

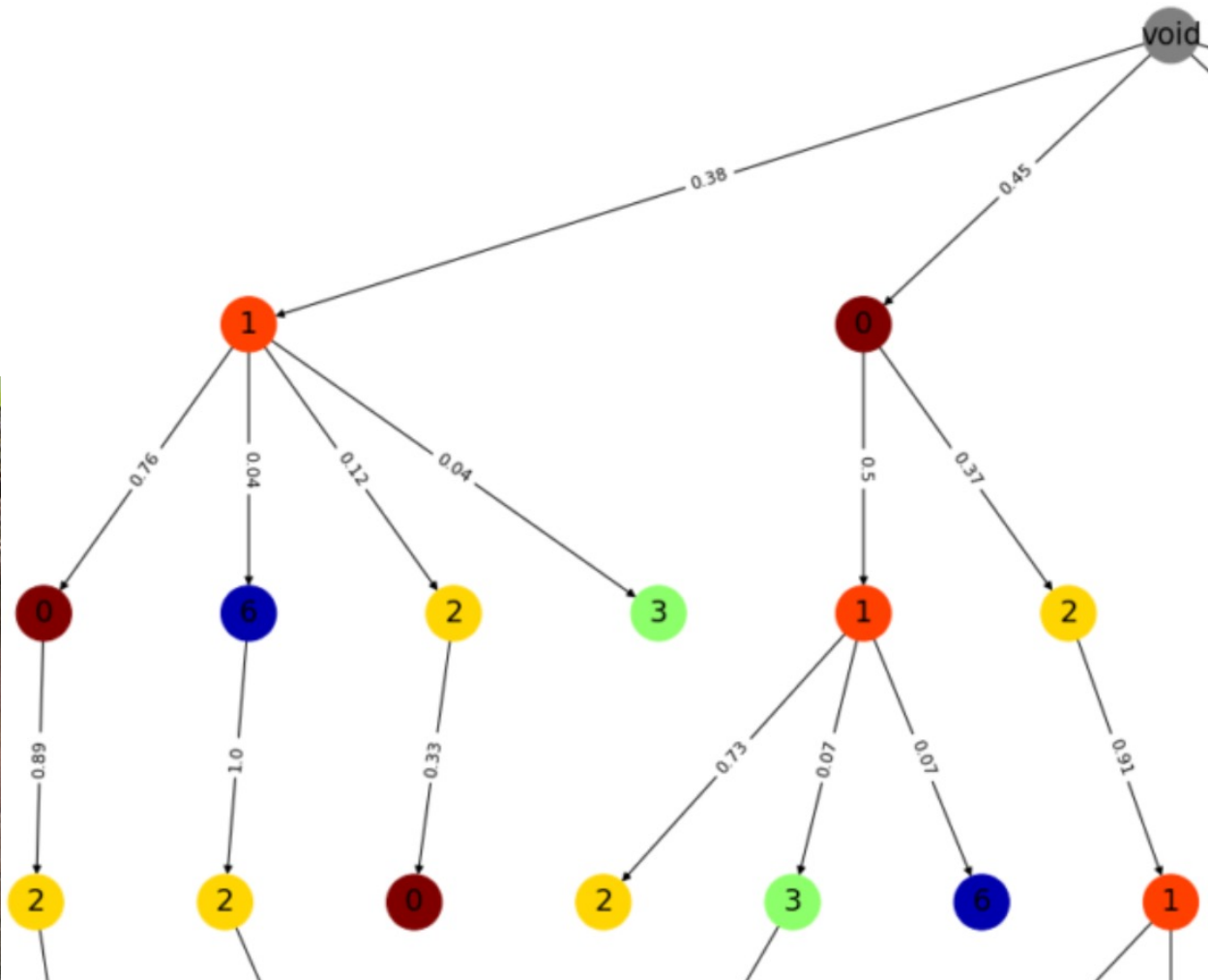
Probabilistic Salient Object Ranking

Guan Huankang

Updated abstract

Salient Object Ranking (SOR) aims to study how humans visually explore complex scenes by predicting an ordered sequence of objects that will attract our attention. Existing SOR approaches typically model this ranking deterministically, assuming a single, fixed ranking of attention. However, such deterministic SOR fails to capture the true nature of human attention. We observe that human attention shifts exhibit variability and stochasticity, showing that the next object of fixation is not a definitive choice but rather a probability distribution. Yet, existing SOR methods and evaluation metrics cannot account for this inherent randomness. To bridge this gap, we propose **ProbSOR**, a novel **Probabilistic Salient Object Ranking** model that explicitly learns the uncertainty of attention shifts by incorporating Group Relative Policy Optimization (GRPO). We leverage a Vision-Language Model (VLM) as the foundation for ProbSOR and extend it with autoregressive mask generation. We also propose a new metric tailored to ProbSOR, as existing SOR metrics only support deterministic rankings. Additionally, we establish a ProbSOR benchmark comprising 15,000 probabilistic SOR samples for evaluation purposes. Extensive experiments demonstrate that ProbSOR achieves strong performance in salient object ranking and mask generation.

