

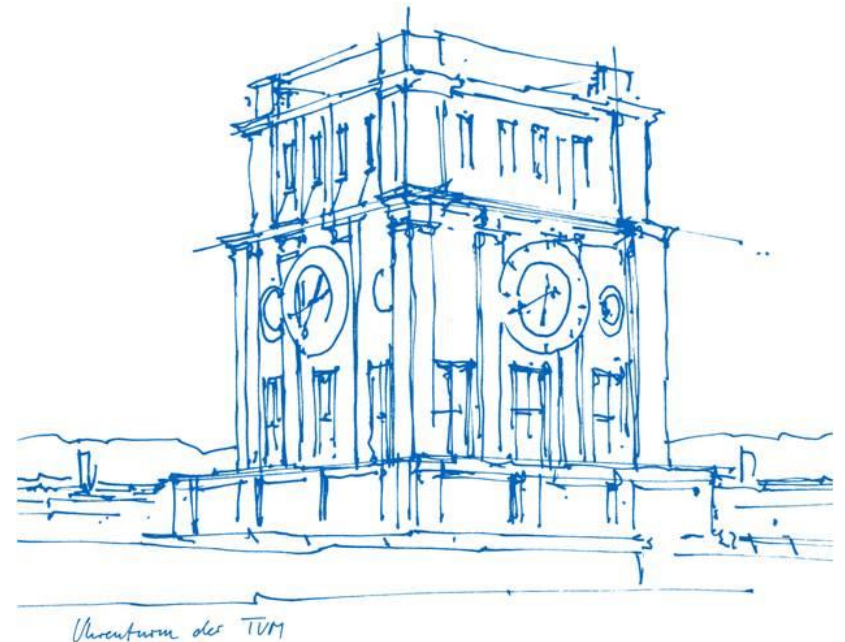
Learning-based Multi-modal Perception

Aleksandar Jevtić

Seminar: Robot Perception and Intelligence

Advisor: Dr. Jaehyung Jung

Munich, 16th of January 2024



Introduction

Learning-based Multi-modal Perception

*“The process of perception involves making useful **models of the environment** from a confusion mass of sensory input data”.*

- Semantic segmentation
- Object detection / tracking
- Pose estimation
- (...)

Introduction

Where do we need Machine Perception?

Many use relevant use cases, e.g.,

- Robotics
- Autonomous Vehicles
- Healthcare



Source: <https://www.cnet.com/home/this-robot-isnt-going-to-replace-your-in-home-nurse-yet/>
<https://venturebeat.com/ai/waymos-autonomous-cars-have-driven-20-million-miles-on-public-roads/>

Introduction

Learning-based Multi-modal Perception

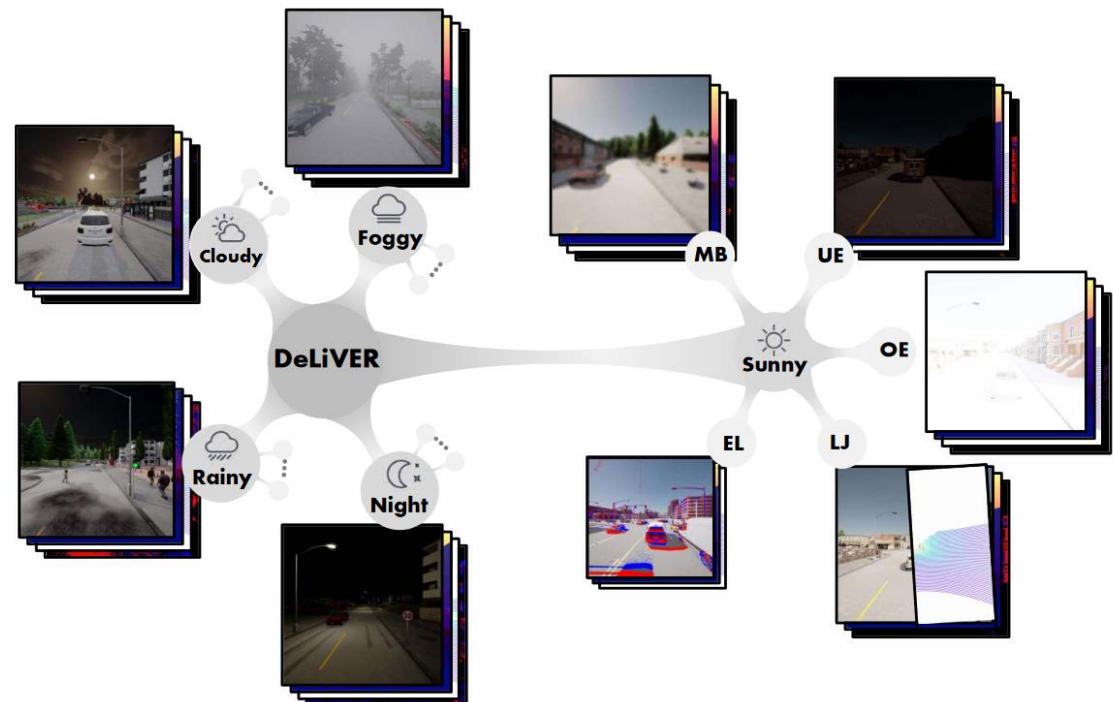
*“The term **multimodality** refers to an individual’s use of **different modes** (i.e. channels of communication) for the purpose of conveying meaning.”*

- RGB images *(non-structural)*
- Depth
 - dense
 - sparse (e.g. LiDAR)
- Thermal imaging
- IMU
- Audio
- Language

Introduction

How does multimodality help?

