



# Spatio-Temporal Data and Foundation Model

Jingtian Ma





#### • Spatio-Temporal Data in Traffic Domain







#### • Visual-Textual Data in Traffic Domain





#### **Captured Image Data**

天安门*(Tian'anmen)*,坐落在中华人民共和国首都北京市的中心、故宫的南端,与天安门广场以及人民英雄纪念碑、毛主 席纪念堂、人民大会堂、中国国家博物馆隔长安街相望,占地面积4800平方米,以杰出的建筑艺术和特殊的政治地位为世人所瞩 目。

天安门是明清两代北京皇城的正门,始建于明朝永乐十五年*(1417年)*,最初名"承天门",寓"承天启运、受命于天"之意。 设计者为明代御用建筑匠师蒯祥。清朝顺治八年*(1651年)*更名为天安门。由城台和城楼两部分组成,有汉白玉石的须弥座,总 高34.7米。天安门城楼长66米、宽37米。城台下有券门五阙,中间的券门最大,位于北京市皇城中轴线上,过去只有皇帝才可以 由此出入。正中门洞上方悬挂着毛泽东画像,两边分别是"中华人民共和国万岁"和"世界人民大团结万岁"的大幅标语。<sup>[1]</sup>

1925年10月10日,故宫博物院成立,天安门开始对民众开放。1949年10月1日,在这里举行了中华人民共和国开国大典,由 此被设计入国徽,并成为中华人民共和国的象征。1961年,中华人民共和国国务院公布为第一批全国重点文物保护单位之一。 [2]





#### **Text Description Data**



#### • Foundation Models

A Foundation Model is "any model that is trained on broad data (generally using self-supervision at scale) that can be adapted (e.g., fine-tuned) to a wide range of downstream tasks", such as BERT [Devlin et al. 2019], GPT-3 [Brown et al. 2020], and CLIP [Radford et al. 2021].

-[On the Opportunities and Risks of Foundation Models 2021 Li Fei-Fei etc.]







#### • Large Language Models (LLMs)



[A Survey of Large Language Models, 2023]



#### • Large Language Models (LLMs)



Feature	BERT	GPT	Т5
Model Type	Bidirectional encoder	Unidirectional decoder	Encoder-decoder
			structure
Task Domain	Mainly used for downstream tasks such as sentiment analysis, named entity recognition, etc.	Mainly used for generative tasks such as text generation, dialogue systems, etc.	Universal sequence-to- sequence tasks, supports various tasks like translation, summarization, etc.
Drotroining Tool	Masked Language	Autoregressive	Sequence-to-Sequence
Pretraining Task	Model (MLM)	Language Model	(seq2seq)
Context Understanding	Bidirectional encoding, understands the entire context information	Unidirectional decoding, understanding of the current token relies on previous tokens	Encoder-decoder structure, capable of processing global context information
Advantages	Strong context understanding, applicable to a variety of tasks	Generates coherent, context-aware text	Universal sequence-to- sequence structure, supports multiple tasks, easy to fine-tune
Disadvantages	High computational resource requirements, sensitive to sequence length	Limited understanding of context, inability to capture global information	Sensitive to sequence length, high computational resource requirements

[Attention Is All You Need, 2017]



#### • Large Vision Models (LVMs)

Vision Transformer (ViT): Addressing tasks in the field of computer vision using the standard Transformer.







### Large Vision Models (LVMs)

- ViT Linear Projection of Flattened Patches
  - □ Split Image into Patches: patch 16 x16
  - $\Box$  Vectorization: If the patches are d1xd2xd3 tensors, then the vectors are d1d2d3×1



[An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale, 2021]



#### • Large Vision Models (LVMs)

- ViT Classification Task
  - □ Add Position Encoding
  - **C**onnect to Transformer



[An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale, 2021]



### Large Vision Models (LVMs)

- ViT Classification Task
  - □ Pretrain the model on Dataset A, fine-tune on Dataset B, then evaluate on Dataset B
  - □ Pretrained on ImageNet (small), ViT is slightly worse than ResNet.
  - □ Pretrained on ImageNet-21K (medium), ViT is comparable to ResNet.
  - □ Pretrained on JFT (large), ViT is slightly better than ResNet.

