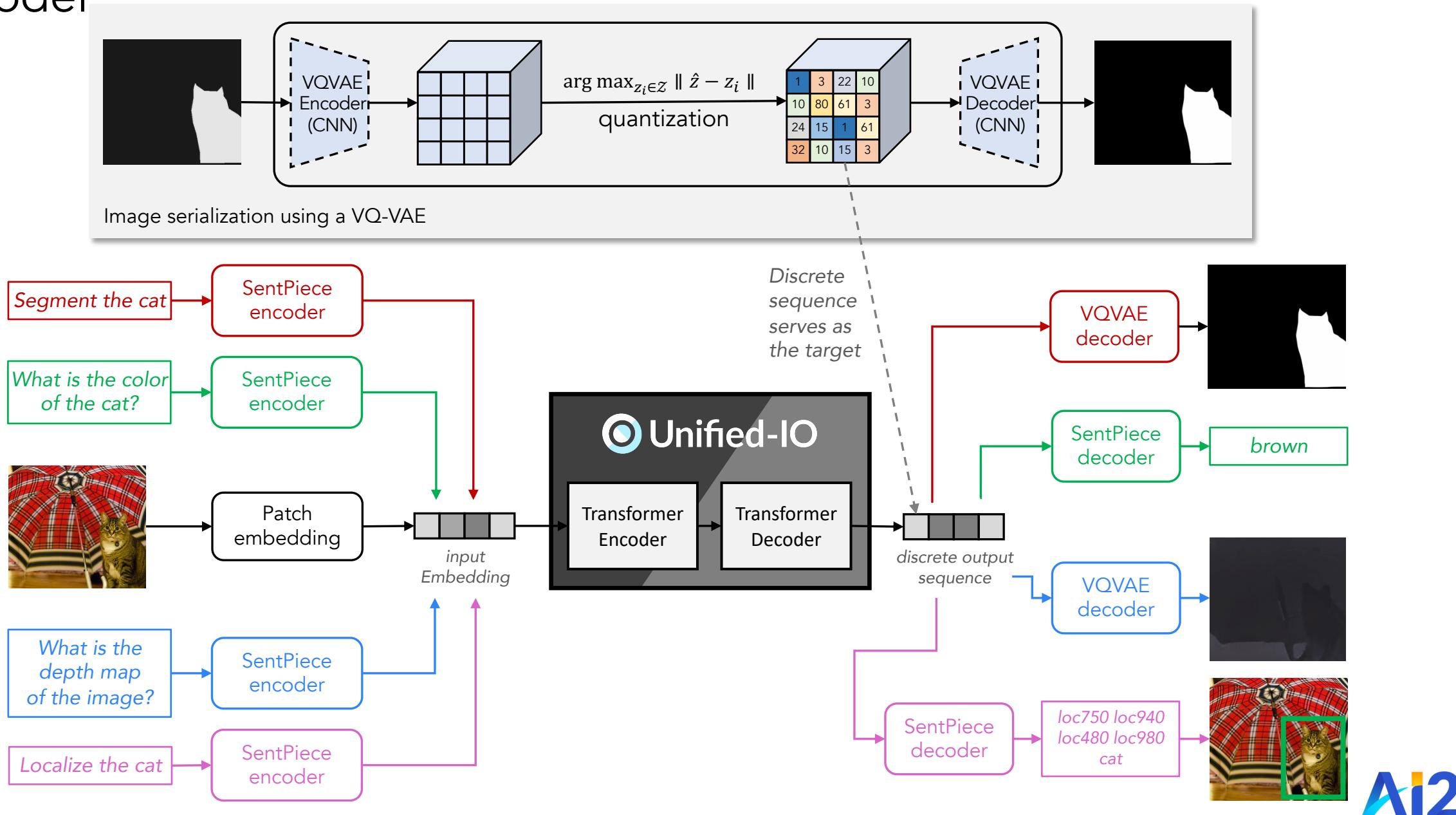


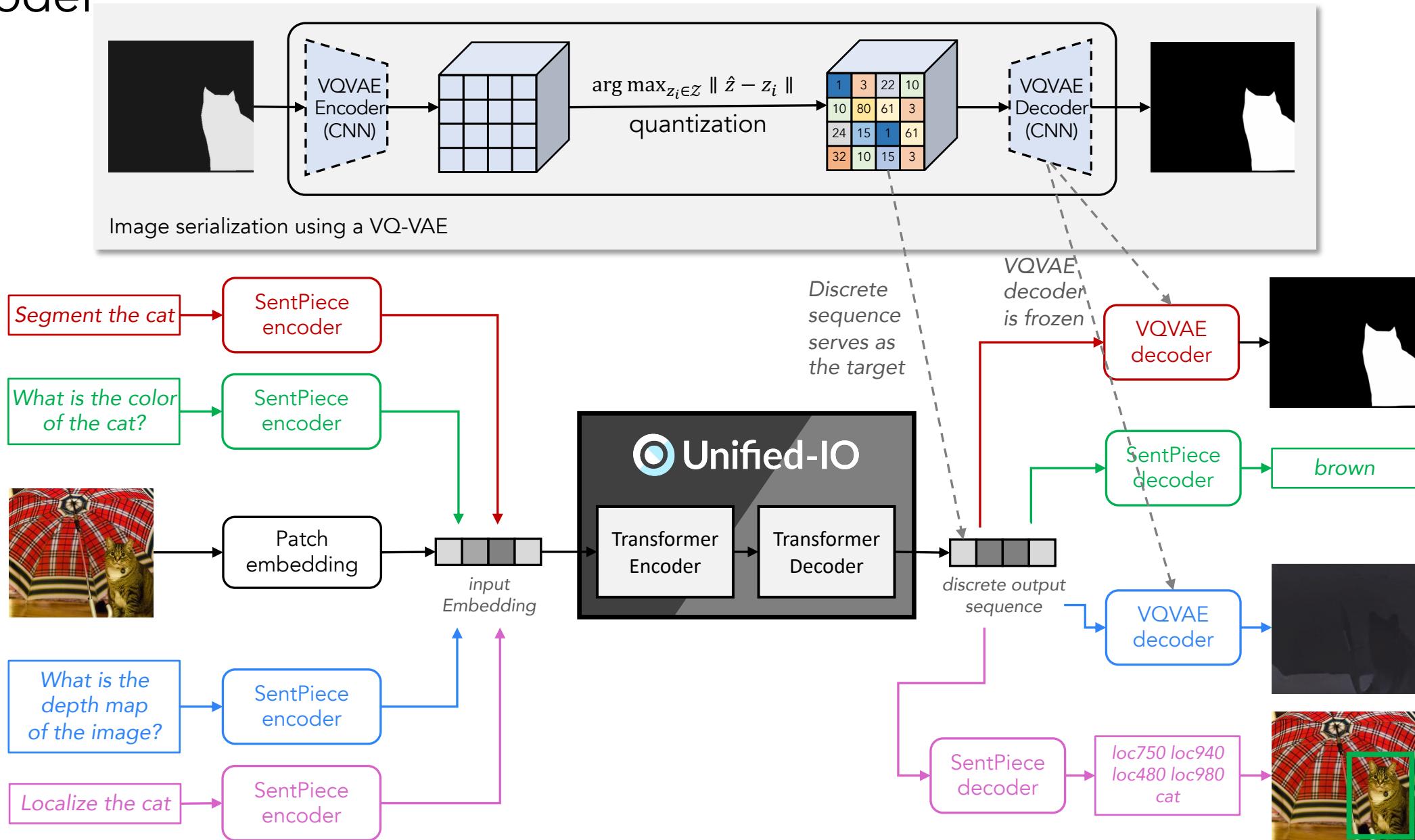
Image serialization using a VQ-VAE



# Model



## Model



# Model Details

- Follow T5 implementation with minimum modification.
- Use pretrained VQ-GAN from Imagenet.
- Use patch encoding for image inputs (no CNN involved)
- Absolute position encoding is important for vision – injecting the position embedding at input
- Relative position encoding within text and image.

```
[ 5,  4,  4,  3,  2,  1,  0,  1,  2,  3,  4,  4,  5,  5],
```



Relative encoding for the text

```
[[45, 44, 44, 43, 42, 41, 40, 41, 42, 43, 44, 44, 45, 45],  
 [37, 36, 36, 35, 34, 33, 32, 33, 34, 35, 36, 36, 37, 37],  
 [37, 36, 36, 35, 34, 33, 32, 33, 34, 35, 36, 36, 37, 37],  
 [29, 28, 28, 27, 26, 25, 24, 25, 26, 27, 28, 28, 29, 29],  
 [21, 20, 20, 19, 18, 17, 16, 17, 18, 19, 20, 20, 21, 21],  
