

# Recurrent Dynamic Embedding for Video Object Segmentation

CVPR 2022



# 提纲

2



Authors



Introduction



Methodology



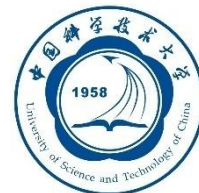
Experiments and Result



Conclusion



Inspiration



# Authors

## □ Authors

### Recurrent Dynamic Embedding for Video Object Segmentation

Mingxing Li<sup>1\*</sup>, Li Hu<sup>2\*</sup>, Zhiwei Xiong<sup>1†</sup>, Bang Zhang<sup>2</sup>, Pan Pan<sup>2</sup>, Dong Liu<sup>1</sup>

<sup>1</sup>University of Science and Technology of China

<sup>2</sup>Alibaba DAMO Academy, Alibaba Group

mxli@mail.ustc.edu.cn {zwxiong, dongeliu}@ustc.edu.cn

{hooks.hl, zhangbang.zb, panpan.pp}@alibaba-inc.com



引用次数

[查看全部](#)

	总计	2017 年至今
引用	10121	9658
h 指数	40	37
i10 指数	115	109



# 提纲

About Article

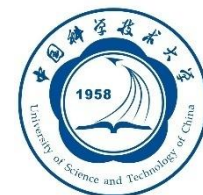
Introduction

Methodology

Experiments and Result

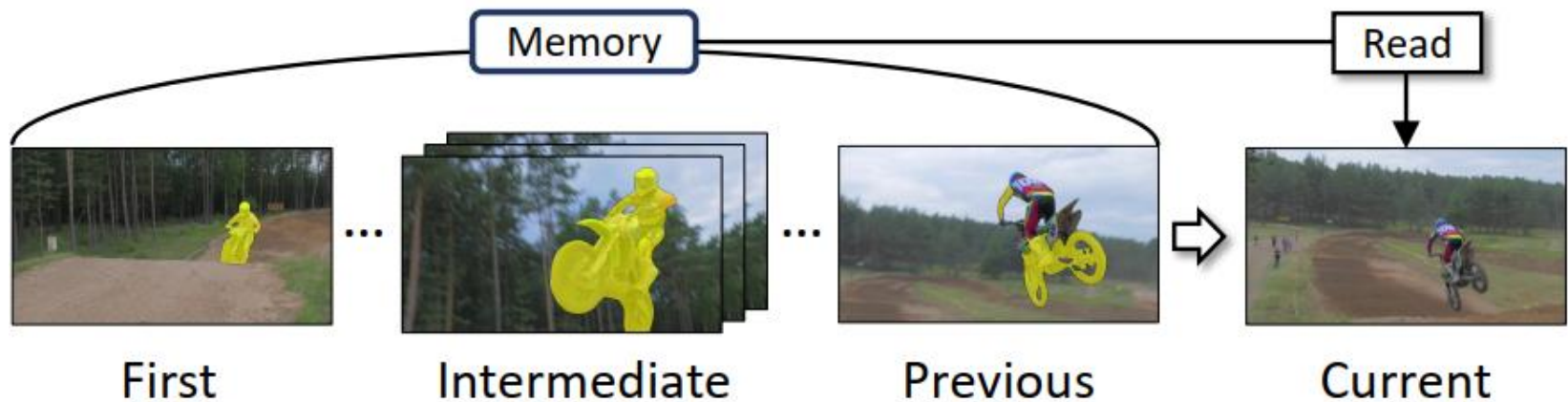
Conclusion

Inspiration



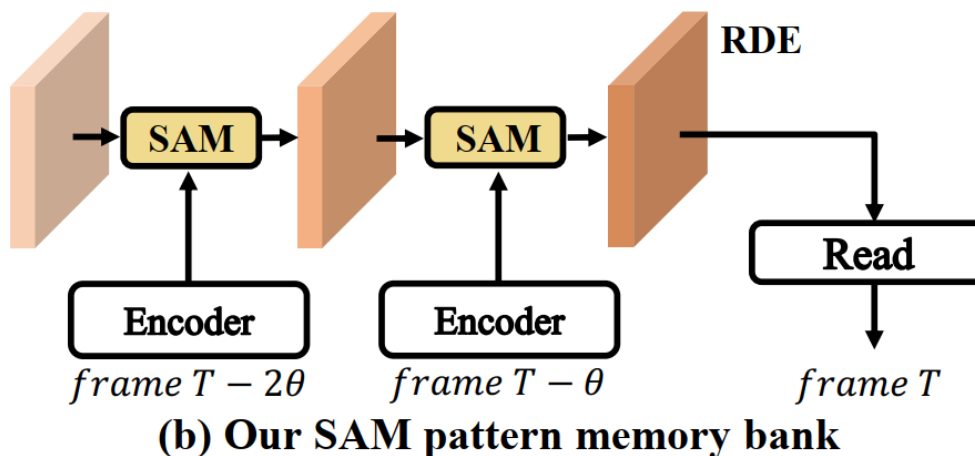
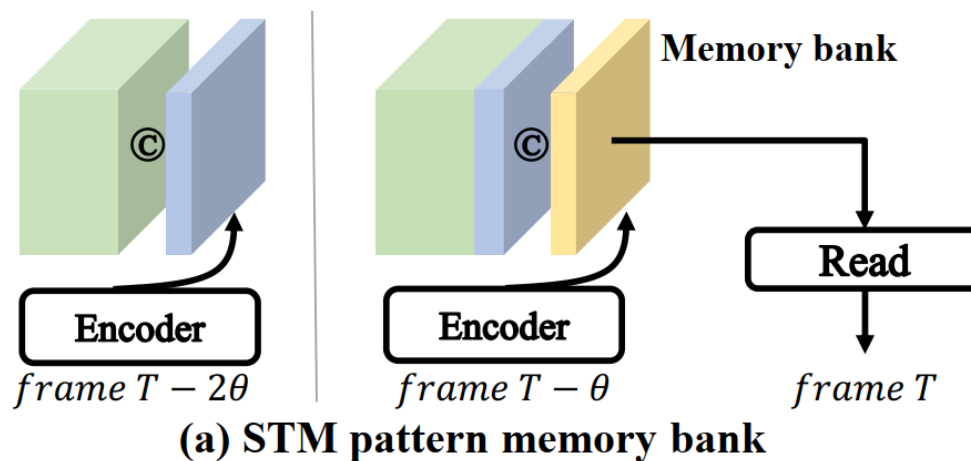
# Introduction

