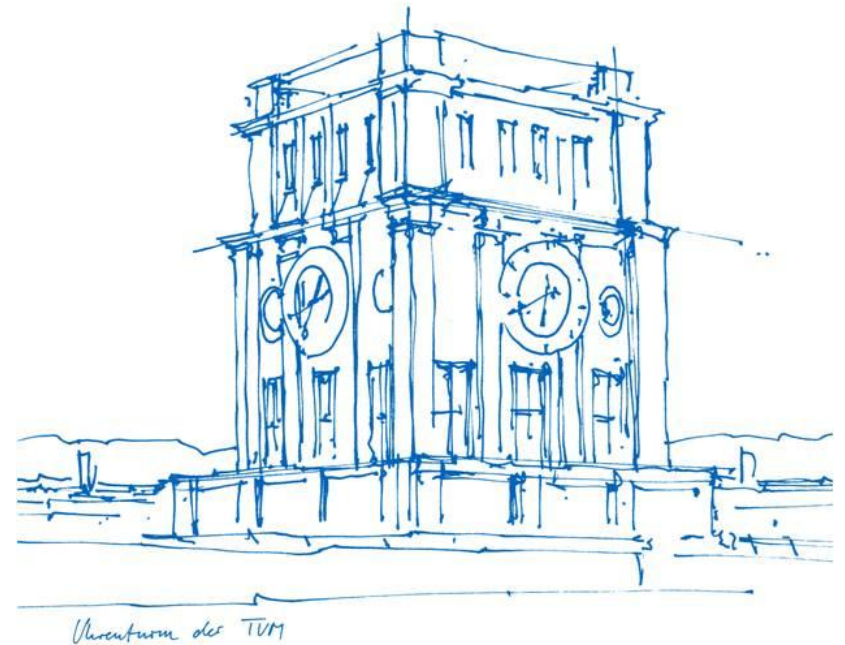# Learning-based Multi-Modal Perception

Timo Class

Technical University Munich

TUM School of Computation, Information and Technology

Smart Robotics Lab
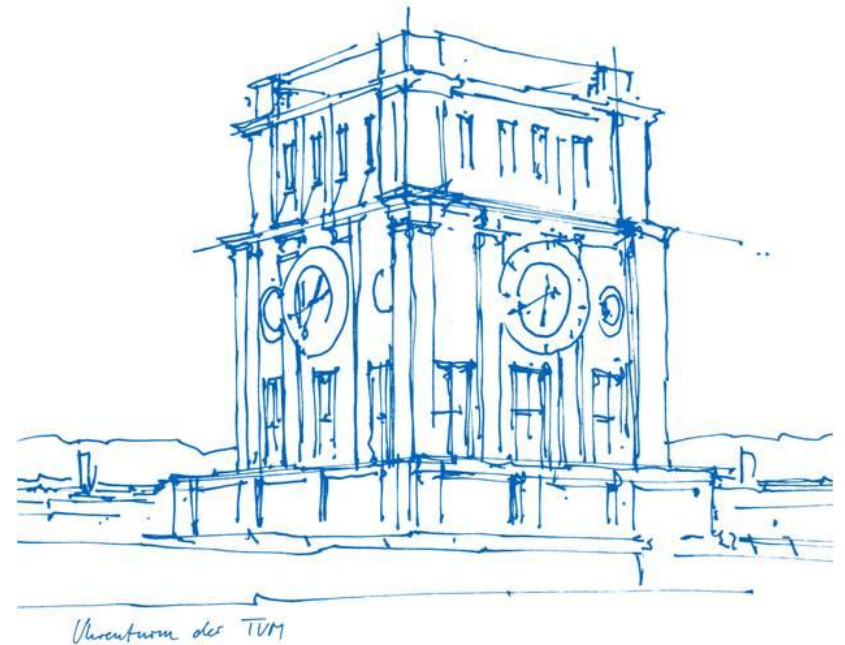
Seminar: Robot Perception & Intelligence

