

Weakly Supervised Instance Segmentation for Videos with Temporal Mask Consistency



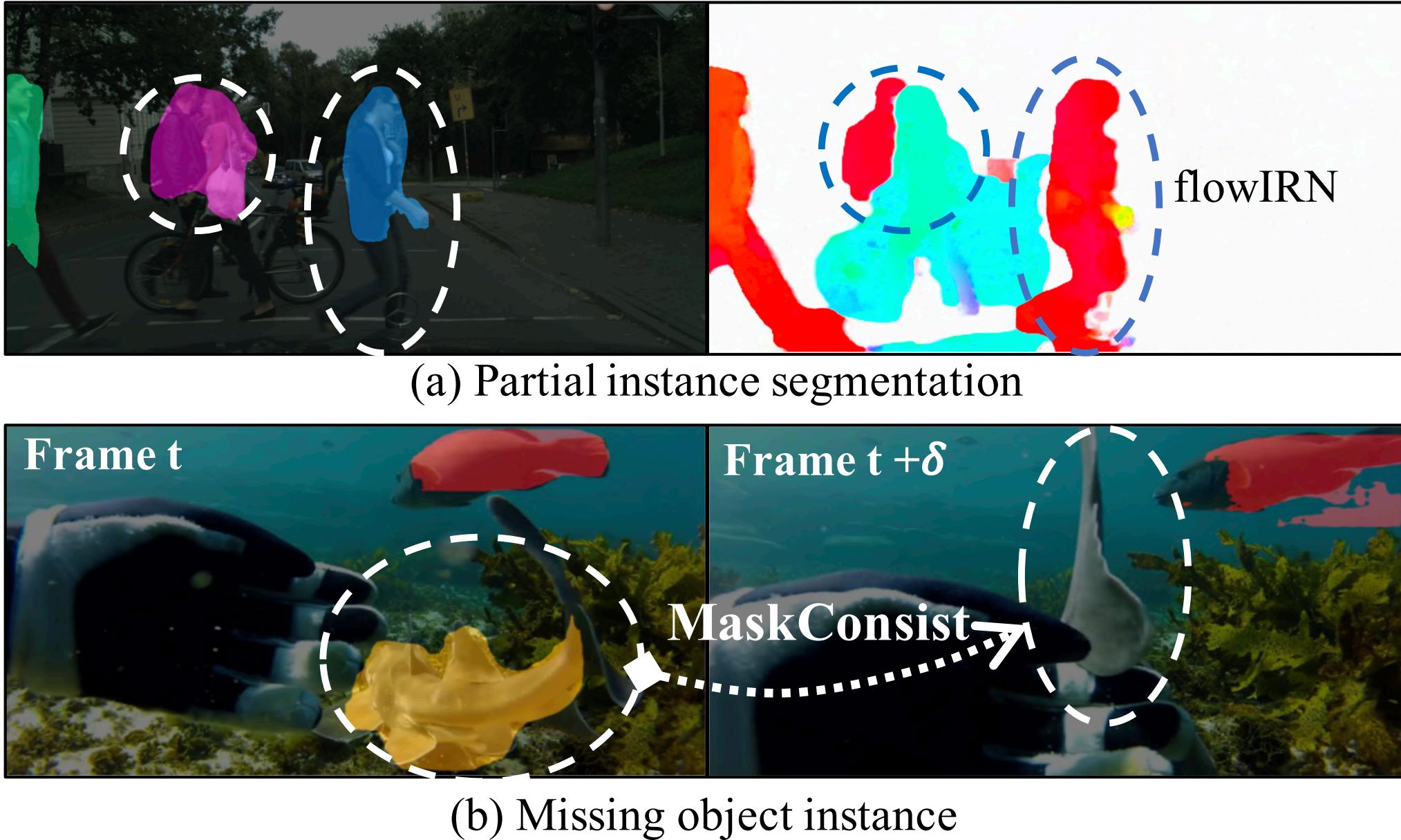
Qing Liu¹ Vignesh Ramanathan² Dhruv Mahajan² Alan Yuille¹ Zhenheng Yang²

