



16TH EUROPEAN CONFERENCE ON
COMPUTER VISION

WWW.ECCV2020.EU



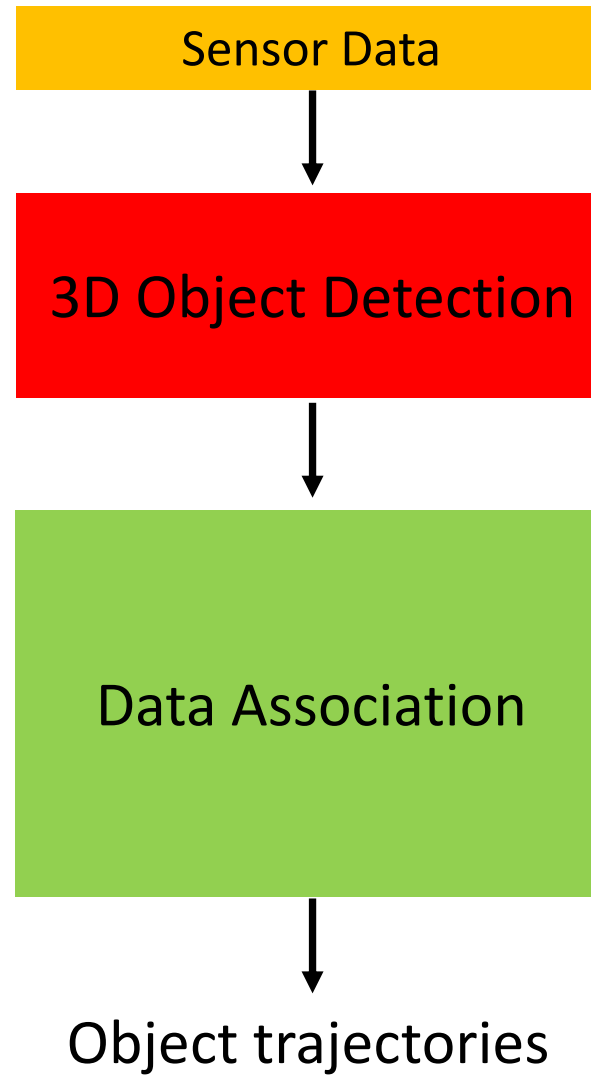
Graph Neural Network for 3D Multi-Object Tracking

Xinshuo Weng, Yongxin Wang, Yunze Man, Kris Kitani
Robotics Institute, Carnegie Mellon University

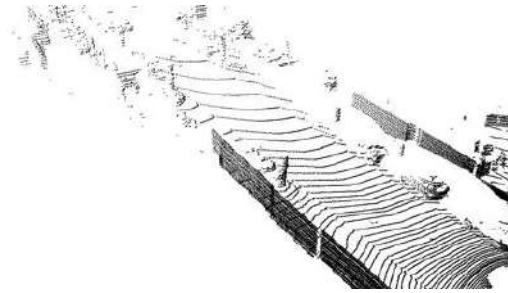
European Conference on Computer Vision (ECCV) Workshops



