

# Imitation with Transformer-Based Sensor Fusion for Autonomous Driving

Kashyap Chitta



University of Tübingen  
MPI for Intelligent Systems  
**Autonomous Vision Group**



# Team



Kashyap Chitta



Aditya Prakash



Bernhard Jaeger



Zehao Yu



Katrin Renz



Andreas Geiger

# Covered Papers

- ▶ **Multi-Modal Fusion Transformer for End-to-End Autonomous Driving**

A. Prakash\*, K. Chitta\* and A. Geiger. CVPR, 2021.

- ▶ **TransFuser: Imitation with Transformer-Based Sensor Fusion**

K. Chitta, A. Prakash, B. Jaeger, Z. Yu, K. Renz and A. Geiger. In Review.

# Evaluating Self-Driving

# Common Task Framework

## Computer Vision



Semantic  
Segmentation

130 benchmarks

2537 papers with code



Image  
Classification

311 benchmarks

2189 papers with code



Object  
Detection

216 benchmarks

1910 papers with code



Image  
Generation

176 benchmarks

839 papers with code



Denoising

103 benchmarks

802 papers with code

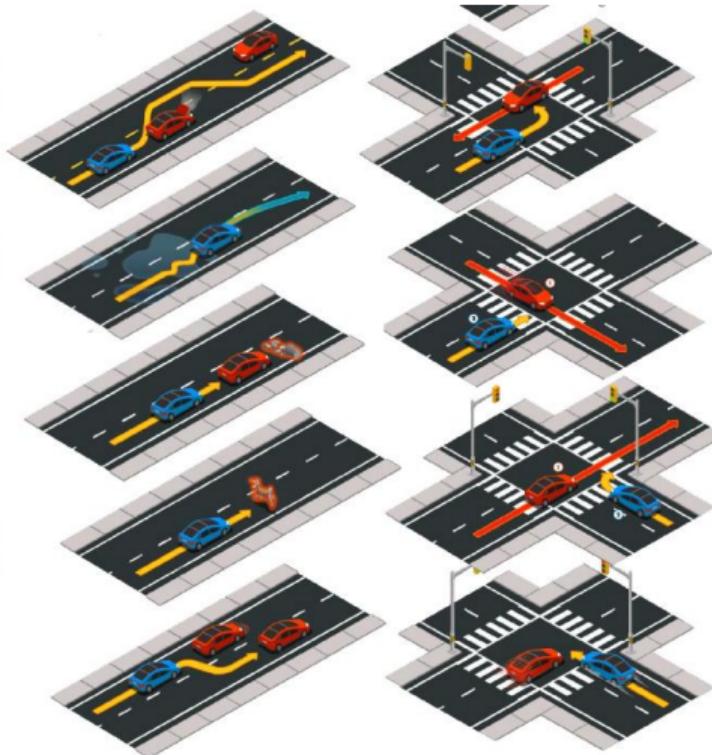
[▶ See all 1235 tasks](#)

- ▶ A **common task framework** accelerates research progress
- ▶ Computer vision: **static benchmarks**
- ▶ How can the community compare **dynamic self-driving** agents?

# CARLA Leaderboard



- 10 routes x 2 weathers x 5 repetitions
- 173 Km of driving experiences



# CARLA Leaderboard Evaluation

$$\frac{1}{n} \sum_{i=1}^n c_i p_i$$

# of routes →

Completion of route  $i$  ↑

Infraction penalty for route  $i$  ↘

$p_i = \prod_{j \in \mathcal{J}} (p^j)^{v_i^j}$

Number of infractions of type  $j$  in route  $i$  ↘

Penalty for infraction of type  $j$  ↗

# Imitation Learning for CARLA

# Imitation Learning

**Motivation:** Hand-designing a sensor-based driving policy is difficult

# Imitation Learning

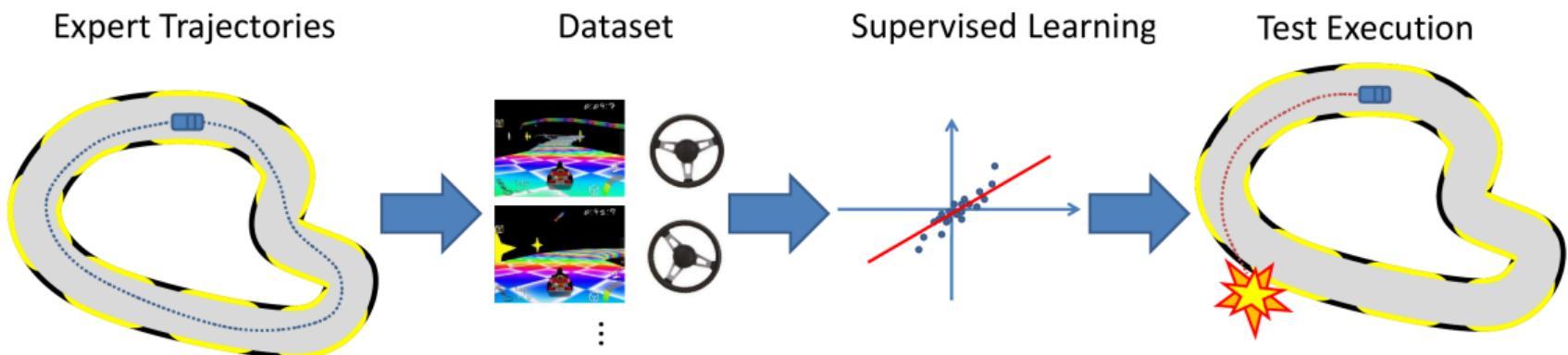
**Motivation:** Hand-designing a sensor-based driving policy is difficult

- ▶ **Step 1:** Hand-design expert which uses privileged information

# Imitation Learning

**Motivation:** Hand-designing a sensor-based driving policy is difficult

- ▶ **Step 1:** Hand-design expert which uses privileged information
- ▶ **Step 2:** Train sensor-based policy to mimic demonstrator