 [13, 12, 12, 11, 10, 9, 8, 9, 10, 11, 12, 12, 13, 13],  
 [ 5,  4,  4,  3,  2,  1,  0,  1,  2,  3,  4,  4,  5,  5],  
 [13, 12, 12, 11, 10, 9, 8, 9, 10, 11, 12, 12, 13, 13],  
 [21, 20, 20, 19, 18, 17, 16, 17, 18, 19, 20, 20, 21, 21],  
 [29, 28, 28, 27, 26, 25, 24, 25, 26, 27, 28, 28, 29, 29],  
 [37, 36, 36, 35, 34, 33, 32, 33, 34, 35, 36, 36, 37, 37],  
 [37, 36, 36, 35, 34, 33, 32, 33, 34, 35, 36, 36, 37, 37],  
 [45, 44, 44, 43, 42, 41, 40, 41, 42, 43, 44, 44, 45, 45],  
 [45, 44, 44, 43, 42, 41, 40, 41, 42, 43, 44, 44, 45, 45]]
```

Relative encoding for the Image

# Objective

- **Pre-Training**

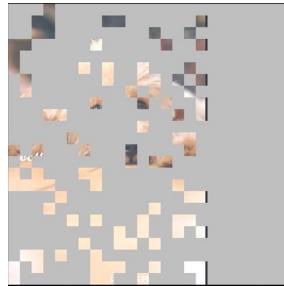
- Text span denoising (15%):

An image of a <M> is lying on the <M>.



<M> dog <M> ground.

- Mask image denoising (75%):



- **Multi-Task Training**

- Jointly train with 80 vision, vision language and language datasets/sets.