¹ Johns Hopkins University ² Facebook



Motivation

- Problems with existing works where video signals can help:



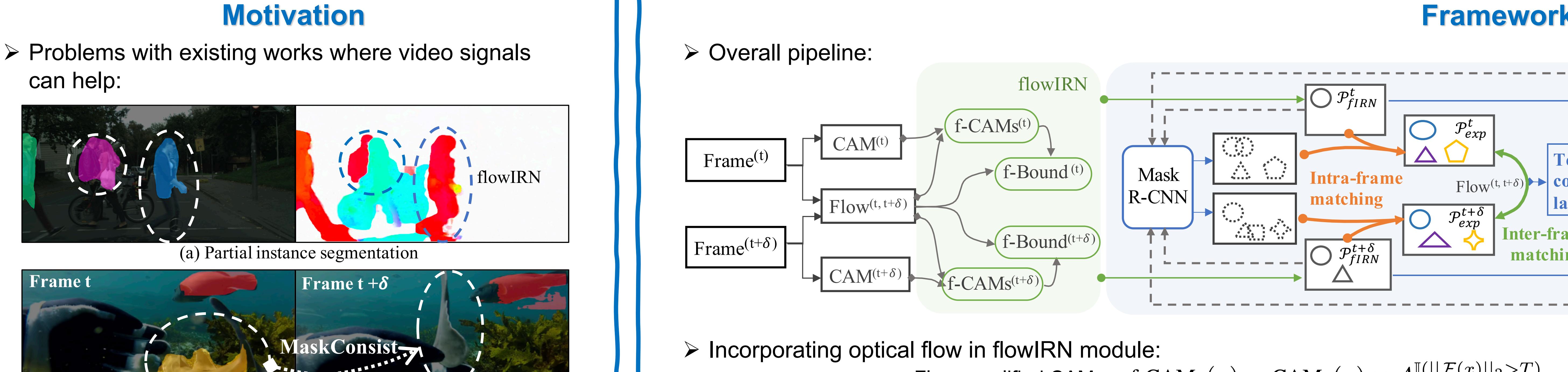
Frame Instance Segmentation Results

Methods	Video Info	Supervision	AP ₅₀
Mask R-CNN [17]	<input checked="" type="checkbox"/>	Mask	78.24
WSIS-BBTP [20]	<input checked="" type="checkbox"/>	Bbox	46.80
WISE [27]	<input checked="" type="checkbox"/>	Class	24.54
F2F [29]+MCG [41]	<input checked="" type="checkbox"/>	Class	26.31
IRN [6]	<input checked="" type="checkbox"/>	Class	29.64
IRN [6]+F2F[29]	<input checked="" type="checkbox"/>	Class	30.27
Ours	<input checked="" type="checkbox"/>	Class	34.66
Ours (self-training)	<input checked="" type="checkbox"/>	Class	36.00

Table 1. Frame-level instance segmentation performance (AP₅₀) on YTVIS train_val split.

