

# NaVid: Video-based VLM Plans the Next Step for Vision-and-Language Navigation

Jiazhao Zhang<sup>\*</sup>, Kunyu Wang<sup>\*</sup>, Rongtao Xu<sup>\*</sup>, Gengze Zhou, Yicong Hong,  
Xiaomeng Fang, Qi Wu, Zhizheng Zhang<sup>†</sup>, He Wang<sup>†</sup>



**BAAI**  
智源研究院

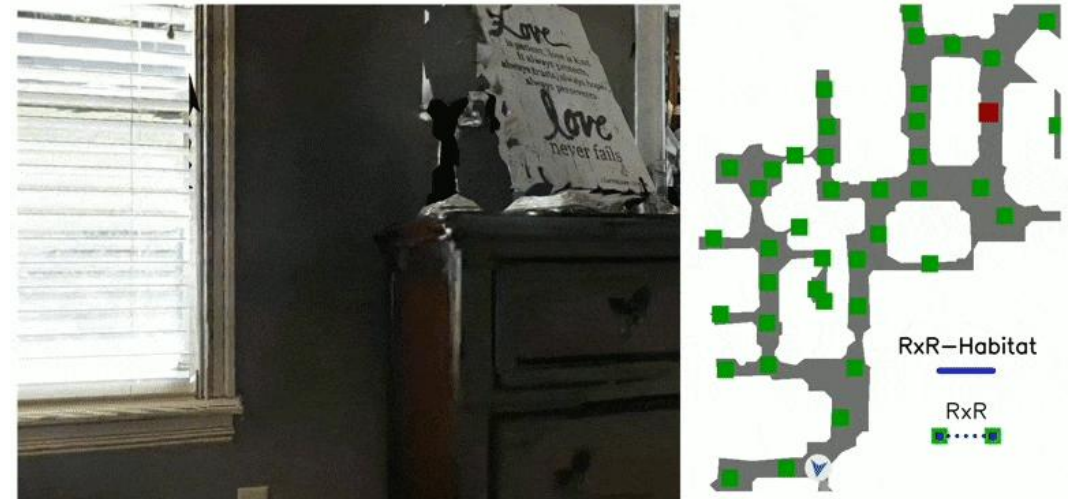
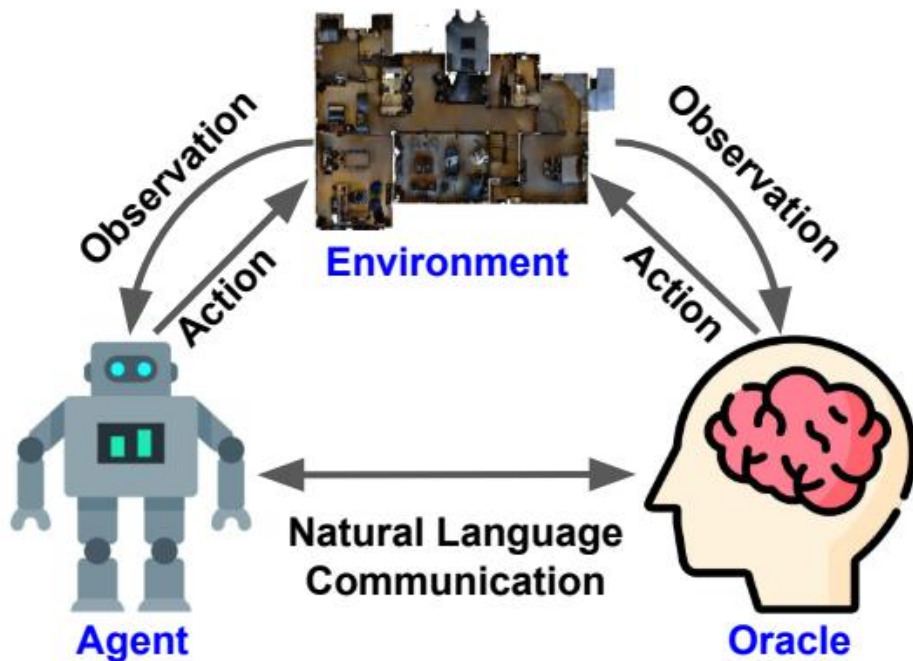


**GALBOT**

# Vision-and-Language Navigation (VLN)

Given **free-form instruction**, the robot is required to follow the instruction to navigate in the **unseen environments**.

