

# Towards Learning a Generalist Model for Embodied Navigation

Duo Zheng<sup>1,2</sup>, Shijia Huang<sup>1</sup>, Lin Zhao<sup>3</sup>, Yiwu Zhong<sup>1</sup>, Liwei Wang<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong    <sup>2</sup>Shanghai AI Laboratory

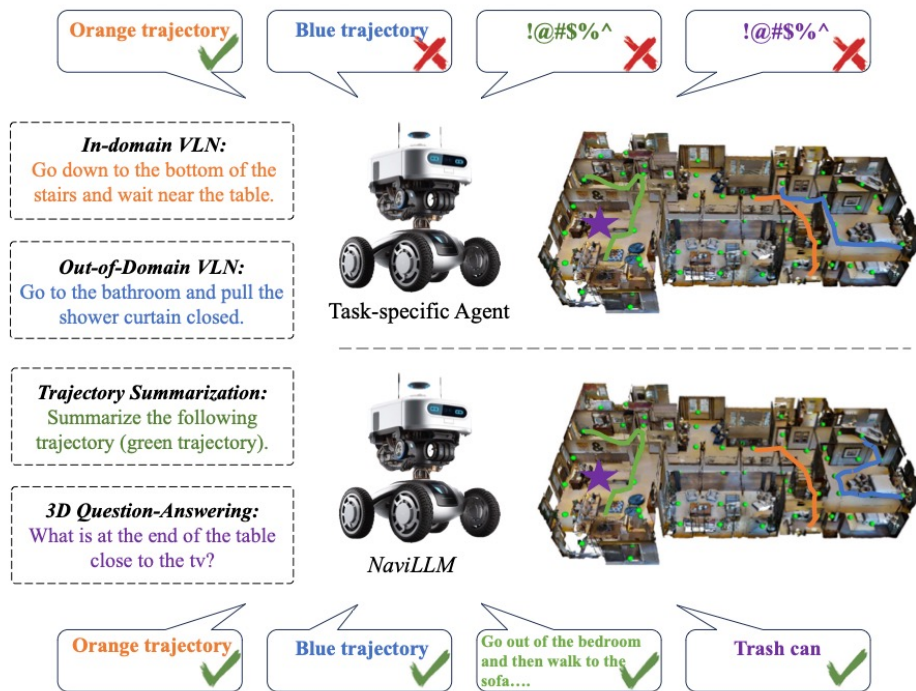
<sup>3</sup>Centre for Perceptual and Interactive Intelligence