## □ Space-Time Memory (STM)



# Introduction

## Memory Bank



# 提纲

About Article

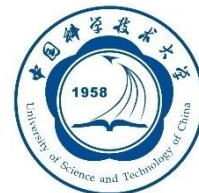
Introduction

Methodology

Experiments and Result

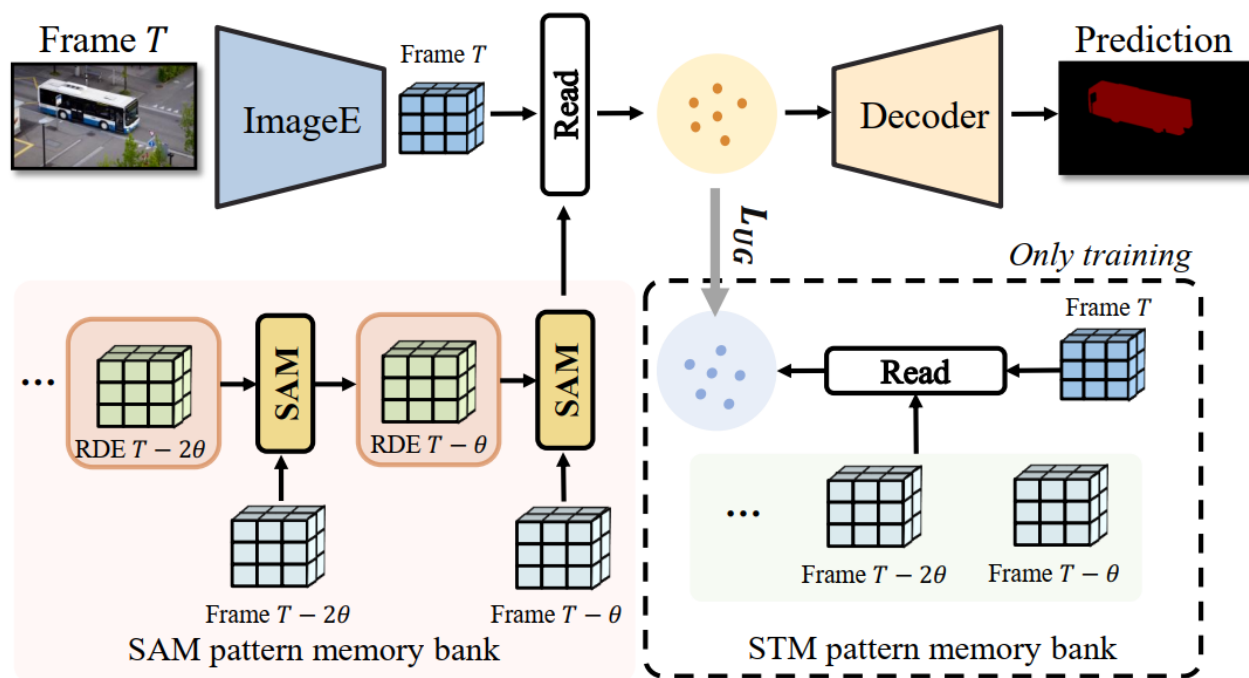
Conclusion

Inspiration

