# Multimodal Perception

**PACCAR**

Fulmi Chang – ajchangwa@gmail.com, Sydney Dukelow - sydneydukelow@gmail.com
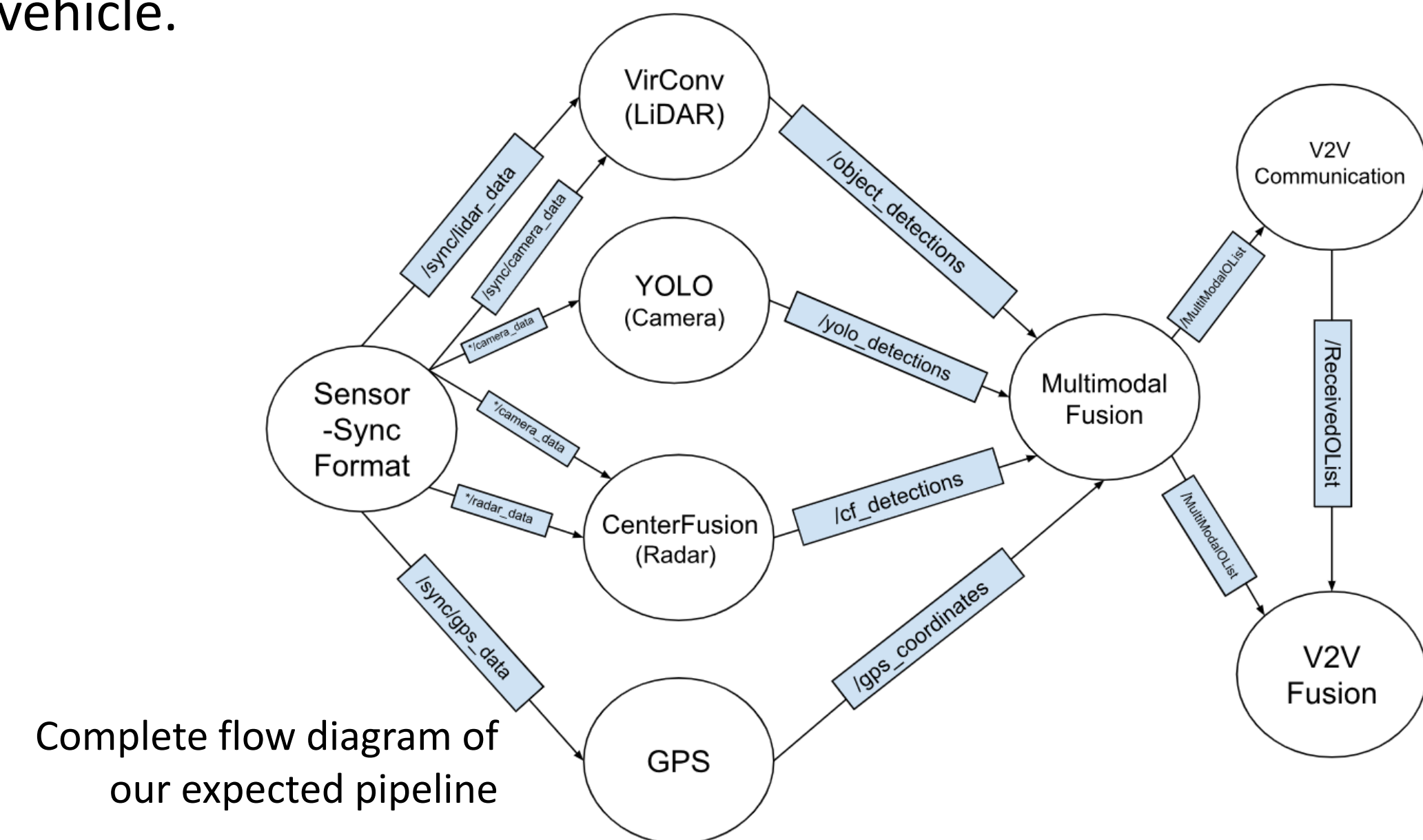
WESTERN WASHINGTON UNIVERSITY

MAKE WAVES.

