Introduction
OO

Methodology
OO
OOOO

Experiment
O
OOO

Related Works
OOOOOOOO

Conclusion
O

# SAM a great Semantic Segmentations
## LLM to generate the reward function

Hantao Zhou

Universitaet Hamburg

January 21, 2024

Introduction
OO

Methodology
OO
OOOO

Experiment
O
OOO

Related Works
OOOOOOOO

Conclusion
O

# Table of Contents

# Multi-modal processing structure



Figure: Multi-modal processing structure

▶ CLIP

Introduction
○●
○○

Methodology
○○
○○○○

Experiment
○○○

Related Works
○○○○○○○○

Conclusion
○

## Prompt-based Techniques



**(A) Pretrain–finetune (BERT, T5)**

Pretrained LM → Finetune on task A → Inference on task A

- Typically requires many task-specific examples
- One specialized model for each task

**(B) Prompting (GPT-3)**

Pretrained LM → Inference on task A

Improve performance via few-shot prompting or prompt engineering

**(C) Instruction tuning (FLAN)**

Pretrained LM → Instruction-tune on many tasks: B, C, D, ... → Inference on task A

Model learns to perform many tasks via natural language instructions

Inference on unseen task

1. Instruct Tunning
2. Prompting

## Motivation

1. to build a good big-modal based image model
2. to harness the capability of zero-shot

Introduction
○○
○●

Methodology
○○
○○○○

Experiment
○
○○○

Related Works
○○○○○○○○

Conclusion
○

## Scientific Questions

1. What task will enable zero-shot generalization?
2. What is the corresponding model architecture?
3. What data can power this task and model?

Introduction
00
00

Methodology
●○
○○○○

Experiment
○
○○○

Related Works
○○○○○○○○

Conclusion
○

## General Methods



(a) **Task**: promptable segmentation

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

(c) **Data**: data engine (top) & dataset (bottom)
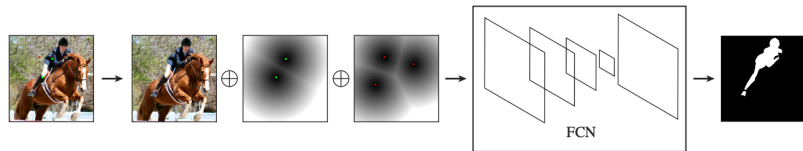
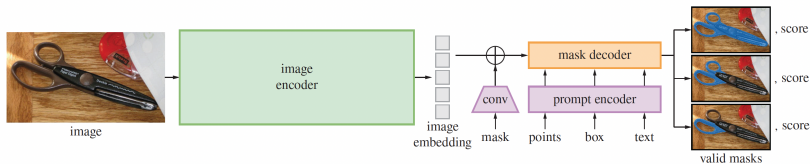1. Promptable Segmentations
2. Encoder-Decoder Architecture
3. Data Engine with Dataset

# Task



1. Translating the idea of Prompting to the task of semantic segmentation
2. Generate mask for any prompt
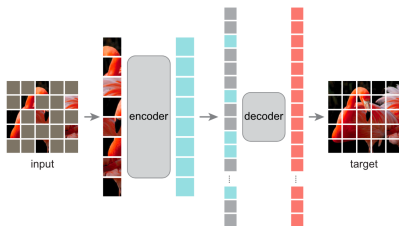3. Leads ot a natural pre-training algorithm

# Pretrain



1. Provide with positive and negative clicks
2. Present the answer of correct mask
3. Unlike the classic interactive semantic segmentation, the annotator can provide the mask for any prompt
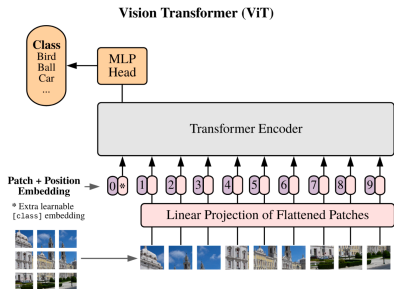
# Model Architecture