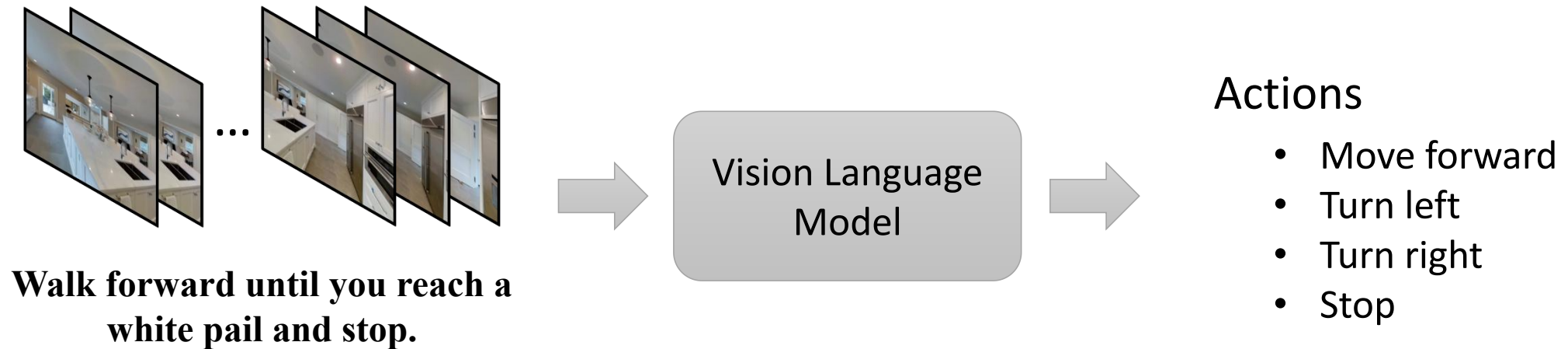
"Leave the bedroom, and enter the kitchen. Walk forward and take a left at the couch. Stop in front of the window"



You are in a bedroom. Turn around to the left until you see a door leading out into a hallway, go through it. Hang a right and walk between the island and the couch on your left. When you are between the second and third chairs for the island stop.

# Key Insight 1 – VLM-Driven Real-World VLN

- Leverage the power of **foundational VLMs** to extend VLN to **real-world** applications, using **pretrained** large model, and **co-tuning** with **web-based data**.



510k Navigation Data + 763k Web-based Data → Total 1.2M training data

# Key Insight 2 – Video-Based VLN Agent

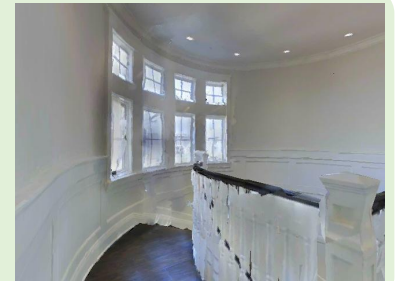
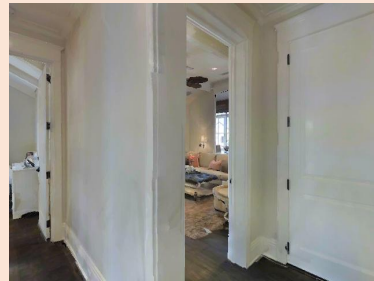
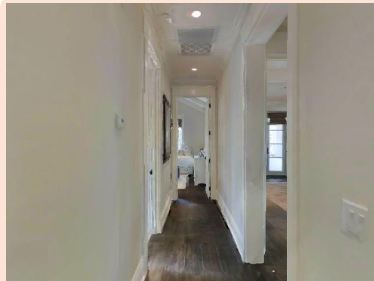
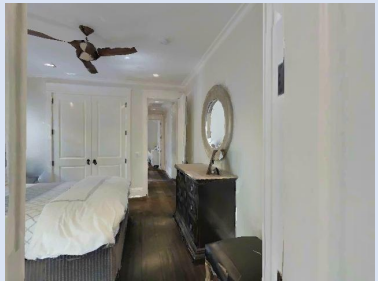
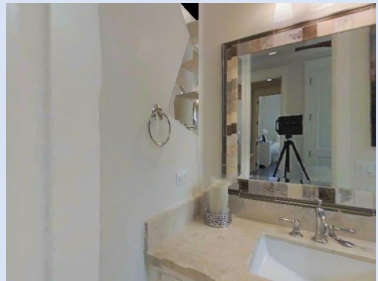
- Navigate in a **human-like** manner, relying **solely** on real-time **video streams** from a monocular camera, **without the need for maps, odometers, or depth inputs.**