# Sensor Fusion

# Sensors

## RGB Camera



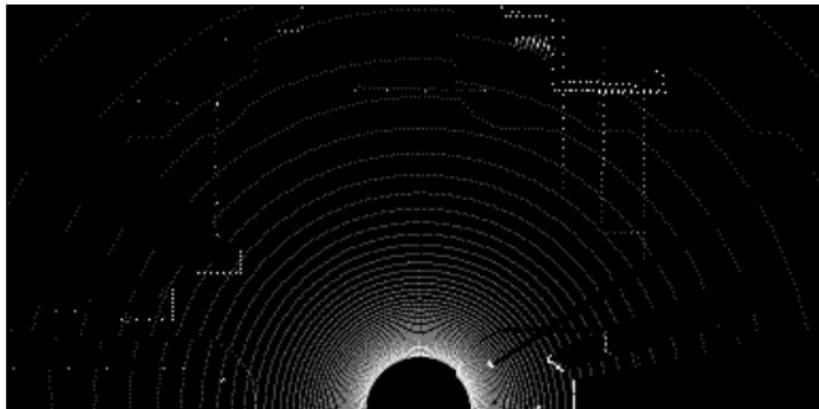
- + Dense RGB input
- Lacks reliable 3D information
- Variation in weather

# Sensors

## RGB Camera



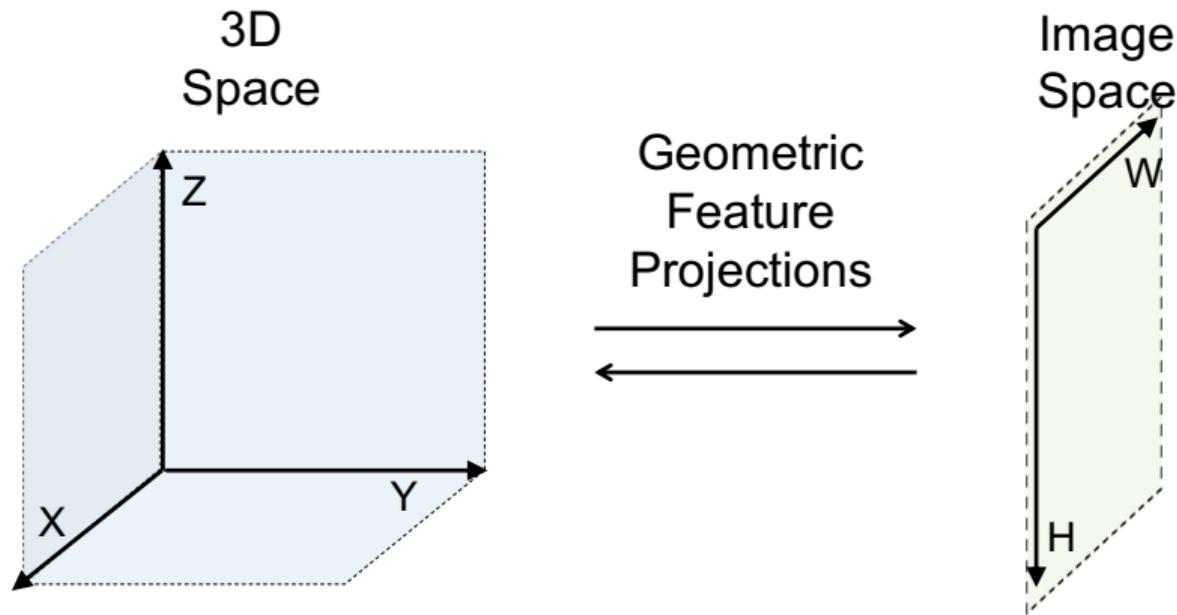
## LiDAR Point Cloud



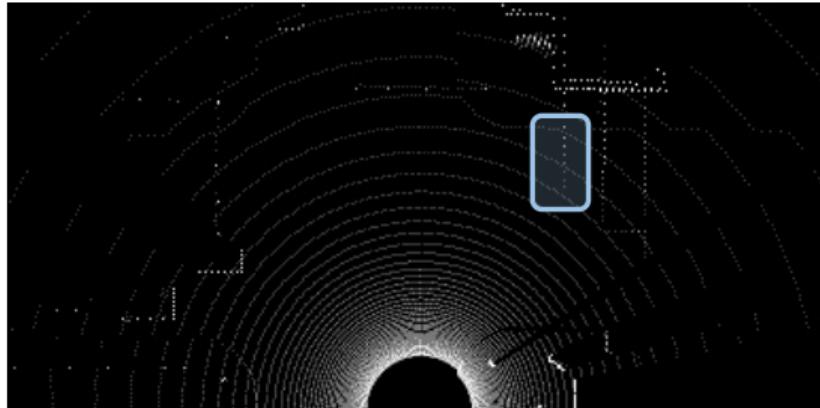
- + Dense RGB input
- Lacks reliable 3D information
- Variation in weather

- + 3D information
- Sparse input
- No traffic light state

# Geometric Fusion

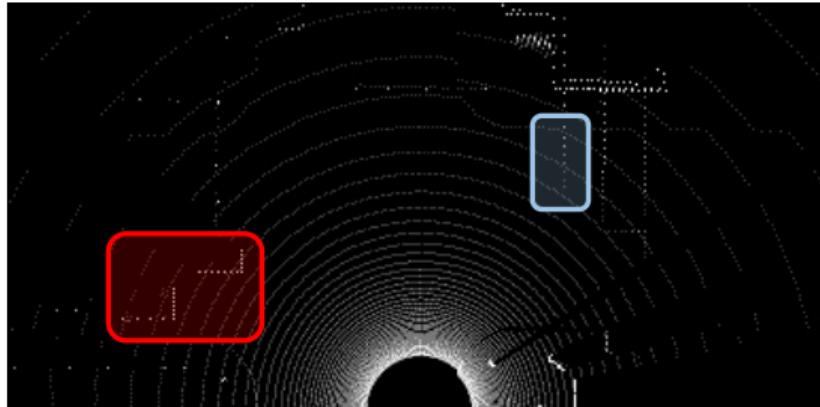


# Geometric Fusion Lacks Global Context



- ▶ From the yellow region, geometric fusion aggregates features to the blue region

# Geometric Fusion Lacks Global Context



- ▶ From the yellow region, geometric fusion aggregates features to the blue region
- ▶ However, for safe navigation, it is useful to aggregate features for the red region since it contains vehicles which are affected by the traffic light

# TransFuser

# Key Idea

Use **attention-based** feature fusion  
to capture the **global context** of the  
scene **across modalities**.



