

Bootstrapping Objectness from Videos by Relaxed Common Fate and Visual Grouping



Long Lian

UC Berkeley



Zhirong Wu

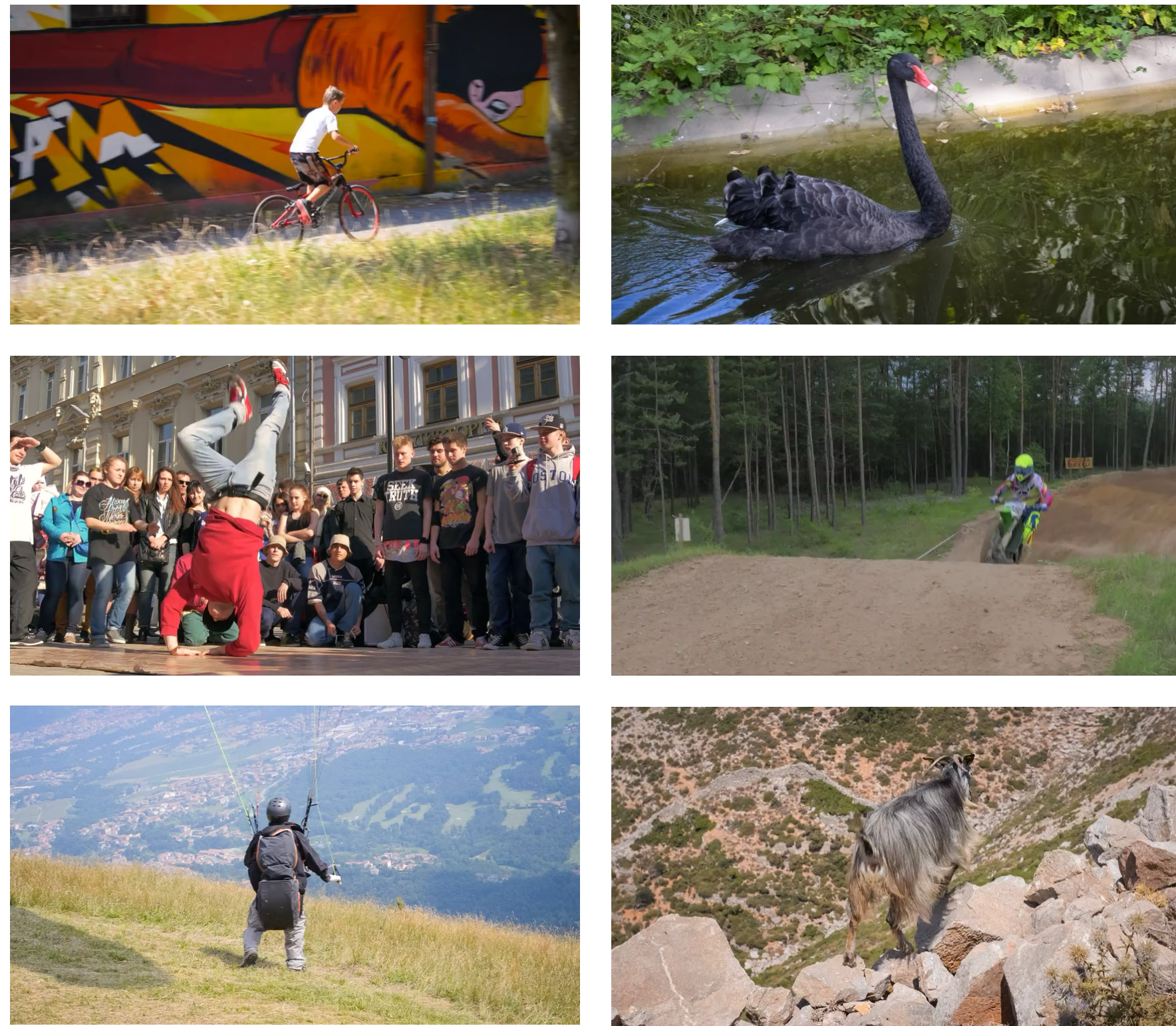
Microsoft Research Asia



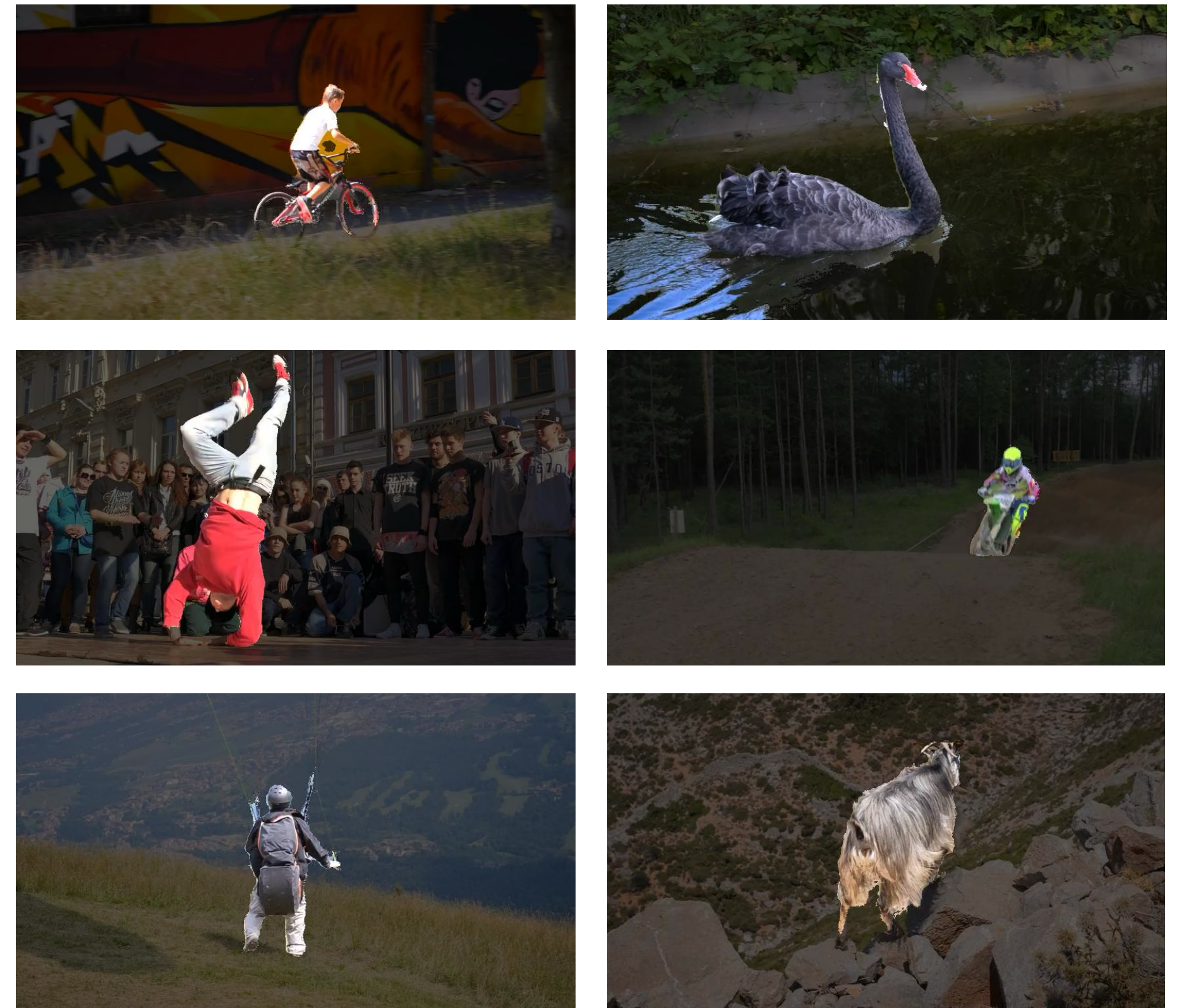
Stella X. Yu

University of Michigan

