

READMem: Robust Embedding Association for a Diverse Memory in Unconstrained Video Object Segmentation

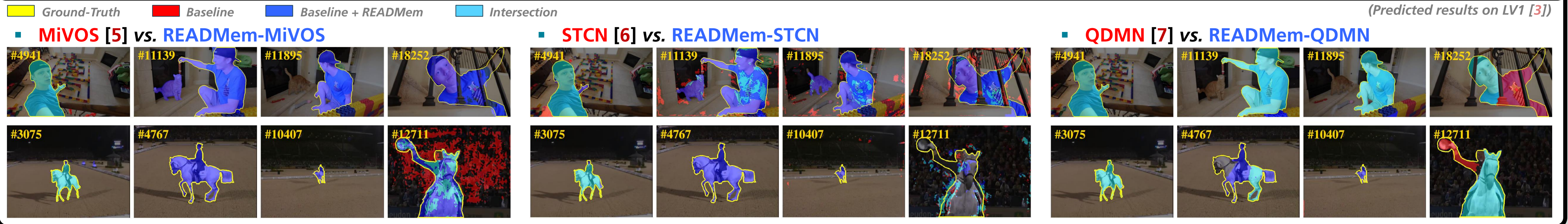
Stéphane Vujasinović¹, Sebastian Bullinger¹, Stefan Becker¹, Norbert Scherer-Negenborn¹, Michael Arens¹ and Rainer Stiefelhagen²



¹Fraunhofer IOSB* (Ettlingen, Germany)

²Karlsruhe Institute of Technology (KIT) (Karlsruhe, Germany)

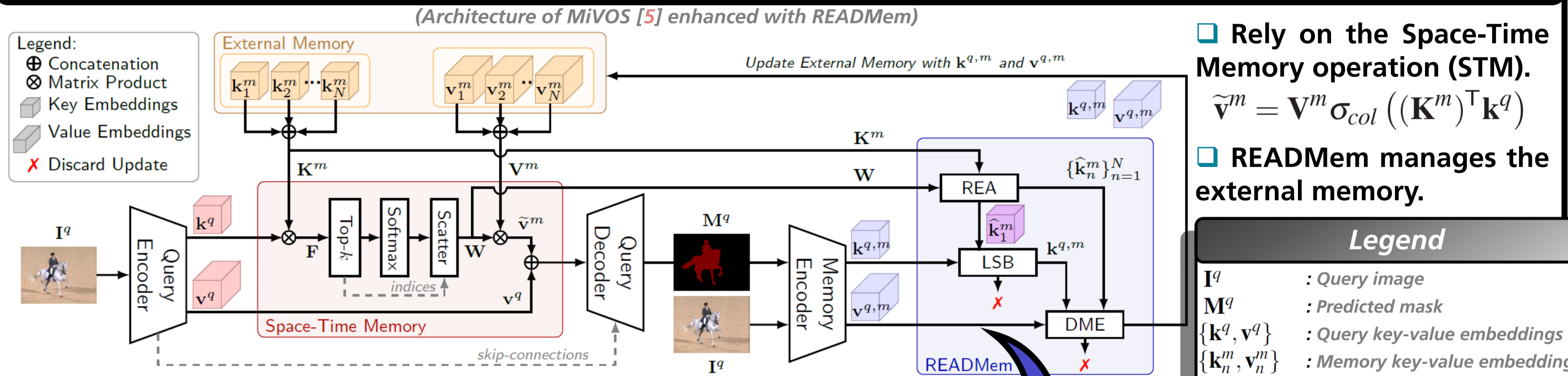
Qualitative Results



TL;DR

- Task:** Semi-automatic Video Object Segmentation (sVOS).
- Objective:** Handle unconstrained videos (arbitrary frame rate, length, object motion and camera motion).
- Approach:** Increase the inter-frame diversity within the memory.
- Implementation:** An extension to manage the memory of STM-like networks [1] available at <https://github.com/Vujas-Eteph/READMem>