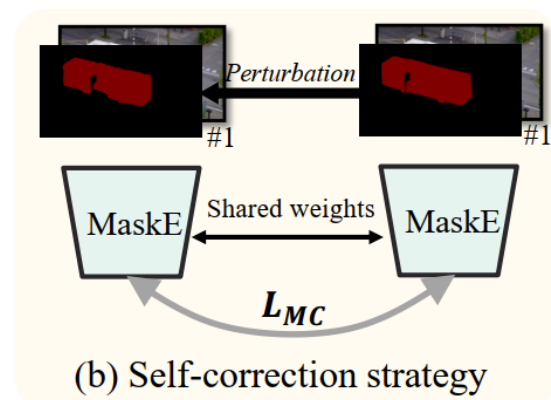


# Methodology

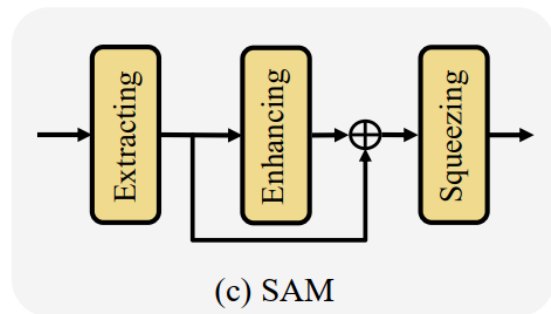
## □ Framework



(a) Main pipeline of our framework



(b) Self-correction strategy

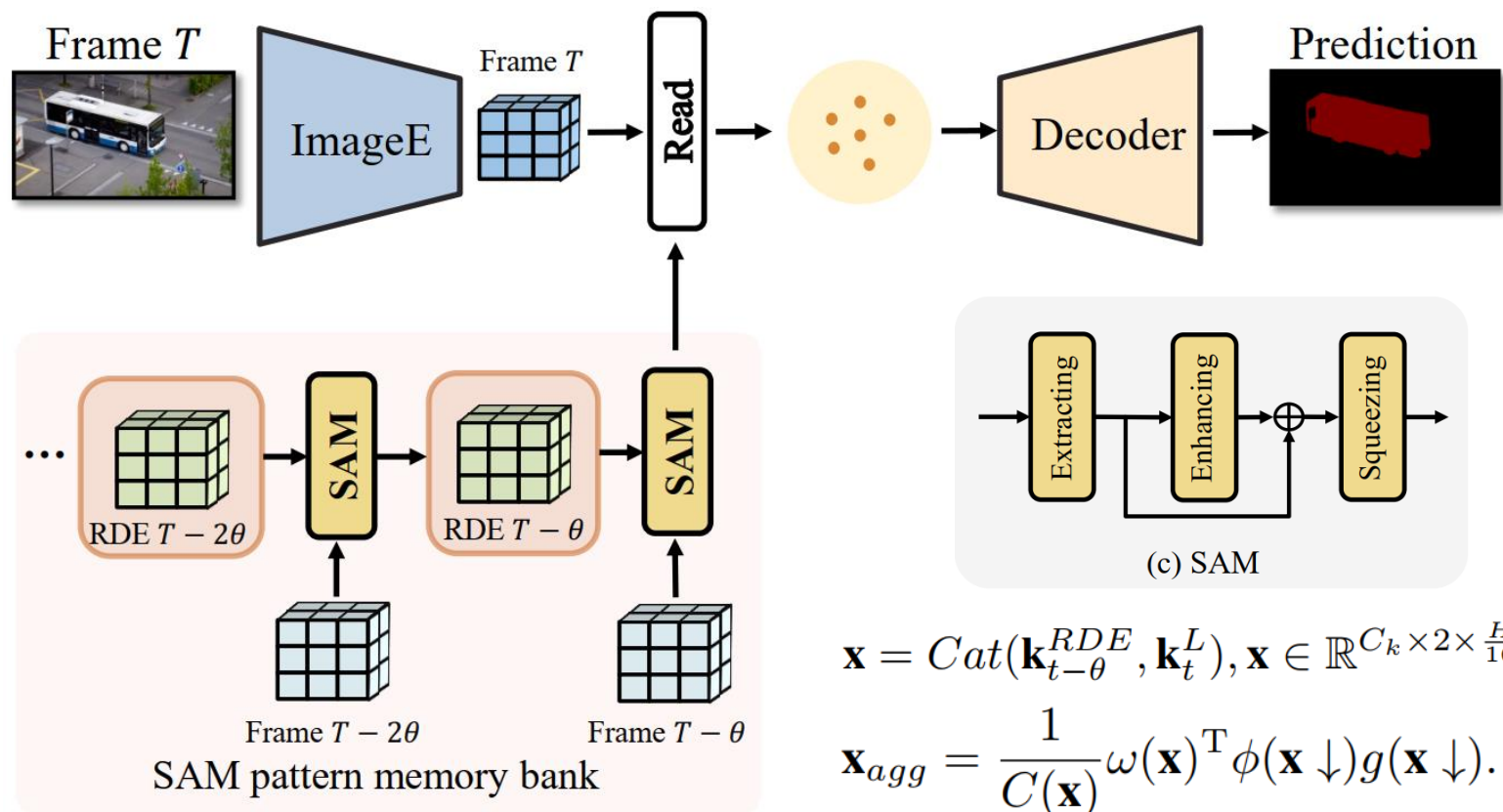


(c) SAM



