

Hochschule München
University of Applied Sciences

Fakultät für Informatik und Mathematik
Department of Computer Sciences and Mathematics

Multimodal Trajectory Prediction in Multi-Agent Scenarios

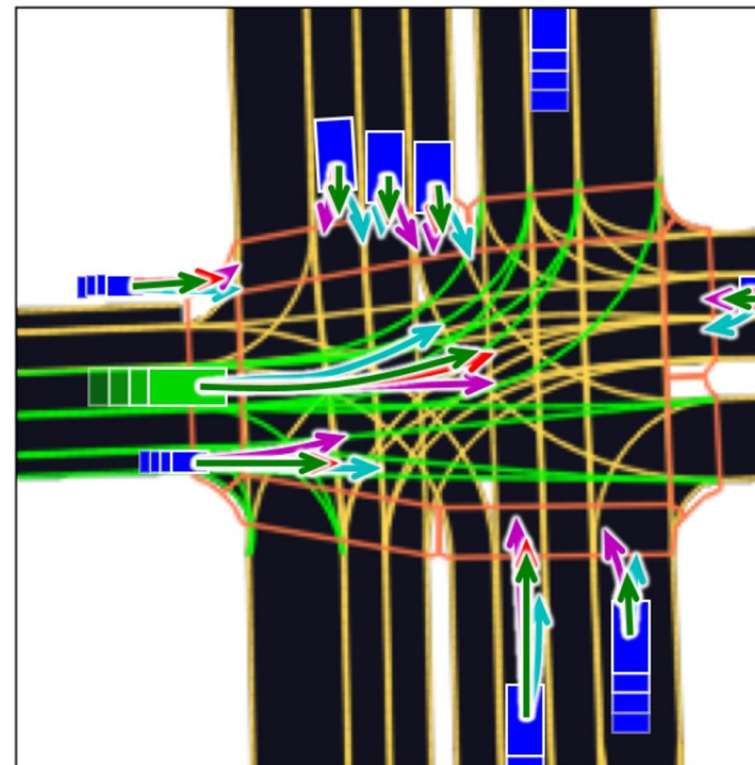
Seminar: Video Analysis
& Object Tracking

16 April 2025
Lukas Röß, Jan Duchscherer



Introduction & Motivation

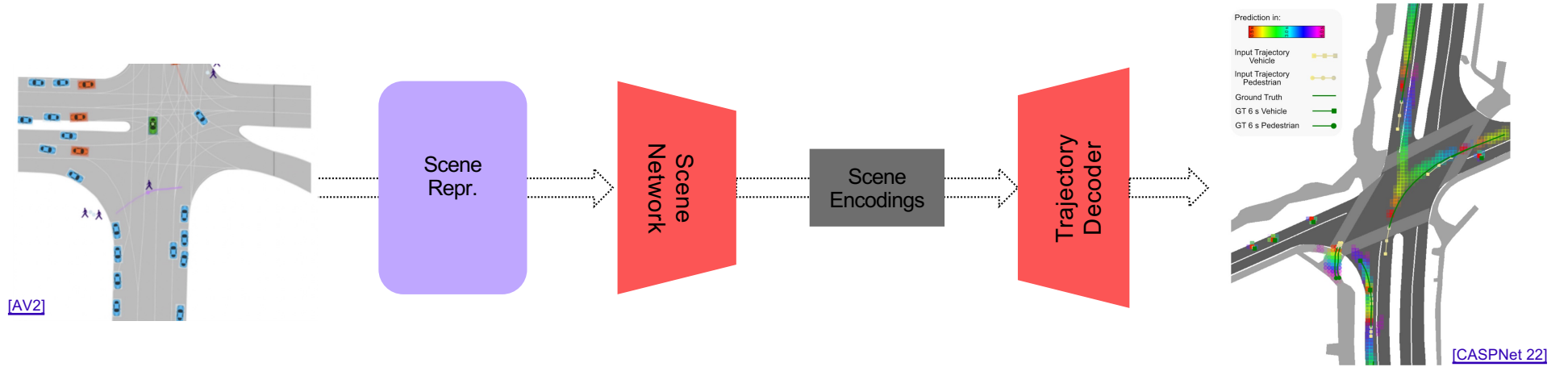
- What is Multimodal Trajectory Prediction?
- Why Multi-Agent Scenarios Matter?
- Challenges



[Explainable multimodal trajectory prediction using attention models 24](#)

From Scene to Prediction

The General Motion Forecasting Pipeline



$$\mathbf{X}_d \in \mathbb{R}^{N \times T_{\text{in}} \times \|F_d\|}$$

$$\mathbf{X}_s \in \mathbb{R}^{K \times \|F_s\|}$$

$$F_d = [x, y, \psi, \dot{x}, \dot{y}, \ddot{x}, \ddot{y}, w, l, c]$$

$$N = \# \text{ agents}$$

$$T_{\text{in}} = \text{past timesteps}$$

$$M = \# \text{ modes}$$

$$F_s = \text{polyline encodings}$$

$$K = \# \text{ map objects}$$

$$T_{\text{out}} = \text{pred horizon}$$

$$\boldsymbol{\pi} \in \mathbb{R}^{N \times M}$$

$$\boldsymbol{\mu}, \mathbf{b} \in \mathbb{R}^{N \times M \times T_{\text{out}} \times 2}$$

