

StreamVLN: Streaming Vision-and-Language Navigation via SlowFast Context Modeling

Meng Wei^{*,1,2} Chenyang Wan^{*,1,3} Xiqian Yu^{*,1} Tai Wang^{*†,1} Yuqiang Yang¹
 Xiaohan Mao^{1,4} Chenming Zhu^{1,2} Wenzhe Cai¹ Hanqing Wang¹ Yilun Chen¹
 Xihui Liu^{‡,2} Jiangmiao Pang^{‡,1}

¹Shanghai AI Laboratory ²The University of Hong Kong
³Zhejiang University ⁴Shanghai Jiao Tong University

Abstract: Vision-and-Language Navigation (VLN) in real-world settings requires agents to process continuous visual streams and generate actions with low latency grounded in language instructions. While Video-based Large Language Models (Video-LLMs) have driven recent progress, current VLN methods based on Video-LLM often face trade-offs among fine-grained visual understanding, long-term context modeling and computational efficiency. We introduce StreamVLN, a streaming VLN framework that employs a hybrid slow-fast context modeling strategy to support multi-modal reasoning over interleaved vision, language and action inputs. The fast-streaming dialogue context facilitates responsive action generation through a sliding-window of active dialogues, while the slow-updating memory context compresses historical visual states using a 3D-aware token pruning strategy. With this slow-fast design, StreamVLN achieves coherent multi-turn dialogue through efficient KV cache reuse, supporting long video streams with bounded context size and inference cost. Experiments on VLN-CE benchmarks demonstrate state-of-the-art performance with stable low latency, ensuring robustness and efficiency in real-world deployment. The project page is: <https://streamvln.github.io/>.

Keywords: Visual-and-Language Navigation, Visual-Language-Action Model

1 Introduction

Vision-and-Language Navigation (VLN) in continuous real-world environments is a critical task in embodied AI, where an agent must ground linguistic cues in visual observations and plan actionable trajectories. However, achieving robust VLN remains challenging due to the need for fine-grained multimodal alignment, long-term sequence reasoning, and generalization to unseen environments. Recent advances in Video MultiModal Large Language Models [1, 2], (Video-LLMs) offer new capabilities for VLN systems. Several research efforts [3, 4, 5] have extended Video-LLMs to vision-language-action models (VLA) for navigation, which integrate visual observation encoding, language understanding, and action prediction in a unified end-to-end framework.

For real-world navigation, VLA models must process continuously incoming video streams, where maintaining long-term context and real-time responsiveness are both crucial. This poses challenges for current Video-LLMs in managing linearly growing visual tokens. Some methods [5, 6] address this by sampling a fixed number of video frames, but the limited temporal resolution may fail to accurately predict low-level actions when fine-grained temporal changes are needed. Other methods [3, 4] compress vision tokens into sparse memory tokens via pooling or token merging, which helps control the token volume but sacrificing temporal and visual details. Furthermore, these methods typically require refreshing the LLM’s dialogue context at every action step. This leads to significant redundant computation during both training and inference, hindering data scalability and real-world deployment.

^{*}Equal contribution.

[†]Project lead.

[‡]Corresponding authors.

In this paper, we propose StreamVLN, a novel streaming vision-and-language navigation framework for low-latency action generation. We extend the Video-LLM into an interleaved vision-language-action model, enabling continuous interaction with a video stream through multi-turn dialogue. To address the challenges of long-horizon context management and computational efficiency, StreamVLN introduces a hybrid strategy that combines a **fast-streaming dialogue context** and a **slow-updating memory context**. Specifically, it employs a sliding-window mechanism to cache the key/value states (KV) of tokens over a fixed number of dialogue turns for highly responsive action decoding. Meanwhile, it leverages the visual context from past windows as memory to enhance long-term temporal reasoning. To control memory growth, StreamVLN applies temporal sampling along with a test-time token pruning strategy that discards redundant tokens based on their 3D spatial proximity. This pruning approach can be performed without altering previously computed tokens, supporting efficient reuse of offloaded KV caches to achieve a stable decoding speed throughout navigation.

In summary, StreamVLN offers an efficient and scalable solution to suit the continuous interaction requirements of real-time navigation. Its slow-fast context modeling design enables the model trained on short clips (e.g., 16 frames), to work effectively on long video streams, without incurring context length growth or compromising inference latency. Experiments on existing VLN-CE benchmarks shows that StreamVLN achieves superior performance while maintaining low latency.

2 Related Work

Vision-and-Language Navigation (VLN). This task requires an agent to follow language instructions while perceiving and acting in environments. Early progress mainly focused on discrete settings, where agents navigate by “teleporting” between predefined nodes of a discrete scene graph [7, 8, 9, 10, 11, 12]. This formulation emphasizes high-level decision-making but ignores the challenges of real-world navigation. More recent work [13, 14, 15, 16, 17] has focused on continuous environments [18], where agents must perform low-level actions in realistic simulators. To address the increased complexity, some methods incorporate a waypoint predictor [16, 19, 20] pretrained in simulators to propose candidate positions, which are then used to guide high-level navigation decisions. Although these approaches have achieved strong performance, the waypoint predictors typically rely heavily on the training scenes and exhibit limited generalization to unseen scenes. Therefore, more flexible and scalable navigation frameworks is needed to generalize better to long-horizon, real-world setting.

Navigation with Multi-Modal Large Language Models (MLLMs). Recent advancements in MLLMs have opened new possibilities for VLN by enabling agents to interpret and reason over natural language instructions in a more generalizable way. Some methods [12, 21, 22, 23] directly use LLMs as planner in a training-free manner within a modular framework. But there’s still a performance gap compared to task-specific models. Other lines of work [3, 4, 5, 6] further fine-tune Video-based LLMs [2, 24, 1] to better capture spatial-temporal information and generate low-level actions in an end-to-end manner, but often face challenges in balancing computational efficiency and long-horizon memory retention. StreamVLN aims to better accommodate streaming video input by introducing an efficient and scalable framework that supports action generation with coherent multi-turn reasoning with low-latency response and bounded memory usage.

3 Method

StreamVLN generates action outputs from continuous video input in an online, multi-turn dialogue manner. Built on LLaVA-Video [2] as the foundational Video-LLM, we extend it for interleaved vision, language, and action modeling. The overall framework of StreamVLN is shown in Figure 1. We briefly introduce the autoregressive generation in continuous multi-turn dialogues for a streaming VLN process (Section 3.1). For both effective context modeling of long sequence and efficient computation for real-time interaction, StreamVLN has: (1) a fast-streaming dialogue context with a sliding-window KV cache (Section 3.2); and (2) a slow-updating memory via token pruning (Section 3.3). Finally, we describe how we curate the navigation data and incorporate diverse multimodal data for multi-task training (Section 3.4).

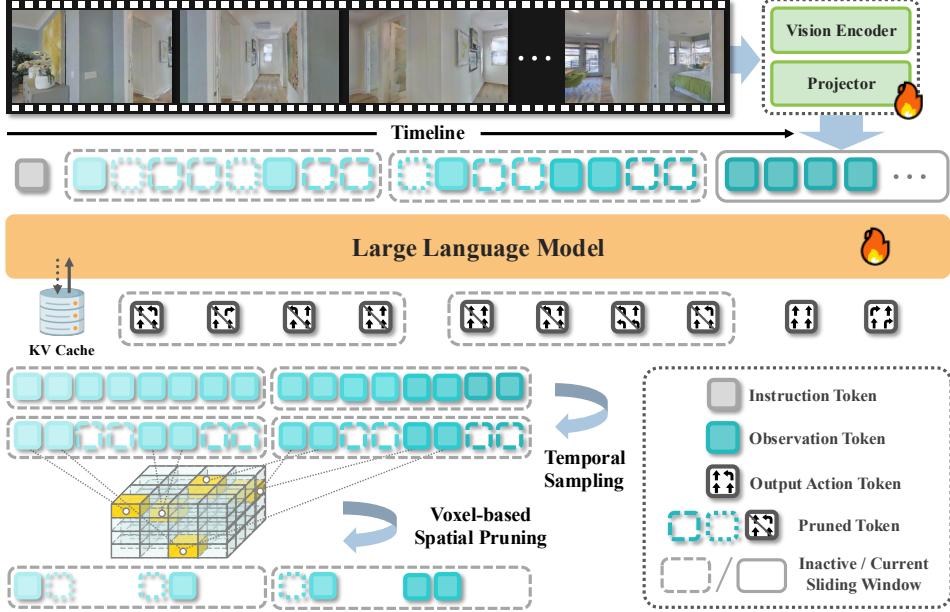


Figure 1: **Framework of StreamVLN.** The input consists of a language instruction and a stream of RGB images. Each navigation episode is framed as a multi-turn dialogue, where the agent continually queries for the next actions. To support long-horizon reasoning while maintaining a manageable context size and low latency, we adopt a fixed-size sliding window to retain recent dialogue history. The context in inactive windows is updated by token pruning to reduce memory overhead.