# Tasks

## Vision Tasks

Image Classification

Object Detection

Semantic Segmentation

Depth Estimation

Image Inpainting

Pose Estimation

Relationship Detection

## V&L Tasks

Image Captioning

Visual Question Answering

Image Generation

Refer Expression

VCR

Region Captioning

Situation Recognition

## Language Tasks

GLUE

Commonsense QA

SNLI

SQuAD

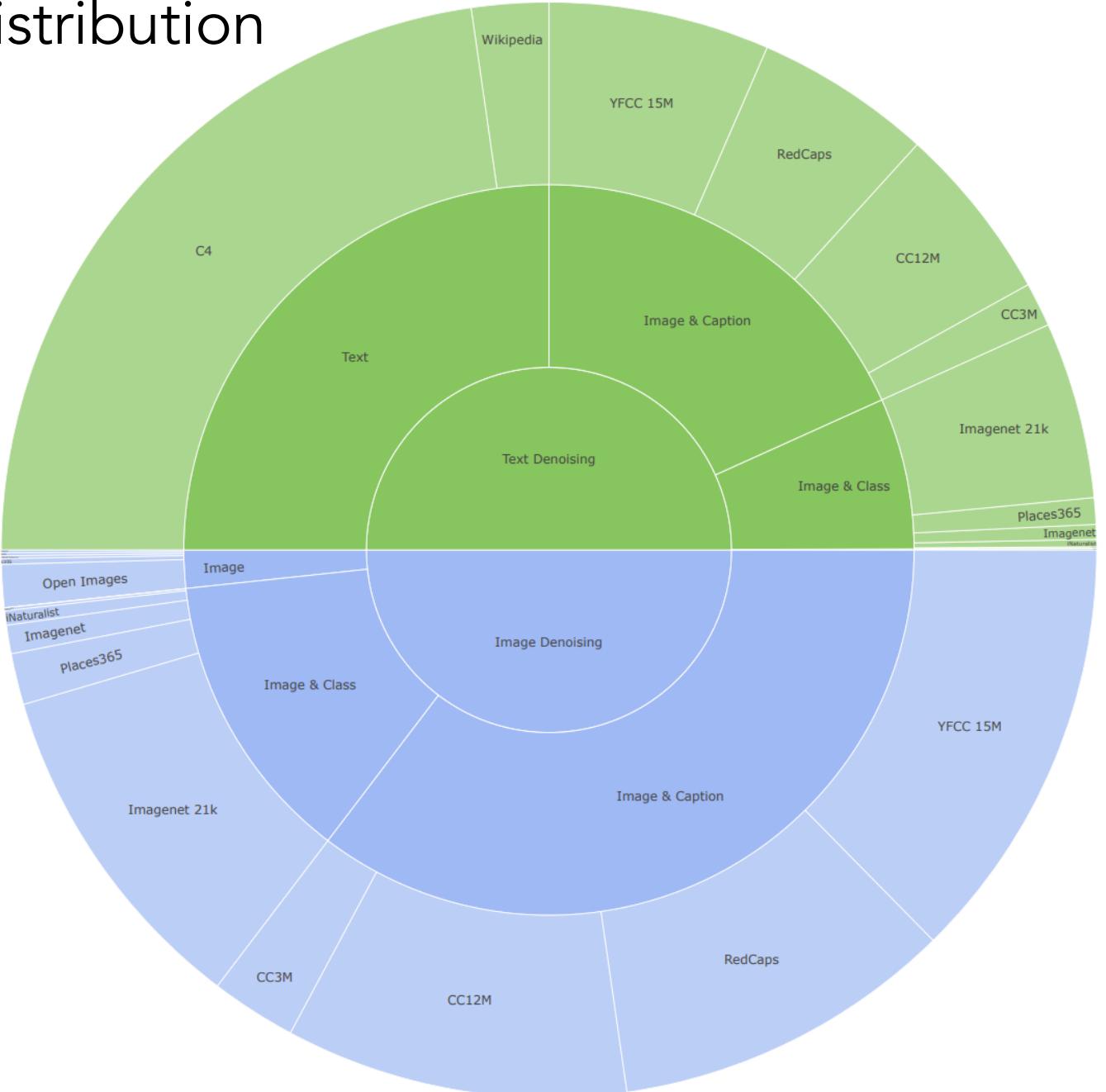
SWAG

Openbook QA

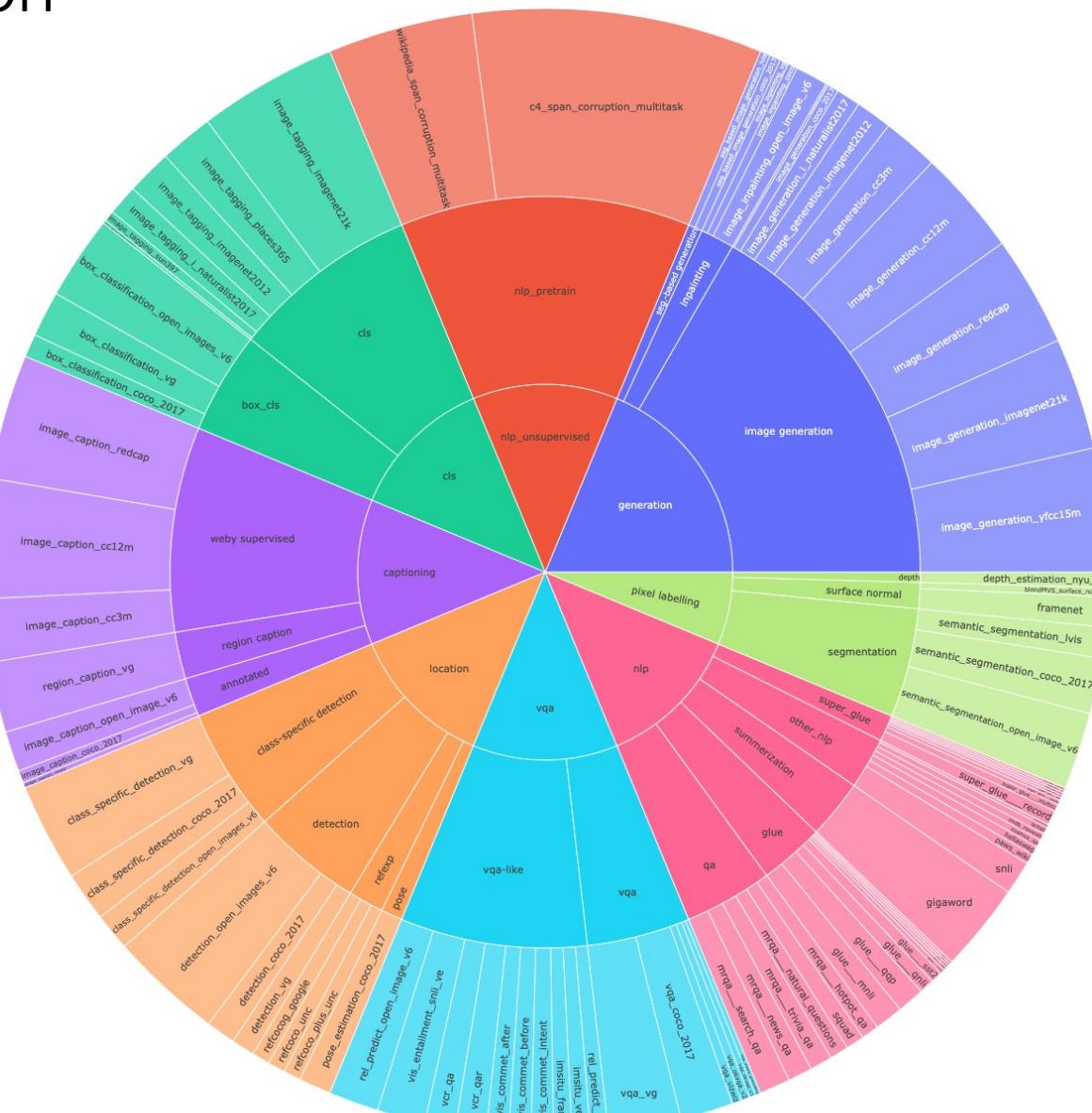
Fever

## Tasks

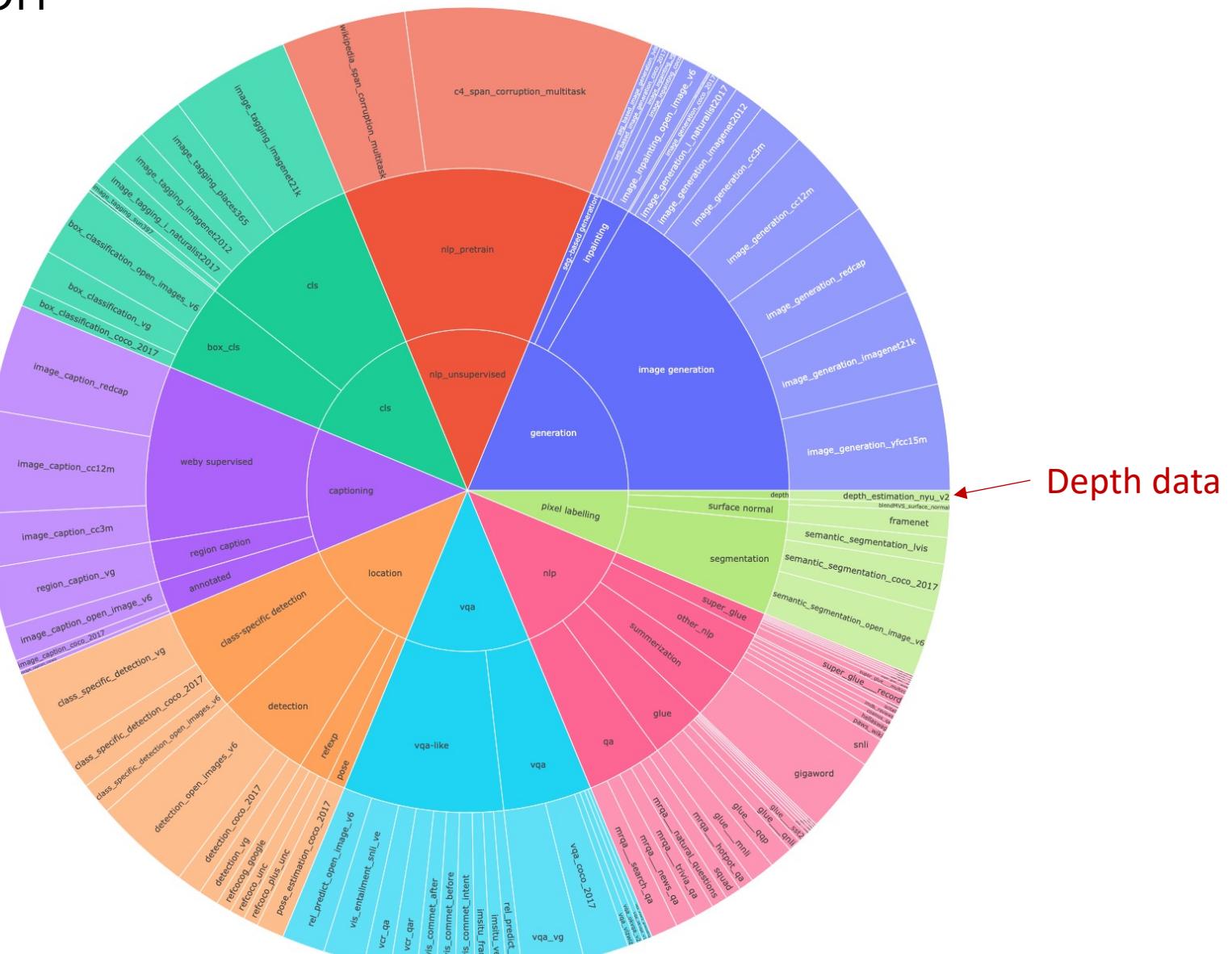
# Pre-training Distribution



# Tasks Distribution



# Tasks Distribution



# Evaluation



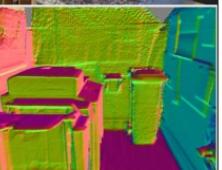
## General Robust Image Task Benchmark

Generality	Robustness	Calibration
Evaluate performance to novel domains and novel concepts for each task	Measure performance degradation with image perturbations	Quantify misinformation and confidence calibration
<b>Tasks</b>		
Object Categorization	Referring Expression Grounding	
Object Localization	Visual Question Answering	
Segmentation	Person Keypoint Detection	
Surface Normal Estimation		

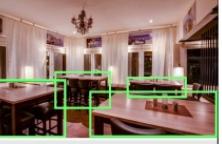
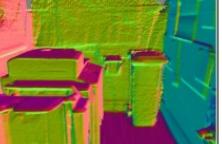
Task	Input Image	Input Query / Options	Output
Categorization		[open_images_categories]	drill
Localization		kitchen & dining room table	
Visual Question Answering		Does this sofa have armrests?	yes
Referring Expressions		man on end black suit	
Segmentation		dolphin	
Pose Keypoints		person	
Surface Normals			

GRIT: General Robust Image Task Benchmark [Gupta et.al. 2022]

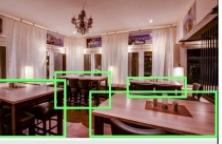
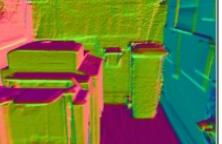
# GRIT requires diverse skills

Task	Input Image	Input Query / Options	Output	Inputs	Outputs
Categorization		[open_images_categories]	drill	→	Bounding box Class name
Localization		kitchen & dining room table		→	Class name Bounding box
Visual Question Answering		Does this sofa have armrests?	yes	→	Questions Answers
Referring Expressions		man on end black suit		→	Phrases Bounding box
Segmentation		dolphin		→	Class name Segmentation
Pose Keypoints		person		→	Joints + visibility