Our Task: Unlabeled Videos \Rightarrow Object Segmentation



Given optical flow detector



(Our RCF Segmentation)

Existing Methods Rely on Common Fate

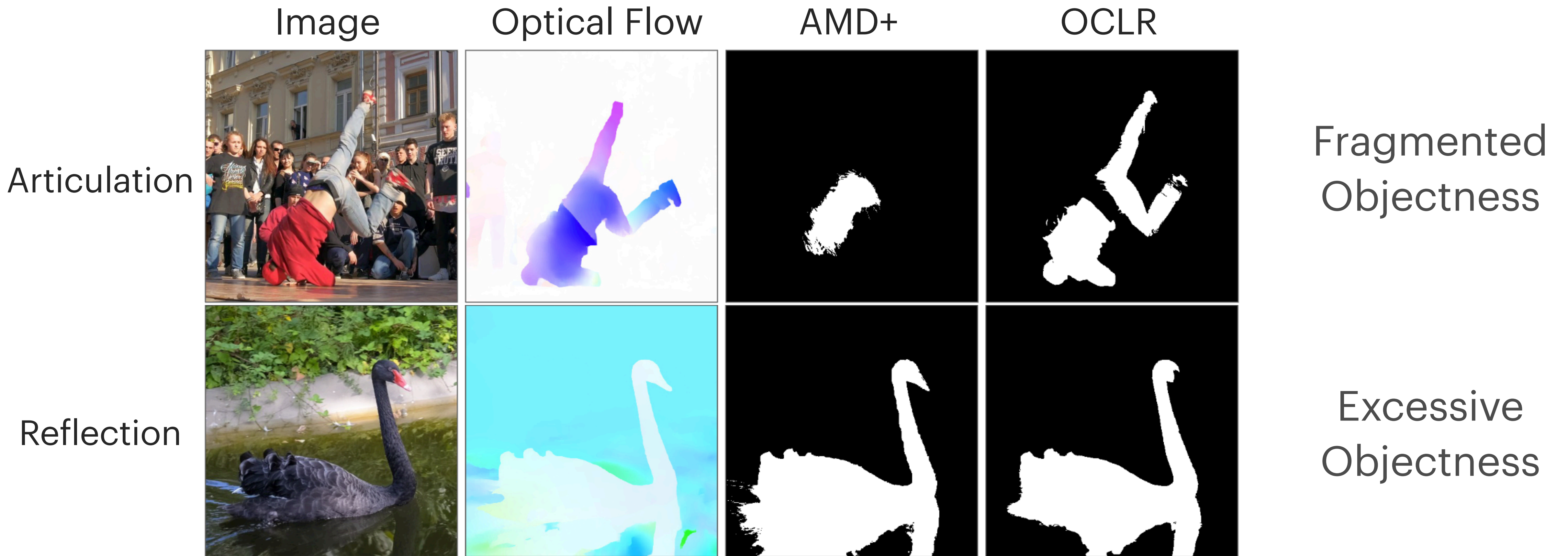


Motion Grouping. ICCV. 2021.
OCLR. NeurIPS 2022.
GWM. BMVC 2022.
AMD. NeurIPS 2021.

