



北京航空航天大學
BEIHANG UNIVERSITY

生成式AI与大模型第5讲

多模态对齐

Beihang University

人工智能研究院
黄雷

01

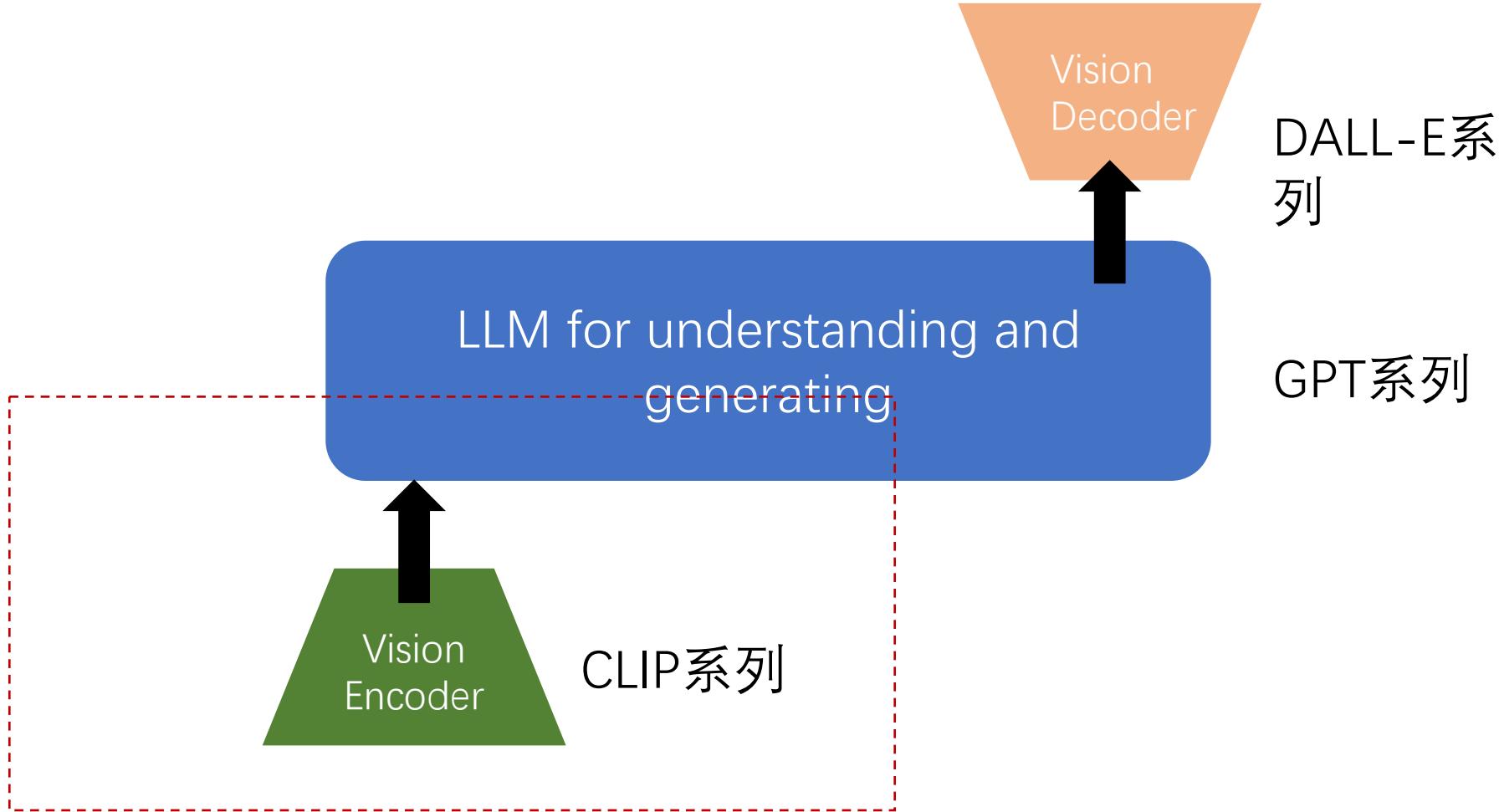
模态对齐

多模态和视觉大模型

图像：

文本：

图像：

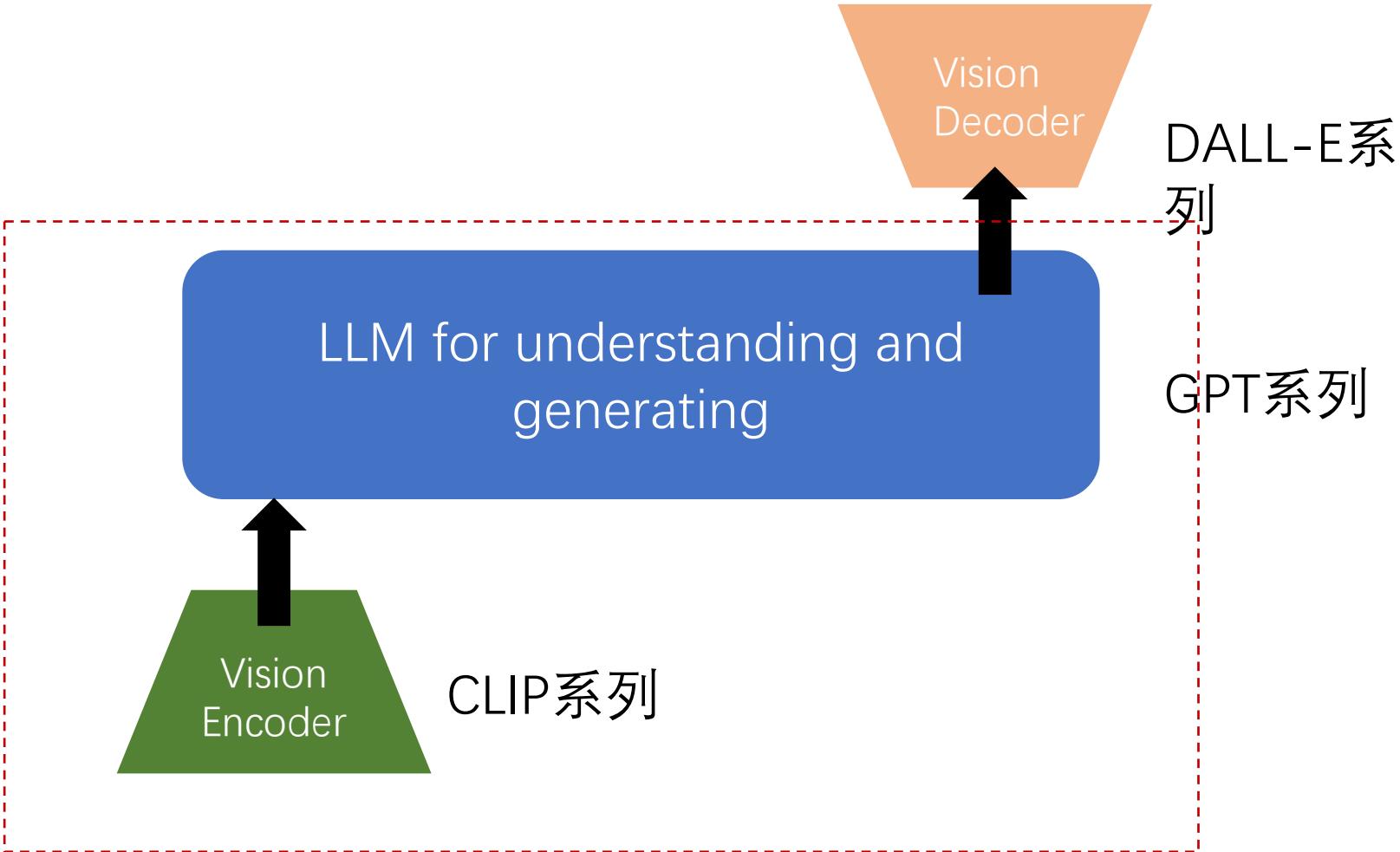


多模态和视觉大模型

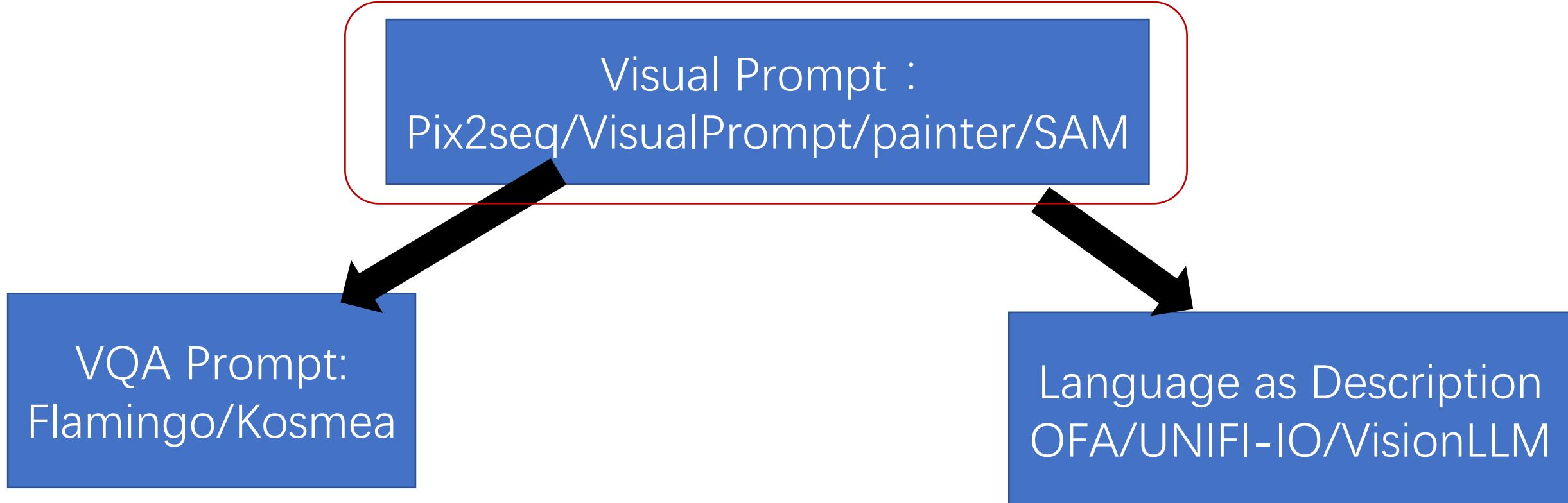
图像：

文本：