Overview

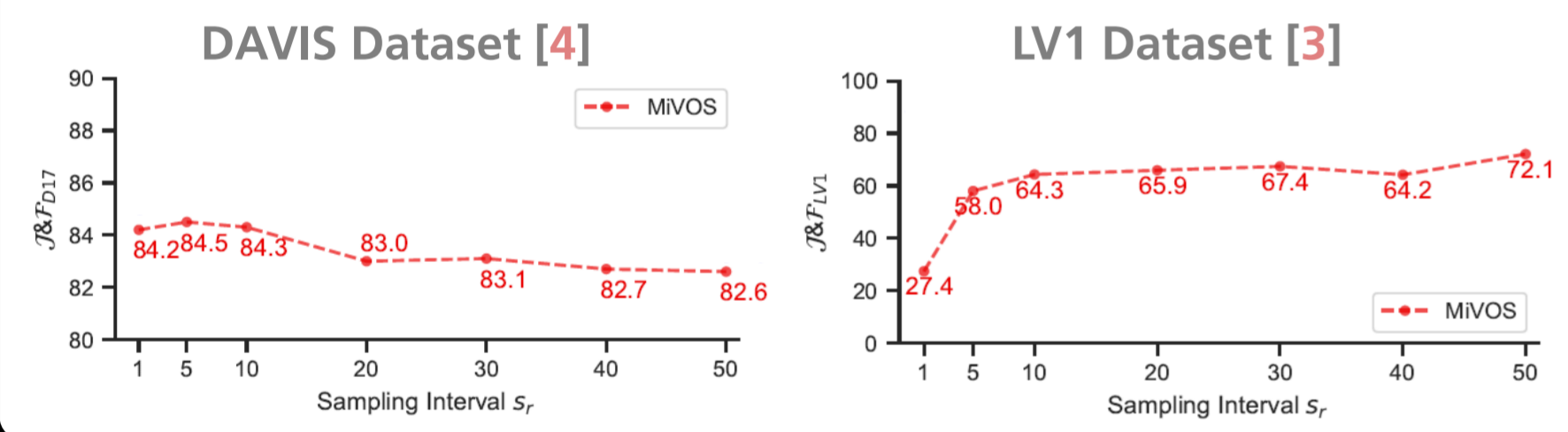


- Rely on the Space-Time Memory operation (STM).
 $\tilde{V}^m = V^m \sigma_{col}((K^m)^T k^q)$
- READMem manages the external memory.

Legend	
I^q	: Query image
M^q	: Predicted mask
$\{k^q, v^q\}$: Query key-value embeddings
$\{k_n^m, v_n^m\}$: Memory key-value embeddings
n	: Index of the memory slot
N	: Total number of memory slots
$\{K^m, V^m\}$: Concatenated embeddings
F	: Similarity Matrix
W	: Affinity Matrix
σ_{col}	: Softmax along the column
$\{k^{q,m}, v^{q,m}\}$: Memory key-value embeddings of the query image
\hat{k}_n^m	: Projected key embeddings
\tilde{v}^m	: Aggregated memory values
T_n	: Transformation Matrix

Problem & Motivation

- Existing sVOS methods are tailored for short-video object segmentation:
 - Rely on an ever-expanding memory.
 - Define a dataset-specific sampling interval s_r , non-generalizable to unseen data.