		200 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100 - 100			
	# of Images	<pre># of Classes</pre>	ResNet is better	ViT is bette	r
ImageNet (Small)	1.3 Million	1 Thousand			# of Images
ImageNet-21K (Medium)	14 Million	21 Thousand			pretraining
JFT (Big)	300 Million	18 Thousand	100M	l l Images 300M I	mages

[An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale, 2021]



#### • Vision Language Models (VLMs)

- Single stream (VisualBERT, VL-BERT, Oscar) Early fusion
- Dual stream (ViLBERT, LXMERT, CLIP) Late fusion





#### • Vision Language Models (VLMs)



Contrastive Language-Image Pre-Training (CLIP)

- Rich pretraining data (400M image-text pairs)
- Mapping images and text to a shared embedding space
- Capability for zero-shot image classification and image-text retrieval

[Learning Transferable Visual Models From Natural Language Supervision, 2021]



#### • Vision Language Models (VLMs)

	ImageNet	
Dataset	ResNet101	CLIP VIT-L
ImageNet	76.2%	76.2%
ImageNet V2	64.3%	70.1%
ImageNet Rendition	37.7%	88.9%
ObjectNet	32.6%	72.3%
ImageNet Sketch	25.2%	60.2%
ImageNet Adversarial		77.1%

[Learning Transferable Visual Models From Natural Language Supervision, 2021]

# How to Combine Spatio-Temporal Data and Foundation Models?





#### • Use Spatio-Temporal Data as Text or Image Formats



**Agent-Driver** 

[A Language Agent for Autonomous Driving, 2023]





**Use Spatio-Temporal Data as Text or Image Formats** 

[A Language Agent for Autonomous Driving, 2023]



#### Use Spatio-Temporal Data as Text or Image Formats



÷. \*\*\*\*\*Environmental information:\*\*\*\*\* Front object detections: Front object detected, object type: pedestrian, object id: 0, position: (-4.32, 13.85), size: (0.76, 0.91) Front object detected, object type: pedestrian, object id: 9, position: (2.23, 19.39), size: (0.57, 0.69) Future trajectories for specific objects: Object type: pedestrian, object id: 0, future waypoint coordinates in 3s: [(-4.44, 14.50), ...] Object type: pedestrian, object id: 9, future waypoint coordinates in 3s: [(1.74, 19.76), ...] Map information (road shoulders): Current ego-vehicle's distance to left shoulder is 0.5m and right shoulder is 4.5m \*\*\*\*\*Common sense:\*\*\*\*\* - Maintain a safe distance from the objects in front of you... \*\*\*\*\*Past driving experience for reference:\*\*\*\*\* Most similar driving experience from memory with similarity score: 0.77: Scenario information: .. The planned trajectory in this scenario for your reference: [(0.04,2.49), ..., (0.30,12.26)] Ô \*\*\*\*\*Chain-of-thoughts reasoning:\*\*\*\*\* - Notable objects: pedestrian at (0.80,18.81), moving to (-2.53,20.89) at 3.0 second - Potential effects: may collide if continue driving at this speed. \*\*\*\*\*Task planning:\*\*\*\*\* Behavior: forward; Speed: deceleration; Driving plan: move forward with a deceleration \*\*\*\*\*Motion planning:\*\*\*\*\* Trajectory: [(-0.03, 2.47), (-0.10, 4.84), (-0.19, 7.10), (-0.29, 9.25), (-0.39, 11.29), (-0.49, 13.22)] \*\*\*\*\*Self-reflection:\*\*\*\*\* \*\*\*\*\*Self-reflection:\*\*\*\*\* No collision. No change to the motion planning result No collision. No change to the motion planning result.





[A Language Agent for Autonomous Driving, 2023]



#### • Use Spatio-Temporal Data as Text or Image Formats





#### • Use Spatio-Temporal Data as Text or Image Formats



□ Limited receptive fields

[VectorNet: Encoding HD Maps and Agent Dynamics from Vectorized Representation, 2020]



#### • Use Spatio-Temporal Data as Text or Image Formats



[Multimodal Motion Prediction with Stacked Transformers, 2021]



### • Employ Specific Spatio-Temporal Models as Tools



Tools



#### Employ Specific Spatio-Temporal Models as Tools



[TrafficGPT: Viewing, Processing and Interacting with Traffic Foundation Models, 2023]



#### • Pre-train a Foundation Model with Spatio-Temporal Data



#### Input: OSM data

Objectives:

- 1) Text-based objective
- 2 Vision-Language
  - multimodal objective
- ③ Road-based objective.



#### • Pre-train a Foundation Model with Spatio-Temporal Data



- Abstract the problem of multi-task decision-making for autonomous driving as a sequence modeling and prediction task.
- Transform the "state-action-reward" tuples from expert data into a token format similar to NLP task to match the input format of the GPT model

Input: 
$$au' = \left(s_1, a_1, g_1, s_2, a_2, g_2, ..., s_T, a_T, g_T\right)$$

- s: Sequence of state
- a: Sequence of action
- g: Sequence of function of trajectory

[MTD-GPT: A Multi-Task Decision-Making GPT Model for Autonomous Driving at Unsignalized Intersections, 2023]

### Others



#### • Prompt Learning





[Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing, 2021]

## Others



### • Parameter-Efficient Fine-Tuning (PEFT)



[LLM-Adapters: An Adapter Family for Parameter-Efficient Fine-Tuning of Large Language Models, 2023]

# **Future Directions**



- Traffic Scene Understanding
- Traffic Description Generation
- Traffic Simulation

. . .

• Trajectory/Road Network Generation





# Thanks for Listening

Jingtian Ma