# TransFuser

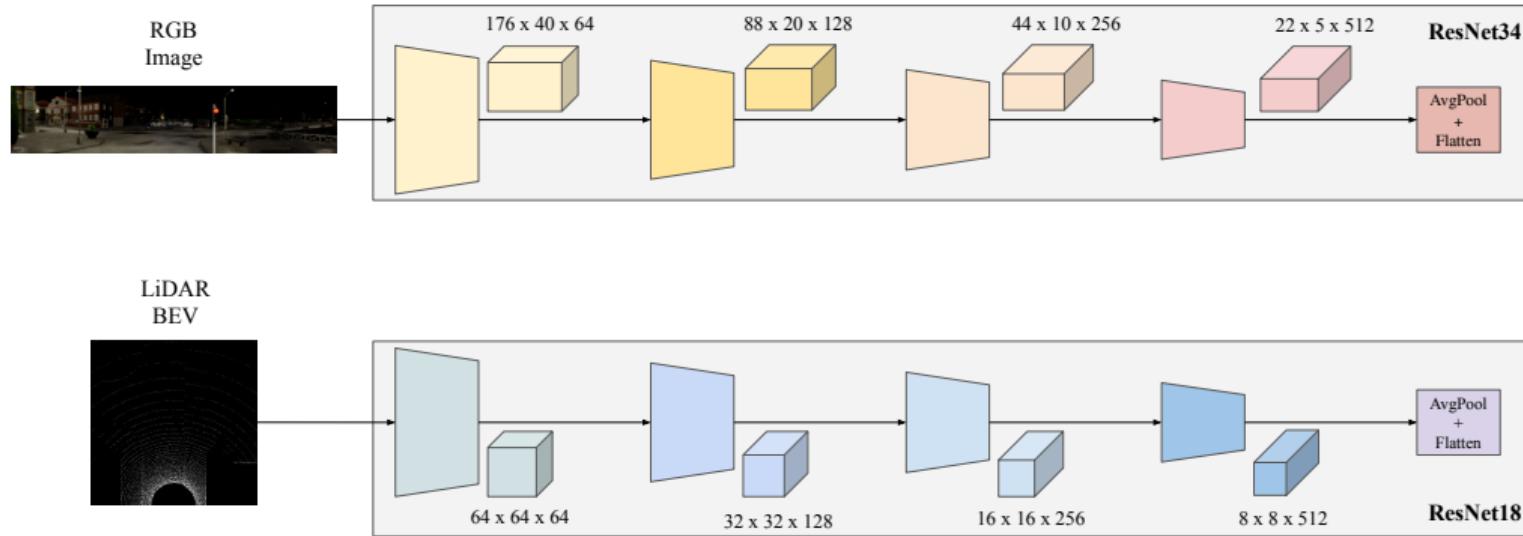
RGB  
Image



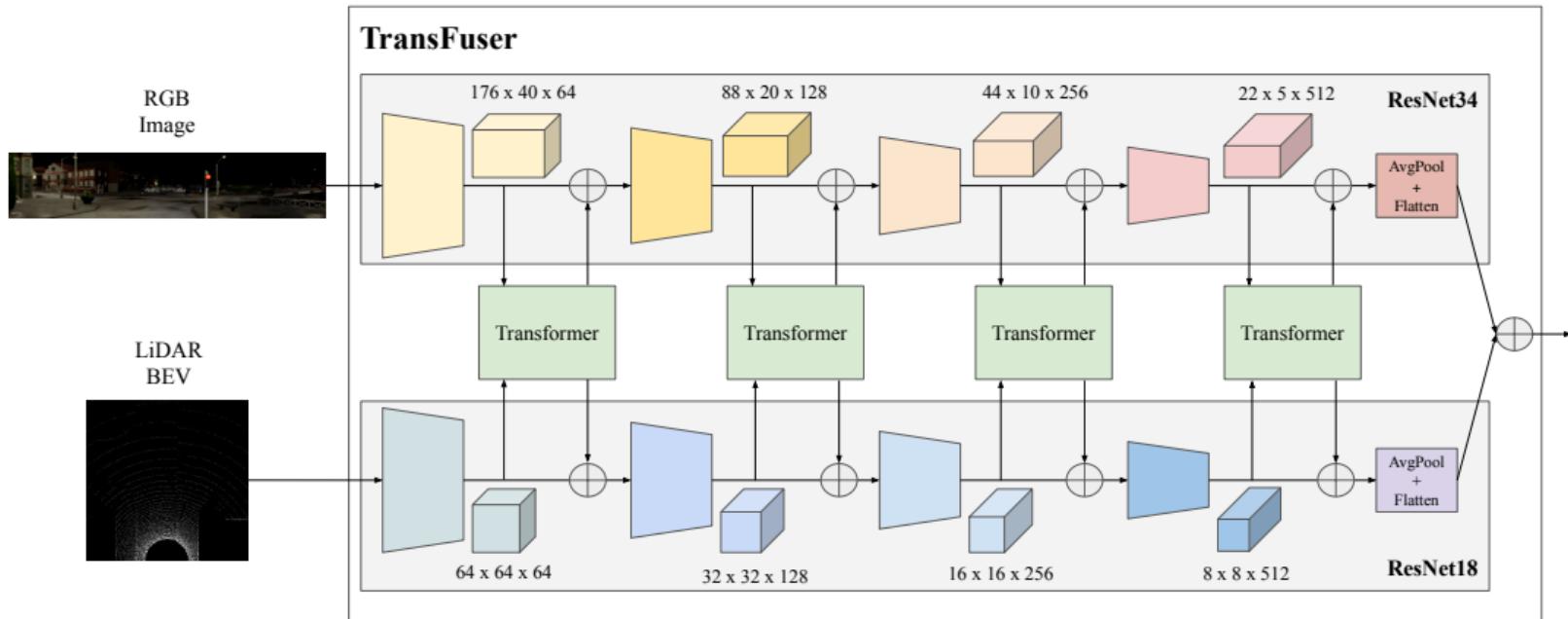
LiDAR  
BEV



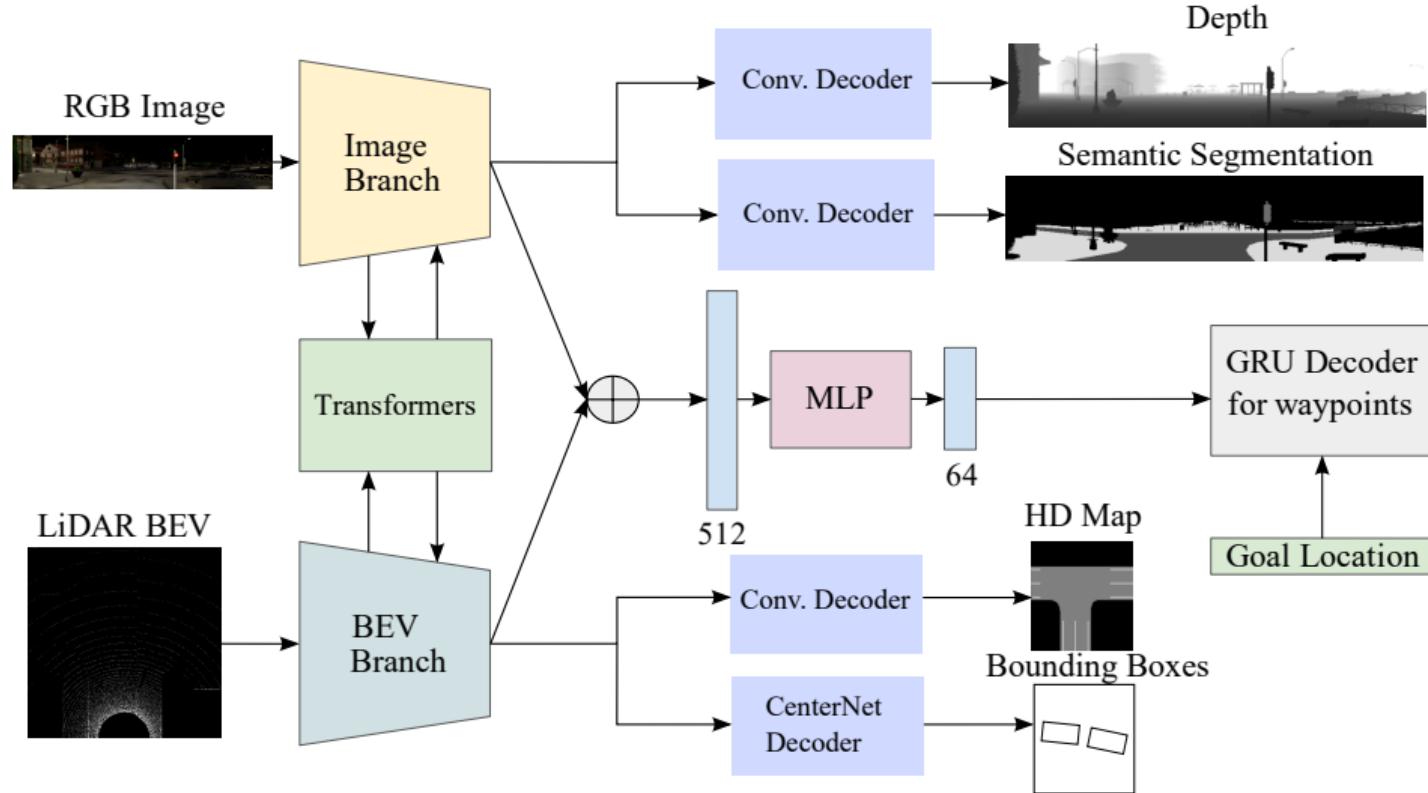
# TransFuser



# TransFuser



# Full Architecture



# Loss Functions

- ▶  $L_1$  loss on waypoints:  $\mathcal{L} = \sum_{t=1}^4 \|\mathbf{w}_t - \mathbf{w}_t^{gt}\|_1$
- ▶ Cross-entropy loss on semantics
- ▶  $L_1$  loss on depth
- ▶ Cross-entropy loss on HD map
- ▶ Focal loss on CenterNet heatmaps
- ▶  $L_1$  loss on CenterNet offsets

# Experiments

## Dataset

- ▶ 8 Towns and randomized weather conditions in CARLA
- ▶ Expert policy based on MPC
- ▶ ~3.5k short routes with hand-crafted scenarios

# Experiments

## Dataset

- ▶ 8 Towns and randomized weather conditions in CARLA
- ▶ Expert policy based on MPC
- ▶ ~3.5k short routes with hand-crafted scenarios

## Sensors

- ▶ RGB cameras:  $704 \times 160$  resolution,  $132^\circ$  FOV
- ▶ LiDAR: 32m range, 64 channels, 10 Hz rotation frequency

# Experiments

## Dataset

- ▶ 8 Towns and randomized weather conditions in CARLA
- ▶ Expert policy based on MPC
- ▶ ~3.5k short routes with hand-crafted scenarios

## Sensors

- ▶ RGB cameras:  $704 \times 160$  resolution,  $132^\circ$  FOV