Standard 3D MOT Pipeline



Standard 3D MOT Pipeline



LiDAR point clouds

Sensor Data

3D Object Detection

Data Association

Object trajectories



RGB frames



Standard 3D MOT Pipeline

Sensor Data



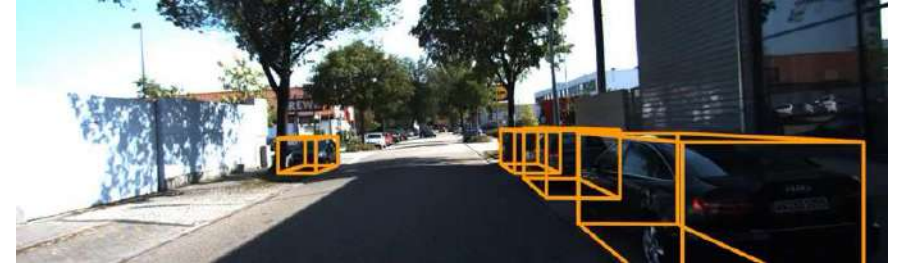
3D Object Detection



Data Association



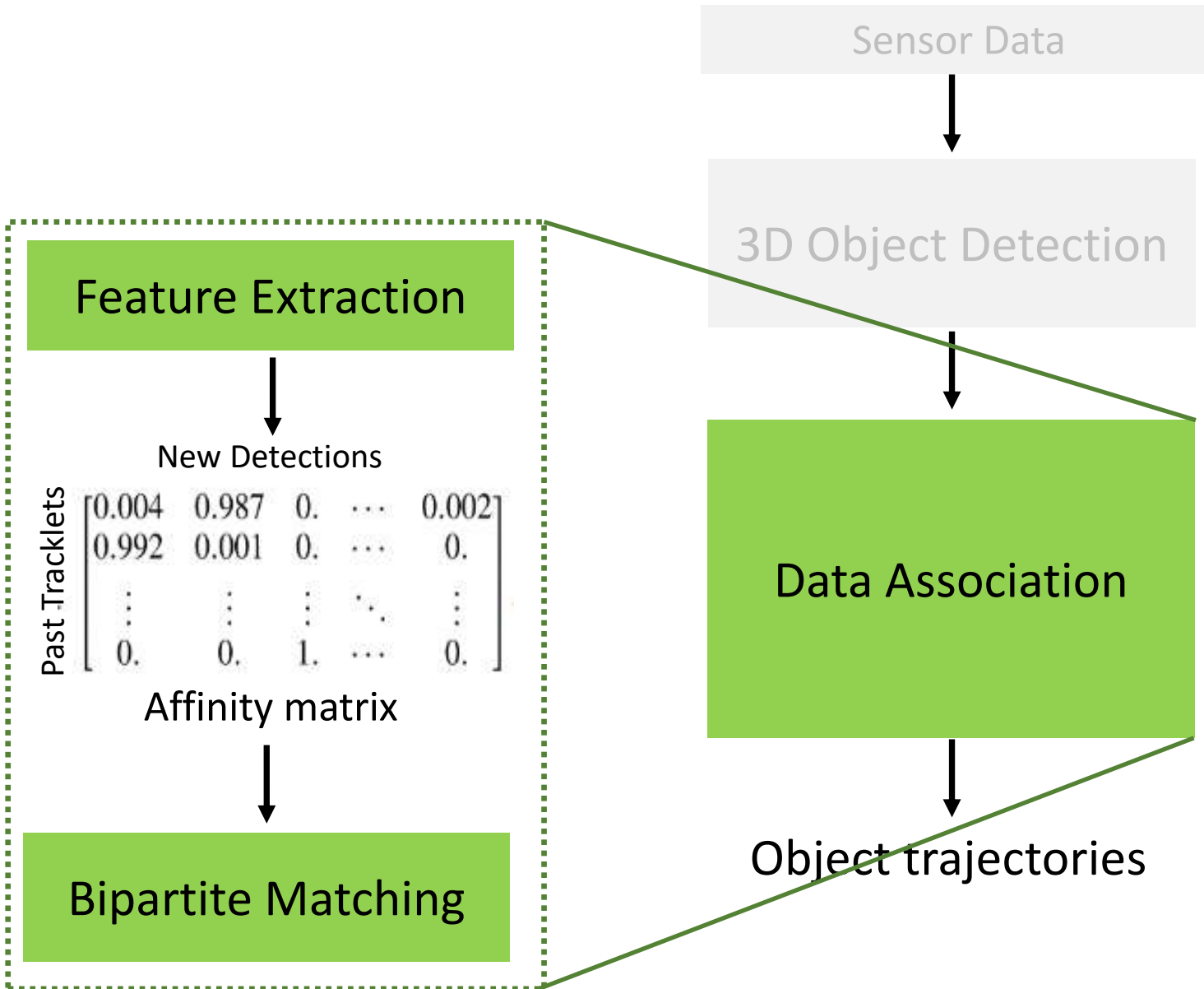
Object trajectories



Detection results



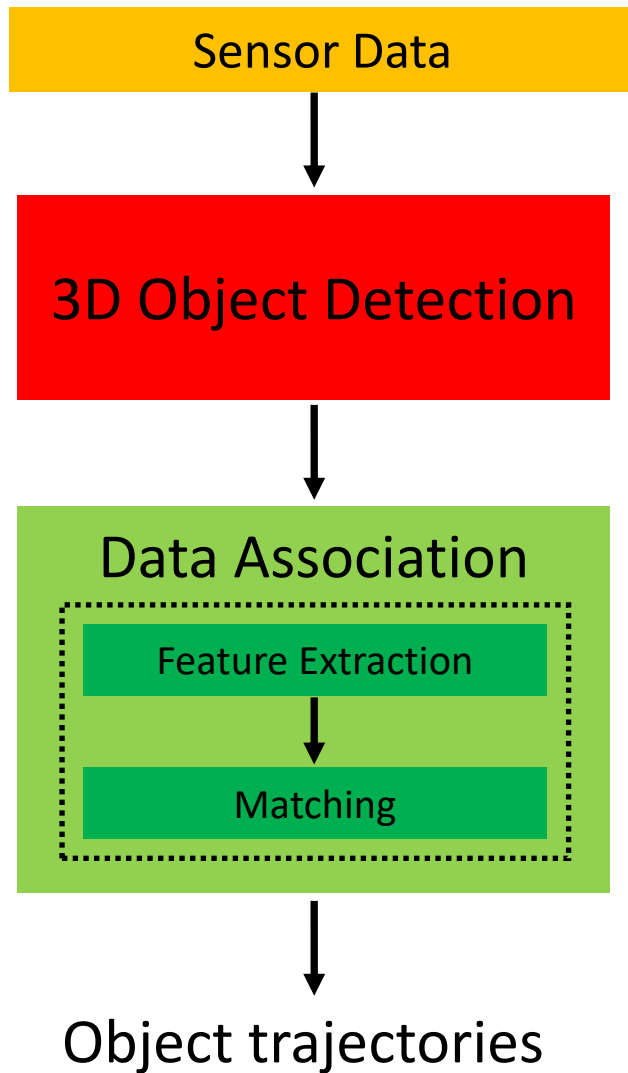
Standard 3D MOT Pipeline



3D MOT results



Limitation of the Prior Work



Limitation

