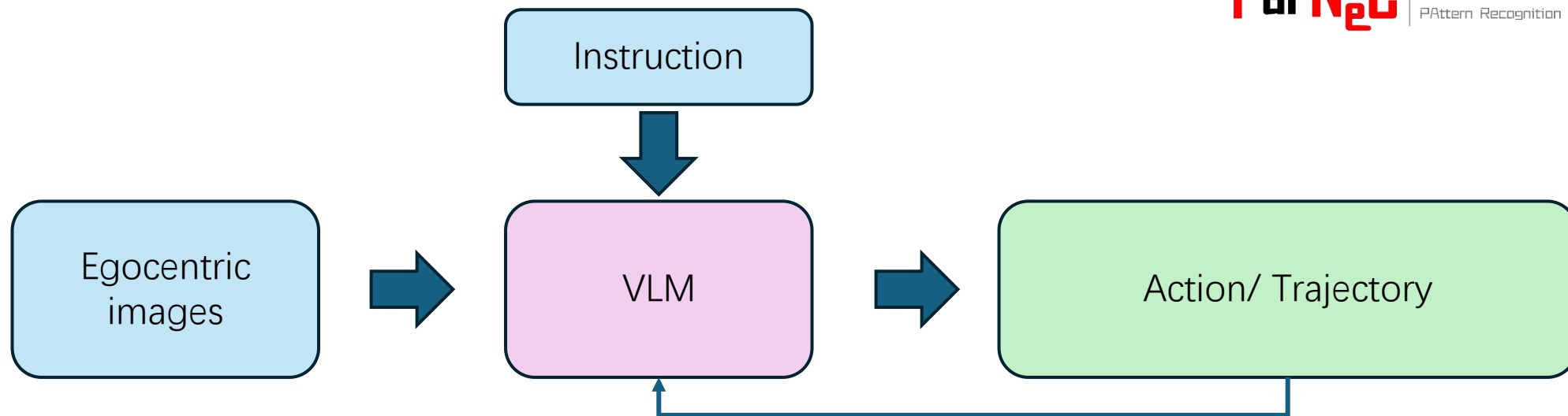


Recent Advances in Vision-and- Language Navigation

欣子豪
2026. 4. 20

1. Existing Paradigms
 1. Single-sys
 2. Dual-sys
 3. Depth-based-sys
2. Our Works
 1. DecoVLN
 2. IDEAL-VLN
 3. AgentVLN
3. Further Works



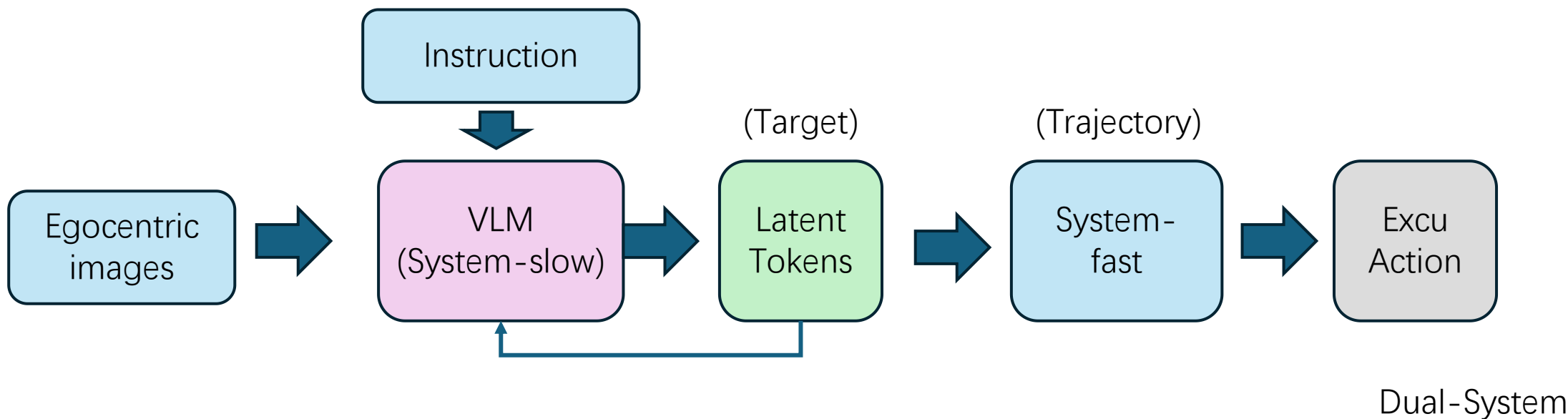
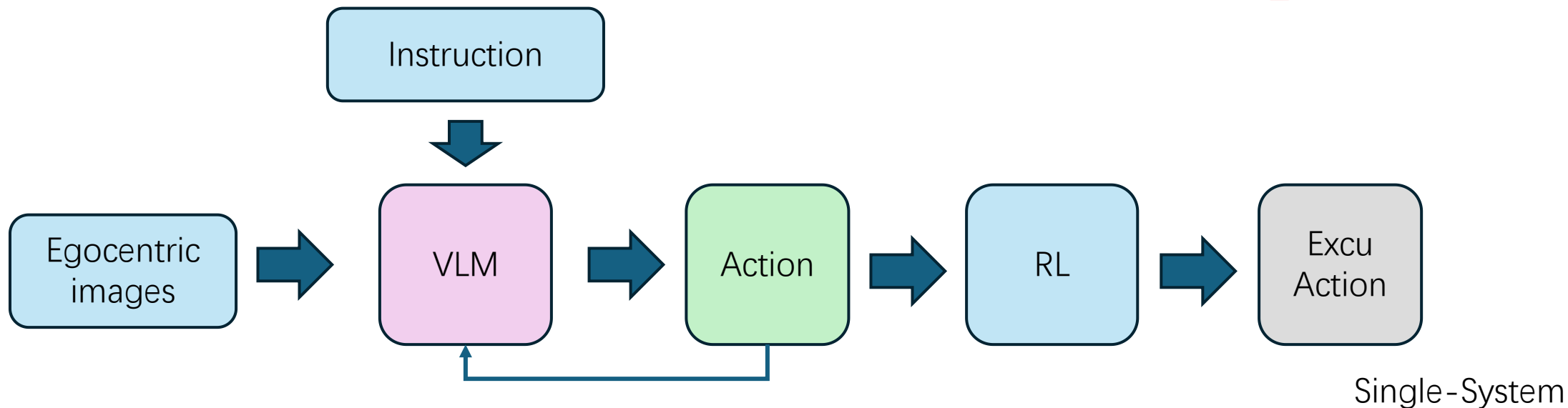
Now you are facing towards the wall, slightly turn to your right, you can see a **open door** to your right, exit room through an open door and move forward, there is a **wall art to the left side of the wall** and few steps in-front of you, walk down the steps and move forward, there are **two black couches** to your right, that's your end point.



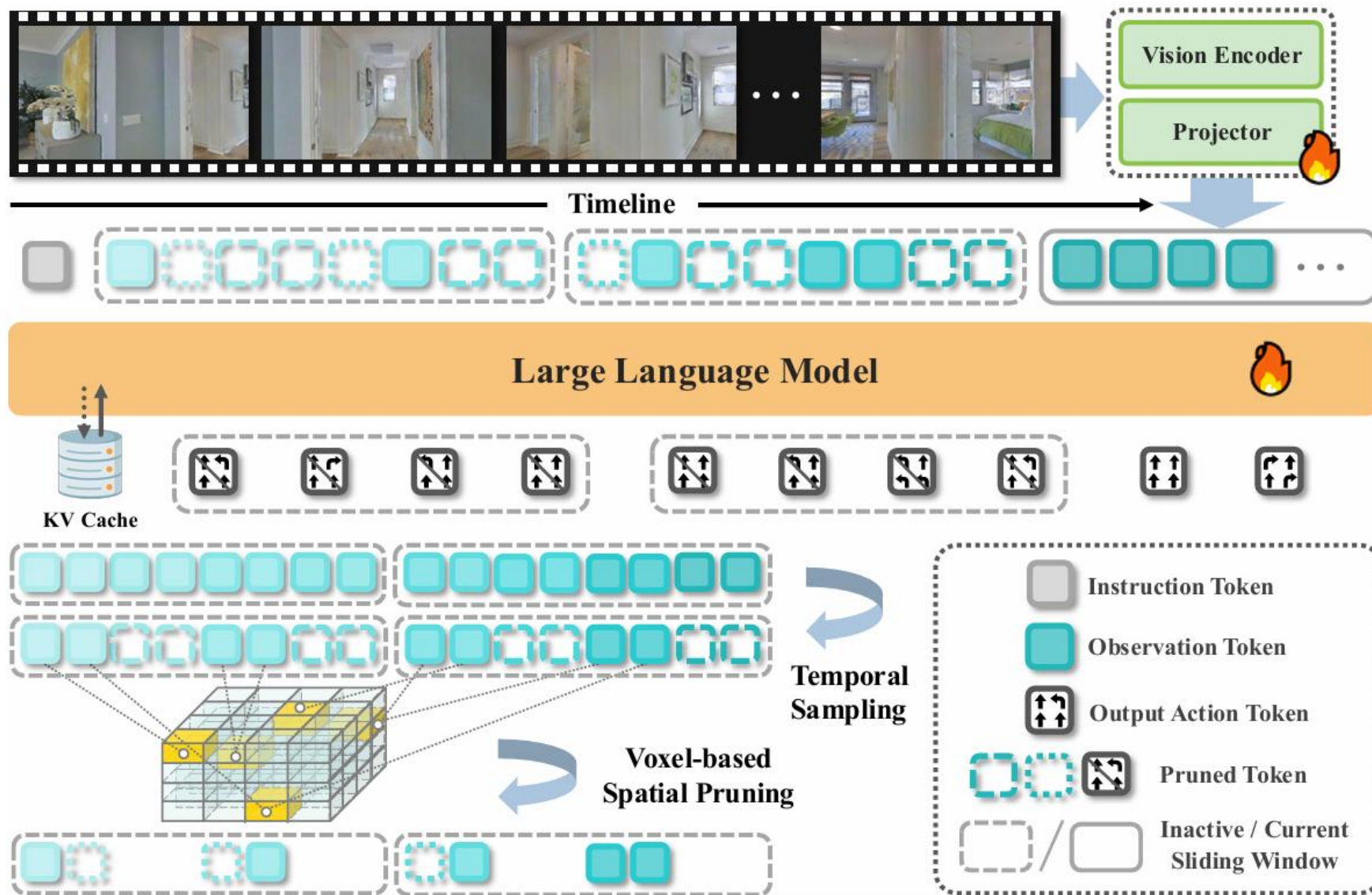
Now you are facing towards the wall, slightly turn to your right, you can see a **open door** to your right, exit room through an open door and move forward, there is a **wall art to the left side of the wall** and few steps in-front of you, walk down the steps and move forward, there are **two black couches** to your right, that's your end point.

What does VLM need to learn?

1. Depth estimation: How far is the wall/door?
2. Contextual semantic analysis: What exactly does "slightly" refer to in terms of perspective?
3. Long and short memory
4. Action Space
- ...

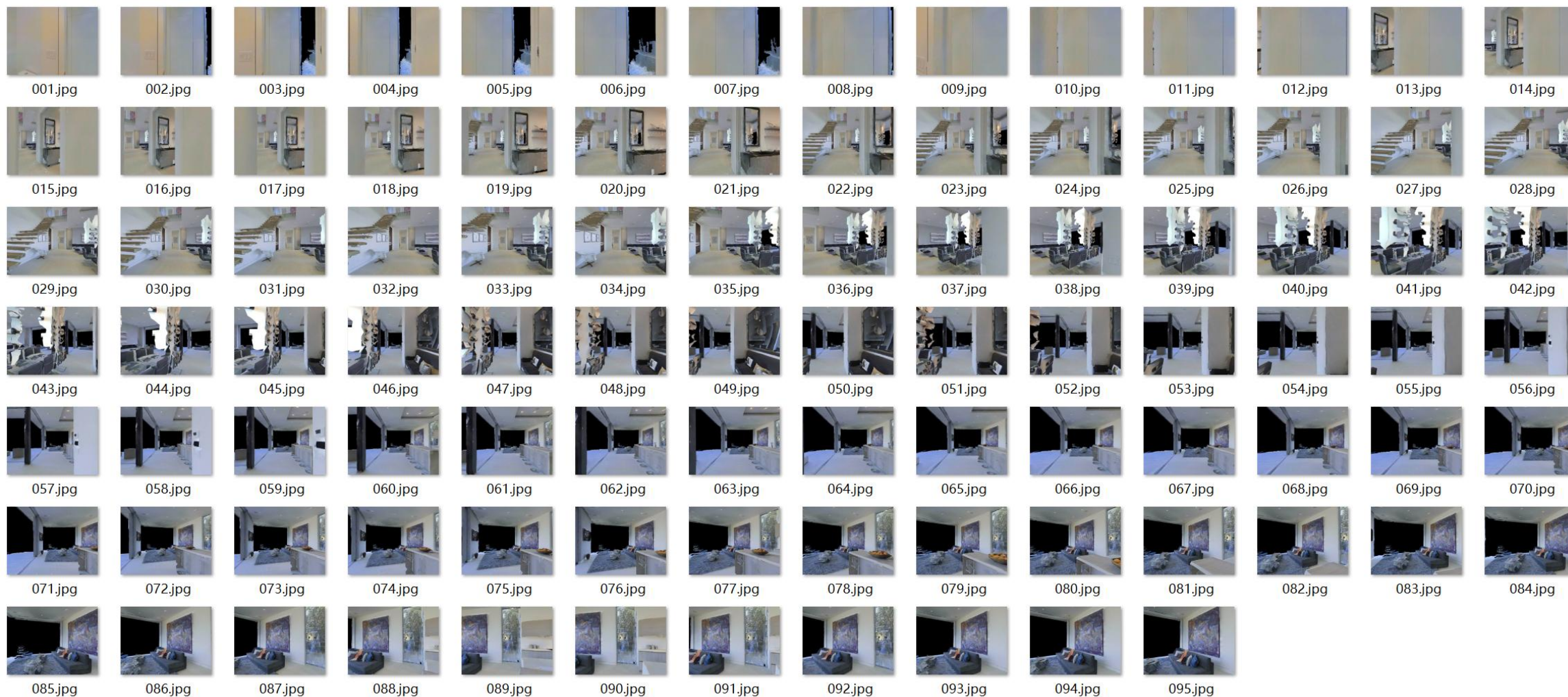


StreamVLN (Single-System, ICRA2026)