- ▶ LiDAR: 32m range, 64 channels, 10 Hz rotation frequency

## Evaluation

- ▶ Long routes ( $\sim 2\text{km}$ ) with dense traffic
- ▶ Ensemble of 3 training runs to reduce variance

## Results: Internal Benchmark

Method	Driving Score ↑	Route Completion ↑	Infraction Score ↑
Late Fusion	$22 \pm 4$	$83 \pm 3$	$0.27 \pm 0.03$
Geometric Fusion	$27 \pm 1$	$91 \pm 1$	$0.30 \pm 0.02$
TransFuser (Ours)	<b><math>47 \pm 6</math></b>	<b><math>93 \pm 1</math></b>	<b><math>0.50 \pm 0.00</math></b>
<i>Privileged Expert</i>	$77 \pm 2$	$89 \pm 1$	$0.86 \pm 0.03$

- ▶ Geometric Fusion, TransFuser and Expert have similar route completion
- ▶ Clear trend in infraction score (Expert > TransFuser > Baselines)

# CARLA Leaderboard

Method	Driving Score ↑	Route Completion ↑	Infraction Score ↑
LAV	61.85	94.46	0.64
TransFuser (Ours)	61.18	<b>86.69</b>	<b>0.71</b>
GRIAD	36.79	61.85	0.60
WOR	31.37	57.65	0.56

- ▶ Simple end-to-end IL (competitors have complex multi-stage training pipelines)
- ▶ Rank 2 at submission (April), with **best infraction score** among top methods
- ▶ Still **gets blocked** more often than LAV
- ▶ DS > 60, rapid overall progress on leaderboard since 2020 (DS < 20)