# Methodology

## Recurrent Dynamic Embedding



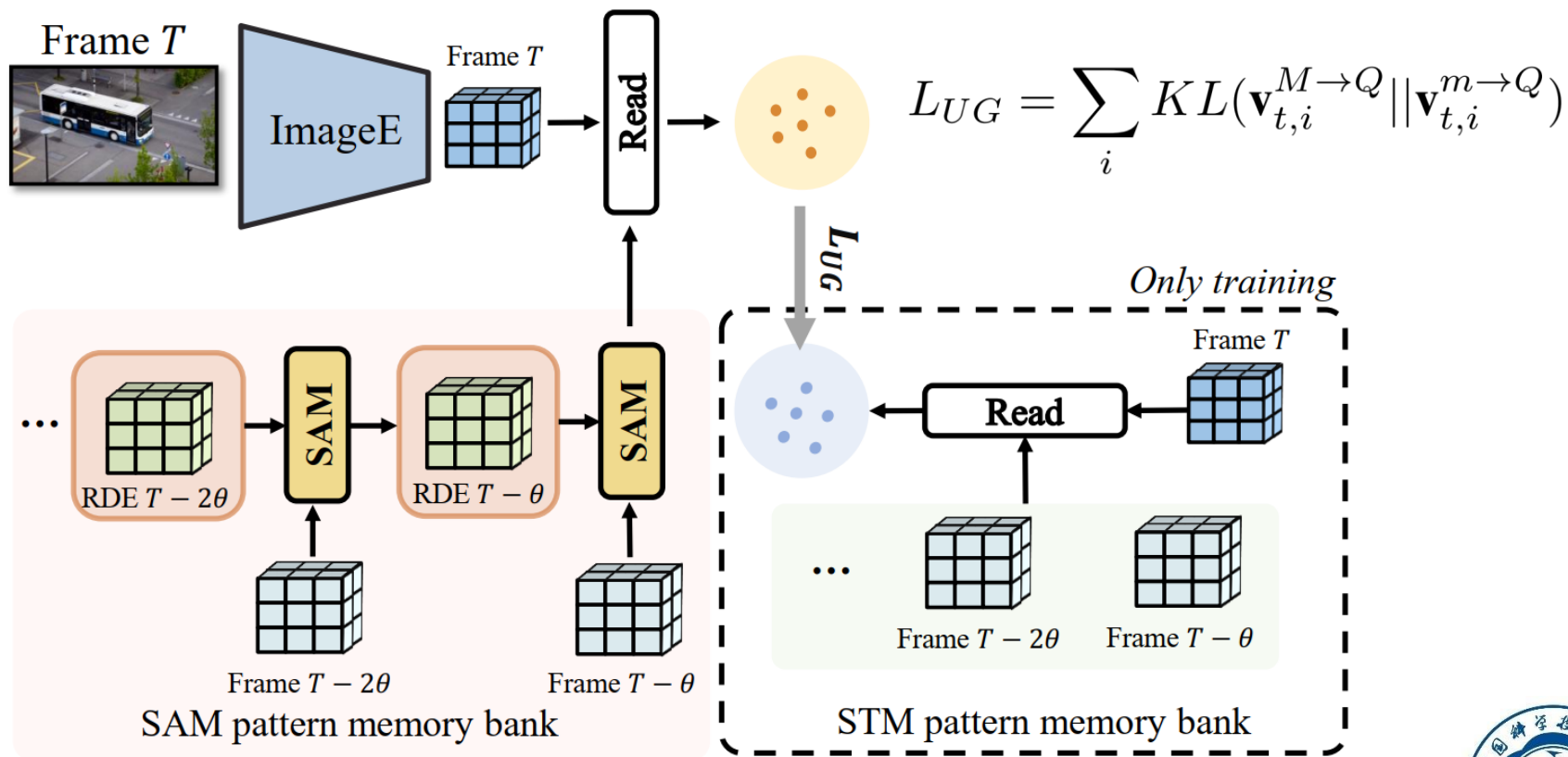
$$\mathbf{x} = \text{Cat}(\mathbf{k}_{t-\theta}^{RDE}, \mathbf{k}_t^L), \mathbf{x} \in \mathbb{R}^{C_k \times 2 \times \frac{H}{16} \times \frac{W}{16}}$$

