



Monocular Multiview Object Tracking with 3D Aspect Parts

Yu Xiang^{1,2*}, Changkyu Song^{2*}, Roozbeh Mottaghi¹ and Silvio Savarese¹
Stanford University¹, University of Michigan at Ann Arbor²

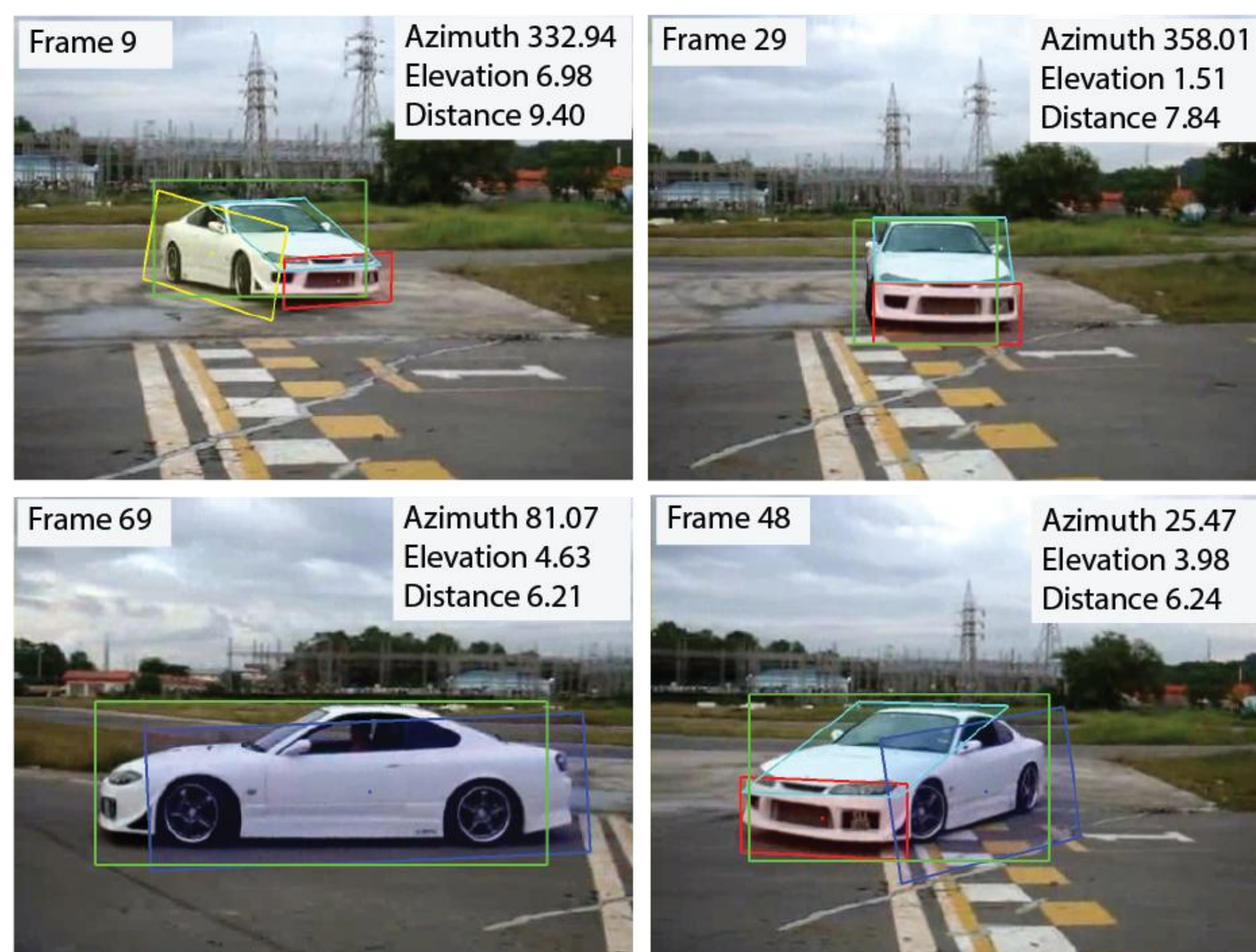


Monocular Multiview Object Tracking

Inputs: video sequences from a single camera

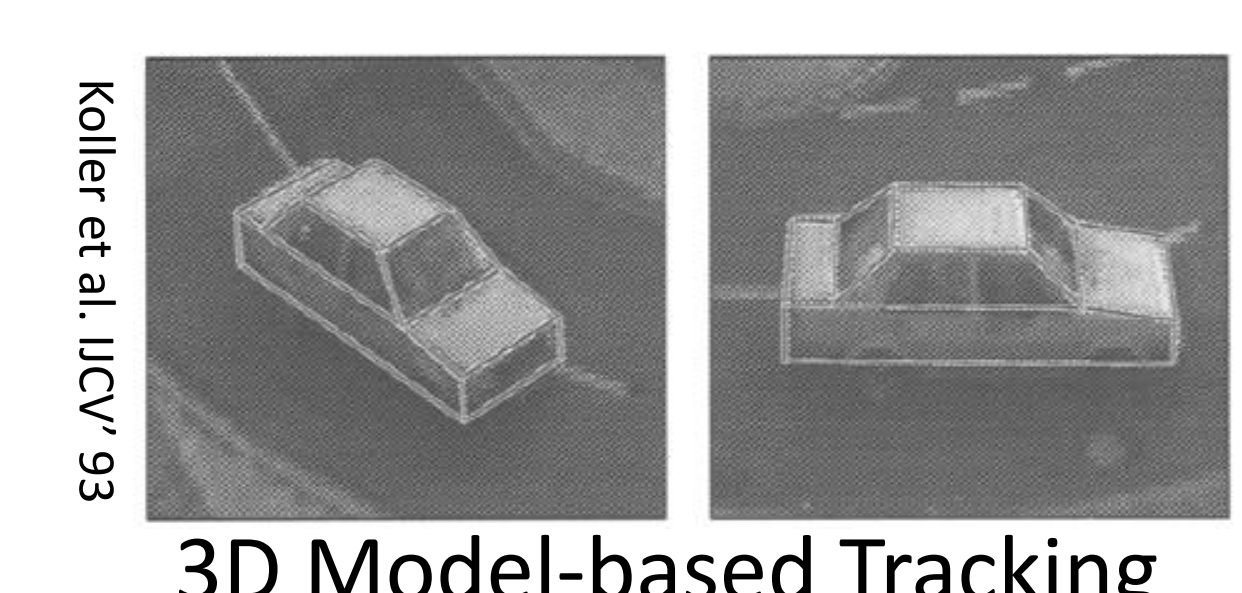
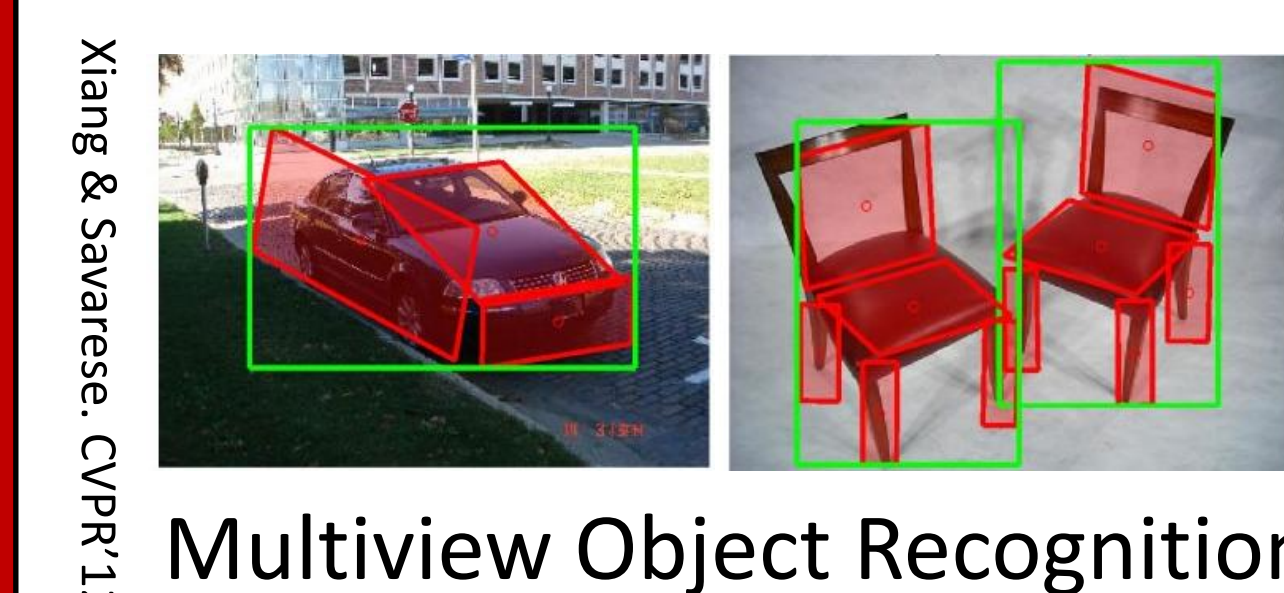
Goals:

- (1) Track the 2D location of the target
- (2) Estimate the 3D pose of the target in time
- (3) Localize the 3D parts of the target in time



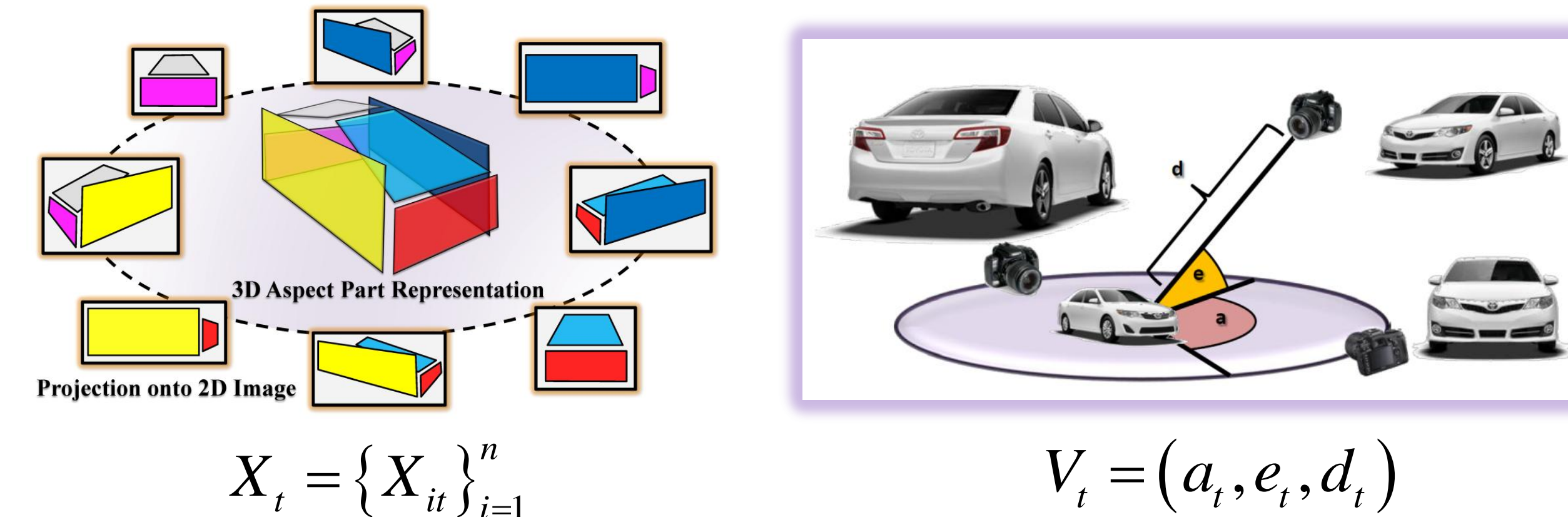
Applications: Autonomous driving, robotics, augmented reality, etc.