Problem



Problem

ior. We observe that human viewing behaviors are influenced by prior viewing trajectory and exhibit variability, meaning that the next object of fixation should be represented as a conditional probabilistic distribution rather than a definitive choice. while previous SOR methods and evaluation metrics fail to account for this inherent randomness. Motivated by this, we propose to reformulate SOR into a

- Our model: **accept multiple answers/ranks**
- Previous model: a single absolute GT/rank
- If we consider top 3 objects
 - GT: [0, 1, 2] SA-SOR: 1.0 probability: 0.16
 - Pred: [1, 0, 2] SA-SOR: 0.5 probability : 0.26
- If we consider top 2 objects
 - GT: [0, 1] SA-SOR: 1.0 probability : 0.23
 - Pred: [1, 0] SA-SOR: **-1.0** probability : 0.29

New Evaluation Metric

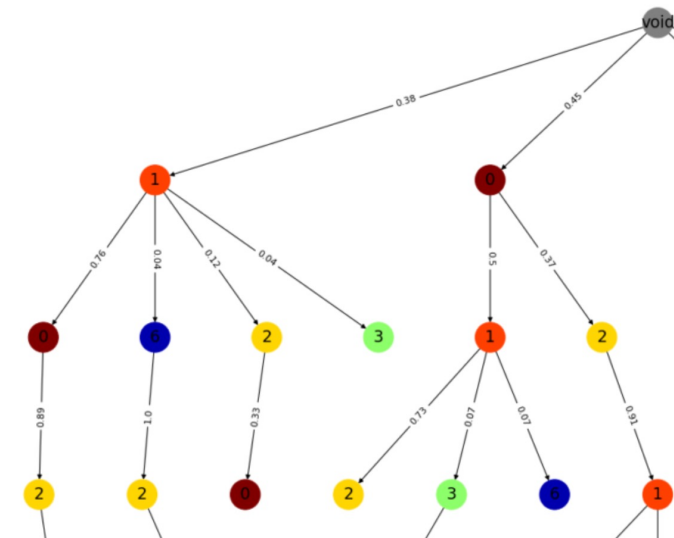
- **Model Input:** RGB image
- **Model Output:** Salient Object Ranking (# dynamic)
- **Metrics:** *Reward-based tree traversal strategy*
 - Starts at root node
 - Go down the tree following the predicted ranking
 - At each node: collects reward = human arrival probability
 - Final score: cumulative rewards along path
 - **Core Concept:** It evaluates SOR by how well they mirror real human attention shifts at each decision point – replacing deterministic 'right/wrong' judgments with plausibility scoring.

Key features of new metrics:

New Evaluation Metric

- Human-Centric Benchmark
 - Rewards = Empirical attention shift probabilities
 - Directly encodes how humans actually allocate attention to objects
- Stochastic Consistency
 - Scores SOR by cumulative attention shift probabilities (not binary or linear association between two sequences.)
 - A rank through high-probability nodes scores well → even if it differs from individual human ranks.
 - Naturally accommodates attention variability across viewers
- Tree Traversal = Simulated Visual Exploration
 - Path following: Mimic attentional shifts between objects
 - Cumulative rewards: Quantifies overall behavioral alignment
- Probabilistic Validation
 - Deterministic metric: "Is object B always viewed after object A?" ❌
 - Our metric: "How probable is this sequence given real human data?" ✅

Previous metrics: Pearson correlation coefficient (measure the strength of a linear association between two variables. It is defined as the covariance of the two variables divided by the product of their standard deviations.)