1. Long and short term memory navigation capability
2. Positioning and planning
3. Exploration and error correction

VLN – Long Term Memory



~15m path = ~100 ego-centric frames

StreamVLN (Single-System, ICRA2026)



0~32

1. Uniform Sampling
2. Voxel-based Cutting

0~64

Memory:
8 history frames
+
4 current frames

0~96

~4k tokens

Qwen2.5-VL-7B (BF16)

Vision Tower: $640 * 480 \text{ image} = 46 * 34 \text{ patch} = 1564 \text{ tokens}$

LLM (after merge): $640 * 480 \text{ image} = 23 * 17 \text{ patch} = 391 \text{ tokens}$

KV Cache = ~23 MB

Peak Inference Mem = Model Weight + **Vision Tower Peak**
+ **LLM Prefill Peak** + KV Cache + Allocator Reserve



Hidden + QKV + MLP $\approx 100\text{-}300 \text{ MB} / \text{frame}$

Qwen2.5-VL-7B (BF16): *10 v.s. 100 frames*

Time of First Token (TTFT)

Vision Tower: 10x

LLM Prefill: ~90x

$$N = \text{num_frames} * 391 + \text{text_tokens}$$

Self-attn: $O(N^2 \cdot d)$

Compute-bound

Time Per Output Token(TPOT)

KV Cache: 10x

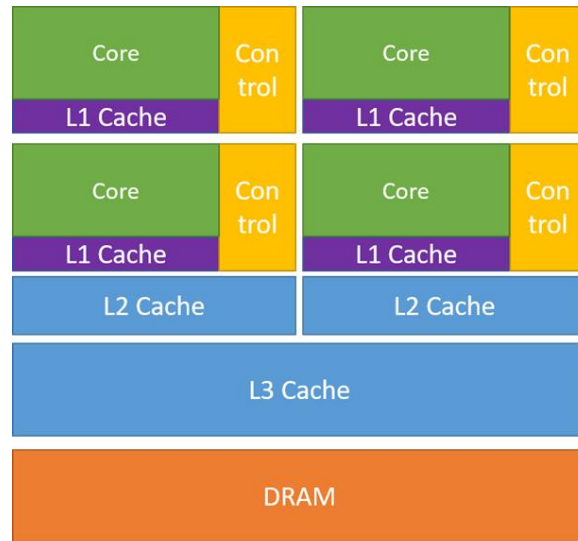
New Token Query (1 x d)

x

KV Cache (d x N)

Memory-bound

GPU Memory Hierarchy

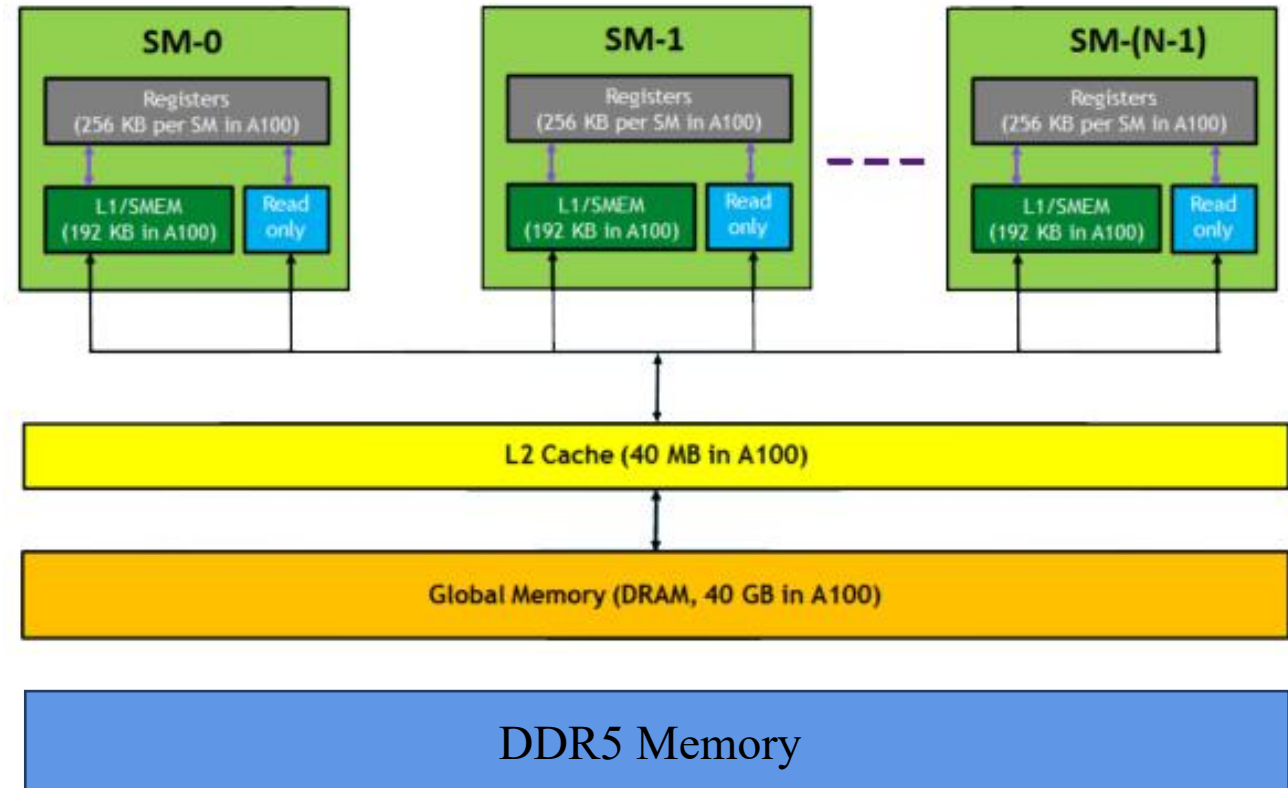


CPU

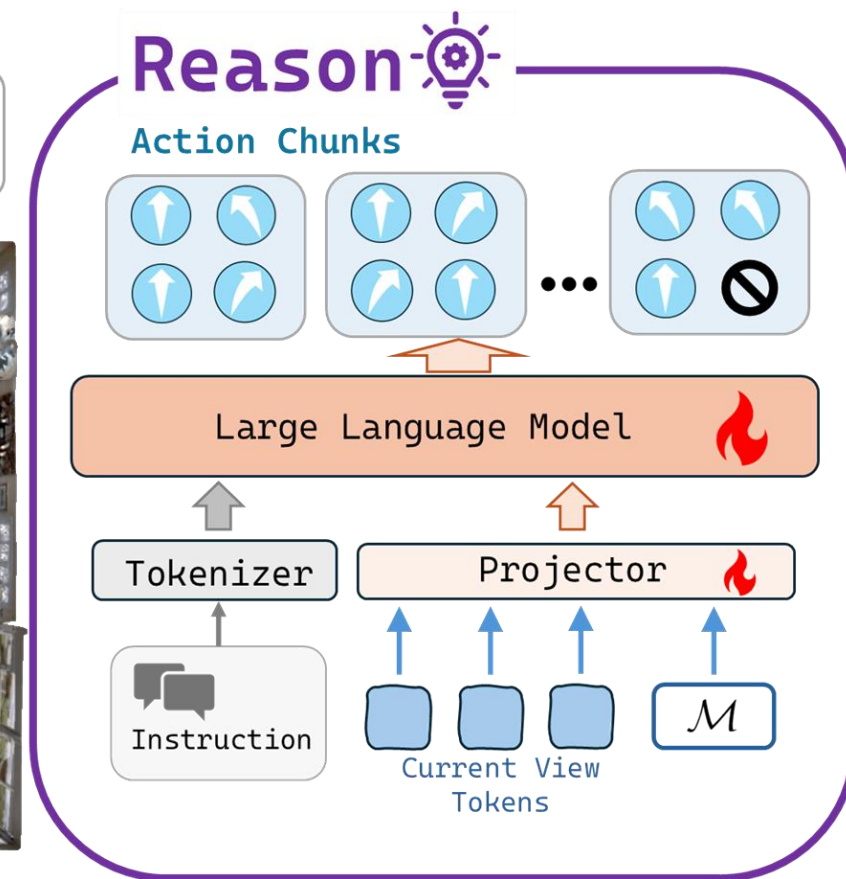
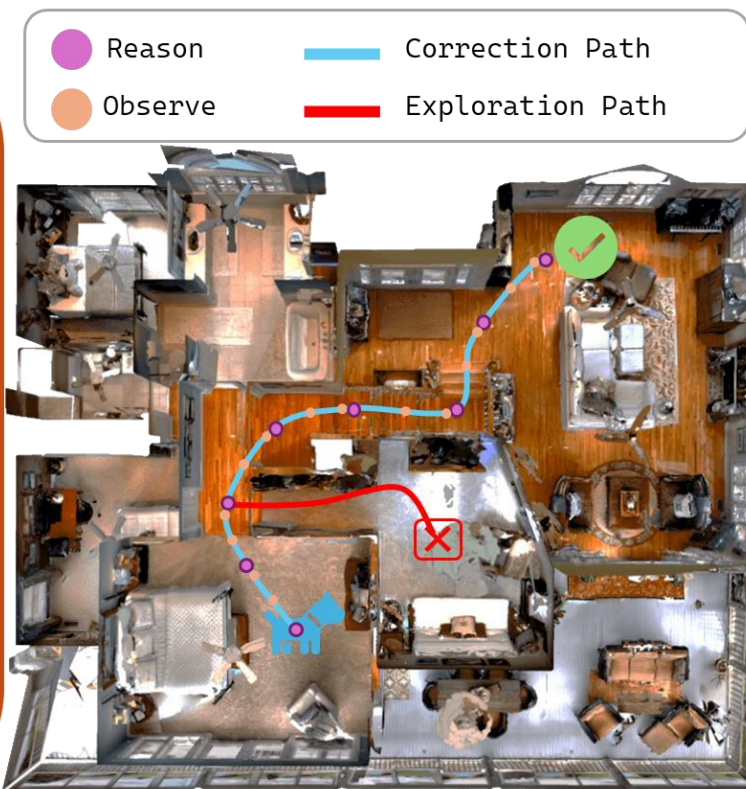
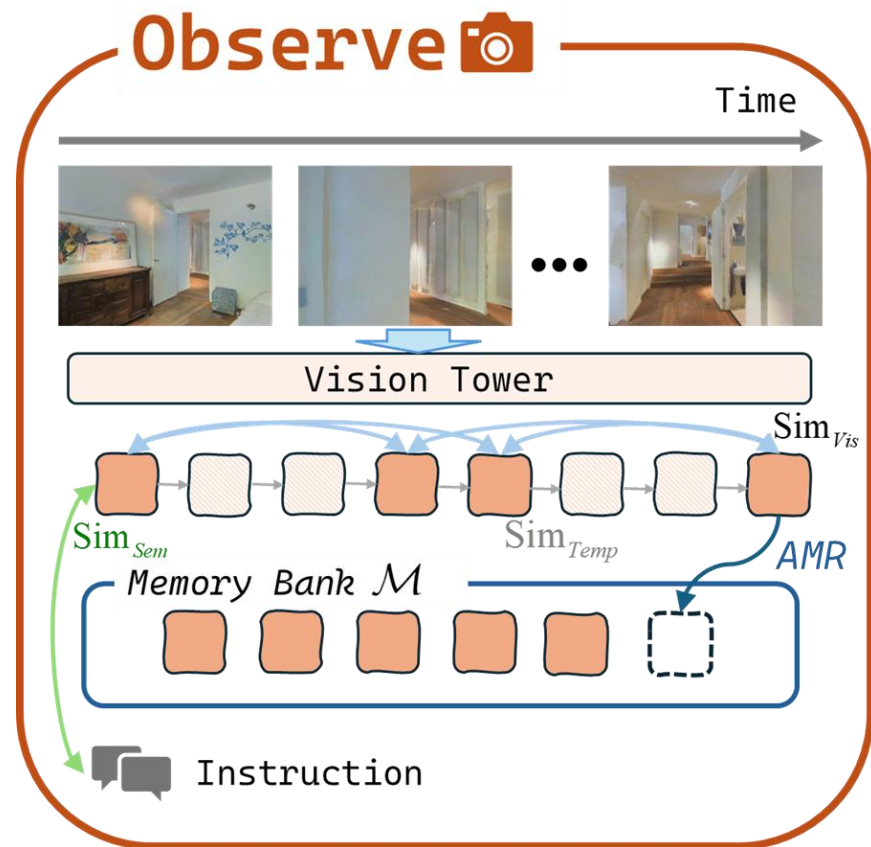


