

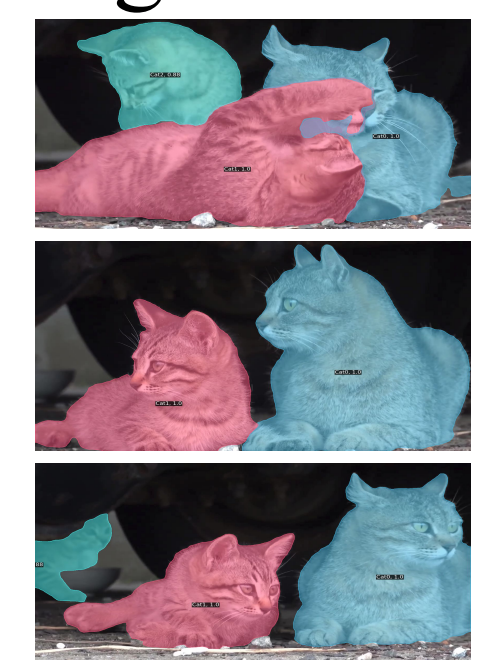


<https://anwesachoudhuri.github.io/ContextAwareRelativeObjectQueries/>

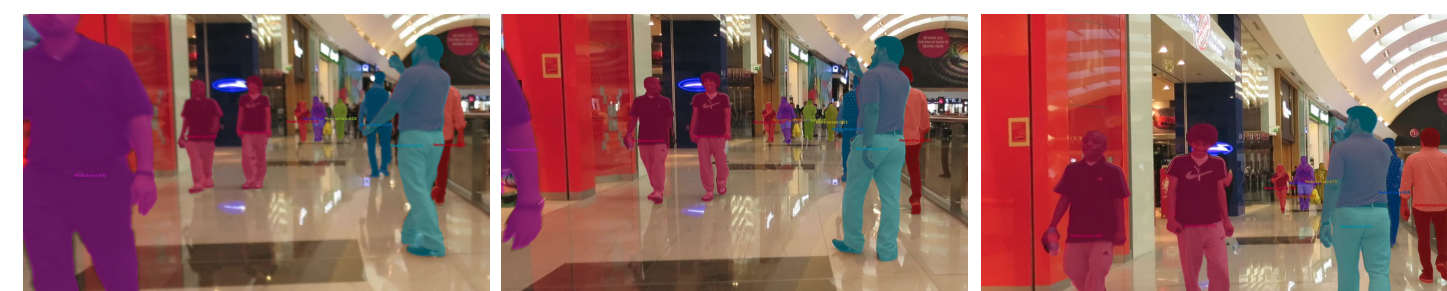
Goal

General and seamless approach for

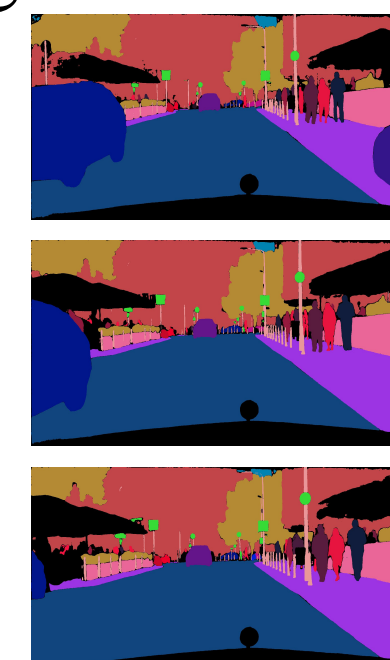
Video Instance Segmentation



Multi-Object Tracking and Segmentation



Video Panoptic Segmentation



Prior Work

- Absolute positional encodings [1, 2, 8]
- Process full video at once [2] or post-processing to merge frames [1]
- Two types of object queries with heuristic post processing [8]