3.1 Preliminary: Continuous Multi-Turn Autoregressive Generation

A multi-turn dialogue session for VLN consists of a sequence of interleaved observations and actions. In each dialogue $d_i = (o_i, a_i)$, the VLN model receives a new observation o_i and produces an action response a_i conditioned on both the current input and the dialogue history. The full input sequence at step i is constructed as: $o_1 a_1 o_2 a_2 \dots o_{i-1} a_{i-1}$. In this streaming setting, new tokens from o_i are appended to the token stream continuously. The response a_i is generated token-by-token via autoregressive decoding. For each dialogue turn, Transformer-based LLMs first perform a **prefill phase** to encode input tokens, caching their key/value (KV) states in attention layers. These cached KV pairs are then used in the **decoding phase** to generate new tokens. If we don't use KV cache across turns, the model will repeat this prefilling process of all previous tokens for a new dialogue.

3.2 Fast-Streaming Dialogue Context

While multi-turn KV cache reuse can eliminate over 99% of prefilling time, it introduces substantial memory overhead. As the number of dialogues increases, the KV cache grows linearly (e.g., 2K tokens can consume around 5GB of memory), making long sessions impractical. In addition, existing Video-LLMs tend to exhibit degraded reasoning performance when processing overly long contexts.

To manage dialogue context, we adopt a sliding window KV cache over continuous dialogues, retaining a fixed number N of recent dialogues in an active window: $W_j = [o_{(i-N+1)} a_{(i-N+1)} \dots o_i a_i]$. When the window reaches capacity, the key/value states are offloaded from the LLM, and the states of non-observation dialogue tokens, such as prompts and generated actions, are immediately discarded. For the new sliding window, the token states from past windows are processed into memory token states $\{\mathcal{M}_0, \dots, \mathcal{M}_j\}$ (as detailed in Section 3.3). Formally, for the latest observation o_i , the decoder generates a_i based on the cached token states and the current window's KV cache:

$$a_i^{W_{j+1}} = \text{Decoder}(o_i, \{\mathcal{M}_0, \dots, \mathcal{M}_j\}, \{k_{(i-N+1)} v_{(i-N+1)}, \dots, k_{(i-1)} v_{(i-1)}\}).$$

3.3 Slow-Updating Memory Context

Balancing temporal resolution and fine-grained spatial perception within a limited context length remains a key challenge for Video-LLMs. Rather than compressing video tokens at the feature level (e.g., through average pooling), which hinders the reuse of the KV cache from previous dialogues, we retain high image resolution while selectively discarding spatially and temporally redundant tokens. We find that this approach better preserves the transferability of Video-LLMs.

To reduce the temporal redundancy, we adopt a simple fixed-number sampling strategy following [5], as varying lengths of memory tokens may induce a temporal duration bias, reduce the model’s robustness across different planning horizons. To further eliminate spatial redundancy across frames, we design a voxel-based spatial pruning strategy. Specifically, we back-project the 2D image patches from the video stream into a shared 3D space using depth information. By discretizing this 3D space into uniform voxels, we can track the voxel indices of the patch tokens over time. If multiple tokens from different frames within a given duration are projected into the same voxel, only the token from the most recent observation is retained, as detailed in Algorithm 1. The voxel pruning mask M is then used to select the preserved token states.

Algorithm 1 Voxel-Based Spatial Pruning

```

1: Voxel map  $V \in \mathbb{Z}^{T \times H \times W}$ , stride  $K$ , threshold  $\theta$ 
2: Pruning Mask  $M \in \{0, 1\}^{T \times H \times W}$ 
3: Initialize  $M \leftarrow \mathbf{0}$ , map  $\text{latest} \leftarrow \emptyset$ 
4: for each token  $(t, x, y)$  with  $V_{t,x,y} \geq 0$  do
5:    $p \leftarrow \lfloor t/K \rfloor$ ,  $v \leftarrow V_{t,x,y}$ 
6:   if  $(p, v)$  not in  $\text{latest}$  or  $t$  is newer then
7:      $\text{latest}[(p, v)] \leftarrow (t, x, y)$ 
8:   end if
9: end for
10: Set  $M_{t,x,y} \leftarrow 1$  for all  $(t, x, y) \in \text{latest}$ 
11: For each  $t$ , if  $\sum_{x,y} M_{t,x,y} < \theta \cdot H \cdot W$ , set  $M_{t,:,:} \leftarrow 0$ 
12: return  $M$ 
```

10: Set $M_{t,x,y} \leftarrow 1$ for all $(t, x, y) \in \text{latest}$
11: For each t , if $\sum_{x,y} M_{t,x,y} < \theta \cdot H \cdot W$, set $M_{t,:,:} \leftarrow 0$
12: **return** M

3.4 Co-Training with Multi-Source Data.

Vision-Language Action Data. We collect navigation-specific training data using the Habitat simulator across multiple public VLN datasets. Specifically, we collect 450K samples (video clips) from 60 Matterport3D [25] (MP3D) environments, sourced from R2R [7], R2R-EnvDrop [26] and RxR [8]. To further improve generalization through increased scene diversity, we incorporate an additional 300K samples from a subset of ScaleVLN [19], spanning 700 Habitat Matterport3D [27] (HM3D) scenes. In addition, we adopt the DAgger [28] algorithm to enhance the model’s robustness and generalization ability in novel scenes and during error recovery. Using Habitat’s shortest-path follower as the expert policy, we collect corrective demonstrations on model rollouts after the initial training stage. These DAgger-collected samples (240K) are then incorporated into the training set for co-training.

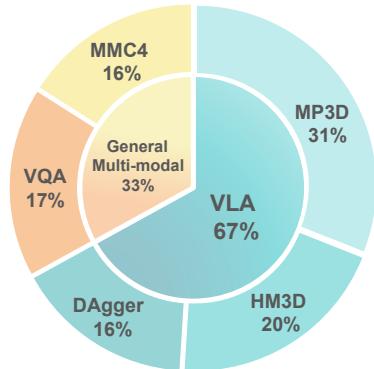


Figure 2: Co-Training Data Recipe of StreamVLN

General Vision-Language Data. To retain the general reasoning capabilities of the pretrained Video-LLM, we incorporate a diverse set of multimodal training data that complements navigation supervision. Specifically, we include 248K video-based visual question-answering (VQA) samples sourced from publicly available datasets LLaVA-Video-178K [29] and ScanQA [30], which combine general video QA with 3D scene understanding to support spatial-temporal and geometric reasoning. To further augment the model’s capacity for multi-turn vision-language interactions, we incorporate 230K interleaved image-text samples from MMC4 [31], which strengthens its ability to parse and generate contextually coherent responses with interleaved visual and textual reasoning.

4 Experiments

4.1 Experimental Setup

Simulation Benchmark Setup. We evaluate our method on two public VLN-CE [18] benchmarks collected from Matterport3D scenes using the Habitat simulator: R2R-CE [7] and RxR-CE [8]. R2R-CE provides 5.6K English trajectories with an average length of 10 meters, while RxR-CE includes 126K multilingual instructions (English, Hindi, Telugu) and features longer, more diverse paths (avg.15 meters). The camera HFOVs are 79° for R2R-CE and RxR-CE. Both benchmarks require realistic indoor navigation under continuous control. As our goal is to assess the generalization ability, we focus on the validation unseen splits of both benchmarks. We report standard VLN metrics, including Navigation Error (NE), Success Rate (SR), Oracle Success Rate (OS), and Success weighted by Path Length (SPL), following prior works.