Related Problems



Our Multiview Tracking Framework

Object and Viewpoint Representation

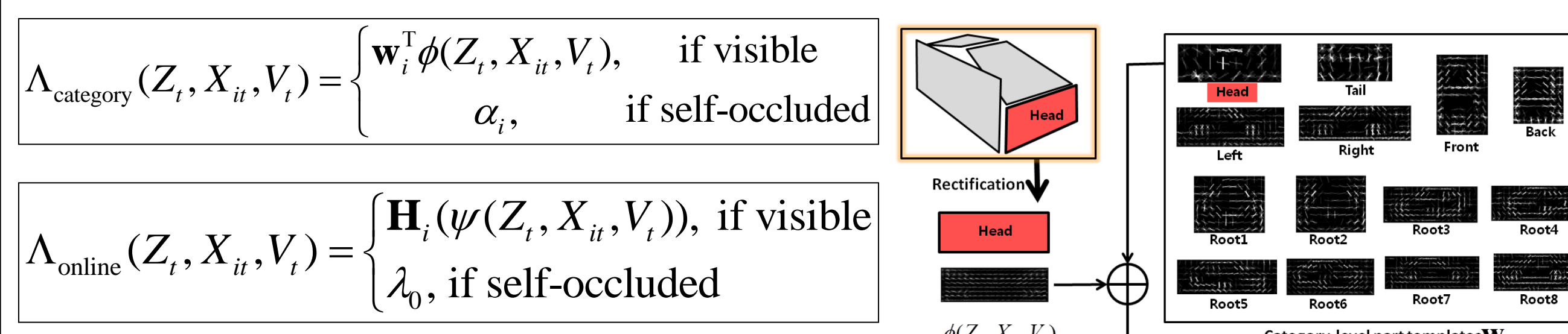


Posterior distribution $P(X_t, V_t | Z_{1:t})$

$$\propto \underbrace{P(Z_t | X_t, V_t)}_{\text{Likelihood}} \underbrace{\int P(X_t, V_t | X_{t-1}, V_{t-1}) P(X_{t-1}, V_{t-1} | Z_{1:t-1}) dX_{t-1} dV_{t-1}}_{\text{Motion Prior Posterior at } t-1}$$

Likelihood $P(Z_t | X_t, V_t) = \prod_{i=1}^n P(Z_t | X_{it}, V_t)$

$$P(Z_t | X_{it}, V_t) \propto \exp(\Lambda_{\text{category}}(Z_t, X_{it}, V_t) + \Lambda_{\text{online}}(Z_t, X_{it}, V_t))$$