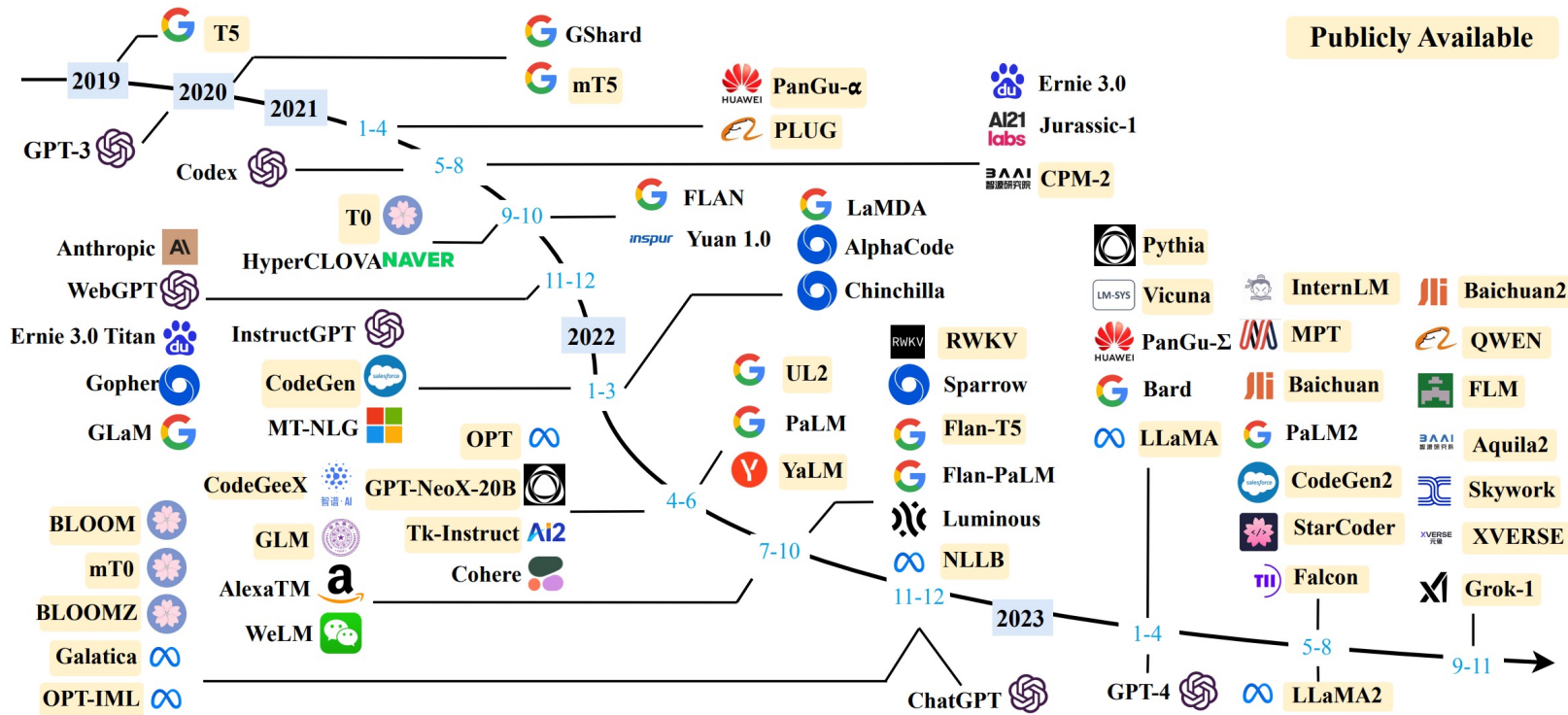


- Background
- Method
- Experimental Results
- Conclusion

- An agent located in 3D environment is required to navigate according to various forms of instructions and provide textual responses based on user queries.

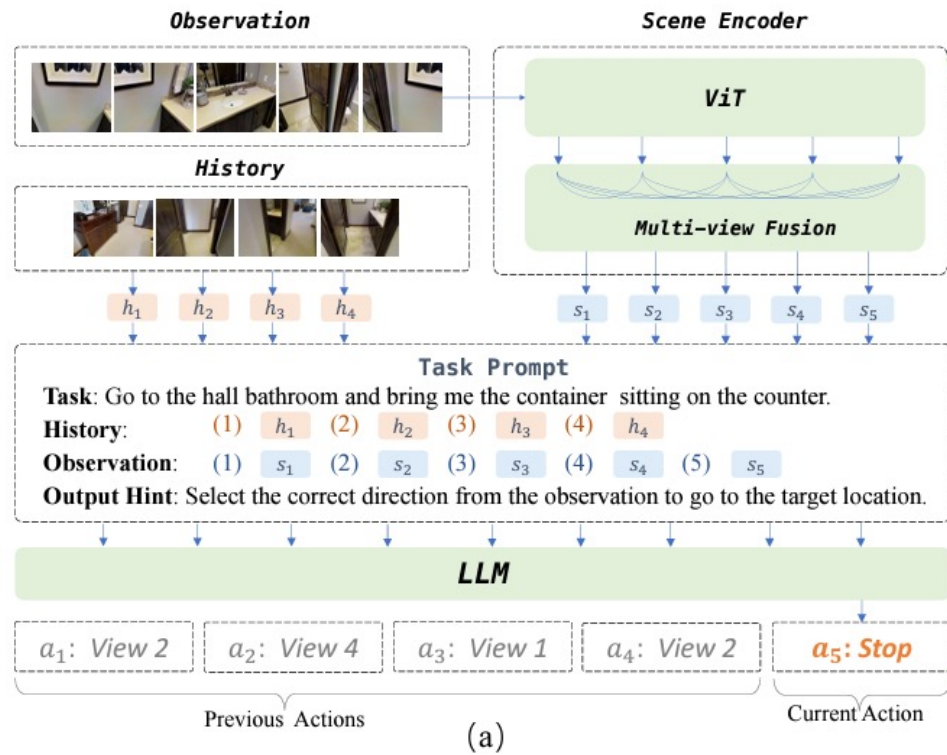


# The Era of Large Language Models



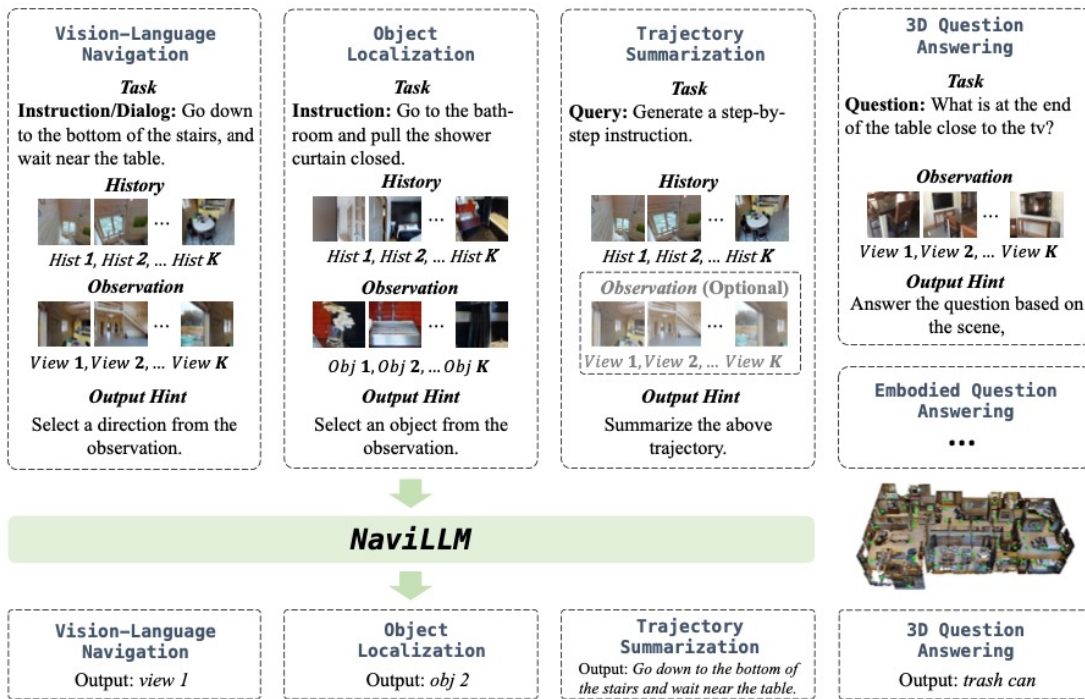
Reference: A Survey of Large Language Models, arXiv 2023