1. Feature representation does not take into account contexts of other objects
2. Feature representation does not fully utilize information from multiple modalities that is complementary



Our Contributions

1. A novel feature interaction mechanism to encode contexts via object interaction
2. A 2D-3D joint feature extractor to learn multi-modal features that are complementary

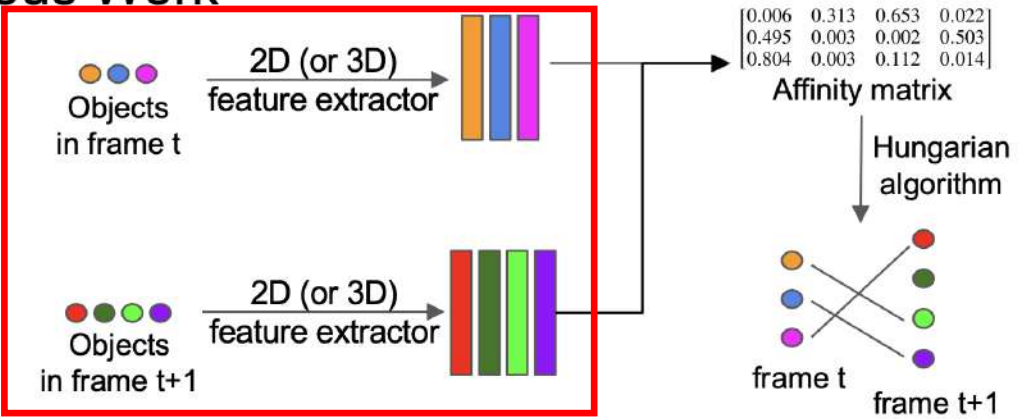


Our Contributions

Prior work

- Feature extraction is independent of each object
- Employs features from one modality (2D or 3D)