GPU

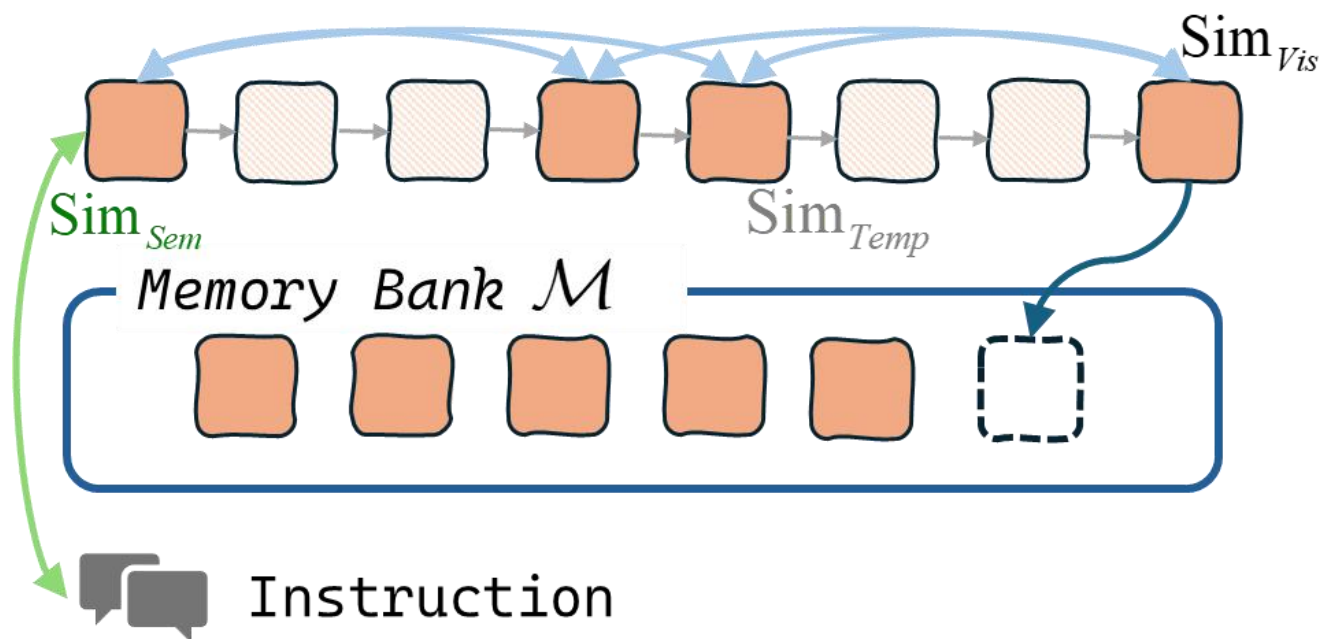
The GPU Devotes More Transistors to Data Processing



DecoVLN (Single-System, CVPR2026)



Consider long-term memory as an optimization problem.



$$\text{Sim}_{Vis}(f, \mathcal{M}) = \max_{m \in \mathcal{M}} \left(\frac{\text{embed}(f) \cdot \text{embed}(m)}{\|\text{embed}(f)\| \|\text{embed}(m)\|} \right)$$

$$\text{Sim}_{Sem}(f, I) = \text{sim}(\mathbf{e}_f, \mathbf{e}_I) = \frac{\mathbf{e}_f \cdot \mathbf{e}_I}{\|\mathbf{e}_f\| \|\mathbf{e}_I\|}$$

$$\text{Sim}_{Temp}(f, \mathcal{M}) = \frac{1}{\min_{m \in \mathcal{M}} |t_f - t_m| + \epsilon}$$

$$f^* = \arg \max_{f \in \mathcal{C} \setminus \mathcal{M}} [\lambda_R \cdot \text{Sim}_{Sem}(f, I) - (1 - \lambda_R) \cdot (w_V \cdot \text{Sim}_{Vis}(f, \mathcal{M}) + w_T \cdot \text{Sim}_{Temp}(f, \mathcal{M}))]$$



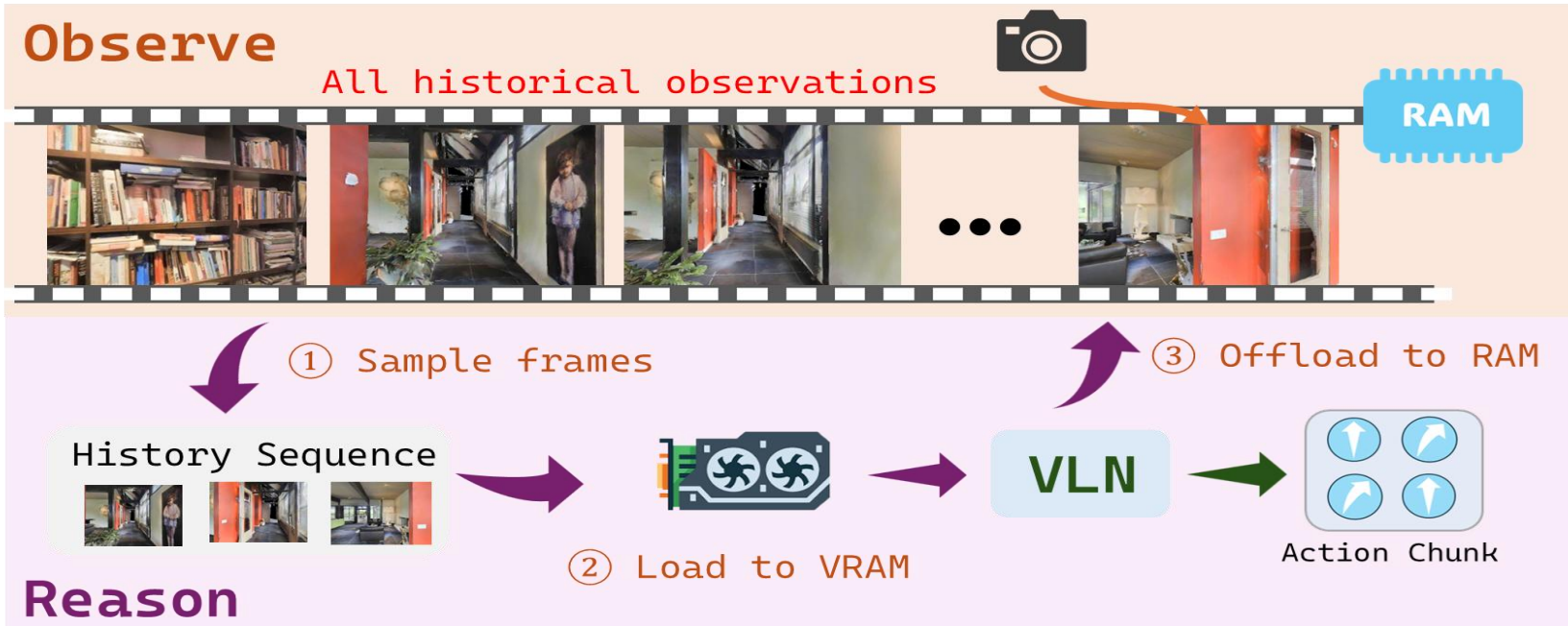
(a) Uniform Sampling



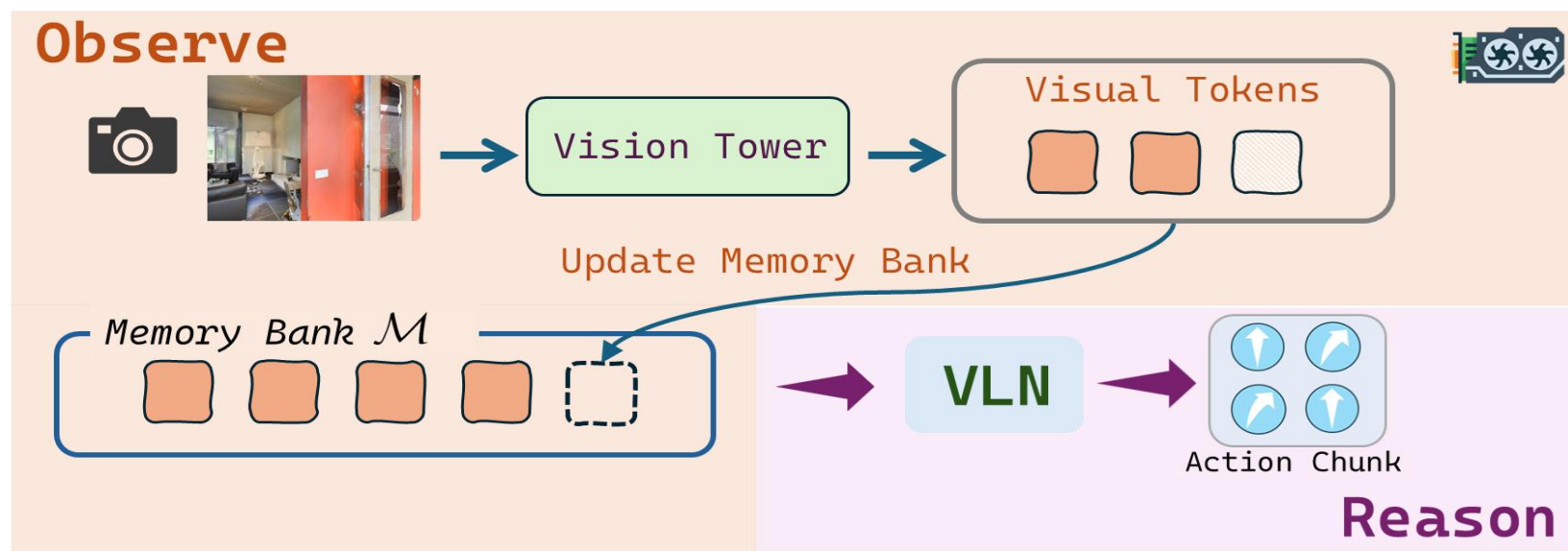
(b) Adaptive Memory Refinement



Now you are facing towards the wall, slightly turn to your right, you can see a **open door** to your right, exit room through an open door and move forward, there is a **wall art to the left side of the wall** and few steps in-front of you, walk down the steps and move forward, there are **two black couches** to your right, that's your end point.



Streaming Paradigm



DecoVLN

StreamVLN



DecoVLN



DecoVLN (Single-System, CVPR2026)

1. Expert Dataset: R2R, RxR (Stage 1)
2. Collect correction data using stage 1 model
3. Multi-Mode Dataset: Correction data + MM Datasets

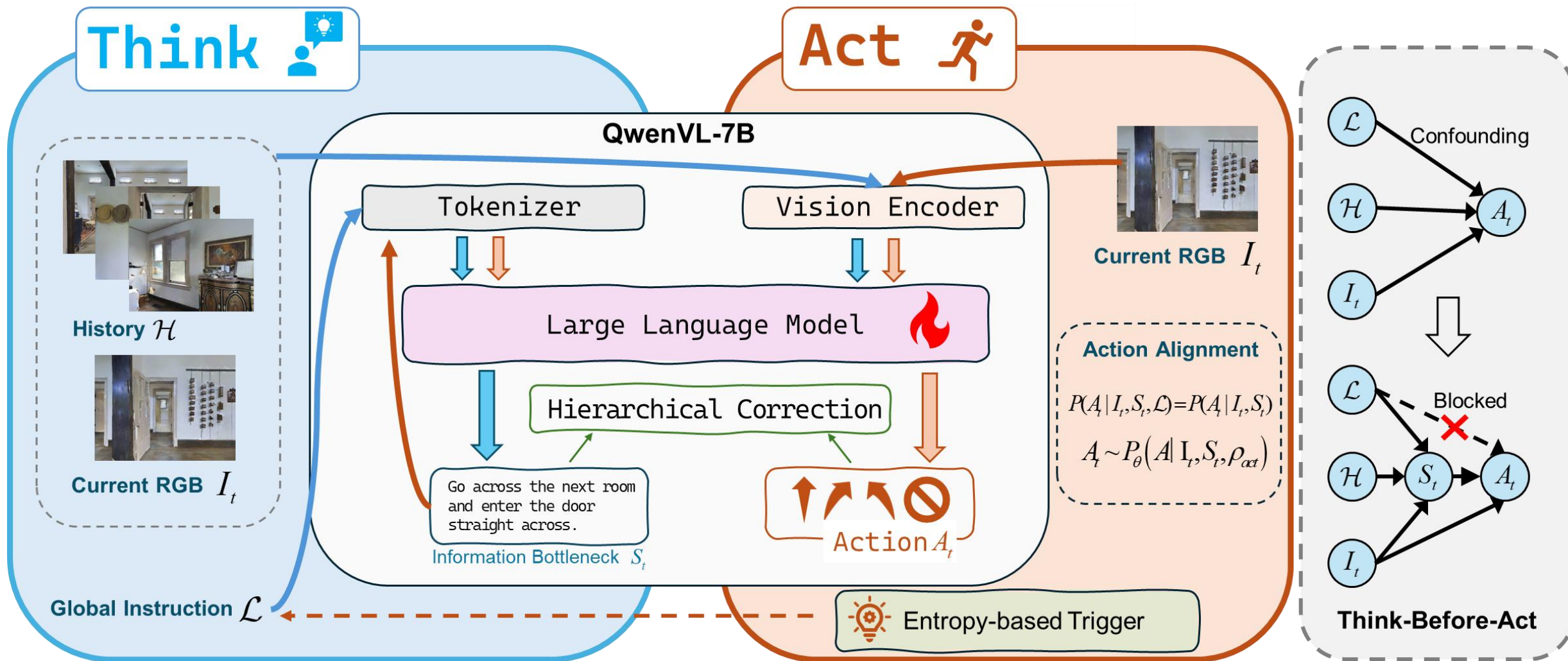
Algorithm 1 Correction Data Collection

Require: Initial policy π_θ , expert policy π^* , dataset $\mathcal{D} = \{(P_{\text{exp}}^i, I^i)\}_{i=1}^M$, deviation threshold τ