Surface Normals				→	Surface normal

# GRIT requires diverse skills

Task	Input Image	Input Query / Options	Output		Inputs	Outputs
Categorization		[open_images_categories]	drill	→	Bounding box	Class name
Localization		kitchen & dining room table		→	Class name	Bounding box
Visual Question Answering		Does this sofa have armrests?	yes	→	Questions	Answers
Referring Expressions		man on end black suit		→	Phrases	Bounding box
Segmentation		dolphin		→	Class name	Segmentation
Pose Keypoints		person		→		Joints + visibility
Surface Normals				→		Surface normal

# GRIT requires diverse skills

Task	Input Image	Input Query / Options	Output		Inputs	Outputs
Categorization		[open_images_categories]	drill	→	Bounding box	Class name
Localization		kitchen & dining room table		→	Class name	Bounding box
Visual Question Answering		Does this sofa have armrests?	yes	→	Questions	Answers
Referring Expressions		man on end black suit		→	Phrases	Bounding box
Segmentation		dolphin		→	Class name	Segmentation
Pose Keypoints		person		→		Joints + visibility
Surface Normals				→		Surface normal

# Results

	Categorization		Localization		VQA		Refexp		Segmentation		Keypoint		Normal		All	
	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test
0 NLL-AngMF [3]	-	-	-	-	-	-	-	-	-	-	-	-	<b>49.6</b>	<b>50.5</b>	7.2	7.1
1 Mask R-CNN [29]	-	-	44.7	45.1	-	-	-	-	26.2	26.2	<b>70.8</b>	<b>70.6</b>	-	-	20.2	20.3
2 GPV-1 [26]	33.2	33.2	42.8	42.7	50.6	49.8	25.8	26.8	-	-	-	-	-	-	21.8	21.8
3 CLIP [56]	48.1	-	-	-	-	-	-	-	-	-	-	-	-	-	6.9	-
4 OFA <sub>LARGE</sub> [73]	22.6	-	-	-	72.4	-	61.7	-	-	-	-	-	-	-	22.4	-
5 GPV-2 [36]	54.7	55.1	53.6	53.6	63.5	63.2	51.5	52.1	-	-	-	-	-	-	31.9	32.0
6 UNIFIED-IO <sub>SMALL</sub>	42.6	-	50.4	-	52.9	-	51.1	-	40.7	-	<b>46.5</b>	-	33.5	-	45.4	-
7 UNIFIED-IO <sub>BASE</sub>	53.1	-	59.7	-	63.0	-	68.3	-	49.3	-	60.2	-	37.5	-	55.9	-
8 UNIFIED-IO <sub>LARGE</sub>	57.0	-	64.2	-	67.4	-	74.1	-	54.0	-	<b>67.6</b>	-	40.2	-	57.0	-
9 UNIFIED-IO <sub>XL</sub>	<b>61.7</b>	<b>60.8</b>	<b>67.0</b>	<b>67.1</b>	<b>74.5</b>	<b>74.5</b>	<b>78.6</b>	<b>78.9</b>	<b>56.3</b>	<b>56.5</b>	68.1	67.7	45.0	44.3	<b>64.5</b>	<b>64.3</b>