- Better accuracy
- Robustness
 - Adverse conditions
 - Failure cases



Overview

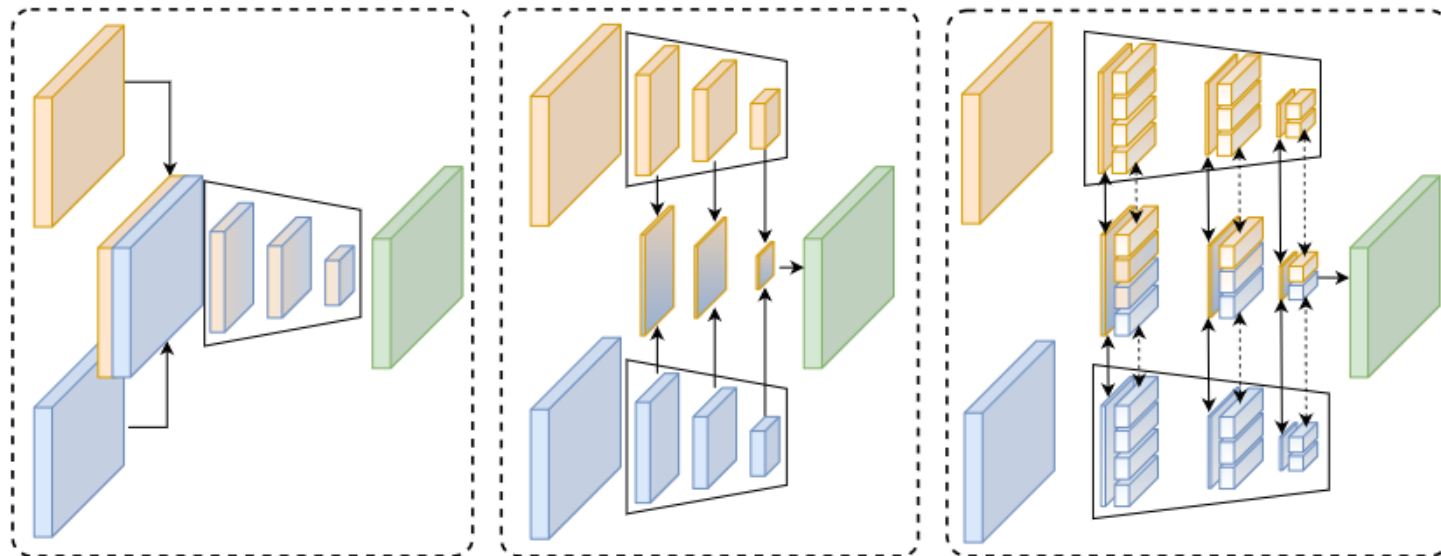
- Introduction and Overview
- Related Work
- Method Descriptions and Results
 - Multi-modal curb detection
 - CMX and CMNeXt
 - Multi-modal knowledge expansion
- Personal Comments
- Future Work

Related Work

Cross-Modal Fusion

(e.g. ShapeConv)

(e.g. SA-Gate)



(a) Input fusion

(b) Feature fusion

(c) Interactive fusion

Source: J. Cao, H. Leng, D. Lischinski, D. Cohen-Or, C. Tu, and Y. Li, "ShapeConv: Shape-Aware Convolutional Layer for Indoor RGB-D Semantic Segmentation," ICCV, 2021.

X. Chen et al., "Bi-directional Cross-Modality Feature Propagation with Separation-and-Aggregation Gate for RGB-D Semantic Segmentation," ECCV, 2020.

Related Work

High-level: Leveraging RGB data and models

Related concept: Semi-supervised Learning

- Consistency regularization
Small input and model perturbations → small output changes
Additional loss term
- Pseudo-labeling
Teacher – Student Architecture
Generation of labels for unlabeled data

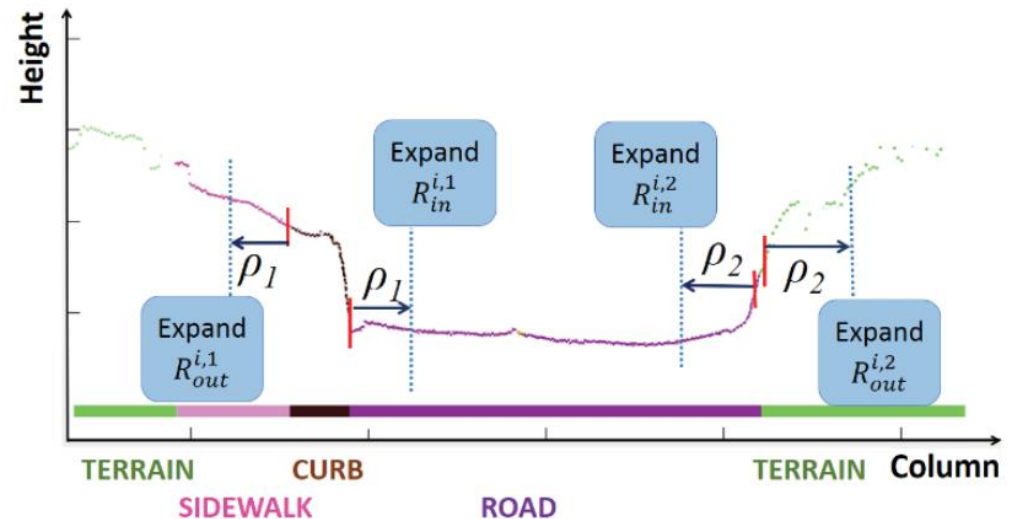
Related Work

Curb Detection Methods with LiDAR

Segmentation in RGB to find Regions of Interest

in these ROIs:
use engineered spatial features

Segmentation is learning-based,
but not fusion!



Multi-modal curb detection and filtering

Sandipan Das^{1,2}, Navid Mahabadi², Saikat Chatterjee¹, Maurice Fallon³

¹ KTH EECS, Sweden. {sandipan, sach}@kth.se

² Scania, Sweden. {sandipan.das, navid.mahabadi}@scania.com

³ Oxford Robotics Institute, UK. mfallon@robots.ox.ac.uk

Multi-modal curb detection and filtering

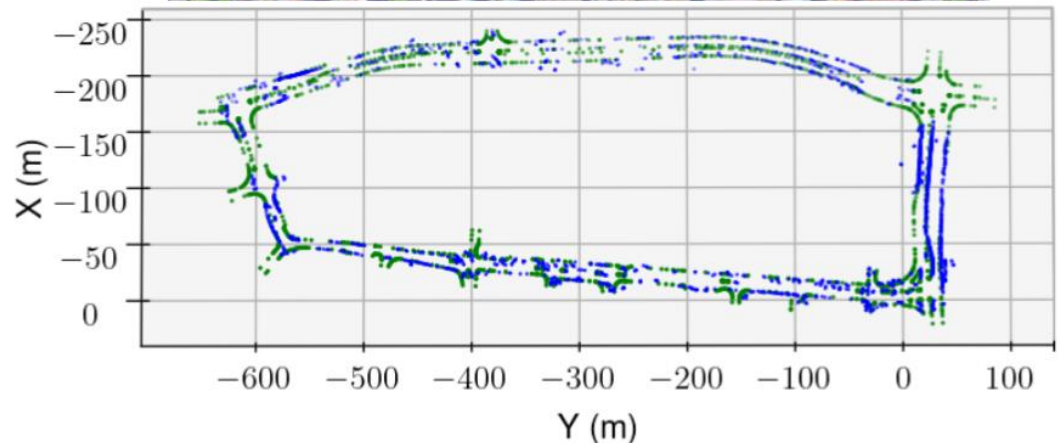
Detection of curb points by
unsupervised clustering

Multi-modal fusion of

- RGB
- LiDAR

Data collection vehicle

- 4 sensors
- varying FoVs



Detected curb features (blue) and ground truth (green)