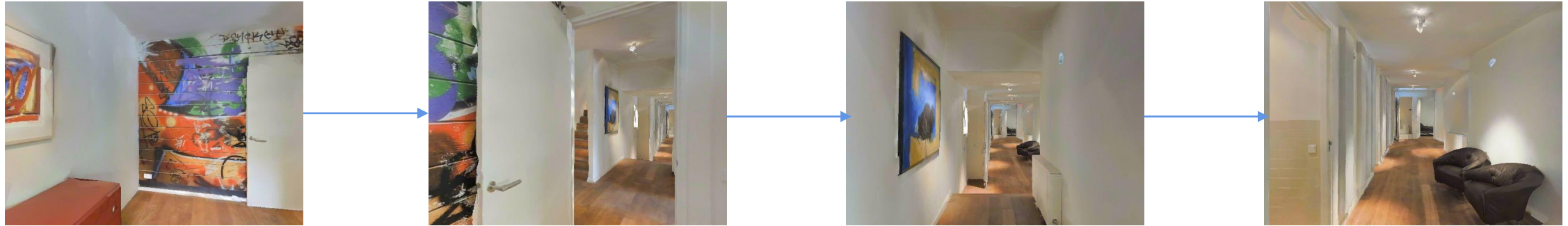
Ensure: Correction dataset \mathcal{D}_c

```
1:  $\mathcal{D}_c \leftarrow \emptyset$  ▷ Initialize correction dataset
2: for all episode  $(P_{\text{exp}}, I) \in \mathcal{D}$  do
3:    $P_{\text{exp}} \leftarrow \{s_0^*, s_1^*, \dots, s_N^*\}$  ▷ Expert reference trajectory
4:    $H \leftarrow \emptyset$  ▷ Initialize memory state
5:   for  $t = 0$  to  $T_{\text{max}}$  do
6:      $DM \leftarrow \min_{s^* \in P_{\text{exp}}} d_g(s_t, s^*)$  ▷ Compute deviation from expert trajectory
7:     if  $0 < DM(s_t) \leq \tau$  then
8:        $a_t^{\text{exp}} \leftarrow \pi^*(s_t, P_{\text{exp}})$ 
9:        $\mathcal{D}_c \leftarrow \mathcal{D}_c \cup \{(s_t, a_t^{\text{exp}}, f_t)\}$  ▷ Store expert correction
10:    else if  $DM(s_t) > \tau$  then
11:      break ▷ Abort episode due to large deviation
12:    end if
13:     $a_t \sim \pi_\theta(f_t, I, H)$  ▷ Roll out current policy
14:     $s_{t+1} \leftarrow \text{Env.Step}(a_t)$  ▷ Environment state update
15:     $f_{t+1} \leftarrow \text{Env.Observe}(s_{t+1})$  ▷ Acquire next observation
16:     $H \leftarrow \text{AMF}(H, f_{t+1})$  ▷ Adaptive memory refinement
17:    if target reached or episode failure then
18:      break ▷ Terminate episode
19:    end if
20:  end for
21: end for
22: return  $\mathcal{D}_c$ 
```

IDEAL-VLN (Single-System)

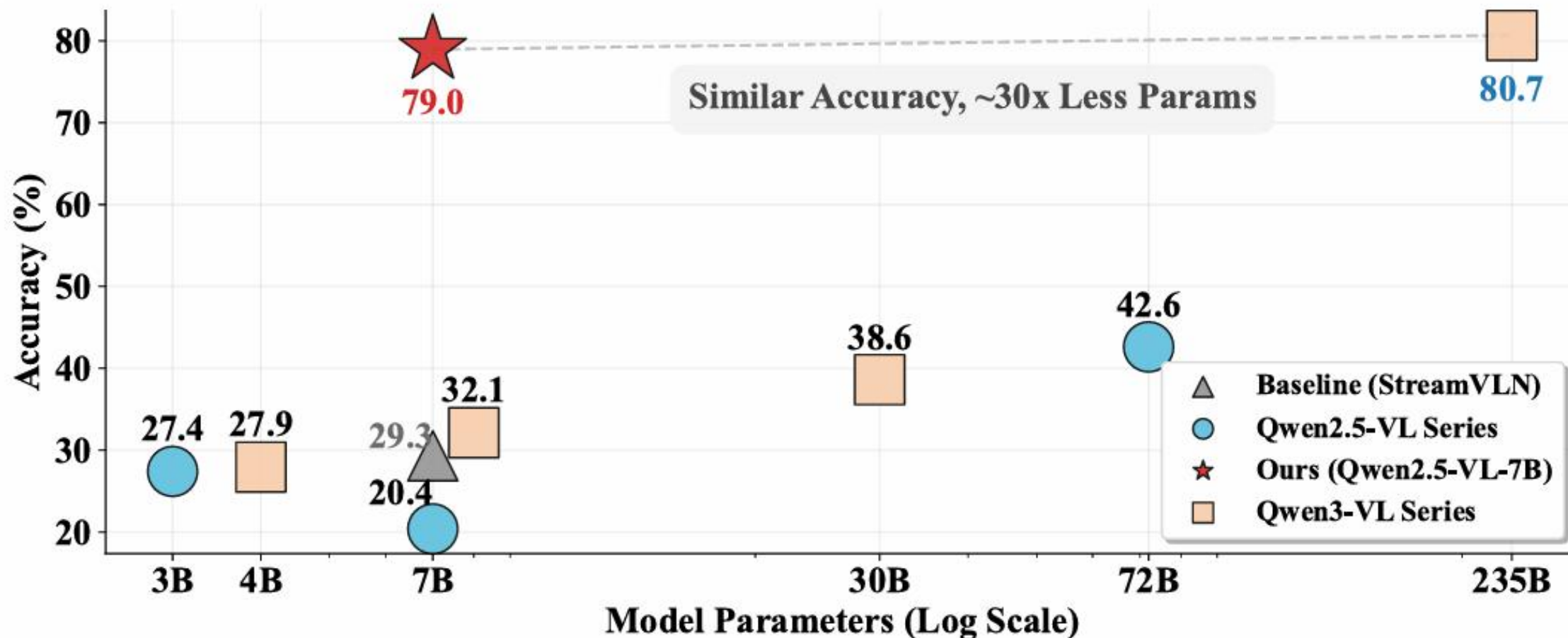


Global Instruction Turn to your left and exit the room. Go across the next room and enter the door straight across. Stop once you enter the bathroom door.

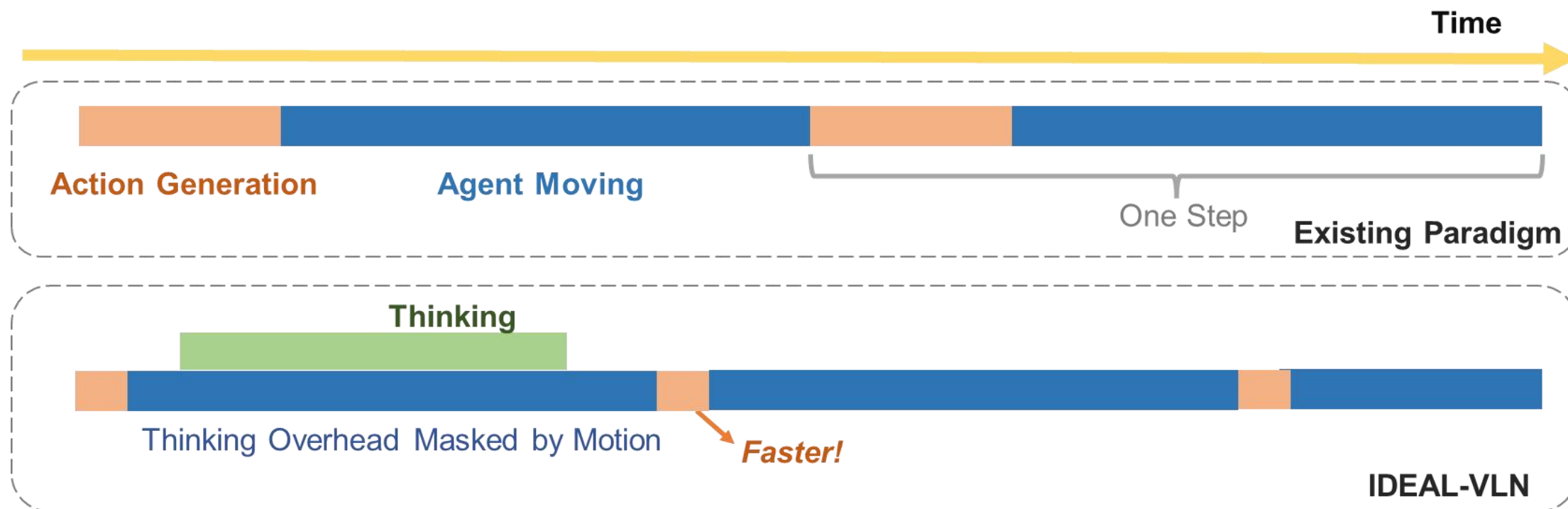


Now you are facing towards the wall, slightly turn to your right, you can see a open door to your right, exit room through an open door and move forward, there is a wall art to the left side of the wall and few steps in-front of you, walk down the steps and move forward, there are two black couches to your right, that's your end point.

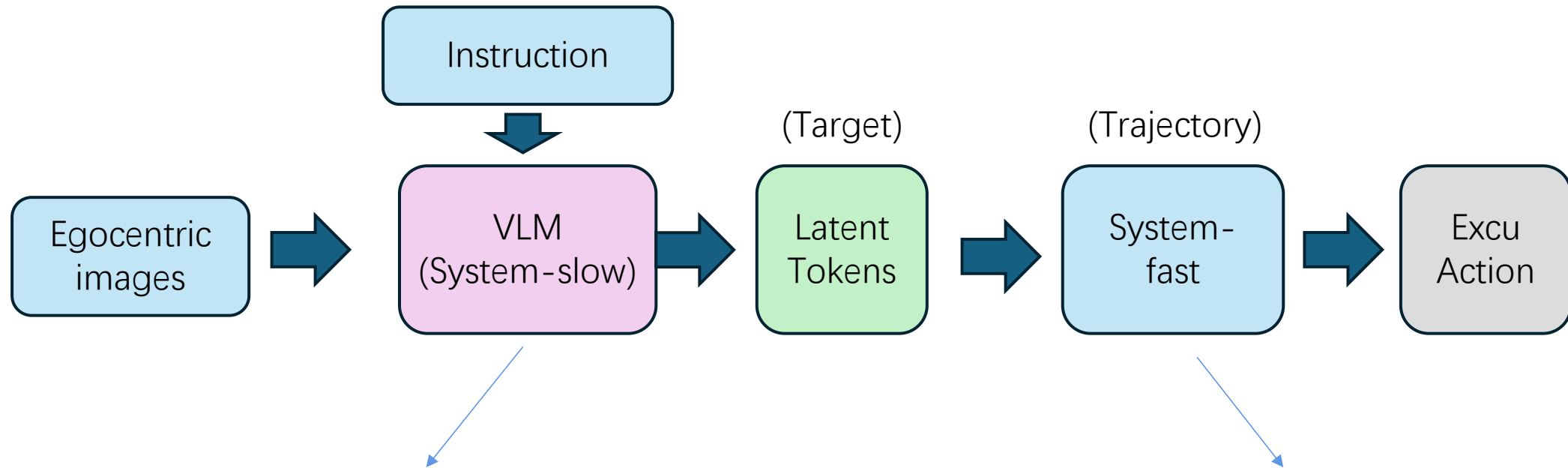
- ✗ Rule-Based Split (word count or period)
- ✗ Text Analysis and Split
- ✗ LLM-based Split



The existing models cannot comprehend the relationship between instructions and ego-centric images.