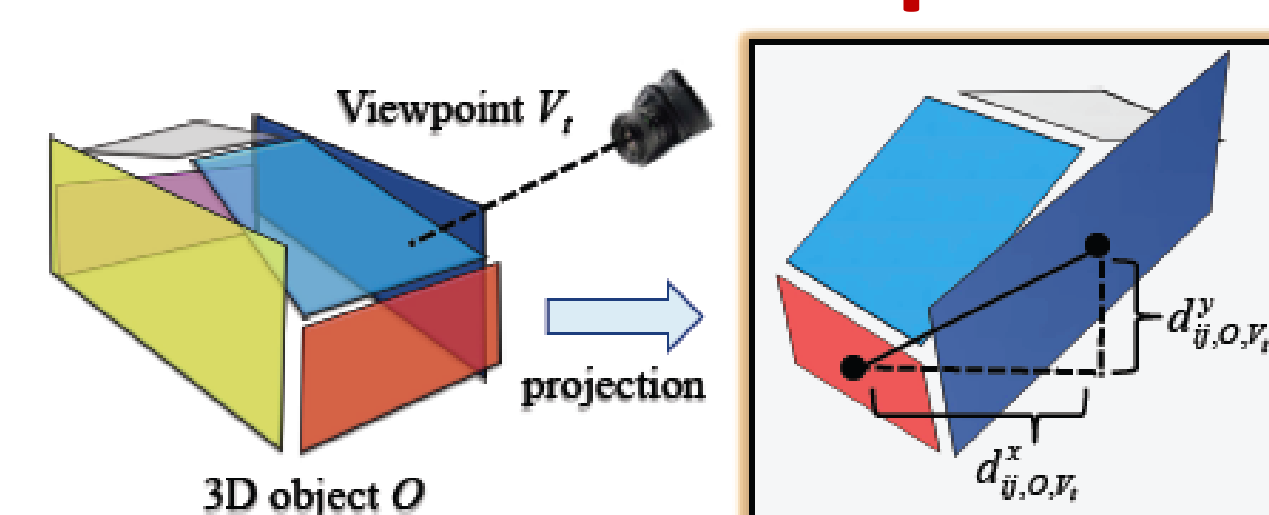


Motion Prior

$$P(X_t, V_t | X_{t-1}, V_{t-1}) = \underbrace{P(X_t | X_{t-1}, V_t)}_{\text{Location}} \underbrace{P(V_t | V_{t-1})}_{\text{Viewpoint}}$$

$$P(X_t | X_{t-1}, V_t) \propto \prod_{i=1}^n P(X_{it} | X_{i(t-1)}) \prod_{(i,j)} \Lambda(X_{it}, X_{jt}, V_t)$$

$$\Lambda(X_{it}, X_{jt}, V_t) = P(\Delta_i(x_i, x_j) | V_t) P(\Delta_j(y_i, y_j) | V_t)$$



Multiview Particle Filtering Object Tracking

References

- [1] B. Babenko, M.H. Yang, S. Belongie. Robust object tracking with online multiple instance learning. TPAMI, 2011.
- [2] C. Bao, Y. Wu, H. Ling, H. Ji. Real time robust l1 tracker using accelerated proximal gradient approach. In CVPR, 2012.
- [3] Z. Kalal, K. Mikolajczyk, J. Matas. Tracking-learning-detection. TPAMI, 2012.
- [4] S. Hare, A. Saari, P.H. Torr. Struck: Structured output tracking with kernels. In ICCV, 2011.
- [5] P.F. Felzenszwalb, R.B. Girshick, D. McAllester, D. Ramanan. Object detection with discriminatively trained part-based models. TPAMI, 2010.
- [6] A. Geiger, P. Lenz, R. Urtasun. Are we ready for autonomous driving? In CVPR, 2012.
- [7] Y. Xiang, S. Savarese. Estimating the aspect layout of object categories. In CVPR, 2012.

Acknowledgement

We acknowledge the support of DARPA UPSIDE grant A13-0895-S002 and NSF CAREER grant N.1054127.

Experiments

2D Object Tracking

Video	MIL [1]	L1 [2]	TLD [3]	Struct [4]	DPM [5]+PF	Category Model	Full Model
YouTube	0.37	0.44	0.38	0.40	0.74	0.74	0.75
KITTI [6]	0.34	0.28	0.29	0.36	0.54	0.55	0.58
06_car [3]	0.19	0.52	0.85	0.48	0.70	0.67	0.70

