# Generative Proto-Sequence: Sequence-Level Decision Making for Long-Horizon Reinforcement Learning

Anonymous authors
Paper under double-blind review

# **Abstract**

Deep reinforcement learning (DRL) methods often face challenges in environments characterized by large state spaces, long action horizons, and sparse rewards, where effective exploration and credit assignment are critical. We introduce Generative Proto-Sequence (GPS), a novel generative DRL approach that produces variable-length discrete action sequences. By generating entire action sequences in a single decision rather than selecting individual actions at each timestep, GPS reduces the temporal decision bottleneck that impedes learning in long-horizon tasks. This sequence-level abstraction provides three key advantages: (1) it facilitates more effective credit assignment by directly connecting state observations with the outcomes of complete behavioral patterns; (2) by committing to coherent multi-step strategies, our approach facilitates better exploration of the state space; and (3) it promotes better generalization by learning macro-behaviors that transfer across similar situations rather than memorizing state-specific responses. Extensive evaluations on mazes of varying sizes and complexities demonstrate that GPS consistently outperforms leading action repetition and temporal methods, where it converges faster and achieves higher success rates across all environments.

# 1 Introduction

Deep reinforcement learning (DRL) has demonstrated impressive performance across diverse applications (Mnih et al., 2015; Silver et al., 2016; Levine et al., 2016). However, significant challenges remain when dealing with environments characterized by large state spaces, long-range tasks, and sparse rewards. In such contexts, traditional DRL methods that select actions sequentially often suffer from inefficient exploration and poor credit assignment (Mesnard et al., 2020; Raileanu & Rocktäschel, 2020; Ecoffet et al., 2021), leading to difficulties in learning effective policies for tasks that require coordinated, multi-step strategies. These challenges are further exacerbated by sparse reward signals, whose limited feedback hinders the agent's ability to discover and reinforce successful behaviors (Arjona-Medina et al., 2019; Hung et al., 2019).

Recent research efforts have attempted to address these challenges using diverse techniques such as hierarchical DRL (Kulkarni et al., 2016; Xu et al., 2022), temporal abstraction (Biedenkapp et al., 2021; Zhang et al., 2022b; Saanum et al., 2023; Patel & Siegelmann, 2024), sequence modeling (Chen et al., 2021; Janner et al., 2021; Giladi & Katz, 2023), and action repetition strategies (Srinivas et al., 2017; Sharma et al., 2017; Dabney et al., 2020). By creating sub-tasks or capturing higher-level behavioral patterns (Rosete-Beas et al., 2023; Vezzani et al., 2022; Wang et al., 2023), these techniques aim to reduce decision frequency and enhance learning efficiency in long-horizon tasks. Although these methods offer promising results, they often require careful sub-task design (Ajay et al., 2023), reward shaping (Liu et al., 2022), or complex training procedures (Seo & Abbeel, 2024b; Rosete-Beas et al., 2023). While there are temporal abstraction methods that generate multi-step action sequences, they often depend on iterative rollouts, autoregressive decoding, or model-based simulations. These solutions introduce computational overhead and restrict the ability to efficiently generate diverse action sequences. To our knowledge, no method supports the generation of coherent, variable-length action sequences directly from state observations in a single decision step.

In this study, we propose Generative Proto-Sequence (GPS), a novel actor-critic architecture capable of producing variable-length action sequences. Instead of actions, our Actor generates a *proto-sequence* embedding, which is then decoded into a discrete action sequence using a Decoder component. The Critic evaluates the state and the entire generated sequence jointly, with gradients flowing from the critic through the Decoder to the Actor, facilitating end-to-end learning of strategic, multi-step action sequences. This design enables the agent to generate and execute complex exploratory behaviors in a single decision, enhancing both generalization and long-horizon credit assignment.

We evaluated GPS on a large set of challenging maze environments with varying sizes and configurations, including rooms, corridors, and randomly generated obstacles. Our results demonstrate that GPS consistently learns more efficiently, generalizes better to novel maze layouts, and significantly outperforms leading baselines in terms of success rate and convergence speed, particularly in large and complex mazes. Our contributions are as follows:

- We introduce a novel architecture that enables end-to-end generation and evaluation of variable-length discrete action sequences, facilitating improved credit assignment and exploration.
- We demonstrate that producing multi-step action sequences in a single shot leads to superior generalization and faster convergence, particularly in large and complex environments.
- We provide extensive empirical results on challenging maze benchmarks, showing significant improvements over top-performing action repetition and temporal methods baselines in metrics such as convergence speed and success rate.