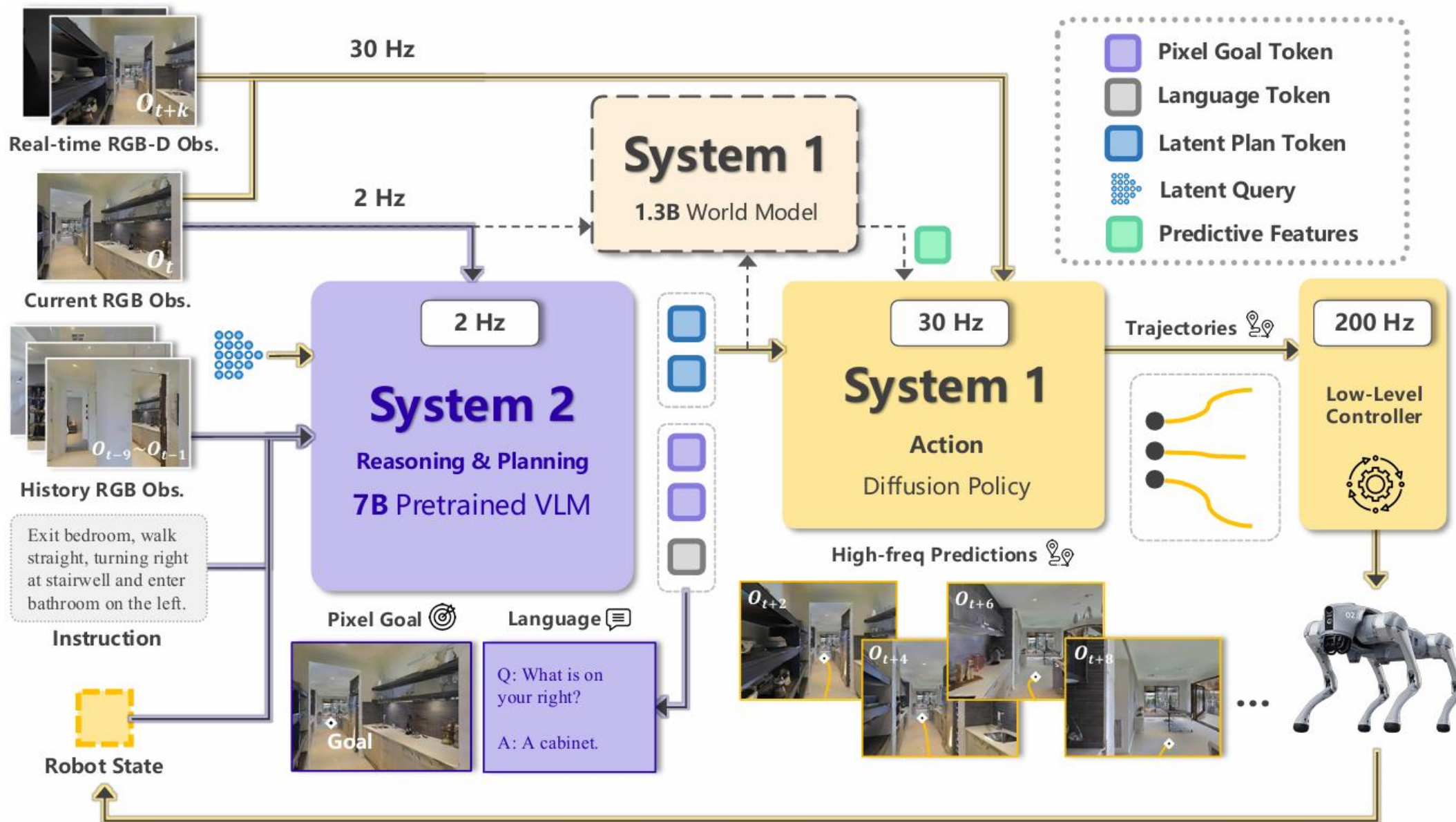
Method	Input tokens	Prefill Time(ms)	Steps/ Episode
Vanilla	7k	1625	21
Ours (Think)	7k	1832	5
Ours (Act)	2k	462	21

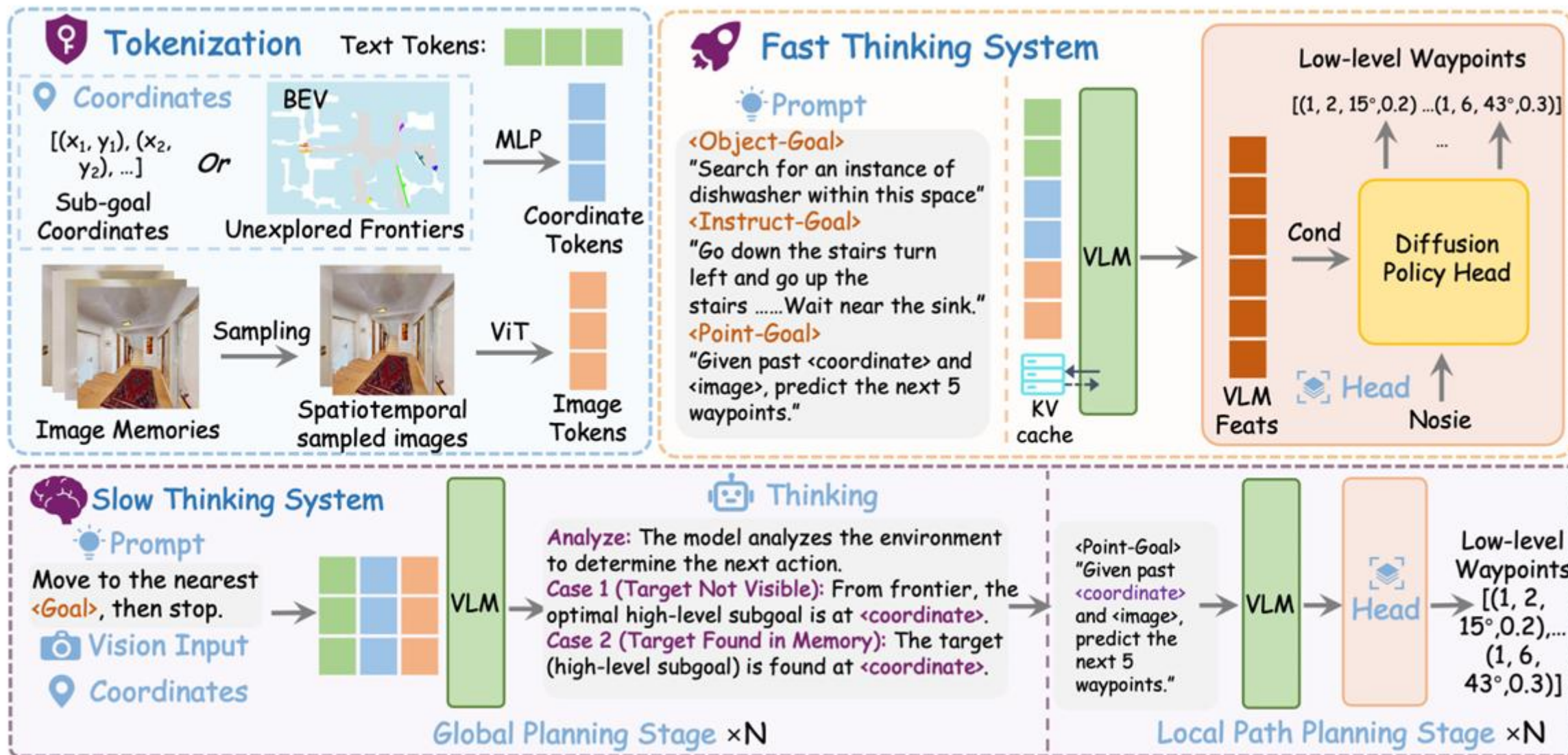


Only responsible for generating navigation points (image coordinate system)

Convert the image coordinate into an executable trajectory

VLM is action-independent





VLM The VLMs in the fast and slow systems share the same parameters.



1. Joint training of the two systems is challenging
2. Incorrect navigation points may cause trajectory errors
3. Only available for specific robots
4. Very, very expensive

A joint training scheme is particularly critical for Vision-and-Language Navigation (VLN), which requires strong general knowledge and semantic understanding, leading to substantial improvements in success rates in open environments. Stage 1 is trained with 96 NVIDIA H20 GPUs for 120 hours, and Stage 2 is trained with 64 NVIDIA H20 GPUs for 48 hours with lower learning rate.

-- OmniNav

18000+ H20 hours

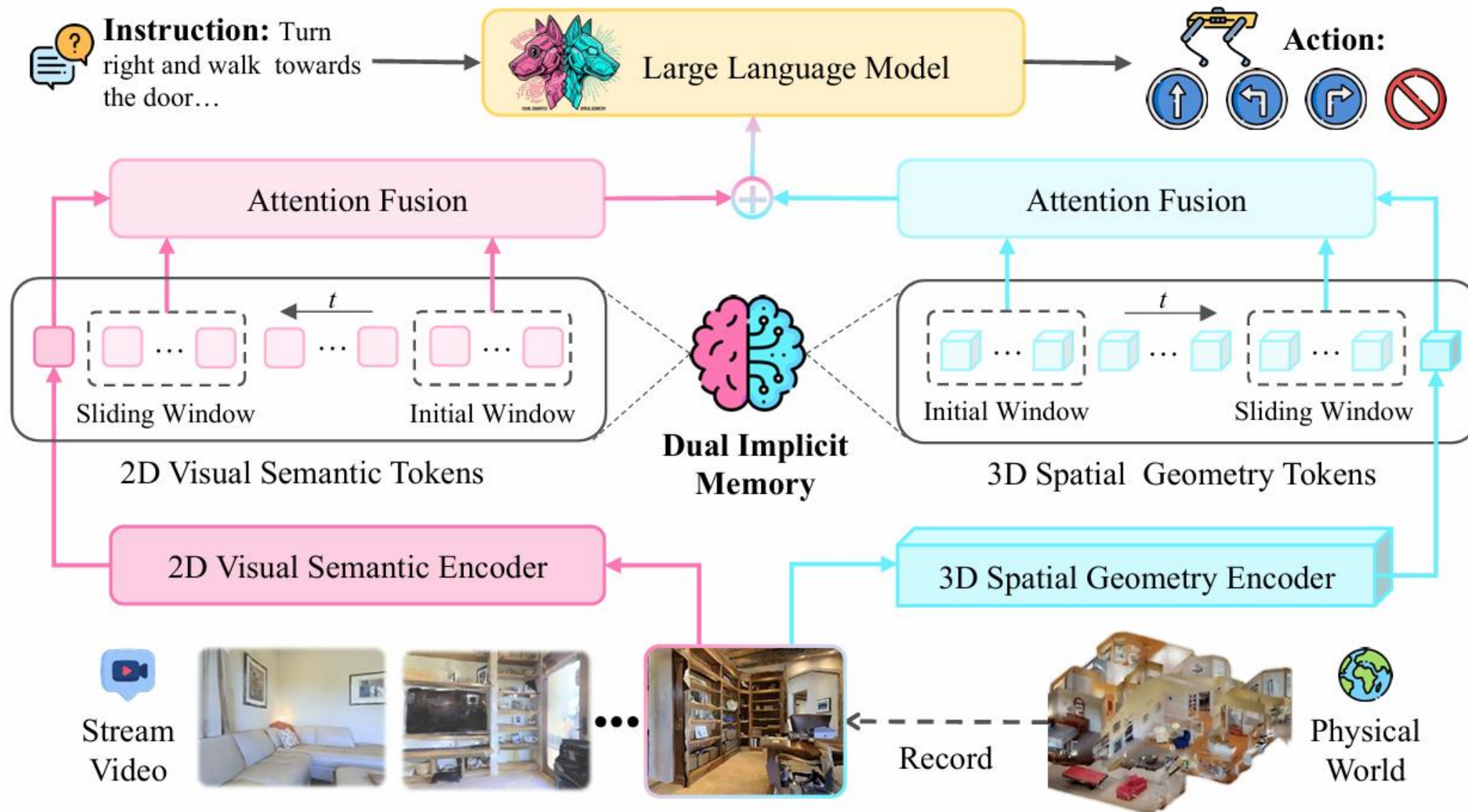
Dual-sys Training Cost: 15W+ RMB

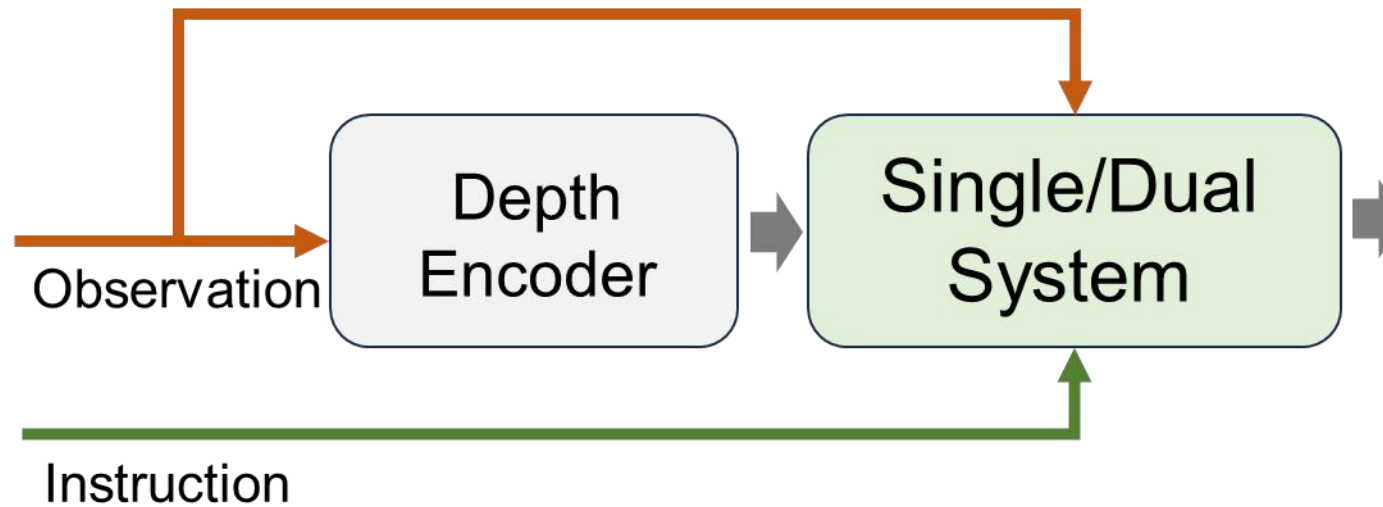
Reproduction is extremely time-consuming

DecoVLN/ IDEAL-VLN: ~ 3k RMB

Depth Encoder-based Sys

JanusVLN (ICLR2026)





Additional depth encoder required

Rely on depth encoder -> poor performance in real-world

Significantly increases computational load and storage overhead

Depth Encoder-based Sys

Why not use real depth map?