*Walk out of the bedroom, turn right, stop before the stairs.*

Walk out of the bedroom

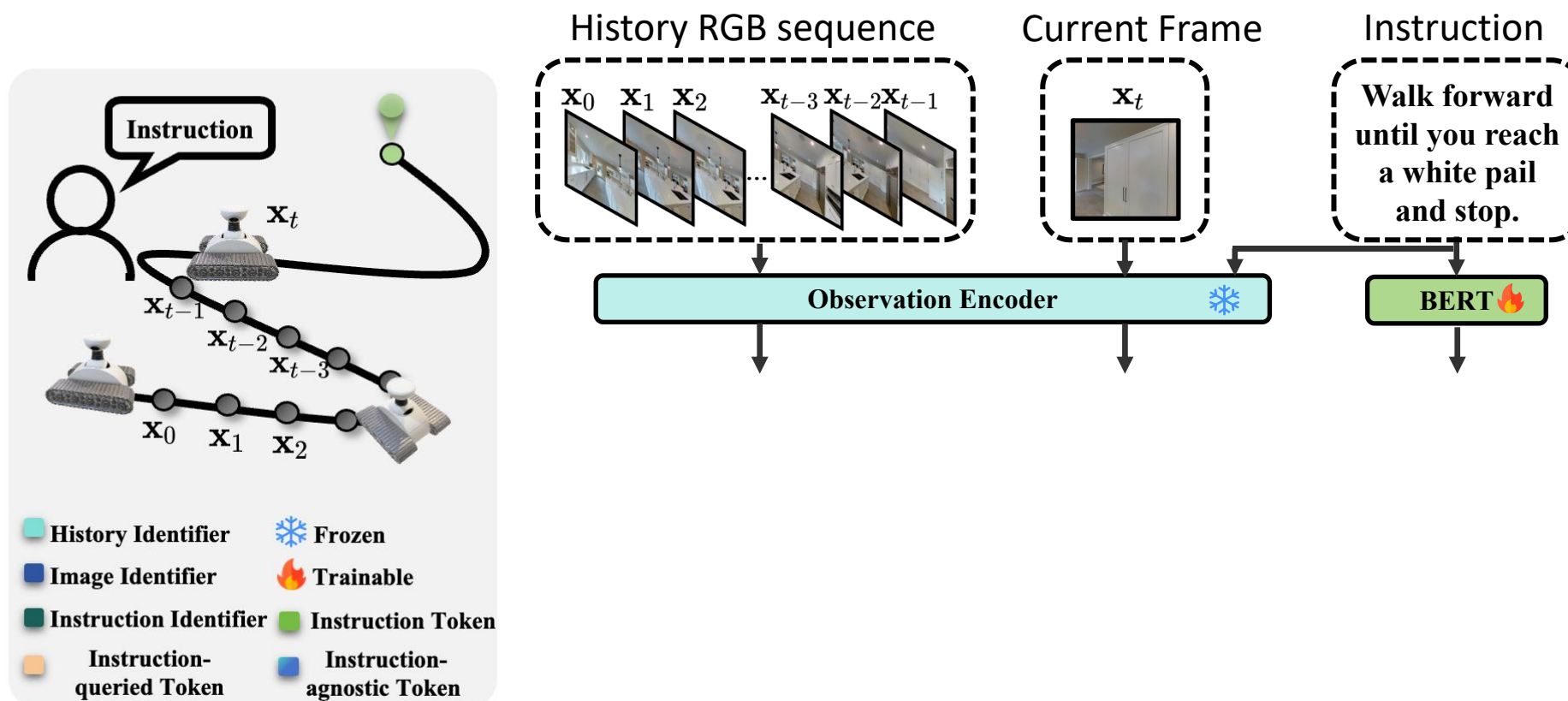
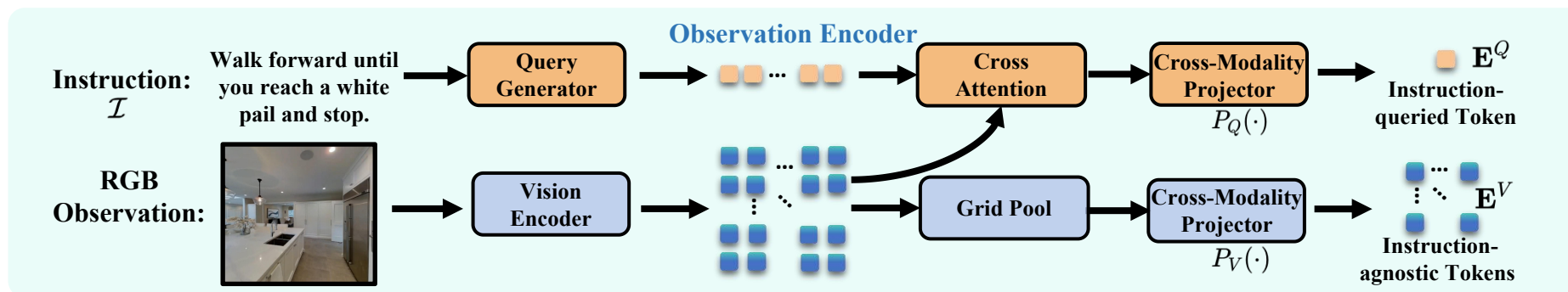
Turn right