Unified IO on GRIT ablation and test benchmark

# Results

	Categorization		Localization		VQA		Refexp		Segmentation		Keypoint		Normal		All	
	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test
0 NLL-AngMF [3]	-	-	-	-	-	-	-	-	-	-	-	-	<b>49.6</b>	<b>50.5</b>	7.2	7.1
1 Mask R-CNN [29]	-	-	44.7	45.1	-	-	-	-	26.2	26.2	<b>70.8</b>	<b>70.6</b>	-	-	20.2	20.3
2 GPV-1 [26]	33.2	33.2	42.8	42.7	50.6	49.8	25.8	26.8	-	-	-	-	-	-	21.8	21.8
3 CLIP [56]	48.1	-	-	-	-	-	-	-	-	-	-	-	-	-	6.9	-
4 OFA <sub>LARGE</sub> [73]	22.6	-	-	-	72.4	-	61.7	-	-	-	-	-	-	-	22.4	-
5 GPV-2 [36]	54.7	55.1	53.6	53.6	63.5	63.2	51.5	52.1	-	-	-	-	-	-	31.9	<b>32.0</b>
6 UNIFIED-IO <sub>SMALL</sub>	42.6	-	50.4	-	52.9	-	51.1	-	40.7	-	<b>46.5</b>	-	33.5	-	45.4	-
7 UNIFIED-IO <sub>BASE</sub>	53.1	-	59.7	-	63.0	-	68.3	-	49.3	-	<b>60.2</b>	-	37.5	-	55.9	-
8 UNIFIED-IO <sub>LARGE</sub>	57.0	-	64.2	-	67.4	-	74.1	-	54.0	-	<b>67.6</b>	-	40.2	-	57.0	-
9 UNIFIED-IO <sub>XL</sub>	<b>61.7</b>	<b>60.8</b>	<b>67.0</b>	<b>67.1</b>	<b>74.5</b>	<b>74.5</b>	<b>78.6</b>	<b>78.9</b>	<b>56.3</b>	<b>56.5</b>	68.1	67.7	45.0	44.3	<b>64.5</b>	<b>64.3</b>

Unified IO on GRIT ablation and test benchmark

# Results

	Categorization		Localization		VQA		Refexp		Segmentation		Keypoint		Normal		All	
	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test
0 NLL-AngMF [3]	-	-	-	-	-	-	-	-	-	-	-	-	<b>49.6</b>	<b>50.5</b>	7.2	7.1
1 Mask R-CNN [29]	-	-	44.7	45.1	-	-	-	-	26.2	26.2	<b>70.8</b>	<b>70.6</b>	-	-	20.2	20.3
2 GPV-1 [26]	33.2	33.2	42.8	42.7	50.6	49.8	25.8	26.8	-	-	-	-	-	-	21.8	21.8
3 CLIP [56]	48.1	-	-	-	-	-	-	-	-	-	-	-	-	-	6.9	-
4 OFA <sub>LARGE</sub> [73]	22.6	-	-	-	72.4	-	61.7	-	-	-	-	-	-	-	22.4	-
5 GPV-2 [36]	54.7	55.1	53.6	53.6	63.5	63.2	51.5	52.1	-	-	-	-	-	-	31.9	32.0
6 UNIFIED-IO <sub>SMALL</sub>	42.6	-	50.4	-	52.9	-	51.1	-	40.7	-	<b>46.5</b>	-	33.5	-	45.4	-
7 UNIFIED-IO <sub>BASE</sub>	53.1	-	59.7	-	63.0	-	68.3	-	49.3	-	60.2	-	37.5	-	55.9	-
8 UNIFIED-IO <sub>LARGE</sub>	57.0	-	64.2	-	67.4	-	74.1	-	54.0	-	67.6	-	40.2	-	57.0	-
9 UNIFIED-IO <sub>XL</sub>	<b>61.7</b>	<b>60.8</b>	<b>67.0</b>	<b>67.1</b>	<b>74.5</b>	<b>74.5</b>	<b>78.6</b>	<b>78.9</b>	<b>56.3</b>	<b>56.5</b>	68.1	67.7	45.0	44.3	<b>64.5</b>	<b>64.3</b>

Unified IO on GRIT ablation and test benchmark

# Results

	Categorization		Localization		VQA		Refexp		Segmentation		Keypoint		Normal		All	
	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test
0 NLL-AngMF [3]	-	-	-	-	-	-	-	-	-	-	-	-	<b>49.6</b>	<b>50.5</b>	7.2	7.1
1 Mask R-CNN [29]	-	-	44.7	45.1	-	-	-	-	26.2	26.2	<b>70.8</b>	<b>70.6</b>	-	-	20.2	20.3
2 GPV-1 [26]	33.2	33.2	42.8	42.7	50.6	49.8	25.8	26.8	-	-	-	-	-	-	21.8	21.8
3 CLIP [56]	48.1	-	-	-	-	-	-	-	-	-	-	-	-	-	6.9	-
4 OFA <sub>LARGE</sub> [73]	22.6	-	-	-	72.4	-	61.7	-	-	-	-	-	-	-	22.4	-
5 GPV-2 [36]	54.7	55.1	53.6	53.6	63.5	63.2	51.5	52.1	-	-	-	-	-	-	31.9	32.0
6 UNIFIED-IO <sub>SMALL</sub>	42.6	-	50.4	-	52.9	-	51.1	-	40.7	-	46.5	-	33.5	-	45.4	-
7 UNIFIED-IO <sub>BASE</sub>	53.1	-	59.7	-	63.0	-	68.3	-	49.3	-	60.2	-	37.5	-	55.9	-
8 UNIFIED-IO <sub>LARGE</sub>	57.0	-	64.2	-	67.4	-	74.1	-	54.0	-	67.6	-	40.2	-	57.0	-
9 UNIFIED-IO <sub>XL</sub>	<b>61.7</b>	<b>60.8</b>	<b>67.0</b>	<b>67.1</b>	<b>74.5</b>	<b>74.5</b>	<b>78.6</b>	<b>78.9</b>	<b>56.3</b>	<b>56.5</b>	68.1	<b>67.7</b>	45.0	44.3	<b>64.5</b>	<b>64.3</b>

