

---

## **CS586 (25 Spring) : Paper Presentation**

---

# **NaVid: Video-based VLM Plans the Next Step for Vision-and-Language Navigation (RSS 24)**

**20244076 Zhaoyan Wang**

**20244050 Xiangchen Liu**



# Review

## Efficient Residual Learning with Mixture-of-Experts for Universal Dexterous Grasping (ICLR 2025)

### Improvements

- **Residual Policy Learning Framework**
  - Efficiently trainable model
- **Geometry-Agnostic Base Policy**
  - High generalization
- **Mixture-of-Experts (MoE)**
  - High diversity and good performance

### Limitation

- **No functional grasping**
- **No experiment on hardware**

### Conclusion

- **Perform better than previous works**
- **zero generalization gap**

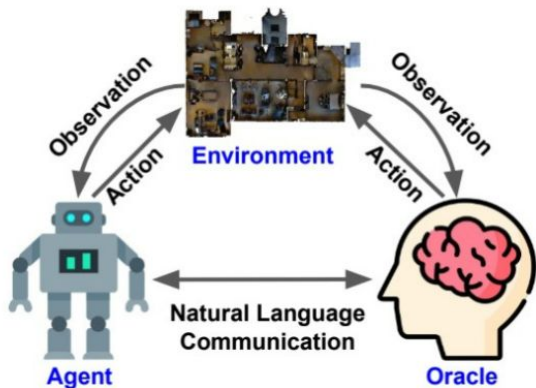
# Table of contents

- Introduction & Motivation
- Related works
- Problem Formulation
- Methodology: THE PROPOSED NAVID AGENT
- Data collection
- Experiments
- Conclusion & Limitations

# Introduction & Motivation

## Vision-and-Language Navigation (VLN)

Goal: Humans communicate with each other using natural language to issue tasks and request help.



[Gu et. al. 2022]



Amazon Astro Robot

A robot that can understand human language and navigate intelligently would significantly benefit human society.

[1] Gu, Jing et al. "Vision-and-Language Navigation: A Survey of Tasks, Methods, and Future Directions." .arXiv preprint arXiv:2203.12667 (2022).

[2] <https://www.cnet.com/home/smart-home/amazon-astro-review>

# Introduction & Motivation

## Vision-and-Language Navigation (VLN)

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

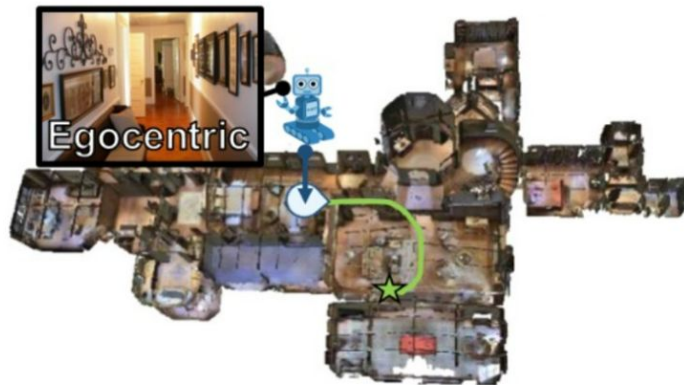
"Leave the bedroom, and enter the kitchen. Walk forward and take a left at the couch. Stop in front of the window"

Observation:

- Egocentric color map
- Egocentric depth map
- Location and orientation

Action:

- Low-level actions  
(Move forward, Turn left, Turn right, Stop)



VLN-CE

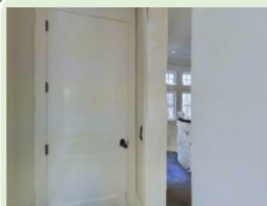
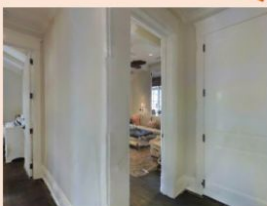
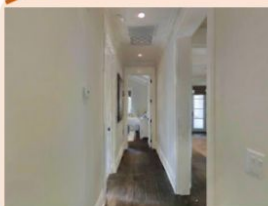
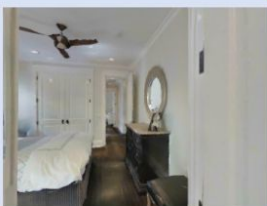
# Introduction & Motivation

## Challenges

High-level understanding:

- Understand long-horizon trajectory with rich visual information.
- Understand free-from text instruction.
- Align the instruction with history trajectory.

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



# Introduction & Motivation

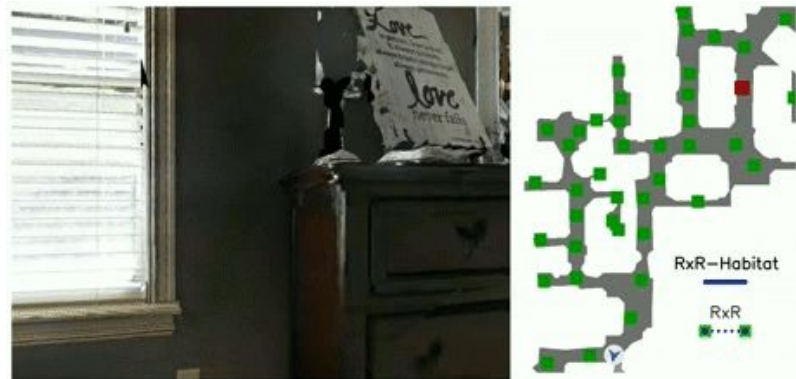
## Challenges

### High-level understanding:

- Understand long-horizon trajectory with rich visual information.
- Understand free-from text instruction.
- Align the instruction with history trajectory.

