# On output representations for end-to-end driving

# Team


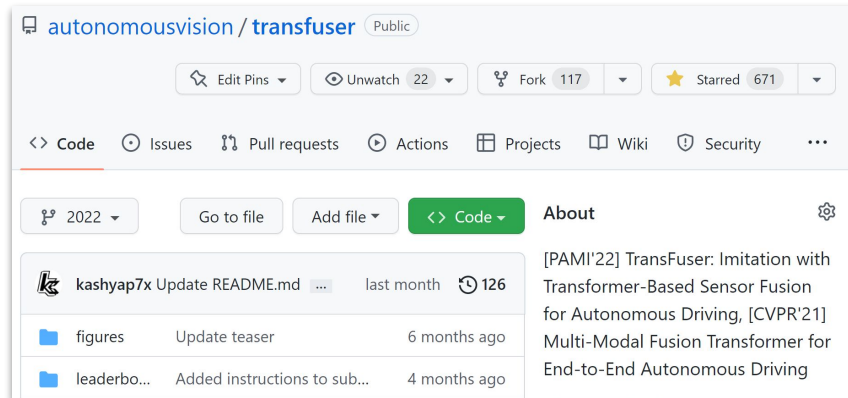Andreas Geiger


Kashyap Chitta


Katrin Renz

# Prior work and challenge

## TransFuser: Imitation with Transformer-Based Sensor Fusion for Autonomous Driving

Kashyap Chitta, Aditya Prakash, Bernhard Jaeger, Zehao Yu, Katrin Renz, and Andreas Geiger

**Abstract**—How should we integrate representations from complementary sensors for autonomous driving? Geometry-based fusion has shown promise for perception (e.g. object detection, motion forecasting). However, in the context of end-to-end driving, we find that imitation learning based on existing sensor fusion methods underperforms in complex driving scenarios with a high density of dynamic agents. Therefore, we propose TransFuser, a