03 December 2024
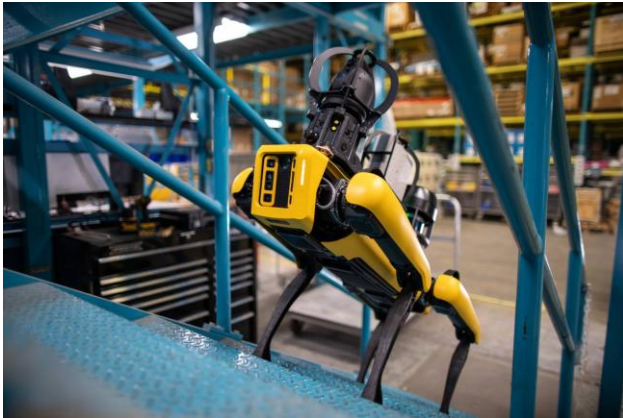
*Uhrenturm der TUM*

# Structure

- Motivation

- Overview

- Focused research

  - Related work

  - Method descriptions

  - Experiments and results

  - Shortcomings and future work

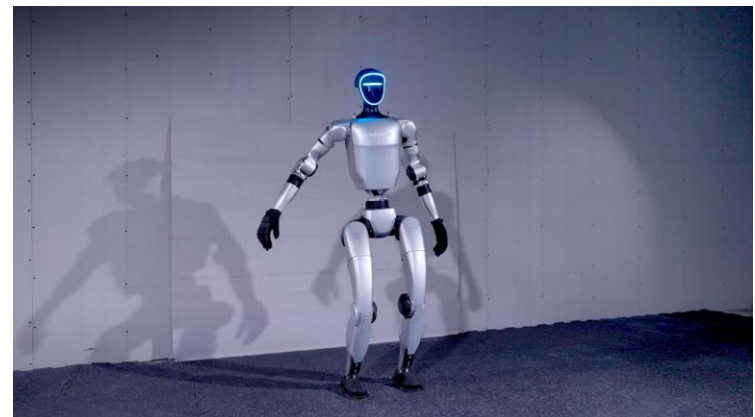- Conclusion

# Motivation


Boston Dynamics, Spot
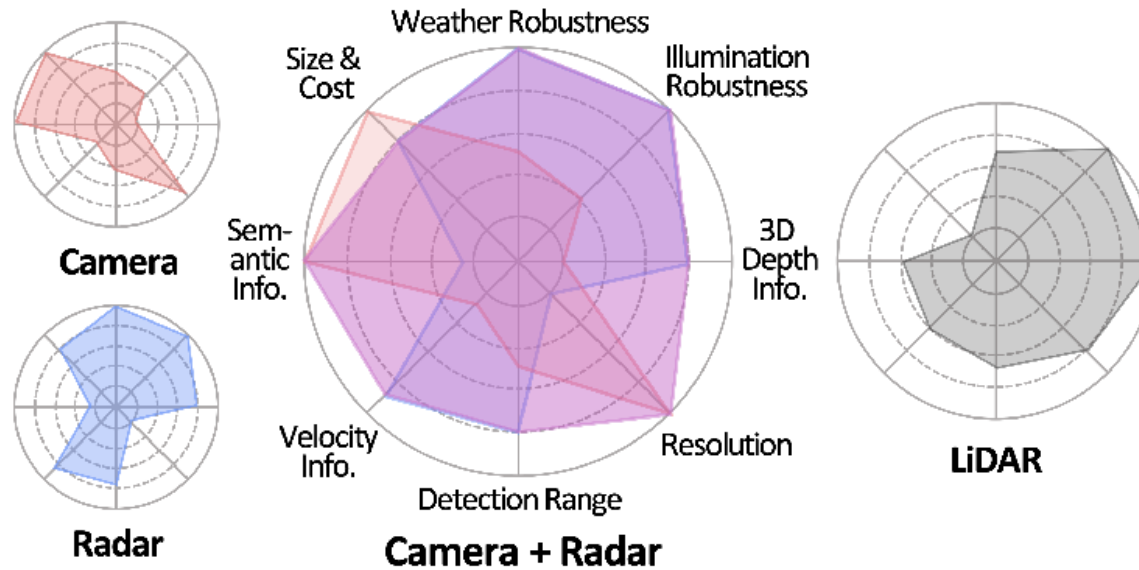

Waymo, Autonomous vehicle


Starship Technologies, Delivery robot


Unitree Robotics, Humanoid Robot G1

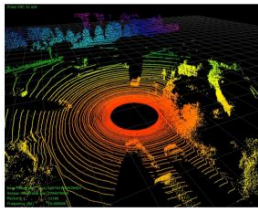# Motivation

Challenges:



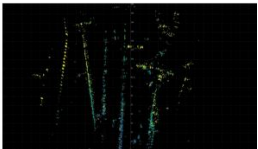Apoorv Singh, Vision-RADAR fusion for Robotics BEV Detections: A Survey

# Motivation



Camera images

LiDAR point cloud

Radar point cloud
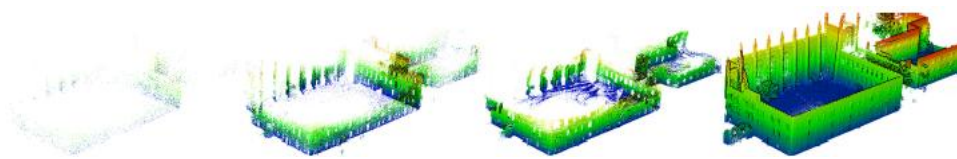
Encoder

Feature space

Decoder

Unified representation / perception task

Encoder-decoder network structure

# Overview

3D LiDAR Reconstruction



RTFNet



BEV Fusion

# Overview

3D LiDAR Reconstruction



RTFNet



BEV Fusion

# 3D Lidar Reconstruction

Method description

- **Problem:**
  - 3D reconstruction with sparse measurements results in incomplete reconstruction
    - → path planning and free-space estimation for autonomous navigation may fail



Legged robot scanning a building



Camera image with 16-channel LiDAR overlay

Sparse reconstruction →



16-channel LiDAR 3D reconstruction

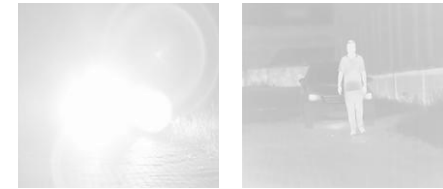# 3D Lidar Reconstruction
## Method description

- **Solution:**
  - Learning-based dense depth completion
  - Incorporate RGB images from three-camera setup (270°)
  - Incorporate learning-based depth uncertainty predictions