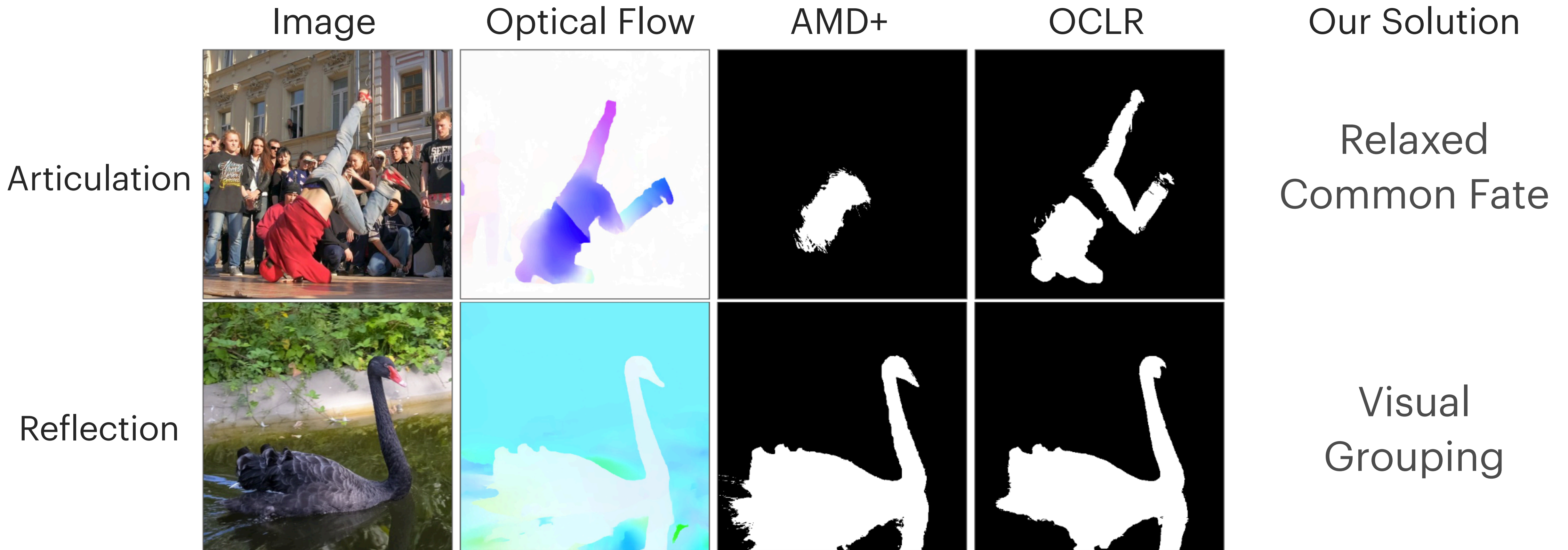
Objectness = What Move Together Belong Together



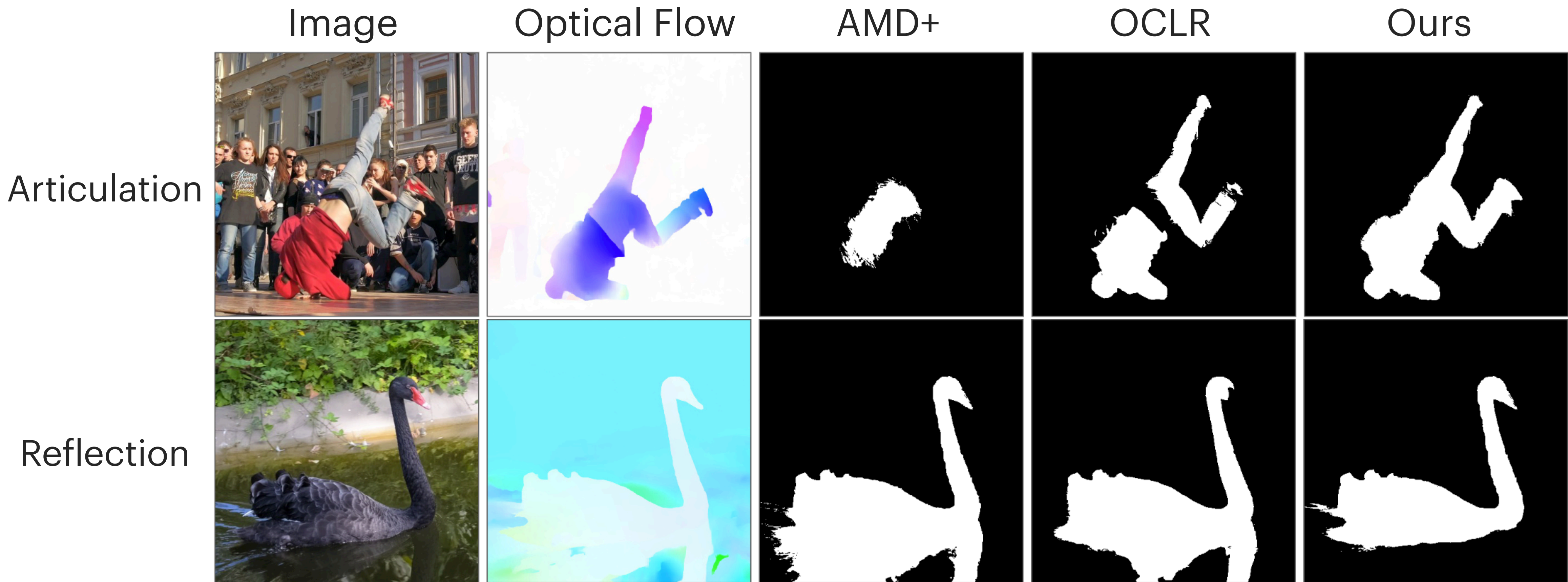
Two Failure Modes From Gestalt Law of Common Fate