## 2 Related Work

# 2.1 Temporal Abstraction Through Action Repetition

Early works in temporal abstraction explored repeating single primitive actions to extend decision horizons. Recent research in DRL has produced various approaches for performing multiple actions as a single block. Earlier works (Srinivas et al., 2017; Sharma et al., 2017) introduced frameworks for dynamic action selection and repetition, though their repetition policies operated independently from chosen actions, limiting strategic development.

DAR (Srinivas et al., 2017) augments discrete action spaces by duplicating each base action with multiple repetition rates. While this expansion can improve learning in environments benefiting from temporal abstraction, it produces an inefficient representation—duplicated actions are treated as unrelated, preventing the agent from exploiting their shared underlying behavior and leading to slower learning and imbalanced trade-offs between coarse and fine control. FiGAR (Sharma et al., 2017) addresses this by decoupling behavior and repetition into two jointly trained policies; however, the repetition policy operates independently from the chosen action, limiting the development of nuanced, action-specific repetition strategies.

The authors of Dabney et al. (2020) proposed an exploration strategy repeating actions for random durations to reduce inefficient dithering. Temporl (Biedenkapp et al., 2021) advanced this by enabling agents to determine both action and repetition duration, improving learning efficiency. However, its hierarchical structure artificially decouples action selection from duration determination. Despite showing promise, these studies share a limitation: temporal abstraction is achieved solely through simple repetition of primitive actions, without generating coherent, variable-length action sequences.

# 2.2 Multi-Step Action Sequence Generation

Beyond single-action repetition, several methods focus on generating and partially committing to multi-step action sequences. The authors of Zhang et al. (2022a) introduced a generative planning method (GPM) that produces multi-step plans. Since GPM is trained by maximizing value, the plans generated from it can be regarded as intentional action sequences to reach high-value states and improve sample efficiency. PrAC (Coad et al., 2022) enables agents to generate n-step plans and commit to them while being predictable, balancing adaptability and control stability. The work of Saanum et al. (2023) incentivizes compressible

action sequences by integrating sequence priors, while Patel & Siegelmann (2024) introduced a model-based sequence RL framework (SRL) reducing decision frequency through action chunking.

Despite recent progress, most existing methods for generating multi-step action sequences still face major limitations. Many rely on heavy processes such as iterative rollouts, autoregressive decoding, or model-based simulation, which can be slow and inflexible (Li et al., 2024; Li, 2023; Zhang et al., 2025). For example, methods like PrAC and SRL use learned environment models for both planning and training, adding extra model-based complexity (Kumar et al., 2024; Luo et al., 2024). To stay adaptable, some approaches also use external switching mechanisms or mid-sequence re-planning, as seen in GPM and PrAC. This treats long-term planning as an add-on to a step-by-step framework rather than as a core design principle. As a result, sequence generation and evaluation are often optimized separately, which can lead to poor credit assignment (Dai et al., 2018). One case is the use of handcrafted regularization, such as rewarding shorter or more "compressible" sequences (Saanum et al., 2023). However, when objectives are split in this way, it becomes unclear whether failures come from a bad plan or from breaking the secondary constraint, making end-to-end training harder and reducing stability during execution.

# 2.3 Temporal Abstraction Using Hierarchies and LLMs

Hierarchical methods have advanced multi-action decision-making through skill discovery and sequencing. TACO-RL (Rosete-Beas et al., 2023) learns latent skills from unstructured data for long-horizon tasks. ASPiRe (Xu et al., 2022) accelerated RL by combining specialized skill priors. The work of Vezzani et al. (2022) introduced a skill scheduler sequencing pretrained skills, while SHRL (Wang et al., 2023) combined high-level policies with low-level skills for visual navigation. These approaches improve temporal abstraction by leveraging reusable skills rather than primitive actions.

Recent works have leveraged large language models and value-based reinforcement learning methods with action discretization for action sequence generation. CQN-AS (Seo & Abbeel, 2024a) proposed a value-based algorithm learning precise value functions from noisy action sequences. AlphaMaze (Dao & Vu, 2025) improved LLMs' spatial reasoning by combining supervised fine-tuning with policy optimization.