$$\mathbf{x}_{agg} = \frac{1}{C(\mathbf{x})} \omega(\mathbf{x})^T \phi(\mathbf{x} \downarrow) g(\mathbf{x} \downarrow).$$

$$\mathbf{k}_t^{RDE} = \text{Squeeze}(\mathbf{x}_{agg} + \text{ASPP}(\mathbf{x}_{agg})).$$

# Methodology

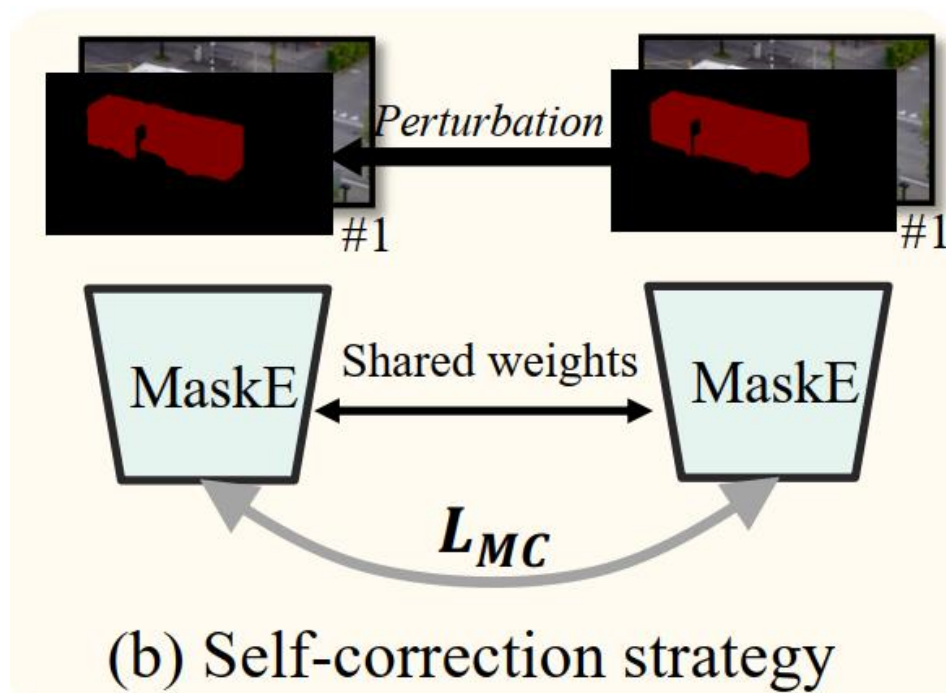
## □ Unbiased Guidance Loss



(a) Main pipeline of our framework

# Methodology

## □ Self-correction Strategy



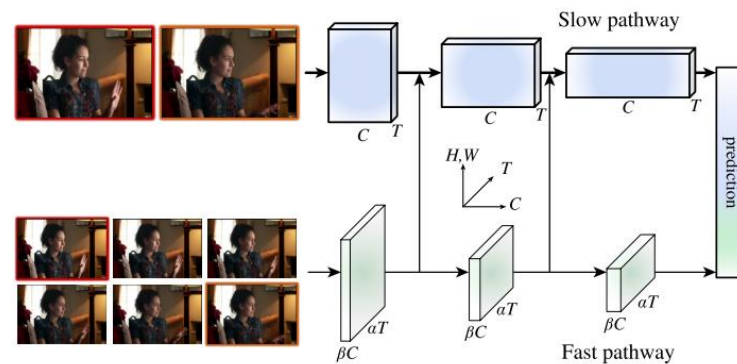
$$L_{MC} = KL(\mathbf{k}_1 || \ddot{\mathbf{k}}_1) + \sum_i KL(\mathbf{v}_{1,i} || \ddot{\mathbf{v}}_{1,i})$$