图像：

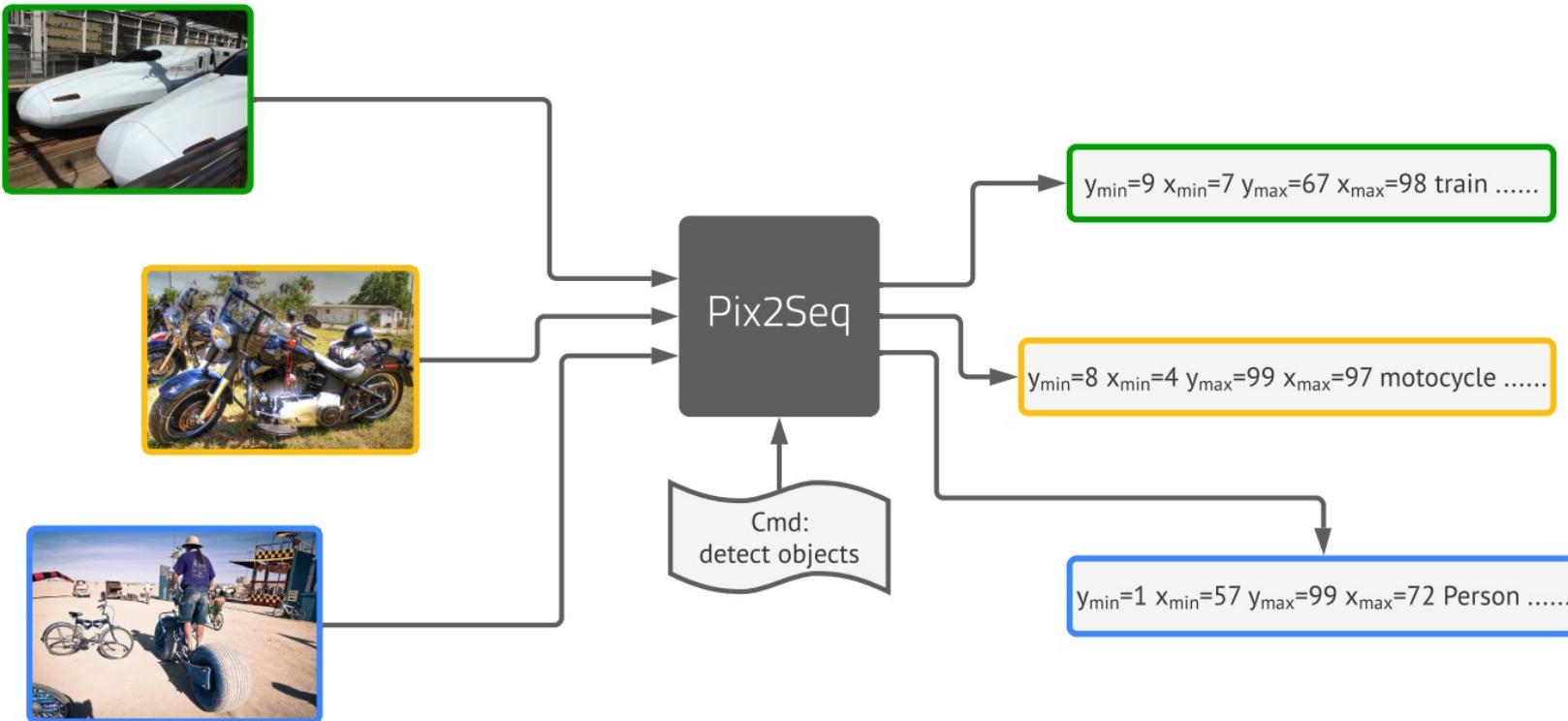


多模态和视觉大模型 (Prompt-like)



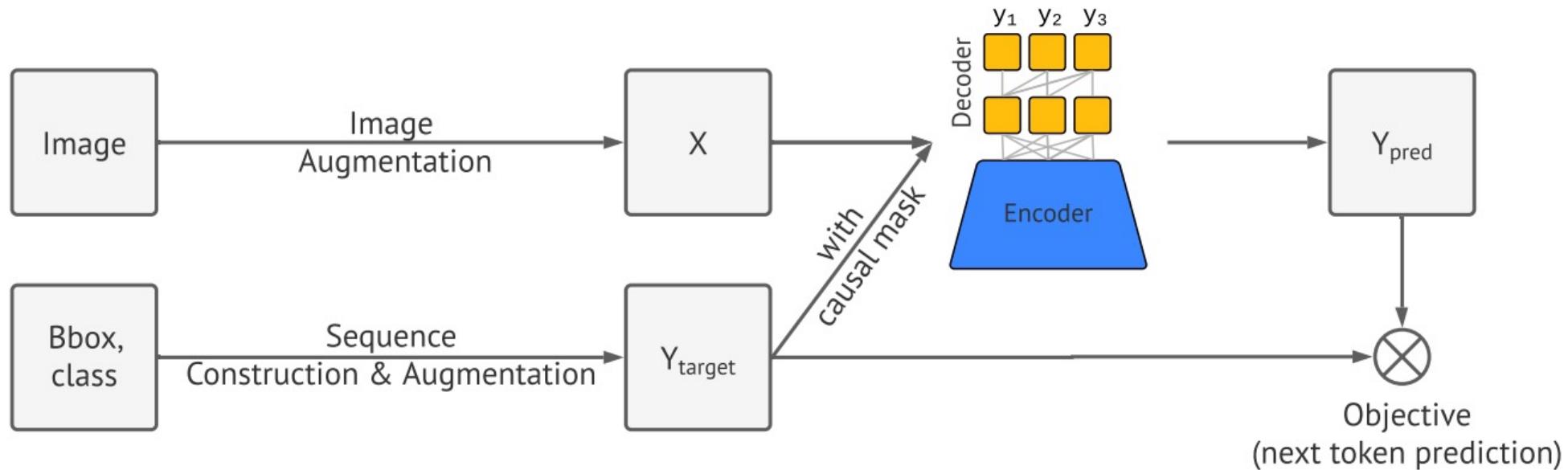
Visual Prompt

- Pix2Seq (Output as token sequences)



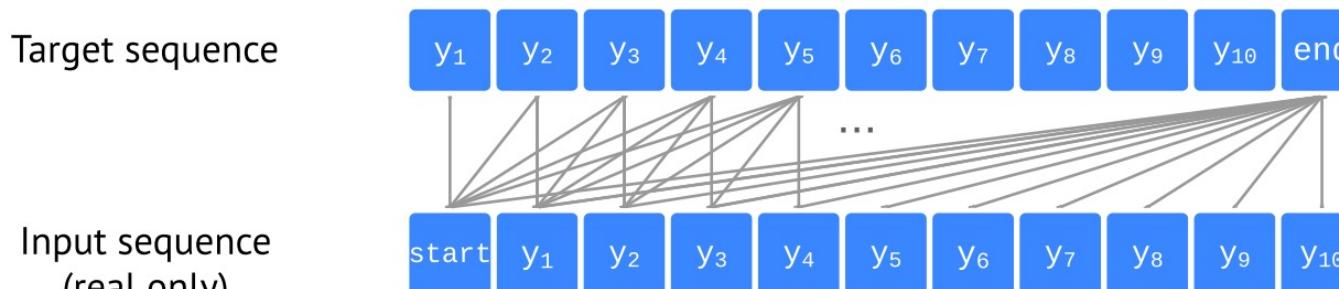
Visual Prompt

- Pix2Seq (Output as token sequences)

