mpii
max planck institut
informatik

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

Renz K.
Koepke A.S.

Chitta K.
Akata Z.

Mercea O.B.
Geiger A.

# PLANT:

## EXPLAINABLE **PLANNING** TRANSFORMERS VIA OBJECT-LEVEL REPRESENTATIONS
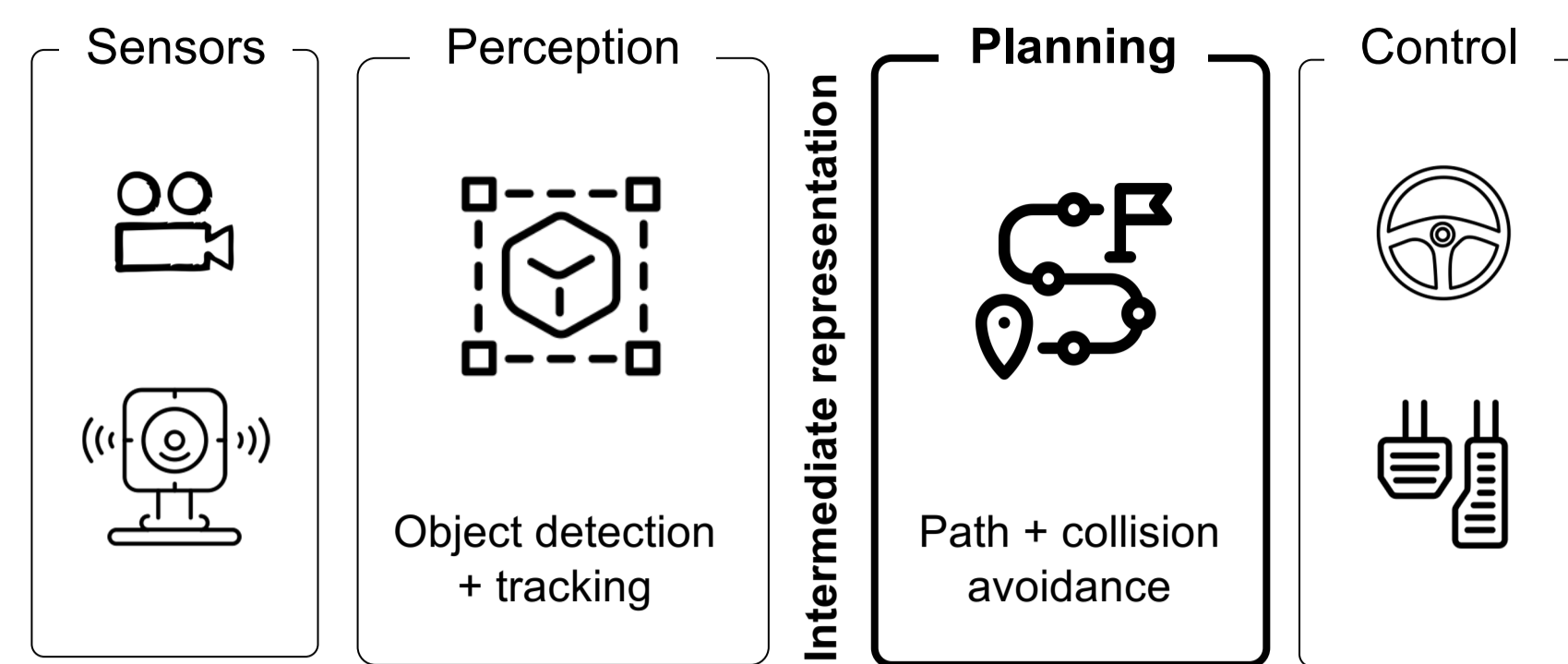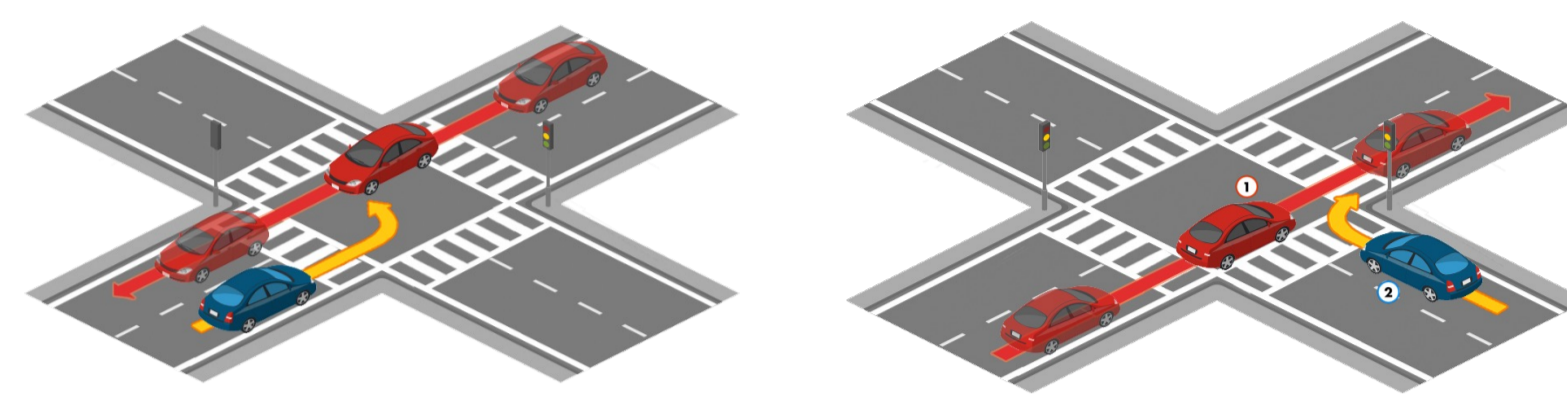
## Abstract