Multi-modal curb detection and filtering

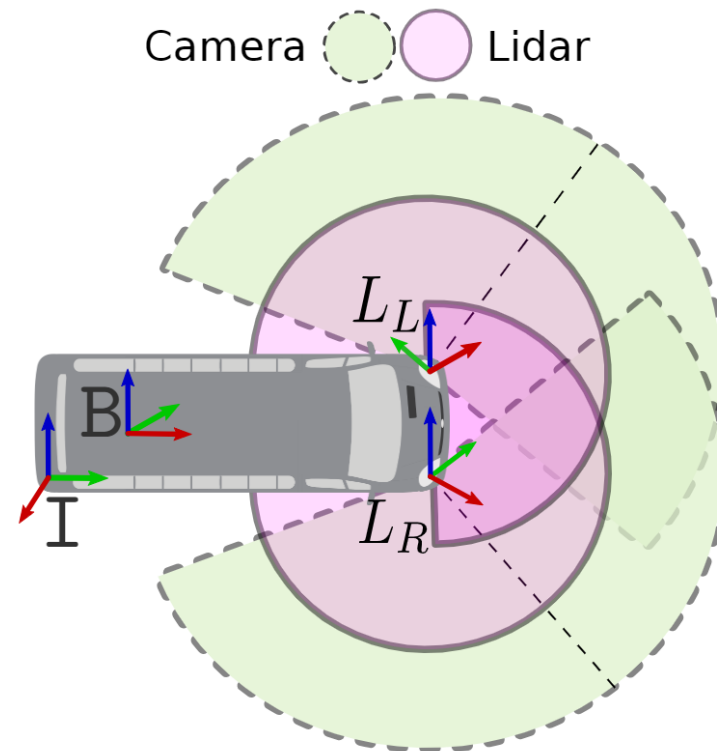
Detection of curb points by
unsupervised clustering

Multi-modal fusion of

- RGB
- LiDAR

Data collection vehicle

- 4 sensors
- varying FoVs



Multi-modal Curb detection - Method

Curb segmentation on RGB

- EfficientNet

Association with LiDAR

Unsupervised clustering

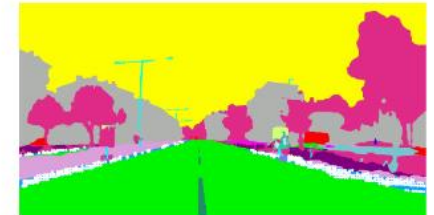
- DBSCAN (density-based)

Filtering

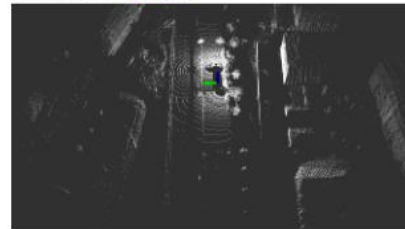
- RANSAC filtering
- Delaunay filtering



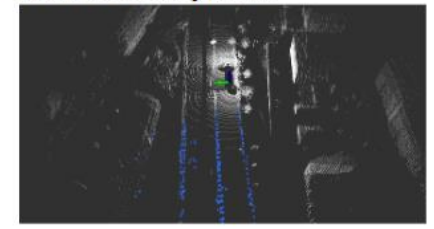
(a) Semantic segmentation results using our modified EfficientNet [18].



(c) Lidar point clouds (white points) overlaid on the segmented curb pixels.



(b) Fused lidar point clouds from lidar sensors.



(d) Curb semantics (blue points) with the fused point cloud.

Multi-modal Curb detection - Method

Curb segmentation on RGB

- EfficientNet

Association with LiDAR

Unsupervised clustering

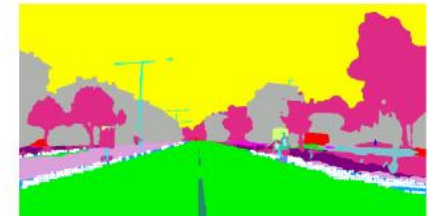
- DBSCAN (density-based)

Filtering

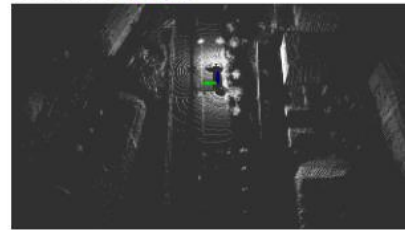
- RANSAC filtering
- Delaunay filtering



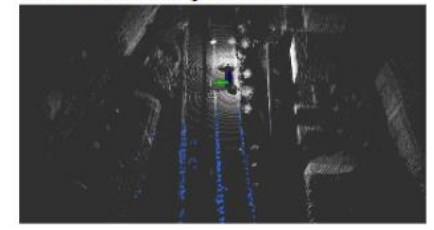
(a) Semantic segmentation results using our modified EfficientNet [18].



(c) Lidar point clouds (white points) overlaid on the segmented curb pixels.