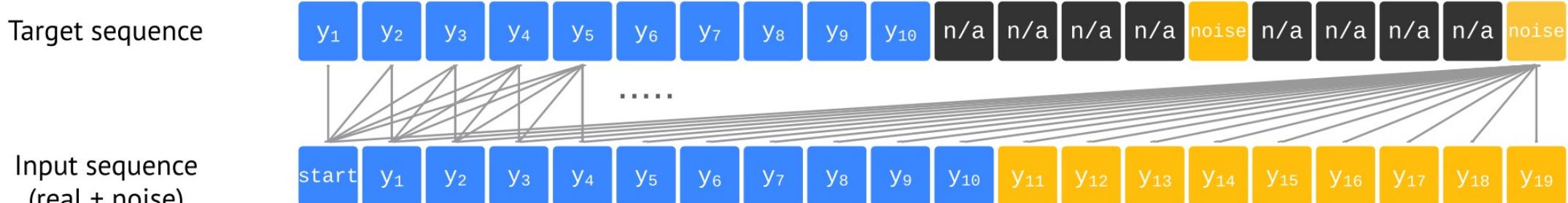


Visual Prompt

- Pix2Seq (Output as token sequences)



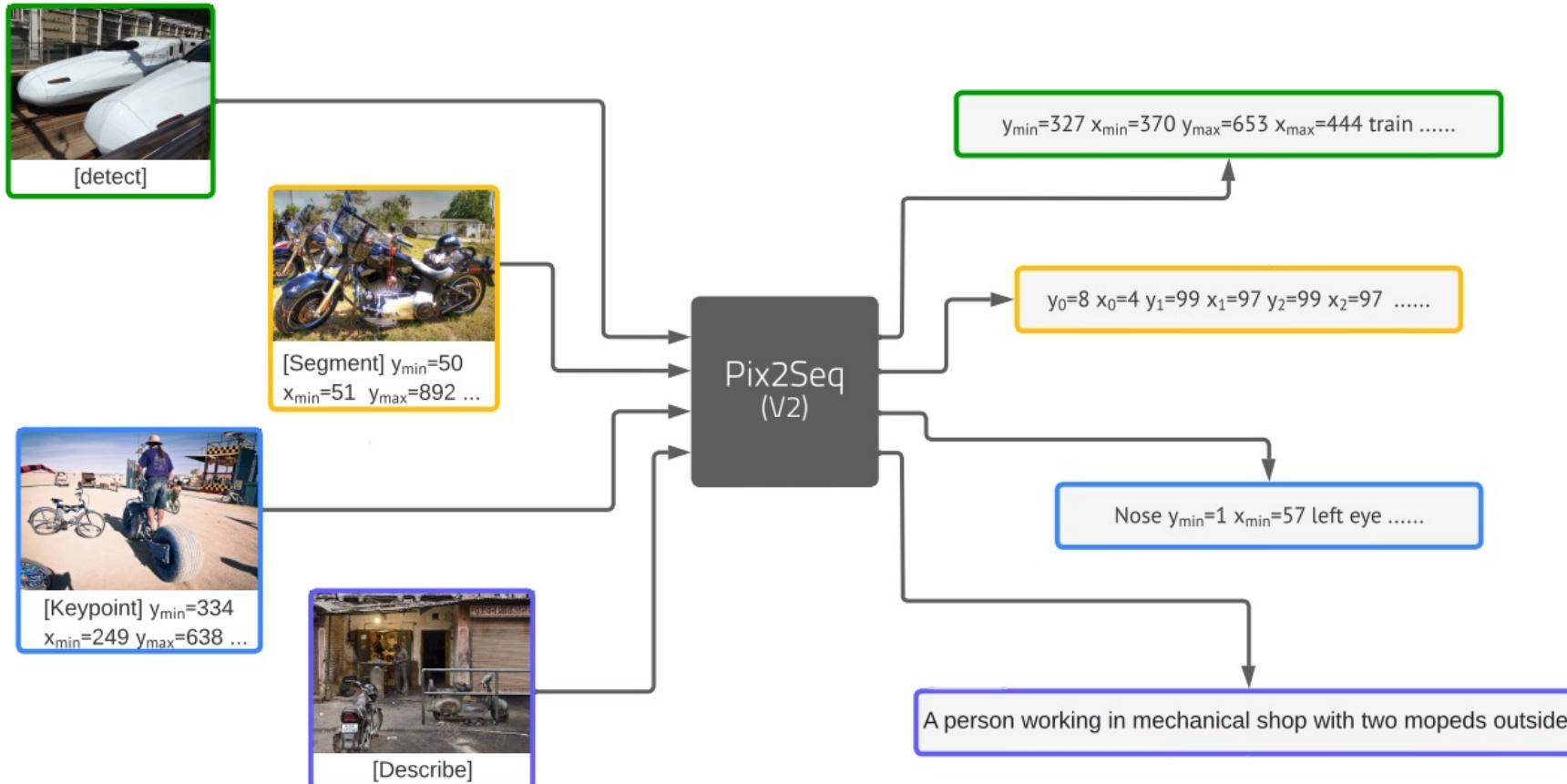
(a) Conventional autoregressive language modeling



(b) Language modeling with sequence augmentation (e.g. adding noise tokens)

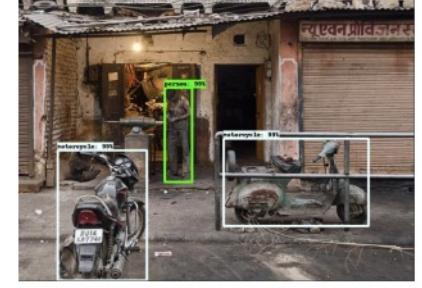
Visual Prompt

- Pix2Seq-v2 (Detection/Segmentation)



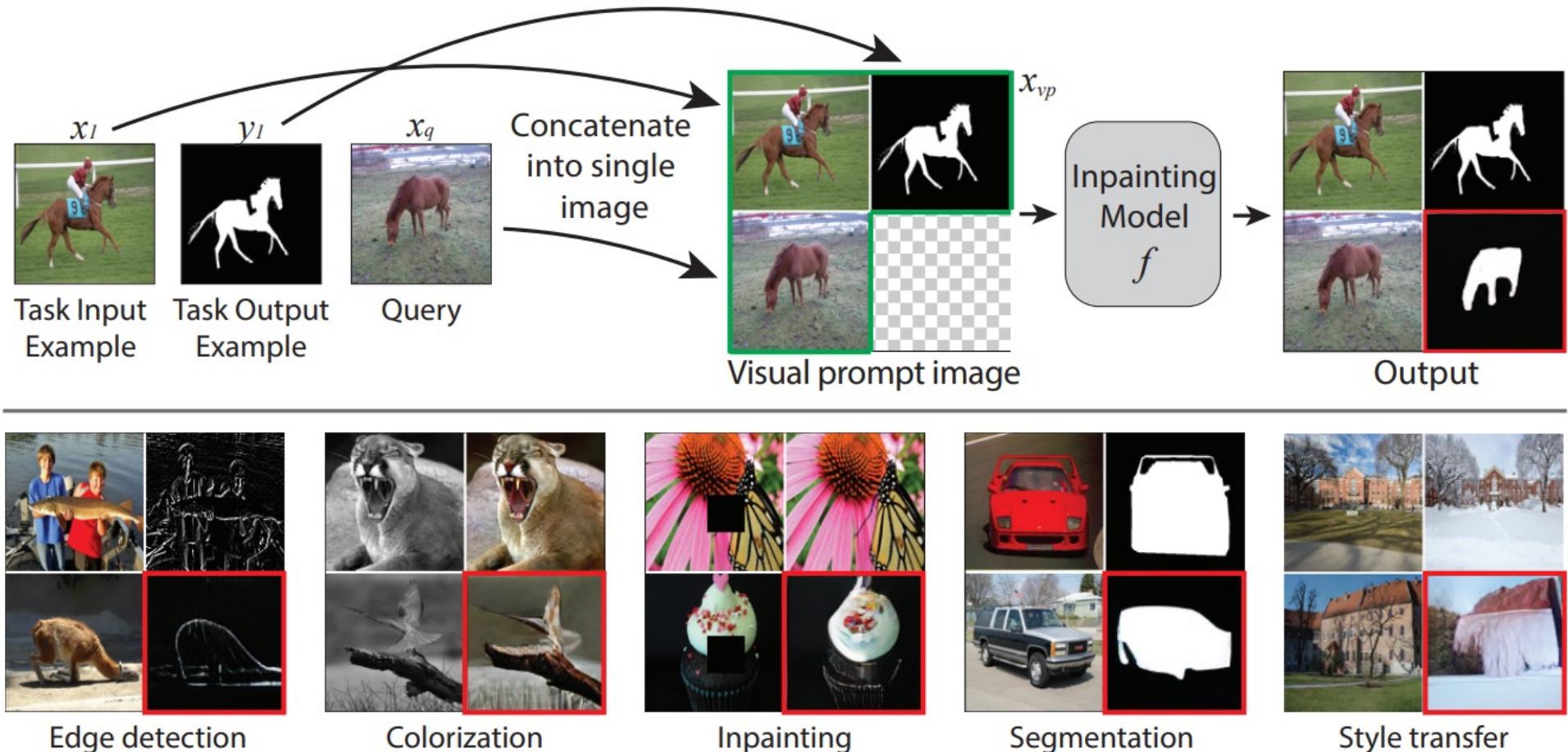
Visual Prompt

- Pix2Seq-v2

	Task prompt	Task output	Output visualization
Input image	[Detect]	$y_{min}=327 \ x_{min}=370$ $y_{max}=653 \ x_{max}=444$ person	  
	[Segment] $y_{min}=503$ $x_{min}=518 \ y_{max}=805$ $y_{max}=892$ Motocycle	$y_0=553 \ x_0=599$ $y_1=788 \ y_1=664$	
	[Keypoint] $y_{min}=327$ $x_{min}=370 \ y_{max}=653$ $x_{max}=444$ person	Nose $y_{min}=1 \ x_{min}=57$ left eye	
	[Describe]	A person working in mechanical shop with two mopeds outside.	