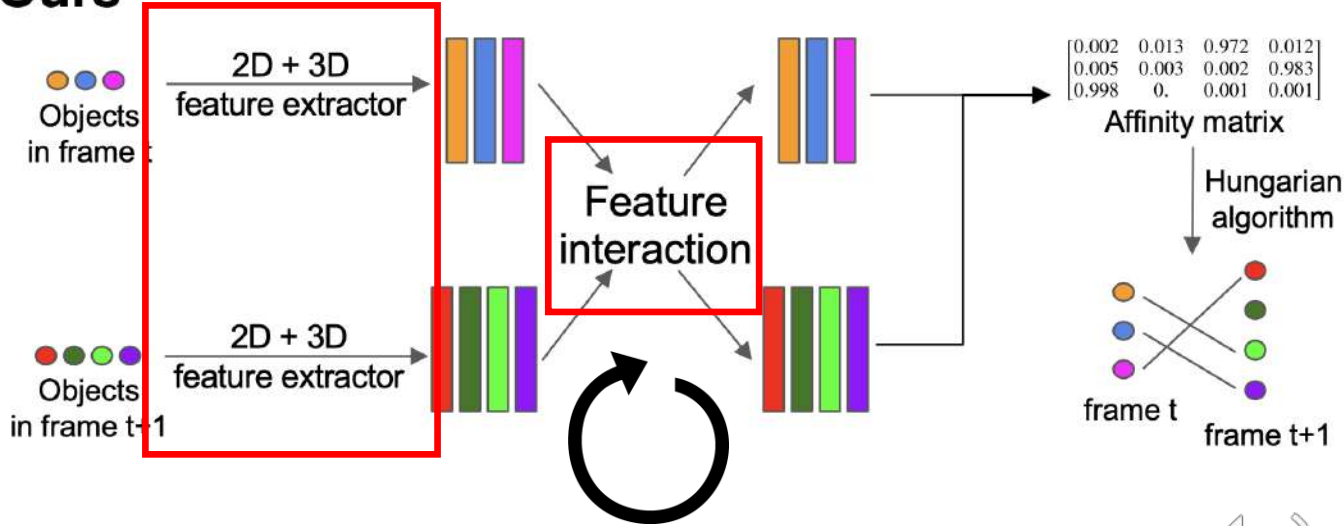
Previous Work



Our Approach

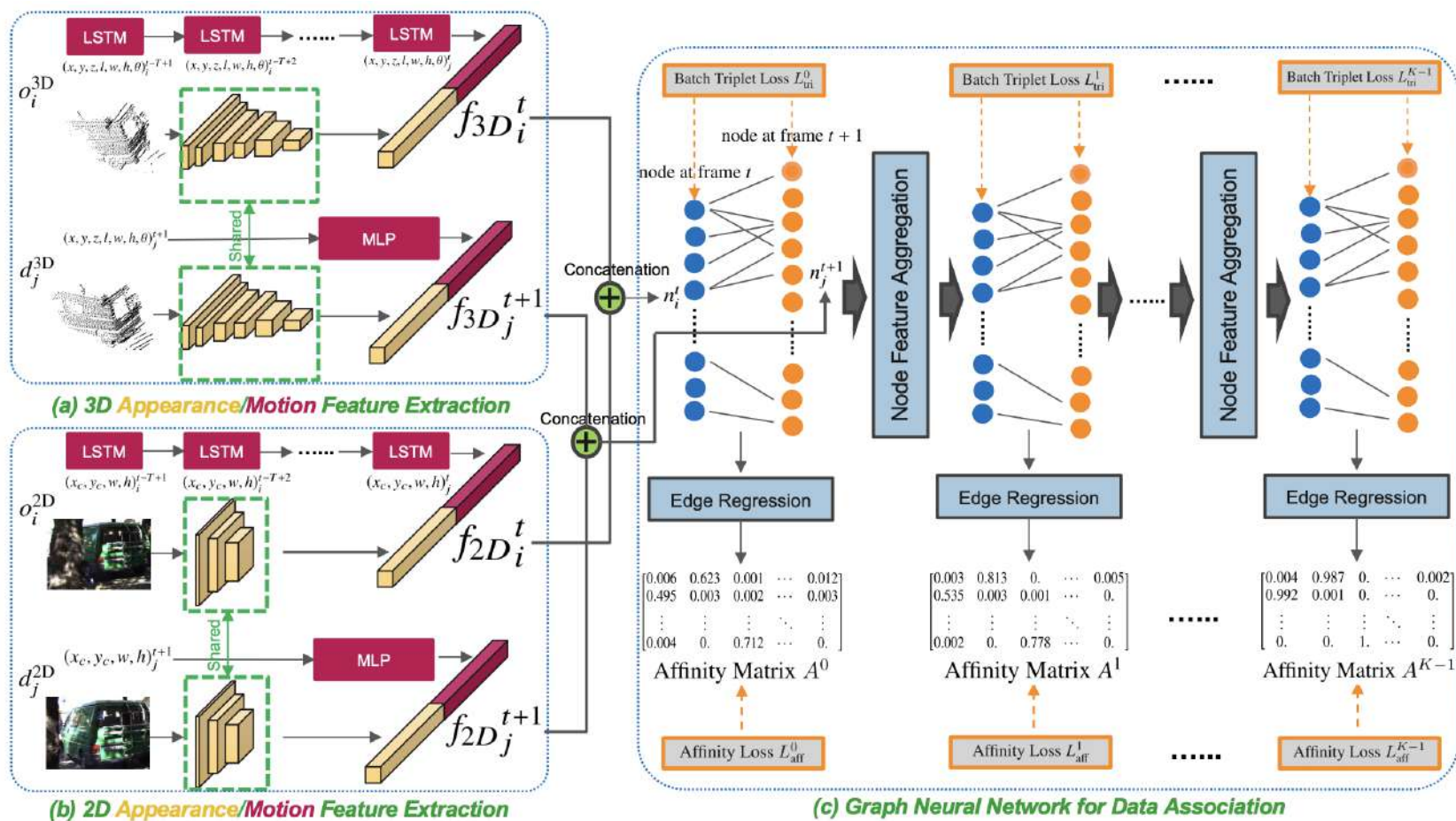
- A joint feature extractor to learn multi-modal features
- A novel feature interaction mechanism to iteratively encode context and improve discriminative feature learning

Ours



Our Approach

- (a) Obtain the appearance / motion features from the 3D space
- (b) Obtain the appearance / motion features from the 2D space
- (c) Learn discriminative object features by encoding context through object feature interaction



Ablation Study



Improve Feature Learning for 3D MOT

- Is encoding the multi-modal features really useful?

Feature Extractor	sAMOTA (%) ↑	AMOTA (%) ↑	AMOTP (%) ↑	MOTA (%) ↑
2D A	88.31	41.62	76.22	79.42
2D M	64.24	23.95	61.13	54.88
3D A	88.27	41.55	76.29	77.38
3D M	88.57	41.62	76.22	81.84