Selected Approaches

CASPFormer (2024):

(SOTA on nuScenes)

- HD Map-Free
- CNN + ConvLSTM Encoder
- Deformable Attention Decoder

MTR (2022):

(MTR v3 SOTA on Waymo)

- Map-Aware
- Intention-Based Queries
- Modular Transformer Design

Project Objective

Main Goal:

- Implement CASPFormer in the UniTraj framework
- Train MTR on dataset within UniTraj if no pre-trained model will be shared

Comparative Analysis:

- Compare models based on algorithmic structure
- Evaluate both models on a shared dataset using key metrics:
 - Average Displacement Error (ADE)
 - Final Displacement Error (FDE)
 - Miss Rate (MR)

Challenges:

- CASPFormer source code not available
 - Use fallback
 - Or implement SmolCASPFormer
- Adapt data processing for the dataset of selected approaches
- Train/Fine-tune MTR Model
- Ensure fair and comparable evaluation, especially if different dataset are used datasets

CASPFomer

BEV (Birds Eye View)

⇒ CNN + ConvLSTM

⇒ Multi-Scale Scene Encodings

⇒ Deformable (recurrent) Transformer

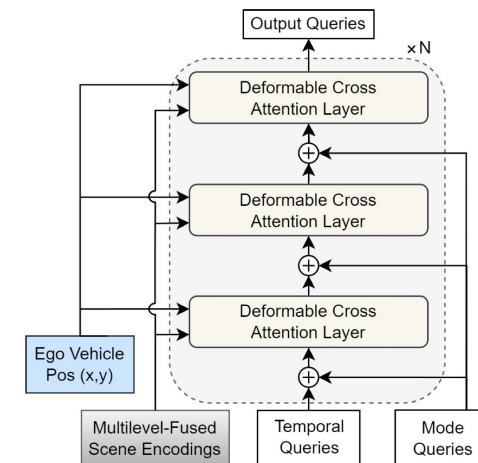
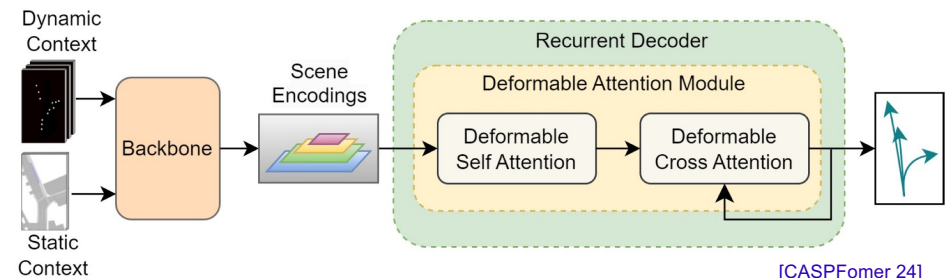
(Mode Queries and Temporal Queries)

⚡ No Code available!

⇒ Try implementing smol version from scratch? 🙄

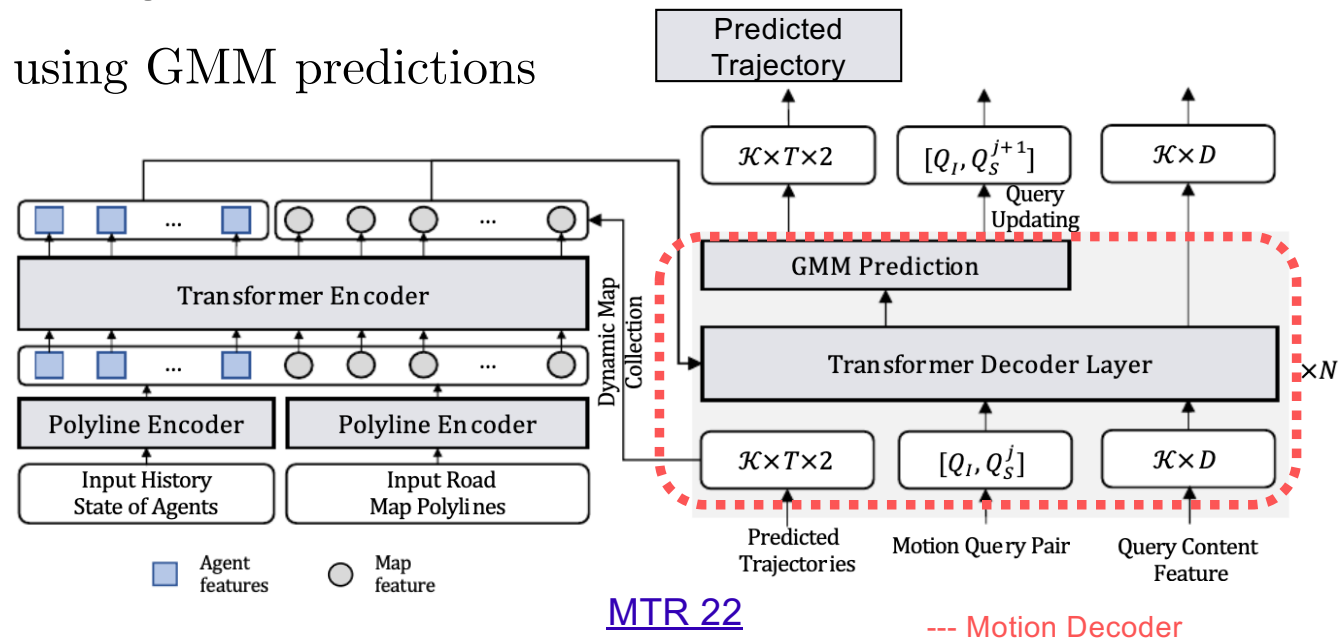
⇒ QCNet as source code blueprint
(loss functions, train pipeline, encodings...)