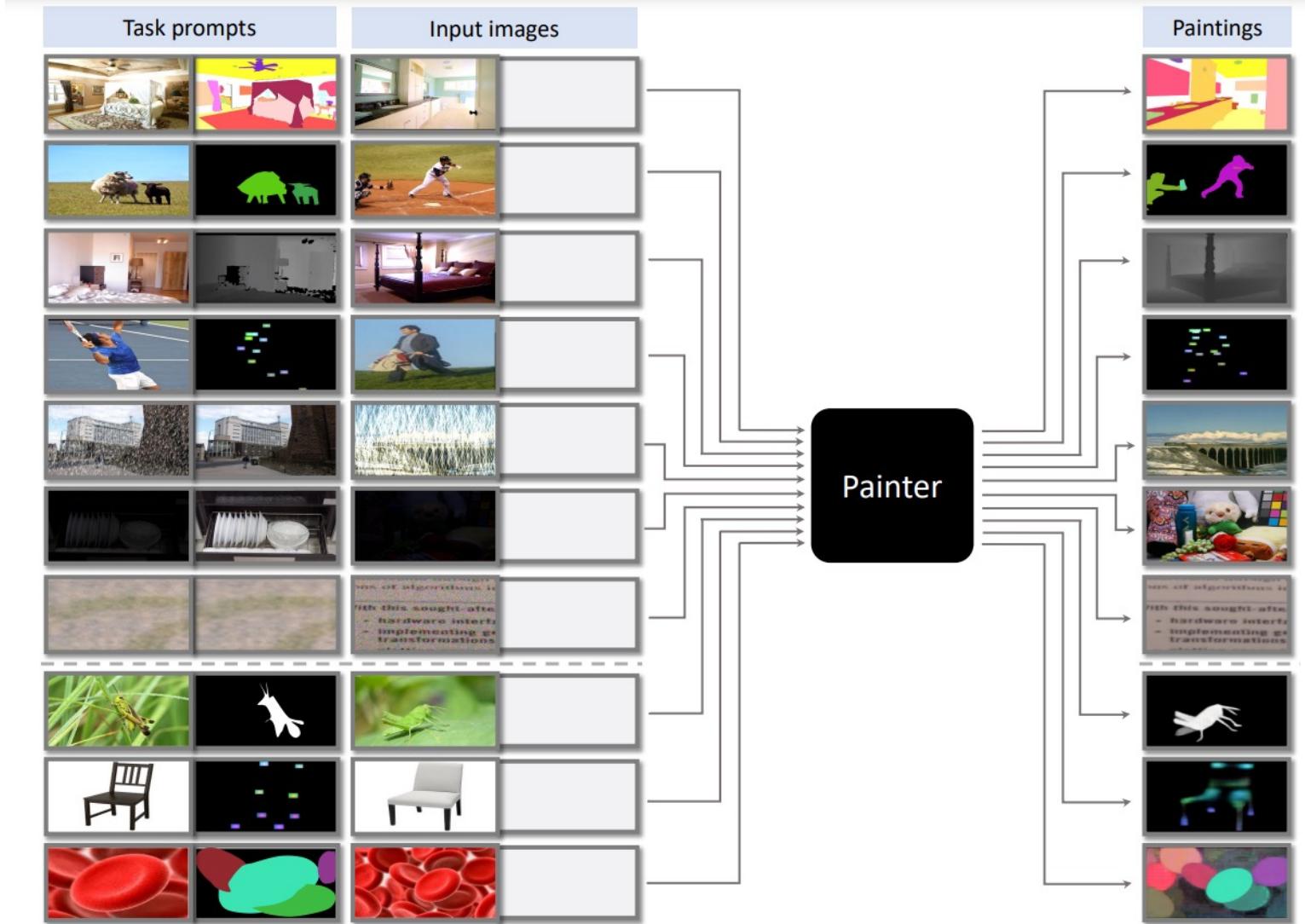
Visual Prompt

- Vision Promter



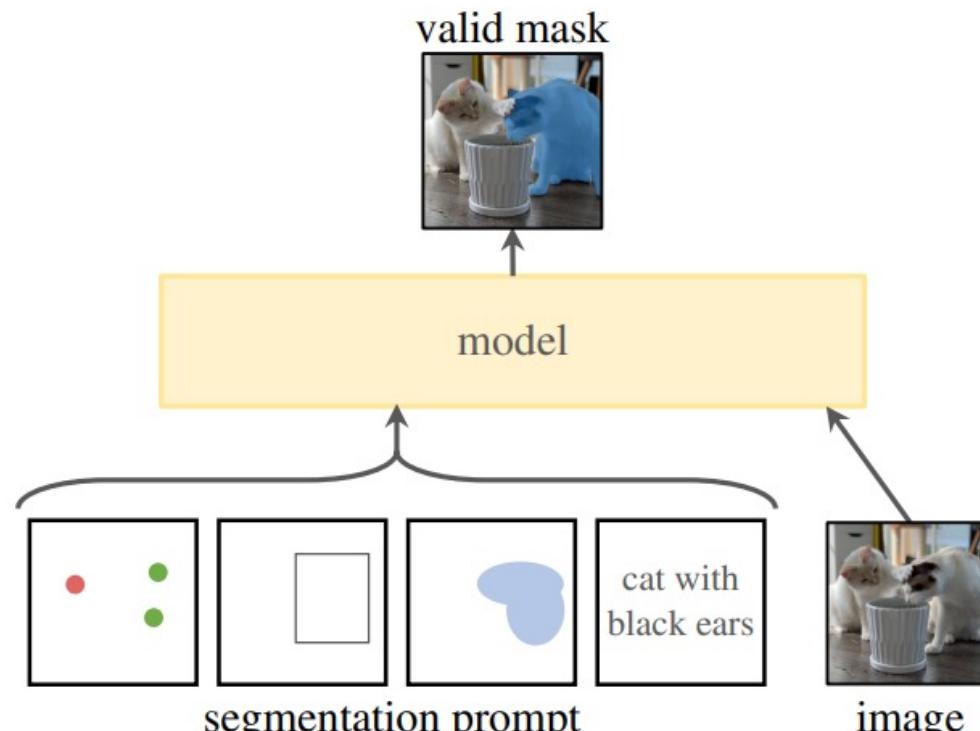
Visual Prompt

- Painter

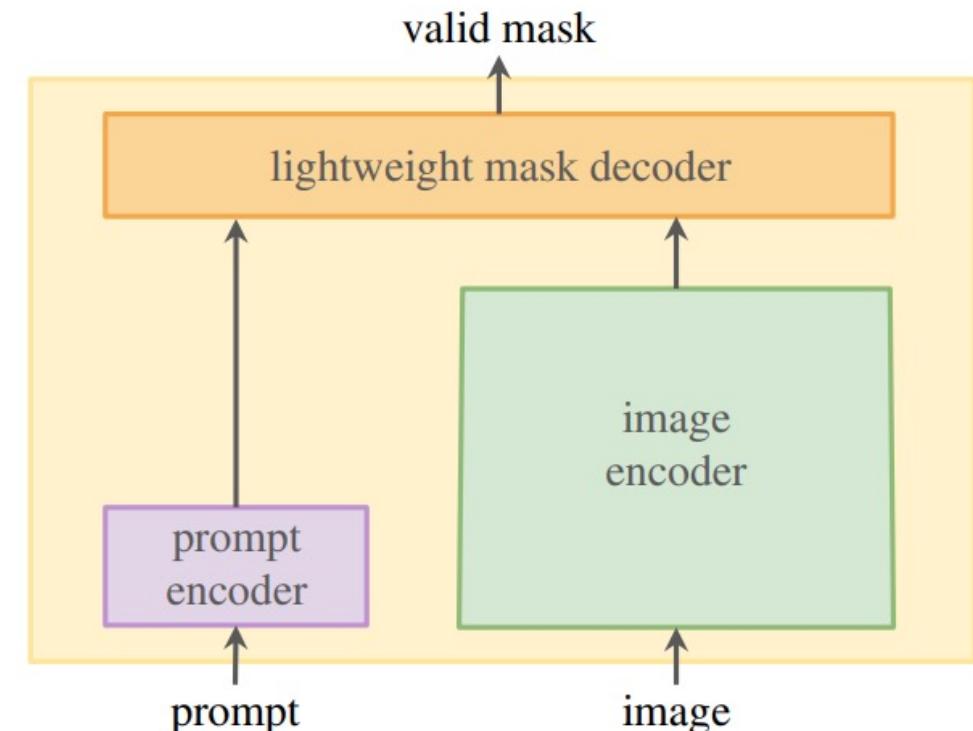


Visual Prompt

- SAM (Segment Anything Model)



(a) **Task:** promptable segmentation



(b) **Model:** Segment Anything Model (**SAM**)

多模态大模型 (Prompt-like)