Real-World Evaluation Setup. We perform real world experiments based on a Unitree Go2 robotic dog. The robot is equipped with a upward facing camera (Intel® RealSense™ D455) for RGB-D observations. We deploy StreamVLN on a remote workstation with an RTX 4090 GPU. The Go2 robot continuously streams visual data to the 4090 server for inference, which returns executable action commands to the robot. The averse inference (0.27s for 4 actions) and communication (0.2s for indoor and 1.0s for outdoor environments) latency enable real-time physical deployment.

4.2 Implementation Details

We build StreamVLN based on the LLaVA-Video [2] 7B model, which uses Qwen2-7B [32] as the language model. Training is conducted in two stages. First, we fine-tune it for one epoch solely on oracle VLN trajectories. Then, we use the model to collect DAgger trajectories and continue training for an additional epoch with a mixture of VLN and general multimodal data. During the warm-up phase, we apply a peak learning rate of 2e-5 for the language model and 5e-6 for the vision encoder. Each training step processes 128 video clips. Training completes in around 1500 A100 GPU hours.

| Method | Observation Encoder | | | | R2R Val-Unseen | | | | RxR Val-Unseen | | | |
|-------------------|---------------------|------|-------|-------|----------------|-------------|-------------|-------------|----------------|-------------|-------------|-------------|
| | Pano. | Odo. | Depth | S.RGB | NE↓ | OS↑ | SR↑ | SPL↑ | NE↓ | SR↑ | SPL↑ | nDTW↑ |
| HPN+DN* [33] | ✓ | ✓ | ✓ | | 6.31 | 40.0 | 36.0 | 34.0 | - | - | - | - |
| CMA* [34] | ✓ | ✓ | ✓ | | 6.20 | 52.0 | 41.0 | 36.0 | 8.76 | 26.5 | 22.1 | 47.0 |
| VLN○BERT* [34] | ✓ | ✓ | ✓ | | 5.74 | 53.0 | 44.0 | 39.0 | 8.98 | 27.0 | 22.6 | 46.7 |
| Sim2Sim* [35] | ✓ | ✓ | ✓ | | 6.07 | 52.0 | 43.0 | 36.0 | - | - | - | - |
| GridMM* [36] | ✓ | ✓ | ✓ | | 5.11 | 61.0 | 49.0 | 41.0 | - | - | - | - |
| ETPNav* [16] | ✓ | ✓ | ✓ | | 4.71 | 65.0 | 57.0 | 49.0 | 5.64 | 54.7 | 44.8 | 61.9 |
| ScaleVLN* [19] | ✓ | ✓ | ✓ | | 4.80 | - | 55.0 | 51.0 | - | - | - | - |
| InstructNav [23] | - | - | - | - | 6.89 | - | 31.0 | 24.0 | - | - | - | - |
| AG-CMTP [37] | ✓ | ✓ | ✓ | | 7.90 | 39.2 | 23.1 | 19.1 | - | - | - | - |
| R2R-CMTP [37] | ✓ | ✓ | ✓ | | 7.90 | 38.0 | 26.4 | 22.7 | - | - | - | - |
| LAW [13] | ✓ | ✓ | ✓ | ✓ | 6.83 | 44.0 | 35.0 | 31.0 | 10.90 | 8.0 | 8.0 | 38.0 |
| CM2 [14] | ✓ | ✓ | ✓ | | 7.02 | 41.5 | 34.3 | 27.6 | - | - | - | - |
| WS-MGMap [15] | ✓ | ✓ | ✓ | | 6.28 | 47.6 | 38.9 | 34.3 | - | - | - | - |
| ETPNav + FF [38] | ✓ | ✓ | ✓ | | 5.95 | 55.8 | 44.9 | 30.4 | 8.79 | 25.5 | 18.1 | - |
| Seq2Seq [39] | ✓ | ✓ | ✓ | | 7.77 | 37.0 | 25.0 | 22.0 | 12.10 | 13.9 | 11.9 | 30.8 |
| CMA [39] | ✓ | ✓ | ✓ | | 7.37 | 40.0 | 32.0 | 30.0 | - | - | - | - |
| NaVid [3] | | | | ✓ | 5.47 | 49.1 | 37.4 | 35.9 | - | - | - | - |
| MapNav [6] | | | | ✓ | 4.93 | 53.0 | 39.7 | 37.2 | - | - | - | - |
| NaVILA [5] | | | | ✓ | 5.37 | 57.6 | 49.7 | 45.5 | - | - | - | - |
| StreamVLN | | | | ✓ | 5.43 | 62.5 | 52.8 | 47.2 | 6.72 | 48.6 | 42.5 | 60.2 |
| NaVILA† [5] | | | | ✓ | 5.22 | 62.5 | 54.0 | 49.0 | 6.77 | 49.3 | 44.0 | 58.8 |
| UniNaVid† [4] | | | | ✓ | 5.58 | 53.3 | 47.0 | 42.7 | 6.24 | 48.7 | 40.9 | - |
| StreamVLN† | | | | ✓ | 4.98 | 64.2 | 56.9 | 51.9 | 6.22 | 52.9 | 46.0 | 61.9 |

Table 1: Comparison with state-of-the-art methods on VLN-CE R2R and RxR Val-Unseen split. * indicates methods using the waypoint predictor from [34]. † denotes methods using additional training data beyond the R2R-CE and RxR-CE benchmarks.

| Method | ScanQA | | | | |
|------------------------------|-------------|-------------|--------------|-------------|-------------|
| | Bleu-4 ↑ | Rouge ↑ | Meteor ↑ | Cider ↑ | EM ↑ |
| ScanRefer [40] | 7.9 | 30.0 | 55.4 | 11.5 | 18.6 |
| ScanQA [30] | 10.1 | 33.3 | 64.9 | 13.1 | 21.0 |
| 3D-VisTA [41] | 10.4 | 35.7 | 69.6 | 13.9 | 22.4 |
| 3D-LLM* [42] | 12.0 | 35.7 | 69.4 | 14.5 | 20.5 |
| LEO [43] | 13.2 | 49.2 | 101.4 | 20.0 | 24.5 |
| ChatScene* [44] | 14.0 | - | 87.6 | - | - |
| Scene-LLM* [45] | 12.0 | 40.0 | 80.0 | 16.6 | 27.2 |
| NaviLLM [46] | 12.0 | 38.4 | 75.9 | 15.4 | 23.0 |
| NaVILA (16 frames) [5] | 15.2 | 48.3 | 99.8 | 19.6 | 27.4 |
| StreamVLN (16 frames) | 15.7 | 48.3 | 100.2 | 19.8 | 28.8 |

Table 2: Comparison on ScanQA [30] Val set. * indicates 3D LLMs with task-specific fine-tuning.

4.3 Comparisons with State-of-the-Arts

Results on VLN-CE benchmark. Table 1 summarizes the performance of our method on the VLN-CE R2R and RxR benchmarks under the Val-Unseen setting, compared with existing VLN-CE methods. Our StreamVLN model achieves state-of-the-art performance among RGB-only methods both without and with extra navigation datasets, reaching 56.9% SR and 51.9% SPL on R2R, and 52.9% SR and 46.0% SPL on RxR. These results highlight the generality and robustness of our approach across both standard and long-horizon navigation tasks. Notably, StreamVLN performs comparably to ETPNav [16], despite not relying on additional panoramic or waypoint supervision. Furthermore, compared to HMAT trained on the large-scale ScaleVLN dataset with 3 million trajectories, StreamVLN surpasses it using only a small subset of ScaleVLN (150k trajectories), demonstrating superior data efficiency.