# Summary

## Conclusions

- Global contextual reasoning is crucial in complex urban scenarios

# Summary

## Conclusions

- ▶ Global contextual reasoning is crucial in complex urban scenarios
- ▶ Attention is effective in aggregating information from multiple modalities

# Summary

## Conclusions

- ▶ Global contextual reasoning is crucial in complex urban scenarios
- ▶ Attention is effective in aggregating information from multiple modalities
- ▶ Driving Score of simple Imitation Learning baseline is competitive (rank 2)

# Summary

## Conclusions

- ▶ Global contextual reasoning is crucial in complex urban scenarios
- ▶ Attention is effective in aggregating information from multiple modalities
- ▶ Driving Score of simple Imitation Learning baseline is competitive (rank 2)

## Code

- ▶ [www.github.com/autonomousvision/transfuser](https://www.github.com/autonomousvision/transfuser)

## Other Work

- ▶ Ohn-Bar et al.: Learning Situational Driving. CVPR, 2020.  
**“Driving in diverse environments is eased by mixture policies.”**

## Other Work

- ▶ Ohn-Bar et al.: Learning Situational Driving. CVPR, 2020.  
**“Driving in diverse environments is eased by mixture policies.”**
- ▶ Prakash et al.: Exploring Data Aggregation in Policy Learning. CVPR, 2020.  
**“Vanilla DAGGER doesn’t work well ⇒ we must sample critical states.”**

## Other Work

- ▶ Ohn-Bar et al.: Learning Situational Driving. CVPR, 2020.  
**“Driving in diverse environments is eased by mixture policies.”**
- ▶ Prakash et al.: Exploring Data Aggregation in Policy Learning. CVPR, 2020.  
**“Vanilla DAGGER doesn’t work well ⇒ we must sample critical states.”**
- ▶ Behl et al.: Label Efficient Visual Abstractions. IROS, 2020.  
**“Visual abstractions help, but annotating less can be more.”**

## Other Work

- ▶ Ohn-Bar et al.: Learning Situational Driving. CVPR, 2020.  
**“Driving in diverse environments is eased by mixture policies.”**
- ▶ Prakash et al.: Exploring Data Aggregation in Policy Learning. CVPR, 2020.  
**“Vanilla DAGGER doesn’t work well ⇒ we must sample critical states.”**
- ▶ Behl et al.: Label Efficient Visual Abstractions. IROS, 2020.  
**“Visual abstractions help, but annotating less can be more.”**
- ▶ Chitta et al.: NEAT: Neural Attention Fields. ICCV, 2021.  
**“BEV predictions from 2D images via neural fields can improve safety.”**

