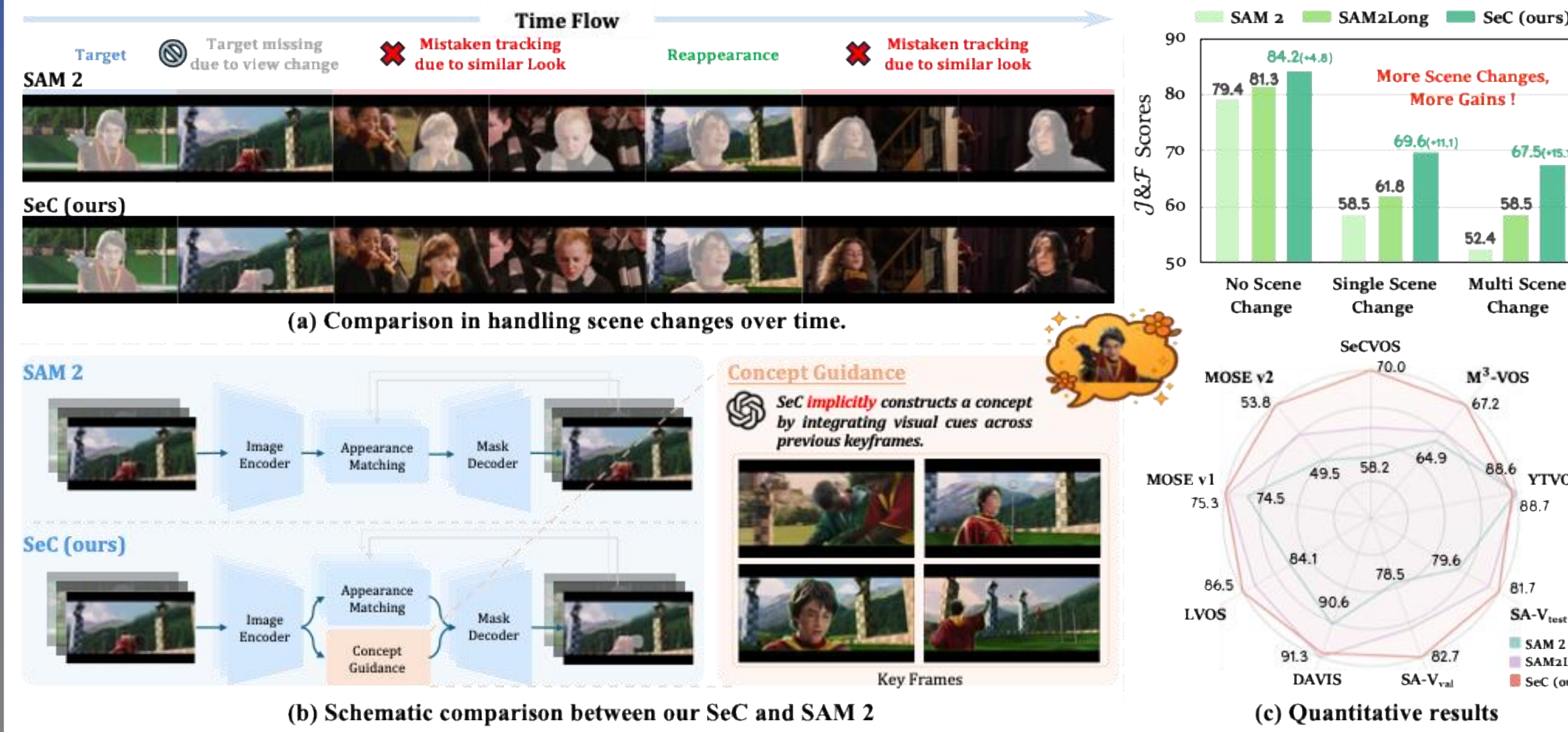


## Motivation

- SAM 2 struggles when facing appearance-similar distractors or abrupt scene transitions.
- If the model understands the object's role (e.g., player or audience) as a concept, can it better resolve SAM 2's limitations in tracking consistency?