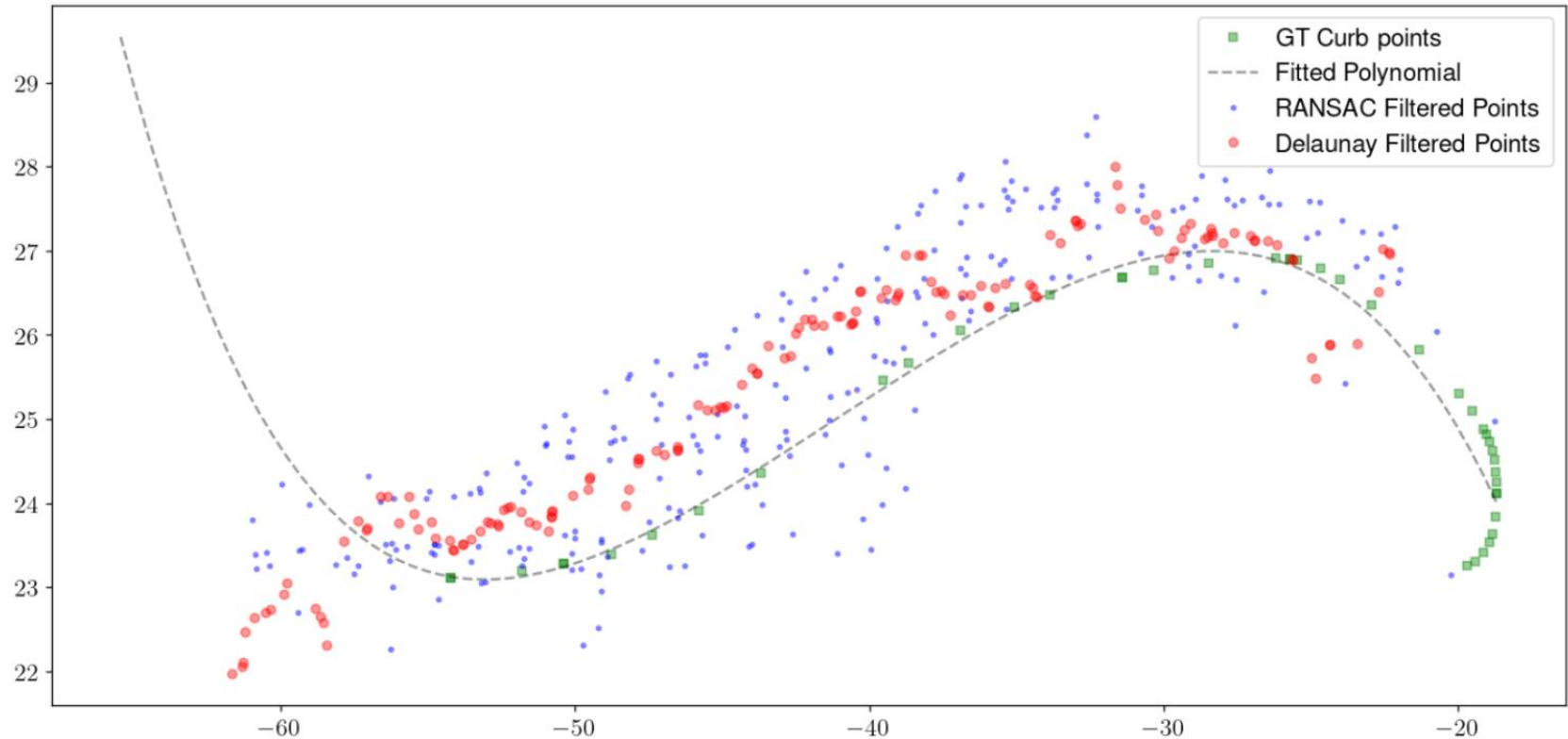


(b) Fused lidar point clouds from lidar sensors.



(d) Curb semantics (blue points) with the fused point cloud.

Multi-modal Curb detection - Results



Source: S. Das, N. Mahabadi, S. Chatterjee, and M. Fallon, "Multi-modal curb detection and filtering," *CoRR*, vol. abs/2205.07096, 2022.

Multi-modal Curb detection - Results

| Manual segment-wise association | | |
|------------------------------------|------------------------|-------------------|
| No Clustering | Normalized L_2 -Norm | # Detected Points |
| RANSAC Filtering | 27.659 | 9578 |
| Delaunay Filtering | 19.947 | 6904 |
| Automatic segment-wise association | | |
| Outlier Removal (RANSAC) | Chamfer Distance | # Detected Points |
| Agglomerative Clustering | 17.427 | 3489 |
| BIRCH | 19.596 | 1351 |
| DBSCAN | 17.220 | 5314 |
| OPTICS | 18.370 | 7446 |
| Outlier Removal (Delaunay) | Chamfer Distance | # Detected Points |
| Agglomerative Clustering | 15.418 | 3924 |
| BIRCH | 16.165 | 3492 |
| DBSCAN | 14.753 | 6678 |
| OPTICS | 15.870 | 4415 |

Source: S. Das, N. Mahabadi, S. Chatterjee, and M. Fallon, "Multi-modal curb detection and filtering," *CoRR*, vol. abs/2205.07096, 2022.

CMX: Cross-Modal Fusion for RGB-X Semantic Segmentation with Transformers

Jiaming Zhang*, Huayao Liu*, Kailun Yang*[†], Xinxin Hu, Ruiping Liu, and Rainer Stiefelhagen

J. Zhang, R. Liu, and R. Stiefelhagen are with Karlsruhe Institute of Technology, 76131 Karlsruhe, Germany.

K. Yang is with Hunan University, Changsha 410082, China.

H. Liu is with NIO, Shanghai 201804, China.

X. Hu is with ByteDance Inc., Hangzhou 310000, China.

*indicates equal contribution.

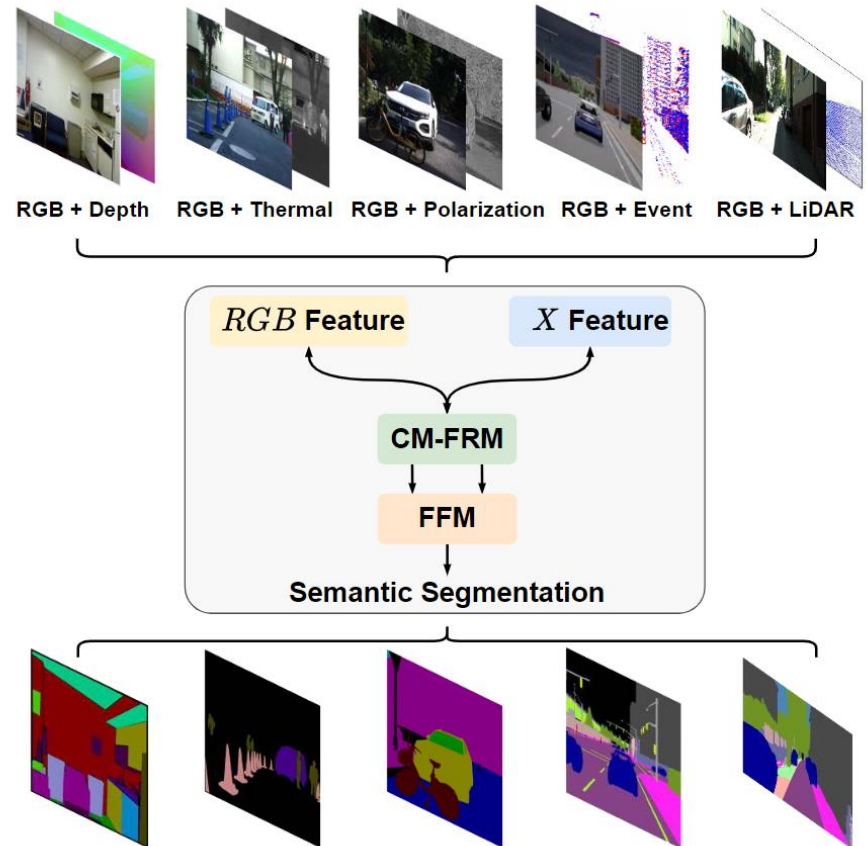
[†]corresponding author. (E-Mail: kailun.yang@hnu.edu.cn.)