Visual Prompt :
Pix2seq/VissualPrompt/painter/SA
M

VQA Prompt:
Flamingo/Kosme
a

Language as Description
OFA/UNIFI-
IO/VisionLLM

VQA Prompt

- Flamingo

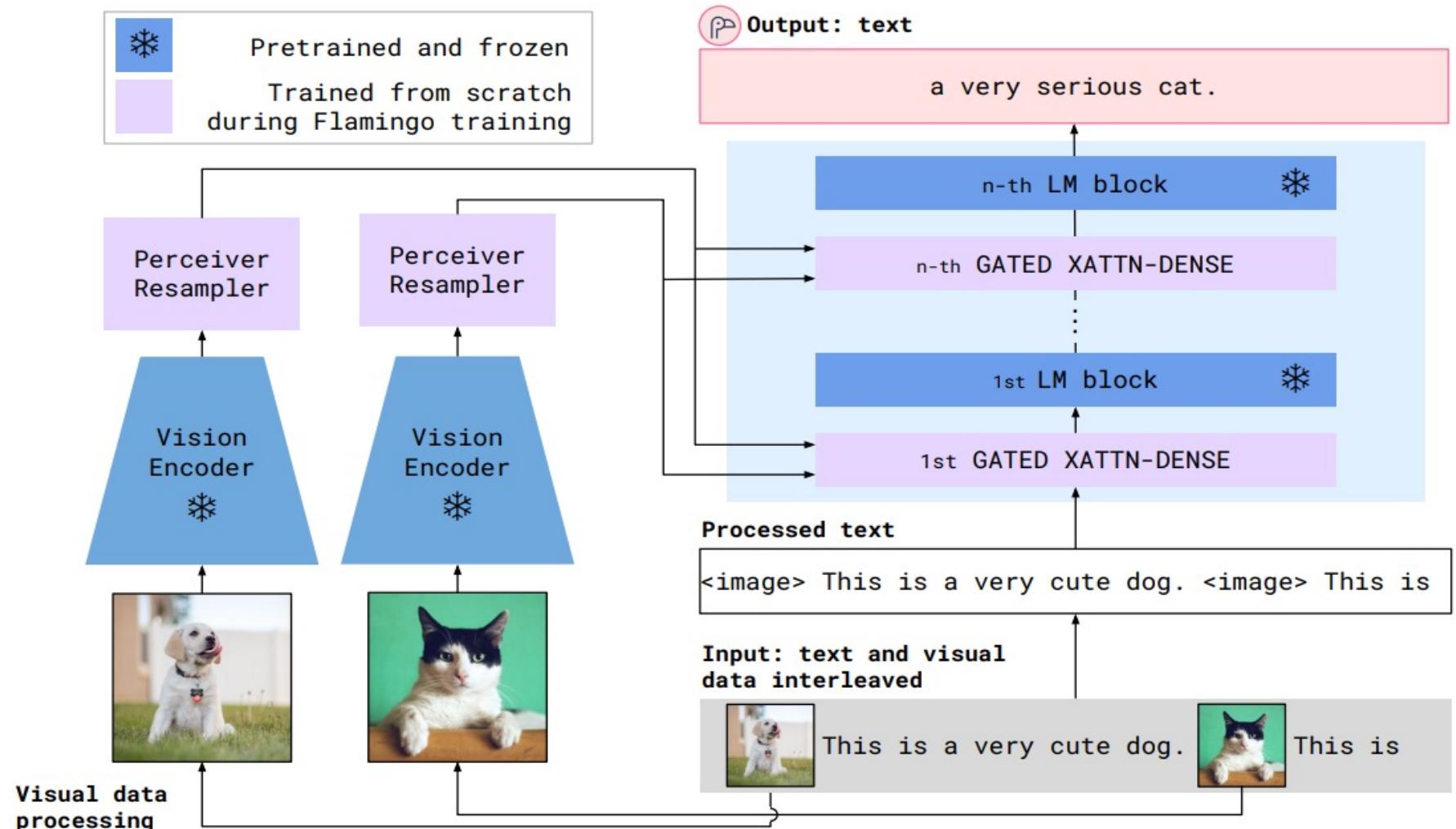
Few-

Input **shot**



VQA Prompt

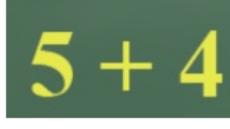
- Flamingo



VQA Prompt

- KOSMOS-1

Zero-shot

<p>Input Prompt</p> 	<p>Question: Explain why this photo is funny? Answer:</p> <p>The cat is wearing a mask that gives the cat a smile.</p> <p>(1)</p>		<p>Question: Why did the little boy cry? Answer:</p> <p>Because his scooter broke.</p> <p>(2)</p>		<p>Question: What is the hairstyle of the blond called? Answer:</p> <p>pony tail</p> <p>(3)</p>		<p>Question: When will the movie be released? Answer:</p> <p>On June 27</p> <p>(4)</p>
<p>Input Prompt</p> 	<p>What is TorchScale?</p> <p>A library that allows transformers to work efficiently and effectively at any scale.</p> <p>(5)</p>		<p>Question: The result is? Answer:</p> <p>$5 + 4 = 9$</p> <p>(6)</p>		<p>Question: What is the heart rate in this picture? Answer:</p> <p>57 bpm</p> <p>(7)</p>		<p>The time now is</p> <p>10:10 on a large clock</p> <p>(8)</p>

VQA Prompt

- LLaVA系列
 - LLaVA, TinyLLaVA等

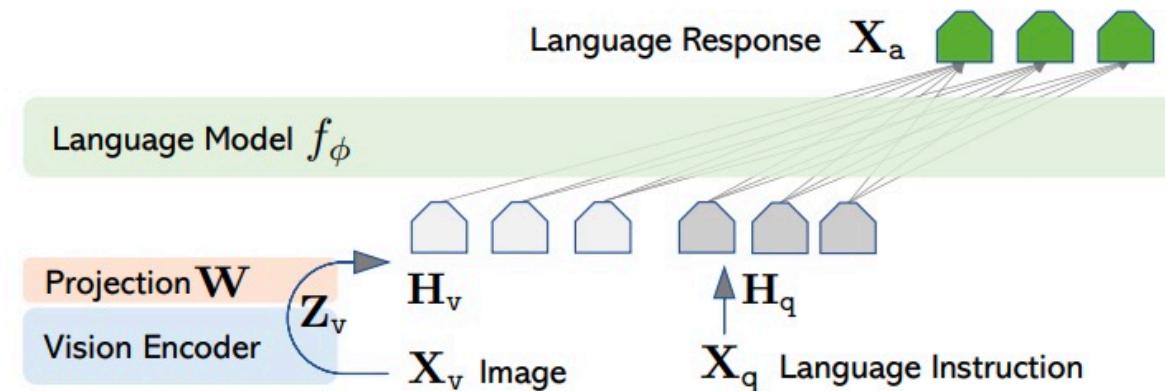
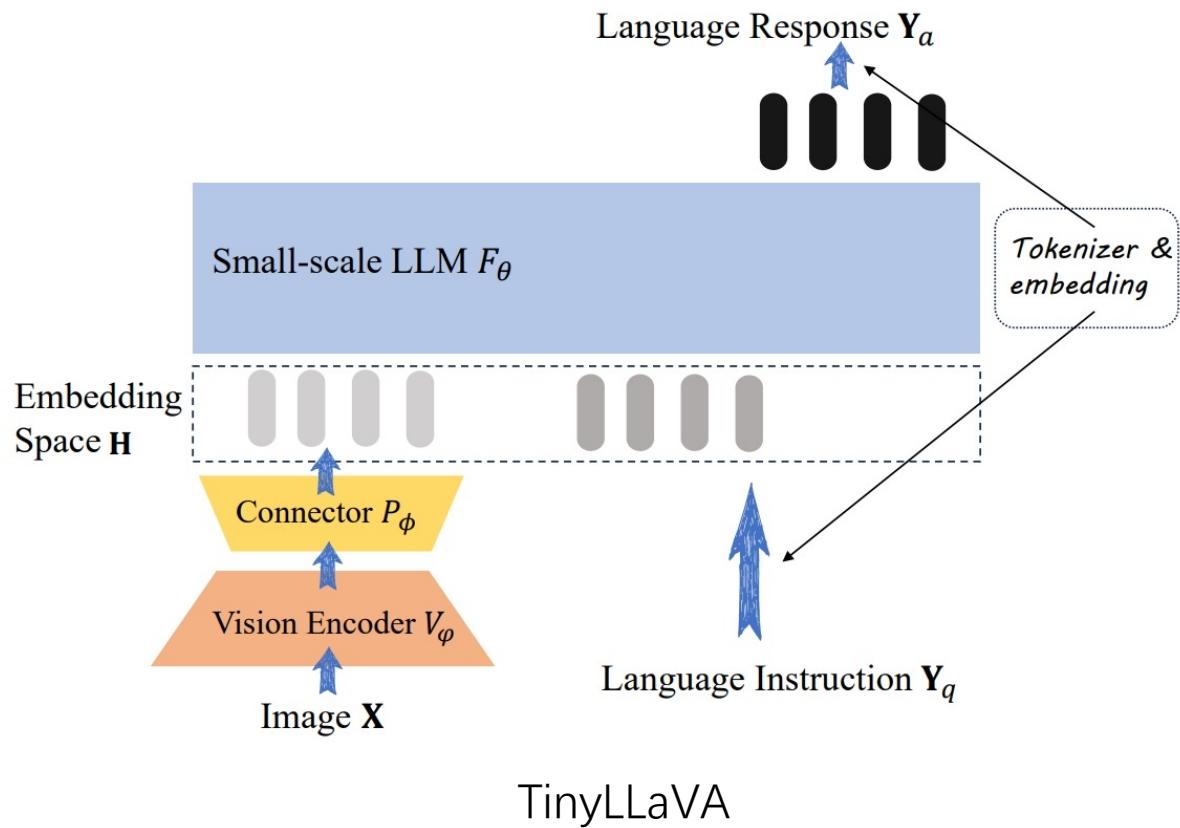