Results on Video Question Answering. To evaluate StreamVLN’s spatial scene understanding capabilities, we conduct experiments on the widely-used ScanQA benchmark for 3D question answering based on real-world scans. StreamVLN answers questions by analyzing 16 multi-view images from each scan. As shown in Table 2, StreamVLN outperforms state-of-the-art generalist navigation models such as NaviLLM [46] and NaVILA [5]. As shown in Figure 3, we observe that the strong VQA capabilities contribute to better generalization to novel navigation instructions. **Real-**

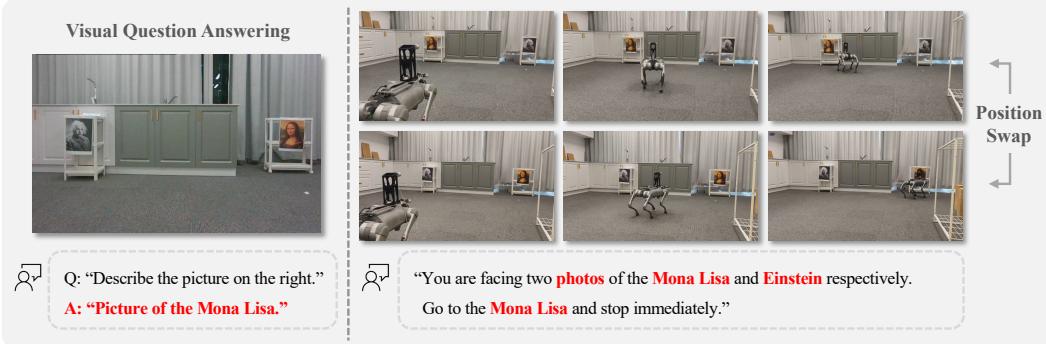


Figure 3: StreamVLN transfers visual reasoning ability to interpreting navigation instructions.

World Qualitative Results. We select several representative real-world settings—**Workspace**, **Mall**, and **Outdoor**. Partial qualitative results are presented in Figure 4. To further assess its robustness, we construct a more challenging task in the **Workspace** environment (first panel), involving multiple landmarks and varying illumination. StreamVLN completes this difficult long-horizon navigation successfully. Moreover, success cases in **Mall** and **Outdoor** environments highlight StreamVLN’s strong generalization to novel scenes and tasks. Throughout all settings, StreamVLN maintains efficient inference, supporting smooth real-world deployment. Please refer to the demo video for full demonstrations <https://streamvln.github.io/>.

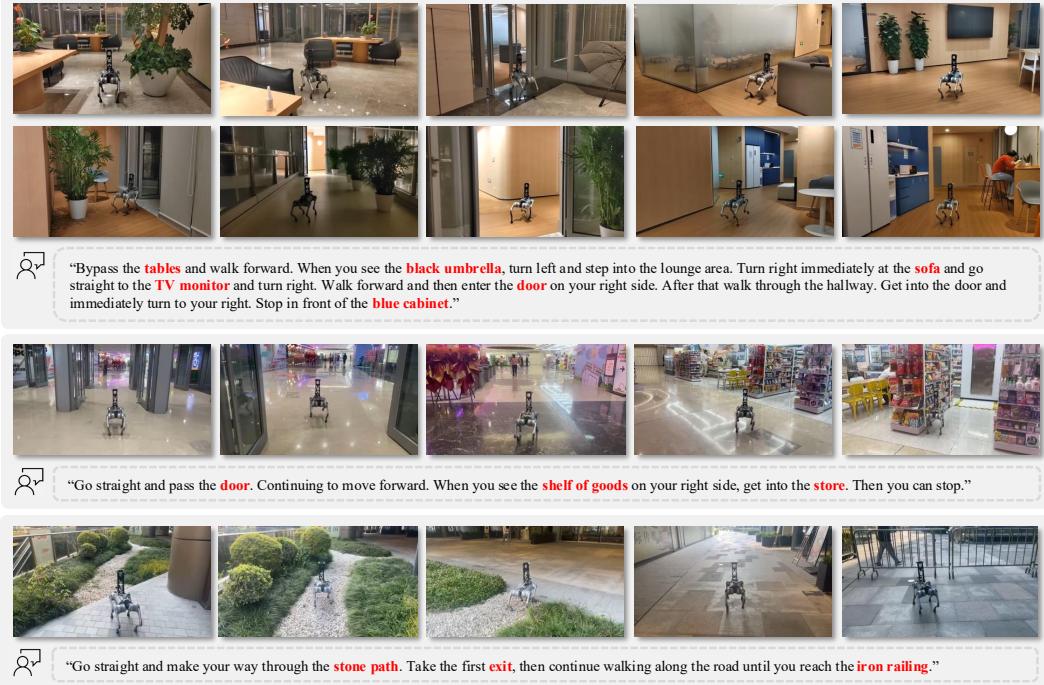


Figure 4: Qualitative results of StreamVLN in several representative real-world environments. From top to bottom are Home, Workspace, Mall and Outdoor. StreamVLN achieves robust performance across diverse VLN scenarios, capable of accurately following complex instructions with various landmarks (marked as red) and handling real-world disturbances.

| R2R | RxR | DAgger | VL Data | ScaleVLN(150K) | NE↓ | OS↑ | SR↑ | SPL↑ |
|-----|-----|--------|--------------|----------------|-------------|-------------|-------------|-------------|
| ✓ | ✓ | | | | 6.05 | 53.8 | 45.5 | 41.6 |
| ✓ | ✓ | ✓ | VideoQA | | 5.47 | 57.8 | 50.8 | 45.7 |
| ✓ | ✓ | ✓ | VideoQA+MMC4 | | 5.43 | 62.5 | 52.8 | 47.2 |
| ✓ | ✓ | ✓ | VideoQA+MMC4 | ✓ | 5.10 | 64.0 | 55.7 | 50.9 |
| ✓ | ✓ | | VideoQA+MMC4 | ✓ | 5.73 | 56.4 | 50.2 | 47.1 |
| ✓ | | ✓ | VideoQA+MMC4 | ✓ | 5.90 | 55.9 | 47.9 | 43.6 |

Table 3: Ablation study of different training data compositions on VLN-CE R2R Val-Unseen split.

4.4 Ablation Studies

Data Ablation. Table 3 presents an ablation study on different training data compositions. All results are reported without using voxel-based spatial pruning. The first row shows the first-stage performance when training with only oracle navigation data. After collecting DAgger data, we co-train oracle data, DAgger data, and vision-language (VL) data. In the second row, we use only VideoQA data as VL data. While the third row mixes VideoQA and MMC4 (M) data in an interleaved image-text format. For a fair comparison, the total number of VL Data is kept the same. We can observe that the second-stage co-training brings significant gains (+5.3 SR / +4.1 SPL) and incorporating MMC4 further improves performance (+2.0 SR / +1.5 SPL). Comparing the third and fourth rows, we see that adding ScaleVLN data brings additional gains (+2.9 SR / +3.7 SPL). To assess the importance of DAgger data, we remove it from the co-training data, as shown in the fifth row. The results show that DAgger data plays a crucial role in boosting performance (+5.5 SR / +3.8 SPL). Furthermore, the last row highlights that incorporating RxR data yields notable performance gains (+7.8 SR / +7.3 SPL).

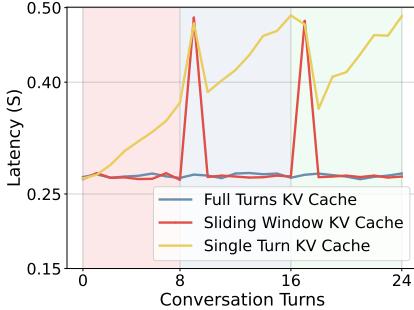


Figure 5: Impact of KV cache reuse in multiple turns. Sliding window size:8.

| Memory | Window | NE↓ | OS↑ | SR↑ | SPL↑ |
|------------|--------|-------------|-------------|-------------|-------------|
| 2*196 | 8 | 6.96 | 48.2 | 37.3 | 34.2 |
| 4*196 | 8 | 6.62 | 49.1 | 38.9 | 35.4 |
| 8*196 | 8 | 6.05 | 53.8 | 45.5 | 41.6 |
| <i>all</i> | 8 | 6.76 | 49.5 | 40.0 | 36.4 |
| 8*196 | 4 | 6.31 | 51.1 | 41.4 | 37.5 |
| 8*196 | 2 | 6.16 | 52.8 | 43.7 | 40.3 |

Table 4: Ablation on the impact of different memory context sizes and sliding window sizes on VLN-CE R2R Val-Unseen split.

| Pruning | R2R | | | | RxR | | | |
|---------|------|------|------|------|------|------|------|-------|
| | NE↓ | OS↑ | SR↑ | SPL↑ | NE↓ | SR↑ | SPL↑ | nDTW↑ |
| X | 5.10 | 64.0 | 55.7 | 50.9 | 6.16 | 51.8 | 45.0 | 62.1 |
| ≈20% | 4.98 | 64.2 | 56.9 | 51.9 | 6.22 | 52.9 | 46.0 | 61.9 |

Table 5: Ablation on voxel-based spatial pruning on R2R, RxR val-unseen split and HM3D val split.

