

A Holistic View of Perception in Intel. Vehicles

Part IV: Key Takeaways and Future Directions

Objectives

Objectives in Part IV

- Takeaway Messages and Key Insights
- Unaddressed Challenges in Perception
 - Context Awareness
 - Embedded Perception
 - V2X Perception
- Future Research Directions
 - Temporal Processing
 - Sensor Processing Architectures
 - Sensors research
 - Infrastructure + AV Datasets

Objectives

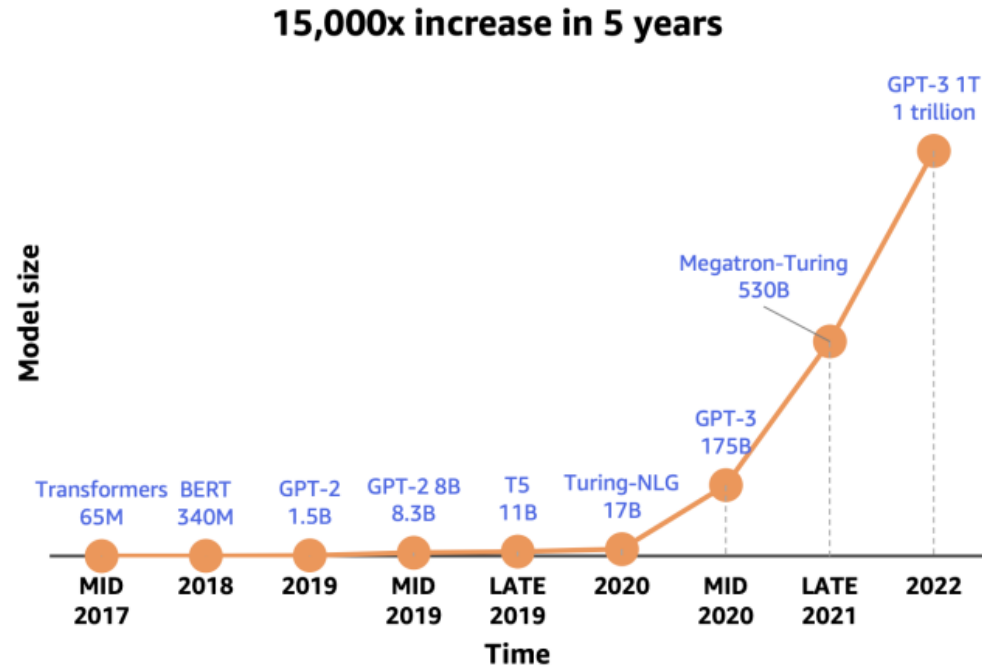
Takeaway Messages and Key Insights

- **Robustness** under challenging conditions, environments, context and surroundings-awareness are **challenges** in AV perception
 - **Deep Learning** provides a **holistic solution** to a number of the above challenges
- **Transfer Learning** and **training at scale** help to create foundation models
 - **Self-supervised Learning** provides a framework for large scale learning on unannotated data
- It is not always clear if aberrant events and challenges must be incorporated in training
 - Instead, **model predictions** must be equipped with **diagnostic tools** at inference
 - These diagnostic tools are **anomaly and uncertainty scores** for decision making and **contextual explainability** for post-hoc stakeholders
 - **Gradients** provide the change induced by an aberrant event in the network and can be used to obtain the required **prediction diagnosis**

Perception in AVs

Unaddressed Technical Challenges for Level 3 Automation

- Challenging weather
- Challenging sensing
- Challenging environments
- **Context awareness**
- **Embedded perception**
- **V2X perception**



- Foundation models are great but the real-time feasibility is an issue
- The inaccuracies from model outputs is dangerous in urban settings

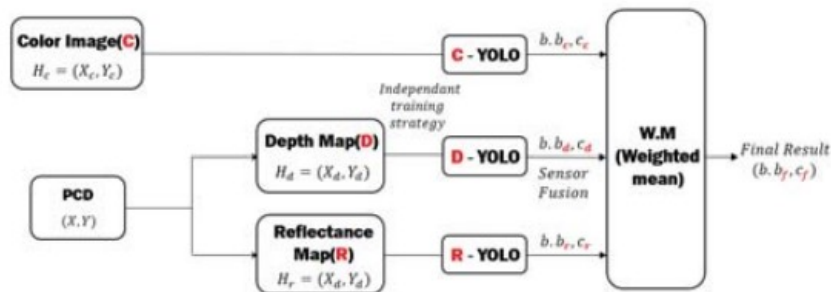
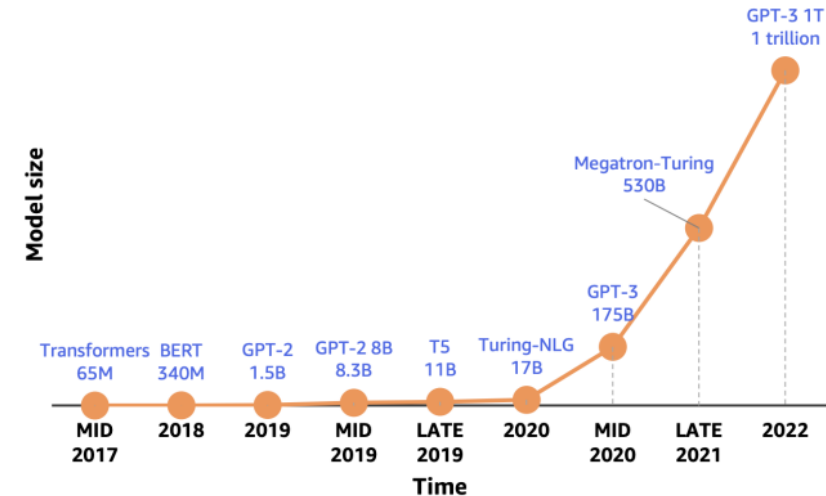
Perception in AVs