READMem

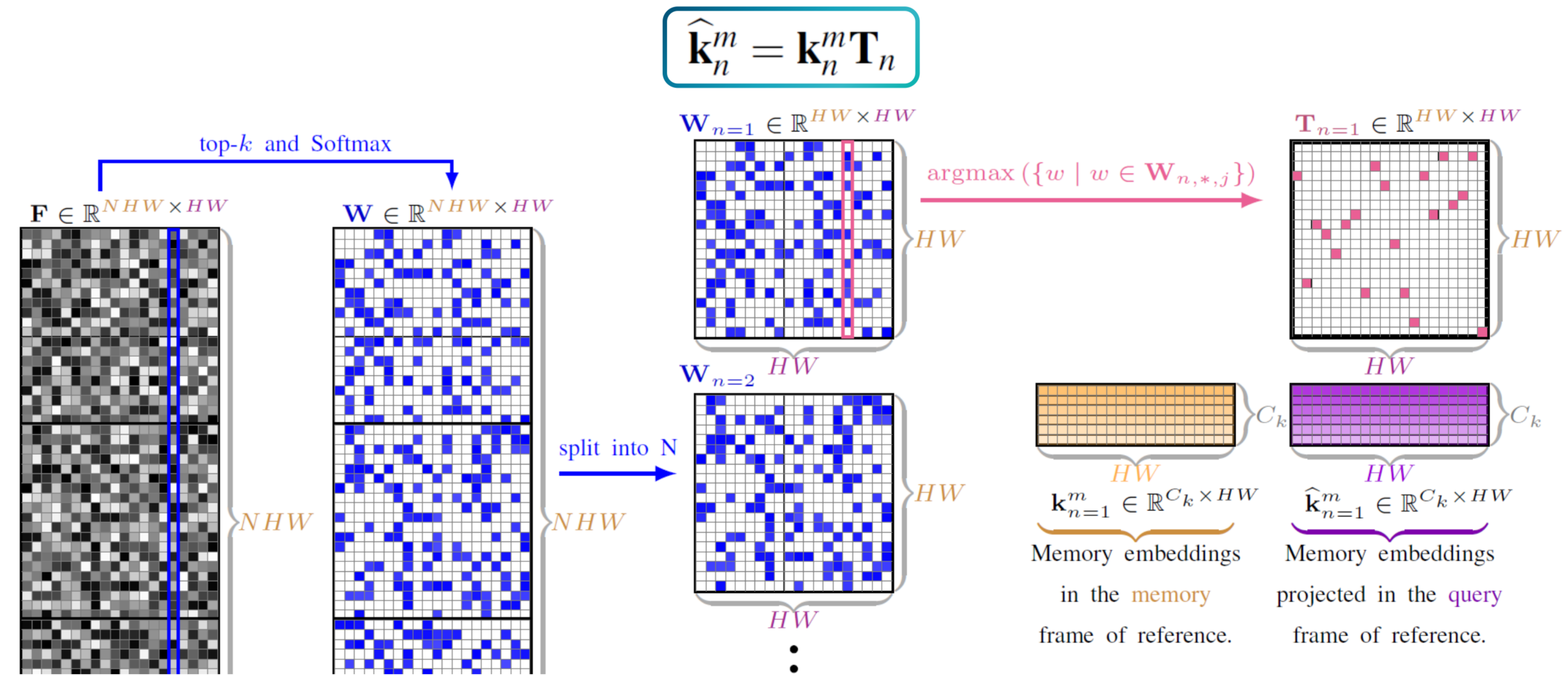
- Diversification of Memory Embeddings (DME):**
 - Motivation:** Increase the diversity of the memory $\mathcal{M}^k = \{k_n^m\}_{n=1}^N$, where $k_n^m \in \mathbb{R}^{C_k \times HW}$
 - Concept:** Represent the diversity of the memory \mathcal{M}^k with the determinant of the associated Gram Matrix by:

$$G \in \mathbb{R}^{N \times N}, \text{ where } X \in \mathbb{R}^{C_k \times HW \times N} \text{ is the flattened representation of } \mathcal{M}^k = \{k_n^m\}_{n=1}^N$$

$$(\text{vol}P(\{k_n^m\}_{n=1}^N))^2 = \det(G), \text{ where } G = X^T X$$

- Process:** Memory keys updated if diversity increases.
- Robust Embedding Association (REA):**
 - Motivation:** Dampen translation and scale variations when computing the similarity.
 - Concept:** Projects the memory key embeddings k_n^m to the query's ($k^{q,m}$) temporal frame of reference by:

$$\hat{k}_n^m = k_n^m T_n$$



- Process:** Compute the similarity between \hat{k}_n^m and $k^{q,m}$ instead of k_n^m .

Contributions

- READMem extends any sVOS method to deal with unconstrained sequences.
- Seamless integration without re-training or fine-tuning.
- Generalizable sampling interval s_r (no need to finetune on the validation set).
- Automatic memory embeddings diversity estimation via Diversification of Memory Embeddings (DME).
- Translation and scale invariant embedding association through Robust Embedding Association (REA).

Quantitative Results

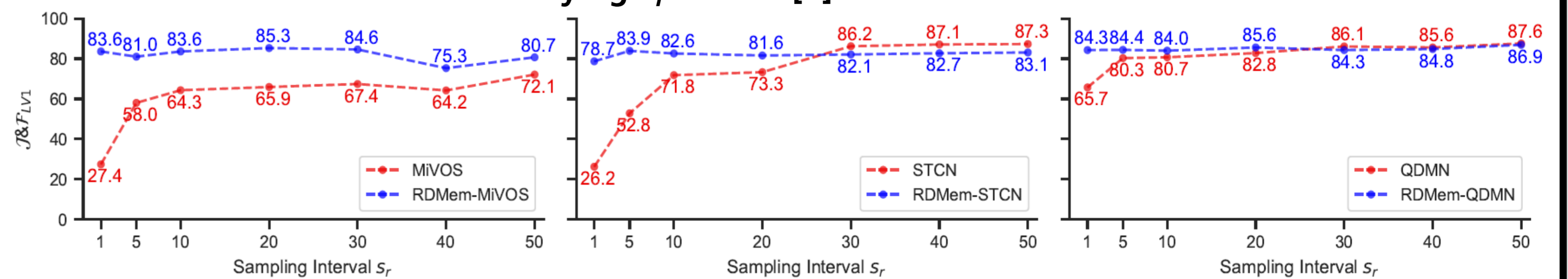
Results on Long-Video [3] and DAVIS [4] with $s_r = 1$.