## Abstract

Our project aims to enhance vehicular environment perception by implementing multimodal and vehicle-to-vehicle (V2V) fusion methods. We use state-of-the-art machine learning algorithms to apply object detection tailored to our datasets. Multimodal fusion integrates data from multiple sensors:
- LiDAR
- Radar
- Camera
- GPS

V2V fusion integrates local multimodal data with another vehicle.

Complete flow diagram of our expected pipeline

## Design Methods and Tools

We utilized a Linux system with the ROS Noetic framework to contain our data processing system. Within the ROS framework, we created nodes for each sensor object detection and fusion. We trained each model to detect cars, cyclists, and pedestrians. We would then send raw data through our data processing pipeline to each model and have it output its respective object lists.

The algorithms used for each node in our system depends on their respective data type. For LiDAR, we used "Virtual Sparse Convolution for Multimodal 3D Object Detection" (VirConv), camera used "You Only Look Once version 8" (YOLOv8), and radar used "Center-based Radar and Camera Fusion for 3D Object Detection" (CenterFusion).

We utilized the KITTI dataset for all modalities except radar, which used the nuScenes dataset.

### Hardware Specifications

|  | Lenovo Legion Pro 7 | Nvidia Jetson AGX Orin |
| --- | --- | --- |
| Manufacturer: | Lenovo | Nvidia |
| Model: | 16IRX8H/16IRX9H | 945-13730-0050-000 |
| System Frequency: | CPU: 5.4GHz GPU: 2040MHz | 2.2GHz |
| Memory: | 64GB | 64GB |
| Storage: | 4TB | 64GB + 1TB |
| Environment: | Windows Subsystem for Linux 2 ROS Noetic | Ubuntu 20.04 LTS ROS Noetic |



Lenovo Legion Pro 7



Nvidia Jetson AGX Orin

### ROS Node Inputs

| Node | Inputs |
| --- | --- |
| VirConv | • LiDAR<br>• Camera |
| YOLO | • Camera |
| CenterFusion | • Radar<br>• Camera |
| GPS | • GPS |
| Multimodal Fusion | • LiDAR<br>• Camera<br>• Radar<br>• GPS |