### Low-level planning:

- Approaching landmarks
- Obstacle avoidance

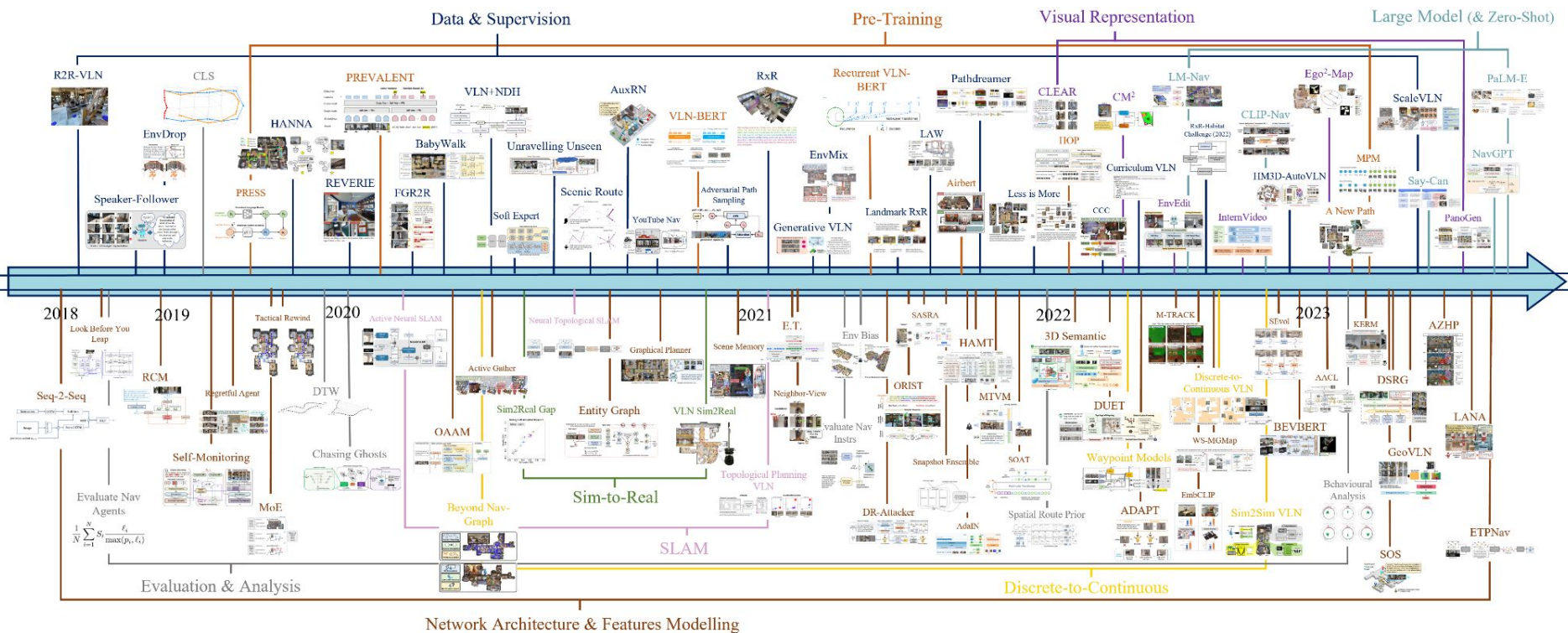


You are in a bedroom. Turn around to the left until you see a door leading out into a hallway, go through it. Hang a right and walk between the island and the couch on your left. When you are between the second and third chairs for the island stop.

VLN-CE RxR



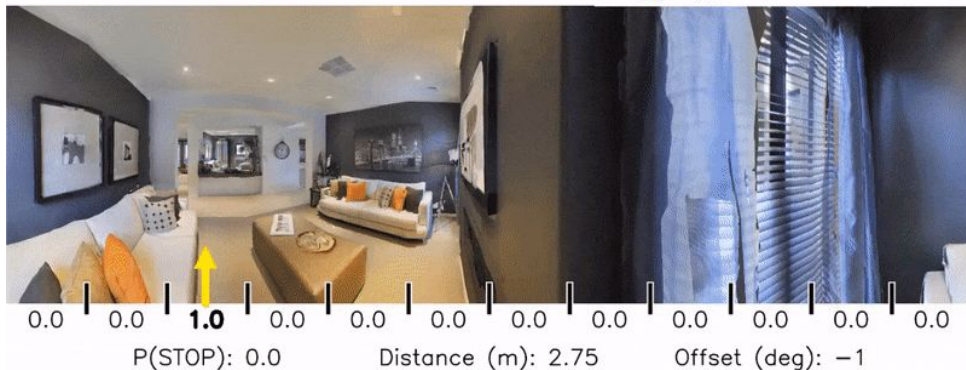
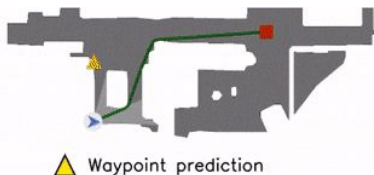
# Introduction & Motivation





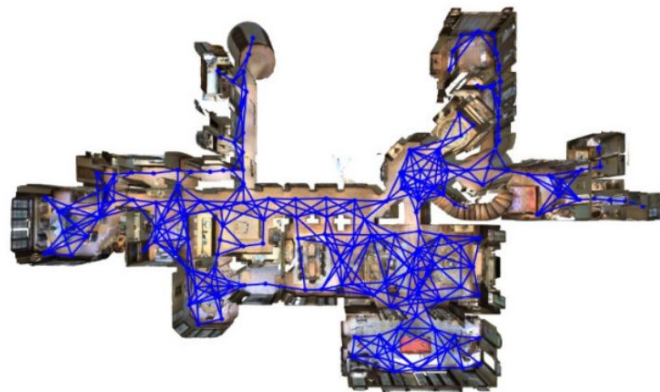
# Related works

Walk around the brown leather ottoman, angling slightly towards the clock on the wall. Turn right at the clock and walk forward. Wait near the dining table.