CMX: Cross-Modal Fusion for RGB-X

Unified fusion framework

RGB-X semantic segmentation

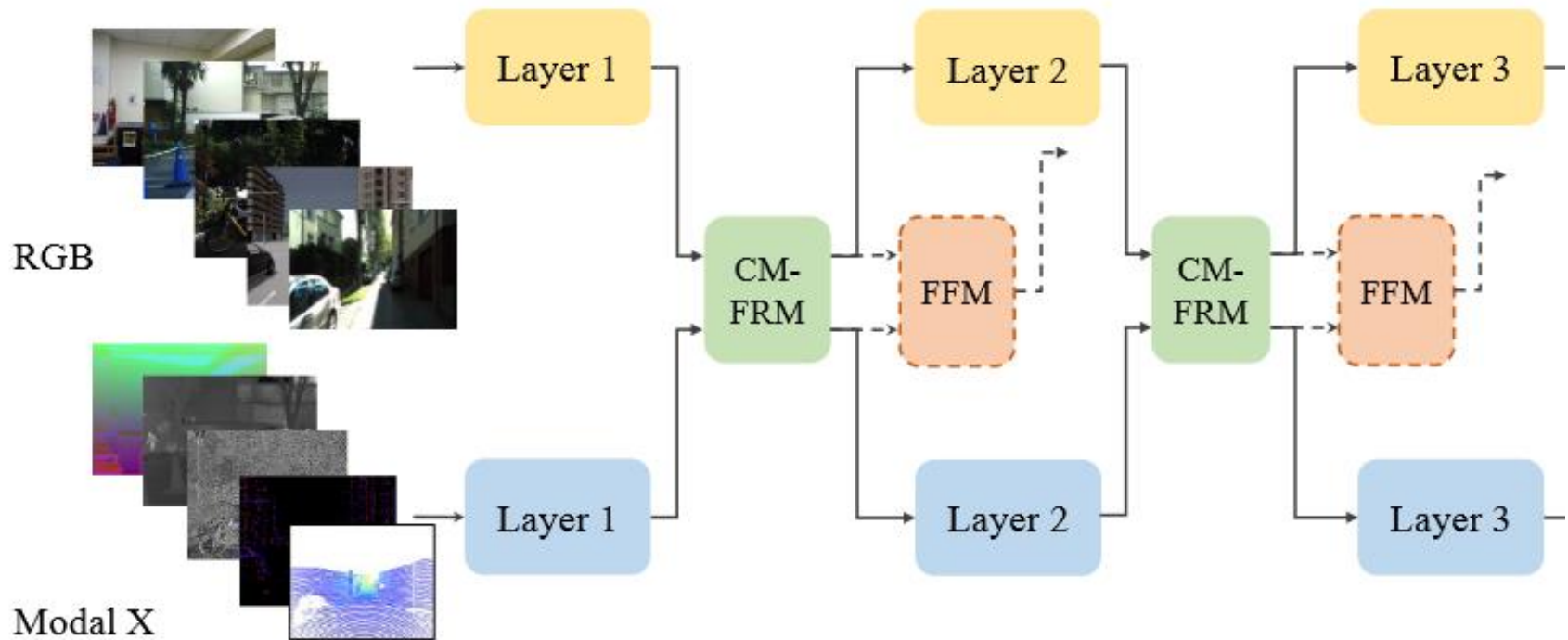
Attention mechanisms enable efficient fusion



CMX - Method

Overall Framework

| | |
|--------|-----------------------------------|
| Layer | Mix Transformer (MiT) |
| CM-FRM | Cross-modal feature rectification |
| FFM | Feature fusion |



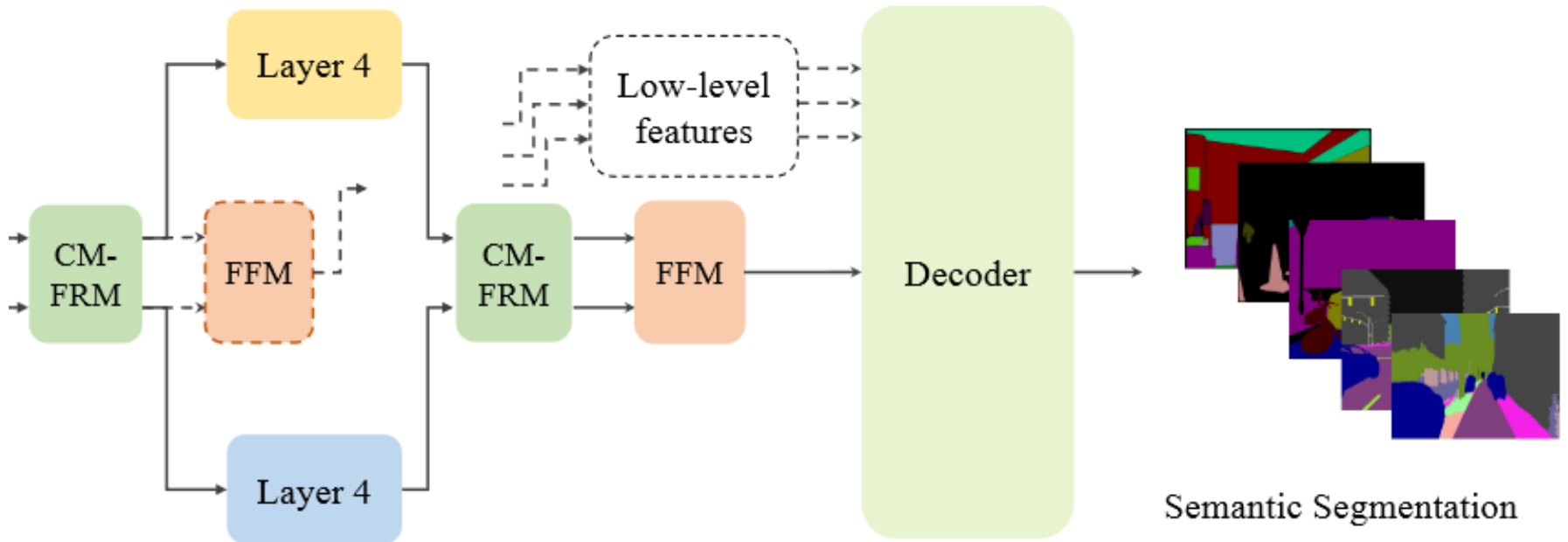
Source: E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "SegFormer: Simple and Efficient Design for Semantic Segmentation with Transformers," NeurIPS, 2021.

J. Zhang, H. Liu, K. Yang, X. Hu, R. Liu, and Rainer Stiefelhagen, "CMX: Cross-Modal Fusion for RGB-X Semantic Segmentation With Transformers," IEEE Transactions on Intelligent Transportation Systems, 2022.

CMX - Method

Overall Framework

| | |
|--------|-----------------------------------|
| Layer | Mix Transformer (MiT) |
| CM-FRM | Cross-modal feature rectification |
| FFM | Feature fusion |