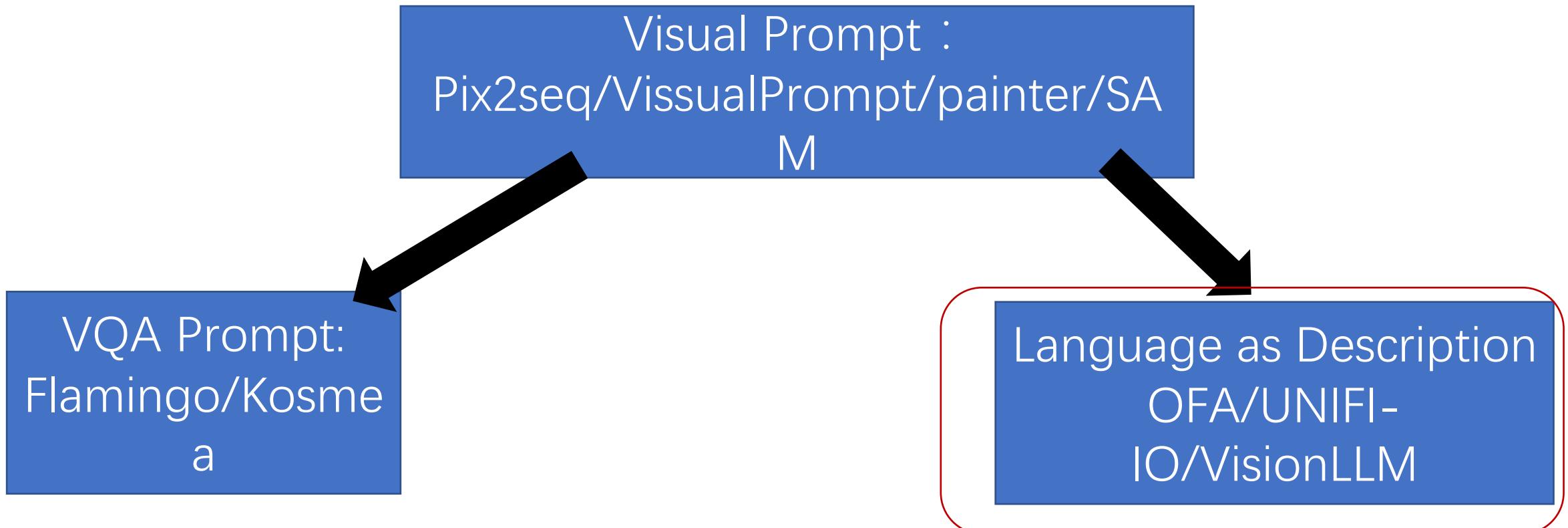
Figure 1: LLaVA network architecture.

LLaVA



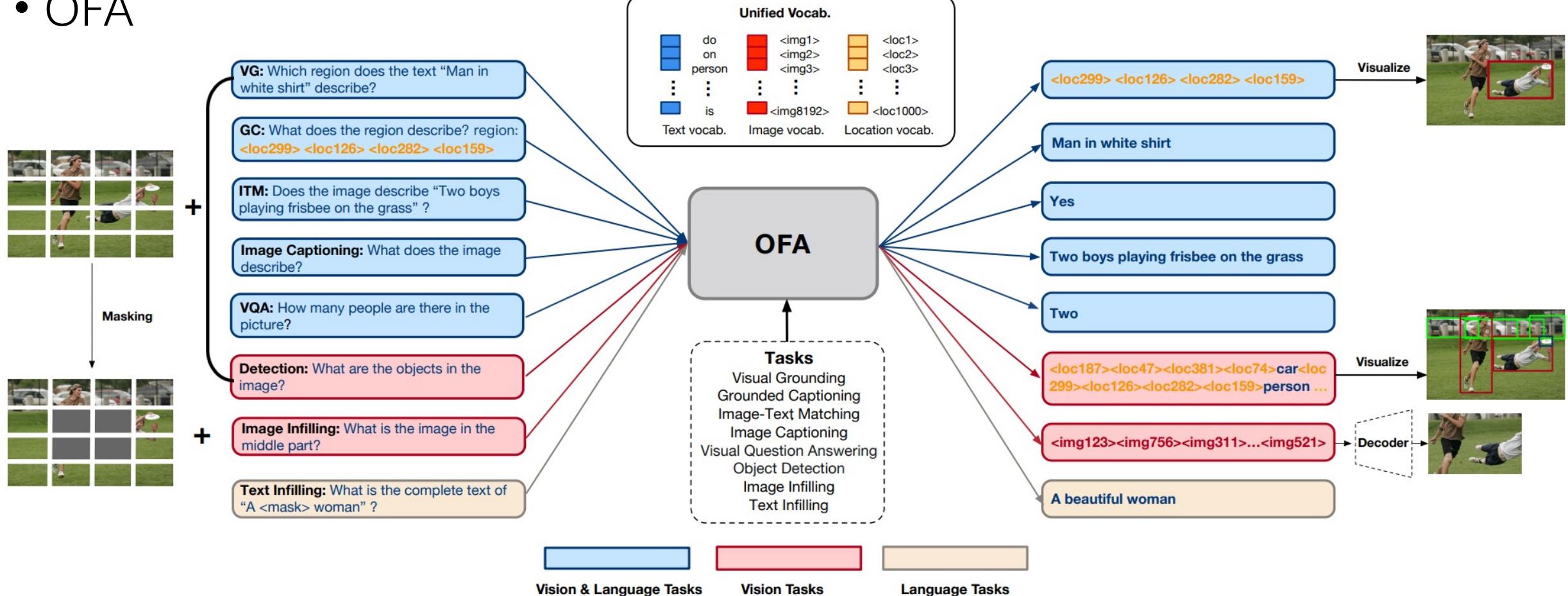
TinyLLaVA

多模态大模型 (Prompt-like)