Two Failure Modes From Gestalt Law of Common Fate



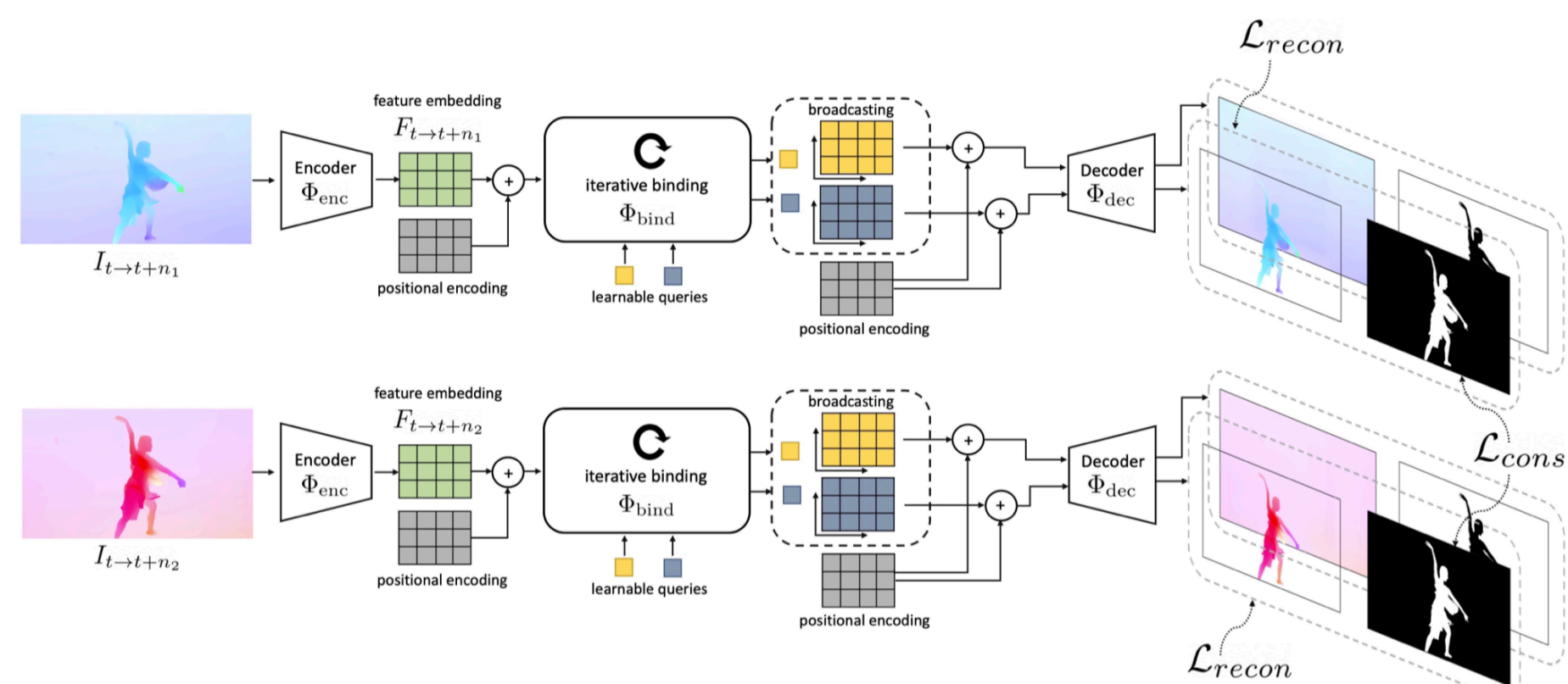
Our Approach Addresses Both Caveats



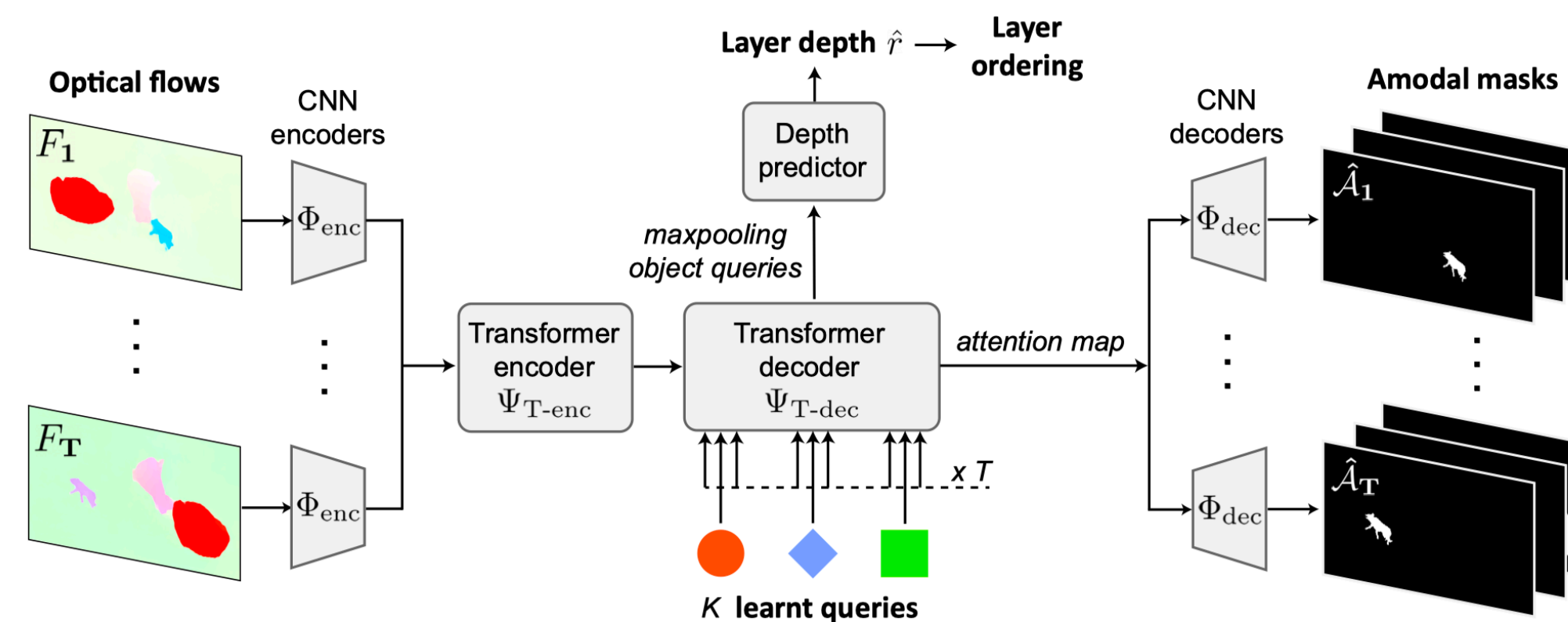
Existing Methods: Three Camps

1. Motion Segmentation

Motion Grouping

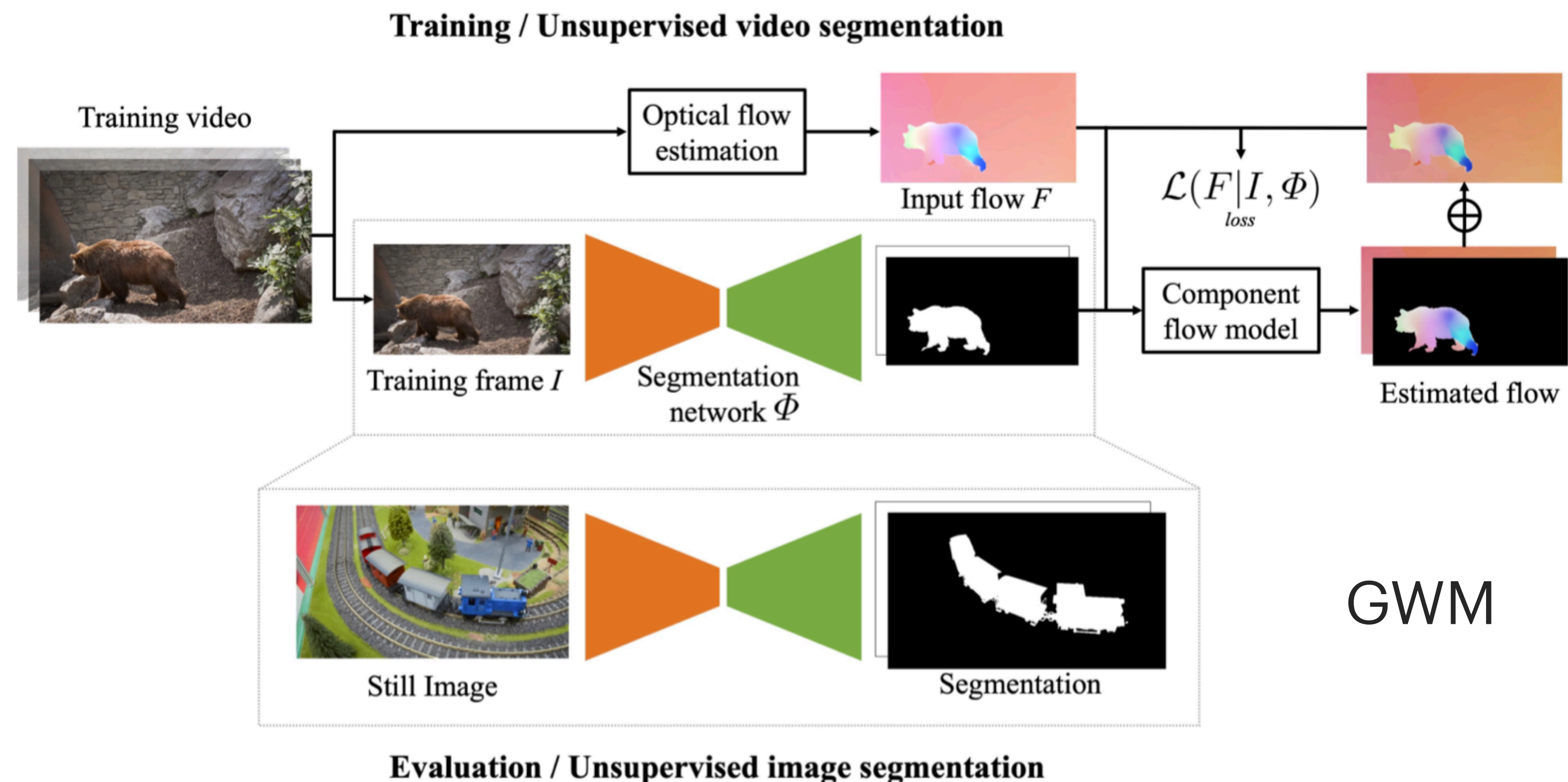


OCLR



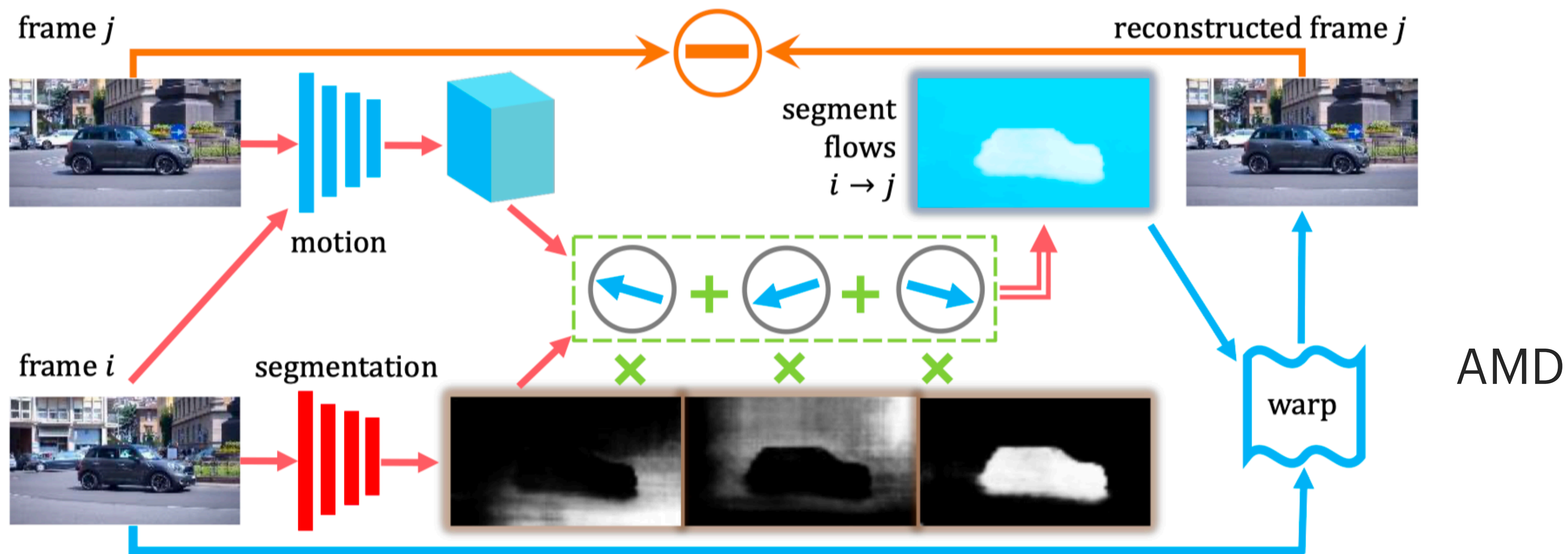
Existing Methods: Three Camps

1. Motion Segmentation
