

ClusterVO: Clustering Moving Instances and Estimating Visual Odometry for Self and Surroundings

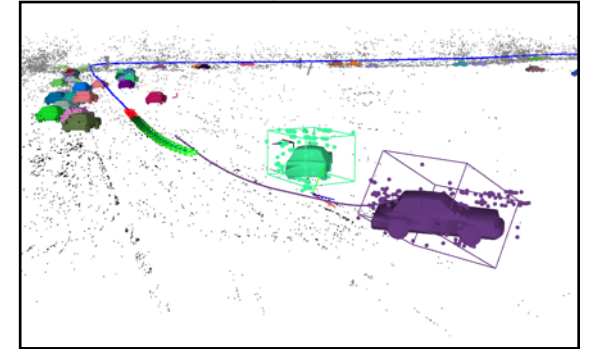
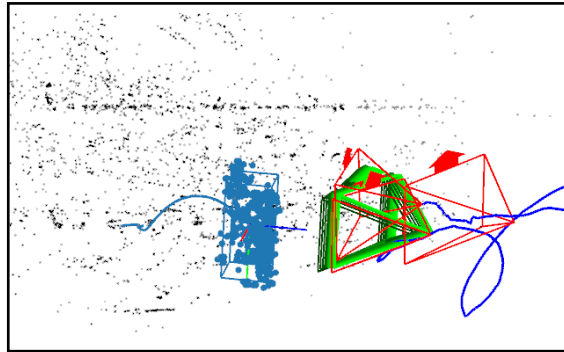
Jiahui Huang¹, Sheng Yang², Tai-Jiang Mu¹, Shi-Min Hu¹

¹BNRist, Department of Computer Science and Technology, Tsinghua University, Beijing, China

²Alibaba Inc., China



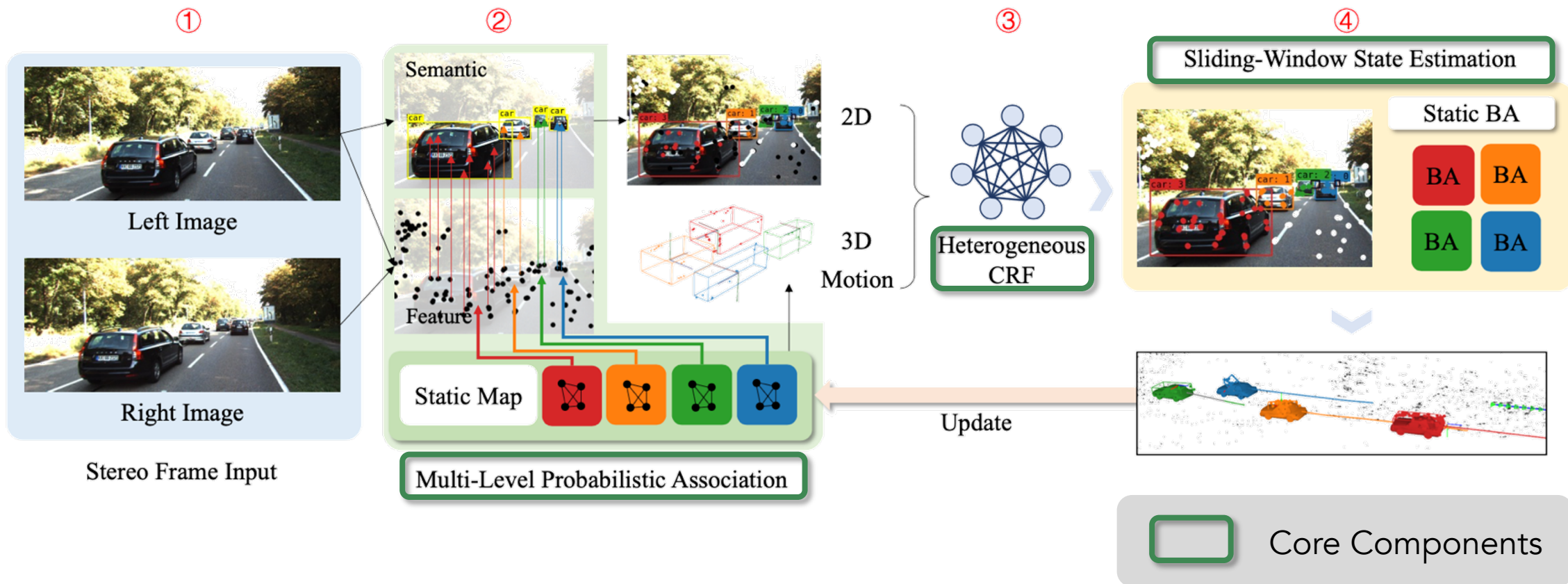
Introduction



ClusterVO is a stereo Visual Odometry which simultaneously clusters and estimates the motion of both ego and surrounding rigid clusters/objects.