Language as Description

- OFA



OFA: UNIFYING ARCHITECTURES, TASKS, AND MODALITIES THROUGH A SIMPLE SEQUENCE-TO-SEQUENCE LEARNING FRAMEWORK. ICML 2022

Language as Description

- Unified-IO

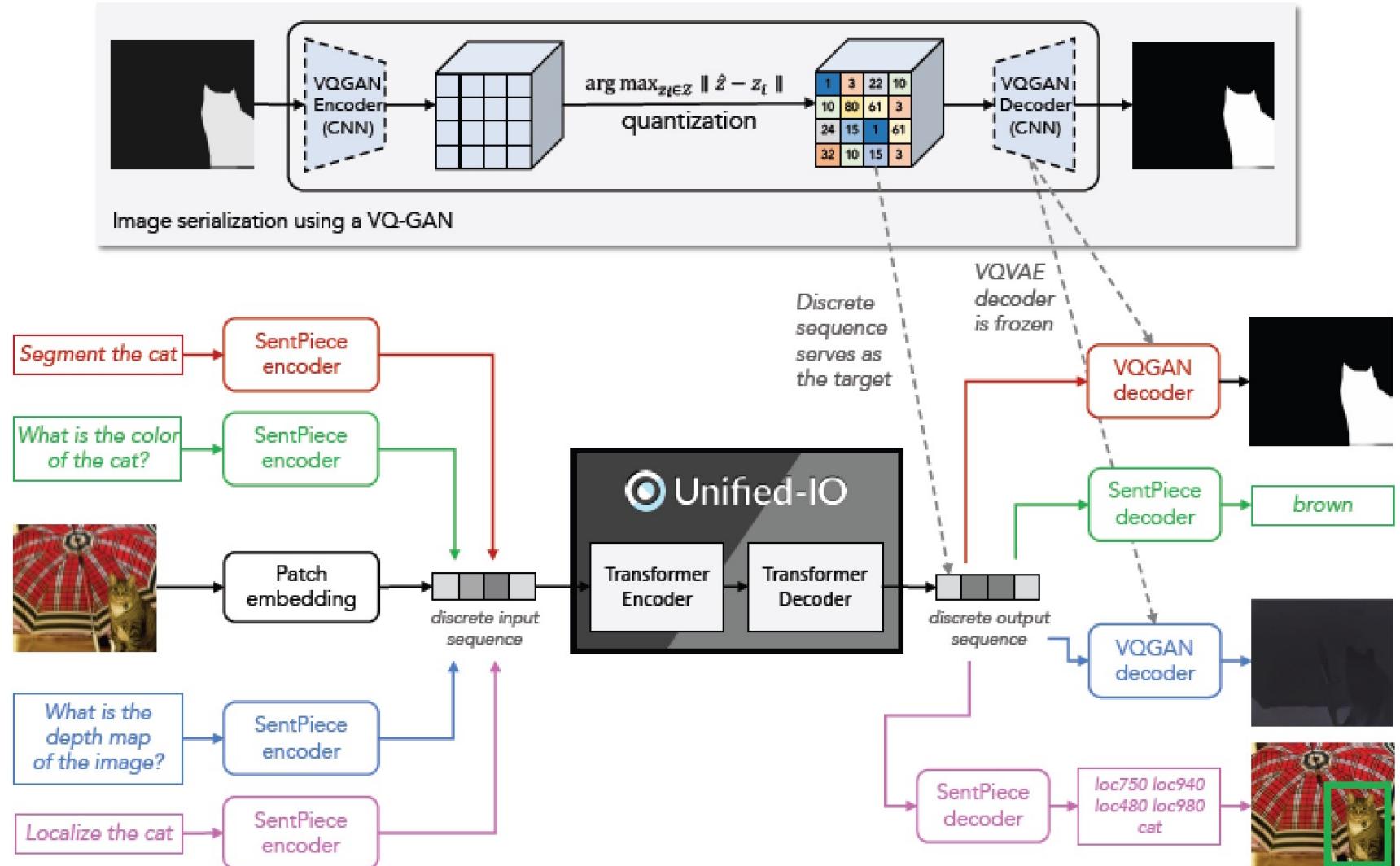


Tasks

Image Classification
Object Detection
Semantic Segmentation
Depth Estimation
Surface Normal Estimation
Segment-based Image Generation
Image Inpainting
Pose Estimation
Relationship Detection
Image Captioning
Visual QA
Referring Expressions
Situation Recognition
Text-based Image Generation
Visual Commonsense
Classification in context
Region Captioning
GLUE Benchmark tasks
Reading comprehension
Natural Language Inference
Grounded Commonsense Inference

Language as Description

- Unified-IO



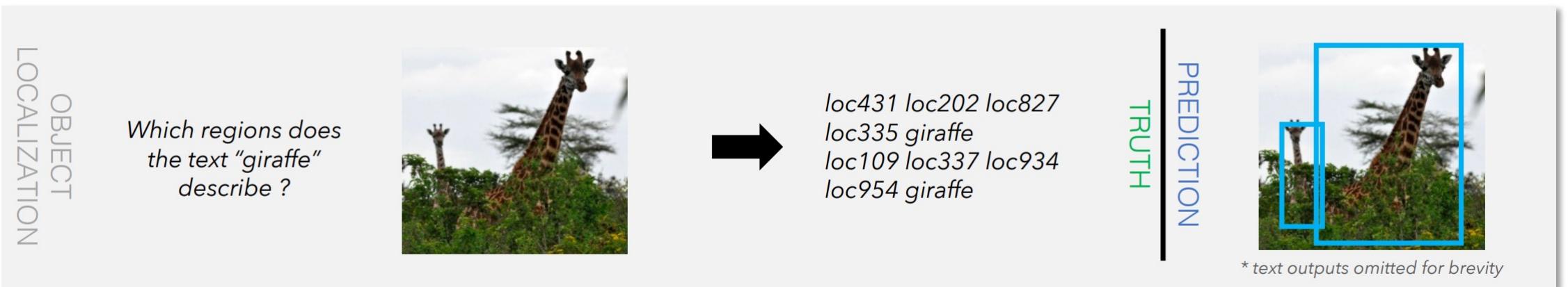
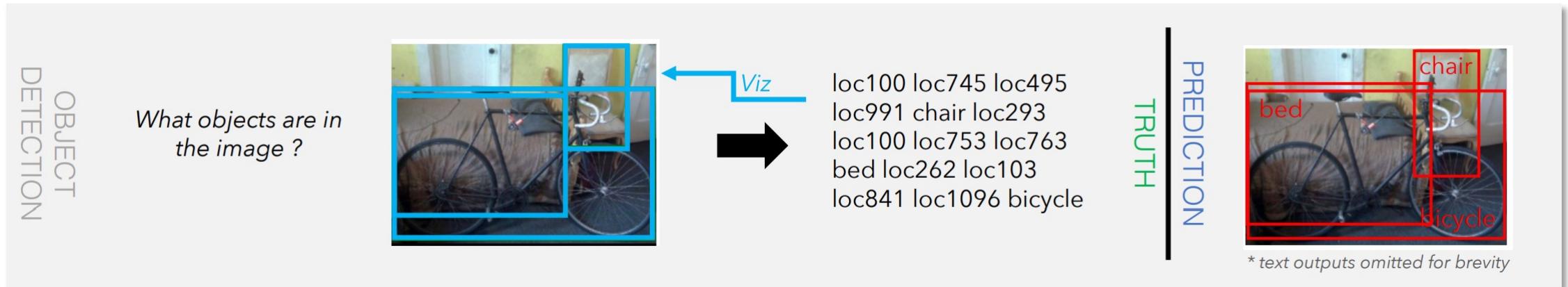
Language as Description

- Unified-IO: Image synthesis



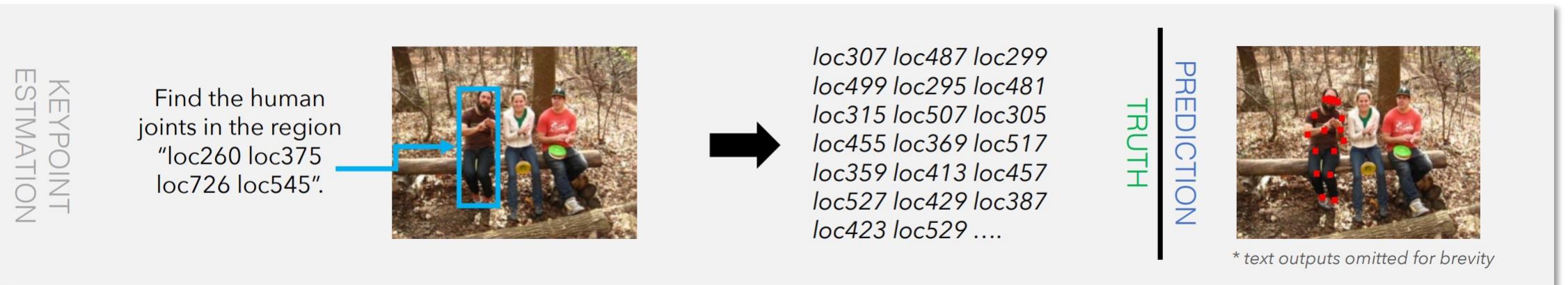
Language as Description

- Unified-IO: Sparse label



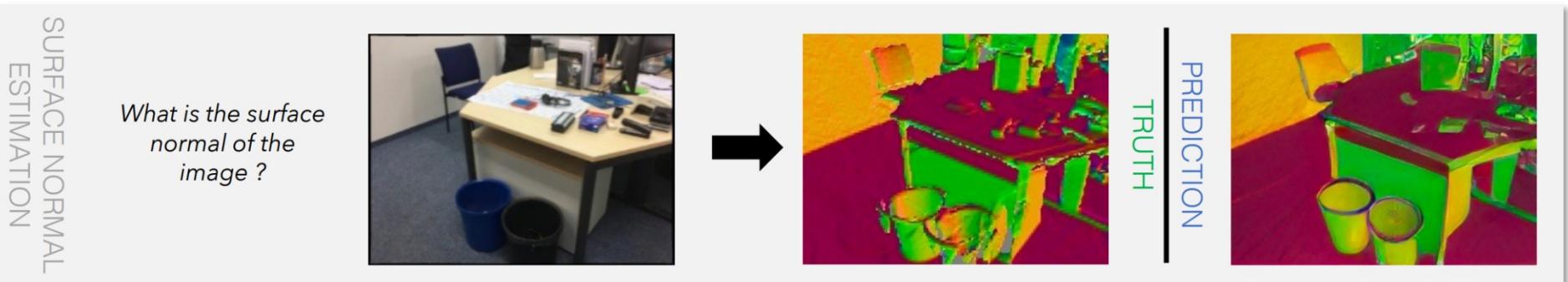
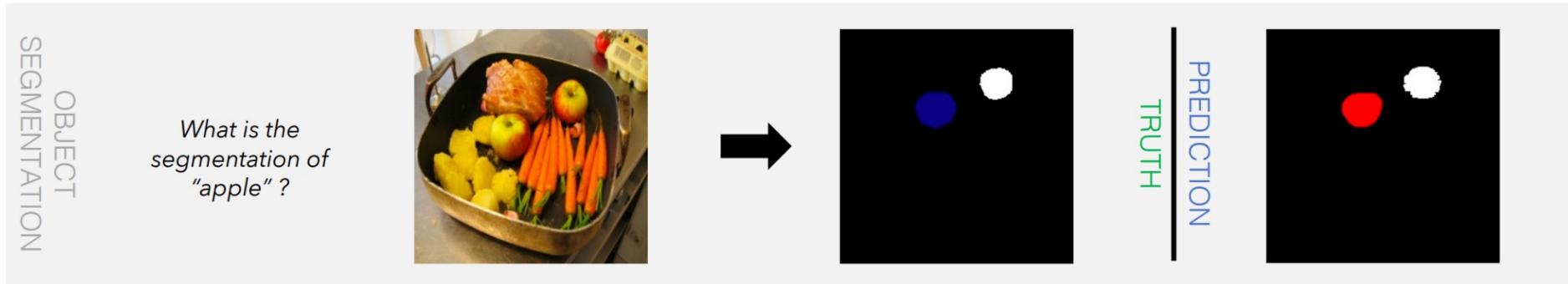
Language as Description

