

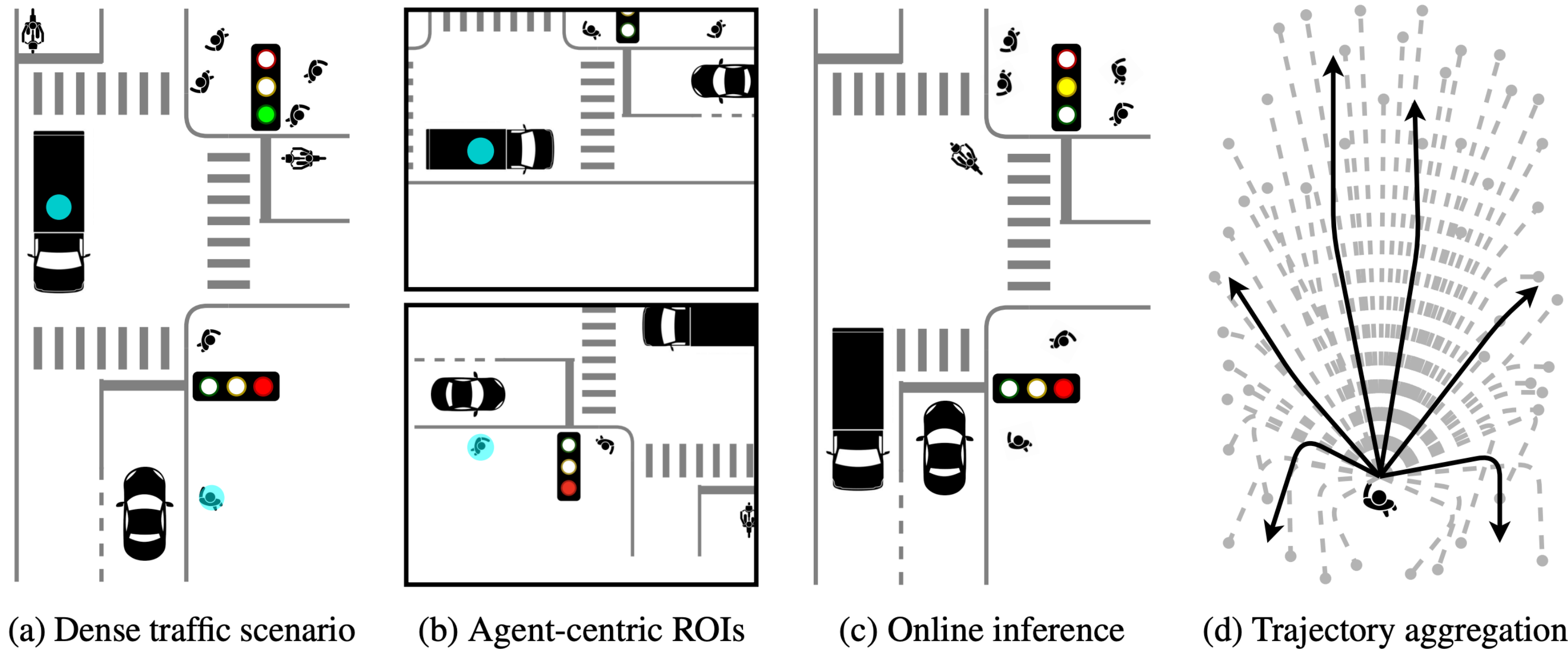
Real-Time Motion Prediction via Heterogeneous Polyline Transformer with Relative Pose Encoding

Zhejun Zhang¹, Alexander Liniger¹, Christos Sakaridis¹, Fisher Yu¹, and Luc Van Gool^{1,2,3}

¹Computer Vision Lab, ETH Zurich, CH. ²PSI, KU Leuven, BE. ³INSAIT, Un. Sofia, BU.