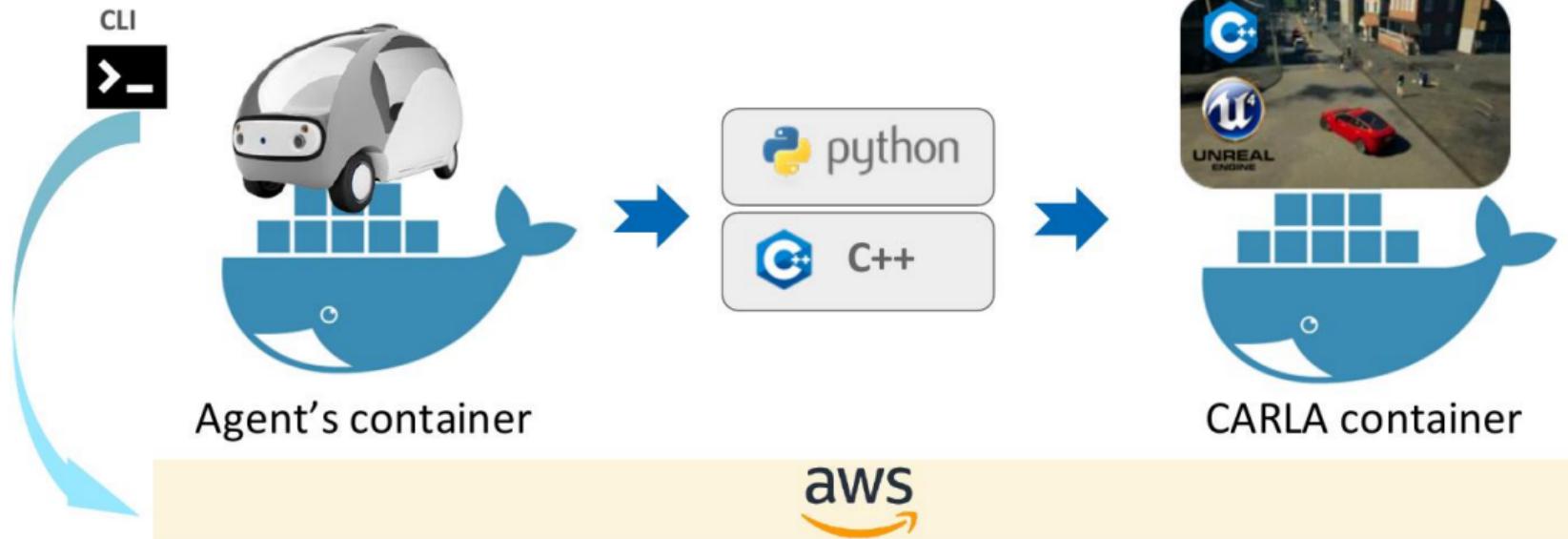
# Extra Slides

# CARLA Leaderboard

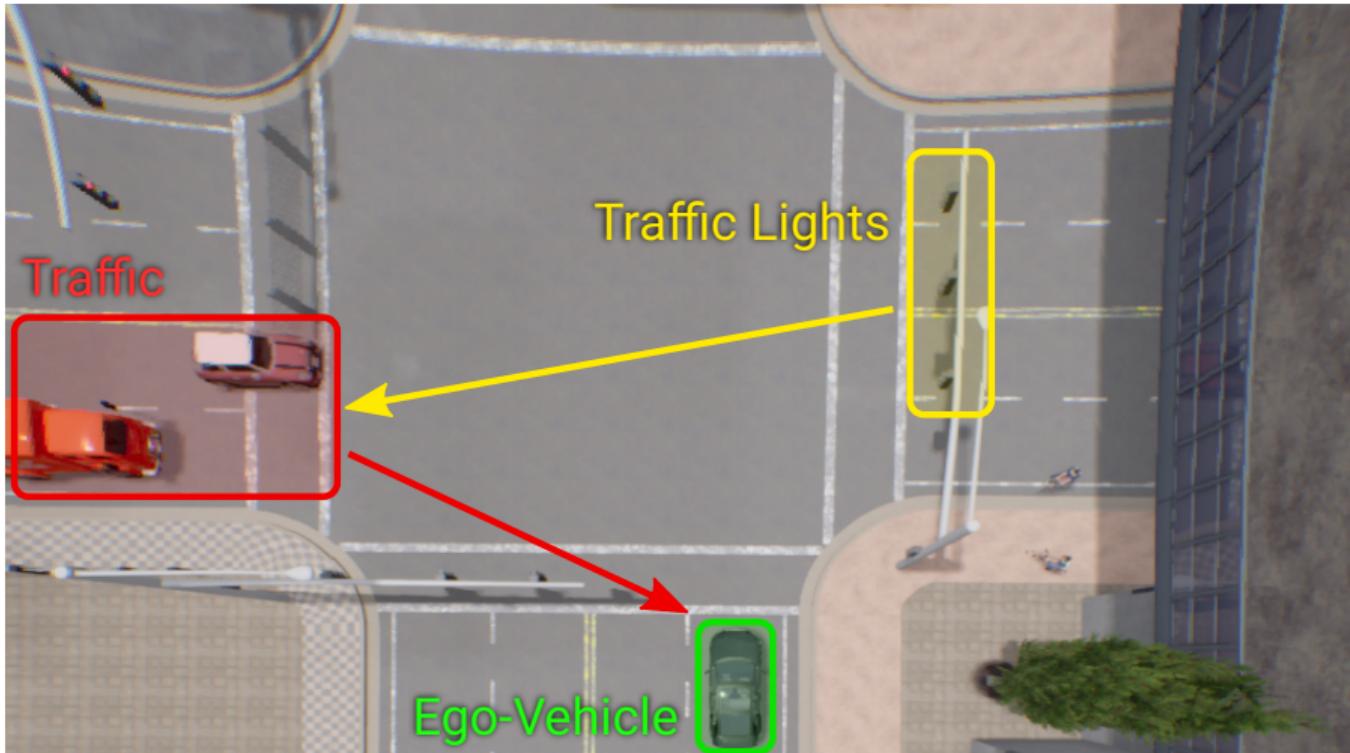
- Open **test bed** to evaluate AD agents for the driving task
- Common maps, situations, and metrics
- Built upon the CARLA simulator
- Aim to **accelerate progress** in the research community



# CARLA Leaderboard Submission



# Motivation



# Research Questions

- ▶ How to integrate representations from multiple modalities?

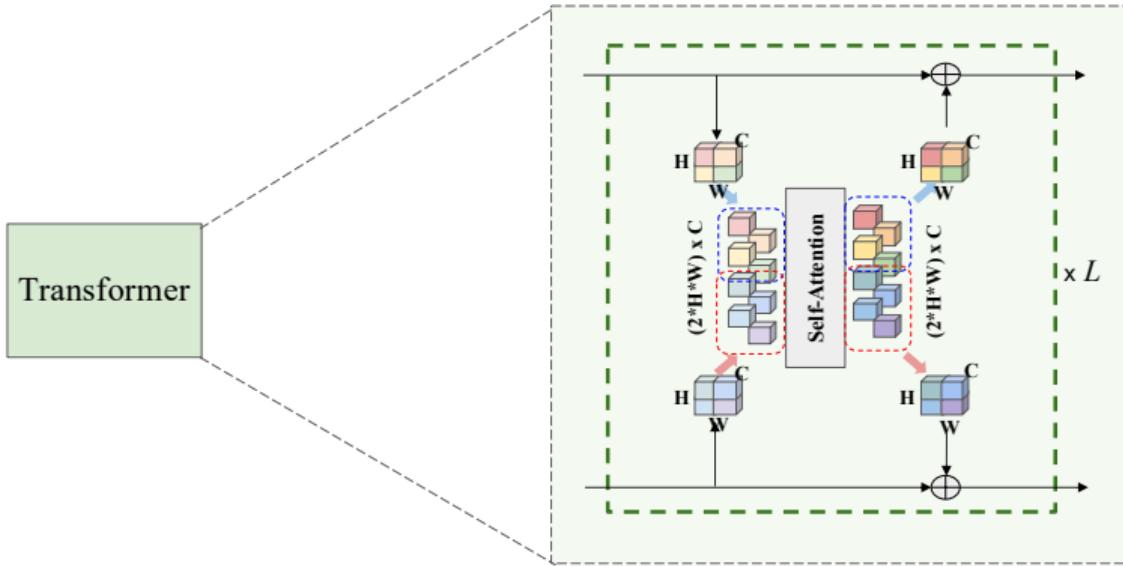
# Research Questions

- ▶ How to integrate representations from multiple modalities?
- ▶ To what extent should the different modalities be processed independently?

# Research Questions

- ▶ How to integrate representations from multiple modalities?
- ▶ To what extent should the different modalities be processed independently?
- ▶ What kind of fusion mechanism to use for maximum performance?

# Attention-based Feature Fusion



- ▶ Consider feature maps as **sets of tokens** (cells of grid = tokens)
- ▶ Pass all tokens to **self-attention** module and reshape back into grid form

# Overall Pipeline

## ► **Step 1 - Privileged Agent (Data Collection)**

- Demonstrator
- Routes
- Sensors

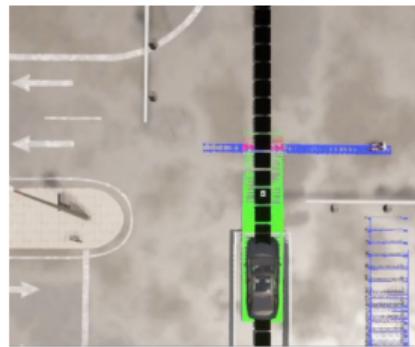
# Overall Pipeline

- ▶ **Step 1 - Privileged Agent (Data Collection)**
  - ▶ Demonstrator
  - ▶ Routes
  - ▶ Sensors
- ▶ **Step 2 - Sensorimotor Agent (Training)**
  - ▶ Architecture
  - ▶ Loss function
  - ▶ Controller