Camera setup on walking robot

# 3D Lidar Reconstruction

## Related work



Fangchang Ma, Sparse-to-Dense



Marija Popović, Volumetric Occupancy Mapping



Fehr Marius, Predicting Unobserved Space For Planning via Depth Map Augmentation

# 3D Lidar Reconstruction

## Method description

Linear sensor model: $\sigma(d)$

Depth filter:
$$\sigma_{pred} < 2\sigma(d)$$





Camera Image

Sparse Depth Image

Encoder

Depth Decoder

Uncertainty Decoder

Dense Depth Image

Uncertainty Image

Mapping System

Dense Reconstruction

System pipeline for 3D reconstruction
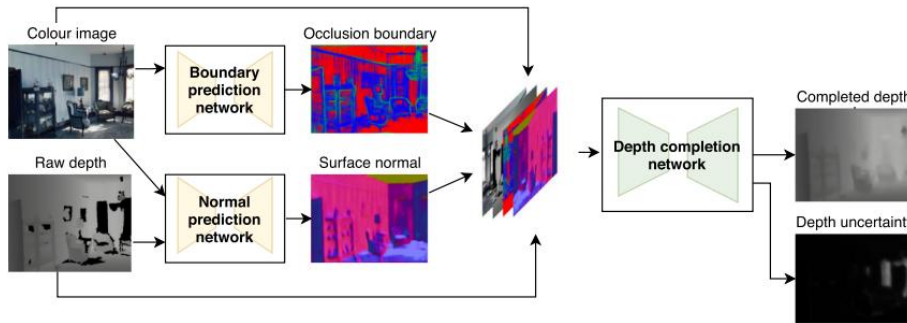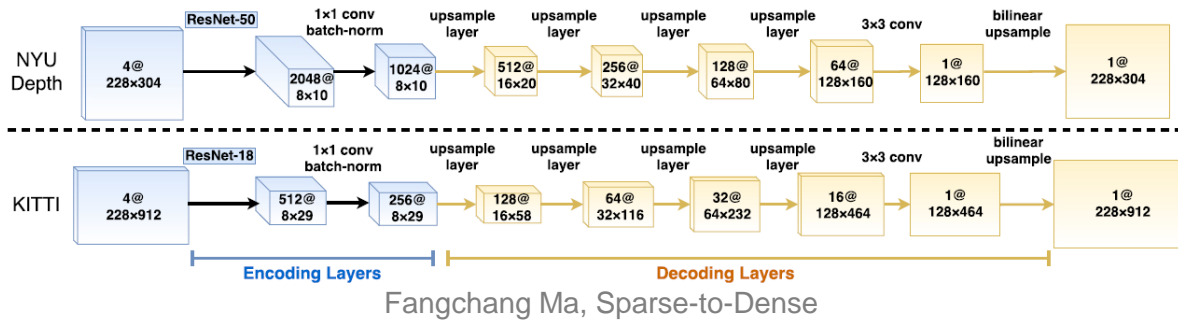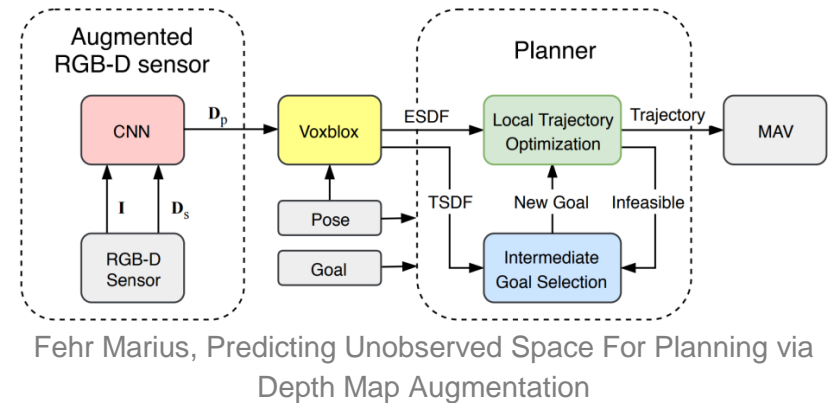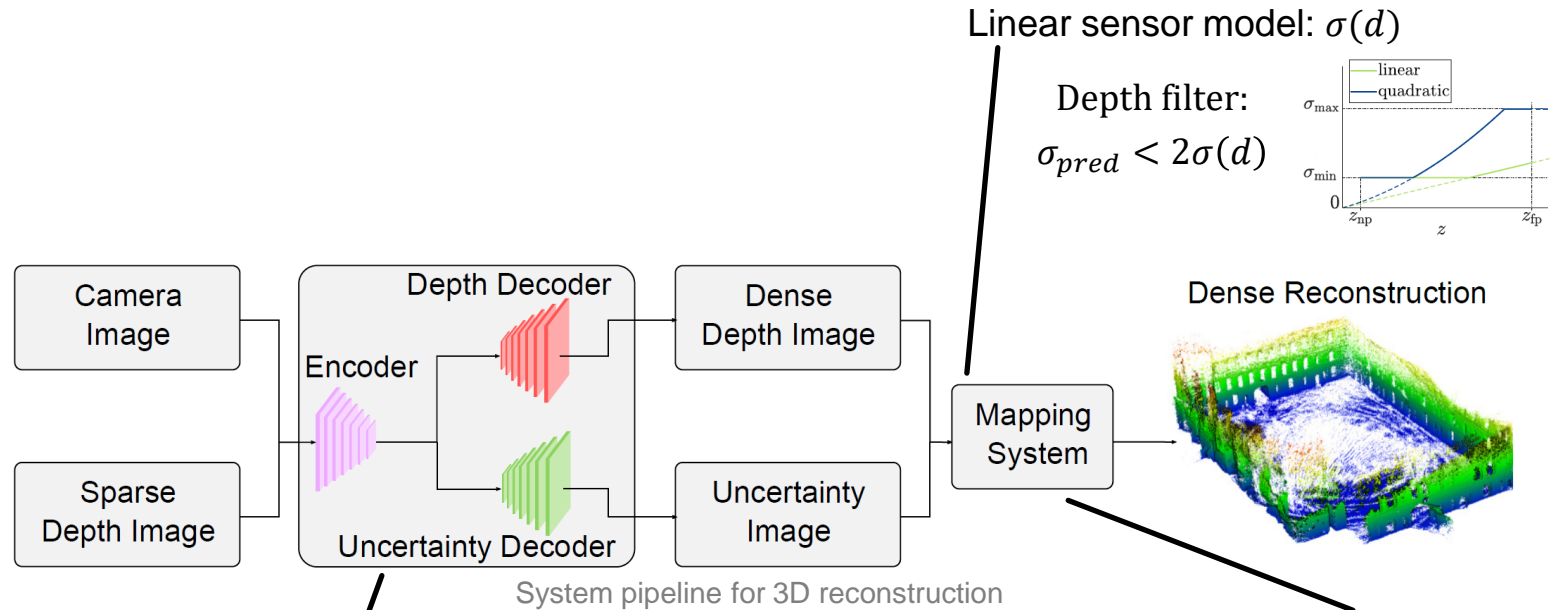


NYU Depth

$4@$ $228{\times}304$ → ResNet-50 → $2048@$ $8{\times}10$ → 1×1 conv batch-norm $1024@$ $8{\times}10$ → upsample layer $512@$ $16{\times}20$ → upsample layer $256@$ $32{\times}40$ → upsample layer $128@$ $64{\times}80$ → upsample layer $64@$ $128{\times}160$ → 3×3 conv $1@$ $128{\times}160$ → bilinear upsample $1@$ $228{\times}304$

KITTI

$4@$ $228{\times}912$ → ResNet-18 → $512@$ $8{\times}29$ → 1×1 conv batch-norm $256@$ $8{\times}29$ → upsample layer $128@$ $16{\times}58$ → upsample layer $64@$ $32{\times}116$ → upsample layer $32@$ $64{\times}232$ → upsample layer $16@$ $128{\times}464$ → 3×3 conv $1@$ $128{\times}464$ → bilinear upsample $1@$ $228{\times}912$

**Encoding Layers**    **Decoding Layers**

Fangchang Ma, Sparse-to-Dense (Depth prediction CNN network S2D)