# Methodology

## □ Training Strategy

$$L_{Seg} = \frac{1}{2} \left( \sum_i \sum_{t=2,4} \underbrace{BCE(\tilde{\mathbf{y}}_{t,i}^M, \mathbf{y}_{t,i})}_{STM \text{ pattern item}} + \sum_i \sum_{t=3,5} \underbrace{BCE(\tilde{\mathbf{y}}_{t,i}^m, \mathbf{y}_{t,i})}_{SAM \text{ pattern item}} \right)$$

$$Loss = L_{Seg} + \mathbb{1}[t = 3, 5] \mu L_{UG} + \gamma L_{MC}$$



# 提纲

About Article

Introduction

Methodology

Experiments and Result

Conclusion

Inspiration



# Experiments and Result

## □ Experiments

Method	CC	$\mathcal{J}\&\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$	FPS
RMNet <sup>†</sup> [37]	×	88.8	88.9	88.7	11.9
STM <sup>†</sup> [25]	×	89.3	88.7	89.9	6.3
KMN <sup>†</sup> [30]	×	90.5	89.5	91.5	8.4
LCM <sup>†</sup> [13]	×	90.7	89.9	91.4	8.5
HMMN <sup>†</sup> [31]	×	90.8	89.6	92.0	10.0
MiVOS <sup>†*</sup> [5]	×	91.0	89.7	92.4	16.9
STCN <sup>†*</sup> [6]	×	<b>91.7</b>	<b>90.4</b>	<b>93.0</b>	<b>26.9</b>
GCNet [17]	✓	86.6	87.6	85.7	25.0
CFBI+ <sup>†</sup> [41]	✓	89.9	88.7	91.1	5.9
SwiftNet <sup>†</sup> [33]	✓	90.4	<b>90.5</b>	90.3	25.0
<b>RDE-VOS<sup>†</sup></b>	✓	91.1	89.7	92.5	<b>35.0</b>
<b>RDE-VOS<sup>†*</sup></b>	✓	<b>91.6</b>	90.0	<b>93.2</b>	<b>35.0</b>

Table 3. Results on the DAVIS 2016 validation set. CC denotes constant cost during the inference.

Method	CC	Overall	$\mathcal{J}_{seen}$	$\mathcal{F}_{seen}$	$\mathcal{J}_{unseen}$	$\mathcal{F}_{unseen}$
STM <sup>†</sup> [25]	×	79.2	79.6	83.6	73.0	80.6
MiVOS <sup>†*</sup> [5]	×	82.4	80.6	84.7	78.2	85.9
STCN <sup>†*</sup> [6]	×	<b>84.2</b>	<b>82.6</b>	<b>87.0</b>	<b>79.4</b>	<b>87.7</b>
CFBI <sup>†</sup> [40]	✓	81.0	80.6	85.1	75.2	83.0
SST <sup>†</sup> [8]	✓	81.8	80.9	-	76.6	-
<b>RDE-VOS<sup>†</sup></b>	✓	81.9	81.1	85.5	76.2	84.8
<b>RDE-VOS<sup>†*</sup></b>	✓	<b>83.3</b>	<b>81.9</b>	<b>86.3</b>	<b>78.0</b>	<b>86.9</b>

Table 4. Results on the YouTube-VOS 2019 validation set.

Method	CC	$\mathcal{J}\&\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$	FPS
STM <sup>†</sup> [25]	×	81.8	79.2	84.3	10.2
KMN <sup>†</sup> [30]	×	82.8	80.0	85.6	<8.4
JOINT <sup>†</sup> [23]	×	83.5	80.8	86.2	4.0
LCM <sup>†</sup> [13]	×	83.5	80.5	86.5	<8.5
RMNet <sup>†</sup> [37]	×	83.5	81.0	86.0	<11.9
MiVOS <sup>†*</sup> [5]	×	84.5	81.7	87.4	11.2
HMMN <sup>†</sup> [31]	×	84.7	81.9	87.5	<10.0
STCN <sup>†*</sup> [6]	×	<b>85.3</b>	<b>82.0</b>	<b>88.6</b>	<b>20.2</b>
GCNet [17]	✓	71.4	69.3	73.5	<25.0
Liang <i>et al.</i> [19]	✓	74.6	73.0	76.1	4.0
G-FRTM <sup>†</sup> [26]	✓	76.4	-	-	18.2
PReMVOS [21]	✓	77.8	73.9	81.7	0.01
SwiftNet <sup>†</sup> [33]	✓	81.1	78.3	83.9	<25.0
SST <sup>†</sup> [8]	✓	82.5	79.9	85.1	-
Ge <i>et al.</i> <sup>†</sup> [10]	✓	82.7	80.2	85.3	6.7
<b>RDE-VOS<sup>†</sup></b>	✓	84.2	80.8	87.5	<b>27.0</b>
<b>RDE-VOS<sup>†*</sup></b>	✓	<b>86.1</b>	<b>82.1</b>	<b>90.0</b>	<b>27.0</b>