Unaddressed Technical Challenges for Levels 4 and 5

Foundation models with multiple sensor modalities

- Challenging weather
- Challenging sensing
- Challenging environments
- **Context awareness**
- **Embedded perception**
- **V2X perception**

15,000x increase in 5 years



- Levels 4 and 5 automation relies on roadside infrastructure to obtain high-resolution predictions
- 10x is the rough estimate of the increase in processing power between levels of automation

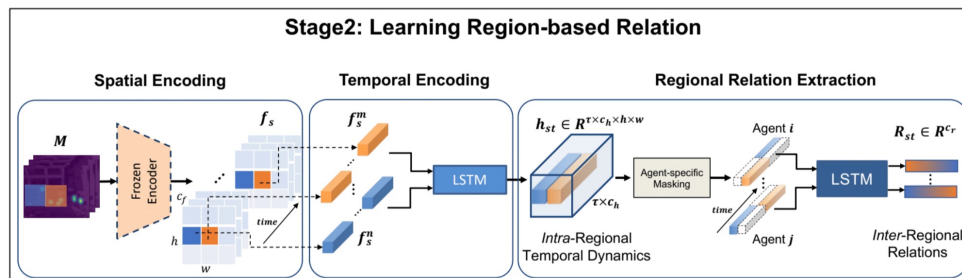
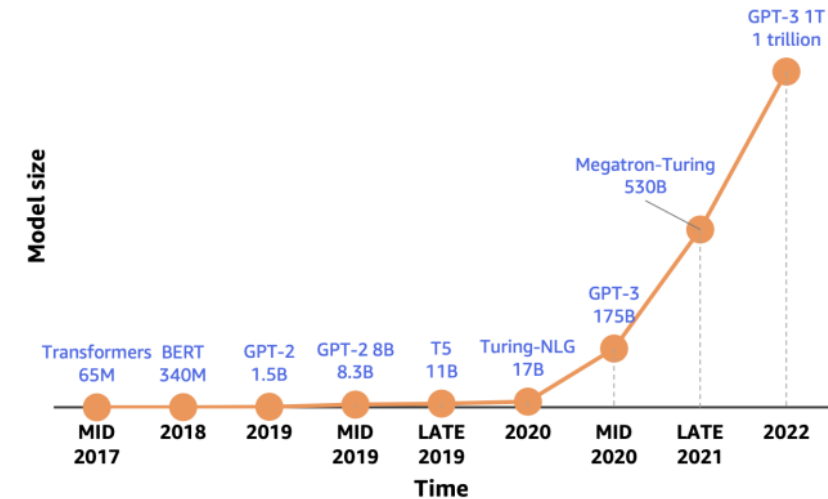
Perception in AVs

Unaddressed Technical Challenges for Levels 4 and 5

Foundation models with multiple sensor modalities and on temporal data

- Challenging weather
- Challenging sensing
- Challenging environments
- **Context awareness**
- **Embedded perception**
- **V2X perception**

15,000x increase in 5 years

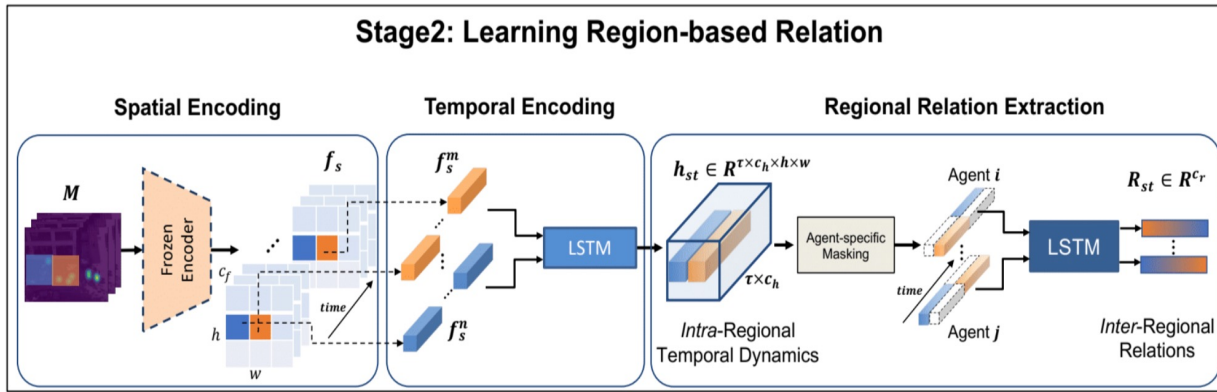


- Levels 4 and 5 automation relies on roadside infrastructure to obtain high-resolution predictions
- 10x is the rough estimate of the increase in processing power between levels of automation
- **Current temporal processing = linear spatial processing in time**