Memory Context Size. We study the impact of the memory context size in the hybrid context modeling strategy. As shown in Table 4 (results shown are from first-stage training using only oracle VLN data), increasing the memory size from $2 * 196$ to $8 * 196$ while keeping the window size fixed at 8 significantly improves navigation performance, with SR rising from 37.3 to 45.5. This indicates the importance of fine-grained memory in supporting long-horizon reasoning. Notably, using the entire visual context (*all*) as memory doesn’t yield the best results, suggesting that an overly long and varied context token sequence may introduce bias in training and hinder generalization at test time.

Sliding Window Size. We also evaluate the effect of the number of dialogue turns retained in the sliding window in Table 4. A smaller window size leads to more frequent shifts, resulting in a significantly larger number of training samples. For example, a window size of 8 yields approximately 450K samples, while sizes of 4 and 2 increase this to 815K and 1.5M respectively. This growth not only raises the training cost linearly but also introduces greater class imbalance, which may affect training stability. We find that retaining 8 continuous dialogue turns achieves the best balance—delivering strong navigation performance while maintaining the lowest training cost.

Effectiveness of KV-Cache Reuse. We evaluate the impact of KV cache reuse on the decoding latency under different settings. As shown in Figure 5, reusing the KV cache across all dialogue turns (Full Turns) achieves consistently low latency, since only the current observation tokens require prefill computation for generating the 8 action tokens. If the KV cache is maintained only within 8 turns (Sliding Window), the decoding latency will increase at the beginning of each sliding window due to the need to prefill the previous window context tokens. Under the Single Turn setting, where the KV cache is not reused across turns (as in prior work), decoding latency steadily increases with the number of turns. Turns 0–8 incur lower latency since no historical context is included, while turns 8–16 and 16–24 has similar latency growth with a fixed memory size.

Effectiveness of Voxel-Based Spatial Pruning. Table 5 shows the effect of applying voxel-based spatial pruning during inference across different tasks. We are able to reduce the number of input tokens by approximately 20% on average across the three datasets, leading to more efficient encoding of the scene information. Despite the reduction in input tokens, voxel pruning consistently improves navigation performance. Specifically, we observe a +1.2% SR and +1.0% SPL gain on R2R, a +1.1% SR and +1.0% SPL gain on RxR, which means proper token pruning can help model focus on relevant tokens and thus improves performance.

5 Conclusion

This paper presents *StreamVLN*, a new streaming vision-language-navigation framework based on Video-LLMs. Compared to previous Video-LLM-based VLN methods that treat each interaction as an independent dialogue and refresh history at every step, StreamVLN can reuse past key/value (KV) states through a hybrid memory design. By maintaining a fast-updating sliding window for immediate responsiveness and a slow-updating long-term memory for temporal reasoning, StreamVLN enables efficient, coherent, and scalable action generation over long video streams. Empirical results on standard VLN-CE benchmarks demonstrate that StreamVLN achieves superior performance with lower latency, paving the way for real-time, memory-efficient, and long context-aware navigation.

6 Limitations

While StreamVLN enables efficient navigation through streaming dialogue with a multi-modal Video LLM, it also faces several limitations. First, directly generating low-level actions from raw visual observations remains less robust to variations in viewpoint and occlusion, potentially leading to suboptimal control in real-world environments. Second, the current hybrid context modeling strategy still encounters challenges in longer-horizon navigation scenarios, where maintaining consistent reasoning over extended sequences is difficult. Third, the reliance on explicit action history as part of the dialogue context introduces additional complexity for asynchronous inference and deployment, as it requires synchronizing past actions to preserve dialogue coherence.

References

- [1] Y. Li, C. Wang, and J. Jia. Llama-vid: An image is worth 2 tokens in large language models. In *European Conference on Computer Vision*, 2024.
- [2] Y. Zhang, J. Wu, W. Li, B. Li, Z. Ma, Z. Liu, and C. Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024.
- [3] J. Zhang, K. Wang, R. Xu, G. Zhou, Y. Hong, X. Fang, Q. Wu, Z. Zhang, and H. Wang. Navid: Video-based vlm plans the next step for vision-and-language navigation. *Robotics: Science and Systems*, 2024.
- [4] J. Zhang, K. Wang, S. Wang, M. Li, H. Liu, S. Wei, Z. Wang, Z. Zhang, and H. Wang. Uni-navid: A video-based vision-language-action model for unifying embodied navigation tasks. *Robotics: Science and Systems*, 2025.
- [5] A.-C. Cheng, Y. Ji, Z. Yang, Z. Gongye, X. Zou, J. Kautz, E. Biyik, H. Yin, S. Liu, and X. Wang. Navila: Legged robot vision-language-action model for navigation. *Robotics: Science and Systems*, 2025.
- [6] L. Zhang, X. Hao, Q. Xu, Q. Zhang, X. Zhang, P. Wang, J. Zhang, Z. Wang, S. Zhang, and R. Xu. Mapnav: A novel memory representation via annotated semantic maps for vlm-based vision-and-language navigation. *arXiv preprint arXiv:2502.13451*, 2025.
- [7] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and A. Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
- [8] A. Ku, P. Anderson, R. Patel, E. Ie, and J. Baldridge. Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. *arXiv preprint arXiv:2010.07954*, 2020.
- [9] S. Chen, P.-L. Guhur, C. Schmid, and I. Laptev. History aware multimodal transformer for vision-and-language navigation. *Advances in Neural Information Processing Systems*, 2021.

- [10] S. Chen, P.-L. Guhur, M. Tapaswi, C. Schmid, and I. Laptev. Think global, act local: Dual-scale graph transformer for vision-and-language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16537–16547, 2022.
- [11] H. Wang, W. Wang, W. Liang, C. Xiong, and J. Shen. Structured scene memory for vision-language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [12] G. Zhou, Y. Hong, and Q. Wu. Navgpt: Explicit reasoning in vision-and-language navigation with large language models. *arXiv preprint arXiv:2305.16986*, 2023.
- [13] S. Raychaudhuri, S. Wani, S. Patel, U. Jain, and A. X. Chang. Language-aligned waypoint (law) supervision for vision-and-language navigation in continuous environments. *arXiv preprint arXiv:2109.15207*, 2021.
- [14] G. Georgakis, K. Schmeckpeper, K. Wanchoo, S. Dan, E. Miltsakaki, D. Roth, and K. Daniilidis. Cross-modal map learning for vision and language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [15] P. Chen, D. Ji, K. Lin, R. Zeng, T. H. Li, M. Tan, and C. Gan. Weakly-supervised multi-granularity map learning for vision-and-language navigation. *arXiv preprint arXiv:2210.07506*, 2022.
- [16] D. An, H. Wang, W. Wang, Z. Wang, Y. Huang, K. He, and L. Wang. Etpnav: Evolving topological planning for vision-language navigation in continuous environments. *arXiv preprint arXiv:2304.03047*, 2023.
- [17] H. Wang, W. Liang, L. Van Gool, and W. Wang. Dreamwalker: Mental planning for continuous vision-language navigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.
- [18] J. Krantz, E. Wijmans, A. Majundar, D. Batra, and S. Lee. Beyond the nav-graph: Vision and language navigation in continuous environments. In *European Conference on Computer Vision (ECCV)*, 2020.
- [19] Z. Wang, J. Li, Y. Hong, Y. Wang, Q. Wu, M. Bansal, S. Gould, H. Tan, and Y. Qiao. Scaling data generation in vision-and-language navigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
- [20] Z. Wang, X. Li, J. Yang, Y. Liu, J. Hu, M. Jiang, and S. Jiang. Lookahead exploration with neural radiance representation for continuous vision-language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [21] Y. Long, X. Li, W. Cai, and H. Dong. Discuss before moving: Visual language navigation via multi-expert discussions. *arXiv preprint arXiv:2309.11382*, 2023.
- [22] P. Chen, X. Sun, H. Zhi, R. Zeng, T. H. Li, G. Liu, M. Tan, and C. Gan. a^2 nav: Action-aware zero-shot robot navigation by exploiting vision-and-language ability of foundation models. *arXiv preprint arXiv:2308.07997*, 2023.
- [23] Y. Long, W. Cai, H. Wang, G. Zhan, and H. Dong. Instructnav: Zero-shot system for generic instruction navigation in unexplored environment. *arXiv preprint arXiv:2406.04882*, 2024.
- [24] J. Lin, H. Yin, W. Ping, P. Molchanov, M. Shoeybi, and S. Han. Vila: On pre-training for visual language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024.
- [25] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang. Matterport3D: Learning from RGB-D data in indoor environments. *International Conference on 3D Vision (3DV)*, 2017.