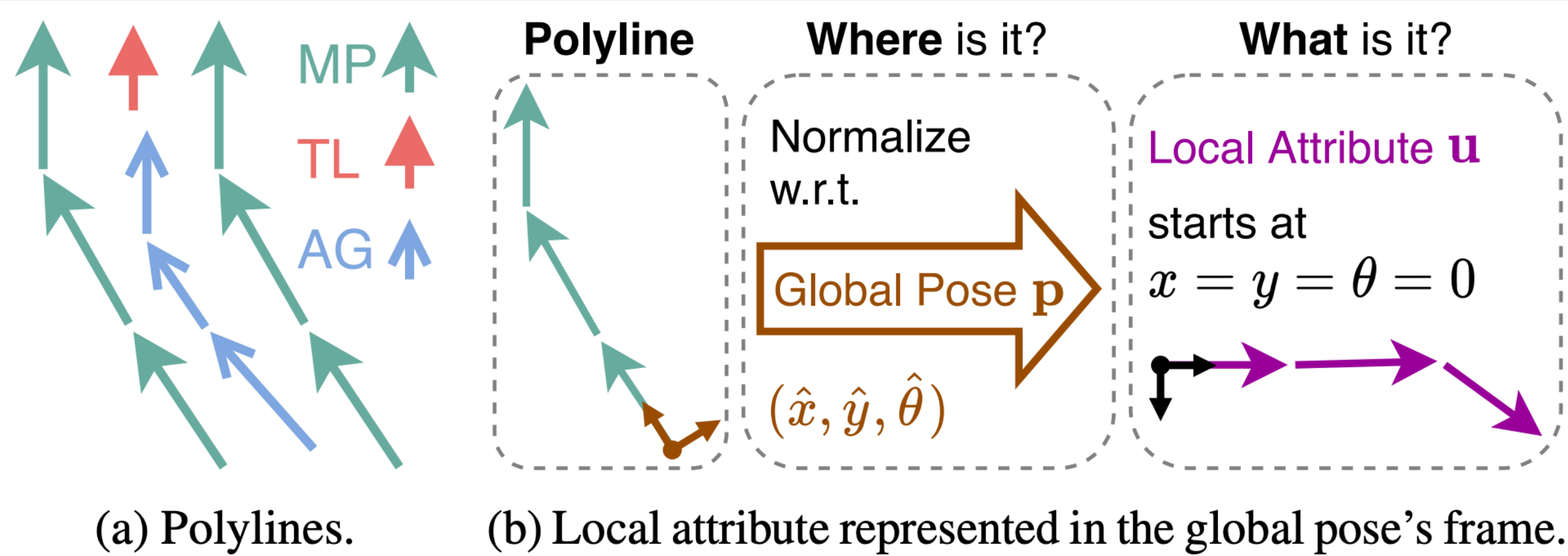


Motivation



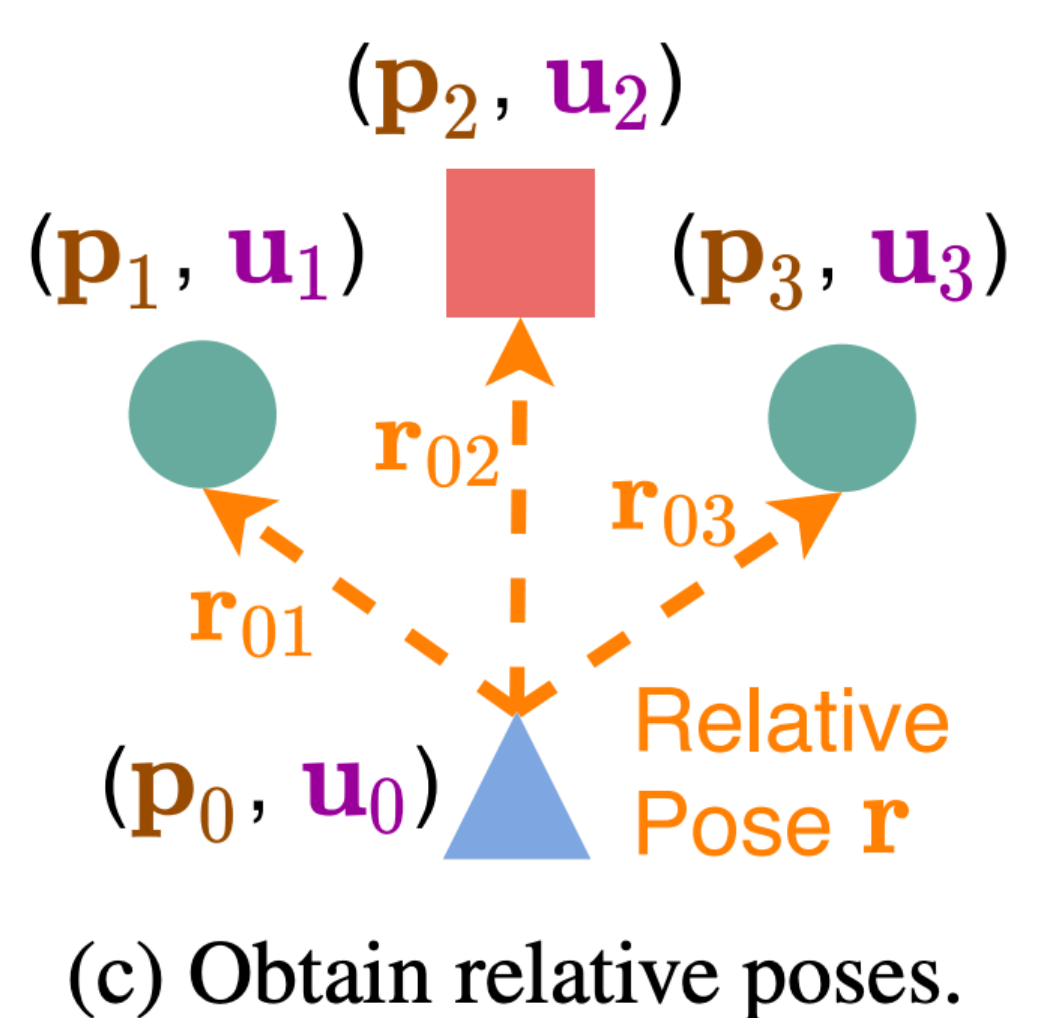
- a) **Real-time and on-board** motion prediction in dense urban scenario.
- b) Agent-centric: **Good** performance. **Bad** scalability.
- c) **Online inference** with streaming inputs.
- d) Expensive post-processing and ensembling.