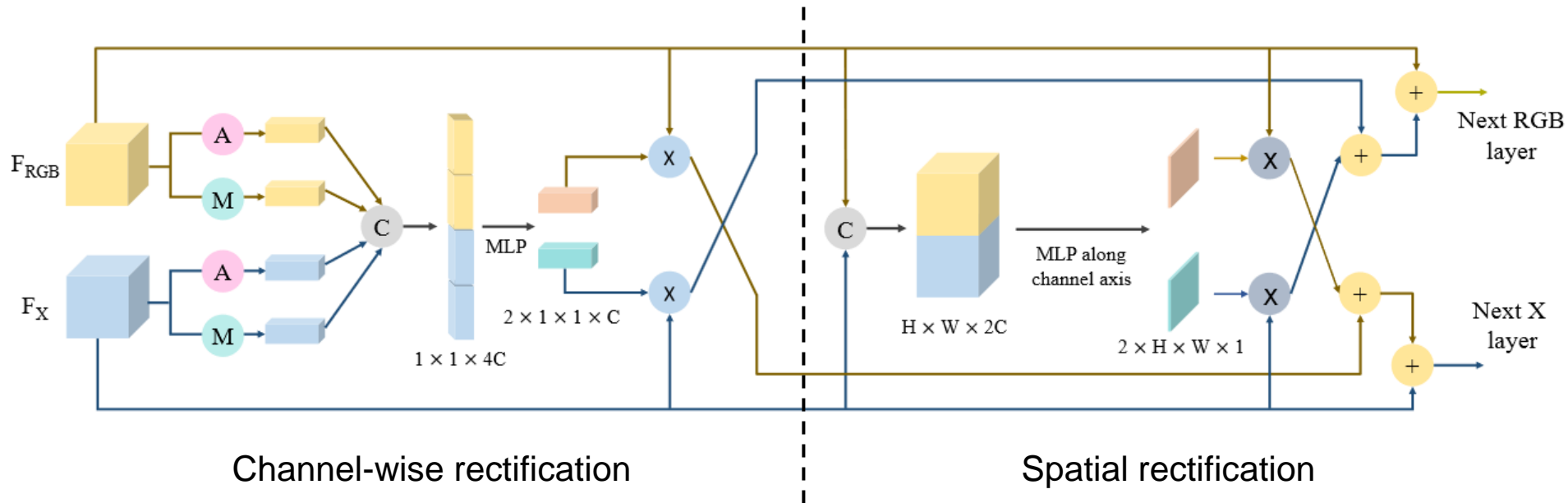


Source: E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "SegFormer: Simple and Efficient Design for Semantic Segmentation with Transformers," NeurIPS, 2021.

J. Zhang, H. Liu, K. Yang, X. Hu, R. Liu, and Rainer Stiefelhagen, "CMX: Cross-Modal Fusion for RGB-X Semantic Segmentation With Transformers," IEEE Transactions on Intelligent Transportation Systems, 2022.

CMX - Method

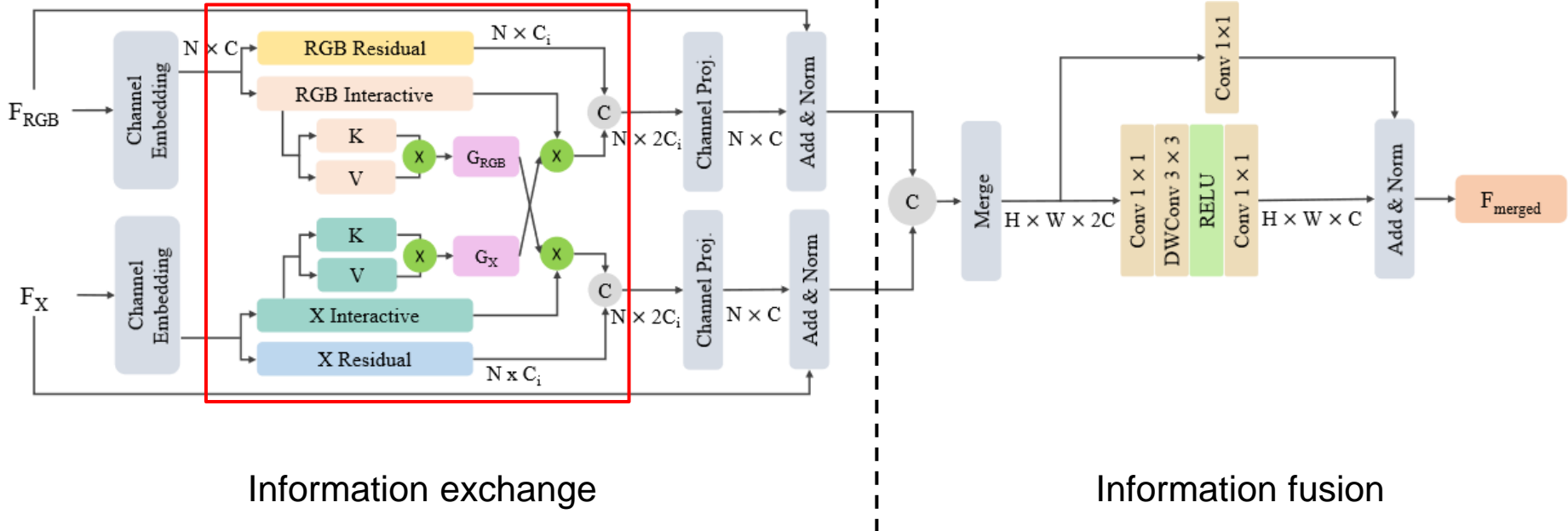
Cross-modal feature rectification module (CM-FRM)



CMX - Method

Feature fusion module (FFM)

Efficient attention across modes



Delivering Arbitrary-Modal Semantic Segmentation

Jiaming Zhang^{1,*}, Ruiping Liu^{1,*}, Hao Shi³, Kailun Yang^{2,†}, Simon Reiß¹,
Kunyu Peng¹, Haodong Fu⁴, Kaiwei Wang³, Rainer Stiefelhagen¹

¹Karlsruhe Institute of Technology, ²Hunan University, ³Zhejiang University, ⁴Beihang University

*Equal contribution.

†Corresponding author (e-mail: kailun.yang@hnu.edu.cn).

¹The DELIVER dataset and our code will be made publicly available
at: <https://jamycheung.github.io/DELIVER.html>.