Planning an optimal route in a complex environment requires efficient reasoning about the surrounding scene. In this paper, we propose PlanT, a novel approach for that uses a standard transformer architecture. PlanT is based on imitation learning with a compact object-level input representation. Combining PlanT with an off-the-shelf perception module provides a sensor-based driving system that is more than 10 points better in terms of driving score than the existing state of the art.
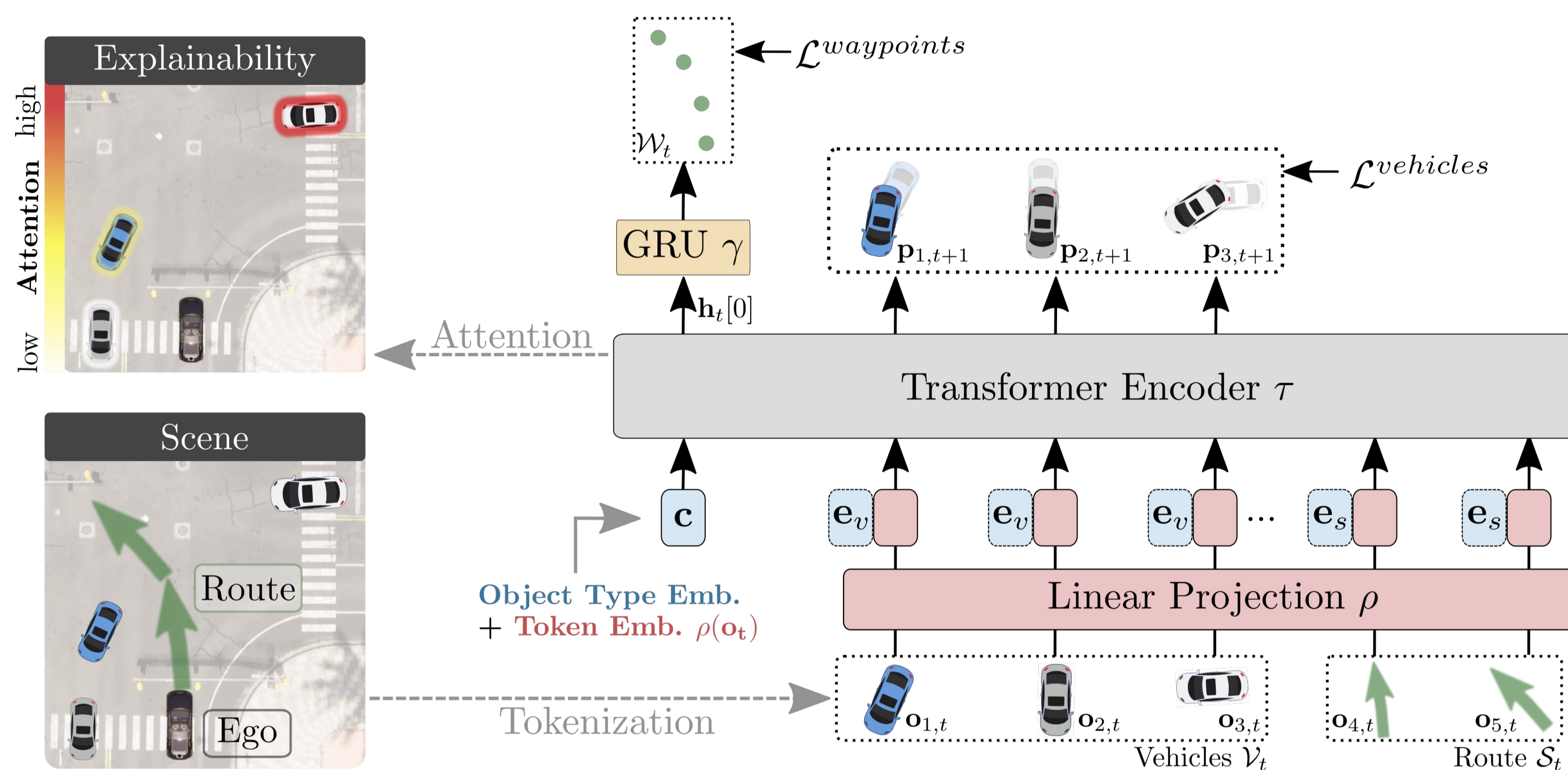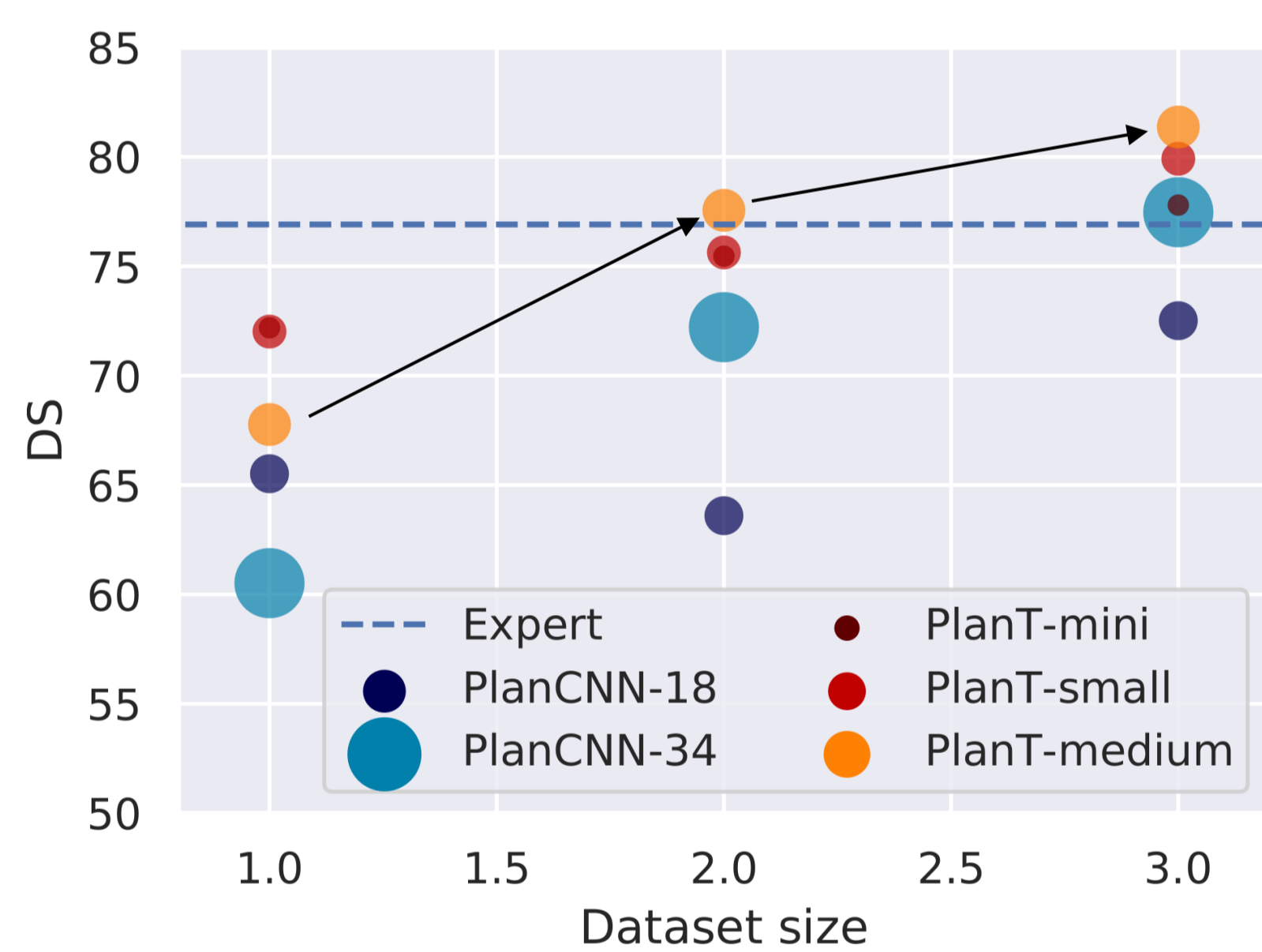
## Task