Pairwise-Relative Representation



Global Pose p: Where is the polyline?

Local Attribute u: What kind of polyline is it?



Input to the Network:

3-dimensional relative pose r.

High-dimensional local attribute u.

Rotation and translation invariance.

Good scalability by sharing **u**.

So far only exploited by GNNs.

KNARPE

K-nearest **N**eighbor **A**ttention with **R**elative **P**ose **E**ncoding.

$$\mathbf{z}_i = \text{KNARPE}(\mathbf{u}_i, \mathbf{u}_j, \mathbf{r}_{ij} \mid j \in \kappa_i^K) = \sum_{j \in \kappa_i^K} \alpha_{ij} (\mathbf{u}_j \mathbf{W}^v + \mathbf{b}^v + \text{RPE}(\mathbf{r}_{ij}) \hat{\mathbf{W}}^v + \hat{\mathbf{b}}^v)$$

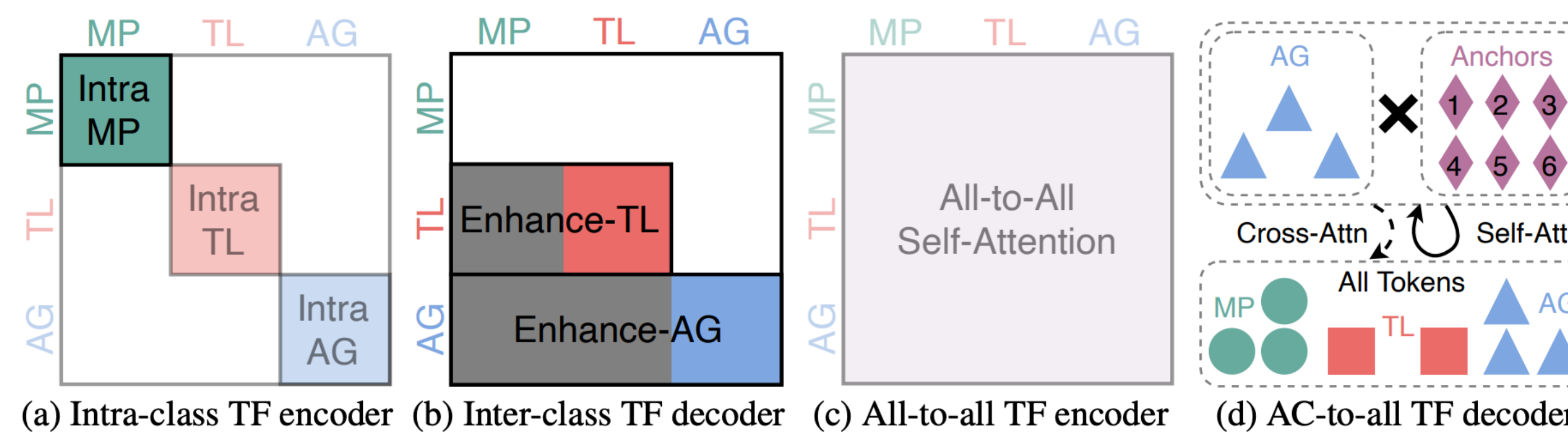
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in \kappa_i^K} \exp(e_{ik})}, \quad e_{ij} = \frac{(\mathbf{u}_i \mathbf{W}^q + \mathbf{b}^q)(\mathbf{u}_j \mathbf{W}^k + \mathbf{b}^k + \text{RPE}(\mathbf{r}_{ij}) \hat{\mathbf{W}}^k + \hat{\mathbf{b}}^k)}{\sqrt{D}}$$

$$\text{RPE}(\mathbf{r}_{ij}) = \text{concat}(\text{PE}(x_{ij}), \text{PE}(y_{ij}), \text{AE}(\theta_{ij})),$$

- Based on multi-head dot-product attention.
- Implemented with **basic matrix operations**: matrix indexing, summation and element-wise multiplication.
- Self-Attention: Local context aggregation like CNN.
- Cross-Attention: Rotated ROI alignment with CNN.