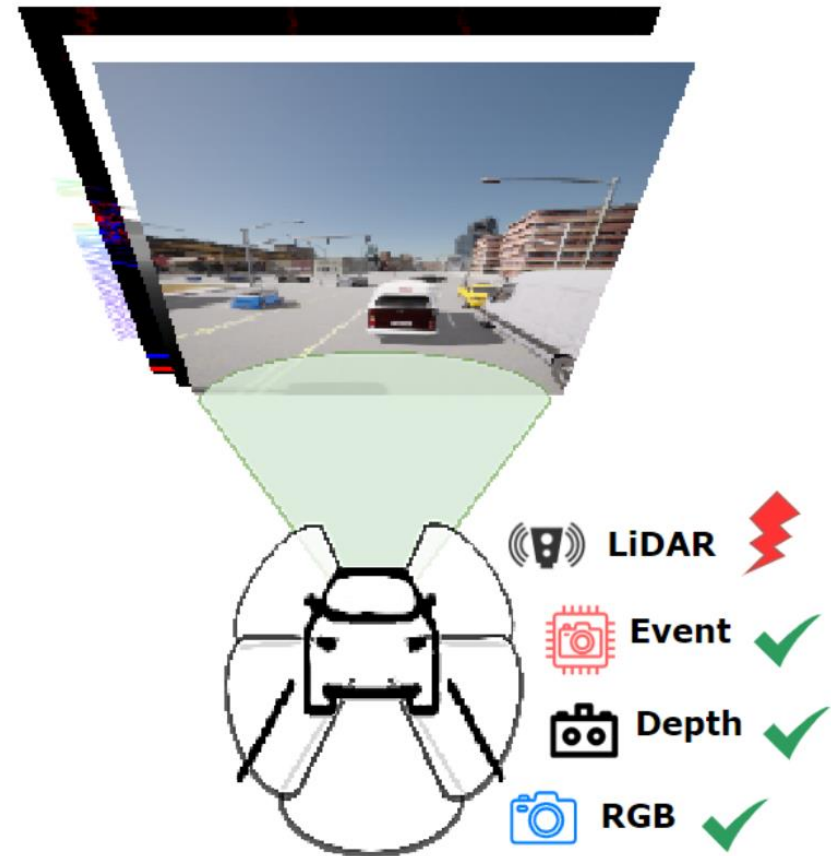
CMNeXt: Arbitrary-Modal Fusion

Extending CMX

- Multiple additional modalities
- Retains two-stream architecture

Synthetic dataset DeLiVER

- Depth
- LiDAR
- Multiple Views
- Event



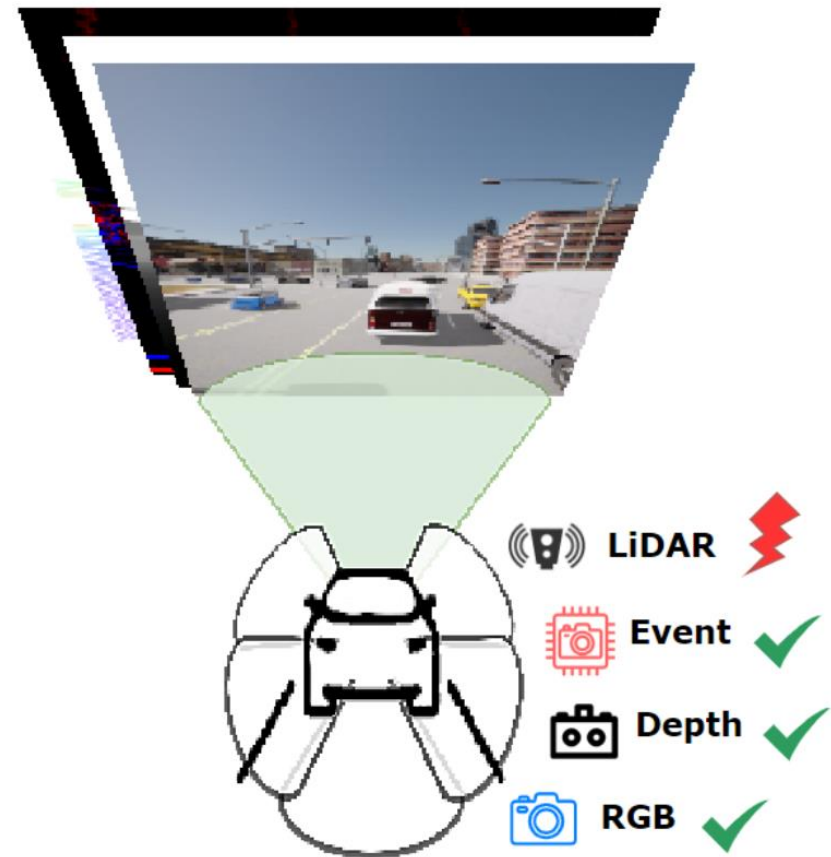
CMNeXt: Arbitrary-Modal Fusion

Extending CMX

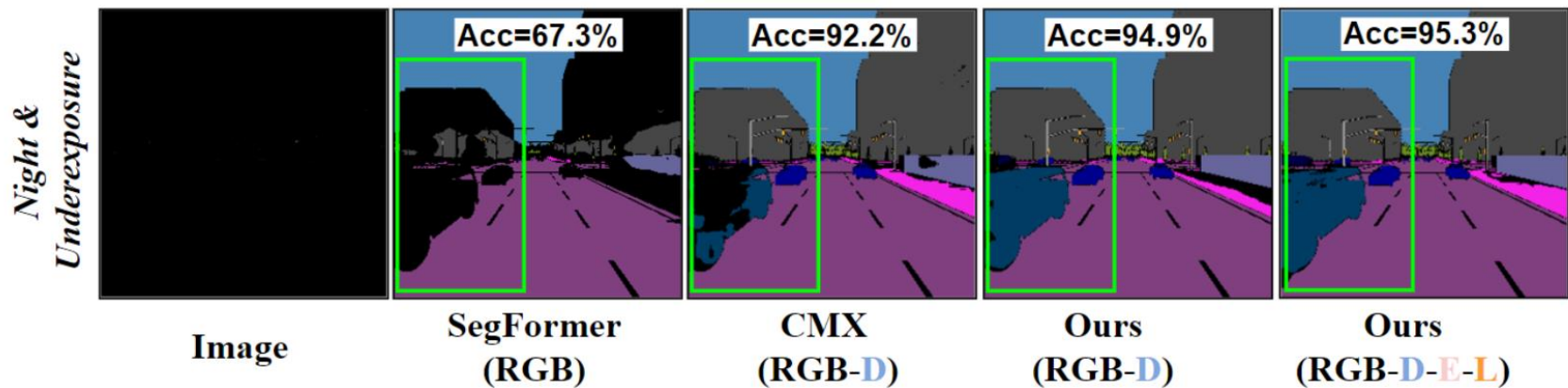
- Multiple additional modalities
- Retains two-stream architecture

Synthetic dataset DeLiVER

- Depth
- LiDAR
- Multiple Views
- Event



CMX and CMNeXt - Results



CMX and CMNeXt - Results

(b) Results on MFNet.

| Method | Modal | mIoU |
|-------------------|-------|-------------|
| SwinT [50] | RGB | 49.0 |
| SegFormer [80] | RGB | 52.0 |
| ACNet [35] | RGB-T | 46.3 |
| FuseSeg [66] | RGB-T | 54.5 |
| ABMDRNet [96] | RGB-T | 54.8 |
| LASNet [41] | RGB-T | 54.9 |
| FEANet [15] | RGB-T | 55.3 |
| MFTNet [101] | RGB-T | 57.3 |
| GMNet [103] | RGB-T | 57.3 |
| DooDLeNet [20] | RGB-T | 57.3 |
| CMX (MiT-B2) [49] | RGB-T | 58.2 |
| CMX (MiT-B4) [49] | RGB-T | 59.7 |
| CMNeXt (MiT-B4) | RGB-T | 59.9 |

(c) Results on NYU Depth V2.

| Method | mIoU |
|--------------------|-------------|
| ACNet [35] | 48.3 |
| SGNet [9] | 51.1 |
| ShapeConv [5] | 51.3 |
| NANet [92] | 52.3 |
| SA-Gate [11] | 52.4 |
| PGDENet [104] | 53.7 |
| TokenFusion [72] | 54.2 |
| TransD-Fusion [78] | 55.5 |
| MultiMAE [2] | 56.0 |
| Omnivore [25] | 56.8 |
| CMX (MiT-B4) [49] | 56.3 |
| CMX (MiT-B5) [49] | 56.9 |
| CMNeXt (MiT-B4) | 56.9 |

Multimodal Knowledge Expansion

Zihui Xue^{1,2}, Sucheng Ren^{1,3}, Zhengqi Gao^{1,4}, and Hang Zhao ^{*5,1}

¹Shanghai Qi Zhi Institute, ²UT Austin

³South China University of Technology

⁴MIT, ⁵Tsinghua University

Multimodal Knowledge Expansion - MKE

Models need to be trained on data!