Unified IO on GRIT ablation and test benchmark

# Results

	Categorization		Localization		VQA		Refexp		Segmentation		Keypoint		Normal		All		
	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	
0 NLL-AngMF [3]	-	-	-	-	-	-	-	-	-	-	-	-	<b>49.6</b>	<b>50.5</b>	7.2	7.1	
1 Mask R-CNN [29]	-	-	44.7	45.1	-	-	-	-	26.2	26.2	<b>70.8</b>	<b>70.6</b>	-	-	20.2	20.3	
2 GPV-1 [26]	33.2	33.2	42.8	42.7	50.6	49.8	25.8	26.8	-	-	-	-	-	-	-	21.8	21.8
3 CLIP [56]	48.1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	6.9	-
4 OFA <sub>LARGE</sub> [73]	22.6	-	-	-	72.4	-	61.7	-	-	-	-	-	-	-	-	22.4	-
5 GPV-2 [36]	54.7	55.1	53.6	53.6	63.5	63.2	51.5	52.1	-	-	-	-	-	-	-	31.9	32.0
6 UNIFIED-IO <sub>SMALL</sub>	42.6	-	50.4	-	52.9	-	51.1	-	40.7	-	46.5	-	33.5	-	45.4	-	
7 UNIFIED-IO <sub>BASE</sub>	53.1	-	59.7	-	63.0	-	68.3	-	49.3	-	60.2	-	37.5	-	55.9	-	
8 UNIFIED-IO <sub>LARGE</sub>	57.0	-	64.2	-	67.4	-	74.1	-	54.0	-	67.6	-	40.2	-	57.0	-	
9 UNIFIED-IO <sub>XL</sub>	<b>61.7</b>	<b>60.8</b>	<b>67.0</b>	<b>67.1</b>	<b>74.5</b>	<b>74.5</b>	<b>78.6</b>	<b>78.9</b>	<b>56.3</b>	<b>56.5</b>	68.1	67.7	45.0	<b>44.3</b>	<b>64.5</b>	<b>64.3</b>	

Unified IO on GRIT ablation and test benchmark

# Results

	Categorization		Localization		VQA		Refexp		Segmentation		Keypoint		Normal		All	
	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test	ablation	test
0 NLL-AngMF [3]	-	-	-	-	-	-	-	-	-	-	-	-	<b>49.6</b>	<b>50.5</b>	7.2	7.1
1 Mask R-CNN [29]	-	-	44.7	45.1	-	-	-	-	26.2	26.2	<b>70.8</b>	<b>70.6</b>	-	-	20.2	20.3
2 GPV-1 [26]	33.2	33.2	42.8	42.7	50.6	49.8	25.8	26.8	-	-	-	-	-	-	21.8	21.8
3 CLIP [56]	48.1	-	-	-	-	-	-	-	-	-	-	-	-	-	6.9	-
4 OFA <sub>LARGE</sub> [73]	22.6	-	-	-	72.4	-	61.7	-	-	-	-	-	-	-	22.4	-
5 GPV-2 [36]	54.7	55.1	53.6	53.6	63.5	63.2	51.5	52.1	-	-	-	-	-	-	31.9	32.0
6 UNIFIED-IO <sub>SMALL</sub>	42.6	-	50.4	-	52.9	-	51.1	-	40.7	-	46.5	-	33.5	-	<b>45.4</b>	-
7 UNIFIED-IO <sub>BASE</sub>	53.1	-	59.7	-	63.0	-	68.3	-	49.3	-	60.2	-	37.5	-	<b>55.9</b>	-
8 UNIFIED-IO <sub>LARGE</sub>	57.0	-	64.2	-	67.4	-	74.1	-	54.0	-	67.6	-	40.2	-	<b>57.0</b>	-
9 UNIFIED-IO <sub>XL</sub>	<b>61.7</b>	<b>60.8</b>	<b>67.0</b>	<b>67.1</b>	<b>74.5</b>	<b>74.5</b>	<b>78.6</b>	<b>78.9</b>	<b>56.3</b>	<b>56.5</b>	68.1	67.7	45.0	44.3	<b>64.5</b>	<b>64.3</b>

Unified IO on GRIT ablation and test benchmark

# Results

	<i>NYUv2</i>	<i>ImageNet</i>	<i>Place365</i>	<i>VQAv2</i>	<i>OkVQA</i>	<i>A-O<sub>k</sub>VQA</i>	<i>VizWizQA</i>	<i>VizWizGround</i>	<i>Swig</i>	<i>SNLI-VE</i>	<i>VisComet</i>	<i>Nocaps</i>	<i>COCO</i>	<i>COCO</i>	<i>MRPC</i>	<i>BoolQ</i>	<i>SciTail</i>	
Split	val	val	val	test-dev	test	test	test-dev	test-std	test	val	val	val	val	test	val	val	test	
Metric	RMSE	Acc.	Acc.	Acc.	Acc.	Acc.	Acc.	IOU	Acc.	Acc.	CIDEr	CIDEr	CIDEr	CIDEr	F1	Acc	Acc	
Unified SOTA	UViM 0.467	- -	- -	- -	Flamingo 57.8	- -	Flamingo 49.8	- -	- -	- -	- -	- -	- -	- -	T5 92.20	PaLM 92.2	- -	
UNIFIED-IO <sub>SMALL</sub>	0.649	42.8	38.2	57.7	31.0	24.3	42.4	35.5	17.3	76.5	-	45.1	80.1	-	84.9	65.9	87.4	
UNIFIED-IO <sub>BASE</sub>	0.469	63.3	43.2	61.8	37.8	28.5	45.8	50.0	29.7	85.6	-	66.9	104.0	-	87.9	70.8	90.8	
UNIFIED-IO <sub>LARGE</sub>	0.402	71.8	50.5	67.8	42.7	33.4	47.7	54.7	40.4	86.1	-	87.2	117.5	-	87.5	73.1	93.1	
UNIFIED-IO <sub>XL</sub>	0.385	79.1	53.2	77.9	54.0	45.2	57.4	65.0	49.8	91.1	21.2	100.0	126.8	122.3	89.2	79.7	95.7	
Single or fine-tuned SOTA	BinsFormer 0.330	CoCa 91.00	MAE 60.3	CoCa 82.3	KAT 54.4	GPV2 38.1	Flamingo 65.7	MAC-Caps 27.3	JSL 39.6	OFA 91.0	SVT 18.3	CoCa 122.4	-	OFA 145.3	Turing 93.8	NLR 92.4	ST-MOE 97.7	DeBERTa

Unified IO on other tasks

# Visualization – image synthesis

IMAGE SYNTHESIS

*What is the complete image? Text*

INPUT

*small personal pizza with bacon and spinach*

*many large kites flying in the sky*

*the train is on the tracks in the station*

*a beach area with black birds flying over it*

PREDICTION

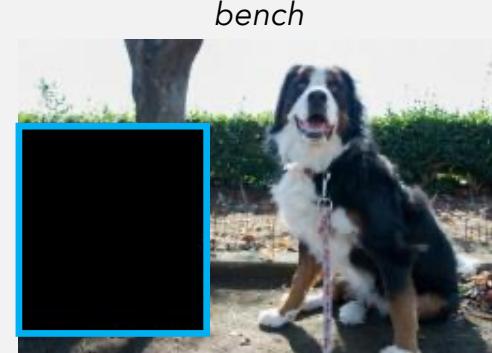
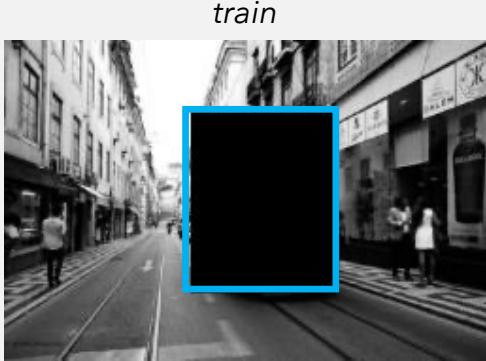


# Visualization – image synthesis

## IMAGE INPAINTING

INPUT

PREDICTION



*Fill in the blank region with this*



# Visualization – image synthesis

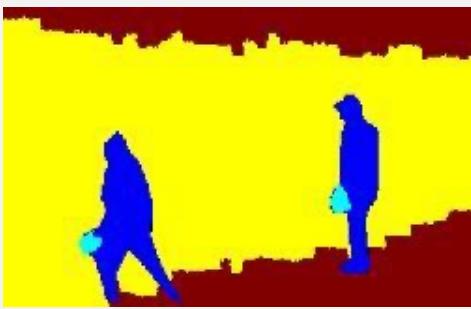
## IMAGE GENERATION FROM SEGMENTATION

fuchsia: tree, lime: dirt....