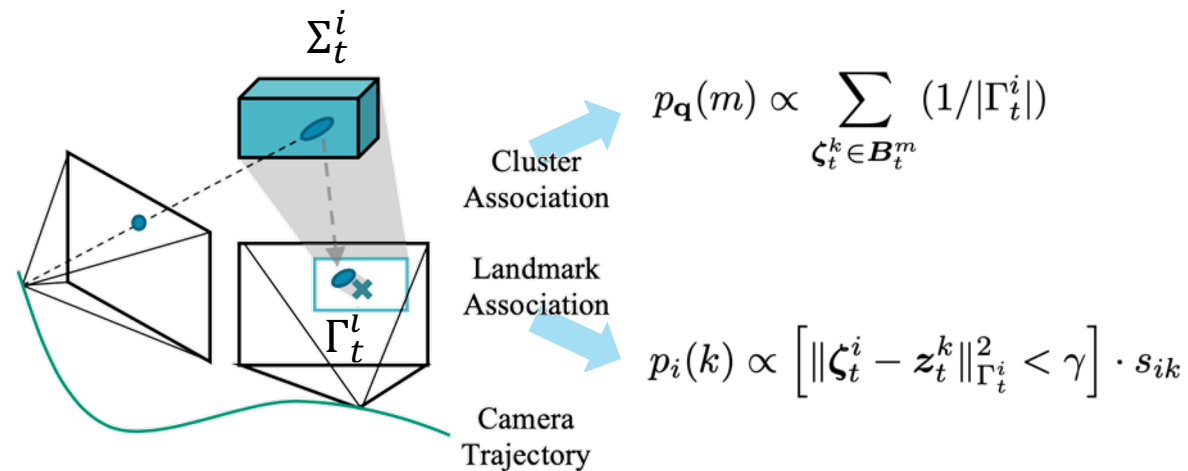
Unlike previous solutions relying on batch input or imposing priors on scene structure or dynamic object models, our method is online, general and applicable in various scenarios including indoor scene understanding and autonomous driving.

System Pipeline



Multi-level Probabilistic Association

- For each new frame, we need to robustly associate detected features and semantic bounding boxes to map landmarks and clusters.
- Association probabilities for the two levels are calculated based on stereo triangulation uncertainty.



$$\Sigma_t^i := \mathbf{R}_{t'}^c \mathbf{z} \Sigma_{t'}^i \mathbf{R}_{t'}^{c\top}, \quad t' := \operatorname{argmin}_{t' < t} |\mathbf{z} \Sigma_{t'}^i|$$

$$\zeta_t^i = \pi(\mathbf{p}_t^i + \mathbf{v}_t^q) \quad \Gamma_t^i = \mathbf{J}_\pi \Sigma_t^i \mathbf{J}_\pi^\top$$

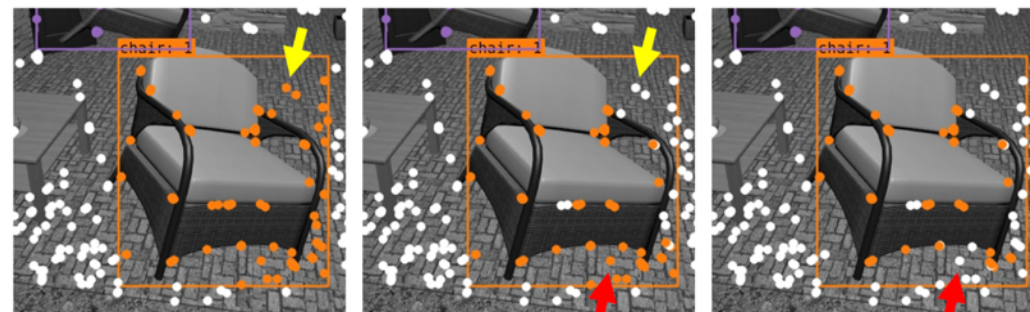
Heterogeneous CRF

- The cluster assignment \mathbf{q}^i of each landmark i observed in the current frame is assigned by a Conditional Random Field model combining semantic (2D), spatial (3D) and motion information, which we call 'Heterogeneous CRF'.