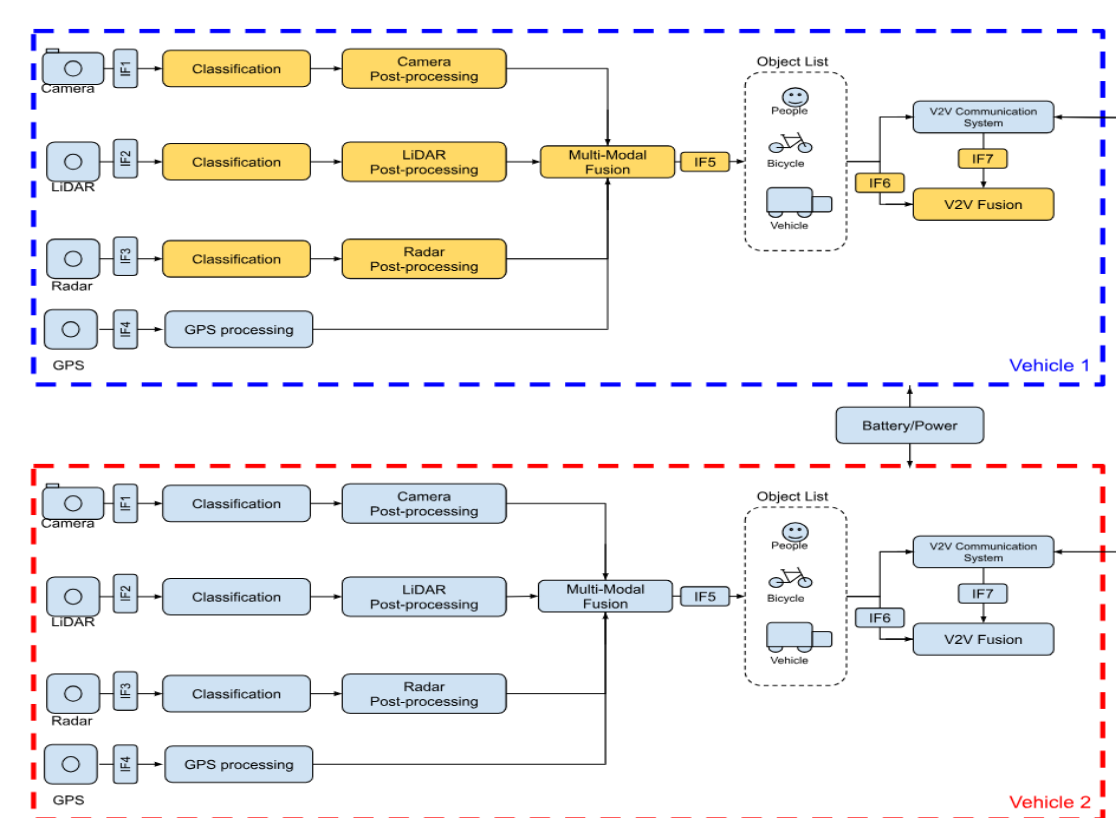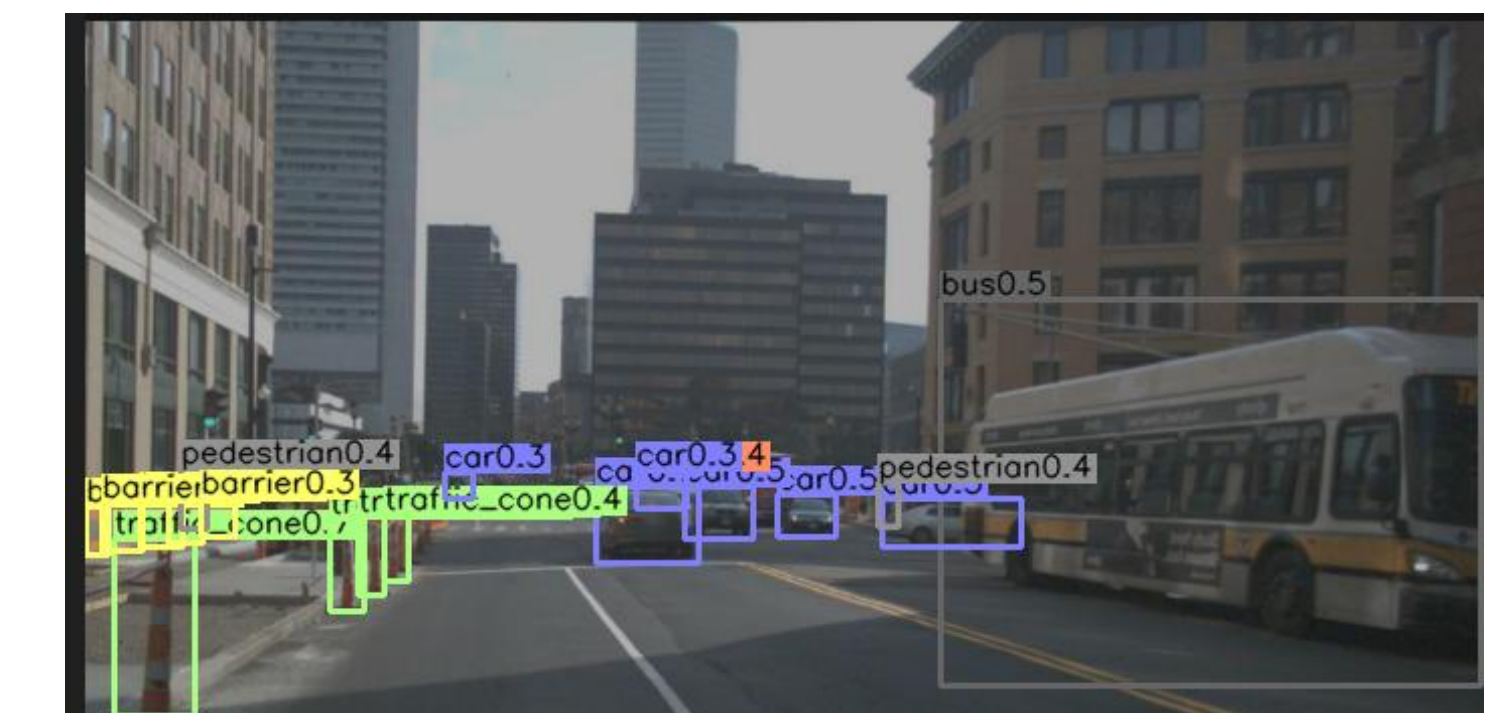CenterFusion Validation Samples