(A) Panoramic images

- Large field of view



(B) Topology node graph

- Known Topology
- Oracle Navigation
- Perfect Location

[1] Krantz, Jacob et al. "Beyond the Nav-Graph: Vision-and-Language Navigation in Continuous Environments.", ECCV 2020

[2] Krantz, Jacob et al. "Waypoint models for instruction-guided navigation in continuous environments.", ICCV 2021

# Related works

## Simplified Settings:

1. Panorama Observation
2. Discrete Connectivity Graph
3. Jump-Point motion

## Problem:

1. Super Large Sim2Real Gap
2. No discrete graph in physical world
3. Pose Estimation is imperfect.



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

[1] Krantz, Jacob et al. "Beyond the Nav-Graph: Vision-and-Language Navigation in Continuous Environments.", ECCV 2020

[2] Krantz, Jacob et al. "Waypoint models for instruction-guided navigation in continuous environments.", ICCV 2021

# Related works

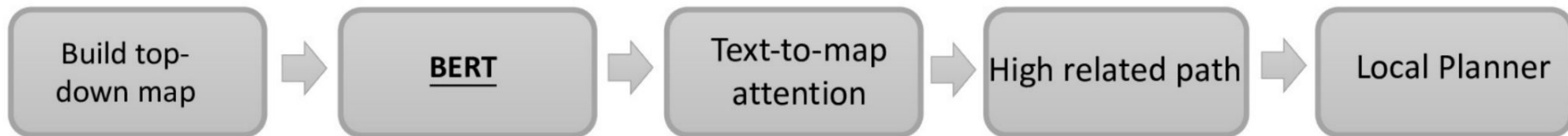
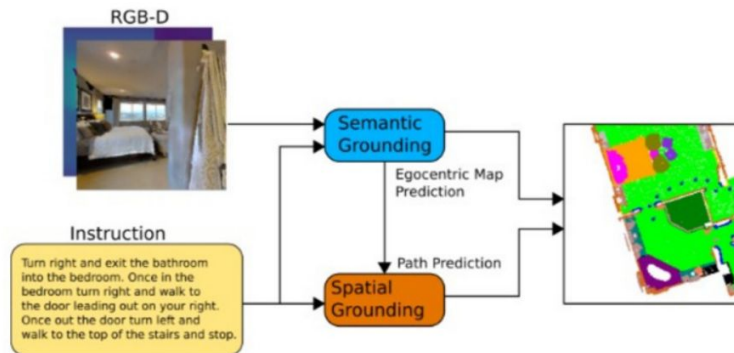
## Map-based VLN method

### Cross-modal Map Learning for Vision and Language Navigation

Georgios Georgakis, Karl Schmeckpeper, Karan Wanchoo, Soham Dan,  
Eleni Miltsakaki, Dan Roth, Kostas Daniilidis  
University of Pennsylvania

{ggeorgak, karls, kwanchoo, sohamdan, elenimi, danroth, kostas}@seas.upenn.edu

Project webpage: <https://ggeorgak11.github.io/CM2-project/>



Noiseless depth image, rotation and translation estimation.

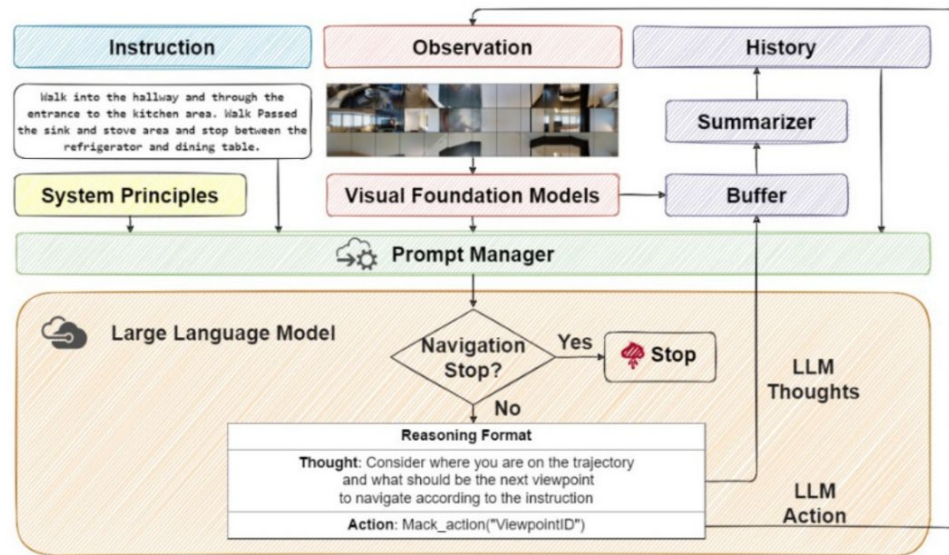
# Related works

## LLM-based VLN method

### NavGPT: Explicit Reasoning in Vision-and-Language Navigation with Large Language Models

Gengze Zhou<sup>1</sup> Yicong Hong<sup>2</sup> Qi Wu<sup>1</sup>

<sup>1</sup>The University of Adelaide <sup>2</sup>The Australian National University  
{gengze.zhou, qi.wu01}@adelaide.edu.au yicong.hong@anu.edu.au  
<https://github.com/GengzeZhou/NavGPT>



Foundation models  
+  
observation images

Text-based  
navigation history

ChatGPT

Node selection

# Related works

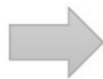
---

Is there a simple way to achieve VLN?