- NaviLLM is an embodied model grounded in LLM, comprising two modules, i.e., a scene encoder and an LLM.





- The schema is designed to be a unified format that can adapt to different data sources and enable flexibility for wide span of tasks.





	CVDN		SOON		R2R		REVERIE		ScanQA	
	Val-U	Test	Val-U	Test	Val-U	Test	Val-U	Test	Val	Test
<i>Separate Model For Each Task</i>										
PREVALENT [25]	3.15	2.44	-	-	53	51	-	-	-	-
HOP [44]	4.41	3.24	-	-	57	59	26.11	24.34	-	-
HAMT [11]	5.13	5.58	-	-	61	<u>60</u>	30.20	26.67	-	-
VLN-BERT [27]	-	-	-	-	57	57	24.90	23.99	-	-
GBE [61]	-	-	13.34	9.23	-	-	-	-	-	-
DUET [12]	-	-	22.58	<u>21.42</u>	60	58	33.73	<u>36.06</u>	-	-
Meta-Explore [31]	-	-	-	<u>25.80</u>	62	<b>61</b>	34.03	-	-	-
AZHP [23]	-	-	-	-	61	<u>60</u>	<b>36.63</b>	35.85	-	-
VLN-SIG [32]	5.52	5.83	-	-	<u>62</u>	<u>60</u>	-	-	-	-
VLN-PETL [45]	<u>5.69</u>	<u>6.13</u>	-	-	60	58	27.67	26.73	-	-
BEV-BERT [1]	-	-	-	-	<b>64</b>	<u>60</u>	<u>36.37</u>	<b>36.41</b>	-	-
3D-LLM [28]	-	-	-	-	-	-	-	-	<u>20.5</u>	<u>19.1</u>
<i>Unified Model For All Tasks</i>										
MT-RCM+Env [54]	4.65	3.91	-	-	49	40	-	-	-	-
<i>NaviLLM</i>	<b>6.16</b>	<b>7.90</b>	<b>29.24</b>	<b>26.26</b>	59	<u>60</u>	35.68	32.33	<b>23.0</b>	<b>26.3</b>

- ❖ Our method delivers SoTA results with a single model.
- ❖ Significant improvement on CVDN can be credited to the utilization of LLM and multi-task training.
- ❖ Better Performance on tasks with complex instructions.



	CVDN		SOON				REVERIE			
	TL	GP $\uparrow$	TL	OSR $\uparrow$	SR $\uparrow$	SPL $\uparrow$	TL	OSR $\uparrow$	SR $\uparrow$	SPL $\uparrow$
DUET (R2R)	21.12	3.38	26.83	7.64	4.66	2.84	7.88	29.11	24.91	20.00
DUET (REVERIE)	76.13	3.30	33.72	20.86	10.24	6.06	-	-	-	-
DUET (SOON)	48.61	2.40	-	-	-	-	38.10	43.45	10.91	3.64
<i>NaviLLM</i>	26.37	<b>4.46</b>	28.66	<b>33.11</b>	<b>19.81</b>	<b>14.29</b>	18.96	<b>51.47</b>	<b>28.10</b>	<b>21.04</b>

Task: Navigate to the object in ‘what color is the stove?’.

Answer the question.



1. The agent is at bathroom  
Walk out of the bathroom

2. Find the kitchen  
Walk to the kitchen



3. Locate the stove  
Answer the question

Output: White

- ❖ Our method generalizes to out-of-domain VLN tasks.
- ❖ *NaviLLM* can combine the learned navigation and question-answering ability to solve more complex tasks, e.g., MP3D-EQA.





- We propose the first generalist model for embodied navigation, enabling a wide spectrum of capabilities required for embodied navigation.
- We unify various tasks in a single model by adapting LLM and introducing schema-based instruction.
- Our single model achieves SoTA results on CVDN, SOON, and ScanQA, with a significant margin of 29% compared to the previous SoTA on CVDN. Furthermore, it also exhibits strong generalizability on unseen tasks.