mechanism to integrate image and LiDAR representations using self-attention. Our approach uses transformer modules at multiple resolutions to fuse perspective view and bird's eye view feature maps. We experimentally validate its efficacy on a challenging new benchmark with long routes and dense traffic, as well as the official leaderboard of the CARLA urban driving simulator. At the time of submission, TransFuser outperforms all prior work on the CARLA leaderboard in terms of driving score by a large margin. Compared to geometry-based fusion, TransFuser reduces the average collisions per kilometer by 48%.
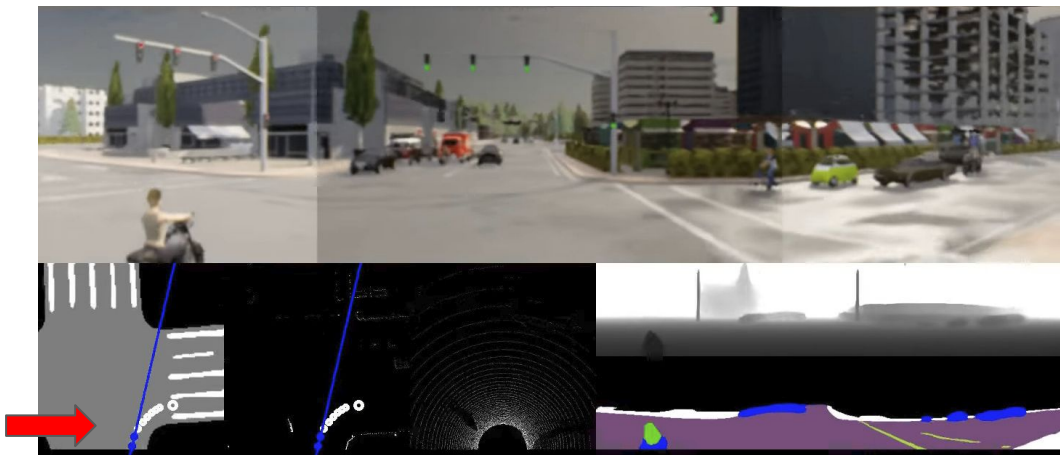
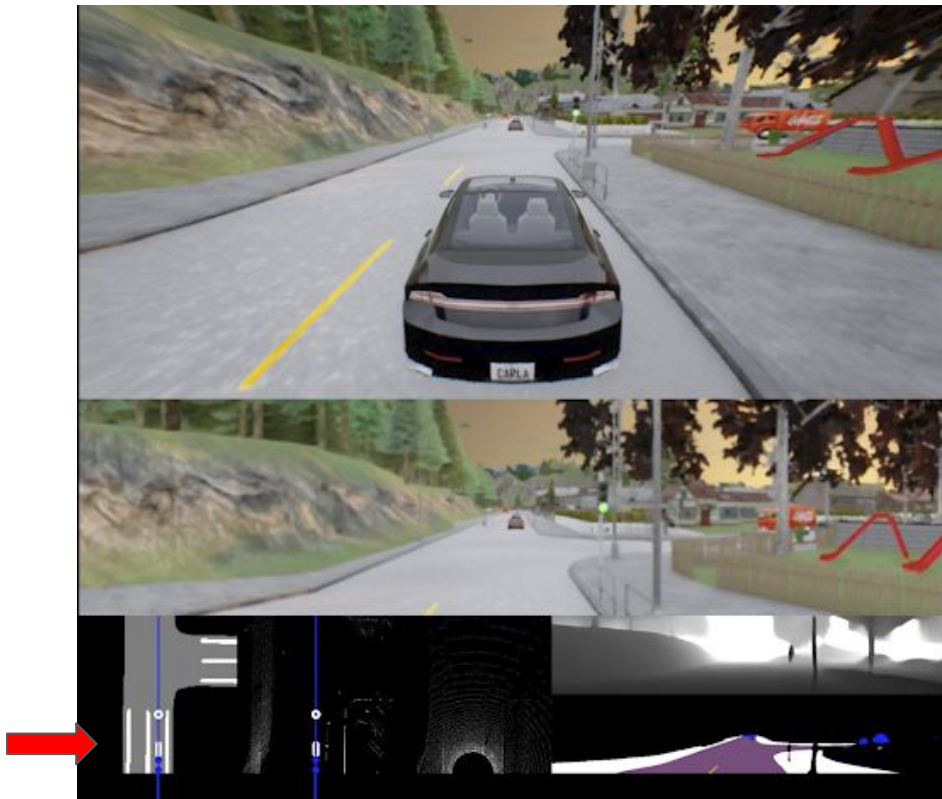**Index Terms**—Autonomous Driving, Imitation Learning, Sensor Fusion, Transformers, Attention.