Stop before the stairs



On-the-fly Video as Input

# Pipeline





Text:

## What is large language model?

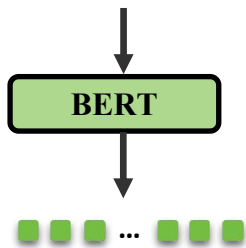
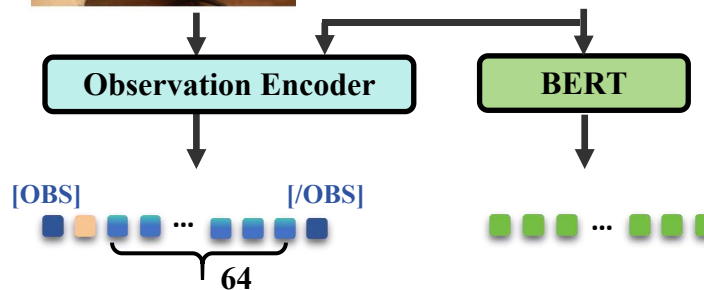


Image:



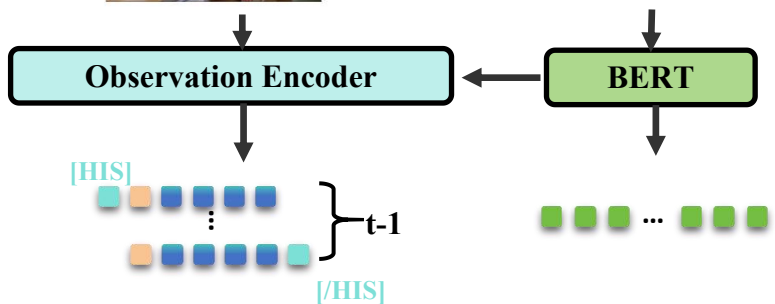
**Suppose you are a detective, what can you infer from the visual clues in the image?**



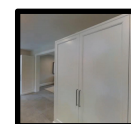
## Video:



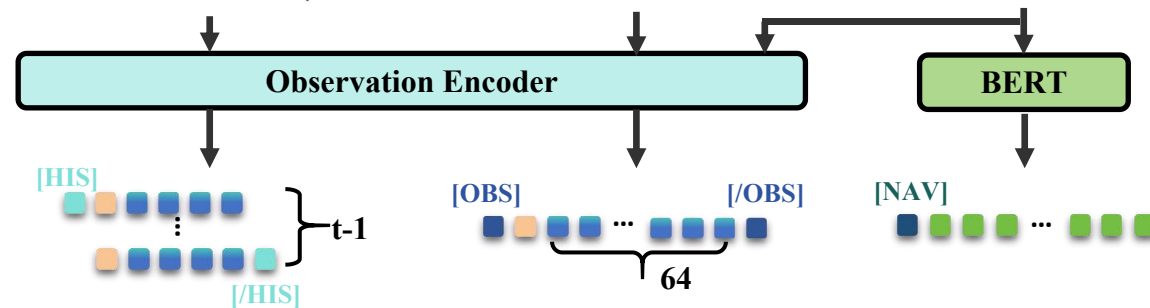
**Please describe  
this video in  
detail.**