Our work draws inspiration from Dulac-Arnold et al. (2015), who generated embedding representations of proto-sequences mapped to discrete actions. In GPS, we propose key improvements: our approach is fully differentiable and trainable end-to-end, unlike Dulac-Arnold et al. (2015) whose k-nearest neighbors mapping broke the computation graph. Additionally, by using a VAE-based decoder instead of clustering, we automatically create sequence representations without manual embedding design. This enables efficient generation of coherent, variable-length action sequences that extend beyond simple repetition or skill sequencing.

# 3 Method

**Overview.** Our proposed approach is presented in Figure 1. GPS consists of three components: Actor,  $Proto-Sequence\ Decoder\ (PSD)$ , and Critic. The Actor receives the current state as input, and produces a  $proto-sequence\ -$  an embedding-based representation of a sequence of actions. The PSD receives the proto-sequence as input, and translates it into a discrete set of actions (e.g.,  $a_t, a_{t+1}, ..., a_{t+L}$ ), which are then executed sequentially by the agent. Finally, our Critic receives the sequence and predicts the expected cumulative reward obtained from its execution.

GPS differs from previous studies in several important aspects. First, unlike previous studies (Dulac-Arnold et al., 2015), it is end-to-end differentiable and does not require training workarounds. Secondly, our VAE decoder produces more diverse and flexible action sequences than autoregressive or model-based approaches, and also does so in an efficient, one-shot manner. Thirdly, sequence generation and evaluation are learned jointly, without regularization or switching mechanisms, thus improving credit assignment. Finally, by committing to the entire sequence (unlike the frequent re-evaluation of Zhang et al. (2022a)) we reduce execution overhead and increase behavioral predictability by forcing GPS to learn robust policies.

Another important aspect of our proposed approach is its ability to generate action sequences that differ from those on which it was trained. By creating novel sequences, GPS does not simply "memorize" a fixed

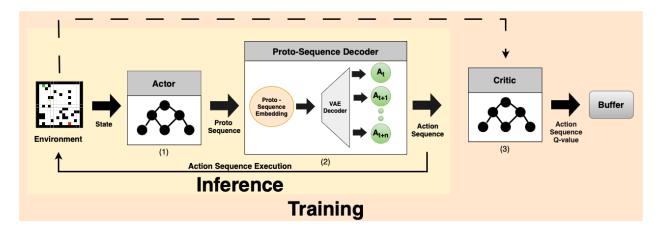


Figure 1: The three components of our proposed approach: (1) The Actor encodes the current state to produce a proto-sequence embedding. (2) The Decoder translates this latent embedding into a variable-length discrete action sequence. (3) The Critic evaluates the state-action-sequence pair and assigns it a Q-value representing the expected cumulative reward (3). During inference, only the actor and decoder components are used.

set of actions, but is able to generalize to larger action spaces. We elaborate on GPS's capacity to produce novel sequences in Section 5.1. Despite using a frozen decoder pre-trained on a limited set of sequences, GPS can generate action sequences that differ from those in the decoder's training data. The ability to generalize is important to our method, since sequence memorization is not applicable to domains with large action spaces. We elaborate on GPS's capacity to produce novel sequences in Section 5.1.

## 3.1 The Actor

The Actor serves as GPS's policy network. Given a state  $s_t$ , the Actor analyzes the input and outputs a proto-sequence embedding  $k = \pi_{\theta^{\pi}}(s_t)$ , where  $\theta^{\pi}$  and  $\pi$  are the parameters of the Actor's neural network and the current policy, respectively. The proto-sequence k is a latent embedding of a sequence of actions, represented as a vector in the embedding space  $k \in \mathbb{R}^d$ , where d is the dimensionality of the embedding space. This representation provides our Actor with significant flexibility, as it can create action sequences of varying length using a fixed-size representation.