We want a straightforward solution for VLN.



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

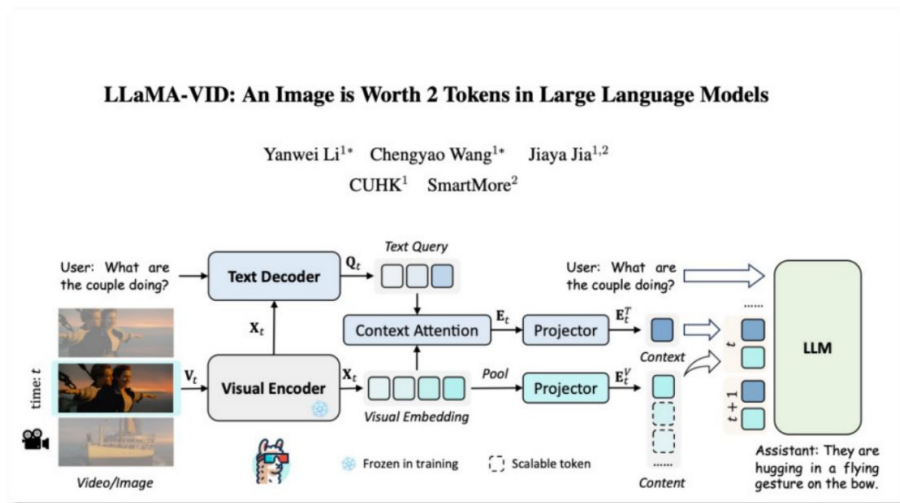


Actions

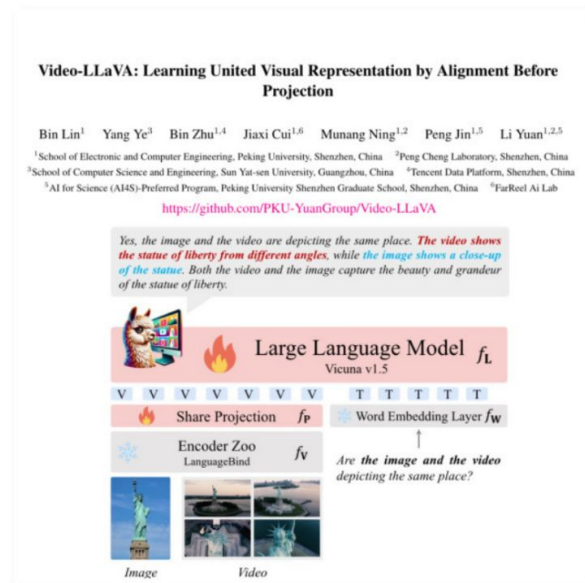
# Related works

## Video-Based Vision-Language Model (VLM)

- Strong performance in understanding video and text.
- Generalizability to novel videos and texts.



LLaMA-VID



VID-LLaVA

[1] Li, Yanwei et al. "LLaMA-VID: An Image is Worth 2 Tokens in Large Language Models." ArXiv abs/2311.17043

[2] Lin, Bin et al. "Video-LLaVA: Learning United Visual Representation by Alignment Before Projection." ArXiv abs/2311.10122



# Problem Formulation

We are focusing on a straightforward and challenging solution: using only **RGB video** as input and directly output **low-level actions** from a **video-based VLM**.



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



Vision Language  
Model



Actions

- Move forward
- Turn left
- Turn right
- Stop

- An intuitive way to drive the agent.
- Eliminates the need for location, orientation, and depth information.

# Problem Formulation

## From Video-based Question & Answer (VQA) to Navigation

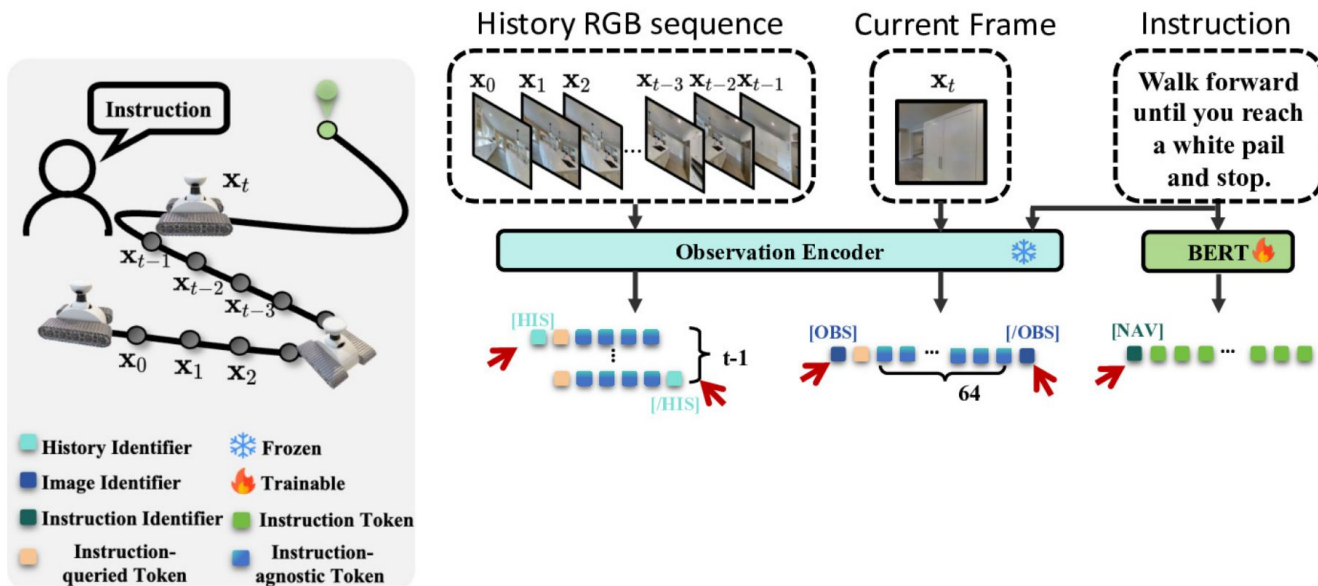
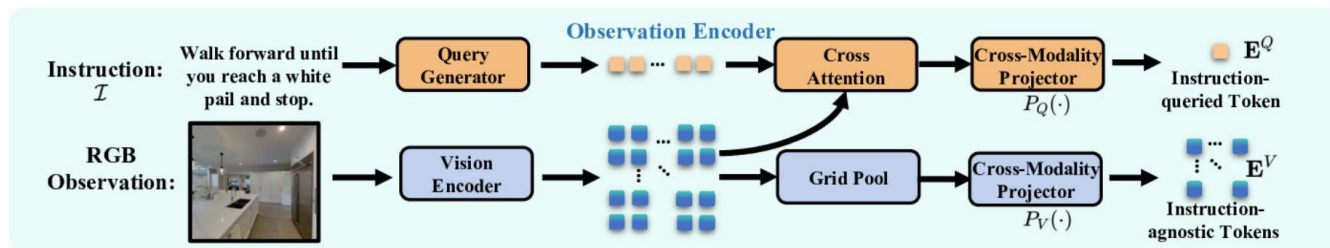