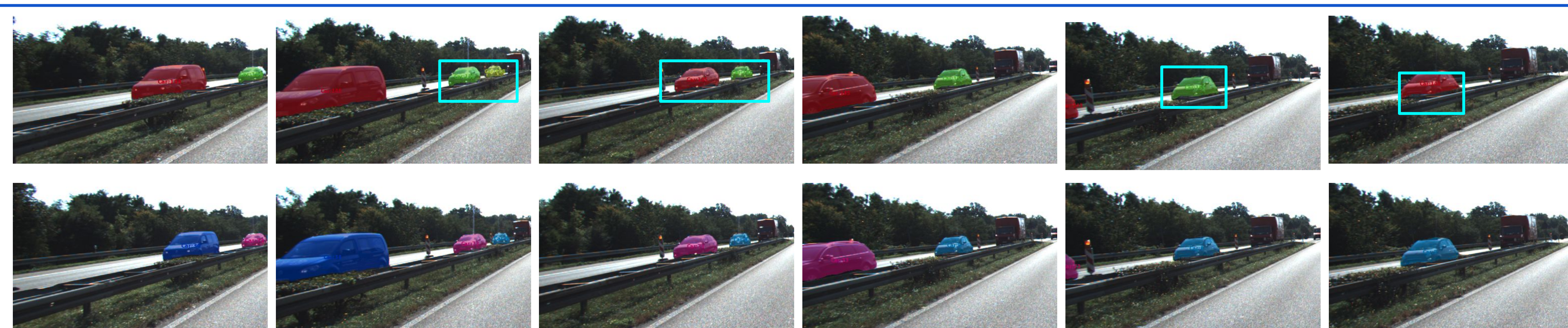
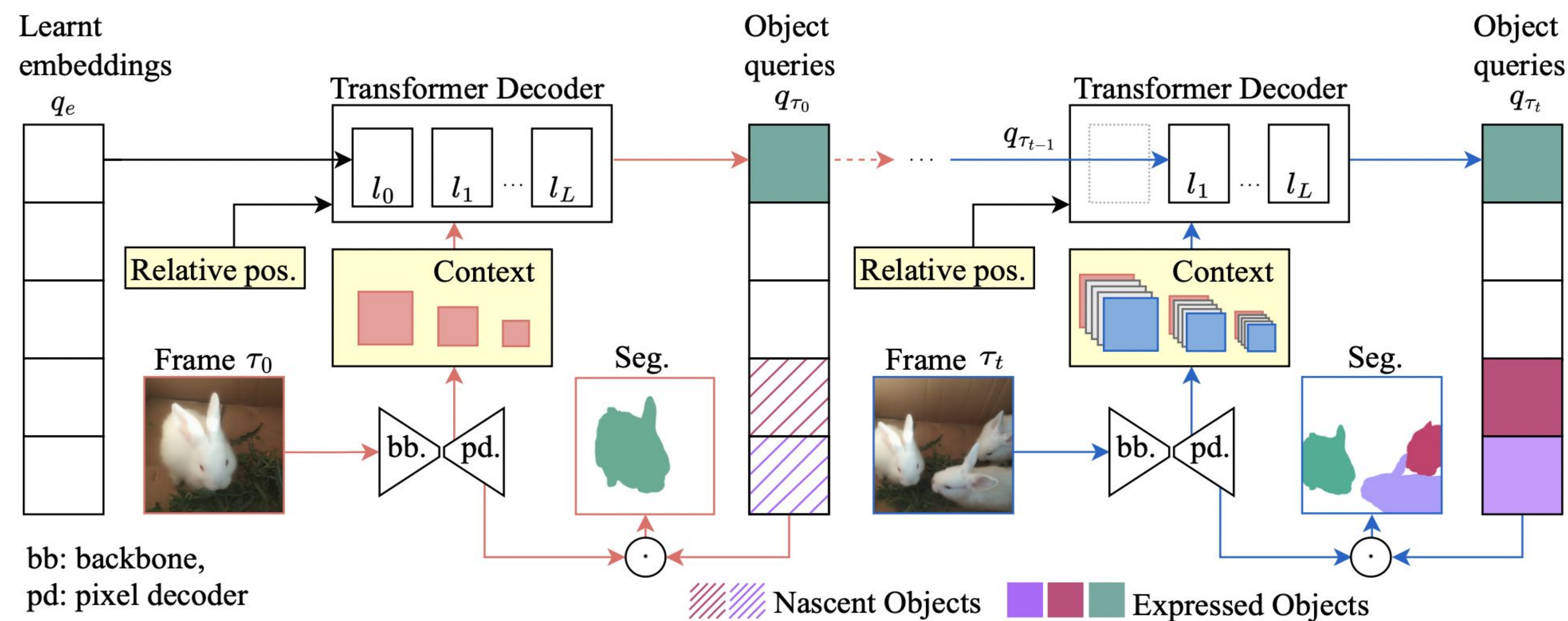
Our Work

- Relative positional encodings
- Frame-by-frame, but no heuristics or post processing
- Only one type of object query