2. Motion Guided Segmentation



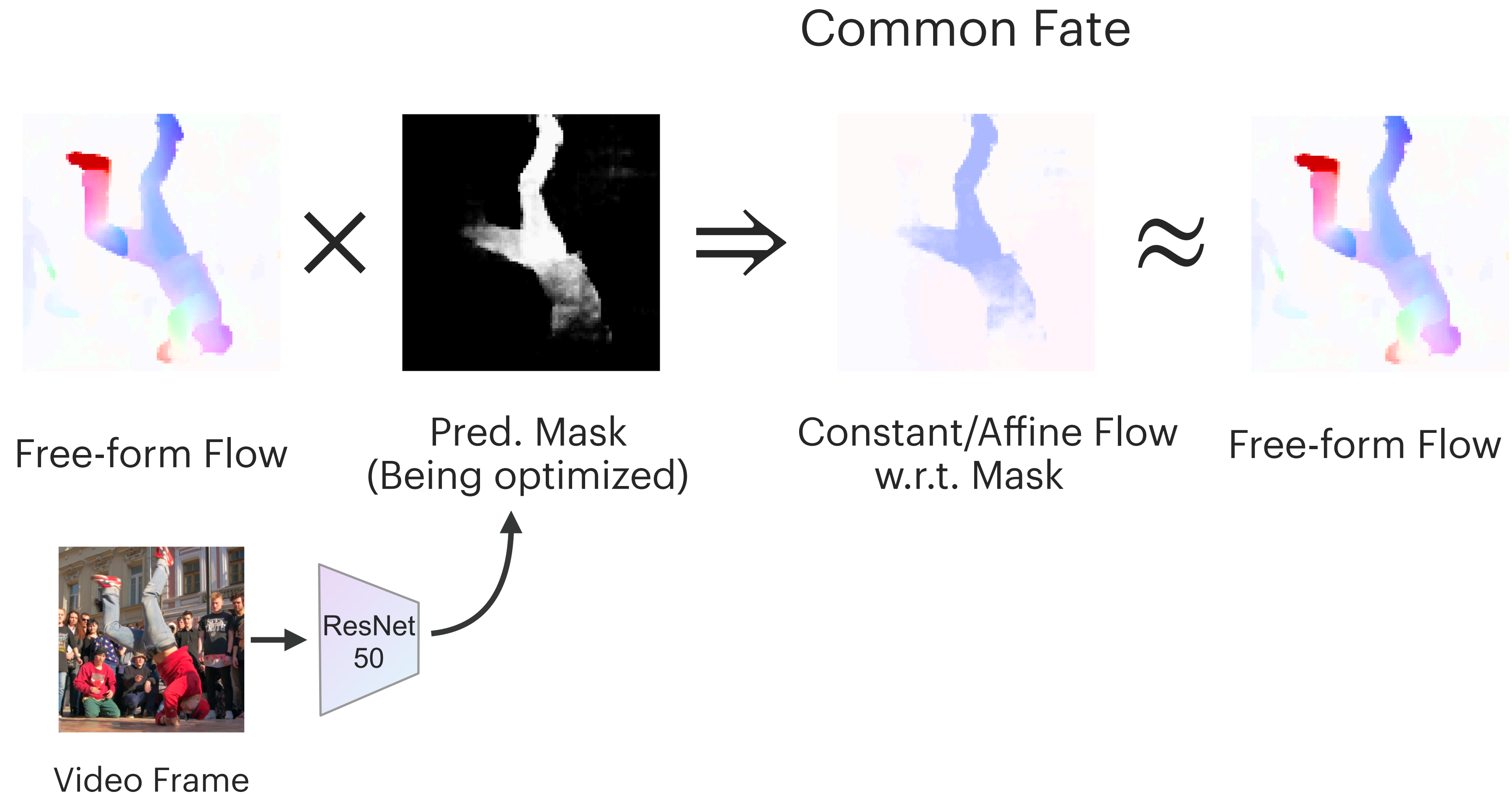
Existing Methods: Three Camps

1. Motion Segmentation
2. Motion Guided Segmentation
3. Motion and Segmentation Jointly Learned



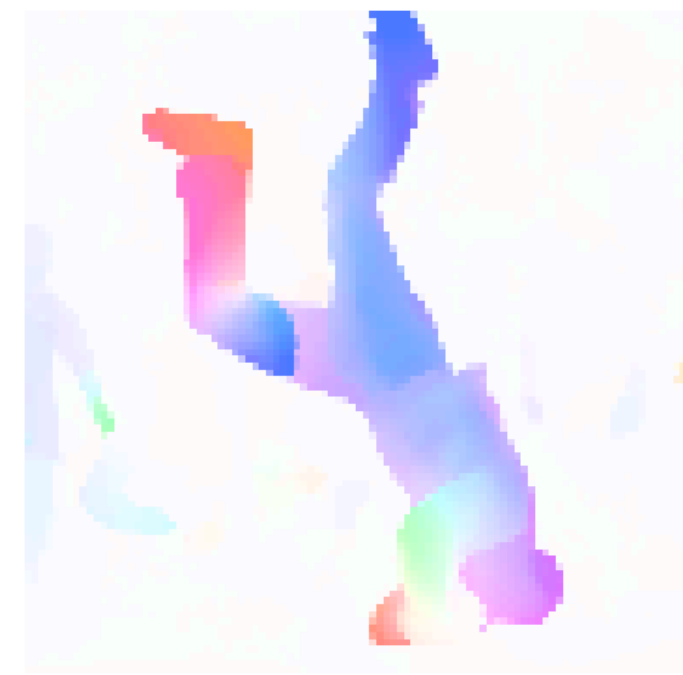
Insight 1: Dealing with Articulation

Previous works:



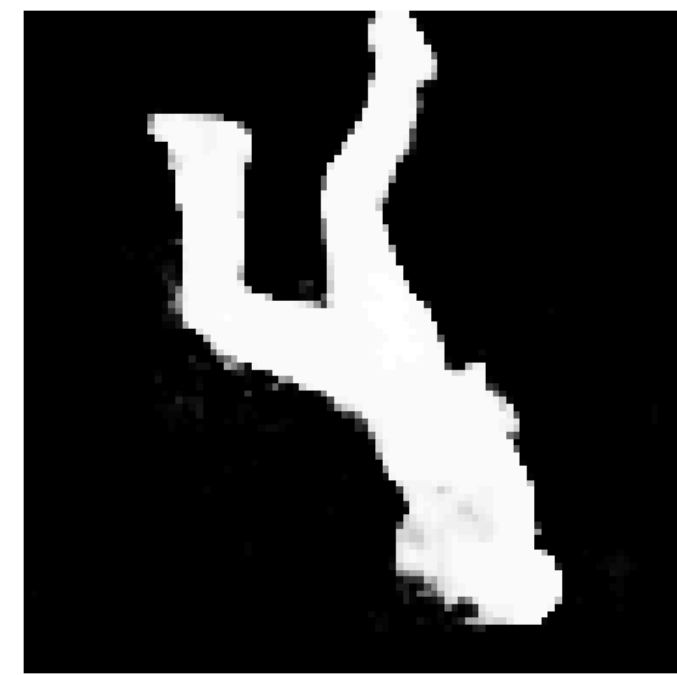
Insight 1: Fitting Flow with *Relaxed Common Fate*

Our work RCF:



Free-form Flow

\times



Pred. Mask
(Being optimized)

\Rightarrow

*Relaxed
Common Fate*



Constant/Affine Flow
w.r.t. Mask

$+$



Residual Flow
(Intra-mask motion)

\approx



Free-form Flow



Video Frame

ResNet
50



Insight 2: Visual Grouping within the Image

Let motion and appearance complement each other for supervision.

Motion Supervision Only



Bootstrapping
from CRF



Motion + Appearance



Insight 2: Visual Grouping Based on Semantics across Images

Let motion and appearance complement each other for supervision.

Motion Supervision Only

Motion + Appearance



Iteratively minimize
the normalized cuts of
DINO feature



Extra Benefits: Label-free Hyperparameter Tuning

Using the normalized cuts of DINO feature as an unsupervised segmentation quality indicator for hyperparameter tuning.