- The modality of VLN is different from the common modalities of large models.

Design a new pipeline of video-based VLM for VLN

- There is a lack of large amounts of high-quality real data for VLN task.

# Methodology: THE PROPOSED NAVID AGENT



# Methodology: THE PROPOSED NAVID AGENT

Text:

What is large language model?

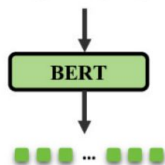
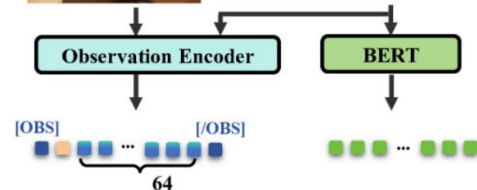


Image:



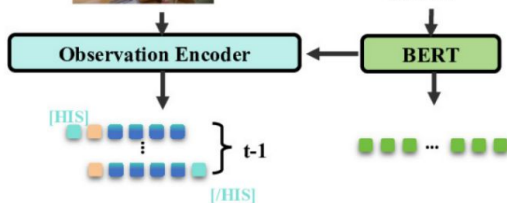
Suppose you are a detective, what can you infer from the visual clues in the image?



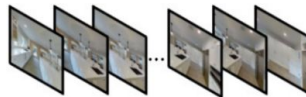
Video:



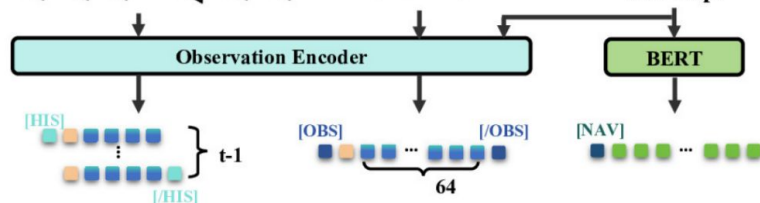
Please describe this video in detail.



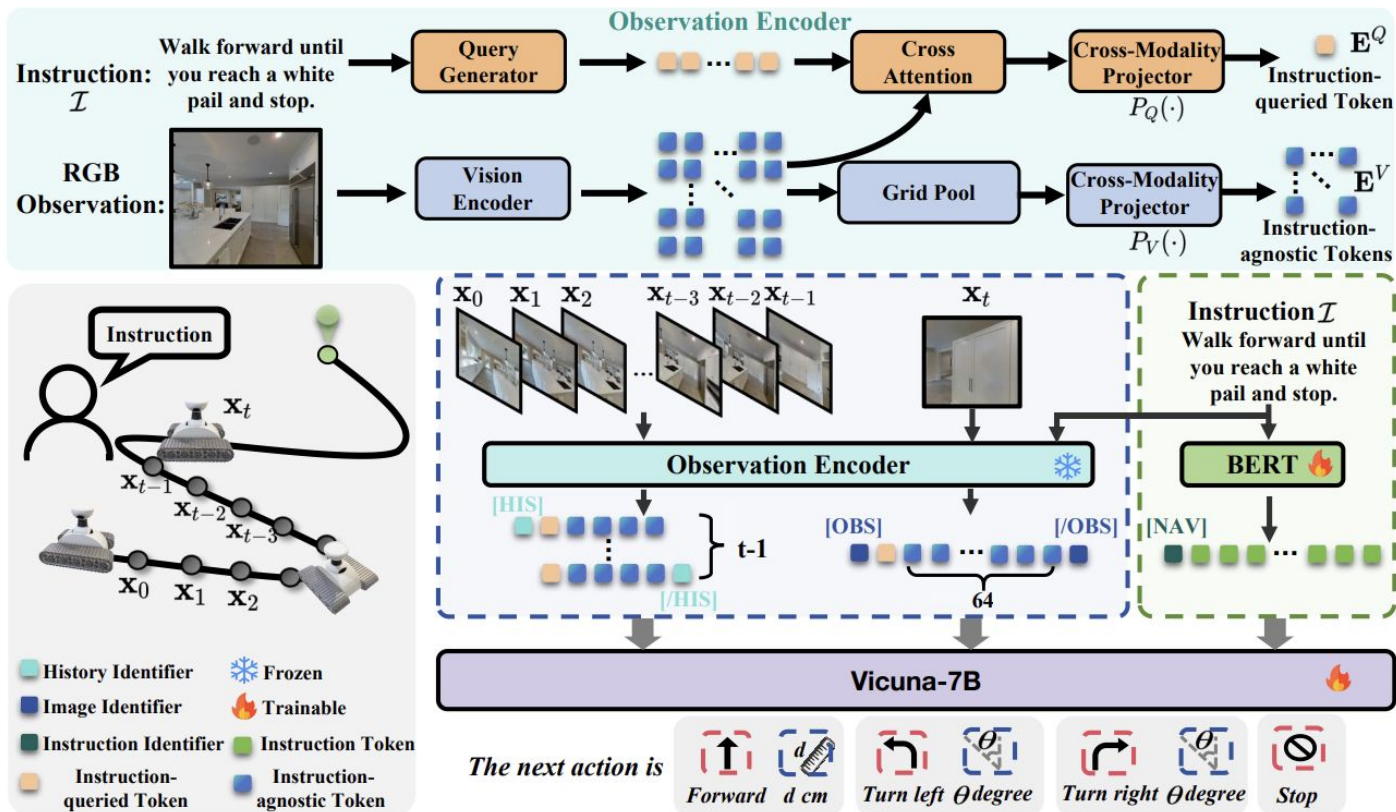
Navigation:



Walk forward until you reach a white pail and stop.



# Methodology: THE PROPOSED NAVID AGENT



# Methodology: THE PROPOSED NAVID AGENT

## From video QA to navigation



The modality of VLN is different from the common modalities of large models.

There is a lack of large amounts of high-quality real data for VLN task.


Collect simulator data for training NaVid




# 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**  


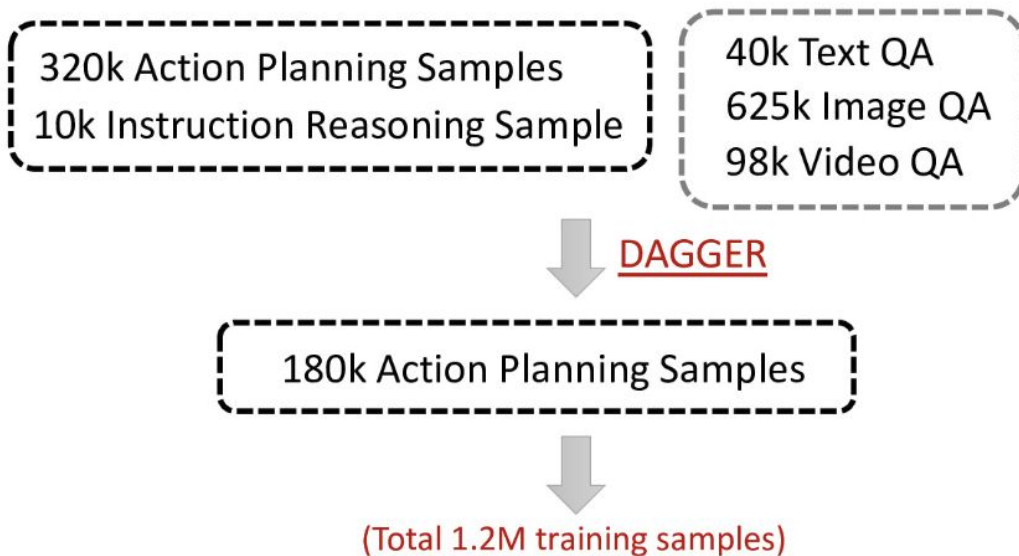
**User:** Assume you are a robot designed for navigation. You are provided with captured images sequences <image>. Based on this image sequence, please describe the navigation trajectory of the robot.

**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)

# Experiments

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



# Experiments

R2R train -> R2R val-unseen



Turn left and go through the kitchen. Turn left and walk past the kitchen island. Turn right and walk past the pottery. Wait inside the room on the left next to the table with the flowers.

The next action is turn left 45 degree.

R2R train -> RxR val-unseen



We start off looking at book shelves. If you slightly look down you see the rug on the floor. Take a step onto the rug on the floor. Now turn to your right. Walk through the open doors. Continue to walk down that pathway. Passing the table on your left, and going straight towards the piano. Now pass the piano keep on walking straight. Once you get to the end of the hallway turn left. Walk through the double door archway. Take one more step forward and turn to your right and you'll now see a bedroom. There is a chair with a circular ottoman. Um take a step towards that circular ottoman and stop there you are done.

The next action is turn right 45 degree.

# Experiments

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>		✓			<b>7.63</b>	<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

# Experiments

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)

Under the same setting, our method demonstrates significant improvements on all metrics.  
(648% improvement on SR and 710% improvement on SPL.)