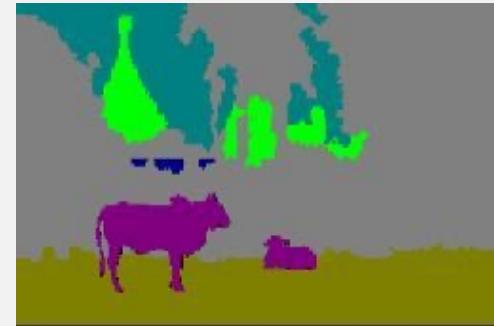


What is the complete image? Segmentation color

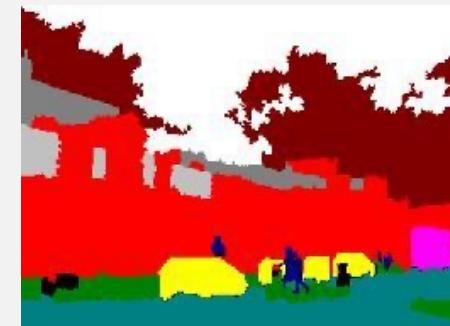
yellow: field, aqua: baseball...



lime: building, navy: wall...



white: tree, red: building...



INPUT

PREDICTION

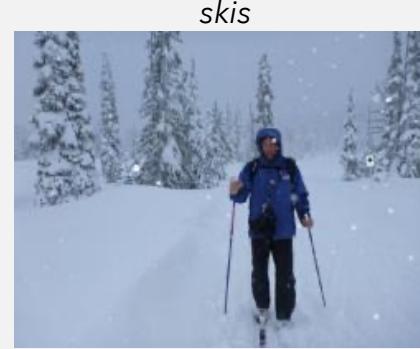


# Visualization – sparse labeling

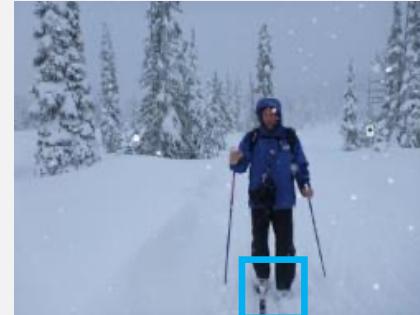
## OBJECT LOCALIZATION

INPUT

PREDICTION



*What region does this describe?*



# Visualization – sparse labeling

## OBJECT DETECTION

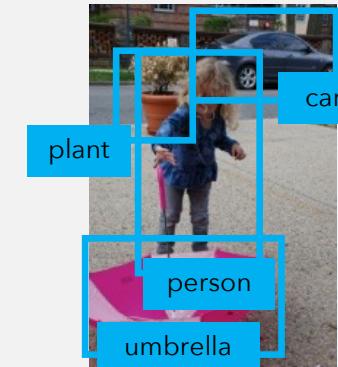
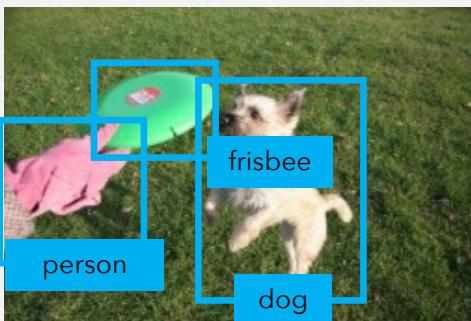
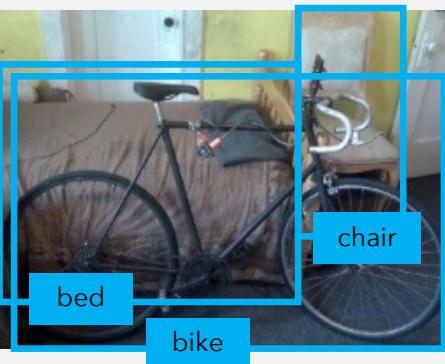
INPUT