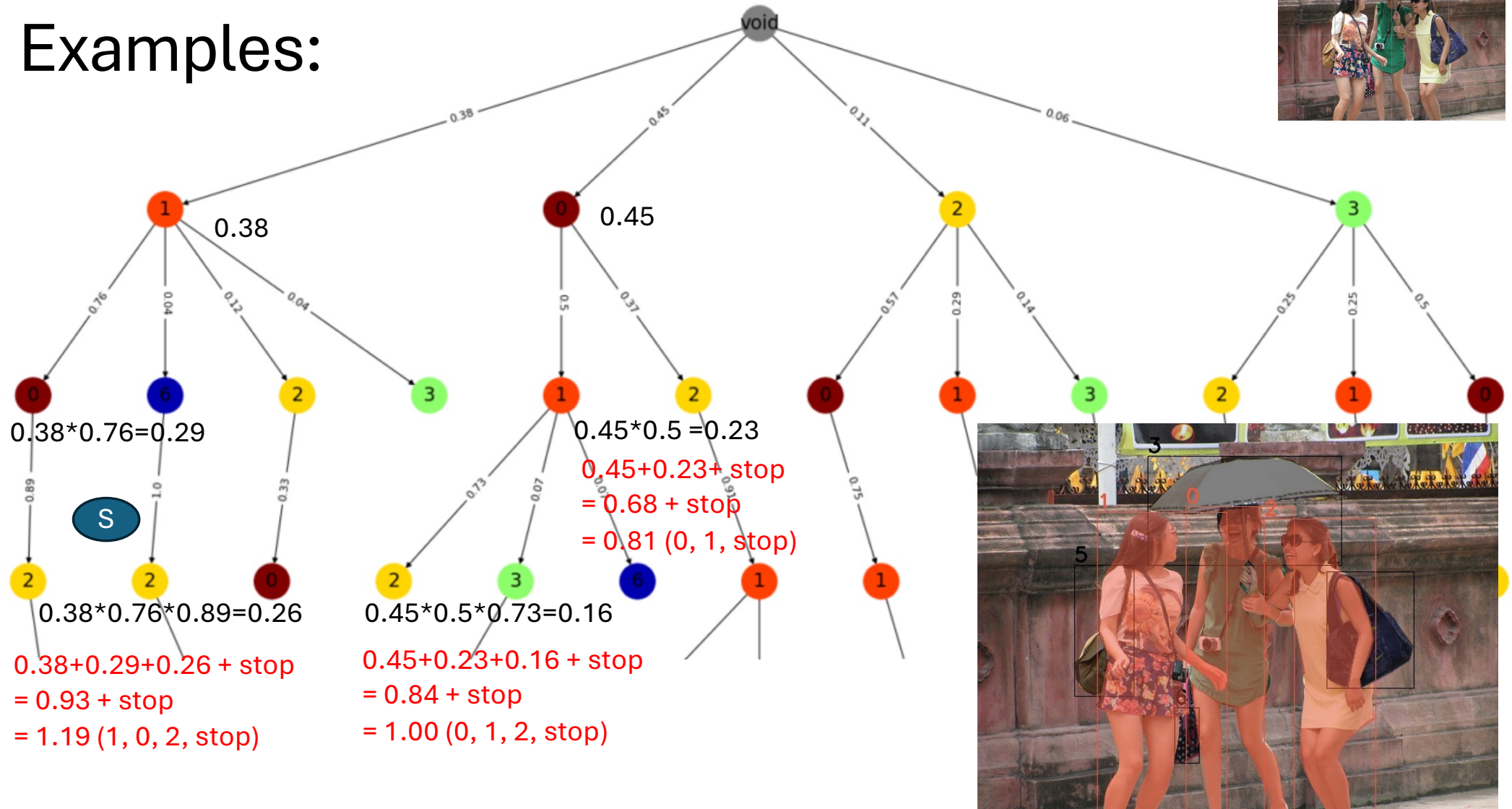
Comparison

Existing Metrics	Our Metrics
Assume one “correct” rank	Accept multiple valid ranks
Penalize deviations	Rewards <i>statistically likely</i> choices
Ignore viewer variability	Embeds population-level attention shift distribution

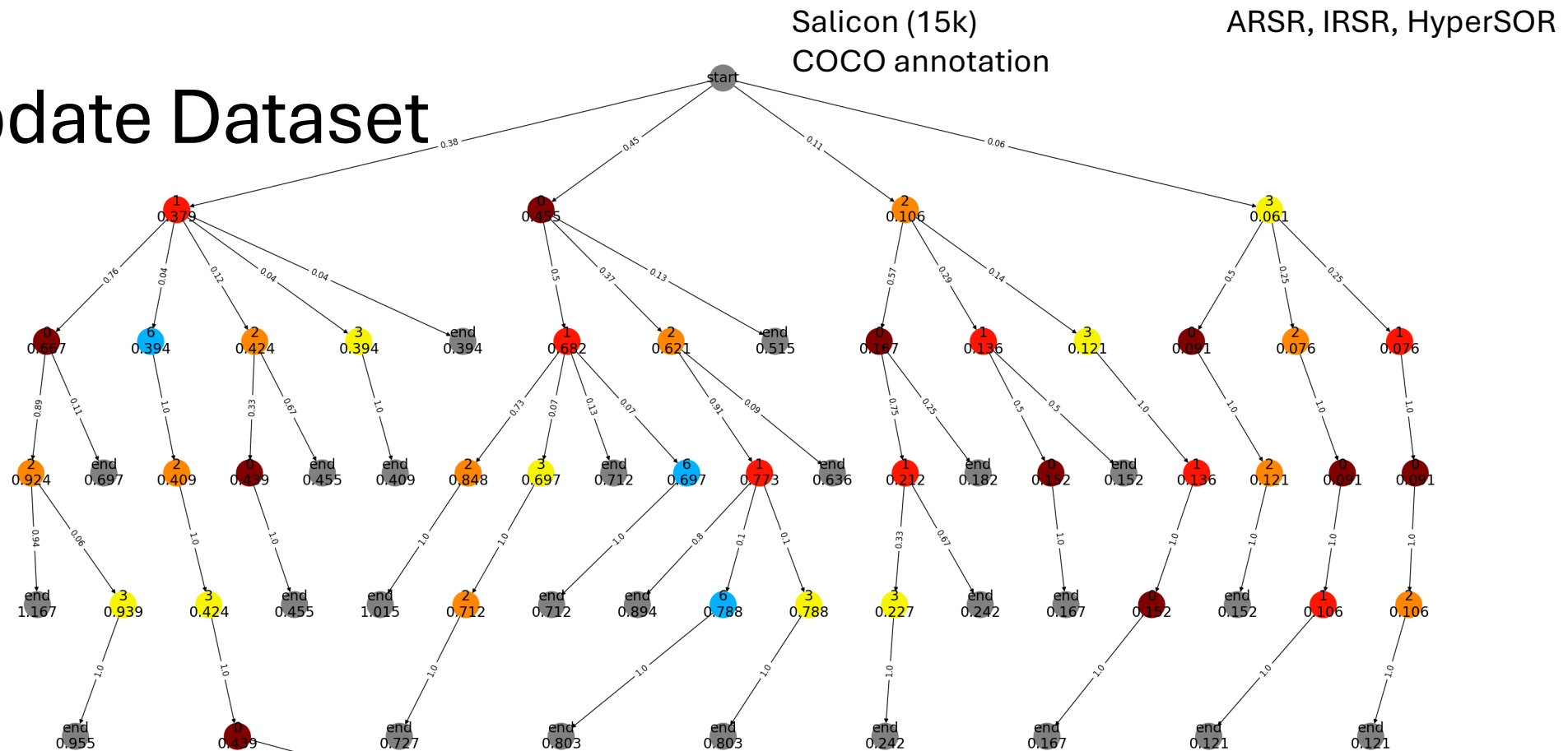
In essence: Instead of asking "Is this the exact path humans take?" (which doesn't exist), we ask "How behaviorally plausible is this path?" – quantified through observable human attention shift probabilities.