- [26] H. Tan, L. Yu, and M. Bansal. Learning to navigate unseen environments: Back translation with environmental dropout. In *Proceedings of NAACL-HLT*, 2019.
- [27] S. K. Ramakrishnan, A. Gokaslan, E. Wijmans, O. Maksymets, A. Clegg, J. Turner, E. Underander, W. Galuba, A. Westbury, A. X. Chang, et al. Habitat-matterport 3d dataset (hm3d): 1000 large-scale 3d environments for embodied ai. *arXiv preprint arXiv:2109.08238*, 2021.
- [28] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, 2011.
- [29] Y. Zhang, J. Wu, W. Li, B. Li, Z. Ma, Z. Liu, and C. Li. Video instruction tuning with synthetic data, 2024. URL <https://arxiv.org/abs/2410.02713>.
- [30] D. Azuma, T. Miyanishi, S. Kurita, and M. Kawanabe. Scanqa: 3d question answering for spatial scene understanding. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022.
- [31] W. Zhu, J. Hessel, A. Awadalla, S. Y. Gadre, J. Dodge, A. Fang, Y. Yu, L. Schmidt, W. Y. Wang, and Y. Choi. Multimodal C4: An open, billion-scale corpus of images interleaved with text. *arXiv preprint arXiv:2304.06939*, 2023.
- [32] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [33] J. Krantz, A. Gokaslan, D. Batra, S. Lee, and O. Maksymets. Waypoint models for instruction-guided navigation in continuous environments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.
- [34] Y. Hong, Z. Wang, Q. Wu, and S. Gould. Bridging the gap between learning in discrete and continuous environments for vision-and-language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [35] J. Krantz and S. Lee. Sim-2-sim transfer for vision-and-language navigation in continuous environments. In *European Conference on Computer Vision*, 2022.
- [36] Z. Wang, X. Li, J. Yang, Y. Liu, and S. Jiang. Gridmm: Grid memory map for vision-and-language navigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
- [37] K. Chen, J. K. Chen, J. Chuang, M. Vázquez, and S. Savarese. Topological planning with transformers for vision-and-language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [38] Z. Wang, X. Li, J. Yang, Y. Liu, and S. Jiang. Sim-to-real transfer via 3d feature fields for vision-and-language navigation. *arXiv preprint arXiv:2406.09798*, 2024.
- [39] J. Krantz, E. Wijmans, A. Majumdar, D. Batra, and S. Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *European Conference on Computer Vision*, 2020.
- [40] Z. Yu, J. Yu, Y. Cui, D. Tao, and Q. Tian. Deep modular co-attention networks for visual question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019.
- [41] Z. Zhu, X. Ma, Y. Chen, Z. Deng, S. Huang, and Q. Li. 3d-vista: Pre-trained transformer for 3d vision and text alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2911–2921, 2023.

- [42] Y. Hong, H. Zhen, P. Chen, S. Zheng, Y. Du, Z. Chen, and C. Gan. 3d-llm: Injecting the 3d world into large language models. *Advances in Neural Information Processing Systems*, 2023.
- [43] J. Huang, S. Yong, X. Ma, X. Linghu, P. Li, Y. Wang, Q. Li, S.-C. Zhu, B. Jia, and S. Huang. An embodied generalist agent in 3d world. *arXiv preprint arXiv:2311.12871*, 2023.
- [44] H. Huang, Y. Chen, Z. Wang, R. Huang, R. Xu, T. Wang, L. Liu, X. Cheng, Y. Zhao, J. Pang, et al. Chat-scene: Bridging 3d scene and large language models with object identifiers. *arXiv preprint arXiv:2312.08168*, 2023.
- [45] R. Fu, J. Liu, X. Chen, Y. Nie, and W. Xiong. Scene-llm: Extending language model for 3d visual understanding and reasoning. *arXiv preprint arXiv:2403.11401*, 2024.
- [46] D. Zheng, S. Huang, L. Zhao, Y. Zhong, and L. Wang. Towards learning a generalist model for embodied navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [47] S. Y. Gadre, M. Wortsman, G. Ilharco, L. Schmidt, and S. Song. Cows on pasture: Baselines and benchmarks for language-driven zero-shot object navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- [48] K. Zhou, K. Zheng, C. Pryor, Y. Shen, H. Jin, L. Getoor, and X. E. Wang. Esc: Exploration with soft commonsense constraints for zero-shot object navigation. *arXiv preprint arXiv:2301.13166*, 2023.
- [49] W. Cai, S. Huang, G. Cheng, Y. Long, P. Gao, C. Sun, and H. Dong. Bridging zero-shot object navigation and foundation models through pixel-guided navigation skill. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024.
- [50] P. Wu, Y. Mu, B. Wu, Y. Hou, J. Ma, S. Zhang, and C. Liu. Voronav: Voronoi-based zero-shot object navigation with large language model, 2024.

StreamVLN: Streaming Vision-and-Language Navigation via SlowFast Context Modeling

Supplementary Material

1 Navigation Instruction Design

2 **Interleaved Observation and Action Format.** StreamVLN is an interleave-VLA framework, in
3 which the system alternates between perceiving the environment and generating actions conditioned
4 on those observations. As demonstrated in the example below, this interleaving of observations
5 and actions ensures that the agent operates in a dynamic, responsive manner. Such a design is
6 critical for simulating real-world decision-making, where an autonomous agent continuously receives
7 environmental feedback and adapts its behavior in real time. Each interaction session consists of 8
8 dialogue turns, resulting in a total of 32 actions executed sequentially without interruption. When a
9 session reaches the maximum number of turns, the observation states are recorded as $\langle\text{memory}\rangle$. A
10 historical prompt is inserted in all subsequent sessions after the first session.

A session of multi-turn dialogue

User: You are an autonomous navigation assistant. Your task is to $\langle\text{Instruction}\rangle$. Devise an action sequence to follow the instruction using the four actions: TURN LEFT (\leftarrow) or TURN RIGHT (\rightarrow) by 15 degrees, MOVE FORWARD (\uparrow) by 25 centimeters, or STOP.

[These are your historical observations: $\langle\text{memory}\rangle$.]

You can spot $\langle\text{image}\rangle$.

Assistant: $\uparrow\uparrow\uparrow$

User: You can see $\langle\text{image}\rangle$.

Assistant: $\leftarrow\leftarrow\uparrow\uparrow$

User: There is $\langle\text{image}\rangle$.

Assistant: $\rightarrow\rightarrow\rightarrow\uparrow$

User: You are toward the $\langle\text{image}\rangle$.

Assistant: $\uparrow\uparrow\leftarrow\uparrow$

User: You can see $\langle\text{image}\rangle$.

Assistant: $\uparrow\uparrow\leftarrow\uparrow$

User: In your sight is $\langle\text{image}\rangle$.

Assistant: $\uparrow\uparrow\uparrow\uparrow$

User: In front of you is $\langle\text{image}\rangle$.

Assistant: $\uparrow\leftarrow\leftarrow\uparrow$

User: Ahead of you is $\langle\text{image}\rangle$.

Assistant: $\uparrow\leftarrow\uparrow\uparrow$

...