Metric: mean bounding box overlap ratio

Continuous Viewpoint Estimation

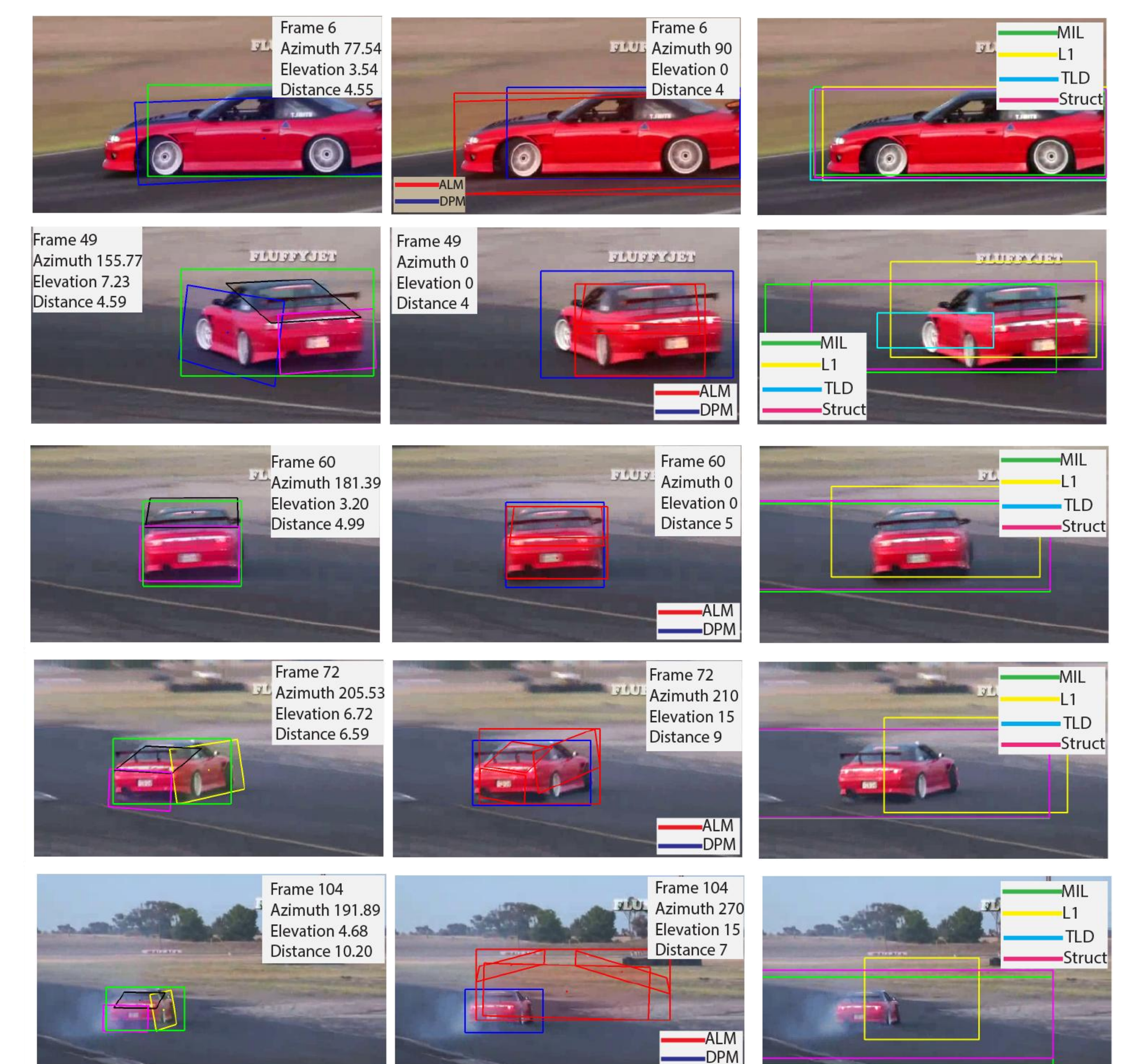
Video	Full Model	Category Model	ALM [7]
YouTube	13.46°	18.38°	47.24°
KITTI [6]	14.66°	23.20°	37.89°

Metric: mean absolute difference in azimuth angle

3D Aspect Part Localization

Video	Full Model	Category Model	ALM [7]
YouTube	0.41	0.40	0.30
KITTI [6]	0.36	0.30	0.26

Metric: mean overlap ratio of part shape



Ours

Object Detection

Online Tracking