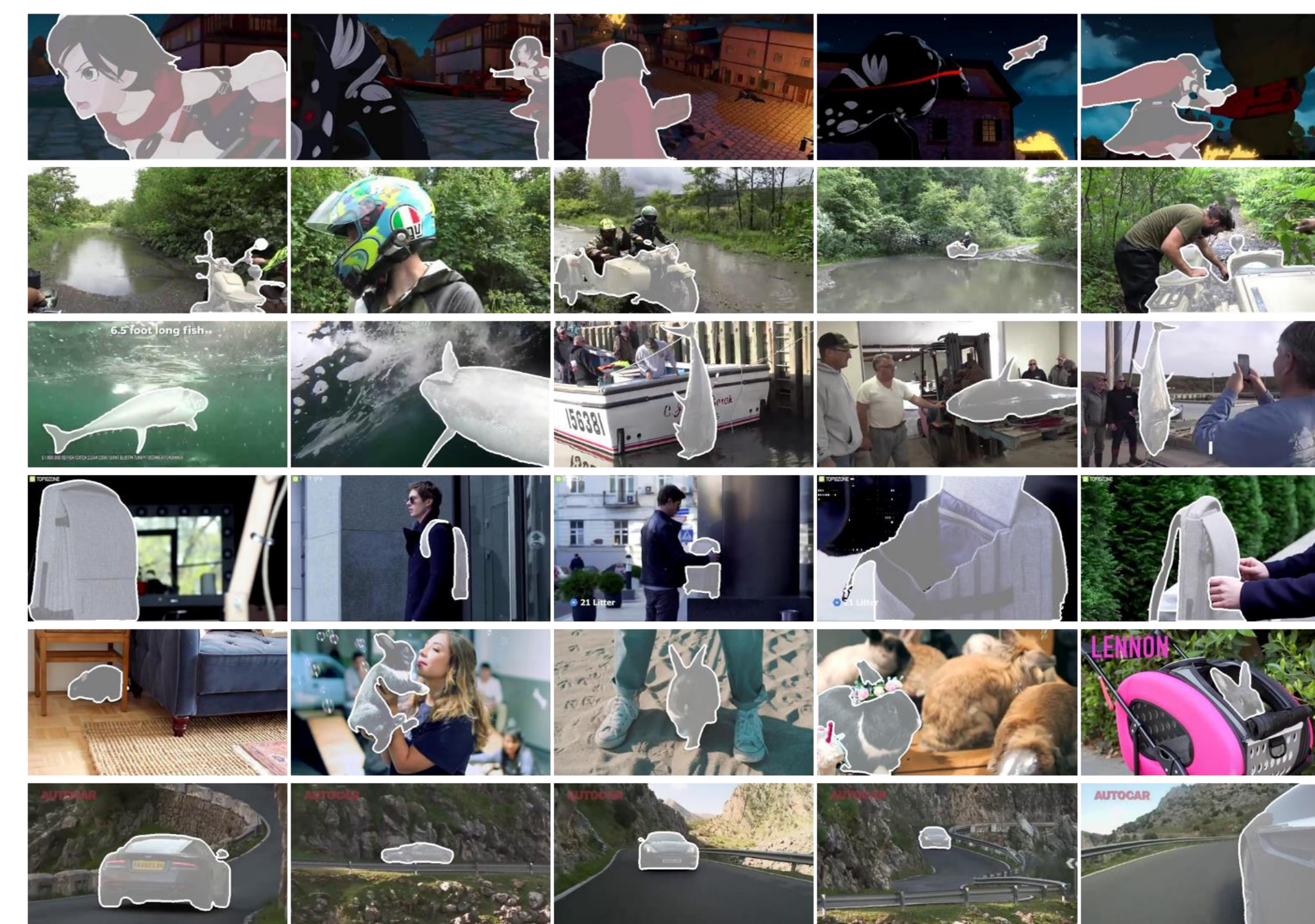


## SeCVOS

### A Benchmark for Semantically Complex VOS

- 160 manually annotated multi-shot video clips;
- Frequent target reappearance and high temporal discontinuity;
- Designed to evaluate performance under high-level semantic reasoning;

### Time Flow



## Experiment Results

- SeC consistently performs better across benchmarks, with more scene changes leading to greater improvements.