Examples:

New Evaluation Metric



Update Dataset



- The tree contains salient and less salient objects
- --> Exclude less salient objects to make the entire tree easier to interpret.

Update Dataset

We can estimate the probability that an object will be noticed or visited as:

$$P(A \text{ is visited}) = \frac{\text{The number of observers who visit } A}{\text{The number of observers}}$$

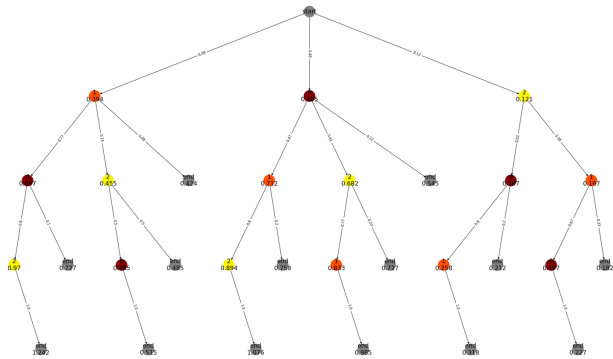
Salient objects are the most visually prominent elements in a scene. Therefore, we can identify salient objects by selecting those with a higher visitation probability P —meaning they are more likely to attract attention.

We simply adopt thresholding method:

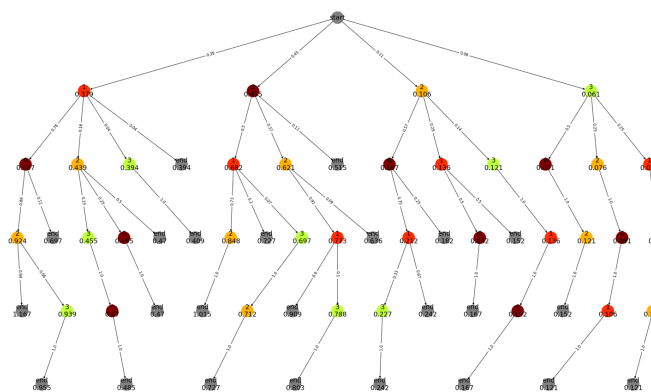
- Assume an **average** of 5 salient objects per image (based on empirical observations), OR Top 5 = 0.12
- Compute the visitation probability threshold using existing benchmark datasets, e.g., ARSR or IRSR. ARSR: 0.18
IRSR: 0.30

Update Dataset

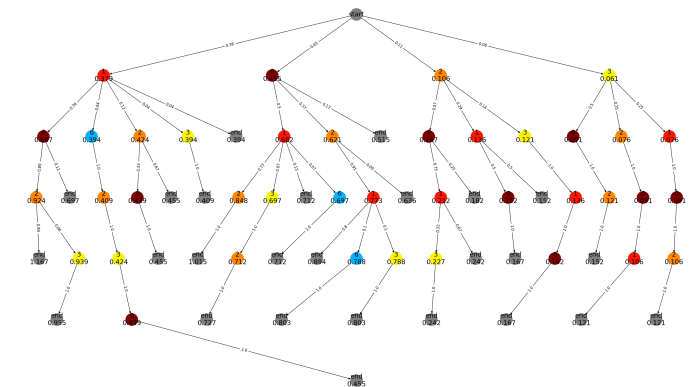
1. Filter out those less salient objects, i.e., $P(\text{OBJ is visited}) < 0.18$
2. Construct the dataset based on previous discussion