Future Direction 1

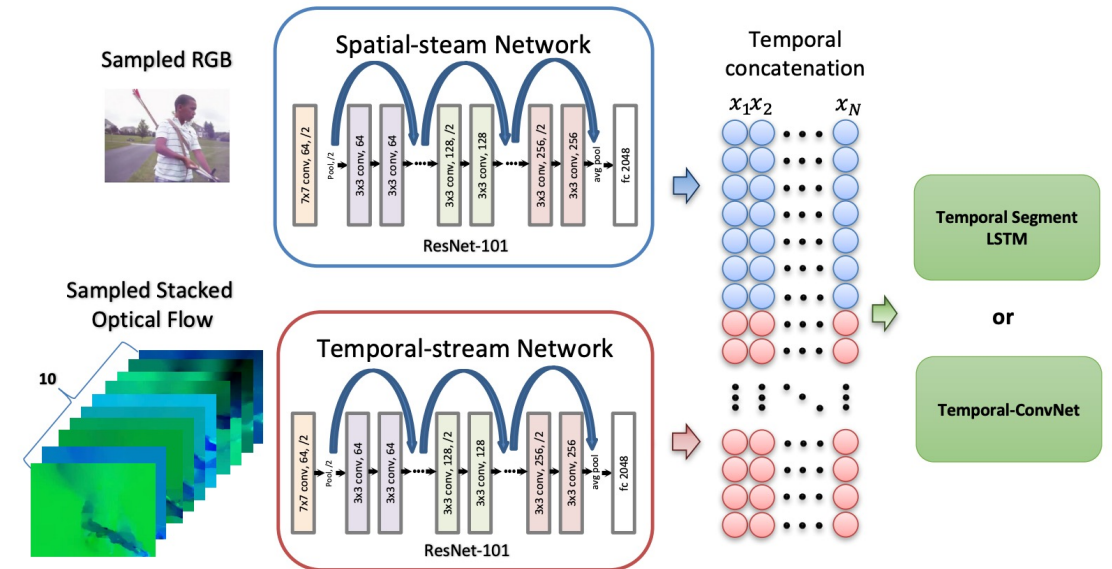
Temporal processing of data

Temporal processing \neq Linear spatial processing



Early temporal fusion: Encode both spatial and temporal information together and fuse them within the network

Late temporal fusion: Encode all spatial data in a time-wise fashion and determine temporal relationships



Future Direction 2

Sensor processing architectures



Vision data processing was revolutionized by CNNs

Language data processing was revolutionized by Transformers

LIDAR data processing is revolutionized by ?

RADAR data processing is revolutionized by ?

...

Future Direction 3

More data with less sensors!

4 Fisheye cameras provide a 360 degree surround view of the car

Results from Zero-shot (i.e. using the trained model out of the box) Segment Anything Model on Woodscape dataset



Important context and objects are not segmented



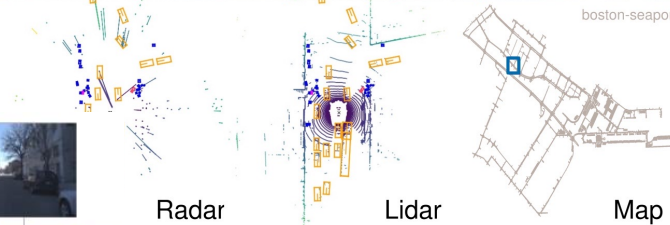
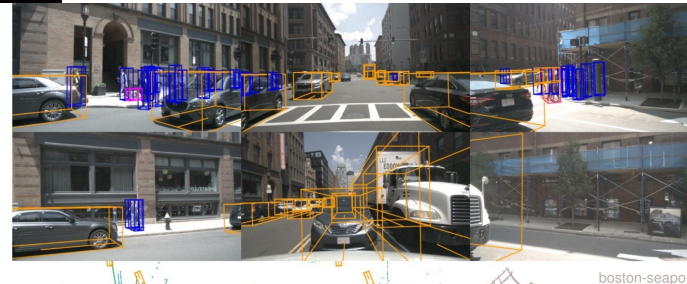
Future Direction 4

Infrastructure + AV Datasets

Abundance of egocentric AV datasets! Dearth of Infrastructure + AV datasets



Argoverse



Radar Lidar Map
"icycle, car makes a u-turn, lane change, peds crossing crosswalk"

NuScenes

- Infrastructure datasets: Stationary sensors at traffic junctures, streets, heavy pedestrian traffic areas etc.
- Infrastructure + AV datasets: Egocentric sensors on vehicles + stationary sensors for the same scenes

Some Memes to Wrap it Up