# Experiments

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

	Observation			VLN-CE RxR Val-Unseen				
	S.RGB	Depth	Odo.	TL	NE↓	OS↑	SR↑	SPL↑
LAW [43]	✓	✓	✓	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)

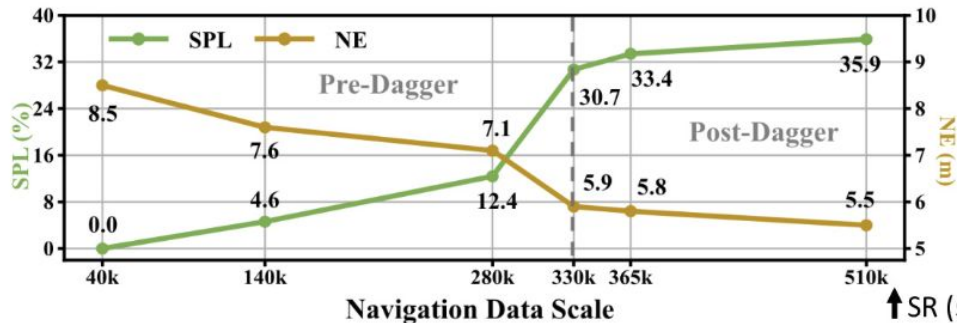
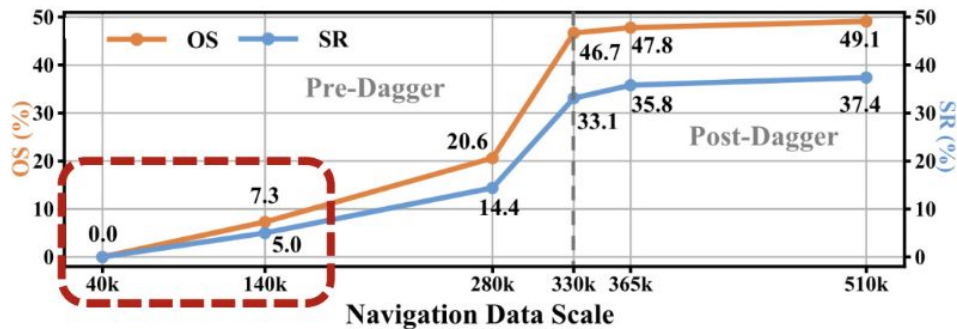
↓ NE (Navigation error)

Our method consistently demonstrates state-of-the-art (SOTA) performance, significantly surpassing baseline metrics.



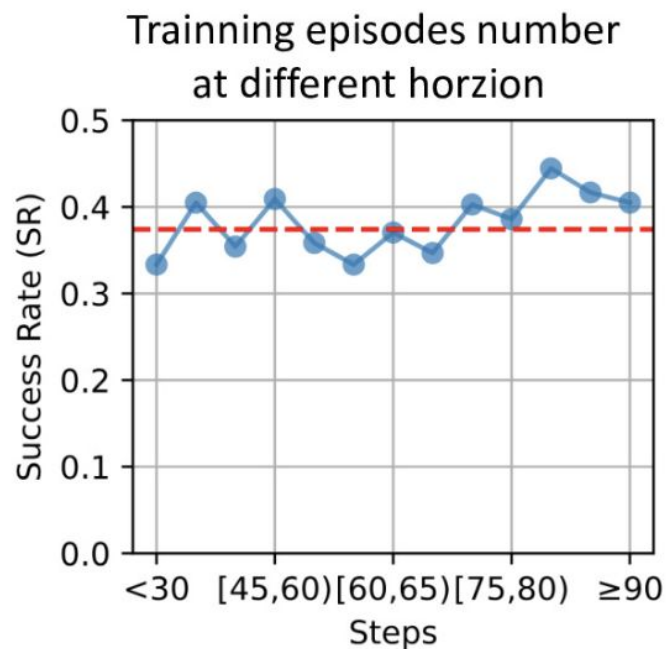
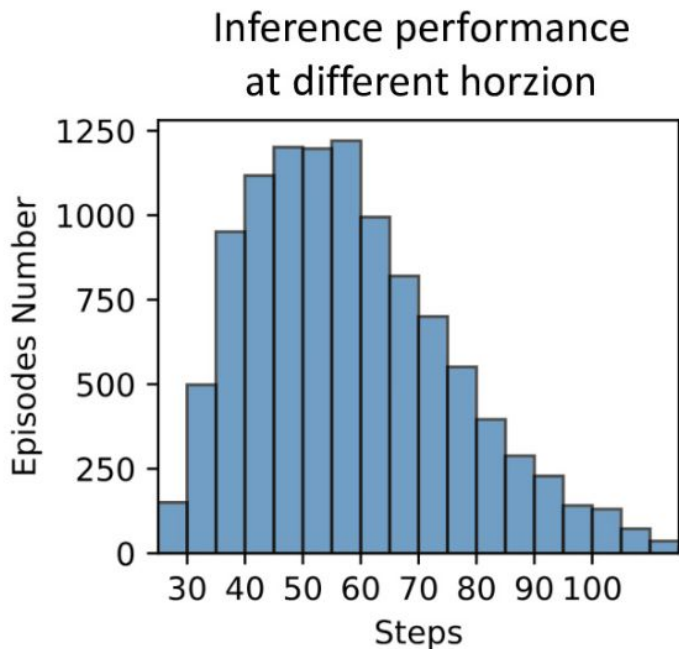
# Experiments

Performance on VLN-CE R2R Val-Unseen



- ↑ SR (success rate)
- ↑ OS (oracle success rate)
- ↑ SPL (success weighted by path length)
- ↓ NE (Navigation error)

# Experiments

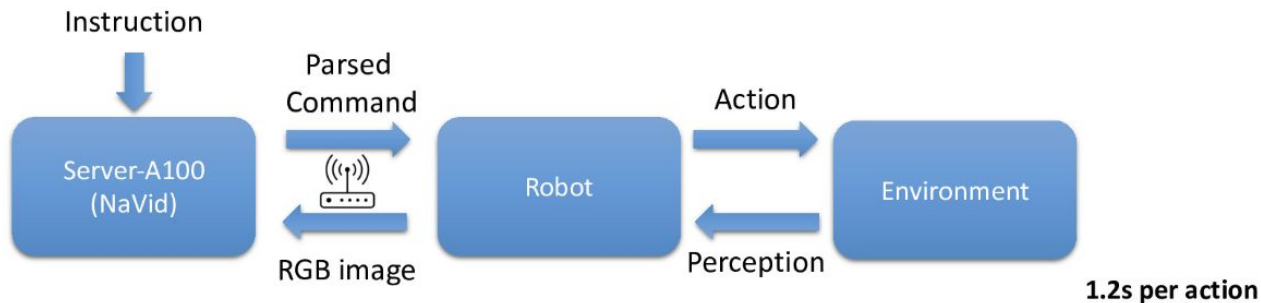


The 'Steps' of the x-axis indicate the oracle actions required by the instructions.