Nils Funk, Multi-Resolution 3D Mapping

# 3D Lidar Reconstruction
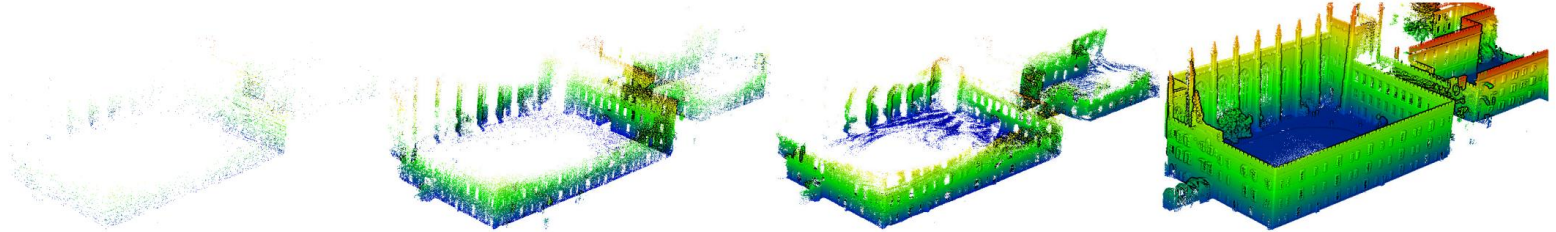
## Experiments and results

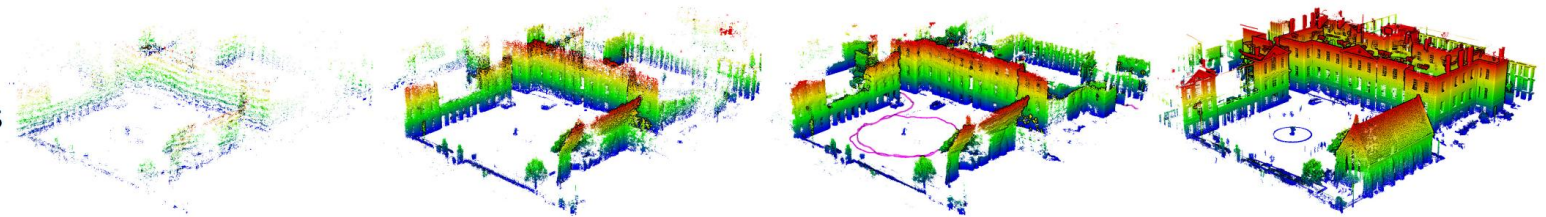Reconstruction & free-space estimation

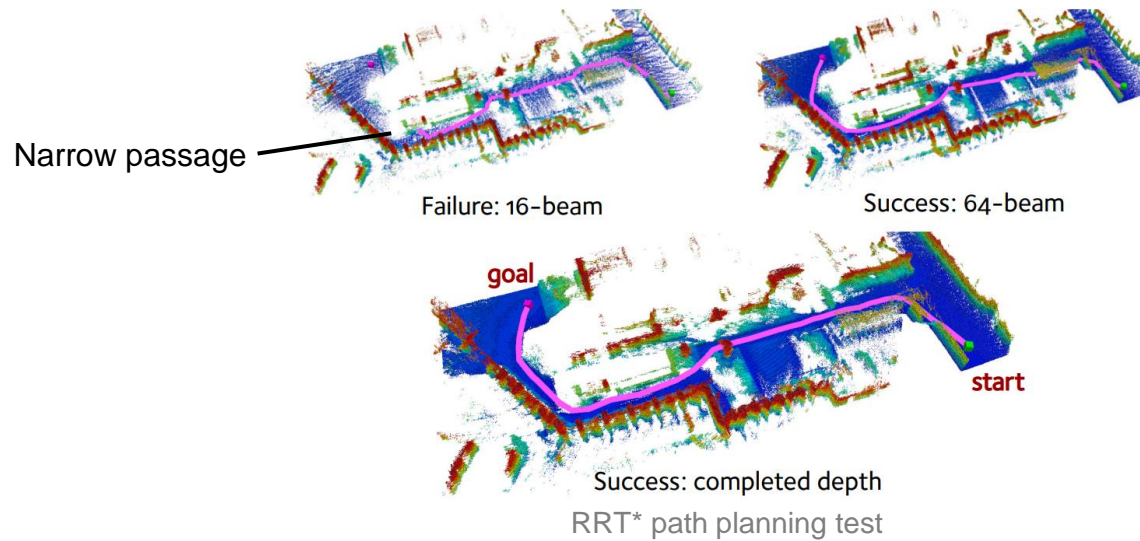Average point-to-point error: 0.2m

NCD

Maths Inst.

3D Reconstruction of NCD & Maths Inst. dataset for (left) 16-channel LiDAR, 64-channel LiDAR, Depth Completion, Ground Truth

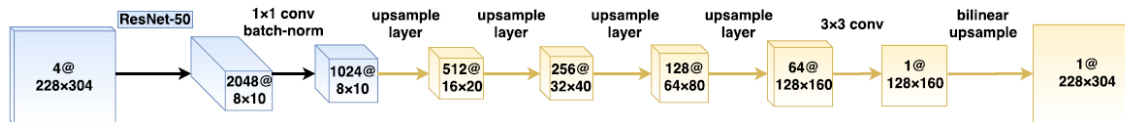# 3D Lidar Reconstruction

Experiments and results



Narrow passage

Failure: 16-beam

Success: 64-beam

goal

start

Success: completed depth

RRT* path planning test
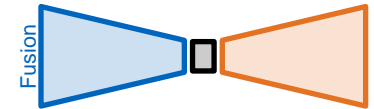
# 3D Lidar Reconstruction
## Shortcomings and future work

- **Shortcoming:**
  - Imbalanced modality representation in Encoder input
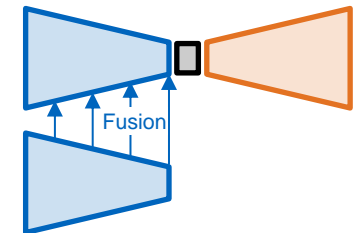  - RGB image potentially dominates input due to higher channel size



Fangchang Ma, Sparse-to-Dense (Depth prediction CNN network S2D)
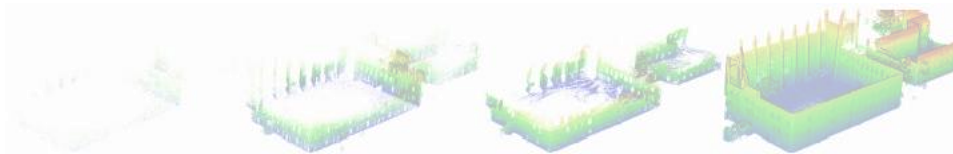
- **Solution:**
  - Add additional encoder
  - Separate feature extraction of RGB and depth image
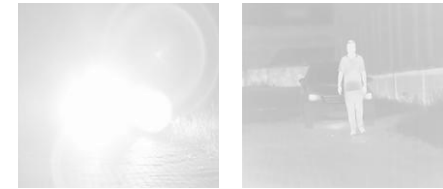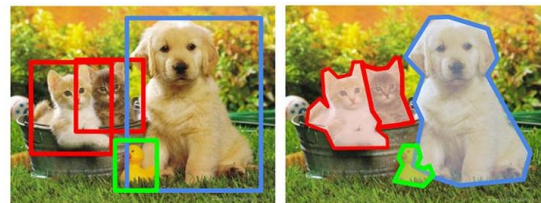  - Fuse features in final layer or simultaneously at each layer

# Overview

3D LiDAR Reconstruction

RTFNet

BEV Fusion

# BEVFusion: Multi-Task Multi-Sensor Fusion

Method description

- **Problem:**
  - Various sensors entail different data modalities
  - Various tasks entail different requirements

- **Solution:**
  - Transformation into unified representation

# BEVFusion: Multi-Task Multi-Sensor Fusion
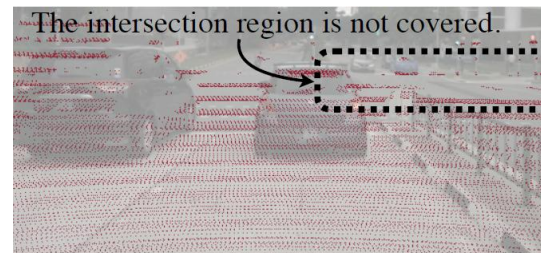## Method description

- **Solution:**
  - To camera → geometric-lossy

  
  close (red) and faraway (blue) points in **3D** are **neighbors** in **2D**

  - To LiDAR → semantic-lossy

  
  The intersection region is not covered.