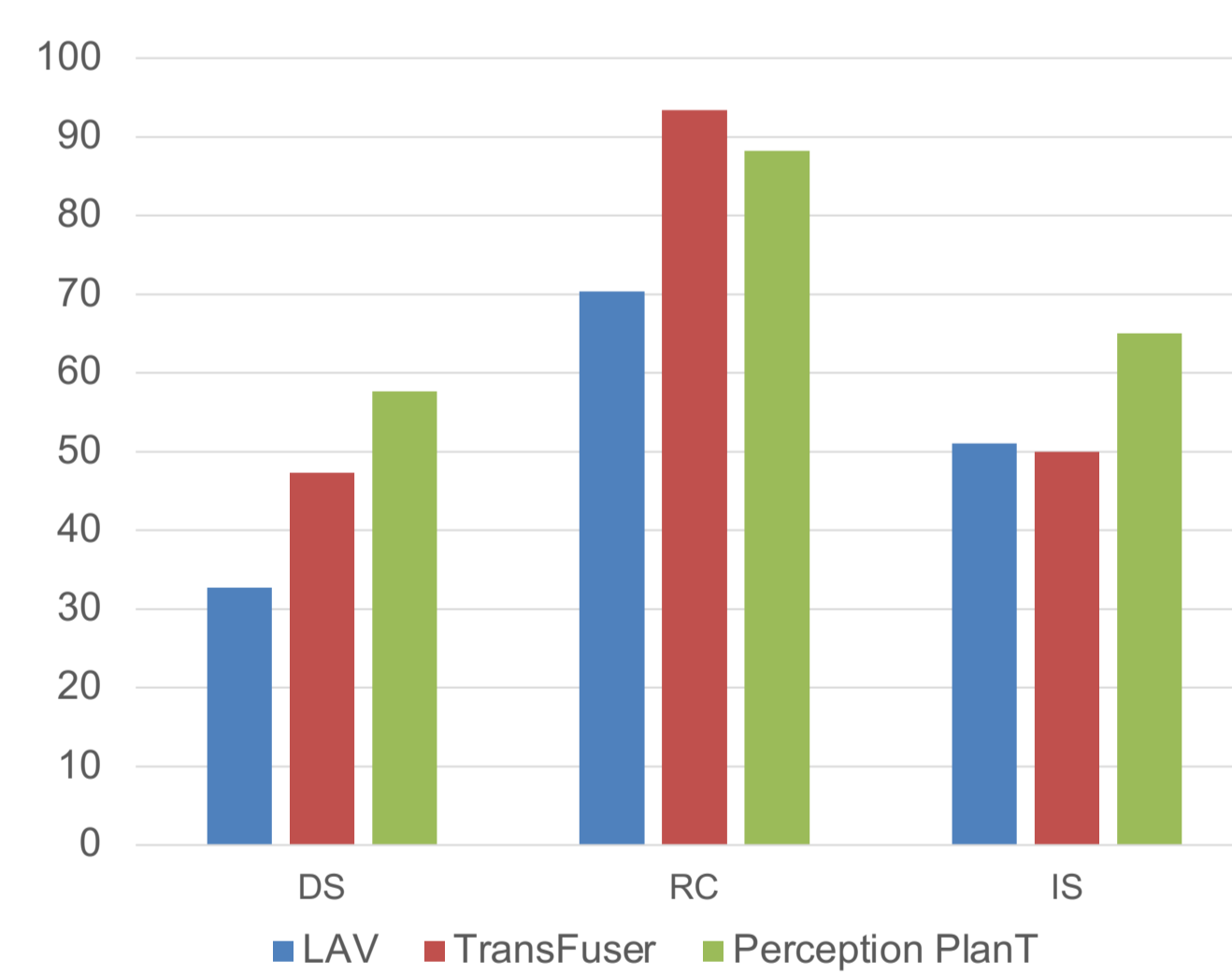

CARLA scenarios



## Architecture



- Train **a standard transformer encoder** from scratch
- Loss on **future positions** of ego vehicle and the other vehicles

## Results



- **Scaling** dataset and model improves performance
- **Expert level performance**

## Perception PlanT



- Adding a perception module
- State of the art on longest 6 benchmark

## Explainability



- Visualization of **attention weights** to show the **most important object**
- **Temporarily** more **consistent** than the CNN-based method + also takes **geometrically distant** objects into account