High Normalized Cuts

(Mask boundary spans coherent features)



Low Normalized Cuts

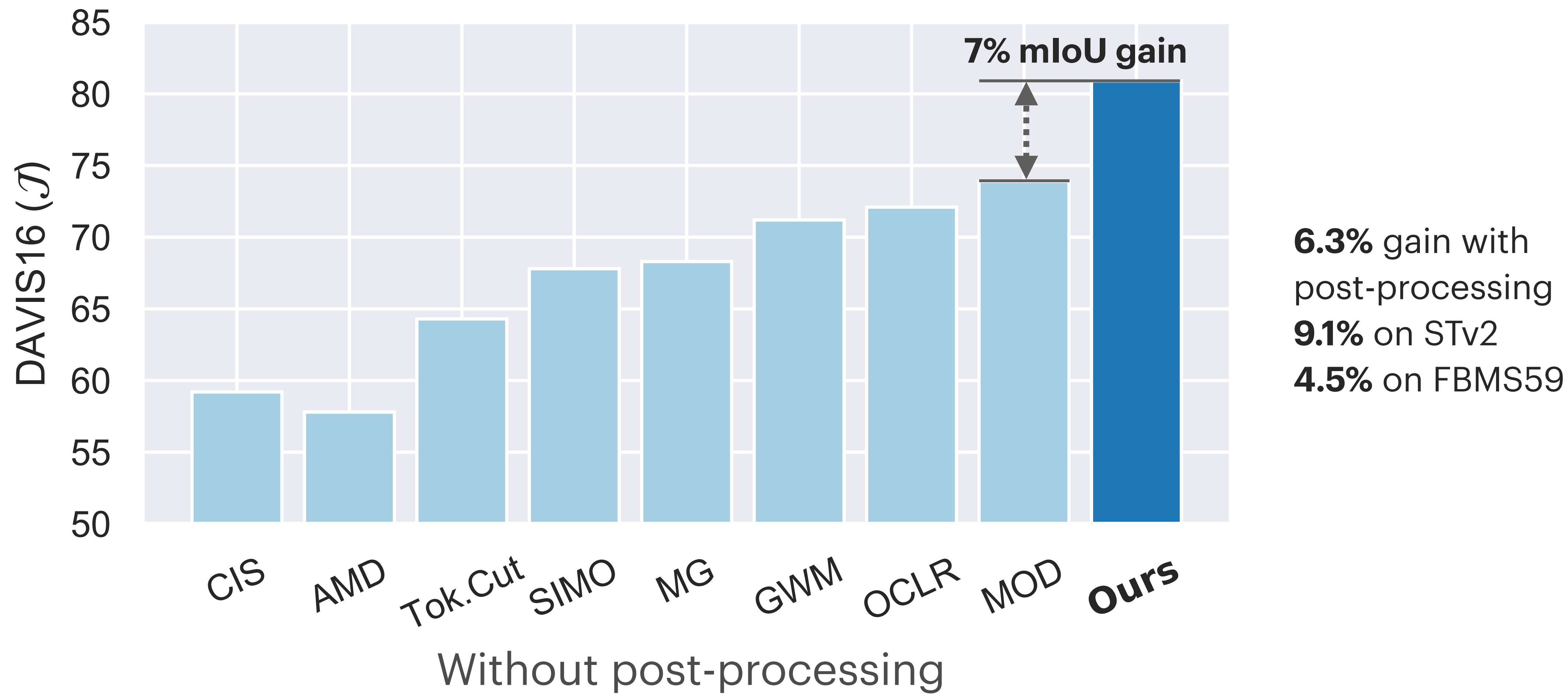
(Mask boundary does not span coherent features)

✓ Use this setting

Advantages of Our RCF to Previous Methods

UVOS Method	Motion Grouping	Emergence of Objectness	Guess What Moves	Our work (RCF)
Sources of supervision	Motion	Motion (Frame Warping)	Motion	Motion + Appearance
Segment stationary objects	✗	✓	✓	✓
Handle articulated/deformable objects	—	✗	✗	✓
Label-free hyperparameter tuning?	✗	✗	✗	✓

Our RCF: SOTA on Unsupervised Object Segmentation



OCLR



RCF (Ours)



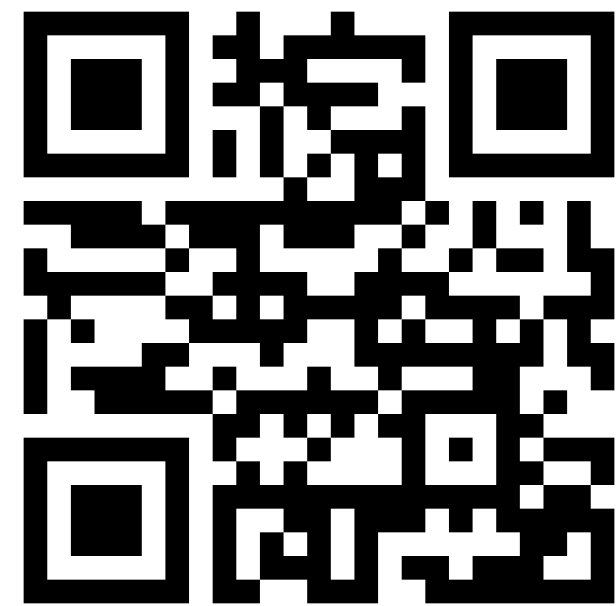
No post-processing applied: results can be further enhanced with post-processing

RCF (Ours)



No post-processing applied: results can be further enhanced with post-processing

Thank you!



Code, Model Zoo,
and Demos Available