# Demonstrator: Components

## Lateral Control

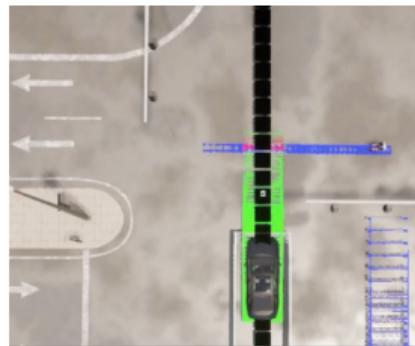
- ▶ Input: HD Map
- ▶ A\* Planner
- ▶ PID controller



# Demonstrator: Components

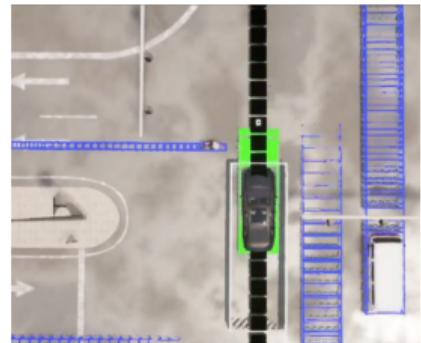
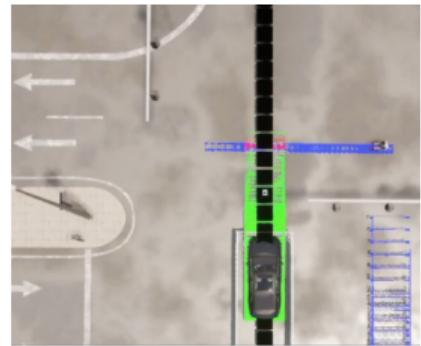
## Longitudinal Control

- ▶ Input: traffic light states
- ▶ Input: nearby actor states
  - ▶ Position
  - ▶ Orientation
  - ▶ Velocity
- ▶ Kinematic bicycle model
- ▶ PID controller



# Demonstrator

- ▶ Simplified version of Model Predictive Control (MPC)
- ▶ 2 candidate trajectories using HD map + PID controllers
  - ▶ Greedy: target speed = 4 m/s
  - ▶ Conservative: target speed = 0 m/s
- ▶ Roll out greedy trajectory with bicycle model
- ▶ Choose conservative trajectory if infraction is detected



## Routes

- ▶ ~ 3000 Junctions (~100m long)
- ▶ ~ 500 Curves (~400m long)
- ▶ 8 CARLA towns (1, 2, 3, 4, 5, 6, 7, 10)
- ▶ 7 CARLA scenarios (1, 3, 4, 7, 8, 9, 10)

## Routes

- ▶ ~ 3000 Junctions (~100m long)
- ▶ ~ 500 Curves (~400m long)
- ▶ 8 CARLA towns (1, 2, 3, 4, 5, 6, 7, 10)
- ▶ 7 CARLA scenarios (1, 3, 4, 7, 8, 9, 10)
- ▶ Time of day: custom distribution around 6 preset values
- ▶ Weathers: 7 CARLA presets
- ▶ Dataset size: 300k frames

# Sensors

## RGB cameras

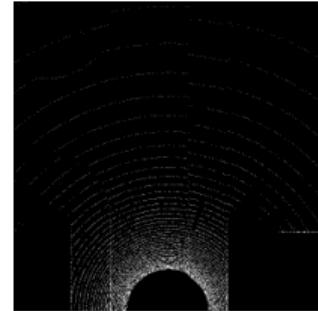
- ▶ 3 cameras: front, 60° left, 60° right
- ▶ Field of view: 60° each
- ▶ Resolution:  $320 \times 160$  pixels each
- ▶ Composited into  $704 \times 160$  input



# Sensors

64 beam LiDAR

- ▶ 10 Hz frequency: use alternate frames
- ▶ Field of view: 180°
- ▶ Rasterized into BEV (256×256, 32m range)
- ▶ 2 channels: ground plane, objects