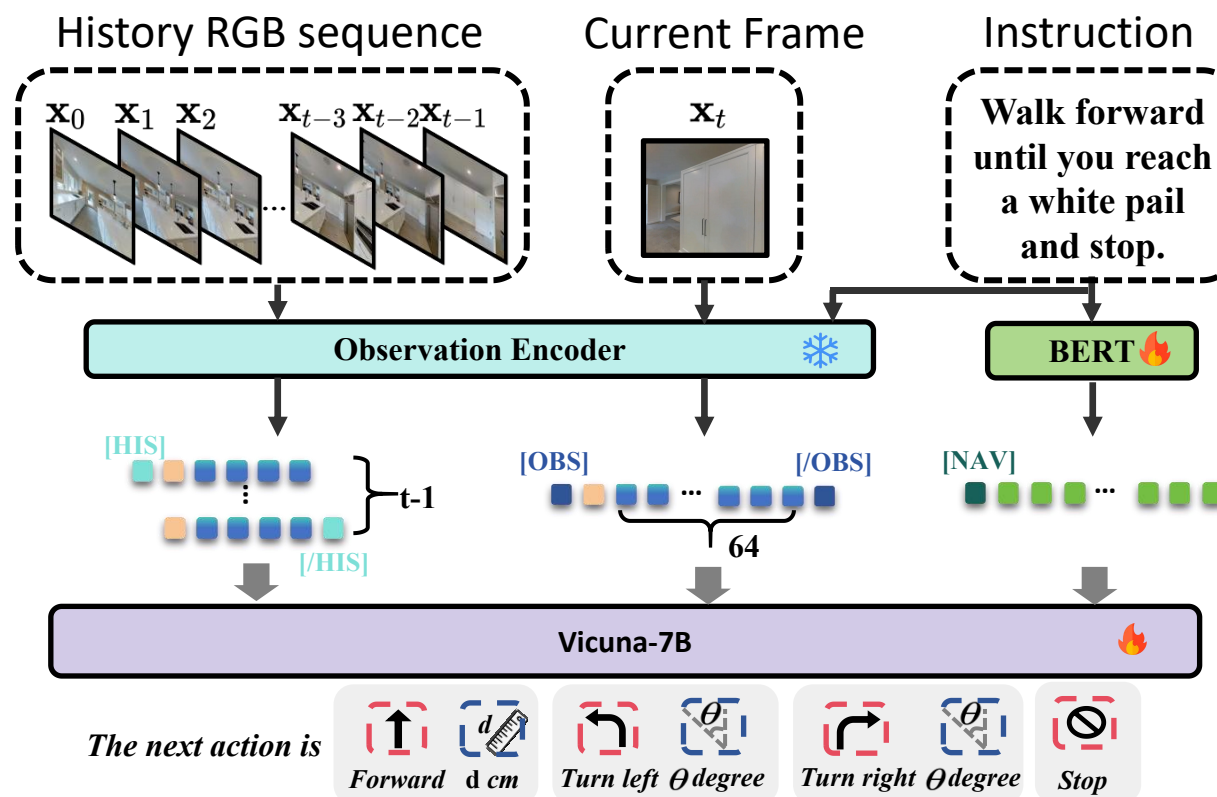
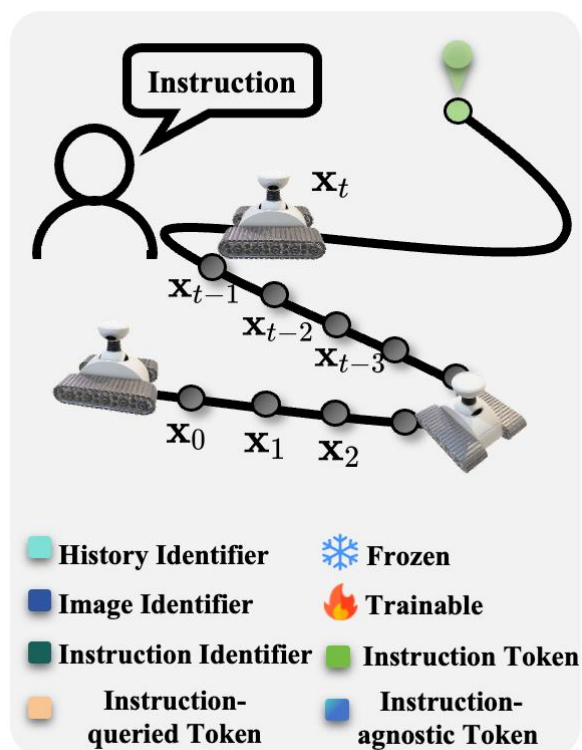
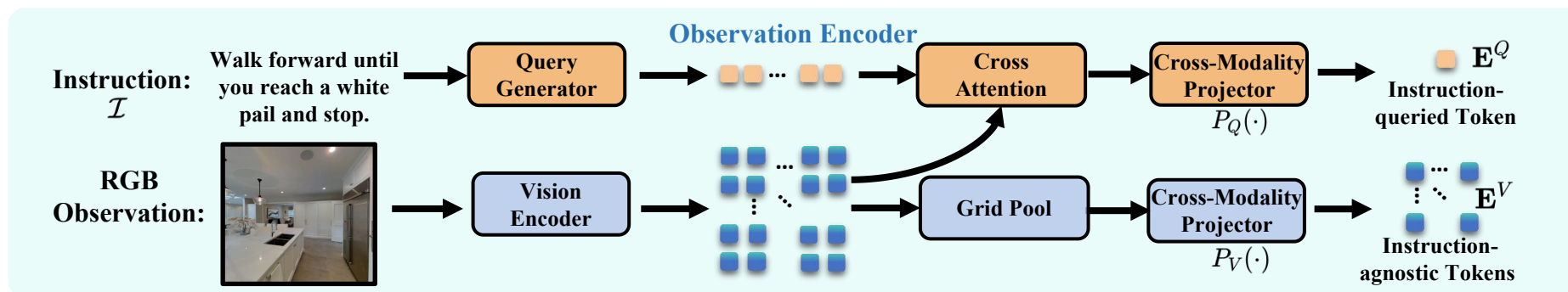
## Navigation:



**Walk forward  
until you reach  
a white pail  
and stop.**



# Pipeline




# Data collection


We collect the navigation data based on **R2R dataset training-split** on VLN-CE simulator:  
**10819 episodes, 61 scenes (MP3D).**

-- Action Planning Sample --

**History Video**



**Current Obs.**




**User:** Imagine you are a robot programmed for navigation tasks. You have been given a video of historical observations and a image of current observation <image>. Your assigned task is: Walk forward into and down the middle of the workspace. Walk forward until you reach a white pail on the floor next to a desk with a black chair and stop. Analyze this series of images to decide your next move, which could involve turning left or right by a specific degree or moving forward a certain distance.

**Assistant:** The next action is move forward 75 cm.

Sample video segment + action  
(**Action Planning Sample**)

-- Instruction reason sampling --

**Trajectory Video**



**User:** Assume you are a robot designed for navigation. You are provided with captured images sequences <image>. Based on this image sequence, please describe the navigation trajectory of the robot.

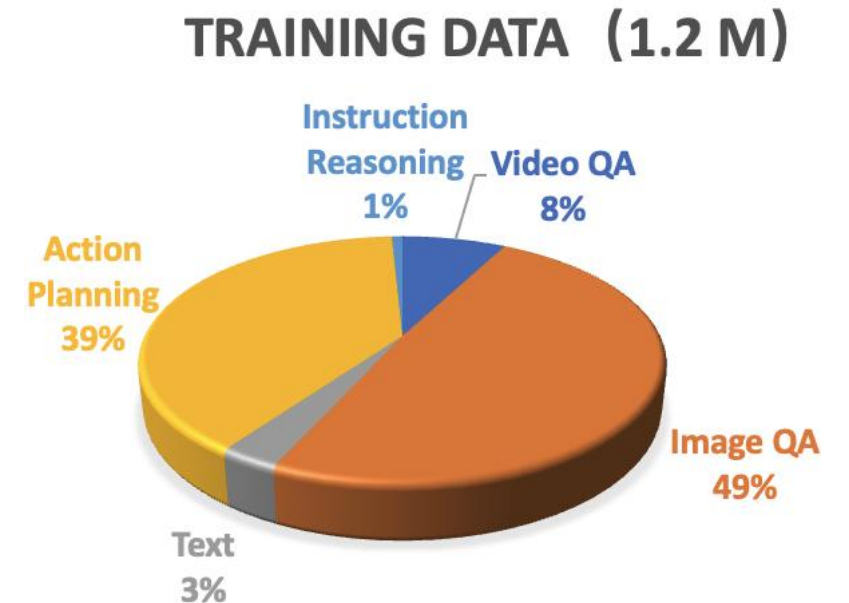