The proto-sequence is next used by the PSD to produce a discrete sequence of actions, and this sequence is evaluated by the Critic (Section 3.3). The parameters  $\theta^{\pi}$  of the Actor are then updated using an actor-critic approach analogous to the Deep Deterministic Policy Gradient (DDPG) algorithm (Lillicrap et al., 2016), leveraging the learning signal provided by the Critic. Specifically, the actor's parameters  $\theta^{\pi}$  are adjusted to produce proto-sequence embeddings k that maximize the expected cumulative reward estimated by the critic,  $Q_{\theta Q}(s_t, \mathbf{a})$ . This optimization is achieved by updating  $\theta^{\pi}$  to minimize the negative Q-value provided by the critic  $-Q_{\theta Q}(s_t, g_{\theta^{\omega}}(\pi_{\theta^{\pi}}(s_t)))$ , using backpropagated gradients from the output of the critic network  $Q_{\theta^{Q}}$ . These gradients pass through the decoder network  $g_{\theta^{\omega}}$  and subsequently through the actor network  $\pi_{\theta^{\pi}}$ , enabling the update of the latter's parameters  $\theta^{\pi}$ .

# 3.2 The Proto-Sequence Decoder

The goal of the PSD is to translate the latent proto-sequence k generated by the Actor into a sequence of executable actions in the original action space  $\{a_t, a_{t+1}, ..., a_{t+L}\} \in A$ . We define the PSD as a function  $g_{\theta^{\omega}}: K \to A'^{L_{max}}$ , parameterized by  $\theta^{\omega}$ , where A' extends A with an EOS token to handle variable-length sequences within a fixed-length format, padding shorter sequences as needed. This function maps from the latent proto-sequence space K to sequences of fixed length  $L_{max}$ .

We use a Variational Autoencoder (VAE) (Kingma & Welling, 2013) as our PSD. We train the architecture on a diverse set of synthetic action sequences of varying lengths. For detailed information on the generation process of these sequences, see Appendix M.1. After training, we discard the VAE's encoder and retain only the learned decoder network  $g_{\theta^{\omega}}$ . The decoder is integrated into our agent architecture, transforming the Actor's latent proto-sequence embeddings into sequences of discrete actions. GPS will then execute the full sequence, without changes or early stopping. We chose VAE for its efficiency, ability to generate complete sequences in a single step, and its structured latent space that enables smooth interpolation and principled probabilistic modeling.

We pre-train the PSD and keep its parameters fixed while jointly training the actor and critic. We had several reasons for this separation: a) simplicity and robustness – separately training the Decoder reduces the number of moving parts in our architecture, enabling faster training. Separate training also prevents a 'moving target' scenario, where multiple components adapt at the same time, compromising training stability Sutton et al. (1999); b) Diversity and prevention of mode collapse: by training the Decoder separately, we have full control of its training set. We use a diverse training set that ensures that the Decoder can generate action sequences of all types. In addition, by freezing the Decoder after its training, we prevent a scenario where this component "forgets" how to generate specific actions during the other components' training; c) Transferability – by ensuring our Decoder can generate action sequences of all types, we can train it once and then use it in all tasks with this particular action space. It should be noted that training GPS end-to-end achieves comparable final results, with some differences in path optimality and efficiency (see Section G in the Appendix).

## 3.3 The Critic

The goal of our Critic is similar to the role of the critic in an actor-critic architecture. The Critic receives the current state  $s_t$  and the one-hot encoded discrete actions sequence  $\mathbf{A} = (a_t, a_{t+1}, \dots, a_{t+L})$  produced by the PSD. It then attempts to predict  $Q_{\theta Q}(s_t, \mathbf{A})$ , which represents the cumulative discounted reward obtained by executing  $\mathbf{A}$  and following the policy after the end of the sequence:

$$Q_{\theta Q}(s_t, \mathbf{A}) \approx \mathbb{E}_{\pi, P} \left[ \sum_{k=0}^{L-1} \gamma^k r_{t+k} + \gamma^L V^{\pi}(s_{t+L}) \right]$$

where  $V^{\pi}$  is the value function under policy  $\pi$ ,  $\theta^{Q}$  are the Critic's parameters, and  $L = eff\_len(\mathbf{A})$  denotes the effective length of the action sequence  $\mathbf{A}$ .