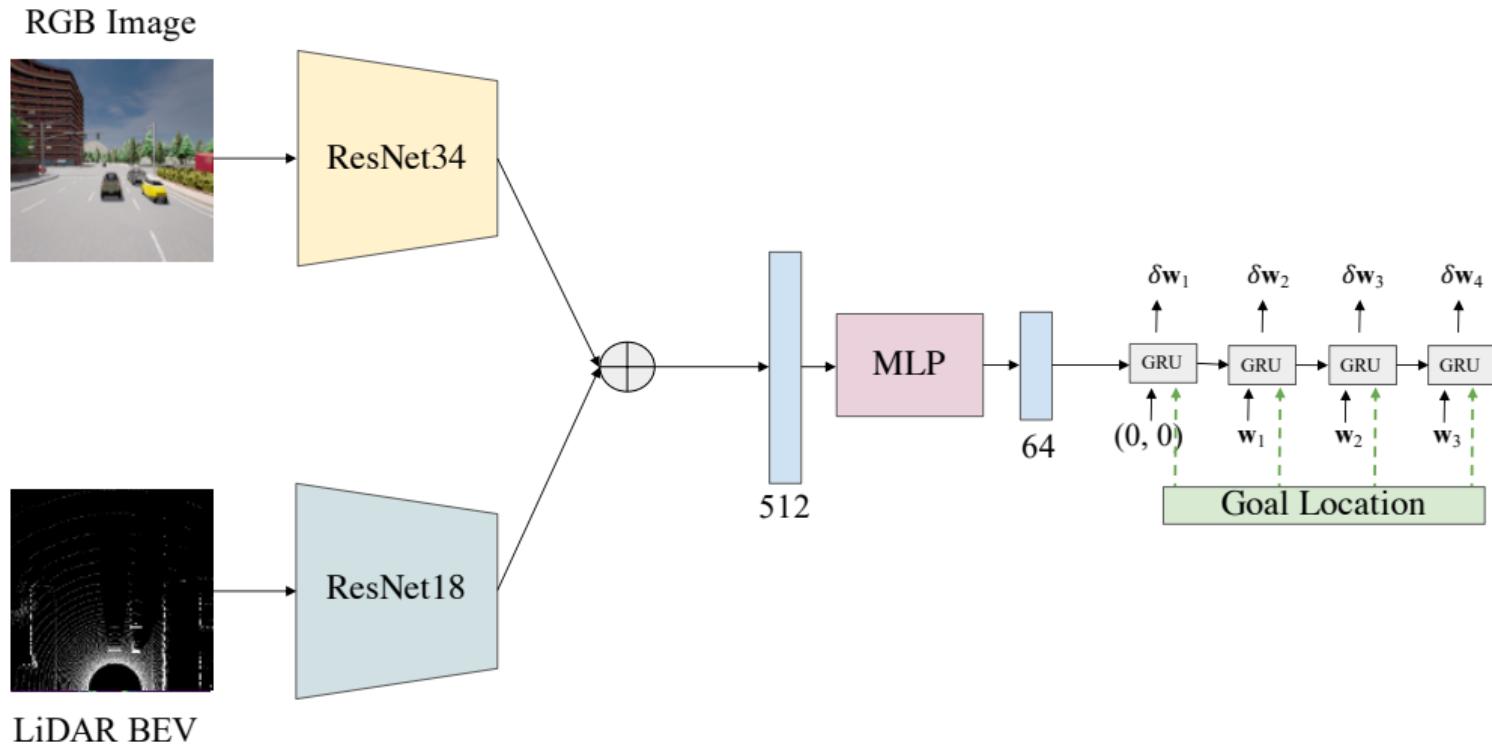
# Sensors

Additional sensors used for auxiliary supervision

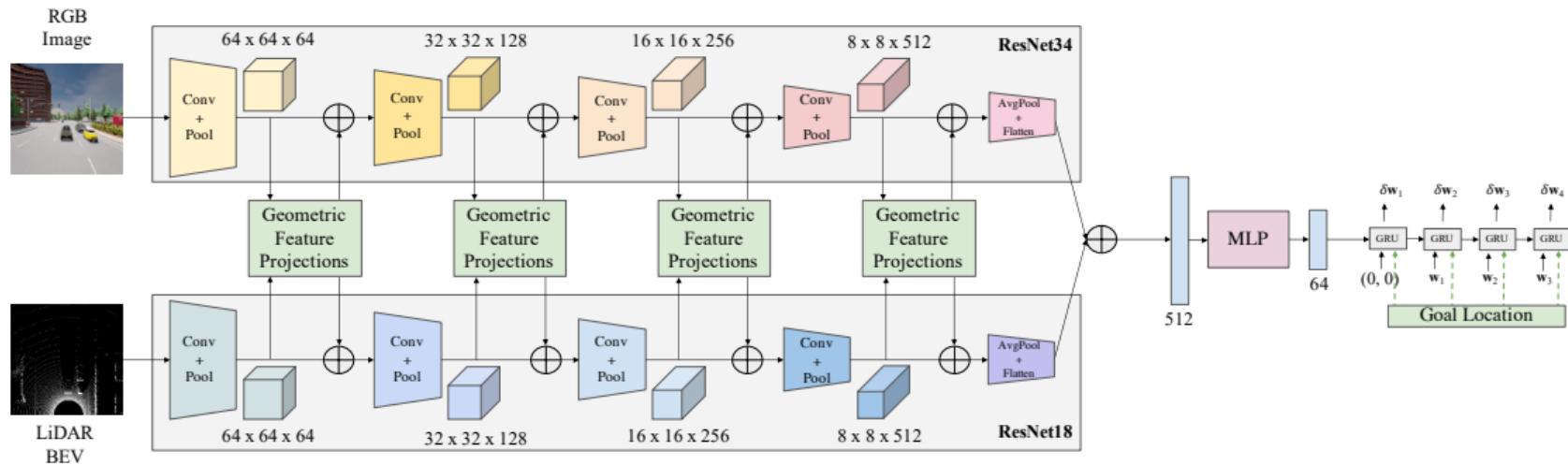
- ▶ Semantic Segmentation
- ▶ Depth
- ▶ HD Map: same coordinate frame as LiDAR



# Baselines - Late Fusion



# Baselines - Geometric Fusion



# Controller

- ▶ Heading and target speed from waypoints
- ▶ PID controllers
- ▶ Inertia problem: creep forward if still for ~1 minute
  - ▶ Safety check: no creeping when LiDAR indicates close proximity

## Runtime

<b>Method</b>	<b>Single Model</b>	<b>Ensemble (3)</b>
Late Fusion (LF)	23.5	46.7
Geometric Fusion (GF)	43.5	69.1
TransFuser (Ours)	27.6	59.6

Table: We show the runtime per frame in ms for each method averaged over all timesteps in a single evaluation route. We measure runtimes for both a single model and an ensemble of three models. A single TransFuser model runs in real-time on an RTX 3090 GPU.

# Auxiliary Tasks

Auxiliary Losses	DS ↑	RC ↑	IS ↑
None	44.29	78.17	0.58
No Depth	55.66	90.87	0.61
No Semantics	53.76	88.40	0.61
No HD Map	50.96	89.52	0.58
All Losses	<b>56.68</b>	<b>92.28</b>	<b>0.62</b>

Table: Results shown are the mean over 3 evaluations on the Longest6 benchmark. HD Map prediction is the most important among the auxiliary tasks.

# Architecture

Parameter	Value	DS ↑	RC ↑	IS ↑
Fusion Scales	1	49.35	84.47	0.57
	2	53.52	91.78	0.59
	3	48.77	85.33	0.60
Attention layers	2	53.49	90.65	0.60
	6	56.24	<b>92.56</b>	0.61
	8	54.00	89.83	<b>0.62</b>
Default Config	-	<b>56.68</b>	92.28	<b>0.62</b>

Table: Results shown are the mean over 3 evaluations on the Longest6 benchmark. The default configuration with 4 fusion scales, and 4 attention layers obtains the best results.

# Model Inputs

Parameter	Value	DS ↑	RC ↑	IS ↑
LiDAR Range	64m × 32m	49.08	91.10	0.54
	64m × 64m	47.55	90.72	0.52
Camera FOV	120°	49.90	90.05	0.56
	90°	42.18	88.49	0.51
No Rasterized Goal	-	54.80	91.63	0.60
No Rotation Aug	-	<b>56.85</b>	<b>92.73</b>	0.61
Default Config	-	56.68	92.28	<b>0.62</b>

Table: Results shown are the mean over 3 evaluations on the Longest6 benchmark. The default configuration with a 32m × 32m LiDAR range and 132° camera FOV obtains the best results.