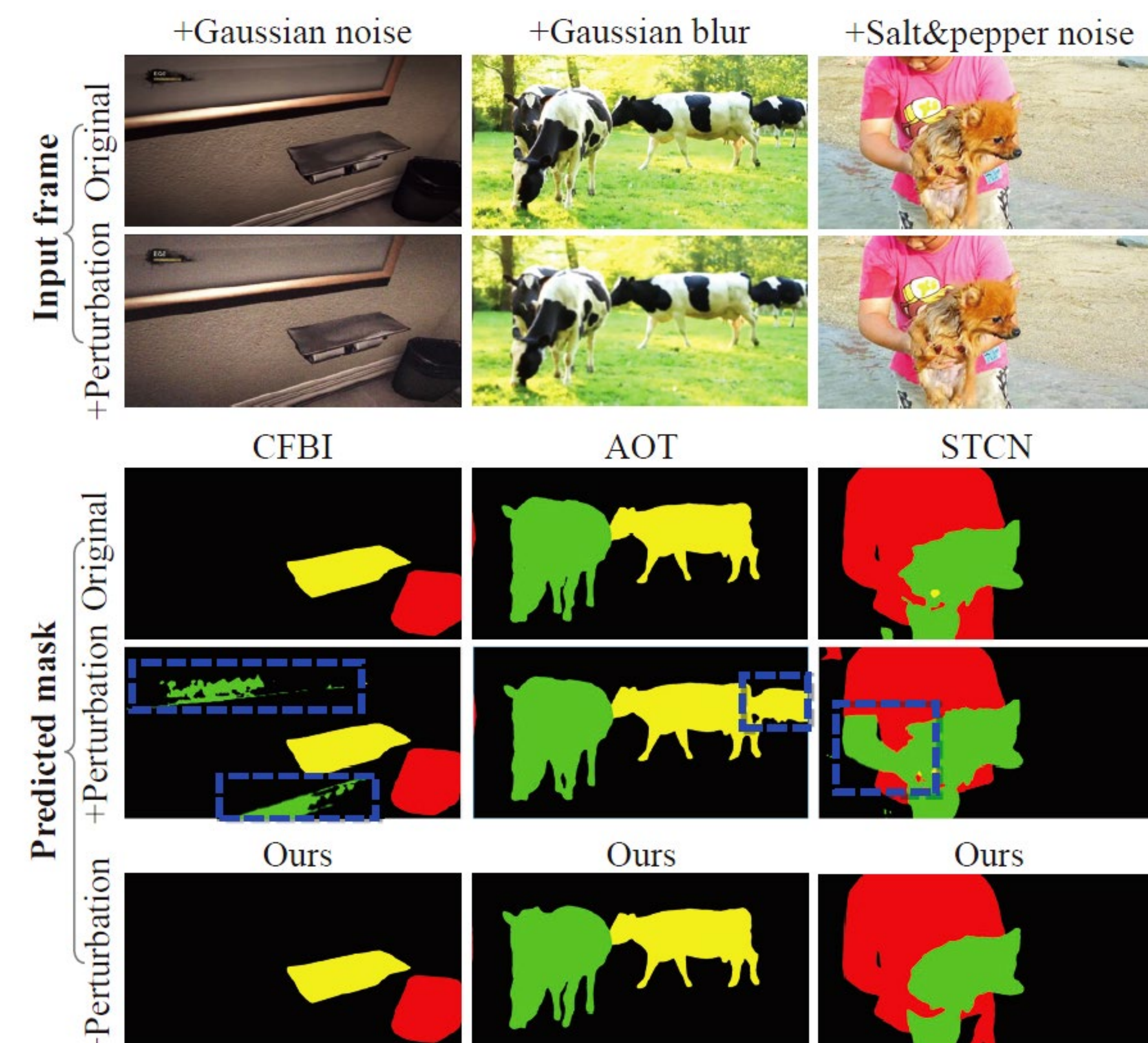


Microsoft  
**Research**  
微软亚洲研究院

## Introduction

- **Task:** Video Object Segmentation (VOS) aims at segmenting target objects in a video clip given the ground truth mask at the reference frame.
- **Motivation:** Our finding indicates that existing VOS models are fragile to natural perturbations.



Our model with adaptive object calibration shows superior robustness against perturbations.