Method	No Scene Change			Single Scene Change			Multi Scene Change			Overall	Pixel-level Association	Concept Guidance	SA-V J&F	SeCVOS J&F
	J&F	J	F	J&F	J	F	J&F	J	F	J&F				
Xmem (Cheng & Schwing, 2022)	71.9	72.0	71.8	47.0	47.9	46.2	41.9	42.4	41.4	48.4	✗	✗	78.6	58.2
DEVA (Cheng et al., 2023)	71.6	71.6	71.5	48.5	48.4	48.6	46.4	46.0	46.8	49.7	✓	✗	82.4	62.2
Cutie-base (Cheng et al., 2024)	72.5	72.2	72.8	53.0	52.9	53.2	48.3	47.8	48.9	52.7	✓	✓	82.7	70.0
SAM2.1 (Ravi et al., 2025)	79.4	79.1	79.7	58.5	58.2	58.8	52.4	52.1	52.6	58.2				
SAMURAI (Yang et al., 2024)	81.8	81.6	81.9	60.6	60.6	60.7	59.3	58.9	59.7	62.2				
SAM2.1Long (Ding et al., 2025d)	81.3	81.0	81.6	61.8	61.6	62.0	58.5	58.1	58.9	62.3				
<b>SeC (Ours)</b>	<b>84.2(+4.8)</b>	<b>83.8</b>	<b>84.5</b>	<b>69.6(+11.1)</b>	<b>69.5</b>	<b>69.7</b>	<b>67.5(+15.1)</b>	<b>67.0</b>	<b>68.0</b>	<b>70.0(+11.8)</b>				

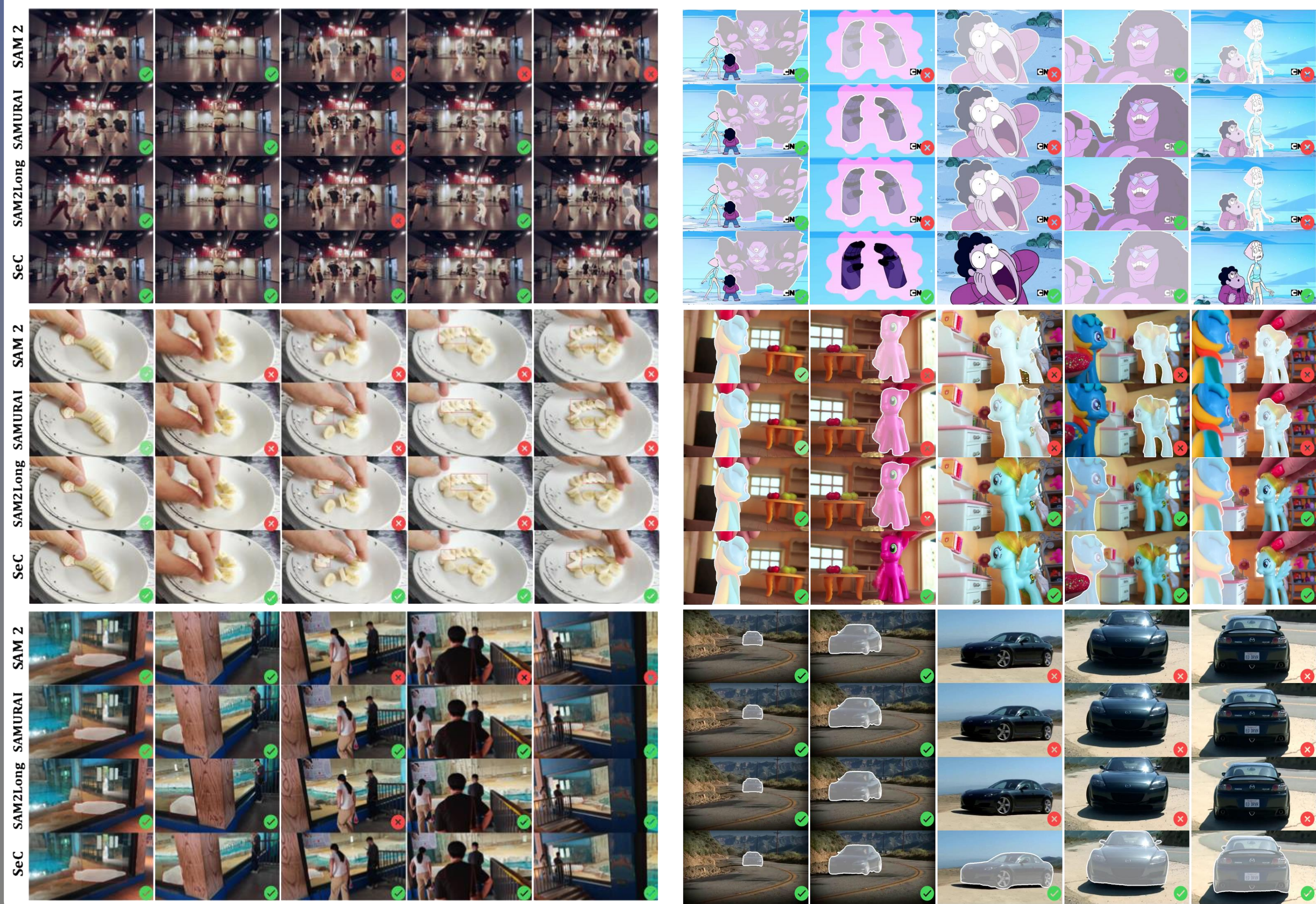
Method	J&F				g	J	J&F
	SA-V val	SA-V test	LVOS v2 val	MOSE v1 val			
STCN (Cheng et al., 2021)	61.0	62.5	60.6	52.5	85.4	82.7	29.7