The Critic's parameters  $\theta^Q$  are updated by minimizing the Mean Squared Error (MSE) loss against a Temporal Difference (TD) target  $y_t$ :

$$L(\theta^Q) = \mathbb{E}_{(s_t, \mathbf{A}, \text{rewards}, s_{\text{next}})} \left[ (Q(s_t, \mathbf{A}; \theta^Q) - y_t)^2 \right]$$

The target  $y_t$  is constructed from the sum of discounted rewards  $R_t(\mathbf{A})$  obtained by executing sequence  $\mathbf{A}$ , and the discounted value of the subsequent state  $s_{t+L}$ , estimated using target Actor (Actor<sub>target</sub>) and target Critic ( $Q_{\text{target}}$ ) networks:

$$y_t = R_t(\mathbf{A}) + \gamma^L Q_{\text{target}}(s_{t+L}, \text{PSD}(\text{Actor}_{\text{target}}(s_{t+L})); \theta^{Q-})$$

This update mechanism, which relies on TD errors and target networks, is characteristic of many actor-critic algorithms, and shares similarities with methods such as DDPG (Lillicrap et al., 2016). While the Critic learns to accurately predict  $Q(s_t, \mathbf{A}; \theta^Q)$ , the Actor is trained to produce proto-sequences that, when decoded by the PSD, maximize this predicted Q-value.

## 3.4 Training Set Augmentation Using Sequence Subsets and Inference

To enhance learning efficiency and improve credit assignment, our training procedure leverages reward information from subsequences of each executed action sequence. For each sequence  $\mathbf{A} = (a_t, \dots, a_{t+L})$  of length L, we extract transitions corresponding to multiple contiguous subsequences  $(a_i, \dots, a_j)$  where

Table 1: The setup and properties of the mazes used in the evaluation.										
Environment	Dist. from	Train	Train Optimal	Val Set	Val Optimal	Test Set	Test Optimal			
Environment	start to goal	Set Size	Avg. Path	Size	Avg. Path	Size	Avg. Path			
8x8	[1 - 14]	100	5.14	100	5.49	1000	5.31			
16x16	[16 - 26]	100	18.04	100	18.0	1000	17.98			
$16x16\_obstacles\_15\%$	[20 - 30]	100	21.02	100	21.35	210	21.31			
$16x16\_obstacles\_25\%$	[20 - 30]	100	21.63	100	21.34	400	21.54			
16x16_rooms	[20 - 30]	100	20.93	100	21.01	585	21.02			
16x16_corridors	[10 - 30]	100	12.84	100	12.76	545	13.14			
24x24	[20 - 30]	100	23.39	100	23.26	1000	23.56			
$24x24\_obstacles\_15\%$	[10 - 20]	100	15.04	100	14.58	1000	14.73			
24x24 obstacles 25%	[10 - 20]	100	15.8	100	15.05	1000	15.11			

Table 1: The setup and properties of the mazes used in the evaluation.

 $t \leq i < j < t + L$ . For each such subsequence starting from an intermediate state  $s_i$ , we calculate the accumulated discounted reward obtained during its execution. This process effectively generates multiple learning samples of varying temporal lengths from a single interaction sequence, enriching the training data.

The subsequence extraction strategies for these state-subsequence-reward tuples, which we add to the replay buffer, include two primary approaches: (1) prefix extraction, which fixes the starting state while varying the end point, and (2) suffix extraction, which fixes the goal state while varying the starting point. This bidirectional approach diversifies the replay buffer with different time scales and enables the Critic to learn value estimates  $Q_{\theta Q}(s_i, (a_i, \dots, a_j))$  for sequences of different lengths concurrently. As shown in our analysis in Section 5.2, these extraction strategies significantly accelerate learning and improve overall performance. It is important to note that during inference (test time), our architecture does not utilize the Critic component, since no training takes place. Instead, the Actor and PSD produce the action sequence, and the latter is executed in full.

# 4 Experiments and Results

## 4.1 Evaluation Environment

Mazes are a foundational benchmark in DRL research, commonly used to evaluate an agent's ability to perform complex sequential decision-making and navigation tasks. Their structured yet variable environments provide a controlled setting for evaluating generalization, exploration, and memory, which are central to DRL performance (Pašukonis et al., 2023). We use four types of mazes in our evaluation:

- Empty. These mazes have no walls or obstacles, except for their boundaries.
- Sparse Obstacles. This setup has randomly placed obstacles in K% of the cells of each maze (e.g., 15%).
- Rooms. This setup consists of four large rooms with small doors between them. We also add randomly placed obstacles in 5% of open cells.
- Corridors. These mazes have only narrow corridors for the agent to navigate.

Our evaluation uses fully observable MDPs where the agent has complete visibility of the entire maze grid, including its position, goal location, and all obstacles.