autonomousvision / **transfuser**  Public

Edit Pins ▾   Unwatch 22 ▾   Fork 117 ▾   Starred 671 ▾

<> Code   ⊙ Issues   ⇄ Pull requests   ⊙ Actions   ⊞ Projects   📖 Wiki   ⊘ Security   ⋯

2022 ▾   Go to file   Add file ▾   <> Code ▾

kashyap7x Update README.md  …   last month  ⏱ 126

figures   Update teaser   6 months ago

leaderbo…   Added instructions to sub…   4 months ago

**About**

[PAMI'22] TransFuser: Imitation with Transformer-Based Sensor Fusion for Autonomous Driving, [CVPR'21] Multi-Modal Fusion Transformer for End-to-End Autonomous Driving

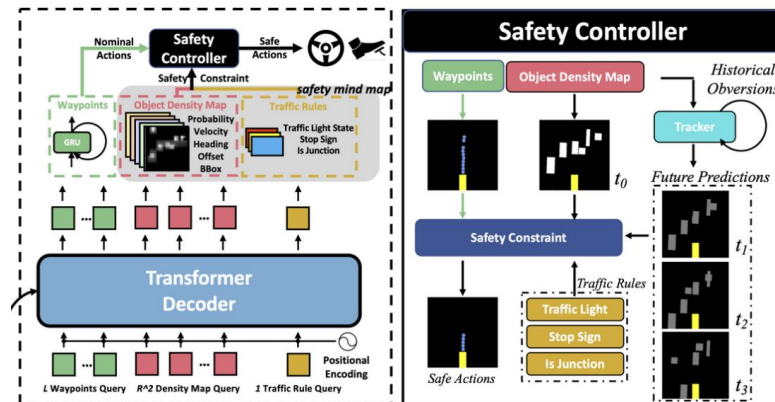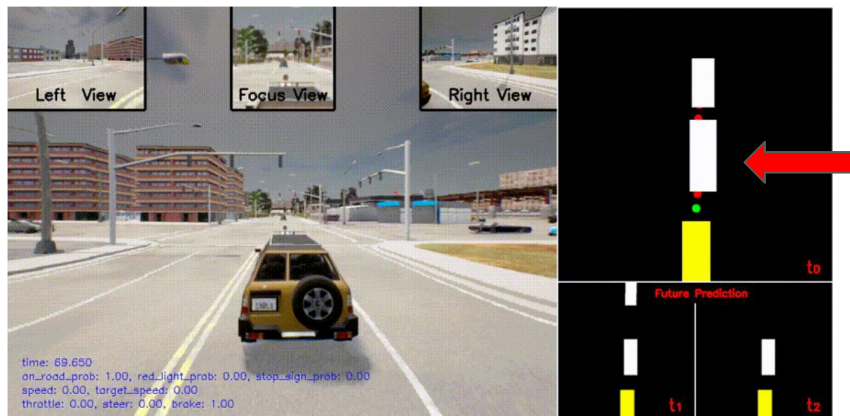| Method (Leaderboard 1.0) | DS ↑ | RC ↑ | IS ↑ | Driving Speed ↑ |
|---|---|---|---|---|
| TransFuser | 61 | 87 | 0.71 | 1x |
| Map TransFuser++ | 61 | 82 | 0.70 | 2x |

3

# The SotA for outputs: Waypoint prediction

- Predicts own future positions
  - location 0.5s, 1.0s, 1.5s, 2.0s … into the future
- Used by TF, LatentTF, LAV, TCP and MMFN
- Entangles path + **future velocities**

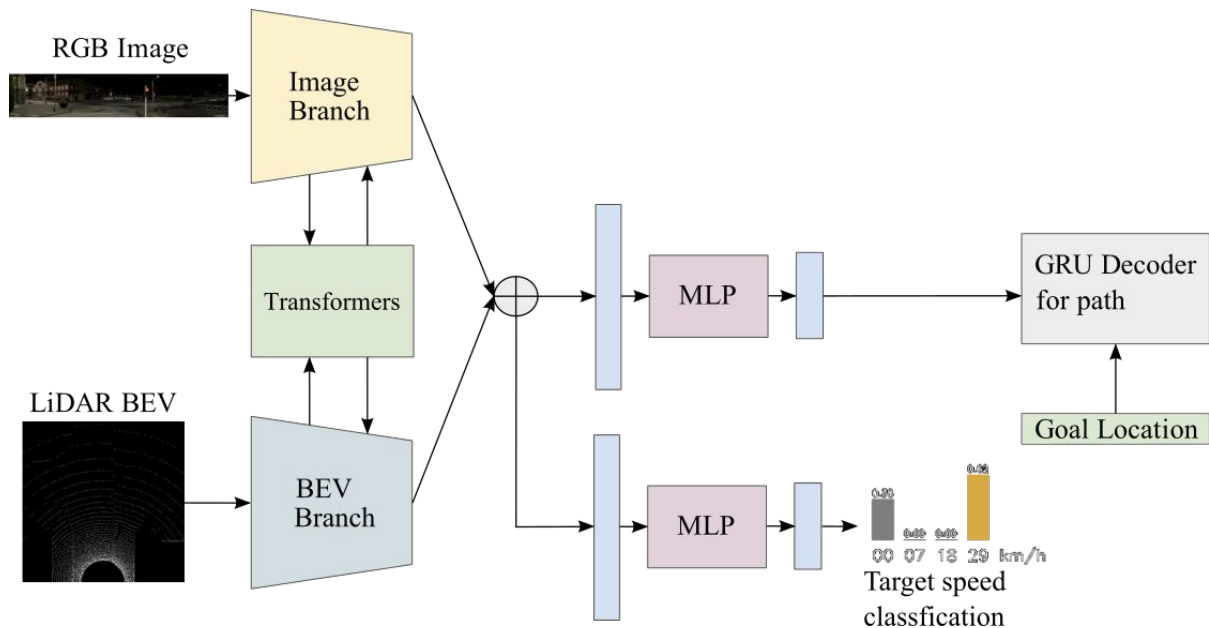# Future velocities are ambiguous!