Threshold = 0.18

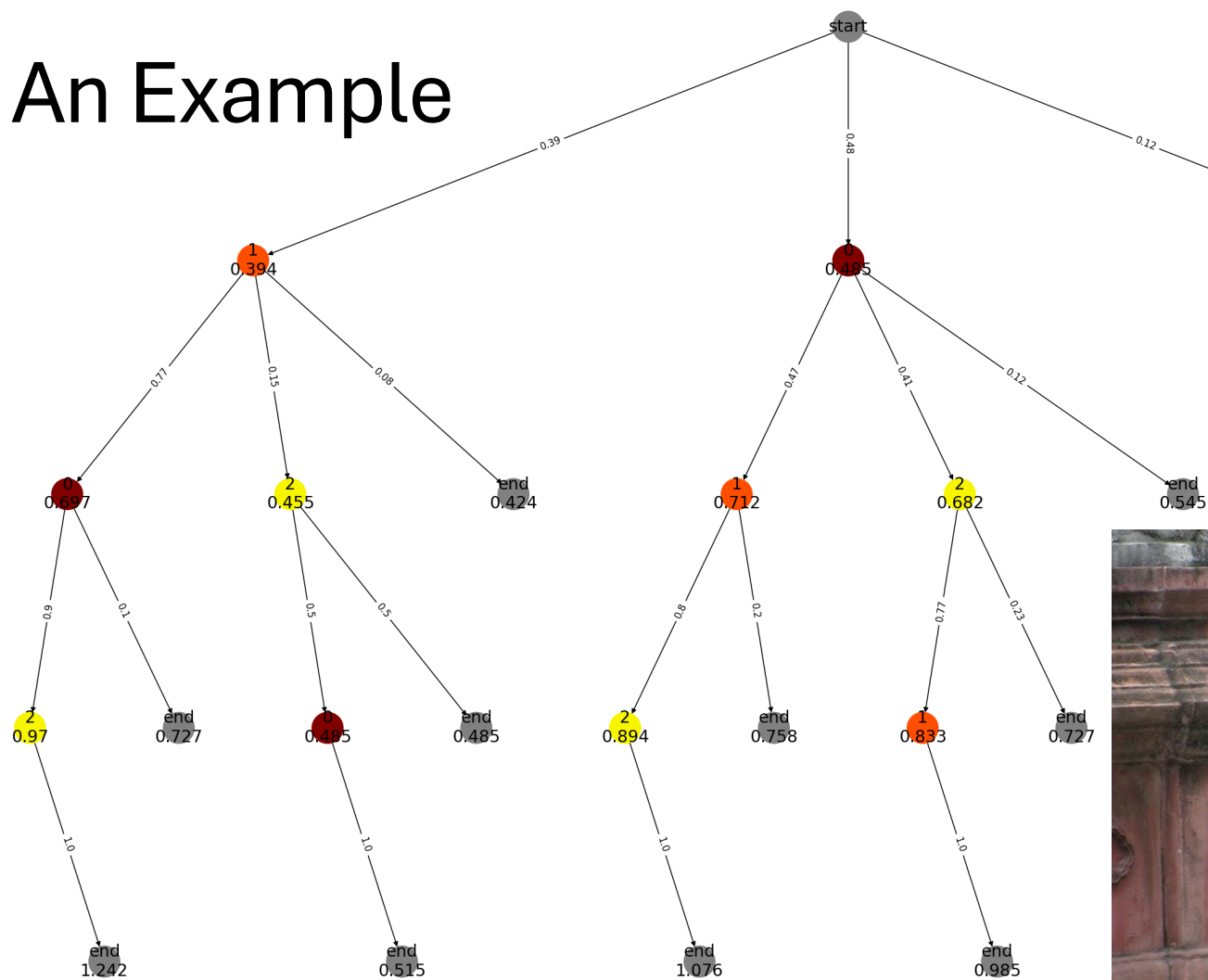


Threshold = 0.12

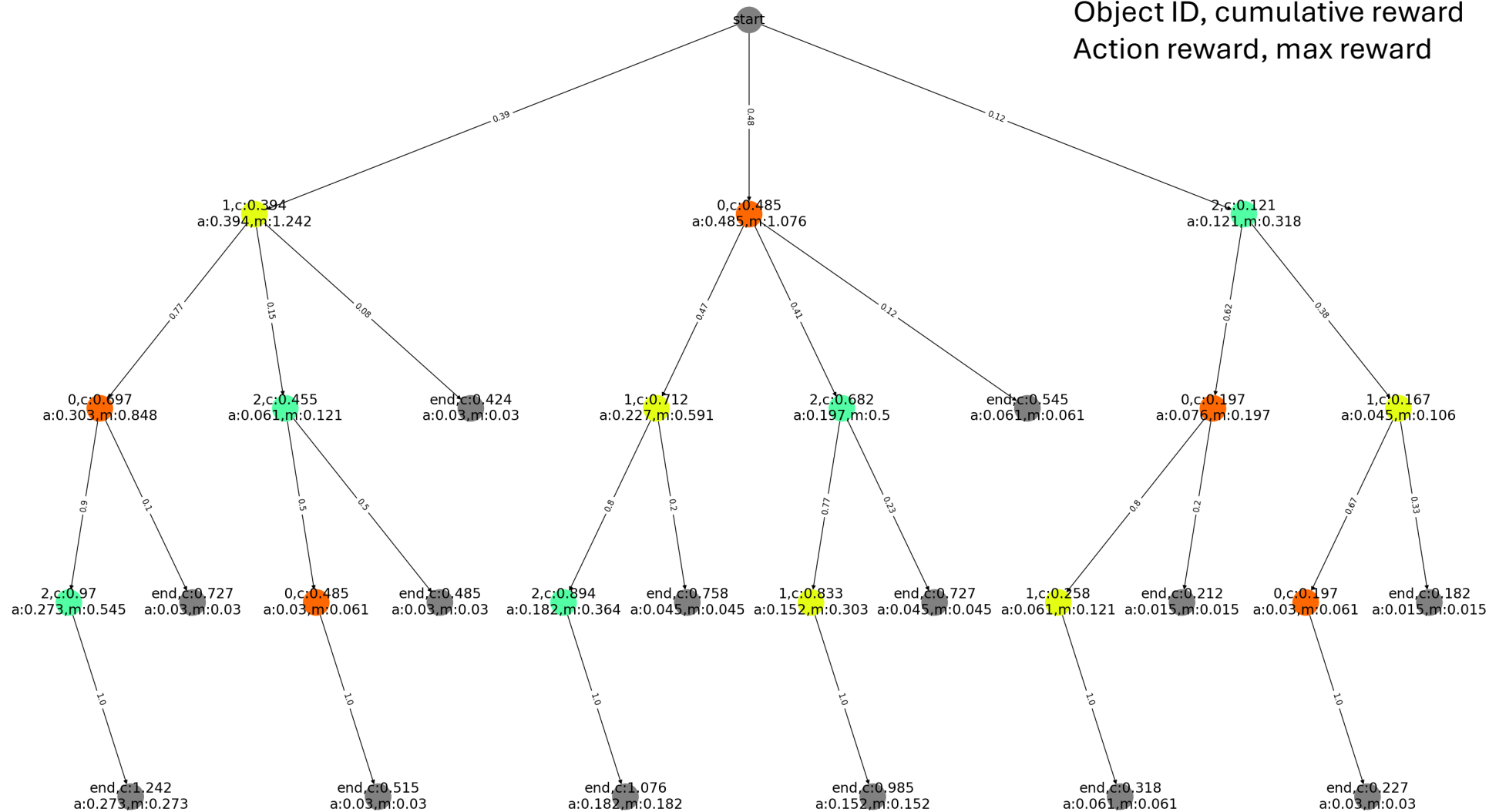


No Threshold

An Example



Object ID, cumulative reward
Action reward, max reward

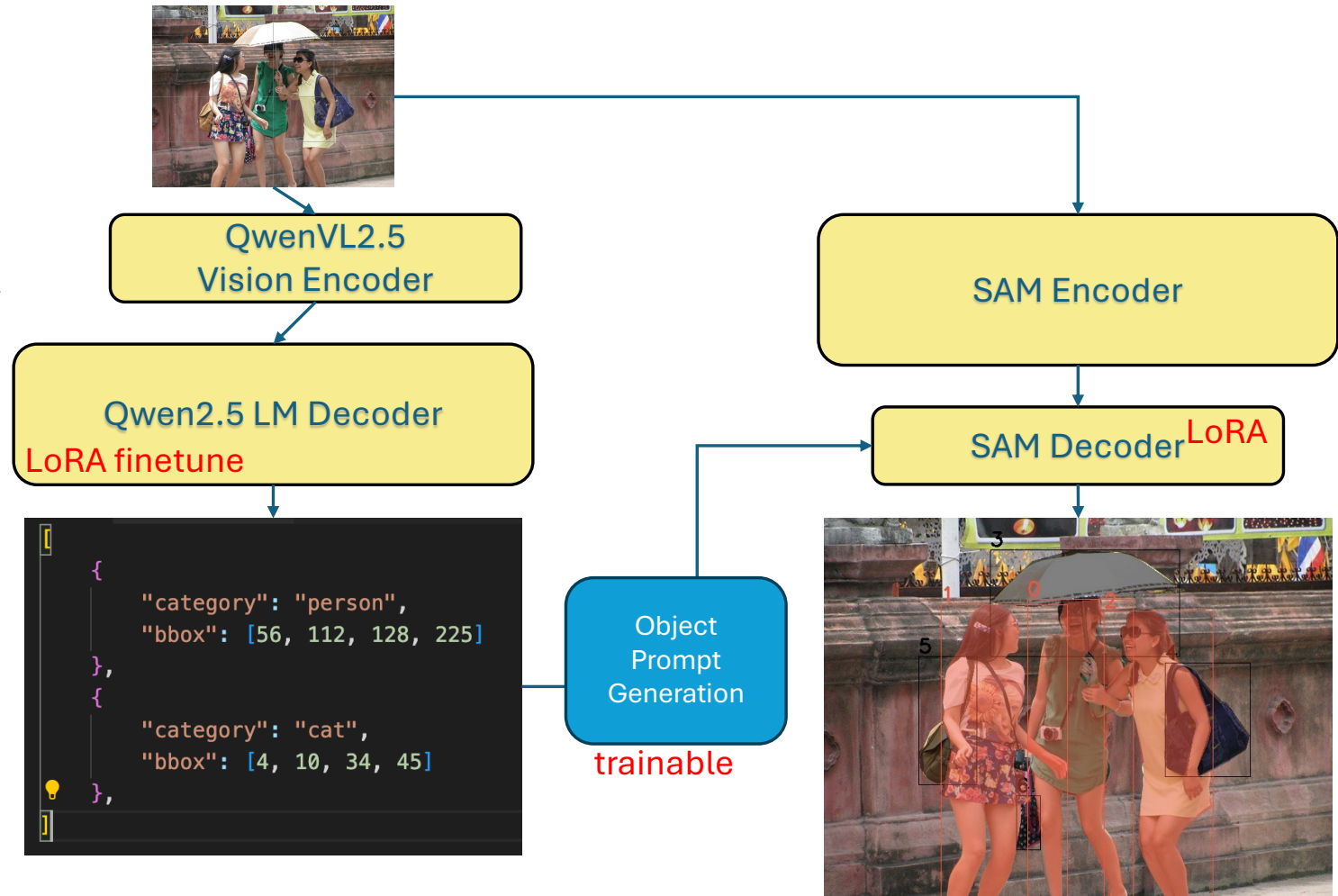


Model Design

1. OPG module: Object Prompt Generation module bridges the QwenVL and SAM for mask prediction.
2. QwenVL is finetuned for SOR prediction.
3. SAM decoder is finetuned to generate object mask.

Three novel things:

1. OPG module
2. Training strategy
3. Metrics



Results: (mask prediction is good enough)

runs.summary["Table-000000301867-1"]

Filter

image



000000301867

data

```
{'name': '000000301867', 'width': 640, 'height': 480, 'records': [{'rank': 1,
