

# MobileVOS: Real-Time Video Object Segmentation

Contrastive Learning meets Knowledge Distillation

Roy Miles, Mehmet Kerim Yucel, Bruno Manganelli, Albert Saa-Garriga

## 1. Background

## 2. Method

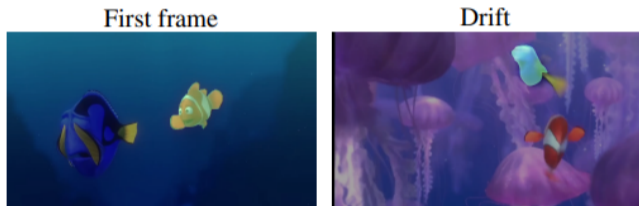
## 3. Experiments

## 4. Conclusion



Figure: Video object segmentation deals with tracking the location of objects over a video sequence. The goal is to generate accurate and consistent pixel-level object masks across the frames.

- This task serves as a foundation for many video understanding tasks such as action recognition, video editing, and augmented reality.



**Figure:** The STM memory model scales poorly for longer video sequences and introduces problems, such as drift, where the model can catastrophically degrade in performance over time.

- The current state-of-the-art relies on **space-time-memory networks** (STM), which relies on densely matching features from previous frames.

# Method - Modified Architecture

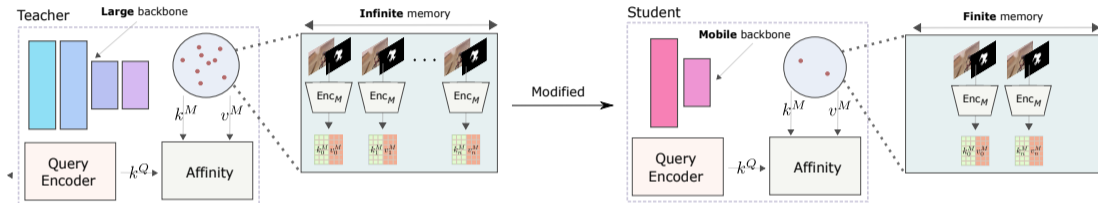


Figure: We restrict our attention to a finite STM network, whereby only the most recent and first frame features are stored in memory.

- To enable real-time performance, we use a much smaller MobileNet backbone.
- To encourage temporally consistent features, we **distill knowledge** a pre-trained infinite memory teacher.

# Method - Knowledge Distillation

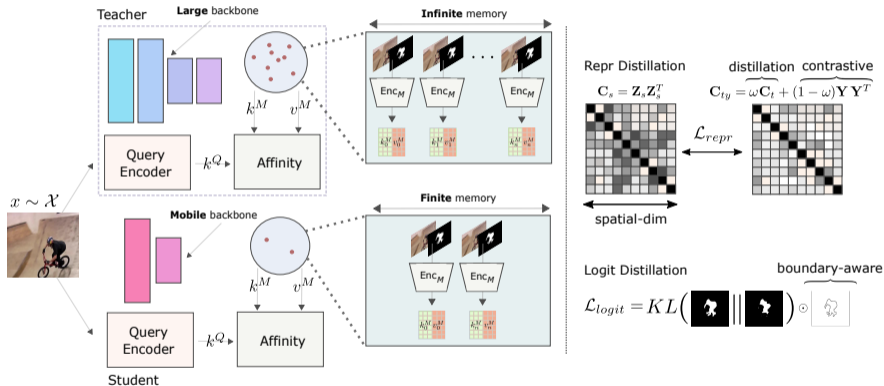


Figure: We jointly introduce **boundary aware sampling** to improve model convergence, and a natural **unification with supervised-contrastive learning** for varying student-teacher capacity gaps.

The proposed distillation loss is given as follows.

$$\mathcal{L}_{repr} = \frac{1}{|\mathbf{C}_s|} \left( \log_2 \|\mathbf{C}_s\|^2 - \log_2 \|\mathbf{C}_s \odot \mathbf{C}_t\|^2 \right) \quad (1)$$

where  $\mathbf{C}_s, \mathbf{C}_t \in \mathbb{R}^{HW \times HW}$  capture the relationship between all pairs of pixels in the student and teacher feature space respectively.

## Boundary-aware sampling

Sampling the boundary pixels not only improves model convergence and addresses observed limitations of SVOS models, but also significantly reduces the memory constraints in constructing these matrices.