Similarly to Dao & Vu (2025), we use an LLM to produce the code used in our maze generation. Our code, as well as the mazes generated for our evaluation, are available in the appendix. All information on our generated mazes is presented in Table 1. For each maze size and type, the table presents: a) the sizes of our training, validation, and test sets, b) the range for the distance between the start and goal positions, and c) the average length of the optimal path.

## 4.2 Baselines & Evaluated Methods

We evaluate two versions of GPS and three discrete-action baselines: DQN, TempoRL, and DAR. Full implementation details of our approach are included in the Appendix.

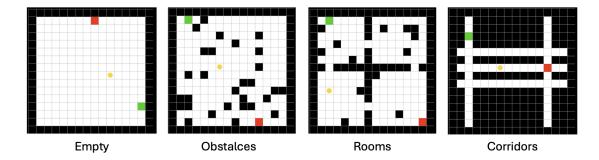


Figure 2: Examples of our generated mazes (16×16). We use four maze environments (left to right): **EMPTY** – open space; **15% Obstacles** – random obstacle placement; **ROOMS** – structured rooms with doorways; **CORRIDORS** – narrow paths requiring precise navigation. Red squares mark start positions, green squares mark goals, and yellow circles show the current agent location.

**GPS:** Our primary approach generates action sequences using a VAE-based decoder with Gumbel-Softmax sampling. This stochastic mechanism applies a temperature-controlled softmax to produce action distributions that maintain differentiability while approximating discrete samples. The Gumbel-Softmax technique creates a relaxation of categorical distributions that preserves gradients for backpropagation, facilitating end-to-end training of our actor-critic architecture.

**GPS-D:** A deterministic variant of our approach that uses argmax operations with a straight-through estimator in the decoder instead of Gumbel-Softmax sampling. This version produces consistent, deterministic action sequences for each proto-sequence embedding.

**DQN:** Deep Q-Network (Mnih et al., 2013) is a foundational model-free DRL algorithm that learns state-action values. DQN utilizes experience replay and a target network to stabilize its learning.

**DAR:** Dynamic Action Repetition (Srinivas et al., 2017) extends discrete action spaces by repeating original actions at varying rates. DAR enables the agent to select different levels of temporal control, allowing for some action abstraction.

**TempoRL:** Temporal Reinforcement Learning (Biedenkapp et al., 2021) introduces a proactive approach, where the agent selects both an action and its duration. TempoRL employs a hierarchical structure with a behavior policy for action selection and a skip policy for duration, enabling more fine-grained temporal abstraction and efficient exploration.

### 4.3 Experimental Setup

State and action representations. We represent the state using a tensor of shape (N, M, 3), where N and M are the height and width of the maze grid. The channels use a similar encoding to that of MiniGridLibrary (Chevalier-Boisvert et al., 2023): a) Object type: identifies all environmental elements including walls, empty spaces, agent position, goal location, and starting position; b) Object color: provides distinguishing colors for the start position, goal location, and current agent position; c) Placeholder channel: consistently set to 0, maintaining compatibility with the MiniGrid format. Our discrete environment supports the four basic actions – up, down, left, right – represented as a four-entry one-hot vector

Reward function. We define the reward function identically across all methods as follows:

$$R = \begin{cases} 1 & \text{if goal is reached} \\ -1/l_{max} & \text{if a valid action is executed} \\ -3/l_{max} & \text{if an invalid action is executed} \end{cases}$$

where  $l_{max}$  is the maximal start-goal distance (see Table 1) acting as a regularizer. When GPS executes a sequence of length L, we sum these Markovian step rewards:  $R_{total} = \sum (r_t)$  for t = 1..L. This reward

structure encourages goal achievement while penalizing excessive steps and invalid actions, ensuring fair comparison across all methods.

Evaluation metrics. We use three evaluation metrics:

- Average success rate (ASR): the percentage of episodes evaluated where the agent navigates successfully from start to goal position within a predefined number of steps.
- Path efficiency ratio (PER): for successfully completed episodes, we calculate the ratio between the episode length and the optimal (minimal) length:

$$PER = \frac{l_{opt}}{l_{episode}}$$

• Sequence Generation Frequency (SGF): This metric reflects how often the agent generates a new action sequence. It is calculated as the average number of times the Actor is invoked per evaluation episode. Lower values suggest the agent relies on longer-term proto-sequences before needing to generate a new sequence. This metric is relevant to GPS, DAR, and TempoRL baselines.