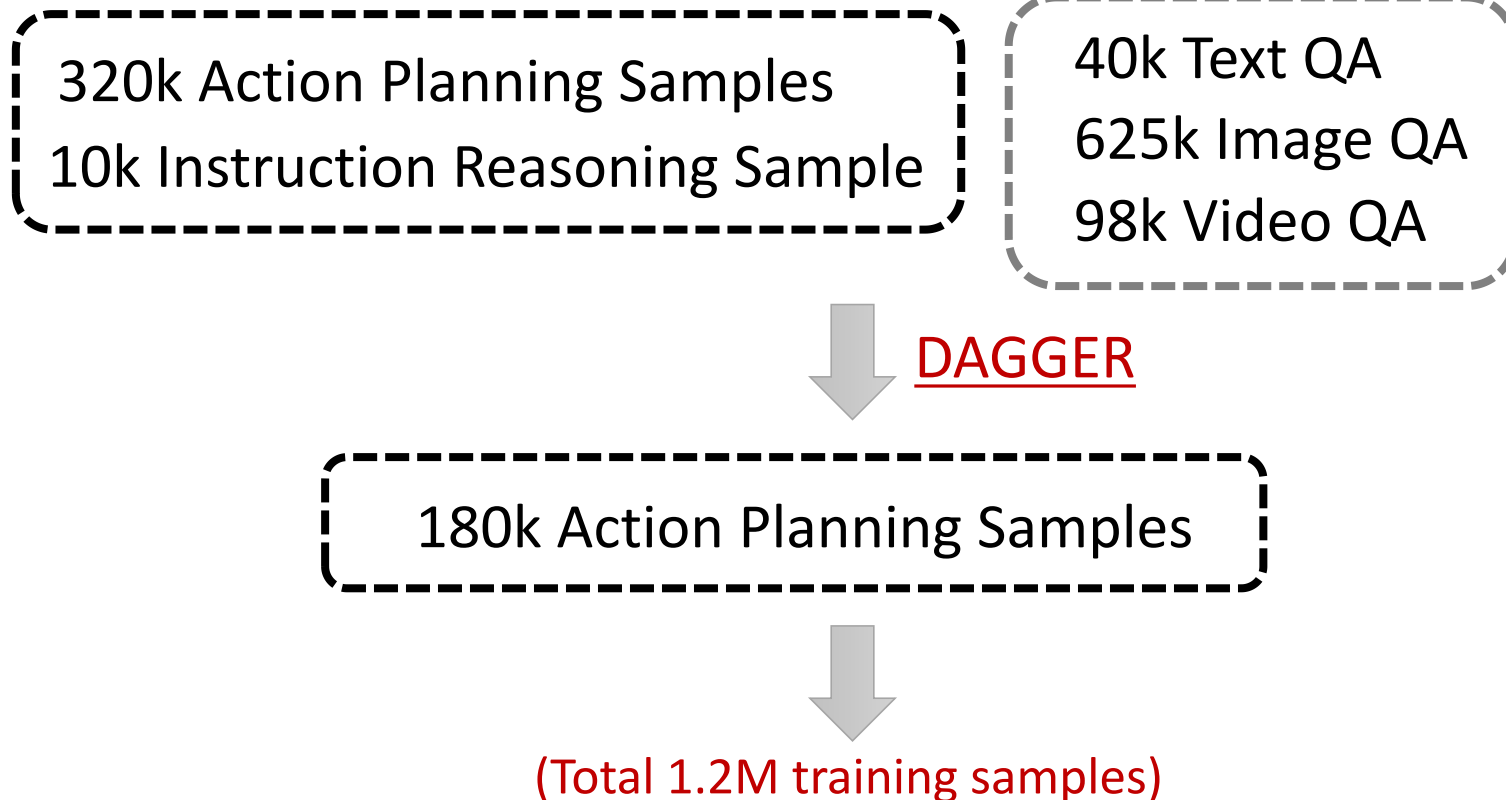
**Assistant:** Walk forward into and down the middle of the workspace. Walk forward until you reach a white pail on the floor next to a desk with a black chair and stop.

Video + instruction  
(**Instruction Reasoning Sample**)



# Data collection

We collect the navigation data based on **R2R dataset training-split** on VLN-CE simulator:  
**10819 episodes, 61 scenes (MP3D).**



# R2R train -> R2R val-unseen (cross split)

	Observation				VLN-CE R2R Val-Unseen				
	Pan.	S.RGB	Depth	Odo.	TL	NE↓	OS↑	SR↑	SPL↑
AG-CMTP [15]	✓		✓	✓	–	7.90	39.2	23.1	19.1
R2R-CMTP [15]	✓		✓	✓	–	7.90	38.0	26.4	22.7
LAW [73]		✓	✓	✓	8.89	6.83	44.0	35.0	31.0
CM2 [29]		✓	✓	✓	11.54	7.02	41.5	34.3	27.6
WS-MGMap [16]		✓	✓	✓	10.00	6.28	47.6	<b>38.9</b>	34.3
Seq2Seq [43]		✓	✓		9.30	7.77	37.0	25.0	22.0
CMA [43]		✓	✓		8.64	7.37	40.0	32.0	30.0
RGB-Seq2Seq		✓			4.86	10.1	8.10	0.00	0.00
RGB-CMA		✓			6.28	9.55	10.8	5.00	4.43
<b>Ours</b>		✓			7.63	<b>5.47</b>	<b>49.1</b>	37.4	<b>35.9</b>