11

12 **Ablation Study on Action Tokens** We conduct an ablation to assess how different action token
13 designs impact navigation performance and inference efficiency. Since textual representations affect
14 token count after tokenization, longer sequences increase latency and computational load. To analyze
15 this trade-off, we compare our design with UniNaVid [4] and NaVILA [5]. UniNaVid uses a compact
16 four-token vocabulary (STOP, FORWARD, LEFT, RIGHT), each mapped to a single token, while
17 NaVILA employs natural language phrases like “move forward 25 cm,” resulting in about 10 tokens.
18 StreamVLN also adopts a one-token-per-action design for efficiency, but uses rare symbols ($\uparrow\leftarrow\rightarrow$)
19 to reduce overfitting to common tokens in the vocabulary. We train StreamVLN’s first stage using
20 each of the three action token schemes. As shown in Table A1, symbolic tokens (StreamVLN)
21 slightly outperform UniNaVid’s textual tokens, but fall short of the performance achieved with
22 NaVILA’s more expressive natural language commands. StreamVLN thus strikes a balance between
23 expressiveness and token efficiency.

| Action Type | Token Number | Generate Time(s) | R2R Val-Unseen | | | |
|---------------------------------|--------------|------------------|----------------|-------------|-------------|-------------|
| | | | NE↓ | OS↑ | SR↑ | SPL↑ |
| forward/left/right/stop | 4 | 0.27 | 6.25 | 52.2 | 44.4 | 41.0 |
| move forward 25cm/turn left ... | 23 | 1.00 | 5.74 | 55.4 | 47.2 | 42.9 |
| ↑←→stop | 4 | 0.27 | 6.05 | 53.8 | 45.5 | 41.6 |

Table A1: Ablation on the impact of different action types on the VLN-CE R2R Val-Unseen split, based on the first-stage training of StreamVLN. The token number denotes the number of generated tokens across the four actions. The generation time refers to the time required for a single call to `model.generate()` to produce the tokens for all four actions (including the system tokens such as `<im_start>`, `<im_end>`, “assistant”).

| Metric | Method | | | | | | |
|--------|----------|----------|---------------|--------------|-----------|------------|-------------|
| | GoW [47] | ESC [48] | PixelNav [49] | VoroNav [50] | Navid [3] | MapNav [6] | StreamVLN |
| SR ↑ | 32.0 | 39.2 | 37.9 | 42.0 | 32.5 | 34.6 | 35.8 |
| SPL ↑ | 18.1 | 22.3 | 20.5 | 26.0 | 21.5 | 25.6 | 25.7 |

Table A2: Comparison with Zero-Shot ObjectNav methods on HM3D validation set.

2 Results on Zero-Shot Object Navigation

We evaluate StreamVLN’s zero-shot generalization on the HM3D-ObjectNav validation set. When the model outputs the stop action at a location where the distance to the goal object is less than 1 meter and the target object is visible, success is achieved. As shown in Table A2, we compare against two types of methods: LLM-based reasoning methods (GoW [47], ESC [48], PixelNav [49], and VoroNav [50]) and end-to-end vision-language-action models (Navid [3], MapNav [6]). StreamVLN demonstrates cross-task and cross-dataset generalization without any ObjectNav-specific training.

3 Voxel-based Spatial Pruning.

At the end of each session, we uniformly sample a fixed number of frames (8) to store in memory. However, uniform sampling introduces temporal and spatial redundancy. To mitigate this, we apply a voxel-based pruning strategy at test time to reduce redundant memory tokens and improve efficiency. We observe that applying pruning during training hurts performance, likely due to reduced input diversity and information. As shown in Figure A1, voxel pruning is more aggressive early in navigation, removing spatially redundant tokens to compress memory, and becomes less aggressive over time, with memory updates relying more on temporal sampling.

4 More Qualitative Results

We present extensive qualitative results of StreamVLN in both real-world and simulated environments to further demonstrate its effectiveness and generalization ability. Specifically, in Figure A2 we showcase six additional successful navigation episodes across the Bedroom, Working Space, and Outdoor environments. These examples highlight the model’s ability to interpret instructions, reason over long-horizon trajectories, and adapt to diverse scene layouts and object configurations. Moreover, we observe that StreamVLN demonstrates notable discriminative capabilities in navigation, effectively distinguishing between various goals and adapting its decisions accordingly, as illustrated in the first case in Figure A2. In this example, the model not only correctly identifies the content depicted in the painting on the right through VQA dialogue, but also selects the appropriate goal location and successfully navigates toward it. Please refer to the corresponding demo video of these real-world qualitative results for more detailed visualizations and analysis. We also provide qualitative results in the simulated environment, as shown in Figure A3 (R2R-CE) and Figure A4 (RxR-CE).



Figure A1: Visualization of voxel pruning on 4 of the 8 sampled frames per session, showing how redundant memory tokens are dynamically removed.

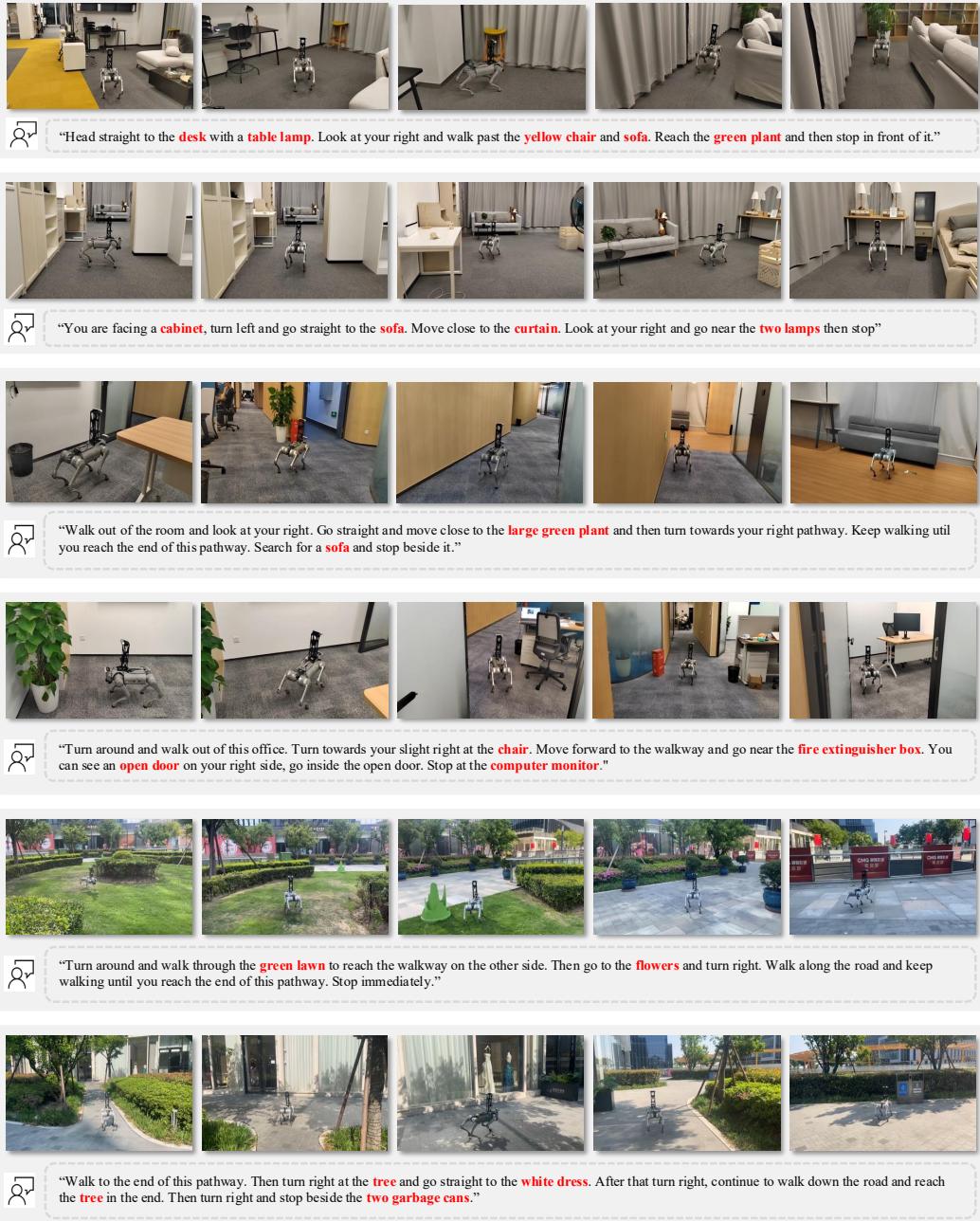


Figure A2: Additional real-world qualitative results of StreamVLN showcasing its strong navigation capabilities across diverse and novel environments. From top to bottom, each row shows: (a) Living room, (b) Bedroom, (c) Workspace, (d) Office, (e) Park, and (f) Walkway.