By introducing known relationships between the pixel-wise features, we can provide a natural scheme to interpolate between knowledge distillation and supervised contrastive learning.

$$\mathbf{C}_{ty} = \omega \mathbf{C}_t + (1 - \omega) \mathbf{Y} \mathbf{Y}^T \quad (2)$$

In the case where  $\omega = 0$ , we arrive at a familiar supervised contrastive setting.

$$\mathcal{L}_{repr} \rightarrow \mathcal{L}_{SupCon} = -\frac{1}{|\mathbf{C}_s|} \log_2 \sum_i \frac{\sum_{j \in \mathcal{P}_i} \text{sim}(\mathbf{Z}_i, \mathbf{Z}_j)}{\sum_k \text{sim}(\mathbf{Z}_i, \mathbf{Z}_k)} \quad (3)$$



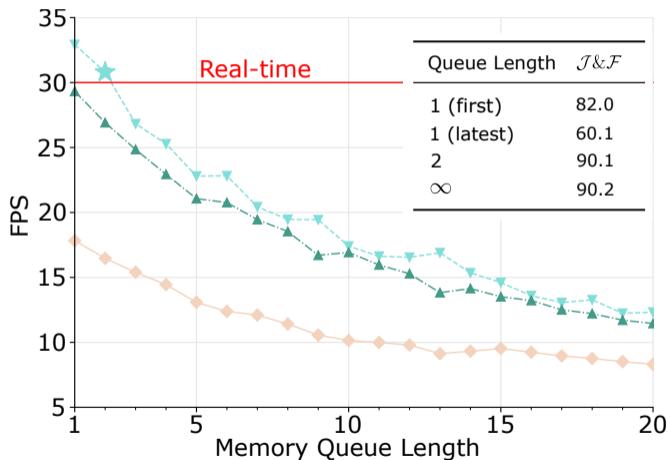
Method	CC	$\mathcal{J}\&\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$	FPS
STM <sup>†</sup>	✗	89.3	88.7	89.9	6.3
MiVOS <sup>†*</sup>	✗	91.0	89.7	92.4	16.9
STCN <sup>†*</sup>	✗	91.7	<u>90.4</u>	93.0	<u>26.9</u>
BATMAN	✗	<b>92.5</b>	<b>90.7</b>	<b>94.2</b>	-
XMem <sup>†*</sup>	✗	92.0	<b>90.7</b>	<u>93.2</u>	<b>29.6</b>
SwiftNet <sup>†</sup>	✓	90.4	<b>90.5</b>	90.3	25.0
RDE-VOS <sup>†*</sup>	✓	<b>91.6</b>	90.0	<b>93.2</b>	35.0
MobileVOS					
ResNet18 <sup>†*</sup>	✓	<u>91.4</u>	<u>90.3</u>	<u>92.6</u>	<b>100.1</b>
MobileNetV2 <sup>†</sup>	✓	90.5	89.5	91.5	81.8
↳ wo/ ASPP <sup>†</sup>	✓	90.1	89.0	91.1	<u>86.0</u>

Method	CC	$\mathcal{J}\&\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$	FPS
STM <sup>†</sup>	✗	81.8	79.2	84.3	10.2
STCN <sup>†*</sup>	✗	85.3	82.0	88.6	<u>20.2</u>
BATMAN	✗	<u>86.2</u>	<u>83.2</u>	89.4	-
XMem <sup>†*</sup>	✗	<b>87.7</b>	<b>84.0</b>	<b>91.4</b>	<b>22.6</b>
SwiftNet <sup>†</sup>	✓	81.1	78.3	83.9	<25.0
RDE-VOS <sup>†*</sup>	✓	<b>86.1</b>	<u>82.1</u>	<b>90.0</b>	27.0
MobileVOS					
ResNet18 <sup>†*</sup>	✓	85.0	81.7	88.3	<b>90.6</b>
MobileNetV2 <sup>†</sup>	✓	82.2	78.7	85.7	79.1
↳ wo/ ASPP <sup>†</sup>	✓	81.8	78.3	85.3	<u>81.3</u>

# Experiments - Mobile Performance

Samsung Research

Method	Params(M)	FPS <i>NVIDIA A40</i>		FPS <i>NVIDIA 1080Ti</i>	
		short	long (10×)	short	long (10×)
STM	38.9	8.9	4.3	6.8	<b>X</b>
GSFM	67.0	18.4	4.2	7.6	<b>X</b>
STCN	54.4	37.4	8.3	18.1	<b>X</b>
RDE-VOS	64.0	32.0	34.2	14.4	14.1
XMem	62.2	38.6	39.9	12.6	12.7
MobileVOS					
ResNet18	8.1	<b>144.7</b>	<b>145.4</b>	<b>76.0</b>	<b>76.3</b>
MobileNetV2	2.5	99.9	99.1	61.6	60.6
↳ wo/ ASPP	<b>1.9</b>	105.1	103.4	66.8	67.4



- Shown that finite memory STM networks are an efficient class of models for mobile device inference.
- Provide a natural unification of supervised contrastive learning and KD.
- Introduce boundary-aware sampling as a task specific trick for improving model convergence and memory constraints for SVOS distillation.