YOLOv8 Validation Samples

## Background

Connected Autonomous Vehicles (CAVs) rely on accurate environmental perception to ensure safe operation.

With a focus on:
- Accuracy
- Energy efficiency
- Timeliness

We developed a data processing pipeline that performs object detection and fusion with onboard sensors.



Hardware level diagram of our system

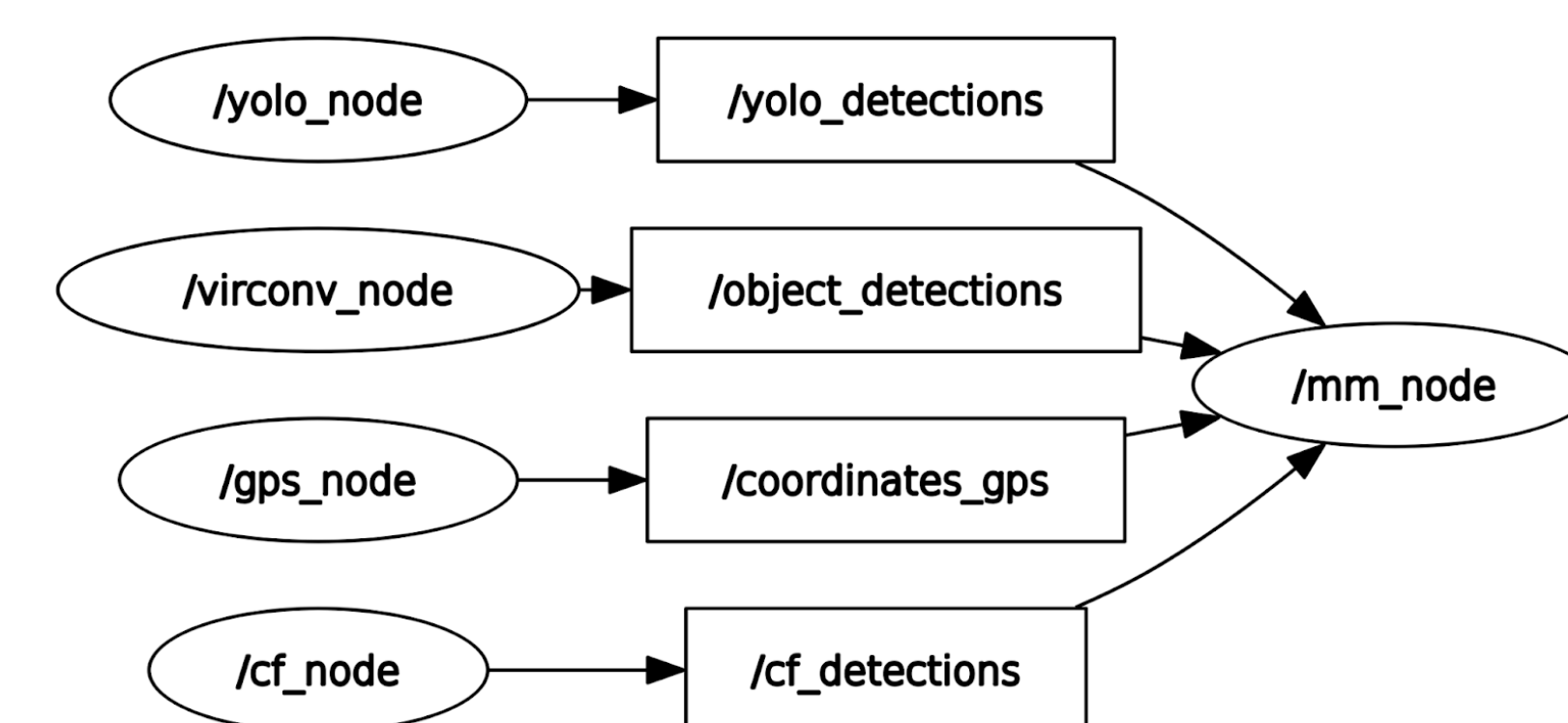## Results

- Successfully created the pipeline that does individual modality object detection
- These tasks are sending their respective object lists to the multimodal fusion node
- Integrated YOLO with Real Time Cooperative Perception Using C-V2X team
- Merged each modality's object list to create a singular comprehensive object list

### Task Metrics

| Node Task | Execution Time per file (milliseconds) | Energy Efficiency (G-FLOPs) |
| --- | --- | --- |
| VirConv | 950ms | 90.85 GFLOPs |
| YOLO | 2ms | 151.01 GFLOPs |
| CenterFusion | 50ms | 306.76 GFLOPs |
| GPS | 50ms | 0 |

### Mean Average Precision Scores per Algorithm

| Algorithm | mAP Score |
| --- | --- |
| VirConv | 0.9623 |
| YOLO | 0.8090 |
| CenterFusion | 0.3298 |



Unprocessed multimodal fusion object list outputs



Node connections as seen through ROS



YOLO Confusion Matrix Output



YOLO Precision-Confidence Validation Plot

## Future Direction

- Train VirConv-S model for better timeliness
- Train a model of CenterFusion on the WWU cluster for better overall performance