  - To Birds-Eye-View
    → preserves geometric information
    → preserves semantic information

  
  **BEV features (camera)**

  
  **BEV features (LiDAR)**

# BEVFusion: Multi-Task Multi-Sensor Fusion
## Related work



Tianwei Yin, Center-based 3D Object Detection and Tracking



Tianwei Yin, Multimodal Virtual Point 3D Detection



Junjie Huang, BEVDet

# BEVFusion: Multi-Task Multi-Sensor Fusion

## Method description



Ze Liu. Swin Transformer

BEVFusion system pipeline

Yin Zhou, VoxelNet

# BEVFusion: Multi-Task Multi-Sensor Fusion

Experiments and results

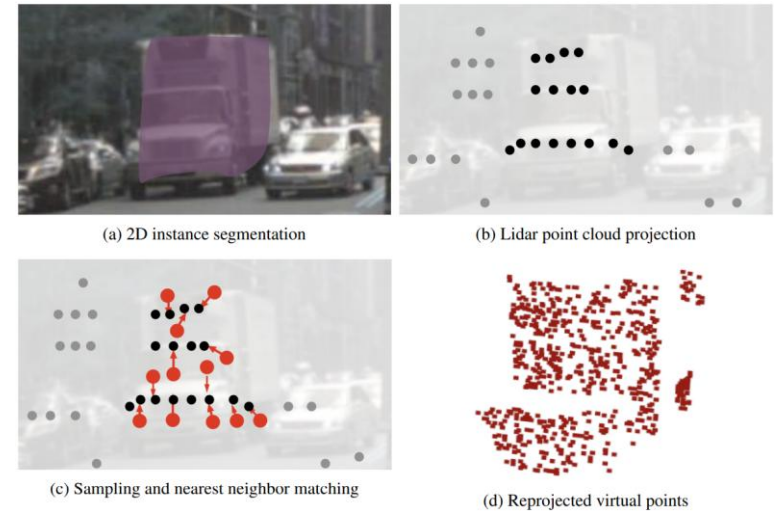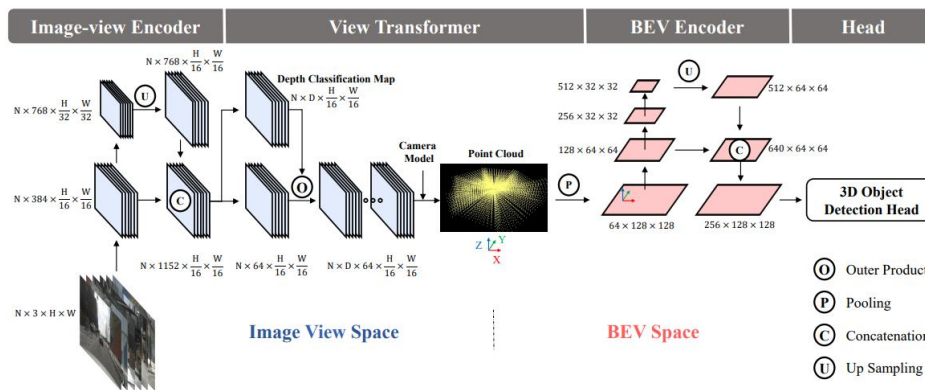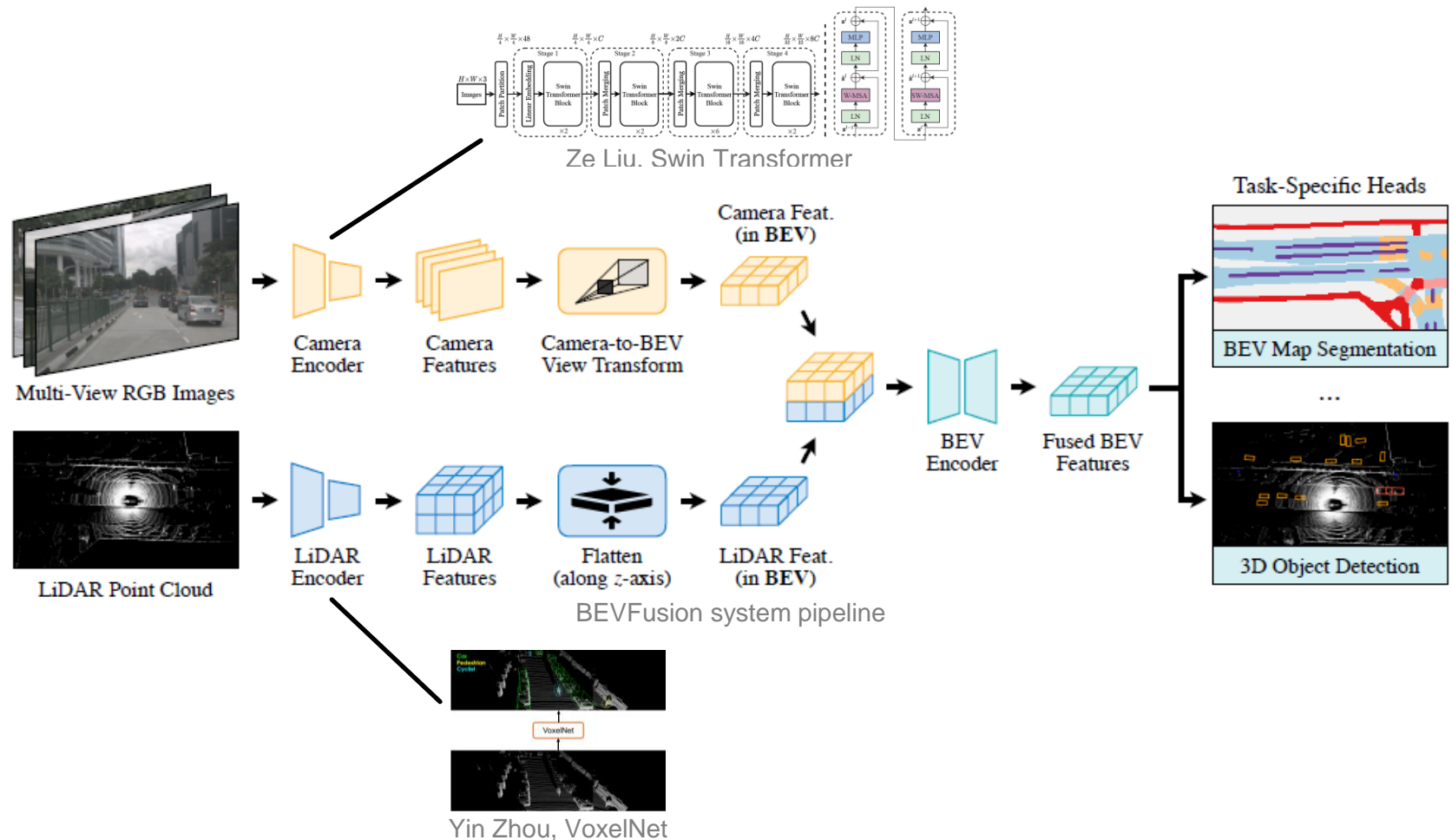| | Modality | Sunny | | Rainy | | Day | | Night | |
|---|---|---|---|---|---|---|---|---|---|
| | | mAP | mIoU | mAP | mIoU | mAP | mIoU | mAP | mIoU |
| CenterPoint | L | 62.9 | 50.7 | 59.2 | 42.3 | 62.8 | 48.9 | 35.4 | 37.0 |
| BEVDet/LSS* | C | 32.9 | 59.0 | 33.7 | 50.5 | 33.7 | 57.4 | 13.5 | 30.8 |
| MVP | C+L | 65.9 (+3.0) | 51.0 (-8.0) | 66.3 (+7.1) | 42.9 (-7.6) | 66.3 (+3.5) | 49.2 (-8.2) | 38.4 (+3.0) | 37.5 (+6.7) |
| BEVFusion | C+L | 68.2 (+5.3) | 65.6 (+6.6) | 69.9 (+10.7) | 55.9 (+5.4) | 68.5 (+5.7) | 63.1 (+5.7) | 42.8 (+7.4) | 43.6 (+12.8) |