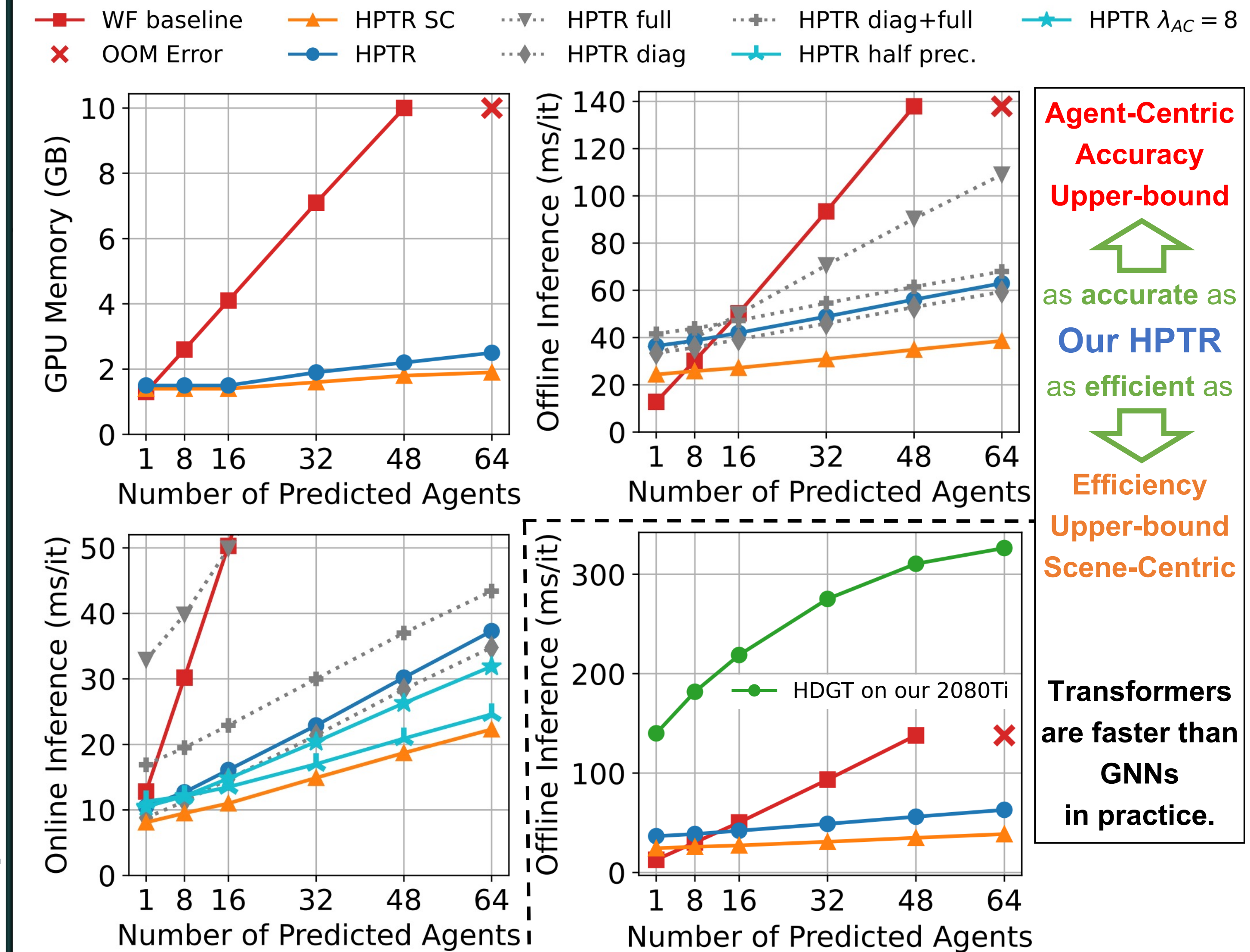
HPTR

Heterogeneous **P**olyline **T**ransformer with **R**elative pose encoding



- Based on KNARPE.
- Inspired by Wayformer.
- Remove redundant attentions.
- Transformers are organized in a hierarchical way.
- Intermediate results can be cached and reused.
- **Asynchronous token update** during online inference.

Performance



Summary

- **KNARPE** allows the pairwise-relative representation to be used by Transformers.
- **HPTR** uses hierarchical architecture to enable asynchronous token update.
- **SoTA** performance among E2E methods: WOMD and AV2 dataset (c.f. paper).
- **Good Performance**: As accurate as agent-centric methods.
- **Good Scalability**: As efficient as scene-centric methods.
- **Real-Time and On-Board** Motion Prediction: **40 FPS** during online inference. 80% reduction on online inference latency and GPU memory.
- **Code**: <https://github.com/zhejz/HPTR>