- Integrate V2V fusion
- Integrate multimodal fusion in real-time
- Improve multimodal object list with a Kalman Filter and Multi-Object Tracking (MOT)
- Produce the energy consumption and flops of the VirConv and CenterFusion algorithms



Diagram of V2V communication

## Acknowledgements

## References

[1] H. Wu, C. Wen, S. Shi, X. Li and C. Wang, "Virtual Sparse Convolution for Multimodal 3D Object Detection," 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Vancouver, BC, Canada, 2023, pp. 21653-21662, doi: 10.1109/CVPR52729.2023.02074.
[2] J. Redmon, S. Divvala, R. Girshick and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, NV, USA, 2016, pp. 779-788, doi: 10.1109/CVPR.2016.91.
[3] R. Nabati and H. Qi, "CenterFusion: Center-based Radar and Camera Fusion for 3D Object Detection," 2021 IEEE Winter Conference on Applications of Computer Vision (WACV), Waikoloa, HI, USA, 2021, pp. 1526-1535, doi: 10.1109/WACV48630.2021.00157.
[4] Geiger A, Lenz P, Stiller C, Urtasun R. Vision meets robotics: The KITTI dataset. The International Journal of Robotics Research. 2013;32(11):1231-1237. doi:10.1177/0278364913491297
[5] H. Caesar et al., "nuScenes: A Multimodal Dataset for Autonomous Driving," 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Seattle, WA, USA, 2020, pp. 11618-11628, doi: 10.1109/CVPR42600.2020.01164.

## WWU

2025 Capstone Project
Electrical & Computer Engineering Program