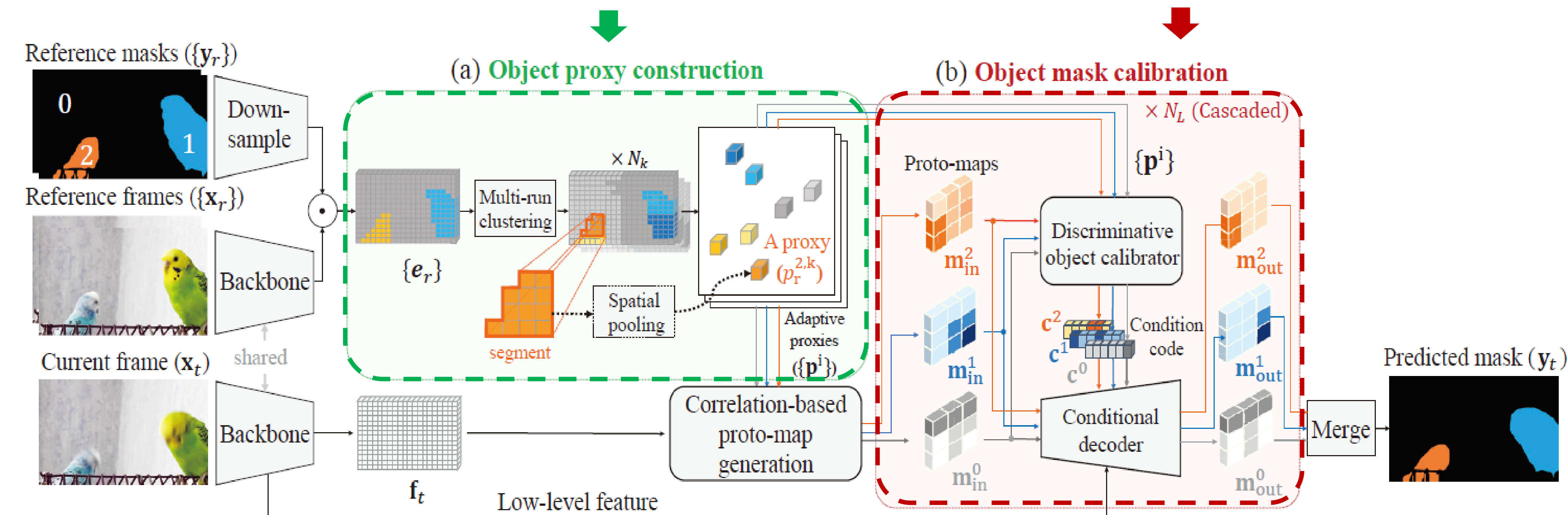
- **Solution:** The key insight is to:
  - introduce an **adaptive object proxy representation** for referenced objects robustly, which **reduces errors incurred by unstable pixel-level matching**.
  - **calibrate the object masks** by updating object representation and masks in an interleaving manner progressively, achieving **discrimination among co-existing objects**.

# Towards Robust Video Object Segmentation with Adaptive Object Calibration

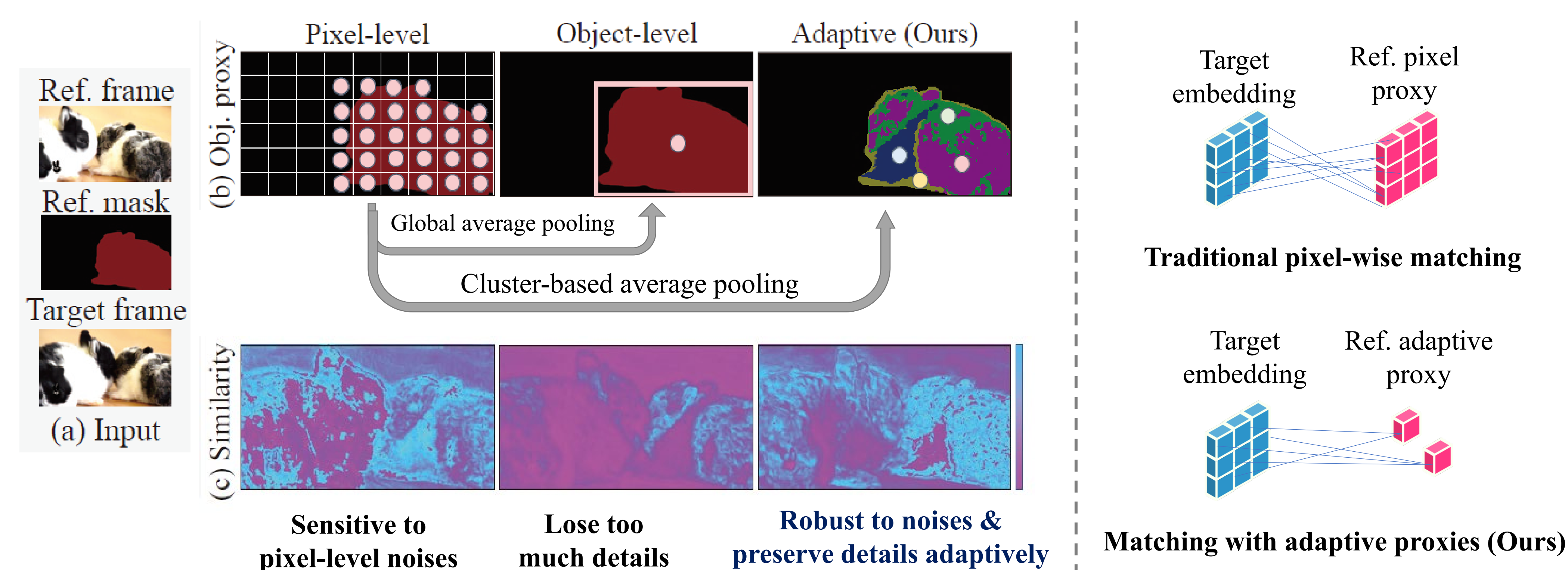
Xiaohao Xu,<sup>1,2</sup> Jinglu Wang,<sup>2</sup> Xiang Ming,<sup>2</sup> Yan Lu<sup>2</sup>  
<sup>1</sup> Huazhong University of Science & Technology  
<sup>2</sup> Microsoft Research Asia