$$E(\{\mathbf{q}^i\}_i) := \sum_i \psi_u(\mathbf{q}^i) + \alpha \sum_{i < j} \psi_p(\mathbf{q}^i, \mathbf{q}^j)$$

$$\psi_u(\mathbf{q}^i) \propto p_{2D}(\mathbf{q}^i) \cdot p_{3D}(\mathbf{q}^i) \cdot p_{mot}(\mathbf{q}^i)$$

$$\psi_p(\mathbf{q}^i, \mathbf{q}^j) := [\mathbf{q}^i \neq \mathbf{q}^j] \cdot \exp(-\|\mathbf{p}_t^i - \mathbf{p}_t^j\|^2)$$



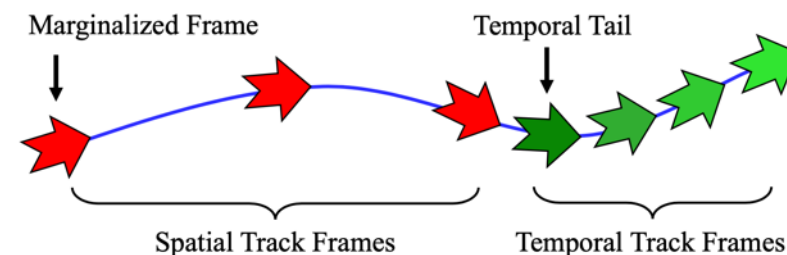
2D Only

2D + 3D

2D + 3D +
Motion

Sliding-Window State Optimization

- We employ a novel double-track frame management design to maintain keyframes in dynamic scenes.
 - Temporal track frames allow for enough observations to track fast-moving clusters.
 - Spatial track frames help create enough parallax for accurate triangulation.
- The energy function for state optimization is defined separately for static and dynamic parts:



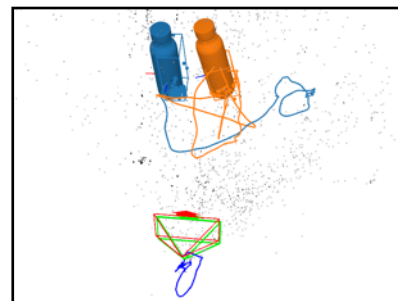
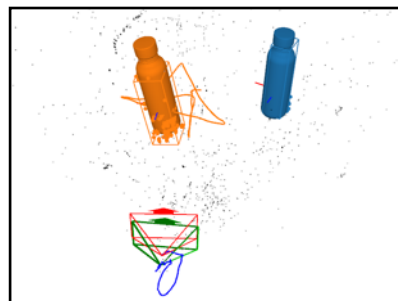
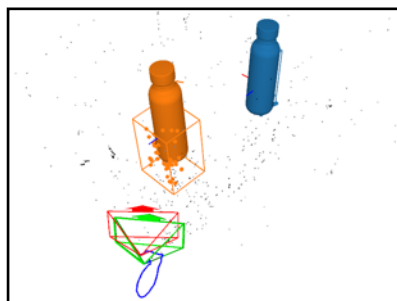
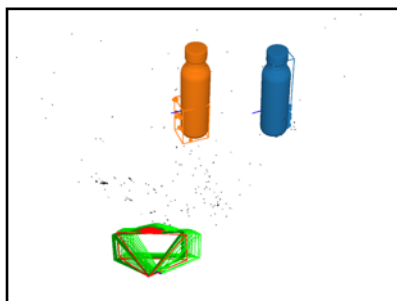
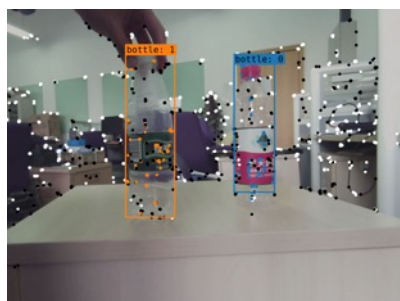
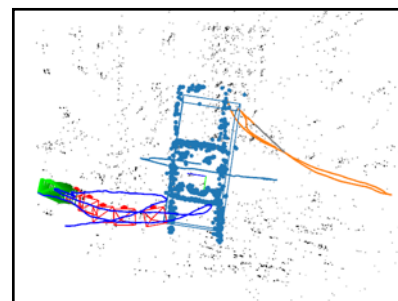
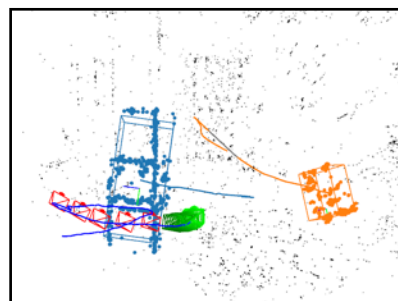
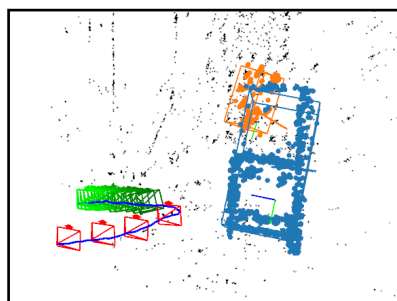
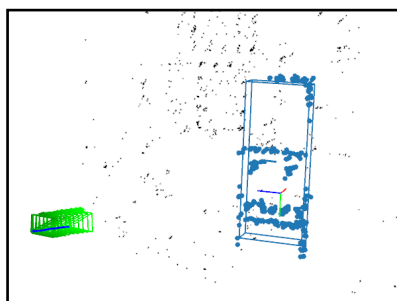
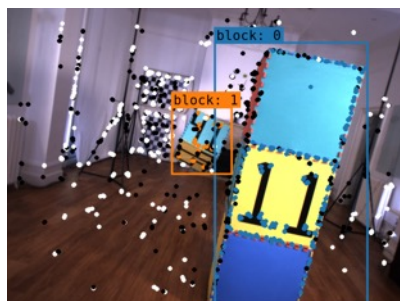
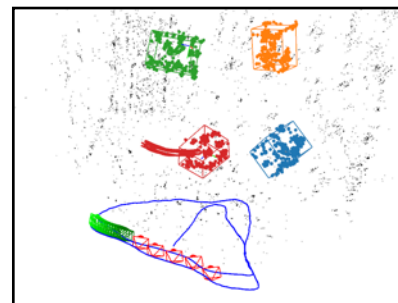
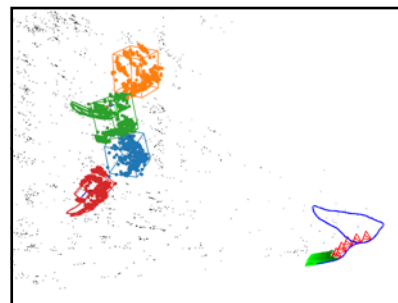
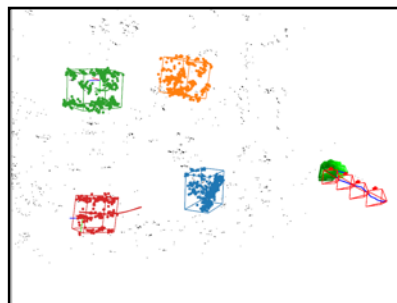
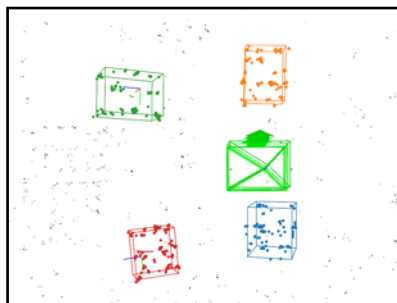
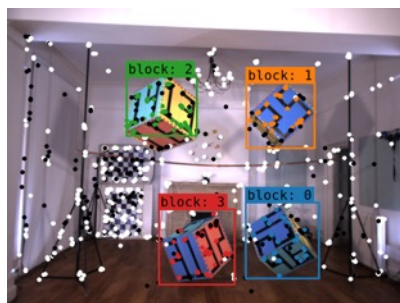
Static scene and Camera pose
BA Term + Marginalization

$$\mathbf{E}_s := \sum_{i \in \mathcal{I}_0, t \in \mathcal{T}_a} \rho(\|z_t^i - \pi((\mathbf{P}_t^c)^{-1} \mathbf{p}_t^i)\|_{z\Sigma}^2) + \sum_{t \in \mathcal{T}_a} \|\delta \mathbf{x}_t^c - \mathbf{H}^{-1} \beta\|_{\mathbf{H}}^2$$