Method	LV1 [3] (Long Sequences)			D17 [4] (Short Sequences)		
	J&F	J	F	J&F	J	F
MiVOS [5] (CVPR 21)	27.4	27.4	27.4	84.3	81.3	87.3
READMem-MiVOS (ours)	83.6 ^{†56.2}	82.3 ^{†54.9}	85.0 ^{†57.6}	84.3	81.4 ^{†0.1}	87.1 ^{†0.2}
STCN [6] [†] (NIPS 21)	26.2	22.9	29.4	83.2	79.9	86.5
READMem-STCN [†] (ours)	80.8 ^{†54.6}	78.4 ^{†55.5}	83.2 ^{†53.8}	83.8 ^{†0.6}	80.4 ^{†0.5}	87.2 ^{†0.7}
QDMN [7] (ECCV 22)	65.7	63.5	67.9	85.1	82.2	88.0
READMem-QDMN (ours)	84.3 ^{†18.6}	81.9 ^{†18.4}	86.7 ^{†18.8}	85.3 ^{†0.2}	82.4 ^{†0.2}	88.1 ^{†0.1}

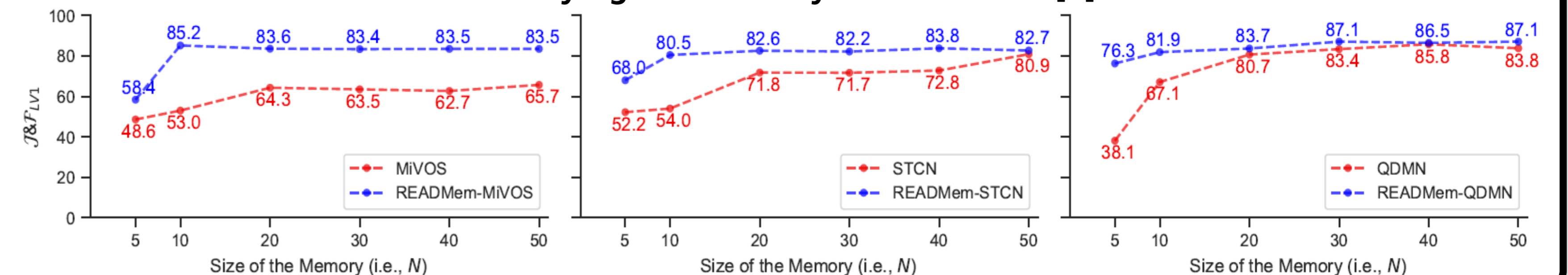
Ablation Study.

Configuration	J&F _{LV1}	J&F _{D17}
MiVOS [5] + adj. frame (baseline)	64.3	84.3
MiVOS [5] + DME (ours)	69.5 ^{†5.2}	81.3 ^{†3.0}
MiVOS [5] + DME + LSB (ours)	75.0 ^{†10.7}	81.3 ^{†3.0}
MiVOS [5] + DME + LSB + adj. frame (ours)	77.4 ^{†13.1}	84.3 ^{†0.0}
MiVOS [5] + DME + LSB + adj. frame + REA (ours)	86.0 ^{†21.7}	84.6 ^{†0.3}

Performance variation when varying s_r on LV1 [3].



Performance variation when varying the memory size N on LV1 [3].



[1] Seoung Wug Oh, et al. Video Object Segmentation using Space-Time Memory Networks. ICCV, 2019.
 [2] Tianfei Zhou, et al. A Survey on Deep Learning Technique for Video Segmentation. TPAMI, 2023.
 [3] Yongqing Liang, et al. Video Object Segmentation with Adaptive Feature Bank and Uncertain-Region Refinement. NeurIPS, 2020.
 [4] Jordi Pont-Tuset, et al. The 2017 DAVIS Challenge on Video Object Segmentation. arXiv, 2017.

[5] Ho Kei Cheng, et al. Modular Interactive Video Object Segmentation: Interaction-to-Mask, Propagation and Difference-Aware Fusion. CVPR, 2021.
 [6] Ho Kei Cheng, et al. Rethinking Space-Time Networks with Improved Memory Coverage for Efficient Video Object Segmentation. NeurIPS, 2021.
 [7] Yong Liu, Ra et al. Learning Quality-Aware Dynamic Memory for Video Object Segmentation. ECCV, 2022.