SwinB-AOT (Yang et al., 2021)	51.1	50.3	-	59.4	85.4	84.5	30.2
SwinB-DeAOT (Yang & Yang, 2022)	61.4	61.8	63.9	59.9	86.2	86.1	32.6
RDE (Li et al., 2022)	51.8	53.9	62.2	46.8	84.2	81.9	32.0
XMem (Cheng & Schwing, 2022)	60.1	62.3	64.5	59.6	86.0	85.6	36.3
SimVOS-B (Wu et al., 2023)	44.2	44.1	-	-	88.0	84.2	-
DEVA (Cheng et al., 2023)	55.4	56.2	-	66.0	87.0	85.4	38.3
ISVOS (Wang et al., 2023)	-	-	-	-	88.2	86.3	-
TarVIS (Athar et al., 2023)	-	-	-	-	85.2	82.7	-
UNINEXT (Yan et al., 2023)	-	-	-	-	81.8	78.6	-
UniVS (Li et al., 2024)	-	-	-	-	76.2	71.5	-
JointFormer (Zhang et al., 2025b)	-	-	-	-	90.1	87.4	37.7
Cutie-base (Cheng et al., 2024)	60.7	62.7	-	69.9	87.9	87.0	42.8
Cutie-base+ (Cheng et al., 2024)	61.3	62.8	-	71.7	88.1	87.5	-
SAM 2.1 (Ravi et al., 2025)	78.6	79.6	84.1	74.5	90.6	88.7	49.5
SAMURAI (Yang et al., 2024)	79.8	80.0	84.2	72.6	89.9	88.3	51.1
SAM2.1Long (Ding et al., 2025d)	81.1	81.2	85.9	75.2	91.4	88.7	51.5
<b>SeC (Ours)</b>	<b>82.7</b>	<b>81.7</b>	<b>86.5</b>	<b>75.3</b>	<b>91.3</b>	<b>88.6</b>	<b>53.8</b>