'category': 'person', 'bbox': {'x1': 150.0, 'y1': 106.07142857142857, 'x2':
280.0, 'y2': 473.5714285714286}, 'mask': None, 'mask_rle': {'size': [480,
640], 'counts': b'^W2c1T=?
PE^NR8i1eG`NU8f1eG`NU8g1cGaNj1oNP3e2PKbNh1POT3c2mJdNg1P0W3c2fJiNj1iN^3e2\\JUNP
{'rank': 2, 'category': 'person', 'bbox': {'x1': 255.71428571428572, 'y1':
98.57142857142857, 'x2': 384.28571428571433, 'y2': 473.5714285714286},
'mask': None, 'mask_rle': {'size': [480, 640], 'counts': b'bfh3f0P>c0@?
C:E;mC[N\\:R2TEhNU:
[3B`0A<_0c0_0i0[0b0@=D7I9G6K6I6K5K5K4L5L3L3M2N3N002N101N1010101mNlEVMV:h2QFRMP
f5X3ZI_Lm0;i5\\4VJdKL2V0aNU5bNfKf2@eNi4dNhK^2LjN\\4hNhKw25P0Q4iNkKR2:U0j3hNnKw
[NA\\KoNb4^1CAkKT0d4W1CBjKW0g4S1@BmK[Of4n0AA\\K@m4f0@W0f2`0Q6K5K4Kanh3'}},
{'rank': 3, 'category': 'person', 'bbox': {'x1': 324.28571428571433, 'y1':
114.64285714285715, 'x2': 485.71428571428567, 'y2': 473.5714285714286},
'mask': None, 'mask_rle': {'size': [480, 640], 'counts':
b'gmi4`0<3Q=S1K4L5L2N3L3N201eGQNb3Q2]LQN`3Q2^LRN`3n1aLSN]3n1dLRNZ3n1gLSNW3n1jL
7DF;;GD8<KA4`0N^0gNCkKo0b5Z0^N0kKh0j5TOYN9iKd0`6aNdMS3\\2\\LbMV3`2gL_M\\3b2aL]M
fLP1Y0]NQ4d0dLo0AVNn3k0`Ln0HPNj3S1\\Ln0Y6R0eIm0^6S0aIm0`6S0^Im0d6S0[Il0f6U0XI\\
'end_of_detection': True, 'output_text': '<|vision_end|>This is the input
image with height = 448 and width = 448.
<|im_end|>\\n<|im_start|>assistant\\n|1|person|105,99,196,442|\\n|2|person|179,92
<|im_end|>'}]}
```