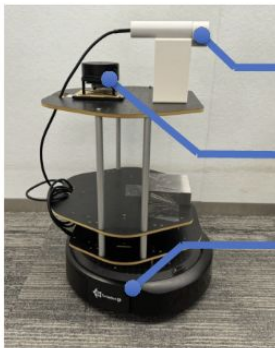
# Experiments

## Real-World Experiments

Pipeline



Robot

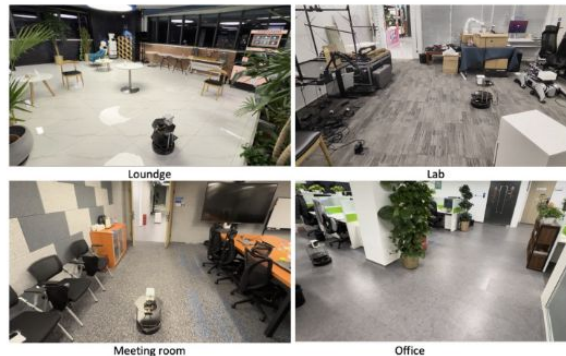


Azure Kinect DK

RPLIDAR A1M8 Lidar

Turtlebot 4

Environment

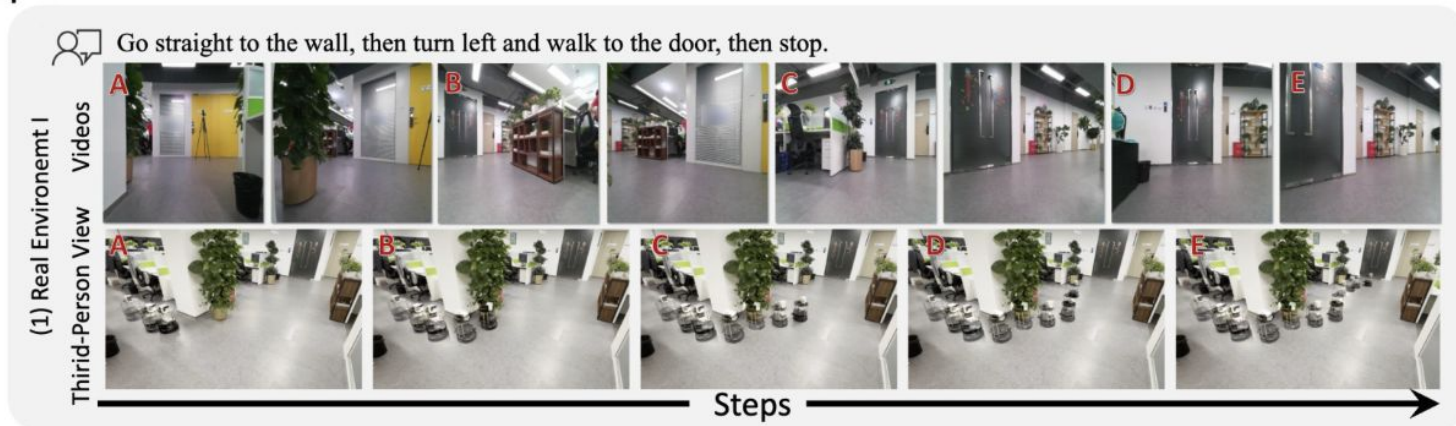


# Experiments

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:



# Experiments

Simple Instruction following

Walk towards the **door** then stop.



Speed up x10

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



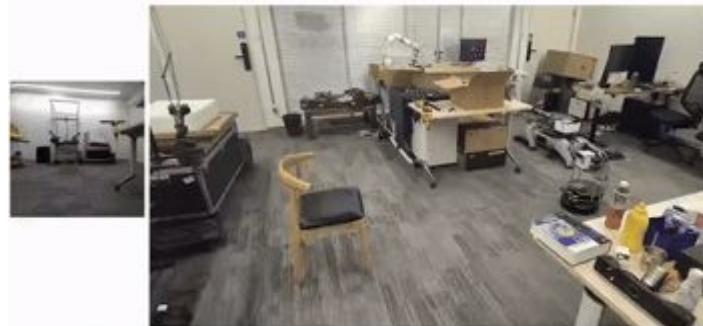


# Experiments

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



Walk towards the chair then turn right, and move to the door, then stop.



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





# Conclusion & Limitations

---

## Takeaway Messages

- NaVid navigates in a human-like manner, requiring solely an on-the-fly video stream from a monocular camera as input, without the need for maps, odometers, or depth inputs.
- We collect 510K VLN video sequences from simulation environments and 763K real-world caption samples to achieve cross-scene generalization.
- With more high-quality data and a better architecture, video-based VLM could be a promising pathway to achieve VLN.

---

# CS586 (25 Spring) : Paper Presentation

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

---

**Thank you.**

20244076 Zhaoyan Wang

20244050 Xiangchen Liu