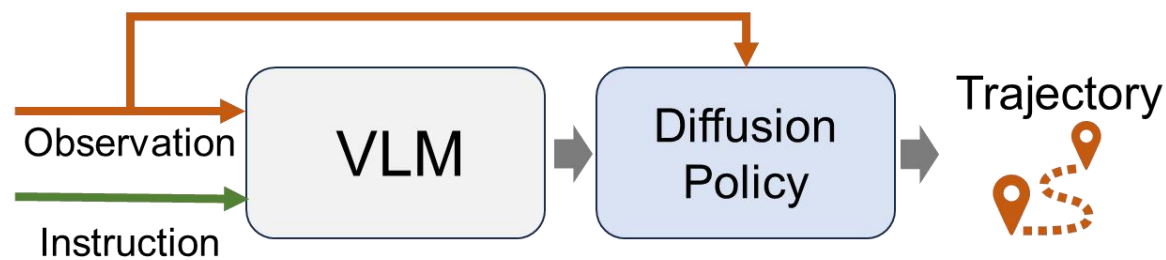
Depth map by realsense camera



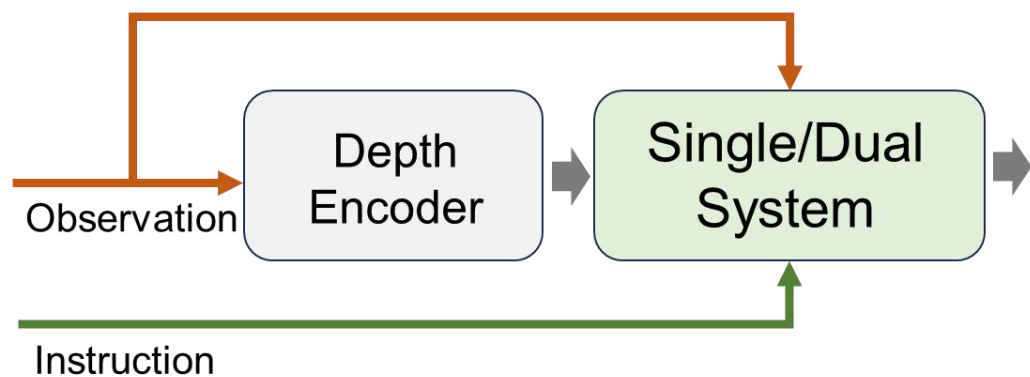
DepthAnything



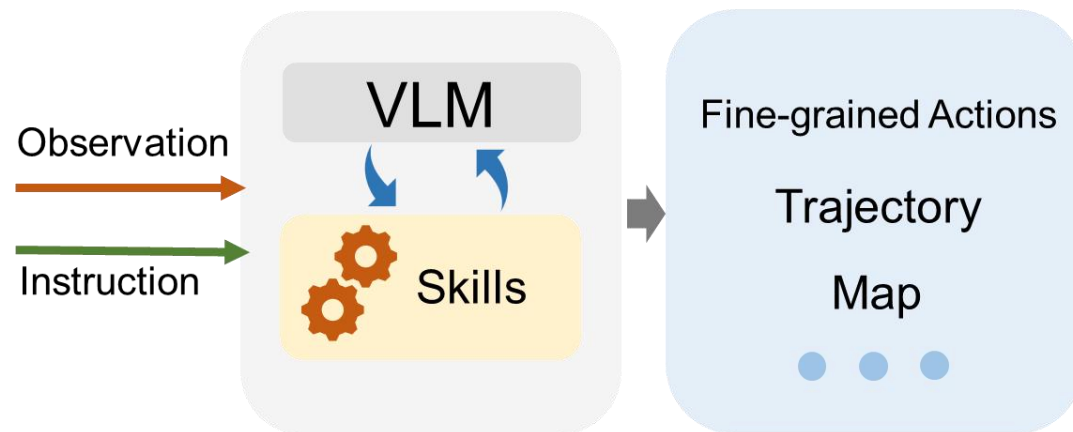
(a) Single-System



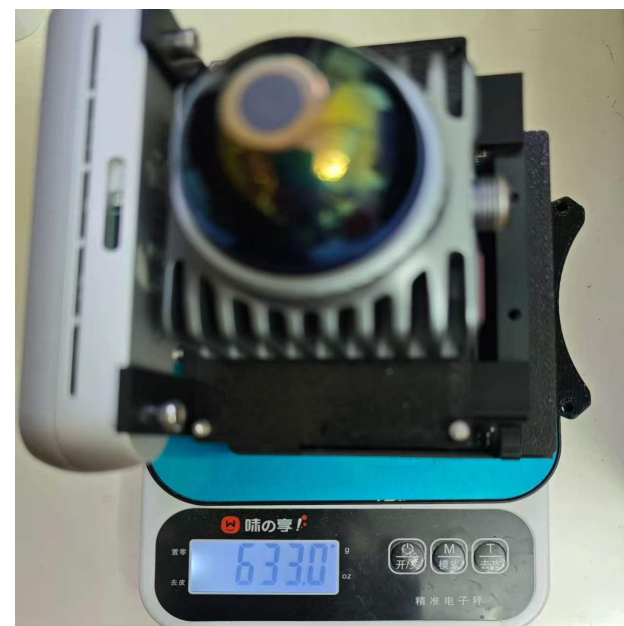
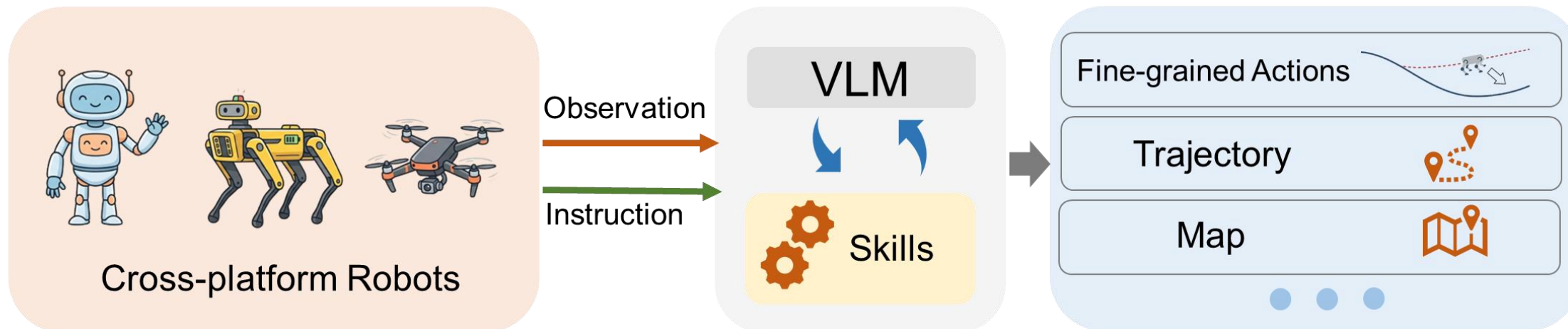
(b) Dual-System

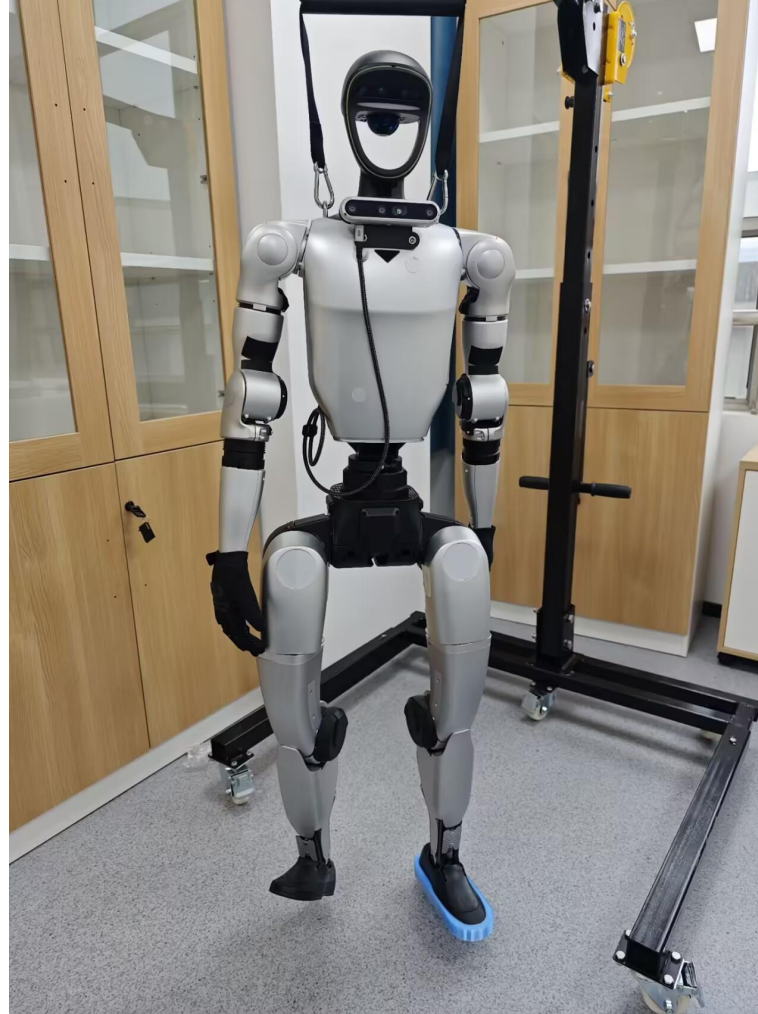
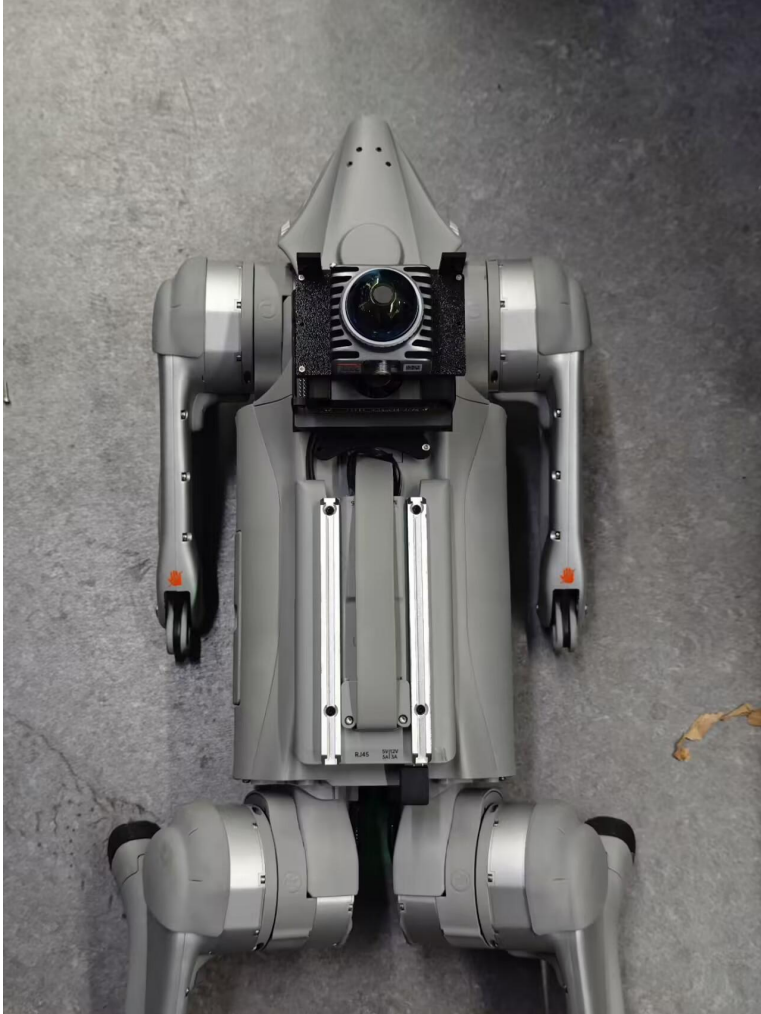


(c) Depth Encoder-based System



(d) AgentVLN



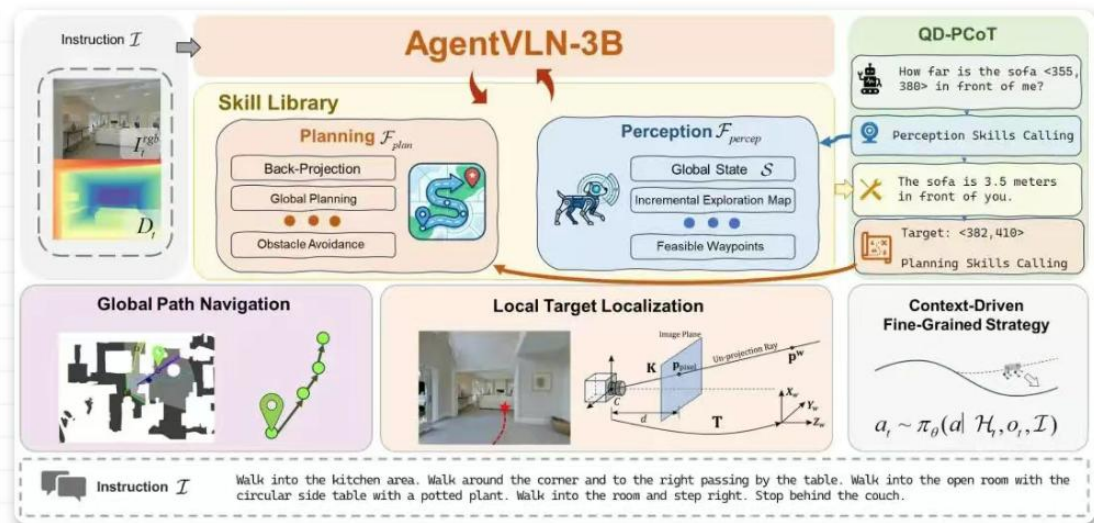


TODO

Quadrotor

小参数量，大能力！AgentVLN：轻量部署的下一代具身导航新范式

VLNer 视觉语言导航 2026年3月25日 21:50 湖南



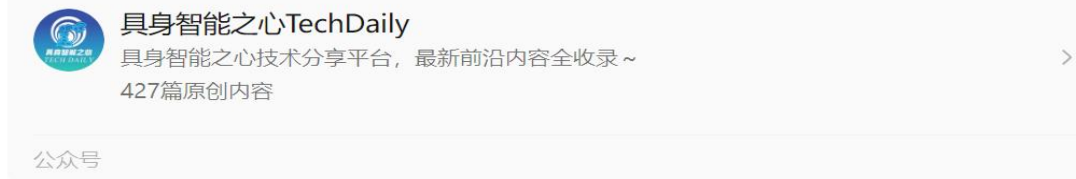
- **作者:** Zihao Xin¹, Wentong Li^{1*}, Yixuan Jiang¹, Ziyuan Huang¹, Bin Wang², Piji Li¹, Jianke Zhu³, Jie Qin¹, Shengjun Huang¹
- **单位:** ¹南京航空航天大学, ²山东大学, ³浙江大学
- **论文标题:** AgentVLN: Towards Agentic Vision-and-Language Navigation
- **论文链接:** <https://arxiv.org/abs/2603.17670v1>
- **代码链接:** <https://github.com/Allenxinn/AgentVLN>