Use feature from
single modality

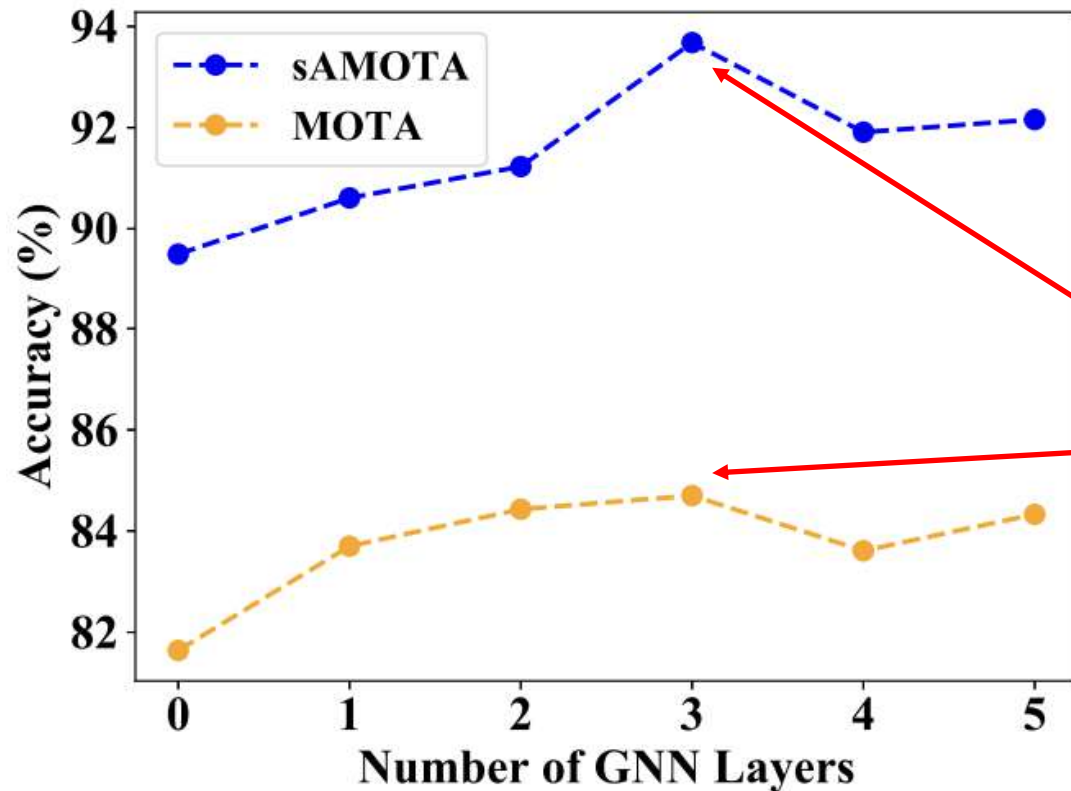
Use feature from
multiple modalities:
Performance
increased!

A: appearance feature, M: motion feature



Improve Feature Learning for 3D MOT

- Is feature interaction using GNNs useful to 3D MOT?



Performance largely increased with GNN layers = 3 v.s. 0 !