*Locate all objects in the image.*



PREDICTION



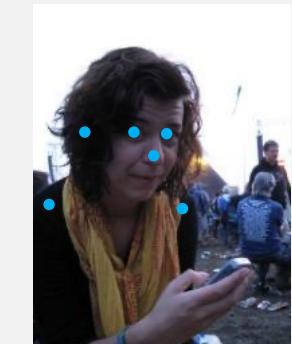
# Visualization – sparse labeling

## KEYPOINT ESTIMATION

INPUT



PREDICTION



*Find the human joints in this region.*

# Visualization – dense labeling

DEPTH ESTIMATION

*What is the depth map of the image?*

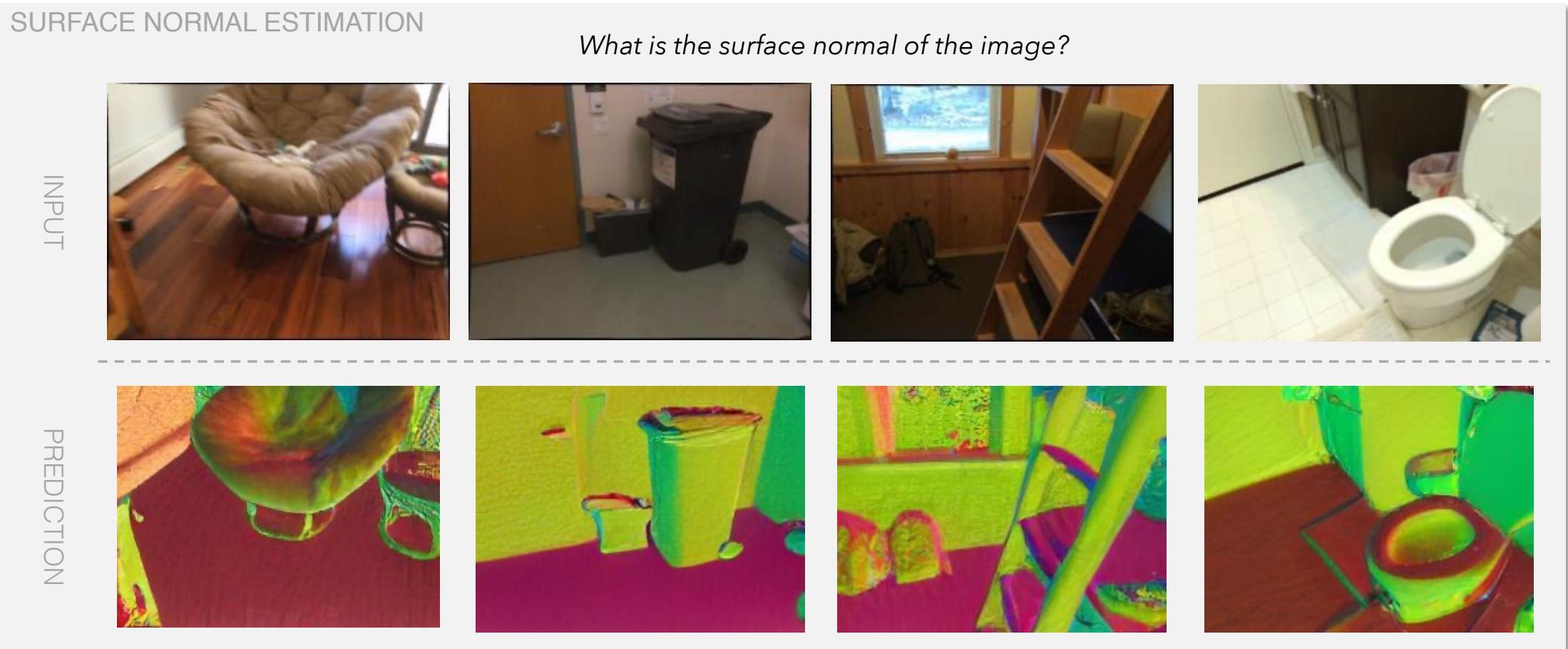
INPUT



PREDICTION



# Visualization – dense labeling



# Visualization – dense labeling

INPUT

OBJECT SEGMENTATION

*What is the segmentation of this?*

*pizza*



*bed*



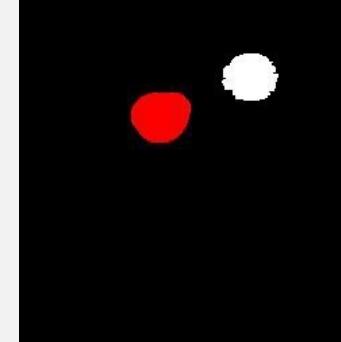
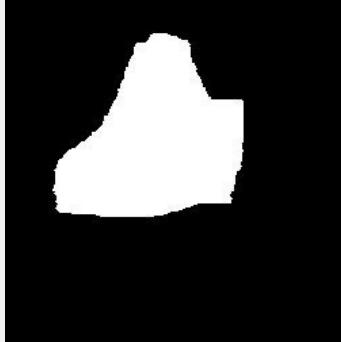
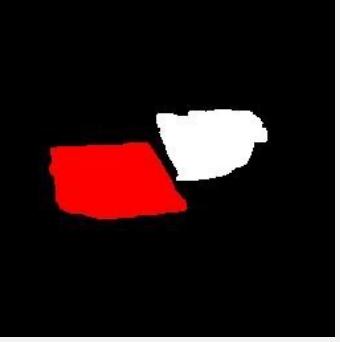
*apple*



*grass*



PREDICTION



# Referring expressions using different prompts

Prompt	Refexp Score
0 Which region does the text “ REFEXP ” describe ?	78.9
1 Which region does the text “REFEXP” describe?	76.7
2 Which region matches the text “ REFEXP ” ?	77.4
3 Locate the “ REFEXP ” .	65.6
4 Which region can be described as “ REFEXP ” ?	64.8
5 Locate the region described by “ REFEXP ” .	43.2
6 Where is the “ REFEXP ” ?	41.5
7 Where is the “REFEXP”?	0.1

# Summary

- Propose unified IO which is the first framework that can handle massive vision, vision – language, and language tasks.
- We treat 2D image tasks as condition image generation tasks.
- We use pre-trained VQ-GAN to convert images into discrete sequences.
- We will release the Code + pre-trained model.

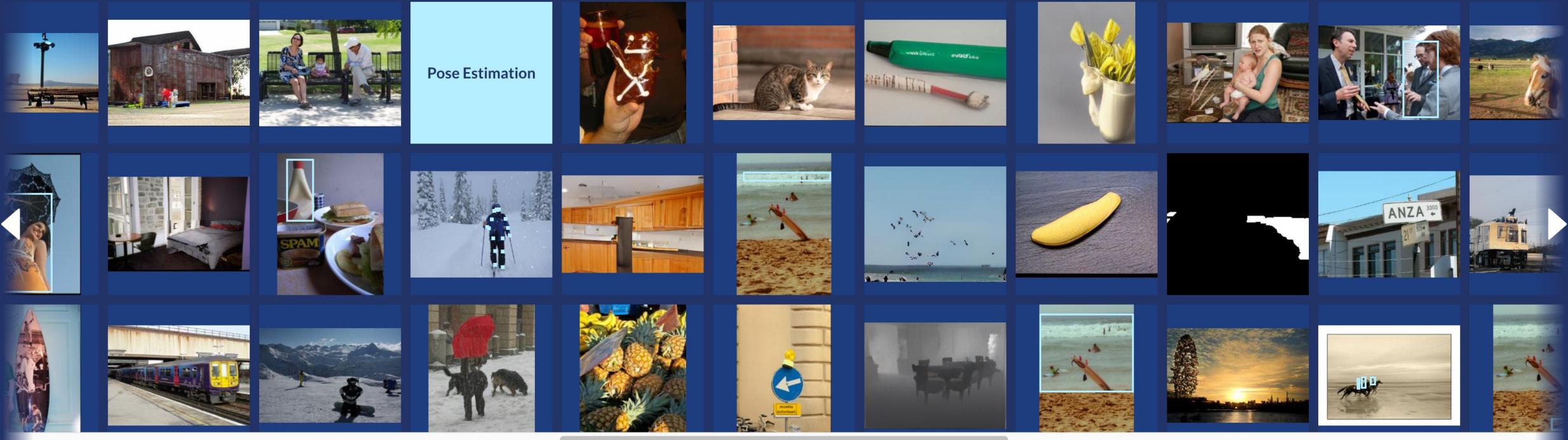
## Acknowledgement

- This research was made possible with cloud TPUs from [Google's TPU Research Cloud \(TRC\)](#).
- Google Jax and T5X team -- amazing libraries.
- AI2 Reviz Team – Great demo pages.



Share to: [Twitter](#) [LinkedIn](#)

A new general-purpose model with **unprecedented breadth**, Unified-IO can perform a wide array of **visual** and **linguistic** tasks.



<https://unified-io.allenai.org>