=> Improved performance at nighttime

=> Improved performance in rainy weather

Performance analysis of BEVFusion under different weather and lightning conditions
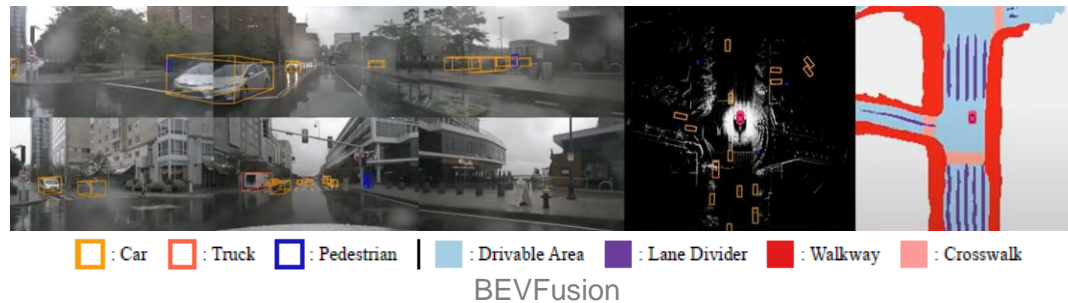
# BEVFusion: Multi-Task Multi-Sensor Fusion

Experiments and results

Missing detection



: Car  : Truck  : Pedestrian | : Drivable Area  : Lane Divider  : Walkway  : Crosswalk

Camera-only baseline

False positive detection

: Car  : Truck  : Pedestrian | : Drivable Area  : Lane Divider  : Walkway  : Crosswalk

LiDAR-only baseline

: Car  : Truck  : Pedestrian | : Drivable Area  : Lane Divider  : Walkway  : Crosswalk

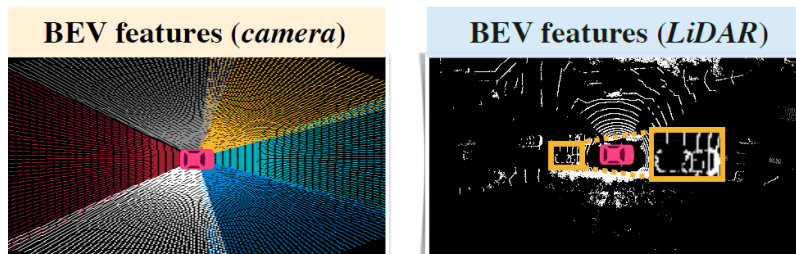BEVFusion

# BEVFusion: Multi-Task Multi-Sensor Fusion

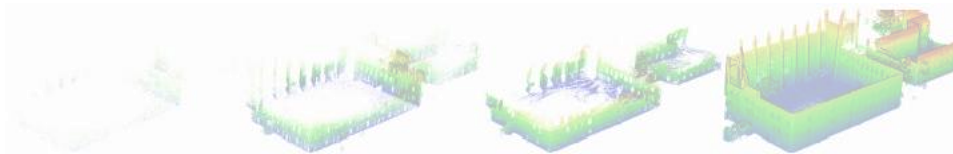Shortcomings and future work

- **Shortcoming:**
  - Loss of z-dimension due to transformation to Bird's-Eye View space



- **Solution:**
  - Multi-Scale BEV representation
  - Hybrid representation (BEV & 3D front-view features)
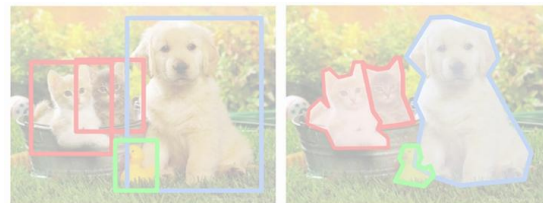  - Volumetric representation (Voxelization)

# Overview

3D LiDAR Reconstruction

RTFNet

BEV Fusion

# RTFNet: RGB-Thermal Fusion Network
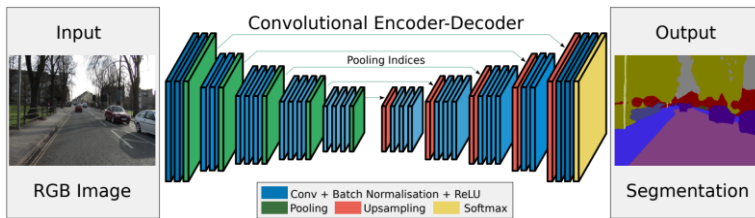## Method description

- **Problem:**
  - RGB camera performance is prone to lightning condition
  - Worse performance in total darkness or glaring situations
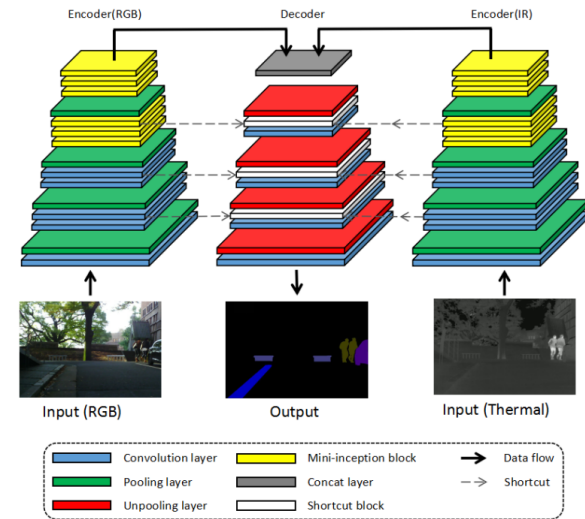
- **Solution:**
  - Incorporate thermal camera image