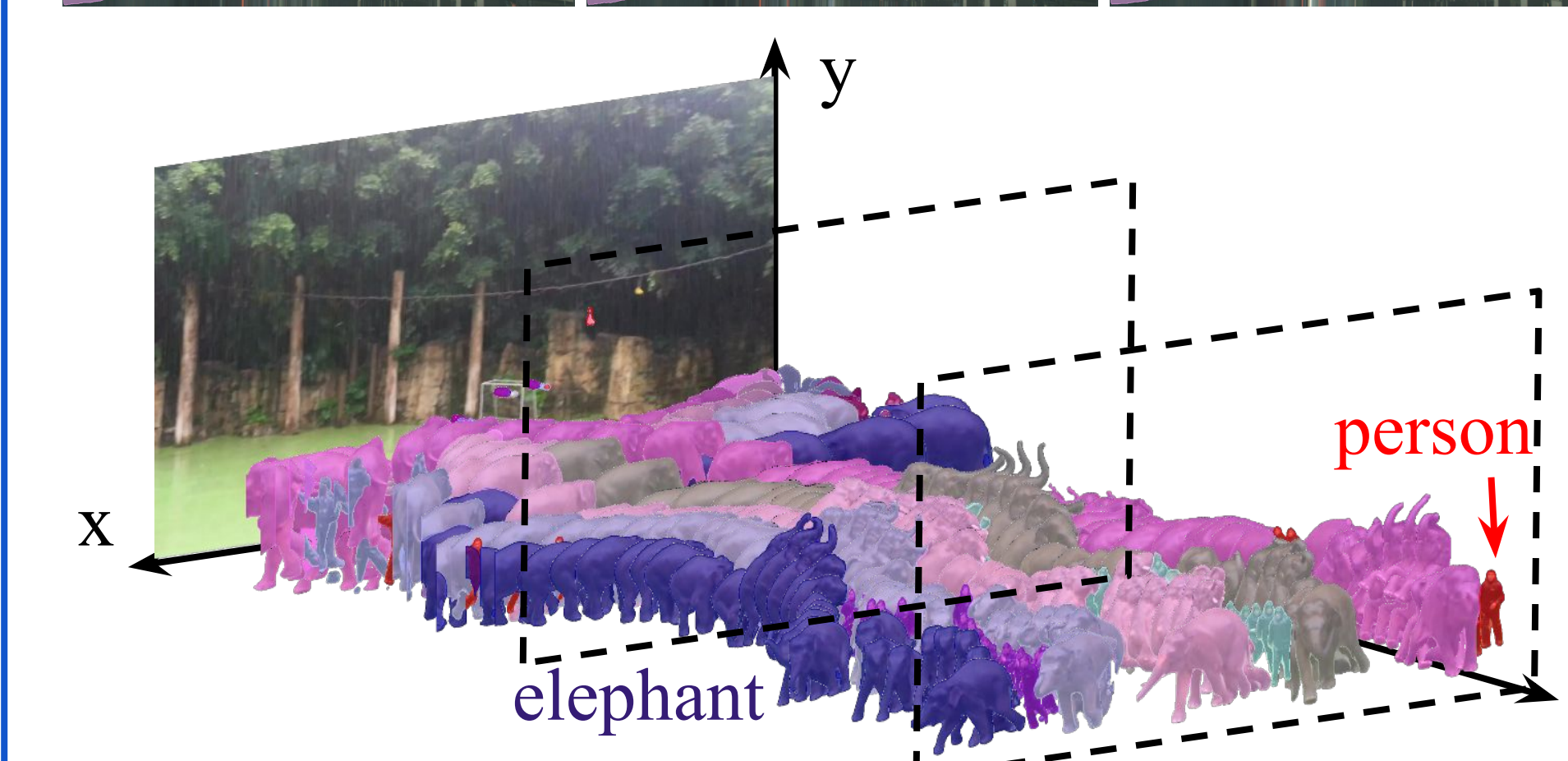
References

- [1] Huang et al., NeurIPS 2022 [5] Choudhuri et al., ICCV 2021
[2] Cheng et al., arXiv 2022 [6] Xu et al., ECCV 2020
[3] Qiao et al., CVPR 2021 [7] Meinhardt et al., CVPR 2022
[4] Kim et al., CVPR 2020

Context-Aware Relative Object Queries



Results



Method	OVIS (AP)	YTVIS 2021 (AP)	YTVIS 2019 (AP)
MinVIS [1]	25.0	44.2	47.4
M2F-VIS [2]	17.3	40.6	46.4
Ours	25.8	43.3	46.7

Cityscapes-VPS				
Method	Depth	VPQ	VPQ _{th}	VPQ _{st}
Vip-DeepLab [3]	✓	63.1	49.5	73.0
VPS-Net [4]		57.5	44.8	66.7
Ours		63.0	48.0	72.8

Ablation

Method	KITTI-MOTS			OVIS		
	Car		Pedestrian	Method	AP	
Ours (Rel. pos.)	83.2	84.5	85.0	64.1	64.4	63.7
Ours (Abs. pos.)	70.0	78.6	62.7	52.0	58.0	46.4

Method	MOTS 2020		KITTI-MOTS	
	sMOTSA		Method	HOTA (car, ped.) / AssA (car, ped.)
Pt.Track [6]	58.1		ASMOTS [5]	68.7, 58.8 / 62.2, 57.6
Tr.Former [7]	58.7		Pt.Track [6]	61.6, 54.4 / 48.8, 48.0
Ours	61.2		Ours	74.5, 62.2 / 64.0, 58.4