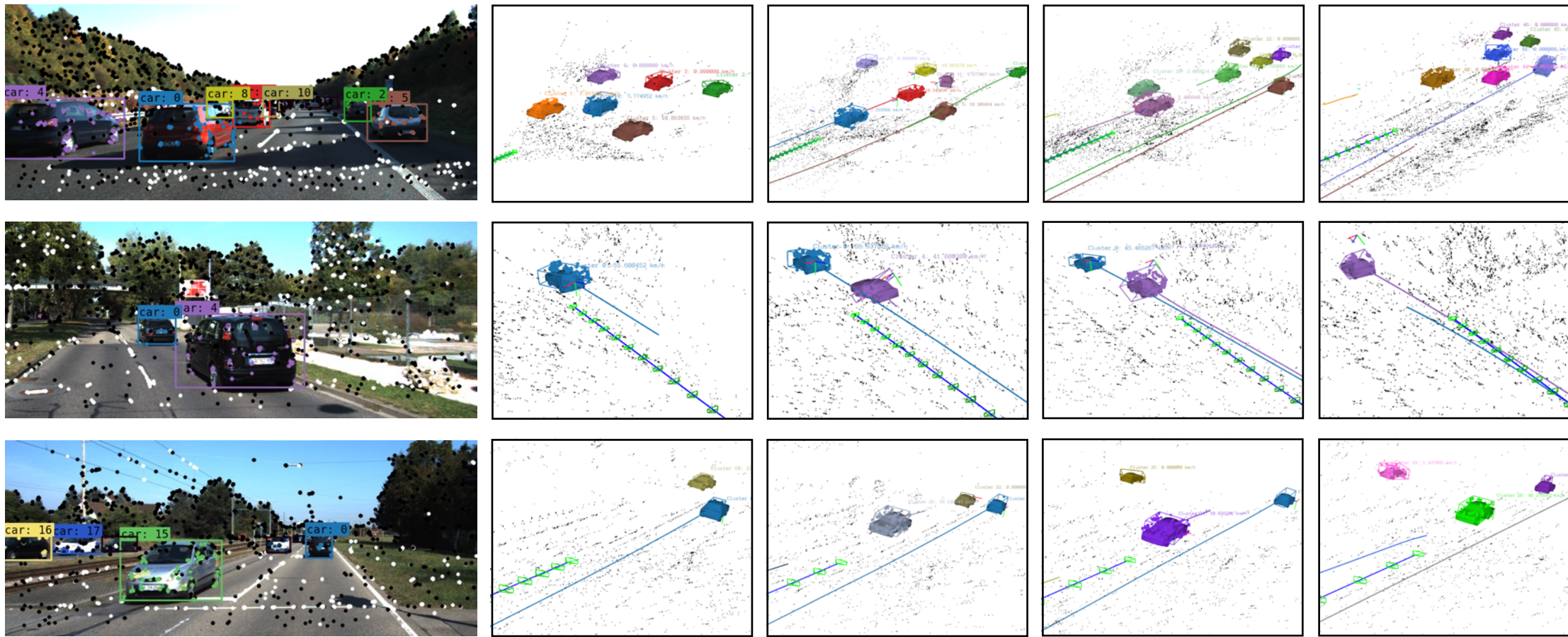
Dynamic clusters
Dynamic BA Term +
Smoothness

$$\mathbf{E}_d := \sum_{t, t+ \in \mathcal{T}_t} \left\| \begin{bmatrix} \mathbf{t}_{t+}^i \\ \mathbf{v}_{t+}^i \end{bmatrix} - \mathbf{A} \begin{bmatrix} \mathbf{t}_t^i \\ \mathbf{v}_t^i \end{bmatrix} \right\|_{\hat{\mathbf{Q}}}^2 + \sum_{i \in \mathcal{I}_q, t \in \mathcal{T}_t} \rho(\|z_t^i - \pi(\mathbf{T}_t^{\text{cqi}} (\mathbf{P}_t^c)^{-1} \mathbf{p}_t^i)\|_{z\Sigma}^2)$$