1. Image Encoder
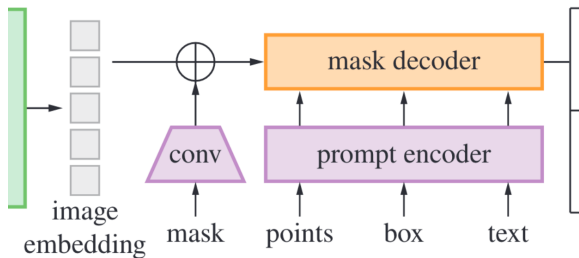2. Prompt Encoder
3. Mask Decoder
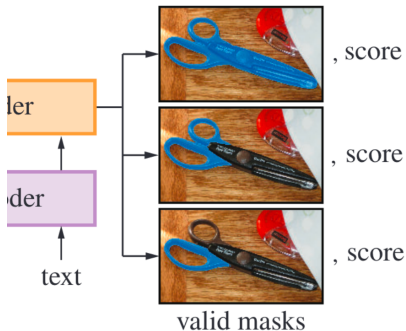4. Resolving Ambiguity

# Image Encoder



1. MAE
2. ViT

## prompt Encoder / Decoder



1. Prompt of Dense and Sparse
2. masks / points, boxes, text
3. Mask encoder map the image embedding, mask and prompts to the result mask

# Resolving Ambiguity

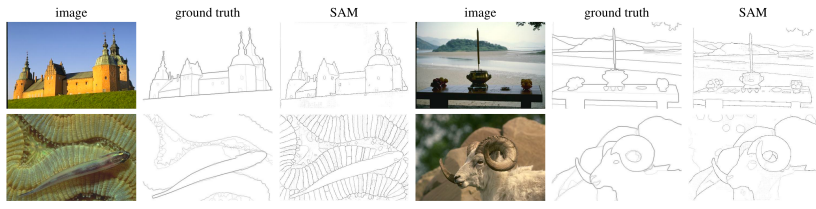

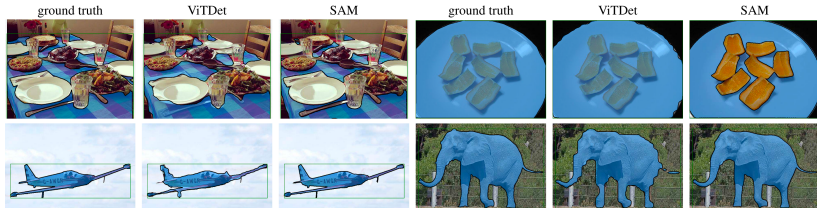1. Three mask is usually sufficient for representing
2. add estimated IoU

# Training

1. Assisted-manual stage
    1.1 like classic interactive semantic segmentation
    1.2 have mechanism for solving granularity problem
    1.3 annotations are based on the models' output
2. Semi-automatic stage
    2.1 Aims to increase the diversity of masks in order timprove the model's generalization ability
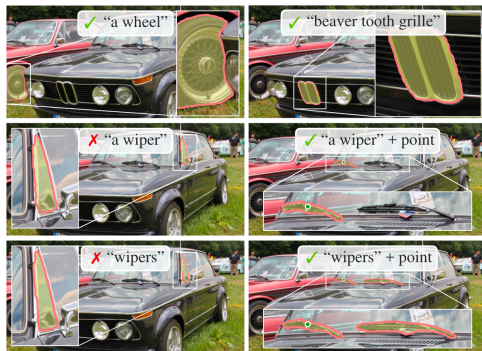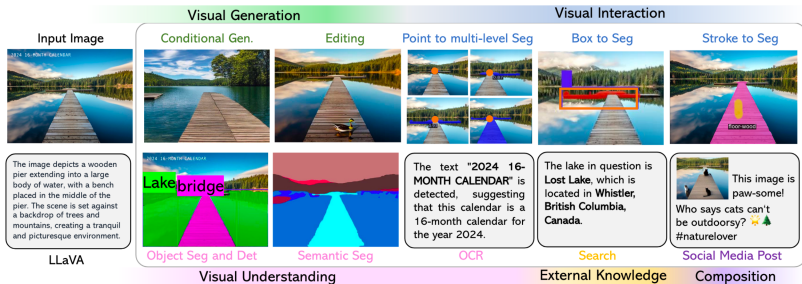    2.2 Ask the annotators to provide different masks
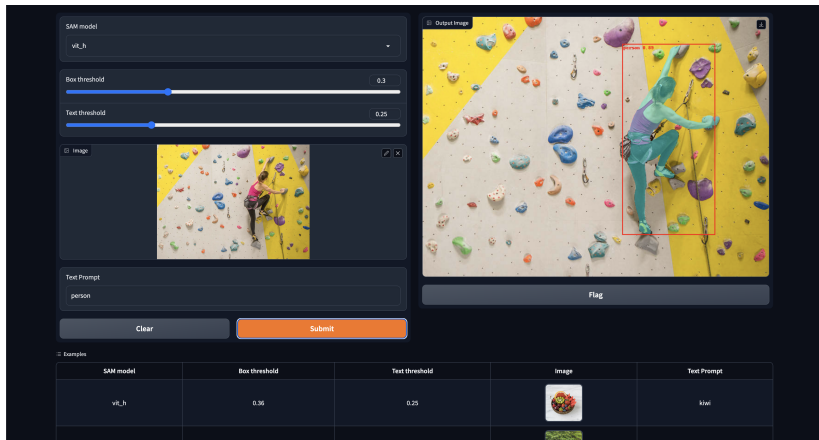3. Fully-automatic stage