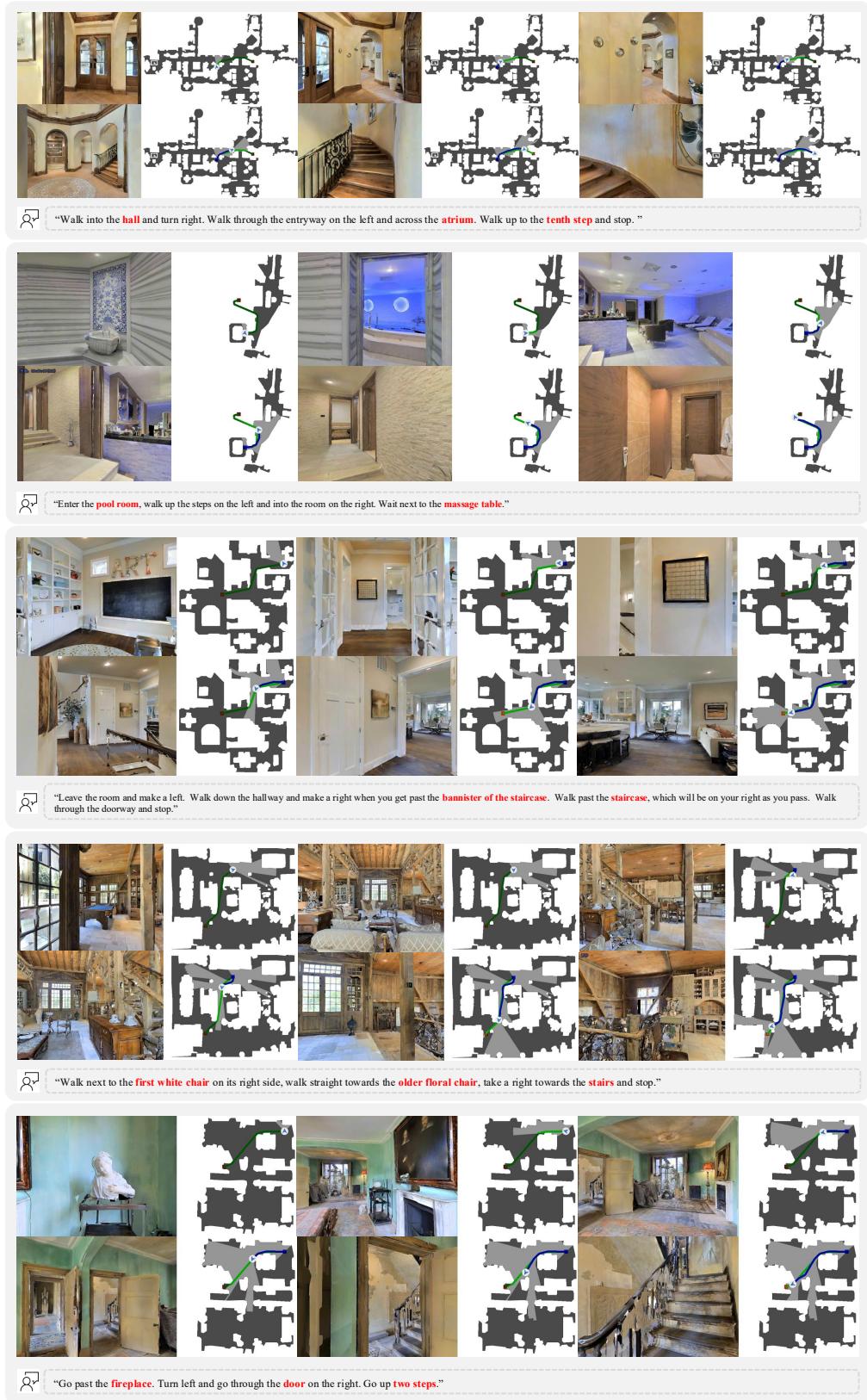


Figure A3: Qualitative results of StreamVLN on R2R-CE.

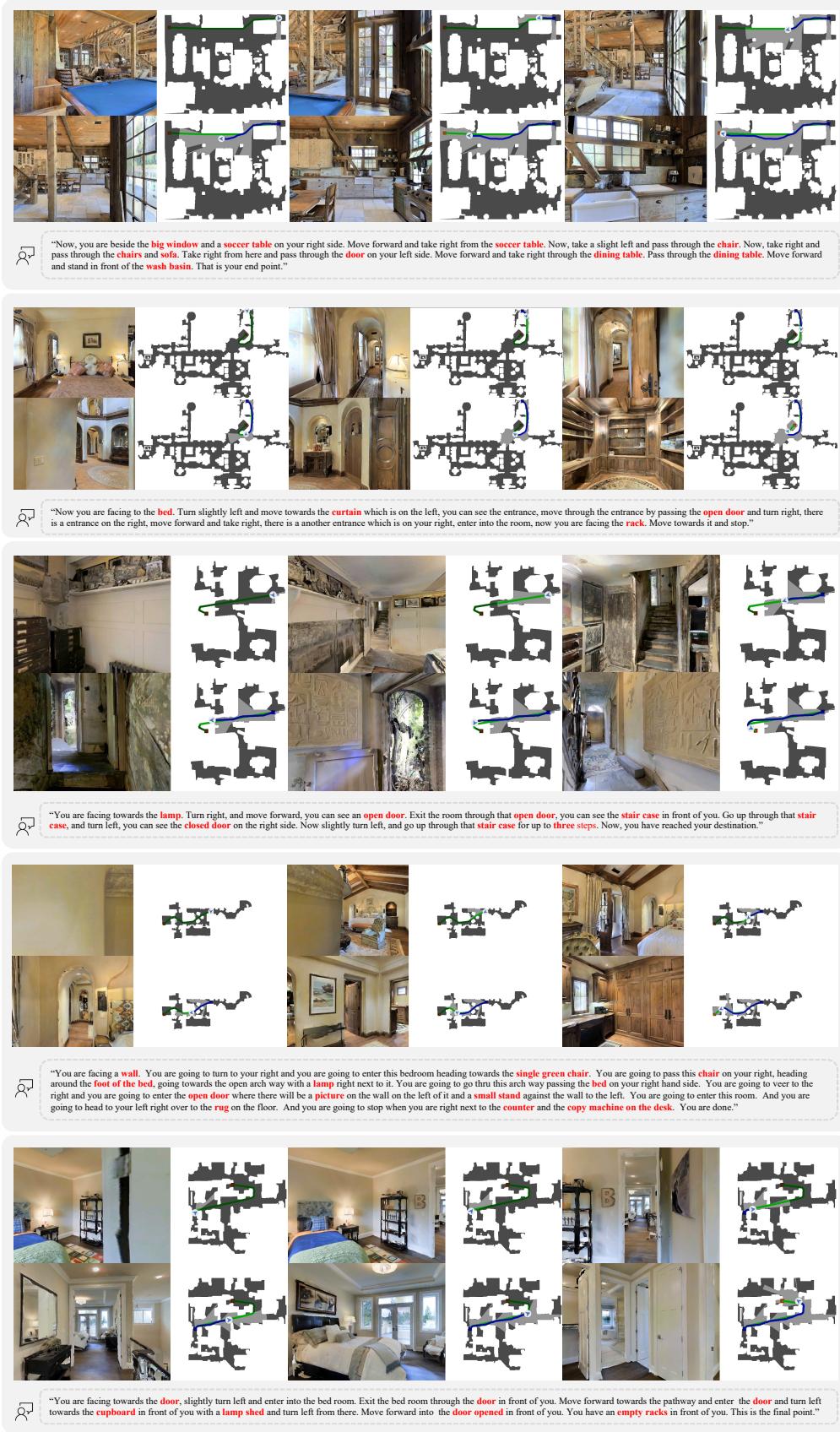


Figure A4: Qualitative results of StreamVLN on RxR-CE.