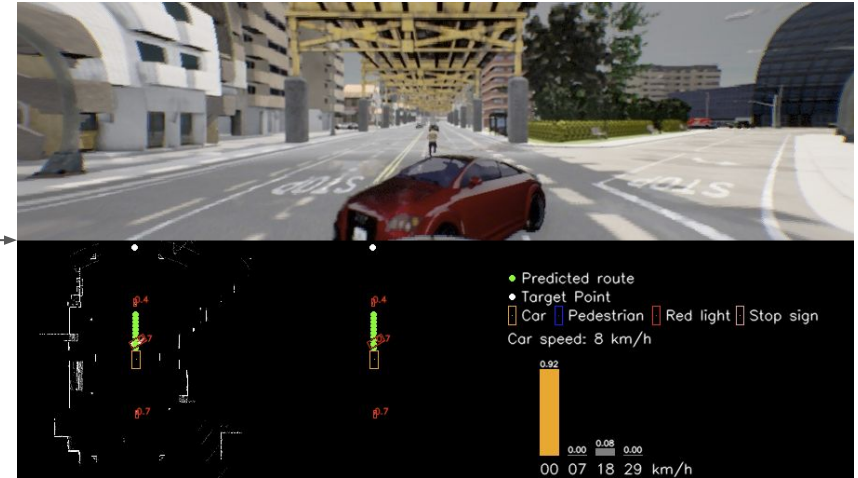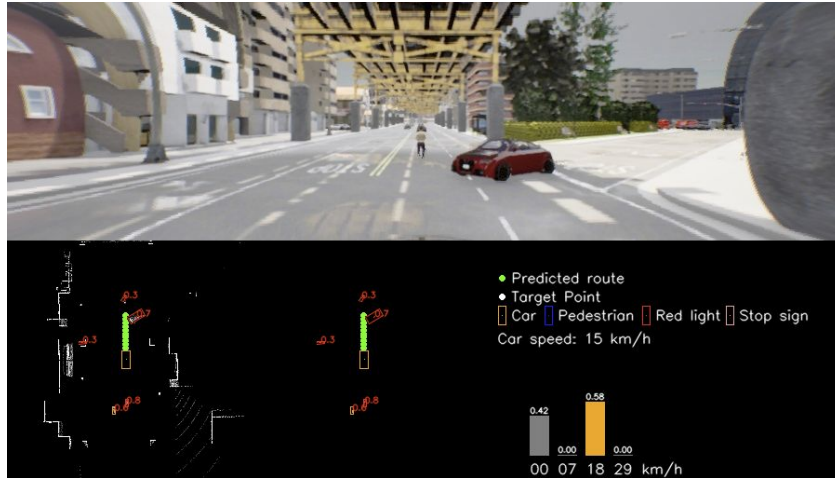# Predict path without velocities



- ● How to do longitudinal control?

- ● Hand-crafted heuristics?

H. Shao, L. Wang, R. Chen, H. Li and Y. Liu "**Safety-Enhanced Autonomous Driving Using Interpretable Sensor Fusion Transformer**" CoRL, 2022