## Method

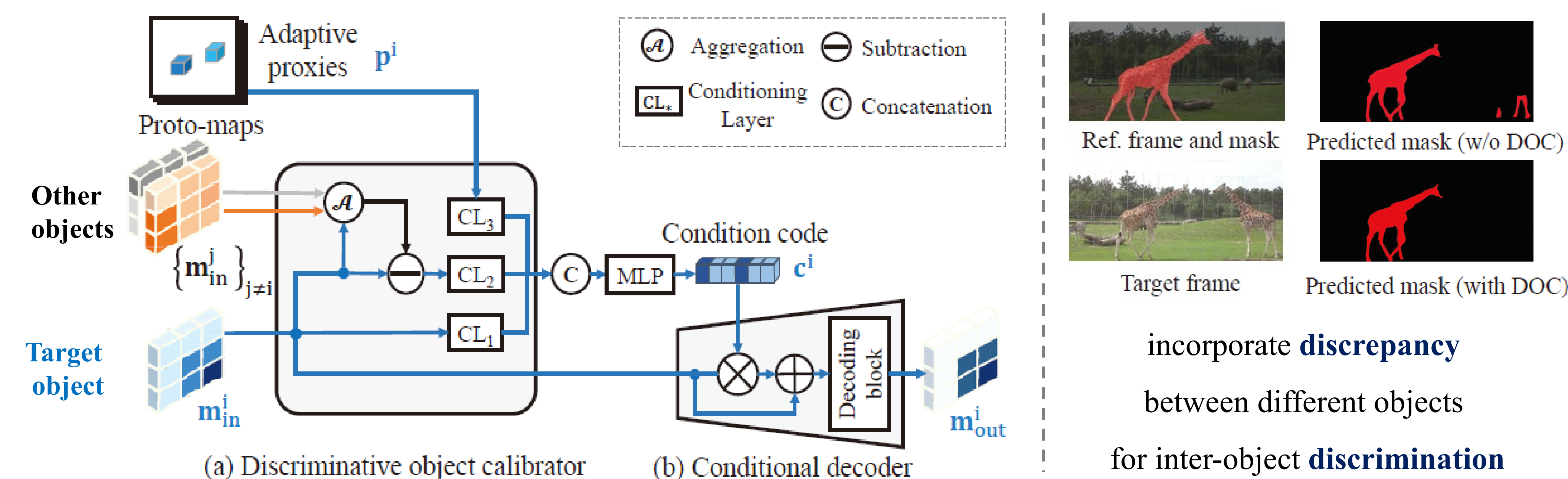
**Robust VOS Framework = Robust object representation + Robust mask decoding**