#Concept Tokens	J&F		
	J&F	J	F
1	70.0	69.7	70.2
2	70.0	69.7	70.3
4	69.9	69.6	70.2

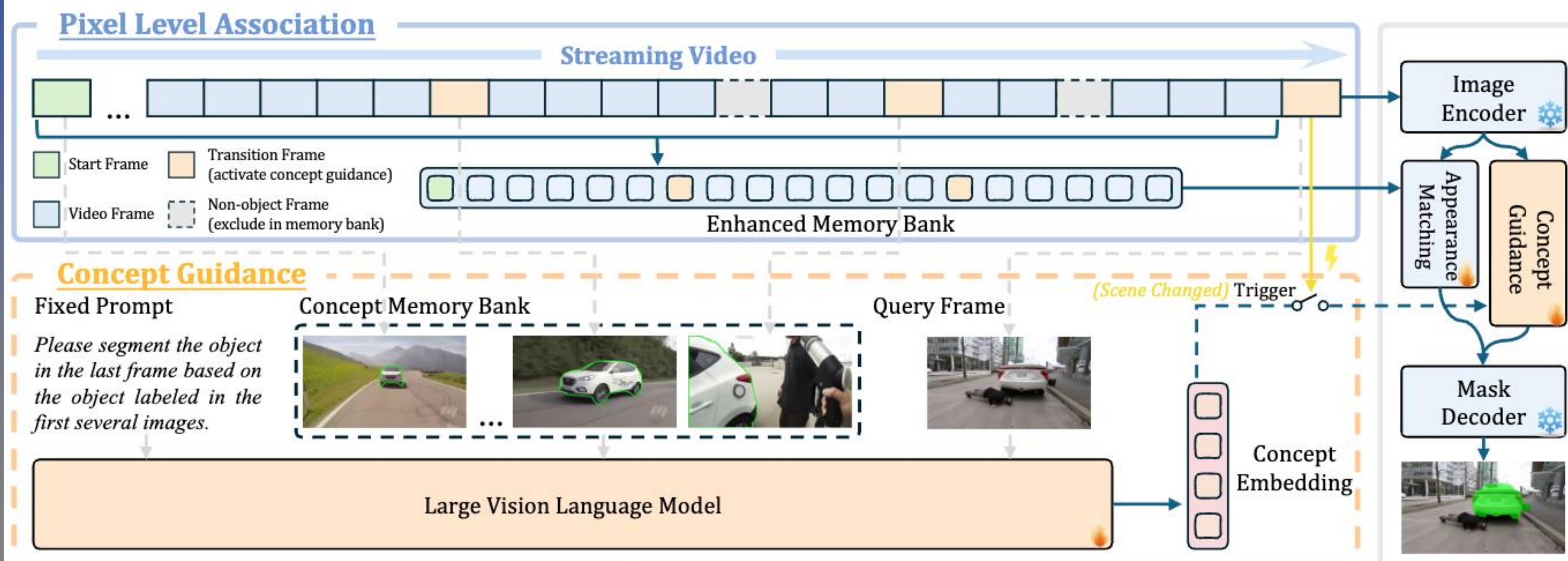
  

Scene Detector	J&F		
	J&F	J	F
ABS DIFF	69.0	68.8	69.3
ORB	68.7	68.5	69.0
SSIM	69.3	69.1	69.6
FLOW	69.5	69.3	69.7
HSV (ours)	70.0	69.7	70.2



## Method: Progressive Concept Construction

- SeC leverage the capability of LLMs to introduce concept-level features.
- Concept Guidance only activate when the scene is changed and concept bank is updated dynamically.



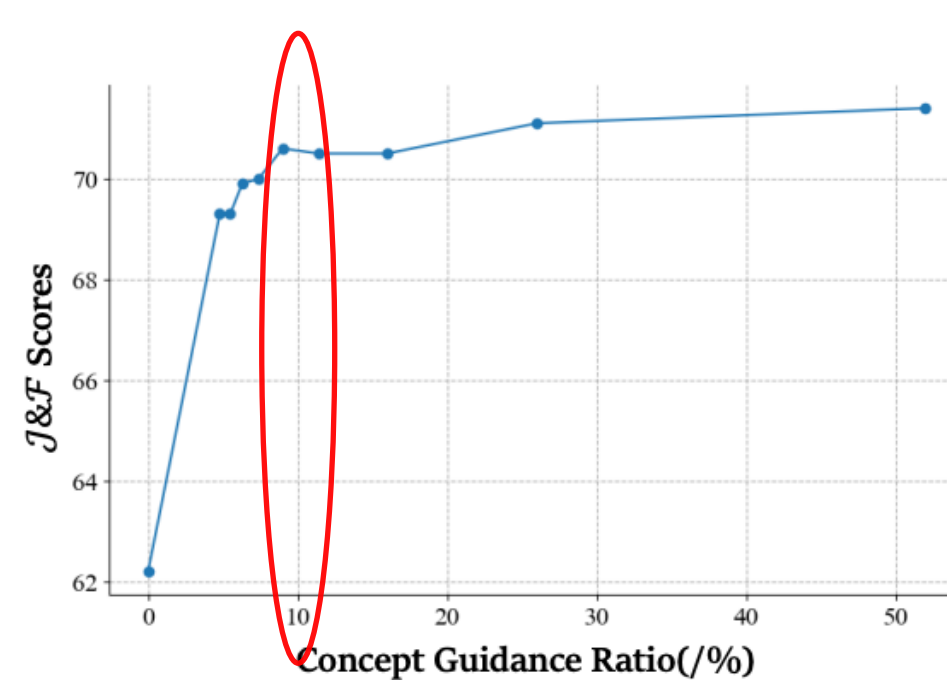
## Discussion

### Does SeC progressively construct concept-level representation?

- SeC learns comprehensive concepts on the fly.
- Offline setting even makes the performance stronger.

### Does SeC require frequent concept guidance?

- Activating less than 10% of frames is sufficient;
- Good trade-off between accuracy and efficiency.



Benchmark	Method	J&F	Con. Guid. Ratio (%)	Throughput (s <sup>-1</sup> )
SeCVOS	SeC	70.0	7.4	14.8
	SAM 2	58.2	N/A	22.0
SA-V	SeC	82.7	1.0	18.1
	SAM 2	78.6	N/A	22.0