# Target speed classification!

- Simple and learnable
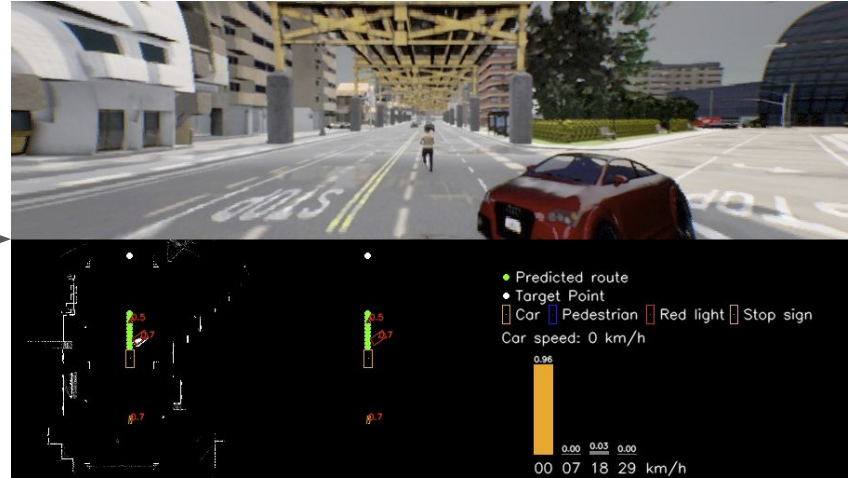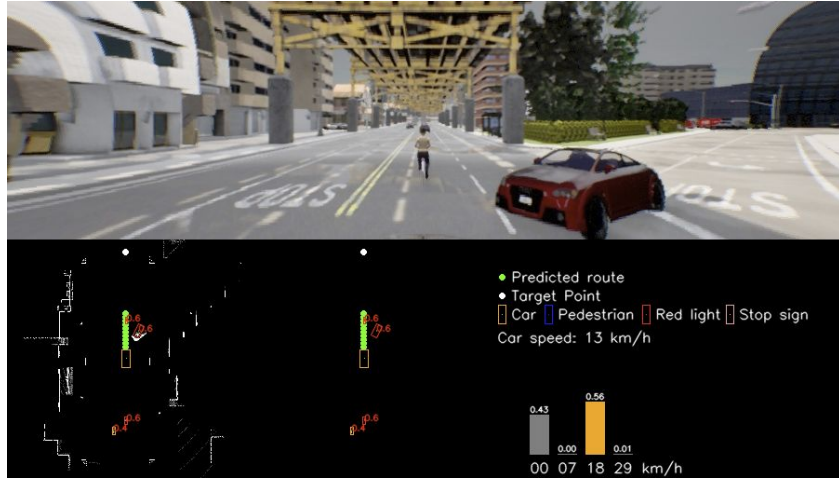  - Provides **confidence**

# Qualitative example: Argmax target speed



- **Idea: Confidence weighted average!**
  - Slow down when uncertain

# Qualitative example: Uncertainty weighted target speed

# Ablations

| Method (longest6 benchmark) | DS ↑ | RC ↑ | IS ↑ | Collisions / km ↓ |
|---|---|---|---|---|
| TransFuser++ Waypoints | 56 | 88 | 0.62 | 0.79 |
| TransFuser++ Disentangled + Argmax | 52 | **97** | 0.53 | 1.31 |
| TransFuser++ Disentangled + Uncertainty | **69** | **96** | **0.71** | **0.63** |

# System comparison

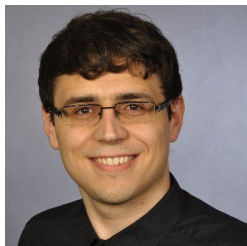| Method (longest6 benchmark) | DS ↑ | RC ↑ | IS ↑ | collisions / km ↓ |
|---|---|---|---|---|
| World on Rails [1] | 21 | 48 | 0.56 | 1.05 |
| TransFuser [2] | 47 | 93 | 0.50 | 2.44 |
| TransFuser++ Waypoints (ours) | 56 | 88 | 0.62 | 0.79 |
| LAV version 2 [3] | 58 | 83 | 0.68 | 0.69 |
| Perception PlanT [4] | 58 | 88 | 0.65 | 0.97 |
| TransFuser++ Disentangled (ours) | **69** | **96** | **0.71** | **0.63** |

[1] D. Chen, V. Koltun, and P. Krähenbühl, "**Learning to drive from a world on rails**" ICCV, 2021
[2] K. Chitta, A. Prakash, B. Jaeger, Z. Yu, K. Renz, and A. Geiger, "**Transfuser: Imitation with transformer-based sensor fusion for autonomous driving**" PAMI, 2022
[3] D. Chen and P. Krähenbühl, "**Learning from all vehicles**." CVPR, 2022
[4] K. Renz, K. Chitta, O. Mercea, A. S. Koepke, Z. Akata and A. Geiger, "**PlanT: Explainable Planning Transformers via Object-Level Representations**" CoRL, 2022
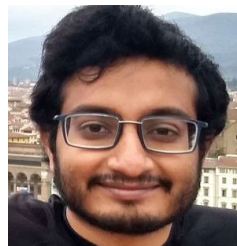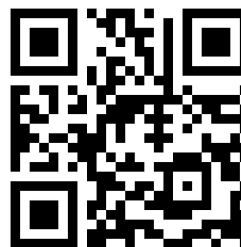
# Thank you!

@Kait0o0

@AutoVisionGroup

@kashyap7x

@KatrinRenz