⇒ Most likely will not yield a working prototype!

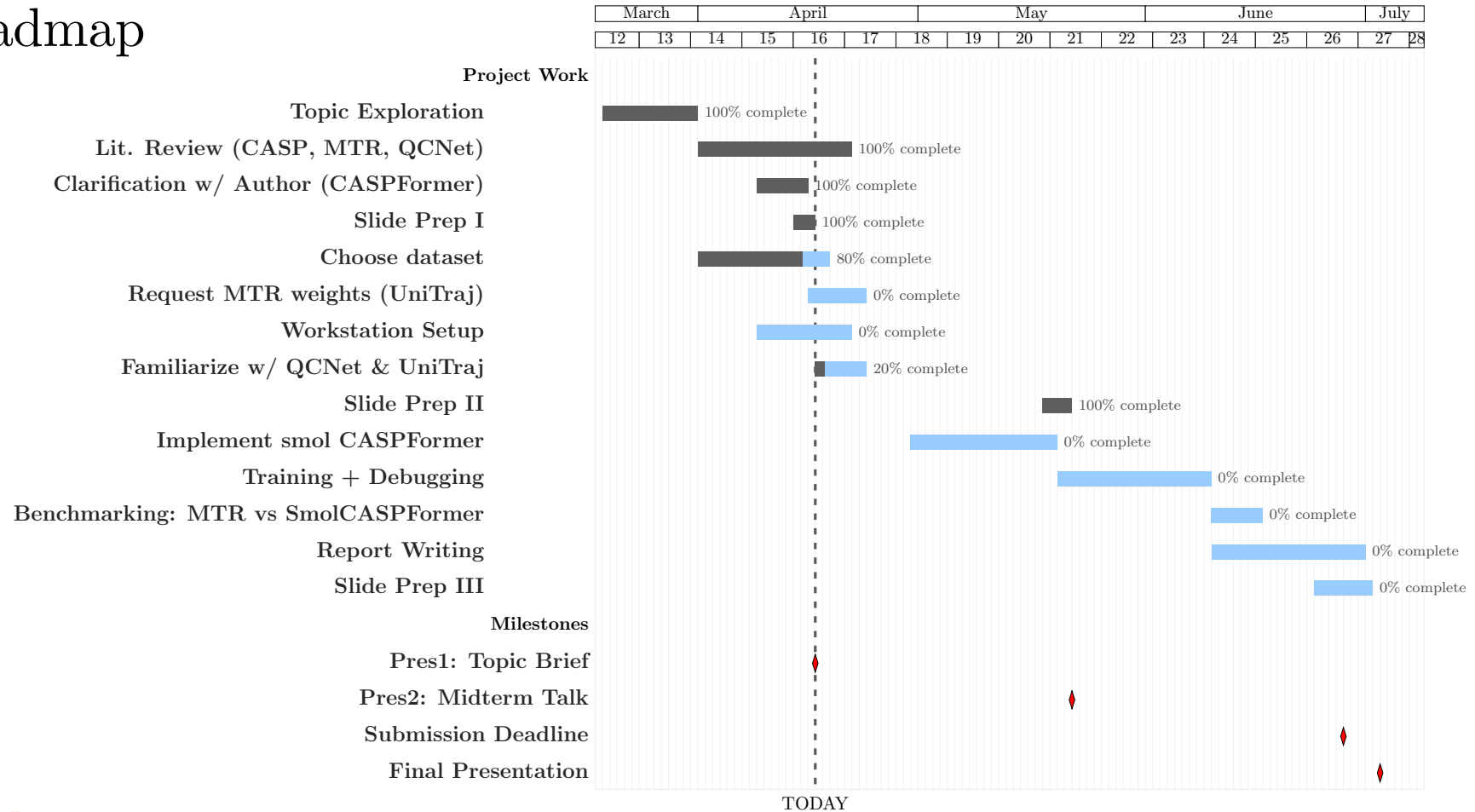


Motion TRansformer (MTR)

- Modular Design
- Goal-conditioned decoding
- Multimodal output using GMM predictions



Roadmap



Discussion