RGB

- Big field of research
- Many datasets
- Well-trained backbones

Multi-modal

- Some labeled datasets, lots of unlabeled data
- Not many pre-trained backbones

Multimodal Knowledge Expansion - MKE

Models need to be trained on data!

RGB

- Big field of research
- Many datasets
- Well-trained backbones

Multi-modal

- Some labeled datasets, lots of unlabeled data
- Not many pre-trained backbones



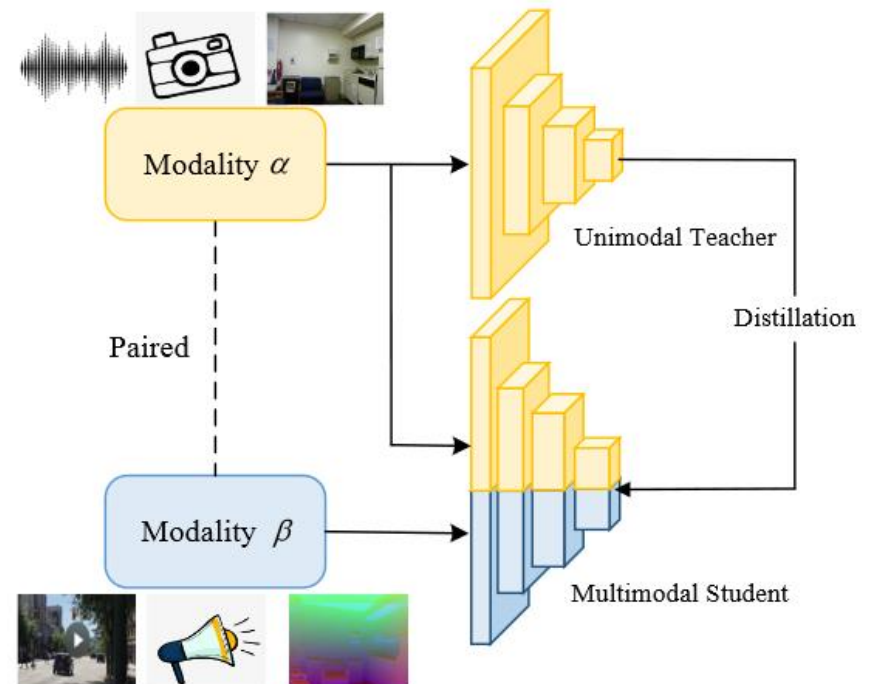
Transfer knowledge to different modes?

MKE - Method

Based on knowledge distillation

Teacher-Student architecture

- Teacher
 - unimodal
 - generates pseudo-labels
- Student
 - multi-modal
 - learns on pseudo-labels



MKE - Method

Confirmation Bias

Student should not strictly confirm to Teacher's pseudo-labels!

Solution → Loss term

(like consistency regularization in SSL)

$$\theta_s^* = \underset{\theta_s}{\operatorname{argmin}} (\mathcal{L}_{pl} + \gamma \mathcal{L}_{reg})$$

$$\mathcal{L}_{pl} = \frac{1}{M} \sum_{i=1}^M l_{cls}(\tilde{\mathbf{y}}_i, \mathbf{f}_s(\mathbf{x}_i^\alpha, \mathbf{x}_i^\beta; \theta_s))$$

$$\mathcal{L}_{reg} = \sum_{i=1}^M l_{reg}[\mathbf{f}_s(\mathbf{x}_i^\alpha, \mathbf{x}_i^\beta; \theta_s), \mathcal{T}(\mathbf{f}_s(\mathbf{x}_i^\alpha, \mathbf{x}_i^\beta; \theta_s))]$$

| | |
|----------------------|--|
| $l_{cls}(\cdot)$ | Cross-entropy loss |
| $l_{reg}(\cdot)$ | Distance metric (L2) |
| $f_s(\cdot)$ | Student model |
| $\mathcal{T}(\cdot)$ | Transformation on student model (i.e. input or model perturbation) |

MKE - Results

| Method | Train data | | | Test mIoU (%) |
|--------------------------|---------------|-------|---------------|------------------|
| | <i>mod</i> | D_l | \tilde{D}_u | |
| UM teacher | <i>rgb</i> | ✓ | | 44.15 |
| Naive student [10] | <i>rgb</i> | | ✓ | 46.13 |
| NOISY student [44] | <i>rgb</i> | ✓ | ✓ | 47.68 |
| Gupta <i>et al.</i> [15] | <i>rgb, d</i> | | ✓ | 45.65 |
| CMKD [49] | <i>rgb, d</i> | | ✓ | 45.25 |
| MM student (no reg) | <i>rgb, d</i> | | ✓ | 46.14 |
| MM student (ours) | <i>rgb, d</i> | | ✓ | 48.88 |

Table 4: Results of semantic segmentation on NYU Depth V2. *rgb* and *d* denote RGB images and depth images.

MKE - Results

| Methods | Train data | | | Accuracy (%) | |
|---------------------|-------------|-------|---------------|--------------|--------------|
| | <i>mod</i> | D_l | \tilde{D}_u | val | test |
| UM teacher | <i>i</i> | ✓ | | 79.67 | 80.33 |
| UM student | <i>i</i> | | ✓ | 79.01 | 77.79 |
| NOISY student [44] | <i>i</i> | ✓ | ✓ | 82.54 | 83.09 |
| MM student (no reg) | <i>i, a</i> | | ✓ | 88.73 | 89.28 |
| MM student (ours) | <i>i, a</i> | | ✓ | 90.61 | 91.38 |
| MM student (sup) | <i>i, a</i> | | ★ | 97.46 | 97.35 |

Table 3: Results of emotion recognition on RAVDESS. *mod*, *i* and *a* denote modality, images and audios, respectively. Data used for training each method is listed. ★ means that the MM student (sup) is trained on true labels instead of pseudo labels in \tilde{D}_u .

Personal Comments

Multi-modal curb detection

- Unimodal Segmentation
- Simple unsupervised fusion in pipeline
- No interaction / end-to-end learning!

CMX and CMNeXt

- Unified fusion framework
- (*close to*) SOTA, even comparing to specialized models
- only image-like formats - no sparse data (LiDAR!)

Personal Comments

Multi-modal knowledge expansion

- Exciting (and surprising) results
- however very theoretical
- still a young field of research

Future Work

Leverage RGB knowledge in SOTA

- Starting point: Multi-modal Knowledge Expansion
- Train SOTA multi-modal model
- Student-teacher architecture
- Improvements?

Future Work

Support of different representations for modalities

- for fully unified multi-modal framework
- Image-like data, point cloud, non-structural (Audio, Language, ...)
- no conversion / loss of structure necessary

even further:

Shared representation between modalities

- Recent advance → ImageBind

Thank you for listening!
Any Questions?