≡ ≡ = -

← < 1 -1 of 1 > →

Export as CSV Columns... Reset table

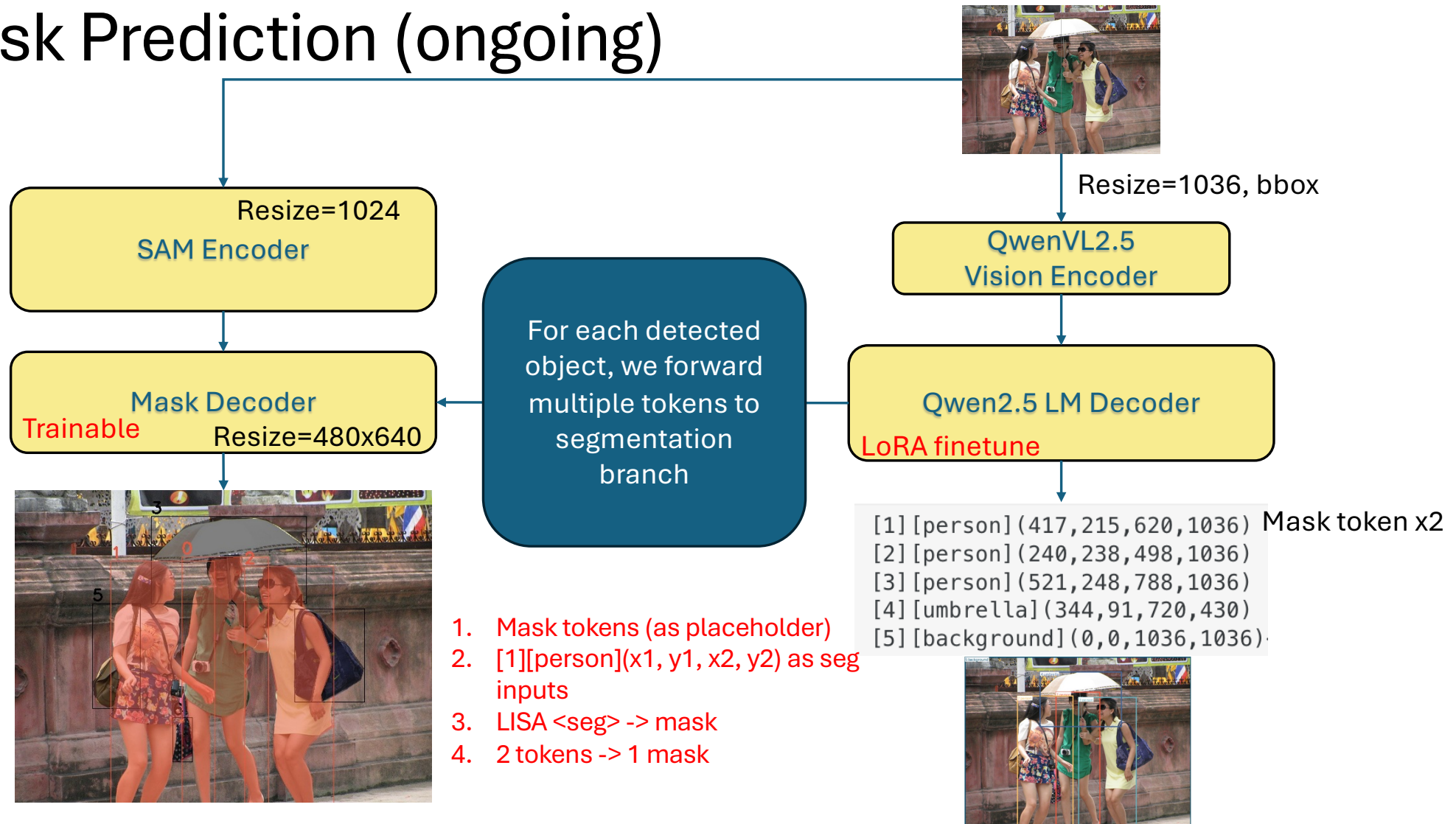
Training Strategy

1. **Exp-1:** Select the best SOR as GT to train the model
 - Global optimal solution: predict the best reward sequence
2. **Exp-2:** Reward-based Learning
 - Image -> Model -> Output ~ Rewards
 - Allow the model to predict different sequences
 - GRPO
3. **Metrics:** reward, normalized reward, iou, mae

Issues

- Inference step is slow, requiring **3 – 9 seconds** for a single image. Our evaluation set contains **5000** images, then it takes up to **8.5 hours** to finish the evaluation.
- Training time is acceptable that it requires **<8 hours** on 1 A100 to finetune the model. Use small image as input (like 448), we can use 1 4090GPU to train.
- Solutions:
 - Use a smaller VLM or smaller image as input
 - there may be technical optimizations that could help improve inference speed

Mask Prediction (ongoing)



Contributions

1. We reformat the SOR to probability-based.
2. The first VLM for SOR, (considering the uncertainty of human attention shift (GRPO), **new mask generation strategy (under revision)**)
3. New Benchmark and new metric

Probabilistic Salient Object Ranking

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Salient Object Ranking (SOR) aims to study how humans visually explore complex
2 scenes by predicting an ordered sequence of objects that will attract our attention.
3 Existing SOR approaches typically model this ranking deterministically, assuming
4 a single, fixed ranking of attention. However, such deterministic SOR fails to
5 capture the true nature of human attention. We observe that human attention shifts
6 exhibit variability and stochasticity, showing that the next object of fixation is not a
7 definitive choice but rather a probability distribution. Yet, existing SOR methods
8 and evaluation metrics cannot account for this inherent randomness. To bridge this
9 gap, we propose **ProbSOR**, a novel **Probabilistic Salient Object Ranking** model
10 that explicitly learns the uncertainty of attention shifts by incorporating Group
11 Relative Policy Optimization (GRPO). We leverage a Vision-Language Model
12 (VLM) as the foundation for ProbSOR and **extend it with autoregressive mask**
13 **generation**. We also propose a new metric tailored to ProbSOR, as existing SOR
14 metrics only support deterministic rankings. Additionally, we establish a ProbSOR
15 benchmark comprising 15,000 probabilistic SOR samples for evaluation purposes.
16 Extensive experiments demonstrate that ProbSOR achieves strong performance in
17 salient object ranking and mask generation.

18 1 Introduction

19 2 Related Works

20 3 Benchmark

21 4 Methodology

22 5 Experiments

23 Metrics: $reward = \sum_{i=1}^n r_i$, where $r_i = p_i$, p_i is the possibility of reaching $node_i$ from the top of
24 the tree. $normalized_reward = \frac{reward}{optimal_reward}$ is used to report the relative goodness of the SOR
25 prediction. We report iou, mae to assess the quality of mask prediction.