↑ SR (success rate)

↑ OS (oracle success rate)

↑ SPL (success weighted by path length)

↓ NE (Navigation error)

**SOTA level performance with  
only RGB video inputs**

# R2R train -> RxR val-unseen (cross dataset)

	Observation			VLN-CE RxR Val-Unseen				
	S.RGB	Depth	Odo.	TL	NE↓	OS↑	SR↑	SPL↑
LAW [73]	✓	✓	✓	4.01	10.87	21.0	8.0	8.0
CM2 [29]	✓	✓	✓	12.29	8.98	25.3	14.4	9.2
WS-MGMap [16]	✓	✓	✓	10.80	9.83	29.8	15.0	12.1
Seq2Seq [43]	✓	✓		1.16	11.8	5.02	3.51	3.43
CMA [43]	✓	✓		5.09	11.7	10.7	4.41	2.47
RGB-Seq2Seq	✓			4.43	11.2	12.2	0.0	0.0
RGB-CMA	✓			13.56	9.55	14.8	0.0	0.0
A <sup>2</sup> Nav [17]	✓			—	—	—	16.8	6.3
<b>Ours</b>	✓			10.59	<b>8.41</b>	<b>34.5</b>	<b>23.8</b>	<b>21.2</b>

↑ SR (success rate)

↑ OS (oracle success rate)

↑ SPL (success weighted by path length)

↓ NE (Navigation error)

**Our method consistently demonstrates SOTA performance, significantly surpassing baseline metrics.**

# R2R train -> Real world (Sim-to-real)

	Meeting Room				Office				Lab				Lounge			
	Simple I.F.		Complex I.F.		Simple I.F.		Complex I.F.		Simple I.F.		Complex I.F.		Simple I.F.		Complex I.F.	
	SR↑	NE↓	SR↑	NE	SR↑	NE↓	SR↑	NE↓	SR↑	NE↓	SR↓	NE↓	SR↑	NE↓	SR↑	NE↓
Seq2Seq [42]	4%	4.45	0%	7.21	0%	4.28	0%	6.92	0%	4.58	0%	6.61	0%	5.95	0%	6.82
CMA [42]	0%	4.27	0%	7.30	8%	4.62	0%	5.71	4%	4.35	0%	5.67	0%	4.63	0%	5.46
WS-MGMap [13]	52%	1.18	24%	2.20	60%	0.96	20%	2.94	44%	1.85	12%	3.18	48%	1.66	32%	2.88
<b>Ours</b>	<b>92%</b>	<b>0.55</b>	<b>56%</b>	<b>0.98</b>	<b>84%</b>	<b>0.63</b>	<b>48%</b>	<b>0.71</b>	<b>76%</b>	<b>0.83</b>	<b>40%</b>	<b>1.89</b>	<b>88%</b>	<b>0.72</b>	<b>44%</b>	<b>1.37</b>

Example:

Go straight to the wall, then turn left and walk to the door, then stop.

(1) Real Environment I

Third-Person View

Videos

A

B

C

D

E

A

B

C

D

E

Steps



# Real-world experiment – Simple Instruction

Simple Instruction following

Speed up x10

Walk towards the **door** then stop.



Walk towards the **white box** then stop.





# Real-world experiment - Outdoor Scenes

Simple Instruction following

Walk forward to the chair then **turn right**, and stop at the stairs.



Speed up x10

Walk forward to the chair then **turn left** and stop at the stairs.

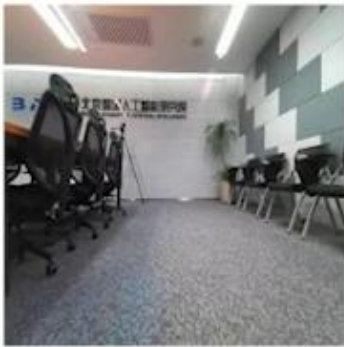


# Real-world experiment – Complex Instruction

Complex Instruction following

Speed up x10

Go straight and move close to the plant, then turn right facing the door, then walk to the door and stop.



# Thanks for your attention

Project page: <https://pku-epic.github.io/NaVid/>



**BAAI**  
智源研究院



**GALBOT**

Contact me at [zhngjizh@gmail.com](mailto:zhngjizh@gmail.com), [hewang@pku.edu.cn](mailto:hewang@pku.edu.cn).