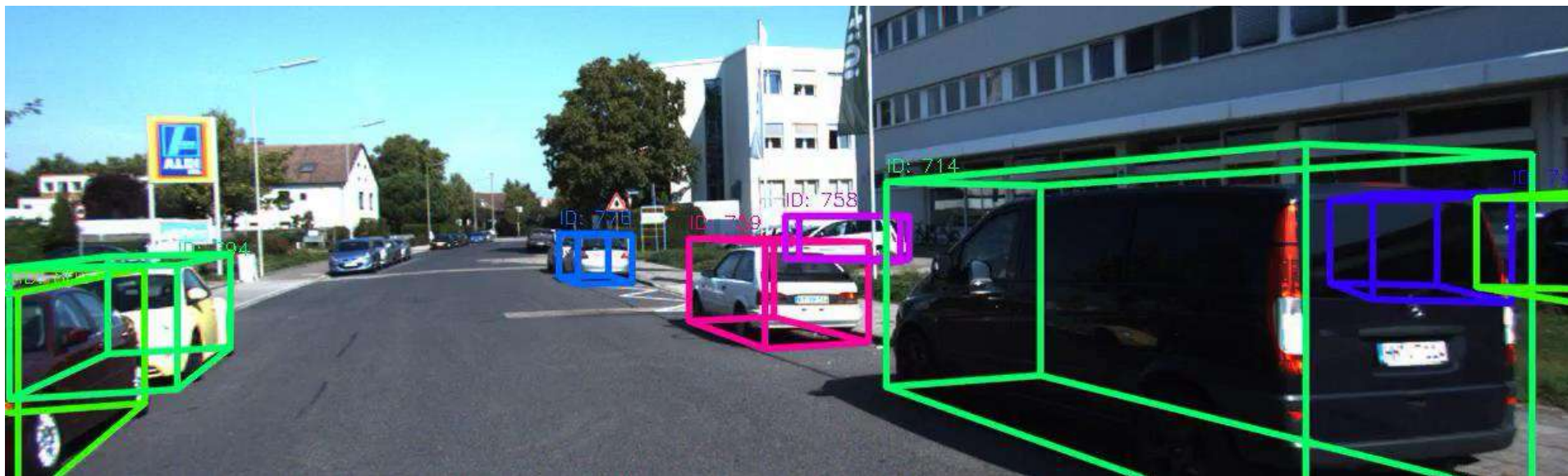
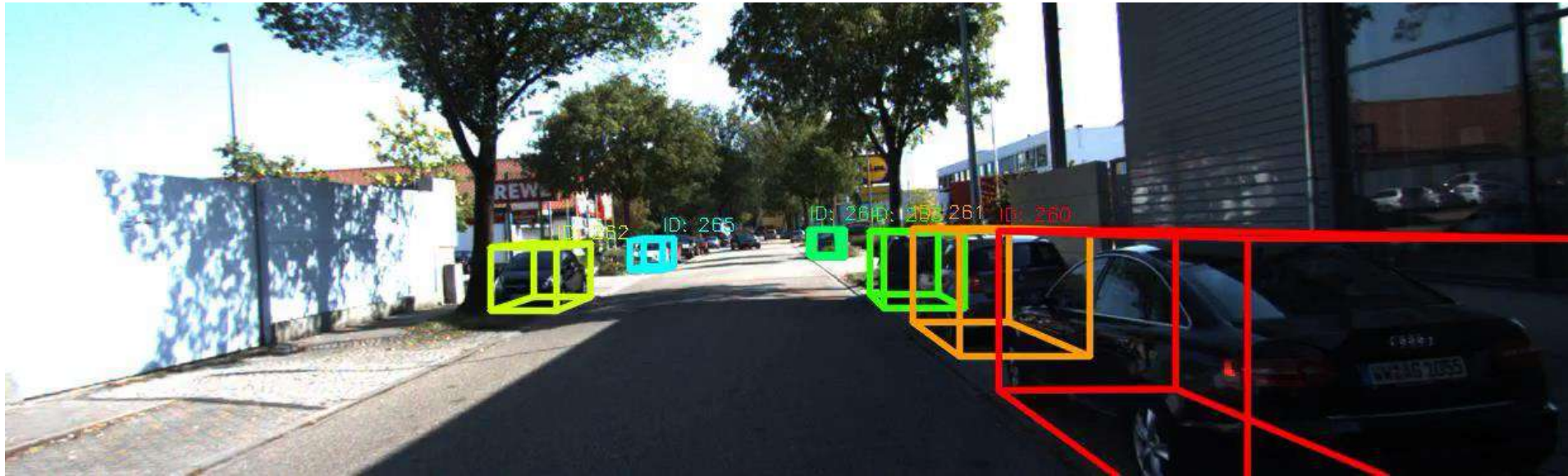
(b) Accuracy v.s. Number of Layers



Qualitative Results



Qualitative Results



Graph Neural Network for 3D Multi-Object Tracking

Xinshuo Weng, Yongxin Wang, Yunze Man, Kris Kitani
Robotics Institute, Carnegie Mellon University

European Conference on Computer Vision (ECCV) Workshops