These metrics are complementary, as they allow us to evaluate the policy's effectiveness, efficiency, and decision frequency, under identical reward optimization across all approaches.

Neural architecture setup. All models use a shared CNN feature extractor followed by method-specific linear layers. DQN outputs Q-values for cardinal directions, DAR expands this for multiple repetition rates, and TempoRL implements a branching architecture for action and skip duration. In GPS, actor and critic networks use separate but identical CNN architectures. The actor produces a 16-dimensional protosequence embedding, which the decoder converts into action sequences through a multi-layer network with normalization. Full details are in the appendix.

**Hyperparameters & Hardware.** Unless otherwise noted for specific ablation studies, experiments were conducted using a common set of key hyperparameters, summarized in Tables 14–19 and 22 in the appendix. We selected the values based on preliminary experiments and common practices. All experiments were conducted on a system running Red Hat 5.14 with x86\_64 architecture. We used an NVIDIA RTX 2080 GPU with 8GB of VRAM.

Training Protocol & Model Selection. We used different training setups based on maze size and type. Detailed step counts are in Table 2. Model selection for final testing used the checkpoint from each run yielding the highest average success rate on a held-out set of validation environments during training. Exploration employed an  $\epsilon$ -greedy strategy, with random sequences being sampled from the same pool of 400 synthetic sequences used to train the PSD (see Appendix M.1 for generation details). This ensures exploratory sequences follow similar structural patterns to those the decoder was trained to generate. While GPS incurs computational overhead per decision, our analysis shows this is often offset by reduced decision frequency and faster convergence (see Appendix F for detailed trade-offs).

### 4.4 Evaluation Results

### 4.4.1 Evaluating the Average Success Rate (ASR).

The results of our evaluation are presented in Table 2. GPS consistently outperforms the baselines in most cases, with several key observations:

Ability to learn, converge quickly, and generalize. GPS demonstrates high sample efficiency and rapid convergence. It achieved an ASR of  $0.70 \pm 0.14$  on the  $16 \times 16$  empty maze with only 100K training steps (compared to DQN's  $0.18 \pm 0.03$ ), and substantially outperformed baselines with an ASR of  $0.82 \pm 0.31$  for  $24 \times 24$  empty maze at 500k steps, while the top baseline TempoRL only reached  $0.29 \pm 0.02$ .

Table 2: ASR Performance at Diff erent Training Steps. Values represent mean  $\pm$  standard deviation over three seeds.