- Unified-IO: Sparse label



Language as Description

- Unified-IO
 - Dense Label



Language as Description

- Unified-IO: Image Classification



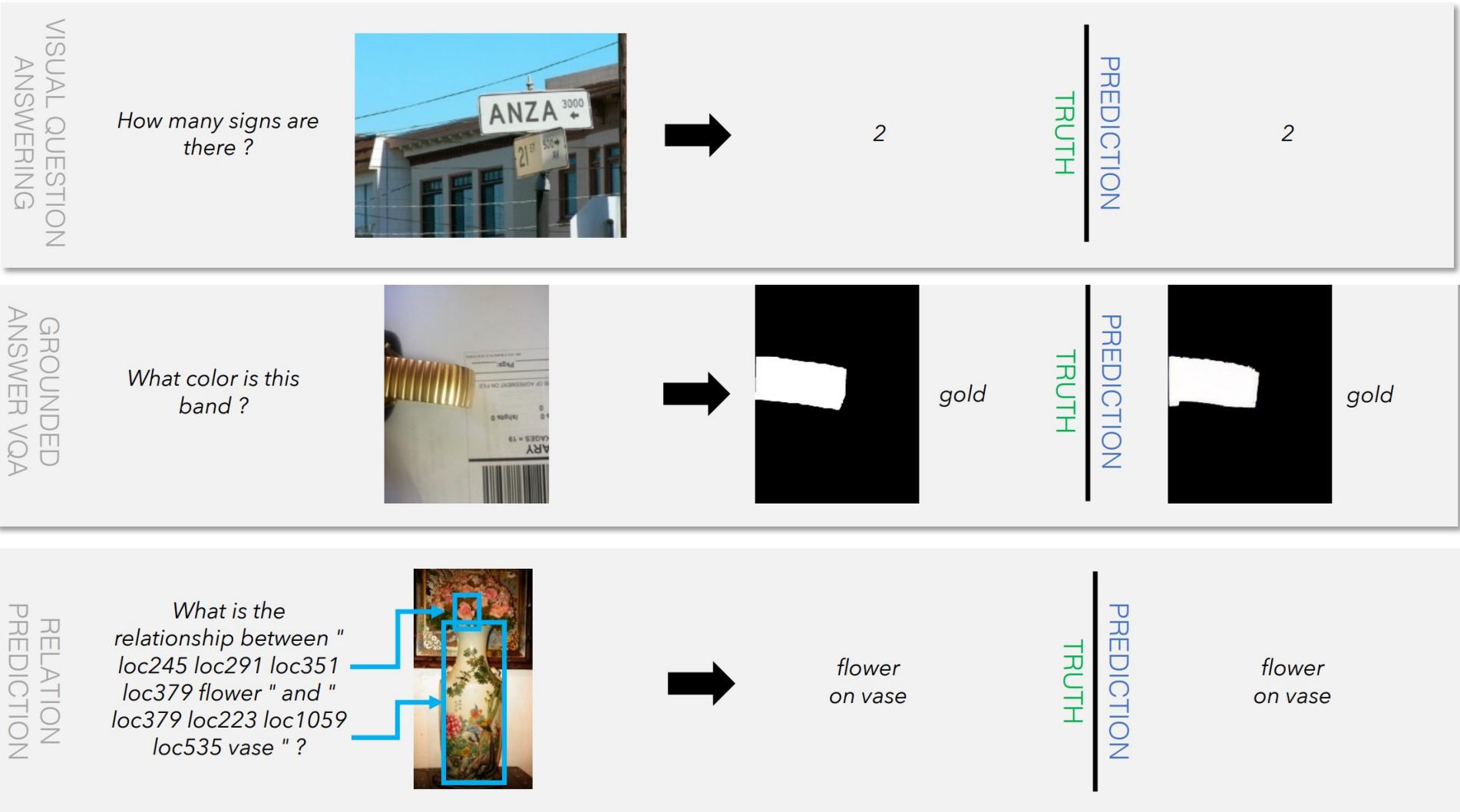
Language as Description

- Unified-IO: Image Captioning Tasks



Language as Description

- Unified-IO
 - V&L Tasks



Language as Description

- Unified-IO : NLP Tasks

TEXT
CLASSIFICATION

context: Swansea striker Lee Trundle has negotiated a lucrative image-rights deal with the League One club. Lee Trundle is in business with the League One club.
question: Does this sentence entail the following sentence?



Yes it entails

PREDICTION
TRUTH

Yes it entails

QUESTION
ANSWERING

context: Uptake of O₂ from the air is the essential purpose of respiration, so oxygen supplementation is used in medicine. Treatment not only increases oxygen levels in the patient's blood....
question: What medical treatment is used to increase oxygen uptake in a patient?



oxygen
supplementation

PREDICTION
TRUTH

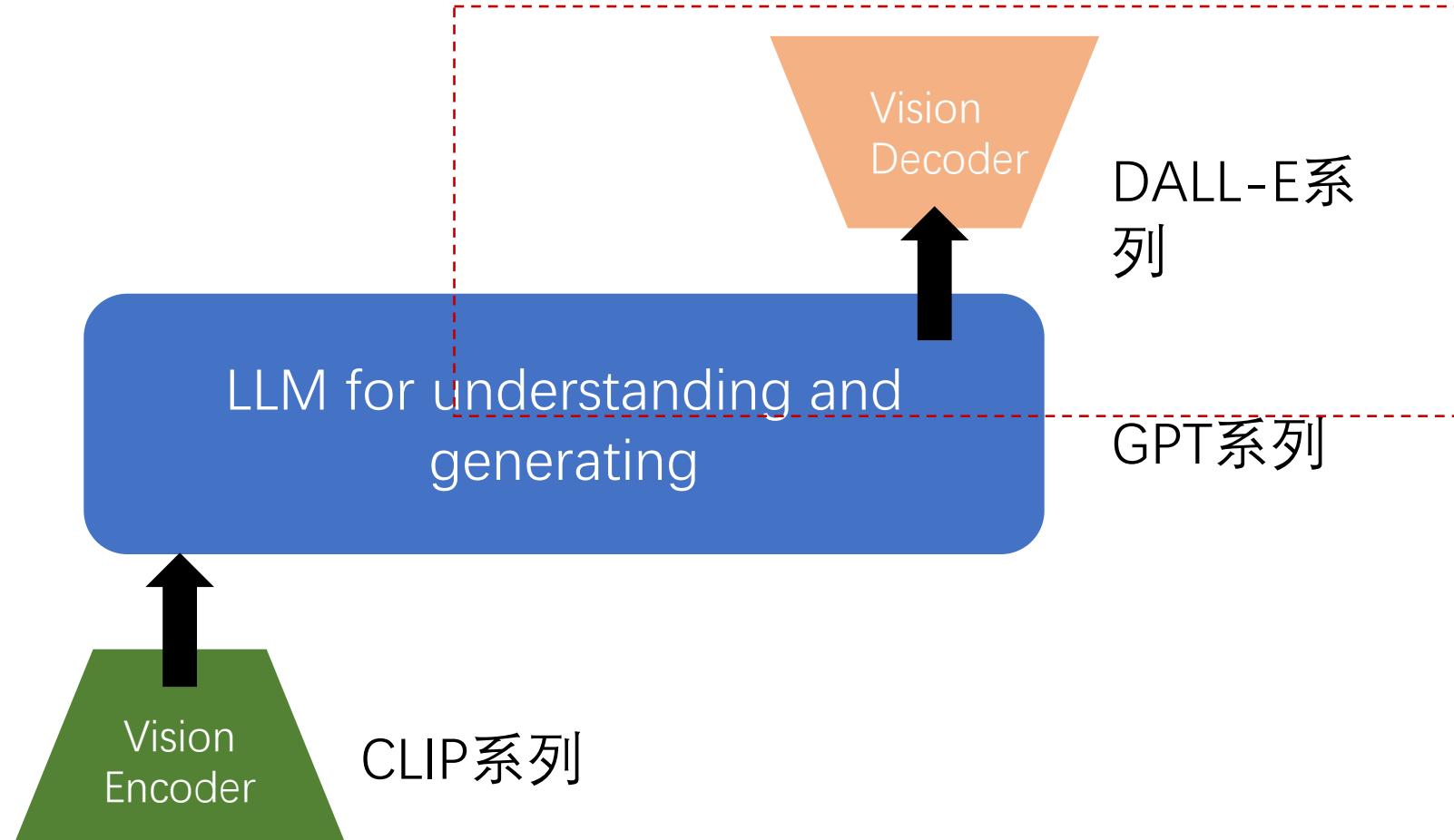
oxygen
supplementation

多模态和视觉大模型

图像：

文本：

图像：



谢谢！

