

Spatial interaction:

- \succ Spatial interaction: relations in Euclidean space. The relative spatial relation matters
- \succ Common and critical. E.g., forecasting behaviors of traffic actors

Approaches to model spatial interaction:

 \succ Graph neural networks (GNNs)

Node: per-actor features

Edge: relative relations (e.g., positions and velocities)

Interaction: aggregate edges and neighboring nodes

- . Have to hand-craft and add relative relations to the edges
- 2. Slower than Convolutional neural networks (CNN)
- \succ How about convolutional layers or CNN?
 - Intuitively, convs can model spatial interactions
 - 2D and 3D conv-layers operate on data in grid forms \rightarrow spatial relations are intrinsically represented
 - \blacksquare Large receptive fields \rightarrow aggregate non-local information
 - Why interaction is ineffectively modeled, even large CNN backbones are widely used?

Questions:

- > How to effectively model interaction using convolutions?
- How effective are convolutions compared to GNNs?

Convolutions for Spatial Interaction Modeling

Su, Chao, David, Carlos, Carl, Nemanja @ Aurora

Experiments

- Test field: forecasting trajectories of traffic actors 2D top-down view
- Test data: large autonomous driving data
- Control experiments: models are almost identical, except for additional light-weight conv-layers and GNNs
- \succ Explicite interaction metrics: overlaps (i.e. collisions)





Voxelized lidar point-cloud at an intersection

Takeaways of the experiments:

- > 3 characteristics to "activate" conv-layers:
 - Large and relevant context as the input to conv-layers
 - Aggregation of per-actor feature maps using a few downsampling conv-layers
 - Overcoming the rotational variance of convolutions
- Conv approach vs. GNN
 - Convs can perform similarly to or better than GNNs
 - Adding the convs considerably improves interaction modeling even when a GNN is used
 - Adding a GNN demonstrates only minor additional gain when the convs is already used

Results: Effective convs







Baseline

Convs (crop size of 60m)



Summary

- interaction



Inference times of additional modules (baseline takes 45.6 ms)

Module	Crop size (m)	Inference (ms)
Convs	0	5.2
Convs	80	8.1
GNN	-	46.9

 \succ Identified 3 characteristics that affects conv-layers in model spatial

 \succ 2D motion forecasting evidences that convs demonstrates comparable or stronger ability than GNNs in modeling interaction \succ Future: generalization to other 2D and 3D tasks with interactions Videos (link) and supplementary (link)