## Adaptive Object Proxy (AOP) for Robust Object Representation



## Discriminative Object Calibration (DOC) for Robust Mask Decoding



paper



repo

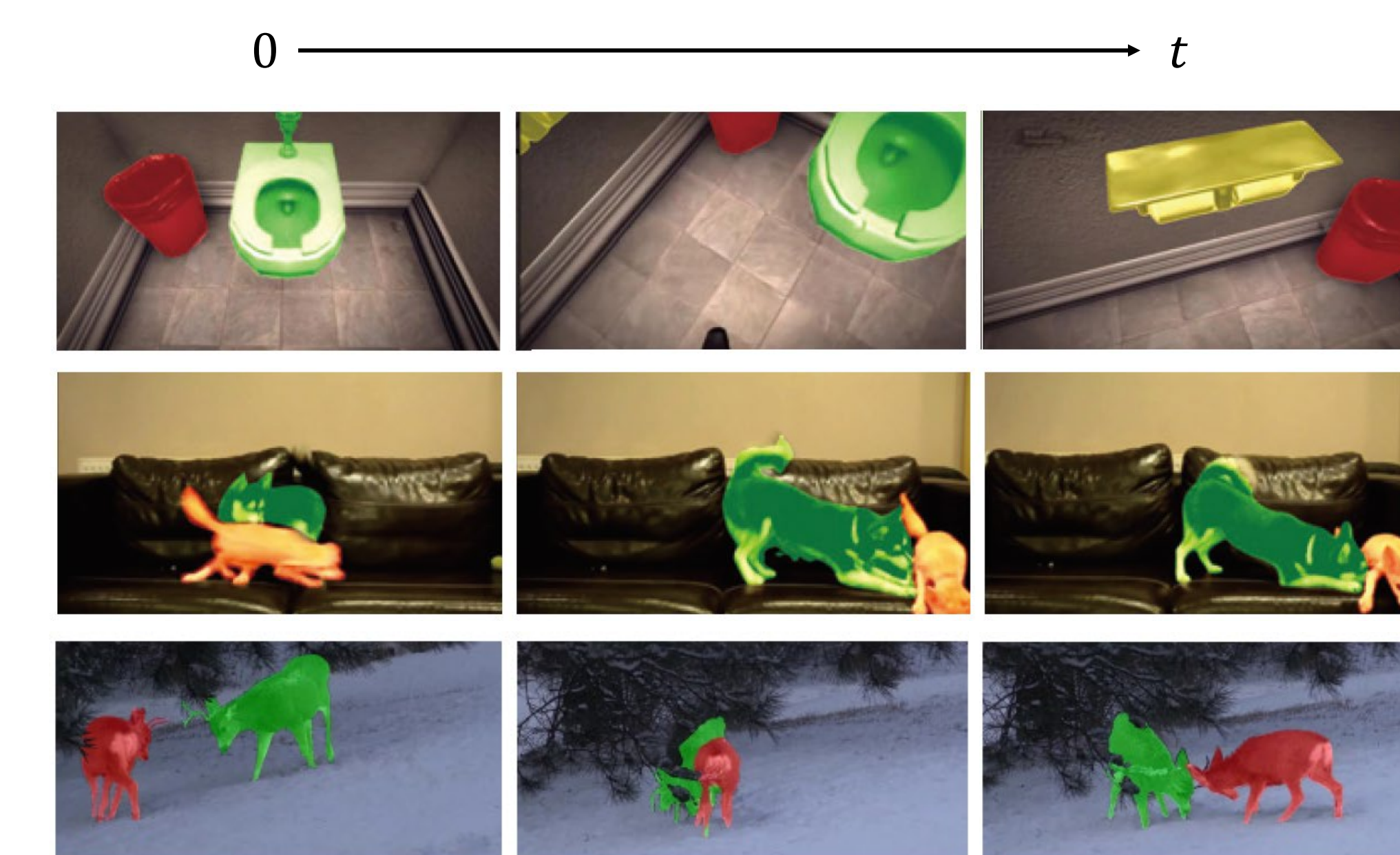
## Results

Method	YouTube-VOS18	DAVIS17 Valid	DAVIS17 Test-dev
STM	79.4	81.8	72.3
CFBI	81.4	81.9	74.8
AOT-B	83.2	82.1	75.5
<b>Ours-Base</b>	<b>83.6</b>	<b>83.1</b>	<b>76.5</b>
STCN	83.0	<b>85.4</b>	76.1
AOT-L	83.7	83.8	78.3
<b>Ours-MF</b>	<b>84.0</b>	<b>83.8</b>	<b>79.3</b>

Comparisons on standard VOS benchmarks  
w.r.t overall performance (J&F).

Method	After-perturbation accuracy (↑)	Perturbation robustness (↓)
CFBI	79.4	1.6
AOT-B	81.6	1.7
<b>Ours-Base</b>	<b>82.3</b>	<b>1.4</b>
STCN	79.7	3.0
AOT-L	81.7	1.9
<b>Ours-MF</b>	<b>82.7</b>	<b>1.4</b>

Pilot study of perturbation robustness for VOS models on the perturbed dataset YouTube-VOS-P.



Qualitative results of our model.