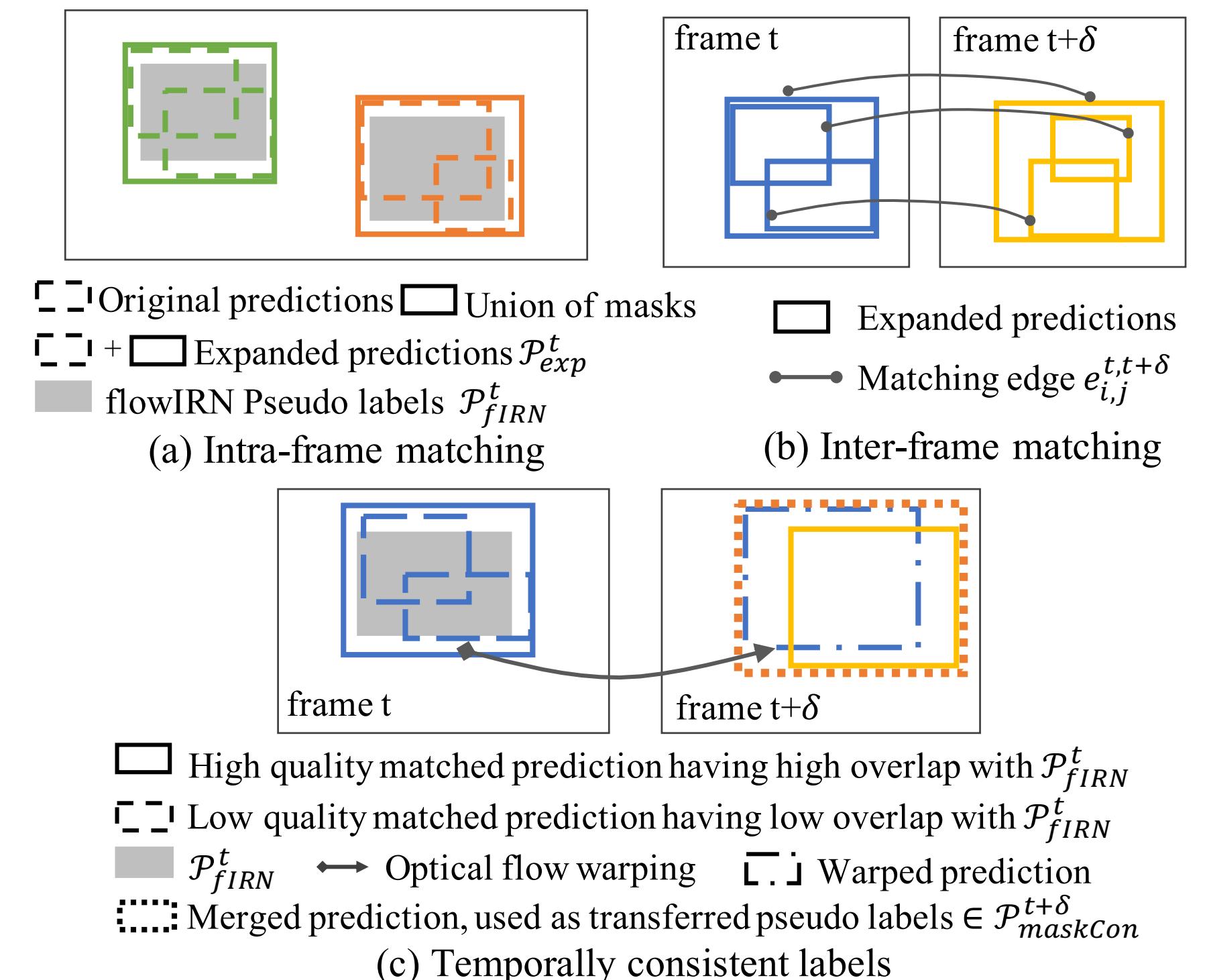
Methods	Supervision	Instance seg	Semantic seg
Mask R-CNN [17]	Mask	38.73	79.23
WISE [27]	Class	10.51	35.82
F2F [29]+MCG [41]	Class	10.73	33.26
IRN [6]	Class	12.33	33.48
IRN [6]+F2F[29]	Class	12.53	34.17
Ours	Class	16.05	39.88
Ours (self-training)	Class	16.82	41.31

Table 2. Frame-level instance segmentation (AP₅₀) and semantic segmentation (IoU) on Cityscapes validation split.



- Incorporating optical flow in flowIRN module:
 - Flow-amplified CAMs: $f\text{-CAM}_c(x) = \text{CAM}_c(x) \times A^{\mathbb{I}(\|\mathcal{F}(x)\|_2 > T)}$
 - Flow-boundary loss: $\mathcal{L}_{\mathcal{F}}^B = \sum_{j \in \mathcal{N}_i} \|\mathcal{F}'(i) - \mathcal{F}'(j)\| \alpha_{i,j} + \lambda |1 - \alpha_{i,j}|$

- Three steps of MaskConsist module:



Video Instance Segmentation Results

Methods	Train_Val Split	Validation Split			
		mAP	AP ₅₀	AP ₇₅	AR ₁
Fully supervised learning methods	IoUTracker+ [58]	-	-	-	-
	DeepSORT [57]	-	-	-	-
	MaskTrack [58]	-	-	-	-
Weakly supervised learning methods	WISE [27]	8.7	22.1	5.5	9.8
	IRN [6]	10.8	26.4	7.7	12.6
	Ours	14.1	34.4	9.4	16.0
	Ours (self-training)	14.1	34.4	9.4	16.0
		mAP	AP ₅₀	AP ₇₅	AR ₁
		23.6	39.2	25.5	26.2
		26.1	42.9	26.1	27.8
		30.3	51.1	32.6	31.0
		30.3	51.1	32.6	35.5

Table 3. Video instance segmentation results on Youtube-VIS dataset.

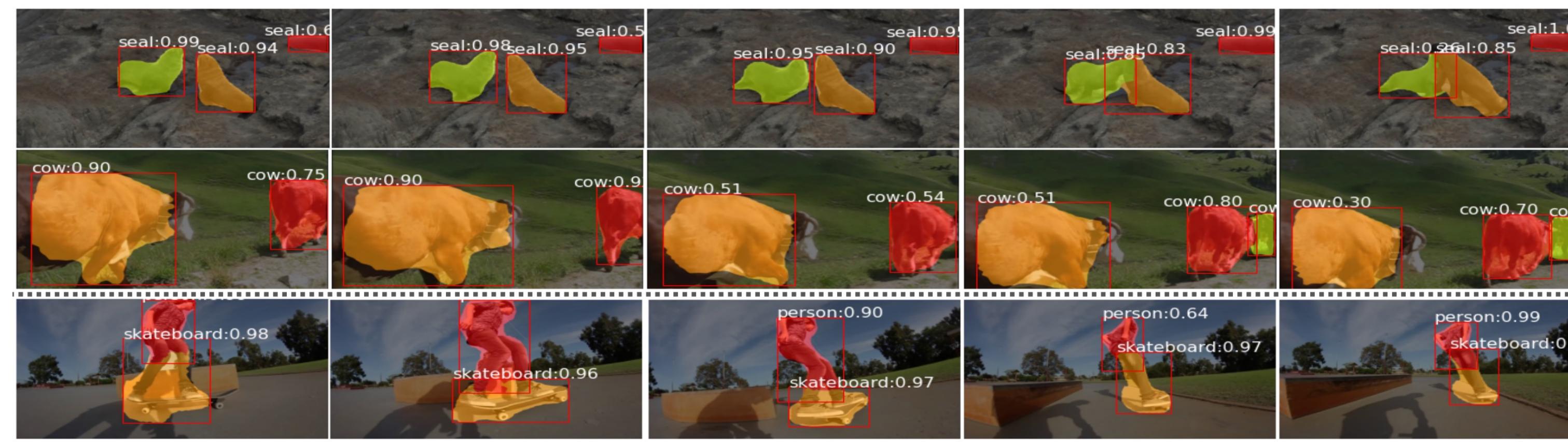
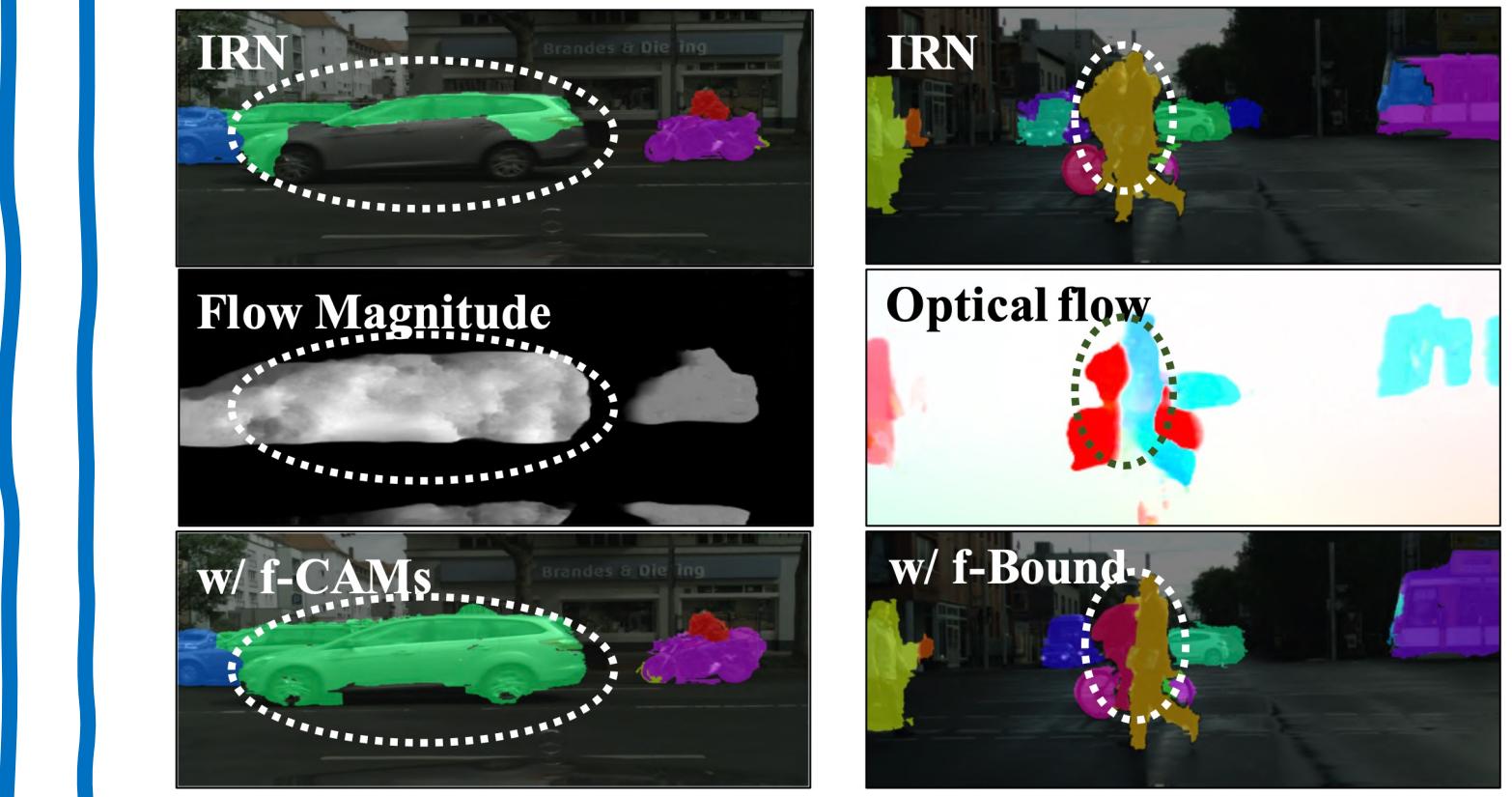


Figure 4. Example Video instance segmentation results from our method on Youtube-VIS dataset.

- Effect of flowIRN:

	YTVIS	Cityscapes
IRN [6]	25.42	8.46
IRN+f-Bound	26.60	9.51
IRN+f-CAMs	27.47	10.55
flowIRN	28.45	10.75



Ablations

- Effect of MaskConsist:

MaskConsist Components			AP ₅₀	
Inter-F	Inter-F	IoM-NMS	YTVIS	Cityscapes
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	31.43	14.66
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	33.75	14.92
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	31.08	14.43
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	33.65	15.27
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	34.66	16.05

Methods	w/ MC	AP ₅₀		TC
		YTVIS	Cityscapes	
WISE [27]	<input checked="" type="checkbox"/>	24.54	10.51	72.08
IRN [6]	<input checked="" type="checkbox"/>	27.03	12.26	80.98
Ours	<input checked="" type="checkbox"/>	31.51	14.72	82.04
FlowIRN Pseudo-label	<input checked="" type="checkbox"/>	31.43	14.66	80.43
Mask R-CNN	<input checked="" type="checkbox"/>	34.66	16.05	84.36
MaskConsist	<input checked="" type="checkbox"/>	34.66	16.05	