Results: Indoor Dataset



Results: KITTI Raw Dataset



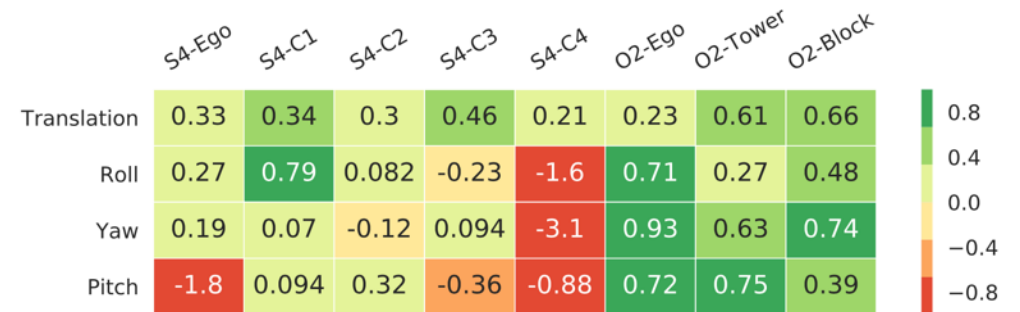
Comparisons

Ego-Motion on KITTI raw Dataset

| Sequence | ORB-SLAM2 [28] | | | DynSLAM [2] | | | Li <i>et al.</i> [24] | ClusterSLAM [15] | | | ClusterVO | | |
|-----------|----------------|-------------|-------------|-------------|-------|--------------|-----------------------|------------------|-------------|-------|-------------|-------------|-------------|
| | ATE | R.RPE | T.RPE | ATE | R.RPE | T.RPE | ATE | ATE | R.RPE | T.RPE | ATE | R.RPE | T.RPE |
| 0926-0009 | 0.91 | 0.01 | 1.89 | 7.51 | 0.06 | 2.17 | 1.14 | 0.92 | 0.03 | 2.34 | 0.79 | 0.03 | 2.98 |
| 0926-0013 | 0.30 | 0.01 | 0.94 | 1.97 | 0.04 | 1.41 | 0.35 | 2.12 | 0.07 | 5.50 | 0.26 | 0.01 | 1.16 |
| 0926-0014 | 0.56 | 0.01 | 1.15 | 5.98 | 0.09 | 2.73 | 0.51 | 0.81 | 0.03 | 2.24 | 0.48 | 0.01 | 1.04 |
| 0926-0051 | 0.37 | 0.00 | 1.10 | 10.95 | 0.10 | 1.65 | 0.76 | 1.19 | 0.03 | 1.44 | 0.81 | 0.02 | 2.74 |
| 0926-0101 | 3.42 | 0.03 | 14.27 | 10.24 | 0.13 | 12.29 | 5.30 | 4.02 | 0.02 | 12.43 | 3.18 | 0.02 | 12.78 |
| 0929-0004 | 0.44 | 0.01 | 1.22 | 2.59 | 0.02 | 2.03 | 0.40 | 1.12 | 0.02 | 2.78 | 0.40 | 0.02 | 1.77 |
| 1003-0047 | 18.87 | 0.05 | 28.32 | 9.31 | 0.05 | 6.58 | 1.03 | 10.21 | 0.06 | 8.94 | 4.79 | 0.05 | 6.54 |

| | AP _{bv} | | | AP _{3D} | | | Time (ms) |
|------------------------|------------------|--------------|--------------|------------------|--------------|--------------|------------|
| | Easy | Moderate | Hard | Easy | Moderate | Hard | |
| Chen <i>et al.</i> [7] | 81.34 | 70.70 | 66.32 | 80.62 | 70.01 | 65.76 | 1200 |
| DynSLAM [2] | 71.83 | 47.16 | 40.30 | 64.51 | 43.70 | 37.66 | 500 |
| ClusterVO | 74.65 | 49.65 | 42.65 | 55.85 | 38.93 | 33.55 | 125 |

Object Detection on KITTI



Trajectory Accuracy on OMD

Thank You!