Comparison of RGB- and thermal image in a bright scene

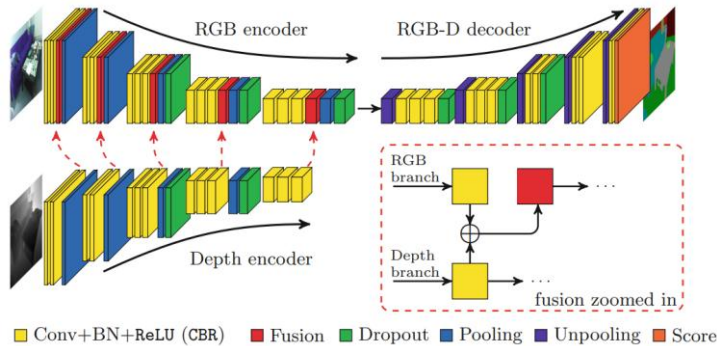# RTFNet: RGB-Thermal Fusion Network

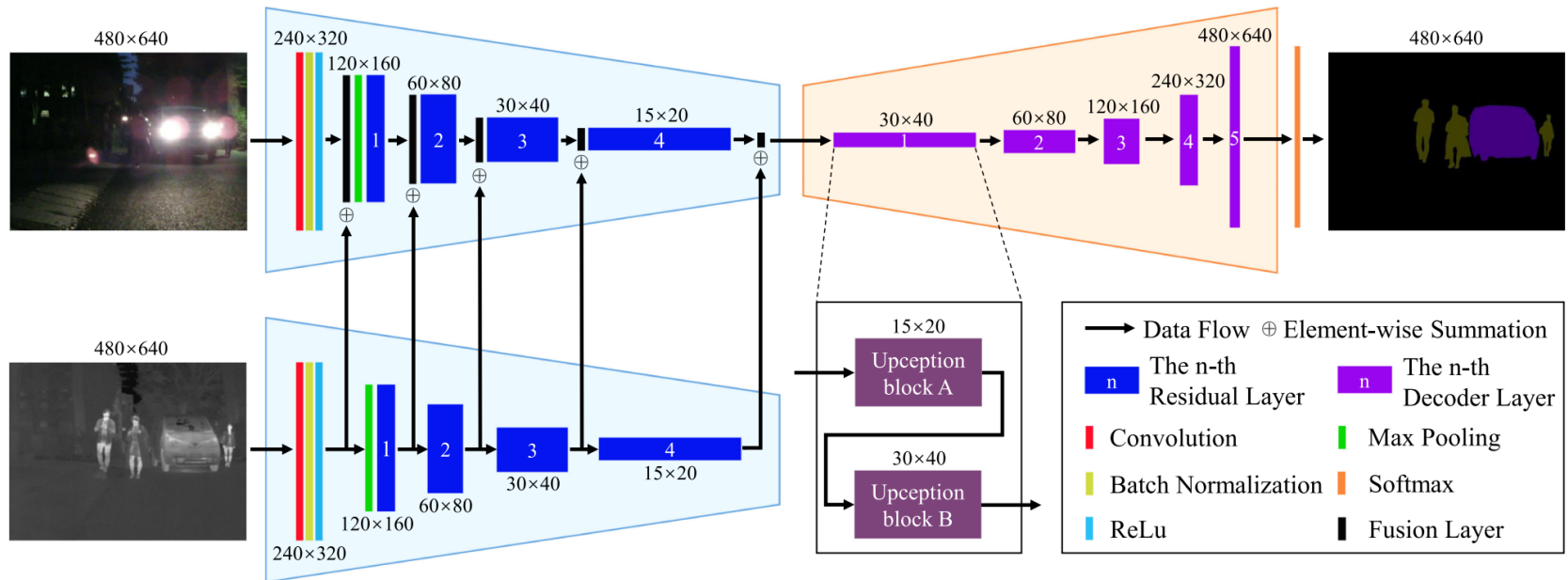## Related work



Vijay Badrinarayanan, SegNet



Caner Hazirbas, FuseNet



Qishen Ha, MFNet
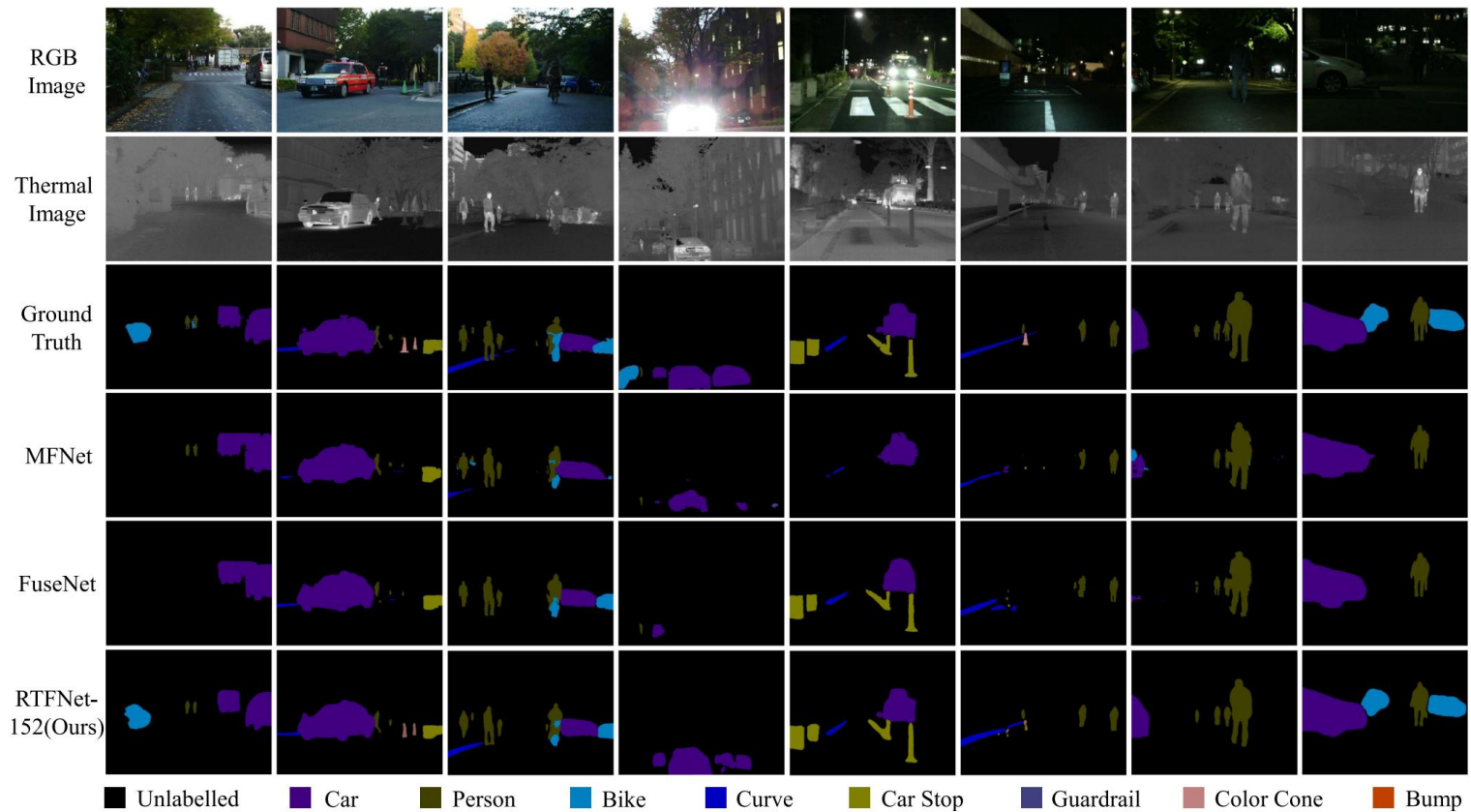
# RTFNet: RGB-Thermal Fusion Network