Method	CC	600p	$\mathcal{J}\&\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$
STM <sup>†</sup> [25]	×	✓	72.2	69.3	75.2
KMN <sup>†</sup> [30]	×	✓	77.2	74.1	80.3
RMNet <sup>†</sup> [37]	×	×	75.0	71.9	78.1
Ge <i>et al.</i> <sup>†</sup> [10]	×	×	75.2	72.0	78.3
STCN <sup>†*</sup> [6]	×	×	77.8	74.3	81.3
MiVOS <sup>†*</sup> [5]	×	×	<b>78.6</b>	<b>74.9</b>	<b>82.2</b>
CFBI <sup>†</sup> [40]	✓	×	74.8	71.1	78.5
Ge <i>et al.</i> <sup>†</sup> [10]	✓	×	75.2	72.0	78.3
CFBI+ <sup>†</sup> [41]	✓	×	75.6	71.6	79.6
<b>RDE-VOS<sup>†</sup></b>	✓	×	77.4	73.6	81.2
<b>RDE-VOS<sup>†*</sup></b>	✓	×	<b>78.9</b>	<b>74.9</b>	<b>82.9</b>



# Experiments and Result

## □ Experiments

Variants	$\mathcal{J} \& \mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$
Strategy permutation			
RDE	81.8	78.0	85.7
First frame	71.6	67.8	75.4
First frame & RDE	85.3	81.6	89.0
Latest frame	80.4	76.9	83.8
Latest frame & RDE	82.2	78.4	86.0
First frame & latest frame	84.6	81.0	88.2
F & L & RDE	85.4	81.6	89.2
First frame $\times 2$ & latest frame	85.1	81.5	88.7
First frame & latest frame $\times 2$	84.0	80.4	87.6
2F & L & RDE	<b>86.1</b>	<b>82.1</b>	<b>90.0</b>
Sampling interval $\theta$			
2F & L & RDE ( $\theta = 2$ )	85.1	81.4	88.9
2F & L & RDE ( $\theta = 3$ )	<b>86.1</b>	<b>82.1</b>	<b>90.0</b>
2F & L & RDE ( $\theta = 4$ )	85.1	81.5	88.8
2F & L & RDE ( $\theta = 5$ )	84.2	80.5	87.9



# Experiments and Result

## □ Experiments

	Ablation Settings	$\mathcal{J} \& \mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$
Loss	w/o $L_{MC}$	83.7	80.5	86.9
	w/o $L_{UG}$	82.9	79.5	86.4
	w/o $L_{MC}$ & $L_{UG}$	82.5	79.1	86.0
	$L_{Seg}$ w/o STM pattern item	83.0	79.4	86.6
	Full	<b>84.2</b>	<b>80.8</b>	<b>87.5</b>

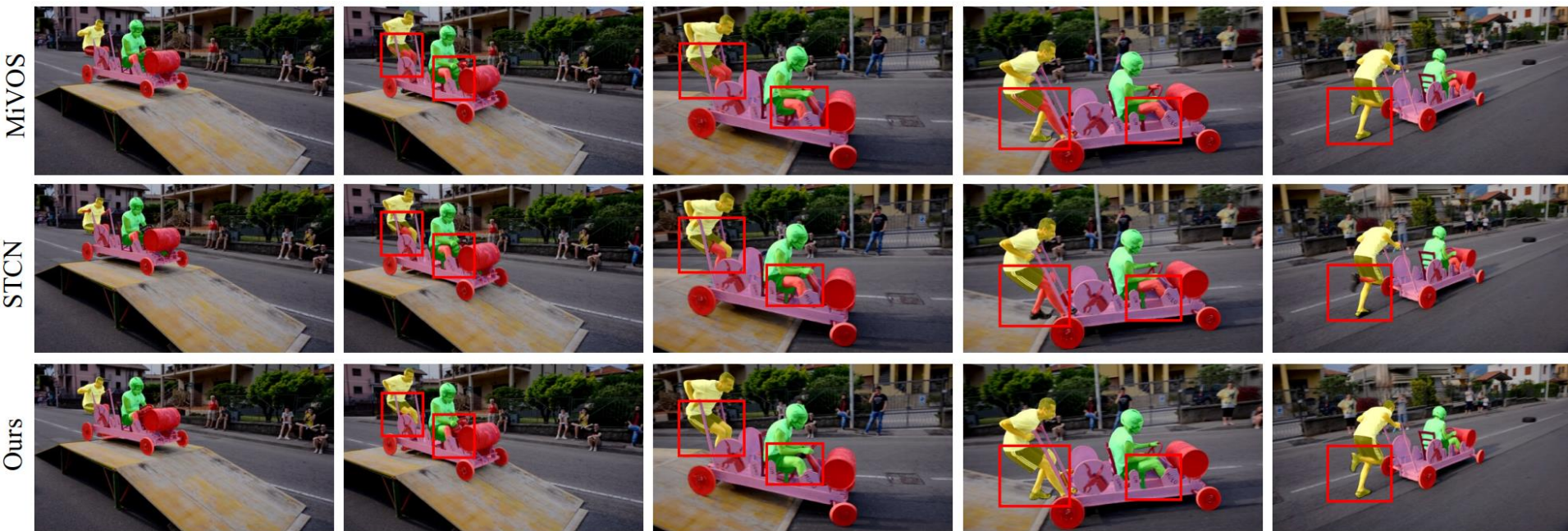
Table 6. Ablation of different loss functions without the BL30K [5] pre-training.