Introduction
00
00

Methodology
00
00
0000

Experiment
0
●00

Related Works
00000000

Conclusion
0

# Edge Detection

## Instance Segmentation

Introduction
OO
OO

Methodology
O
OO
OOOO

Experiment
O
OO●

Related Works
OOOOOOOO

Conclusion
O

# Text-to-Mask

Introduction
OO
OO

Methodology
OO
OOOO

Experiment
O
OOO

Related Works
●OOOOOOO

Conclusion
O

# LLAVA-PLUS

# Language Segment-Anything

Introduction
○○

Methodology
○○
○○○○

Experiment
○○○

Related Works
○○●○○○○○○

Conclusion
○

# Transfer



| Method | Seen | Unseen | New background | More distractors | Average |
|--------|------|--------|----------------|------------------|---------|
| Ours | 82.5 | 80.0 | 65.0 | 75.0 | **75.625** |
| -replace mask with bbox | 50.0 | 40.0 | 25.0 | 30.0 | 36.25 |
| -w/o tracking | 70.0 | 50.0 | 55.0 | 70.0 | 61.25 |
| -single view | 65.0 | 80.0 | 20.0 | 70.0 | 58.75 |
| -RGB-M only | 85.0 | 70.0 | 50.0 | 70.0 | 68.75 |

Introduction
00
00

Methodology
00
0000

Experiment
000

Related Works
0000●000

Conclusion
0

# Grasp Anything



Fig. 2. **Dataset creation pipeline.**



Centroid
Point of contact
Grasp pose
Grasp line

# FLICAR

Introduction
00

Methodology
00
0000

Experiment
000

**Related Works**
000000●00

Conclusion
0

# Instruct2Act

Introduction
○○

Methodology
○○
○○○○

Experiment
○○○

**Related Works**
○○○○○○●○

Conclusion
○

# Agriculture Robots



Fig. 1: Overview of the robot platform architecture showing its components and relations

Introduction
○○
○○

Methodology
○○
○○○○

Experiment
○○○
○○○

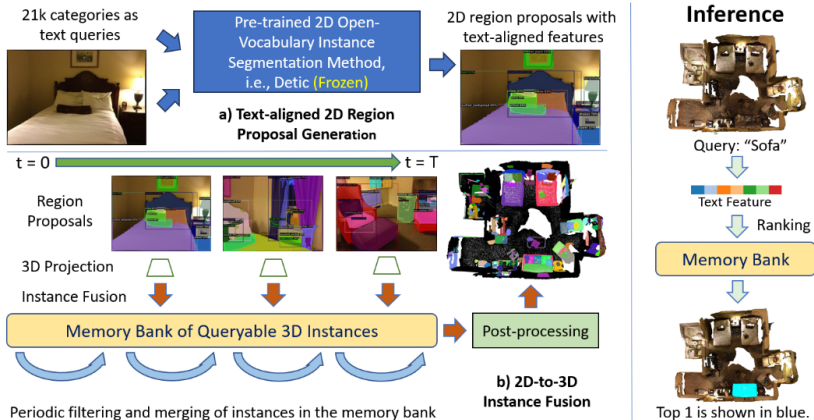Related Works
○○○○○○○●

Conclusion
○

# OVIR-3D



Figure 2: **Pipeline of the proposed method.**

## Takeaways

1. A good semantic segmentation model
2. Encoporating human interaction like Prompting can give more possbliities
3. An existing experiment pattern can achieve great result when combined with new emerging techniques