10k + views, 1k+ shares

AgentVLN：学会“自主导航”，轻量化具身视觉语言导航新范式！

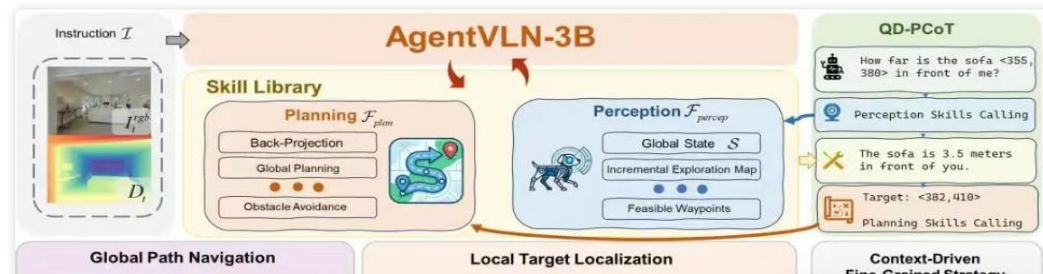
原创 第一具身范式 具身智能之心TechDaily 2026年3月20日 17:32 上海

点击下方卡片，关注“第一具身范式”公众号
第一时间获取具身智能



行业探讨：自主语言导航会终结室内SLAM，还是补齐SLAM最后的短板？

在具身智能的研究版图中，视觉语言导航（VLN）是**极具挑战性**(ps: 好发文章)的方向之一——它要求智能体依托第一视角的动态视觉观测，遵循自然语言指令，在未知的三维物理环境中完成长时程的自主导航。视觉语言模型（VLM^Q）虽在二维语义理解上展现出强大能力，却始终难以突破三维空间感知不足、2D-3D表征错位、单目尺度模糊的瓶颈，让VLN系统的实际落地困难重重。

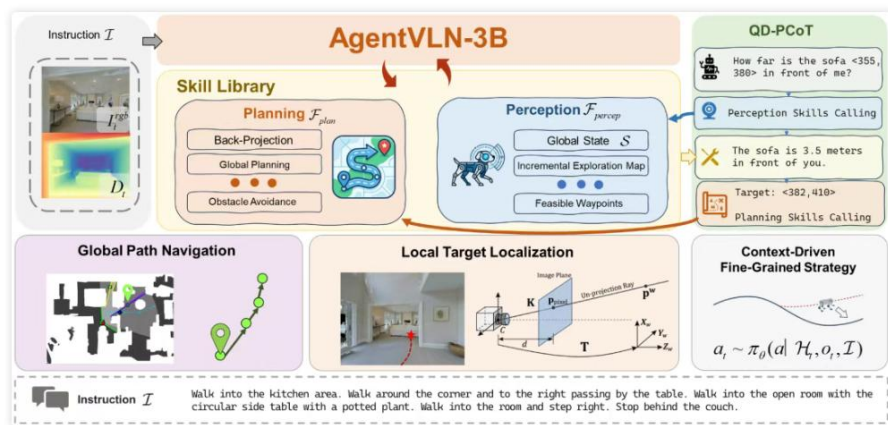


智能体也能“看图说话”自主导航? AgentVLN^Q突破视觉语言导航新边界

基本信息

- 作者: Zihao Xin, Wentong Li, Yixuan Jiang, Ziyuan Huang, Bin Wang, Piji Li, Jianke Zhu, Jie Qin, Shengjun Huang
- 单位: 南京航空航天大学, 山东大学, 浙江大学
- 论文标题: AgentVLN: Towards Agentic Vision-and-Language Navigation
- 论文链接: <https://arxiv.org/abs/2603.17670v1>
- 代码链接: <https://github.com/Allenxinn/AgentVLN>

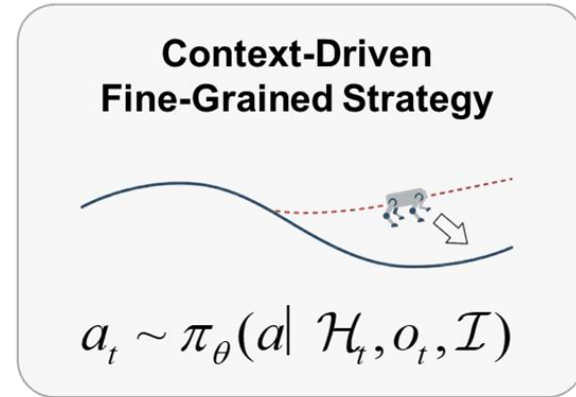
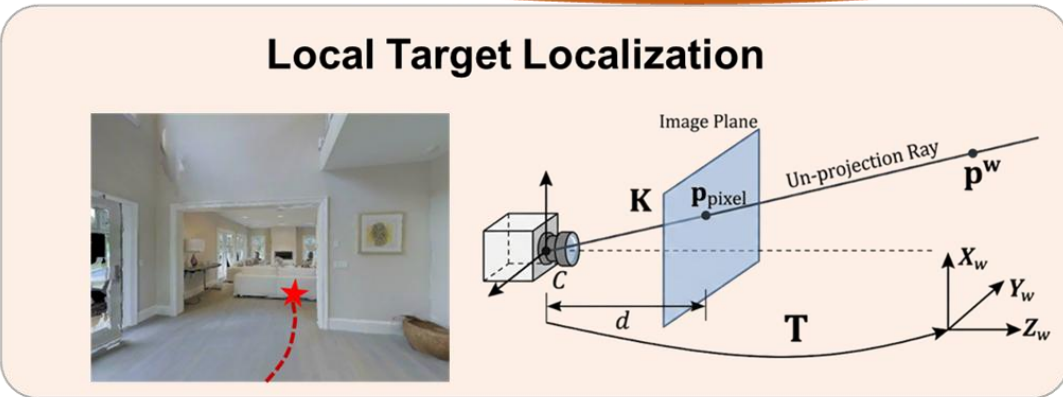
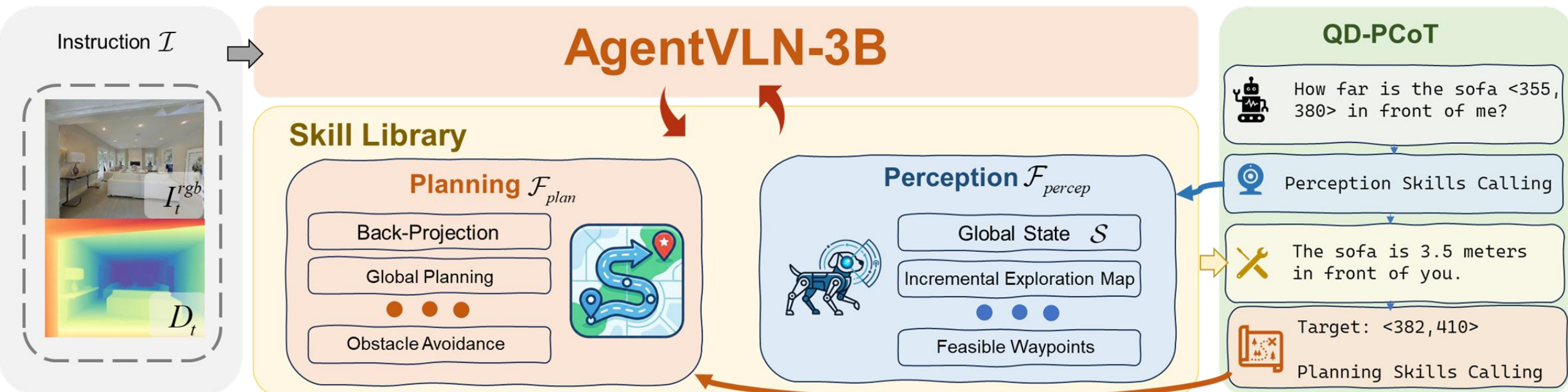
简要介绍



视觉语言导航 (VLN) 要求智能体根据复杂语言指令在未知环境中进行长距离导航。然而, 现有系统常受

The screenshot shows a dataset card for AgentVLN-3B. It includes buttons for **Duplicate** and **Data Studio**. The card displays **Downloads last month: 830** and **Total file size: 215 GB**. There is a link to the **Paper for allenxinn/AgentVLN-Instruct** and a title **AgentVLN: Towards Agentic Vision-and-Language ...** with a link to the paper (2603.17670) published on Mar 18.

Below the screenshot, there is a partial view of a text description: "on and planning via a" and "s of 3D physical".



Instruction \mathcal{I}

Walk into the kitchen area. Walk around the corner and to the right passing by the table. Walk into the open room with the circular side table with a potted plant. Walk into the room and step right. Stop behind the couch.

Duple Camera Sys

$$o_t = \{I_t^{rgb}, D_t\}$$

$$T_t = [R_t \mid \mathbf{t}_t] \in SE(3).$$

$$\mathbf{P}^w = R_t(d \cdot K^{-1} \mathbf{p}_{img}) + \mathbf{t}_t$$

$$s \cdot \mathbf{p}_{path}^{img} = K R_t^{-1} (\mathbf{P}_{path}^w - \mathbf{t}_t)$$

$$\mathbf{P}_{target}^c = d_{target} \cdot K^{-1} \mathbf{p}_{target}^{img} = \begin{bmatrix} \frac{(u_{target} - c_x) \cdot d_{target}}{f_x} \\ \frac{(v_{target} - c_y) \cdot d_{target}}{f_y} \\ d_{target} \end{bmatrix}$$

$$\mathbf{P}_{target}^w = R_t \mathbf{P}_{target}^c + \mathbf{t}_t$$

Lidar + Monocular Camera Sys

$$\mathbf{P}_{L,i} = R_W^L \mathbf{P}_{w,i} + \mathbf{t}_W^L$$

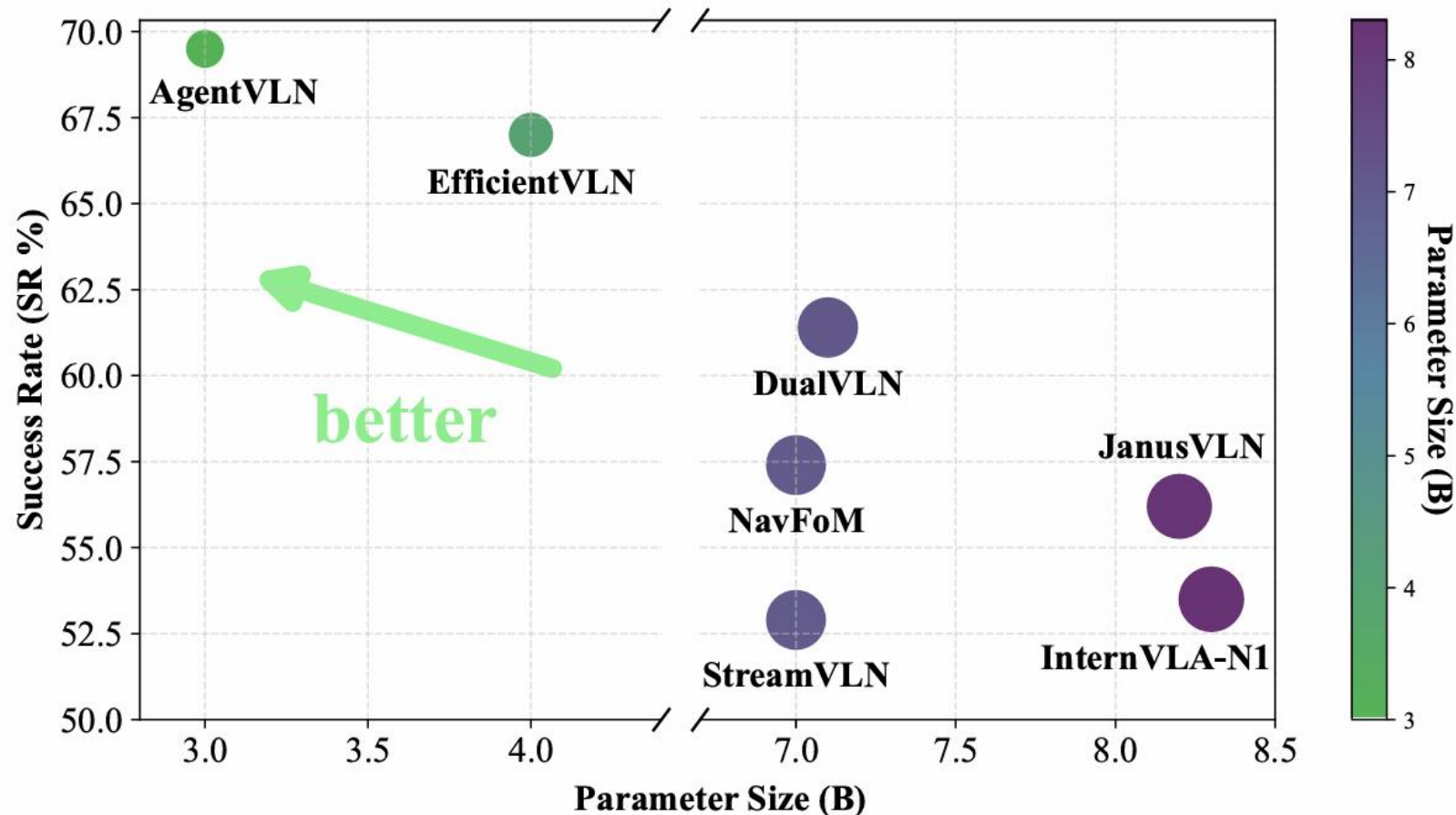
$$\mathbf{P}_{L,i} = R_W^L \mathbf{P}_{w,i} + \mathbf{t}_W^L$$

$$\mathbf{P}_{C,i} = R_L^C \mathbf{P}_{L,i} + \mathbf{t}_L^C$$

$$\hat{d} = \frac{\sum_{m=1}^k \omega_m \cdot d_m}{\sum_{m=1}^k \omega_m}$$

$$\omega_m = \frac{1}{\|\mathbf{P}_{target}^{img} - \mathbf{P}_m^{img}\|_2}$$

$$\mathbf{P}_{target}^C = \hat{d} \cdot K^{-1} \begin{bmatrix} u^* \\ v^* \\ 1 \end{bmatrix}$$



Inference peak memory:

AgentVLN: ~9G

DualVLN: ~20G

StreamVLN: ~18G

The first VLN model deployable on embedded devices.