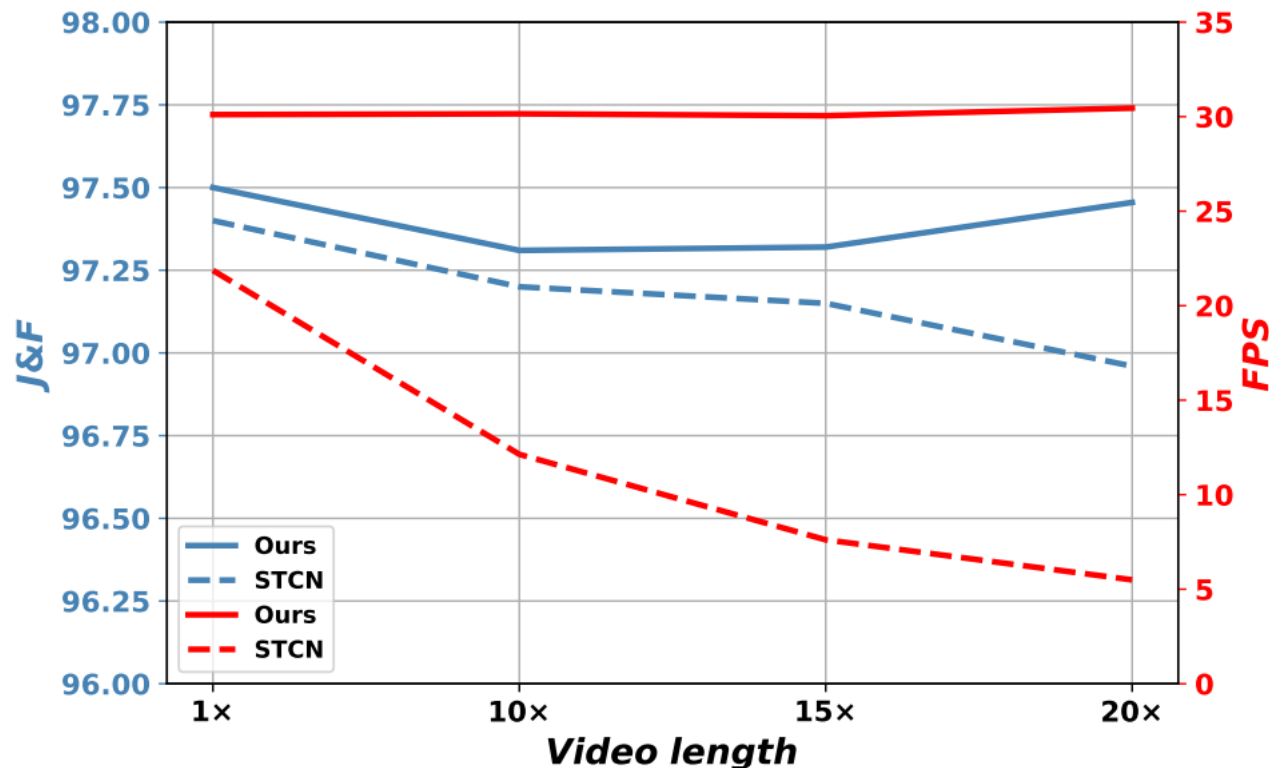
# Experiments and Result

## □ Experiments



# Experiments and Result

## □ Experiments



# 提纲

About Article

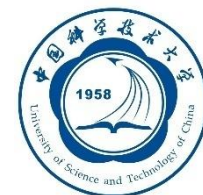
Introduction

Methodology

Experiments and Result

Conclusion

Inspiration



# Conclusion

- Conclusion:
- Constant memory cost
- Lack of Intuitive Interpretation for memory update