Method description



RTFNet network architecture

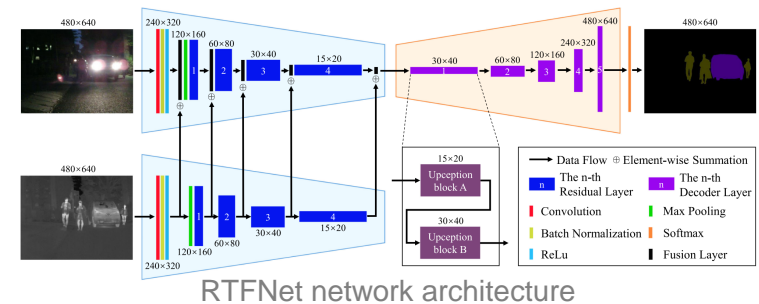# RTFNet: RGB-Thermal Fusion Network

Experiments and results



Qualitative comparison of data-fusion networks and RTFNet in different lightning conditions

# RTFNet: RGB-Thermal Fusion Network
## Shortcomings and future work
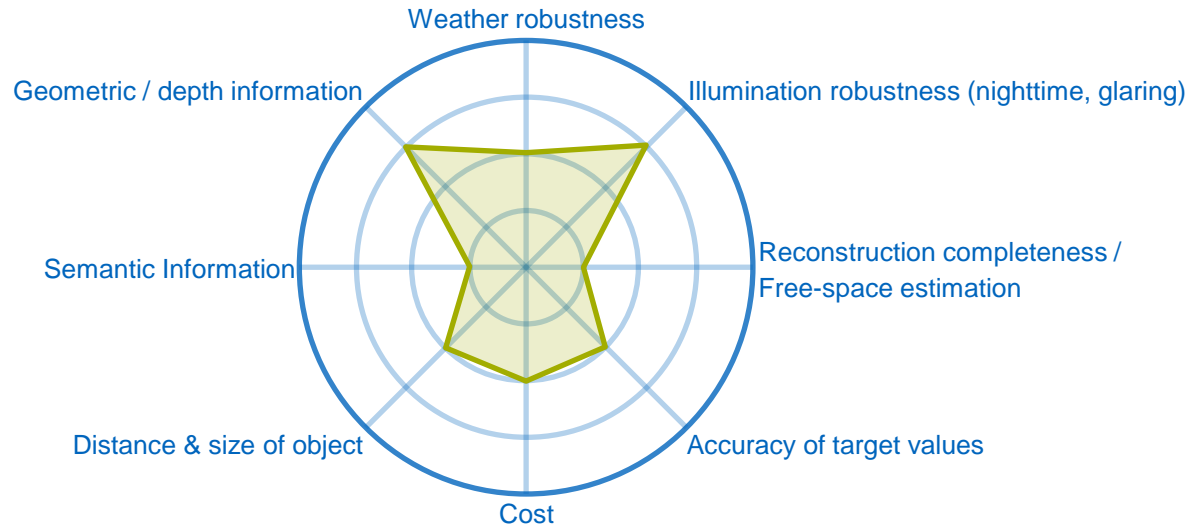
- **1st Shortcoming:**
  - Slow inference speed on embedded platforms
  - Large Encoder network

- **1st Solution:**
  - Reduce network size (especially Encoder)

- **2nd Shortcoming:**
  - Thermal images less informative when near objects share similar temperature

- **2nd Solution:**
  - Develop mechanisms to identify data that is more informative



RTFNet network architecture

# Conclusion



16-channel LiDAR

Radar chart axes:
- Weather robustness
- Illumination robustness (nighttime, glaring)
- Reconstruction completeness / Free-space estimation
- Accuracy of target values
- Cost
- Distance & size of object
- Semantic Information
- Geometric / depth information

# Conclusion



16-channel LiDAR
64-channel LiDAR

- Weather robustness
- Illumination robustness (nighttime, glaring)
- Reconstruction completeness / Free-space estimation
- Accuracy of target values
- Cost
- Distance & size of object
- Semantic Information
- Geometric / depth information

# Conclusion



16-channel LiDAR
64-channel LiDAR
RGB Camera

Weather robustness
Geometric / depth information
Illumination robustness (nighttime, glaring)
Reconstruction completeness / Free-space estimation
Semantic Information
Accuracy of target values
Distance & size of object
Cost

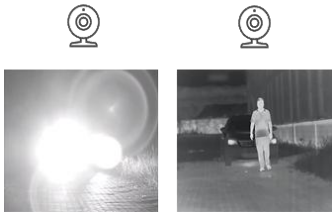# Conclusion

16-channel LiDAR
64-channel LiDAR
RGB Camera
Thermal camera

# Conclusion



RTFNet

16-channel LiDAR
64-channel LiDAR
RGB Camera
Thermal camera
Fusion approaches



Weather robustness

Illumination robustness (nighttime, glaring)

Geometric / depth information

Reconstruction completeness /
Free-space estimation

Semantic Information

Accuracy of target values

Distance & size of object

Cost

3D LiDAR Reconstruction

BEV Fusion