Outstanding load balancing

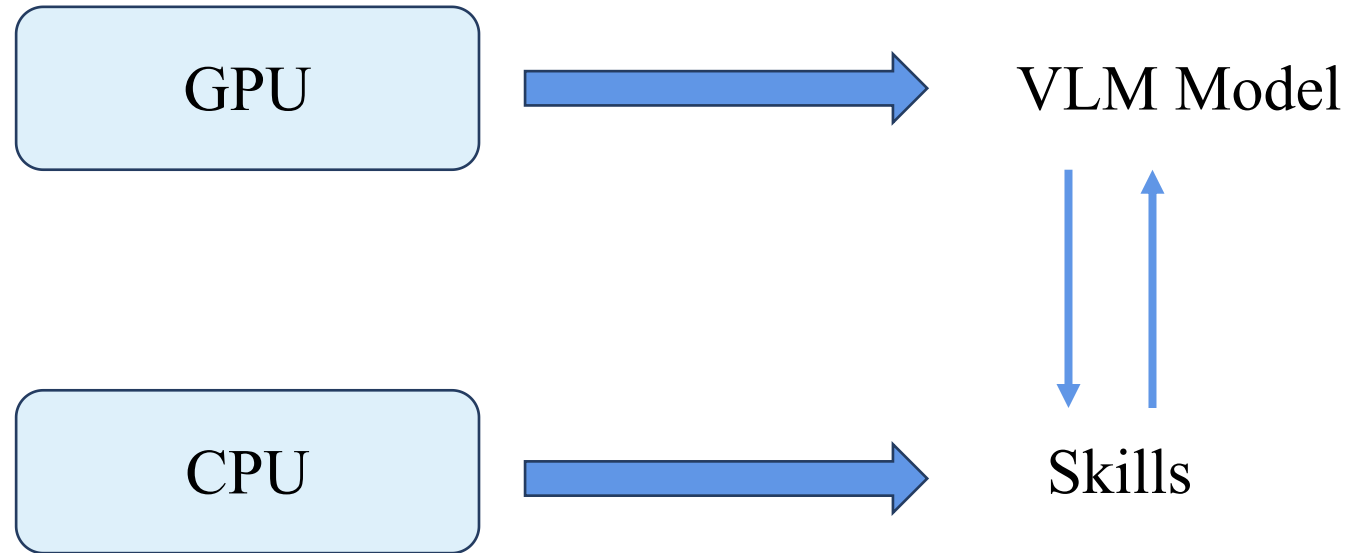




Table 1: Comparison results with SOTA methods on the Val-Unseen dataset for R2R-CE. Our method outperforms other approaches on the same benchmarks, even having fewer parameters.

Method	Observation			R2R-CE				
	S.RGB	Pano.	Depth	Odo.	NE ↓	OS ↑	SR ↑	SPL ↑
HPN+DN [20]	✓	✓	✓		6.31	40.0	36.0	34.0
VLN BERT [17]	✓	✓	✓		5.74	53.0	44.0	39.0
CMA [17]	✓	✓	✓		6.20	52.0	41.0	36.0
Reborn [2]	✓	✓	✓		5.40	57.0	50.0	46.0
Ego ² -Map [18]	✓	✓	✓		5.54	56.0	47.0	41.0
DreamWalker [39]	✓	✓	✓		5.53	59.0	49.0	44.0
ETPNav [1]	✓	✓	✓		4.71	65.0	57.0	49.0
Seq2Seq [21]	✓		✓		7.77	37.0	25.0	22.0
RGB-CMA [21]	✓				9.55	10.0	5.0	4.0
AG-CMTP [8]		✓	✓		7.90	39.0	23.0	19.0
R2R-CMTP [8]		✓	✓		7.90	38.0	26.0	22.0
LAW [35]	✓		✓	✓	6.83	44.0	35.0	31.0
CM2 [14]	✓		✓		7.02	41.0	34.0	27.0
WS-MGMap [9]	✓		✓		6.28	47.0	38.0	34.0
ETPNav+FF [42]	✓		✓		5.95	55.8	44.9	30.4
AO-Planner [7]		✓	✓		5.55	59.0	47.0	33.0
NaVid-7B [51]	✓				5.47	49.0	37.0	35.0
Uni-NaVid-7B [50]	✓				5.58	53.3	47.0	42.7
NaVILA-7B [12]	✓				5.22	62.5	54.0	49.0
StreamVLN-7B [44]	✓				5.10	64.0	55.7	50.9
NavFoM-7B [49]	✓				5.01	64.9	56.2	51.2
JanusVLN-8.2B [48]	✓				4.78	65.2	60.5	56.8
EfficientVLN-4B [53]	✓				4.18	73.7	64.2	55.9
DualVLN-7.1B [43]	✓				4.05	70.7	64.3	58.5
StreamVLN-7B [44]	✓		✓		4.98	64.2	56.9	51.9
InternVLA-N1-8.3B [38]	✓		✓		4.83	63.3	58.2	54.0
AgentVLN-3B	✓		✓		3.88	73.5	67.2	64.7

Table 2: Comparison with SOTA methods on RxR-CE Val-Unseen split.

Method	S.RGB	Pano.	Depth	Odo.	NE ↓	SR ↑	SPL ↑	nDTW ↑
VLN BERT [17]		✓	✓	✓	8.98	27.0	22.6	46.7
CMA [17]		✓	✓	✓	8.76	26.5	22.1	47.0
Reborn [2]		✓	✓	✓	5.98	48.6	42.0	63.3
ETPNav [1]		✓	✓	✓	5.64	54.7	44.8	61.9
Seq2Seq [21]	✓		✓		12.10	13.9	11.9	30.8
LAW [35]	✓		✓	✓	10.90	8.0	8.0	38.0
ETPNav+FF [42]	✓		✓	✓	8.79	25.5	18.1	-
AO-Planner [7]		✓	✓		7.06	43.3	30.5	50.1
Uni-NaVid-7B [50]	✓				6.24	48.7	40.9	-
NaVILA-7B [12]	✓				6.77	49.3	44.0	58.8
NavFoM-7B [49]	✓				5.51	57.4	49.4	60.2
JanusVLN-8.2B [47]	✓				6.06	56.2	47.5	62.1
StreamVLN-7B [44]	✓				6.16	51.8	45.0	61.9
InternVLA-N1-7B [13]	✓				6.41	49.5	41.8	62.6
EfficientVLN-4B [53]	✓				3.88	67.0	54.3	68.4
DualVLN-7.1B [43]	✓				4.58	61.4	51.8	70.0
StreamVLN-7B [44]	✓		✓		6.22	52.9	46.0	61.9
InternVLA-N1-8.3B [38]	✓		✓		5.91	53.5	46.1	65.3
AgentVLN-3B	✓		✓		3.92	69.5	61.3	74.6

Future Work: Infinite Memory

Model architecture specifically designed for VLN tasks.

Static Memory: $H^I = \text{EncodeOnce}(I), \quad M_I = K_I^{(l)}, V_I^{(l)}$

Dynamic Memory: $h_t^{loc} = \text{SWA}(o_{t-K:t}, a_{t-K:t-1}, p_{t-K:t})$

Memory Update: $e_t = \text{WriteMLP}(h_t^{loc}, a_{t-1}, p_t)$

$$S_t = \text{GatedDeltaNetWithDecay}(S_{t-1}, e_t)$$

$$h_t^{hist} = \text{Read}(S_t, q_t^{hist})$$

$$S_t = \rho_t \odot S_{t-1} (I - \beta_t k_t k_t^\top) + \beta_t v_t k_t^\top$$

Model architecture specifically designed for VLN tasks.

```
{ "episode_id": "2", "scene_id": "7y3sRwLe3Va", "step": 0, "phase": "decide", "parsed_task_type": "target", "parsed_coordinate": [306.0, 305.0], "parse_001012802124023, -1.448159098625183, -2.412627935409546], "map_update_ms": 2.46, "obs_prepare_ms": 13.466, "observe_step_ms": 176.639, "obs_total_ms": "decision_prepare_ms": 4.279, "generate_ms": 455.303, "parse_ms": 0.075, "pixel_to_world_ms": 0.26, "path_plan_ms": 0.337, "decision_total_ms": 465.57 "cuda_reserved_mb": 5654.0, "cuda_max_allocated_mb": 5551.118, "cuda_max_reserved_mb": 5654.0, "dist_to_goal": 5.2577, "step_total_ms": 679.601} {"episode_id": "2", "scene_id": "7y3sRwLe3Va", "step": 1, "phase": "observe", "parsed_task_type": null, "parsed_coordinate": null, "parsed_actions": r -1.448159098625183, -2.412627935409546], "map_update_ms": 3.301, "obs_prepare_ms": 13.595, "observe_step_ms": 175.784, "obs_total_ms": 192.761, "path_ "cuda_reserved_mb": 5654.0, "cuda_max_allocated_mb": 5551.12, "cuda_max_reserved_mb": 5654.0, "dist_to_goal": 5.4273, "step_total_ms": 196.336} {"episode_id": "2", "scene_id": "7y3sRwLe3Va", "step": 2, "phase": "decide", "parsed_task_type": "target", "parsed_coordinate": [306.0, 305.0], "parse_001012802124023, -1.448159098625183, -2.412627935409546], "map_update_ms": 2.47, "obs_prepare_ms": 11.079, "observe_step_ms": 175.178, "obs_total_ms": "decision_prepare_ms": 3.944, "generate_ms": 451.436, "parse_ms": 0.05, "pixel_to_world_ms": 0.21, "path_plan_ms": 0.313, "decision_total_ms": 461.165 "cuda_reserved_mb": 5654.0, "cuda_max_allocated_mb": 5551.12, "cuda_max_reserved_mb": 5654.0, "dist_to_goal": 5.6029, "step_total_ms": 669.09} {"episode_id": "2", "scene_id": "7y3sRwLe3Va", "step": 3, "phase": "observe", "parsed_task_type": null, "parsed_coordinate": null, "parsed_actions": r -1.448159098625183, -2.412627935409546], "map_update_ms": 2.354, "obs_prepare_ms": 11.537, "observe_step_ms": 176.91, "obs_total_ms": 190.877, "path_p "cuda_reserved_mb": 5654.0, "cuda_max_allocated_mb": 5551.12, "cuda_max_reserved_mb": 5654.0, "dist_to_goal": 5.6977, "step_total_ms": 194.531} {"episode_id": "2", "scene_id": "7y3sRwLe3Va", "step": 4, "phase": "decide", "parsed_task_type": "target", "parsed_coordinate": [306.0, 305.0], "parse_001012802124023, -1.448159098625183, -2.412627935409546], "map_update_ms": 3.092, "obs_prepare_ms": 11.846, "observe_step_ms": 179.001, "obs_total_ms" "decision_prepare_ms": 4.126, "generate_ms": 455.409, "parse_ms": 0.048, "pixel_to_world_ms": 0.225, "path_plan_ms": 0.321, "decision_total_ms": 465.0 "cuda_reserved_mb": 5654.0, "cuda_max_allocated_mb": 5551.12, "cuda_max_reserved_mb": 5654.0, "dist_to_goal": 5.718, "step_total_ms": 678.243} {"episode_id": "2", "scene_id": "7y3sRwLe3Va", "step": 5, "phase": "stop", "parsed_task_type": "stop", "parsed_coordinate": null, "parsed_actions": nu "map_update_ms": 2.299, "obs_prepare_ms": 12.811, "observe_step_ms": 176.479, "obs_total_ms": 191.665, "path_plan_ms": 0.426, "cuda_allocated_mb": 540
```

Sliding Attention + GDN Attention

Peak memory usage: 5.5G

Achieve 10 FPS

Inference Delay: 0.45s

No inference optimization was performed!