Environment			$100 \mathrm{k} \; \mathrm{Steps}$					500 k Steps					1M Steps					$1.5 \mathrm{M}$ Steps		
	DQN	GPS	GPS-D	TempoRL	DAR	DQN	GPS	GPS-D	TempoRL	DAR	DQN	GPS	GPS-D	TempoRL	DAR	DQN	GPS	GPS-D	TempoRL	DAR
8x8	$0.77 \pm 0.08$	$0.99 \pm 0.01$	$0.93 \pm 0.02$	$0.84 \pm 0.07$	$0.57 \pm 0.05$	$0.86 \pm 0.08$	$1.00 \pm 0.0$	$1.00 \pm 0.01$	$0.95 \pm 0.04$	$0.74 \pm 0.03$	$0.86 \pm 0.08$	$1.00 \pm 0.0$	$1.00 \pm 0.0$	$0.95 \pm 0.04$	$0.75 \pm 0.04$	-	-	-	-	-
16x16	$0.18 \pm 0.03$	$0.70\pm0.14$	$0.50\pm0.07$	$0.31 \pm 0.06$	$0.30\pm0.04$	$0.72 \pm 0.10$	$1.00\pm0.01$	$0.95 \pm 0.01$	$0.83\pm0.06$	$0.62\pm0.04$	$0.75\pm0.05$	$1.00 \pm 0.0$	$0.97 \pm 0.01$	$0.86 \pm 0.05$	$0.67\pm0.02$	$0.76 \pm 0.06$	$1.00 \pm 0.0$	$0.99 \pm 0.01$	$0.87 \pm 0.06$	$0.67 \pm 0.02$
16x16_obst_15%	$0.10\pm0.05$	$0.38\pm0.05$	$0.31\pm0.02$	$0.07\pm0.02$	$0.31\pm0.02$	$0.76\pm0.01$	$0.90\pm0.02$	$0.66\pm0.07$	$0.69\pm0.04$	$0.66\pm0.06$	$0.82\pm0.02$	$0.95\pm0.03$	$0.80 \pm 0.09$	$0.77\pm0.03$	$0.80\pm0.09$	$0.85\pm0.03$	$0.95\pm0.03$	$0.83\pm0.06$	$0.80\pm0.01$	$0.84\pm0.06$
16x16_obst_25%	$0.03 \pm 0.01$	$0.08\pm0.02$	$0.04 \pm 0.02$	$0.03 \pm 0.01$	$0.03\pm0.01$	$0.54\pm0.10$	$0.44\pm0.15$	$0.27 \pm 0.13$	$0.13\pm0.06$	$0.08\pm0.06$	$0.70 \pm 0.07$	$0.59 \pm 0.12$	$0.48 \pm 0.08$	$0.62\pm0.07$	$0.11\pm0.07$	$0.73 \pm 0.08$	$0.64 \pm 0.17$	$0.51\pm0.14$	$0.74 \pm 0.09$	$0.15 \pm 0.04$
16x16_rooms	$0.04 \pm 0.03$	$0.13\pm0.07$	$0.06\pm0.04$	$0.03 \pm 0.02$	$0.02\pm0.01$	$0.55 \pm 0.02$	$0.56\pm0.15$	$0.39 \pm 0.13$	$0.42\pm0.10$	$0.04\pm0.02$	$0.66\pm0.01$	$0.86\pm0.06$	$0.67 \pm 0.06$	$0.60\pm0.04$	$0.09\pm0.02$	$0.66\pm0.01$	$0.89\pm0.07$	$0.72 \pm 0.06$	$0.68\pm0.06$	$0.14 \pm 0.05$
16x16_corr	$0.48 \pm 0.09$	$0.96\pm0.03$	$0.84\pm0.07$	$0.60 \pm 0.07$	$0.24\pm0.14$	$0.81\pm0.02$	$0.99\pm0.01$	$0.92 \pm 0.01$	$0.90\pm0.02$	$0.46\pm0.04$	$0.81\pm0.02$	$0.99\pm0.01$	$0.92\pm0.01$	$0.91\pm0.03$	$0.60\pm0.02$				-	-
24x24	$0.04 \pm 0.01$	$0.25 \pm 0.08$	$0.17 \pm 0.07$	$0.02 \pm 0.01$	$0.05\pm0.01$	$0.11 \pm 0.02$	$0.82\pm0.31$	$0.73 \pm 0.41$	$0.29 \pm 0.02$	$0.20 \pm 0.04$	$0.25\pm0.03$	$0.90 \pm 0.17$	$0.81 \pm 0.28$	$0.51 \pm 0.04$	$0.26\pm0.01$	$0.26 \pm 0.03$	$0.97 \pm 0.06$	$0.87 \pm 0.21$	$0.54 \pm 0.06$	$0.26 \pm 0.01$
24x24_obst_15%	$0.02 \pm 0.01$	$0.20\pm0.04$	$0.09 \pm 0.01$	-	$0.04\pm0.01$	$0.10 \pm 0.03$	$0.67\pm0.07$	$0.36 \pm 0.06$		$0.06\pm0.01$	$0.12\pm0.01$	$0.77\pm0.06$	$0.42 \pm 0.05$	-	$0.10\pm0.01$	$0.14 \pm 0.01$	$0.82\pm0.01$	$0.45 \pm 0.03$	-	$0.12 \pm 0.02$
24x24 obst 25%			-				-	-	-	-	-	-	-		-	-				

Note: The ASR for each algorithm at specific training step intervals. A gray background indicates the highest ASR achieved for that environment across all steps and algorithms. A yellow background indicates the highest ASR within that specific step interval (excluding any cell already marked gray). '-' indicates unavailable data.

This supports our hypothesis that modeling action sequences rather than individual actions enables more strategic exploration. GPS primarily learns to generate sequences that move the agent in the correct general direction toward goals, allowing progress in unseen environments even without perfectly optimized paths. The deterministic variant, GPS-D, also shows strong performance, supporting the robustness of the protosequence concept.

Scalability and superior ability to solve complex environments. The performance gap widens in larger environments. In the empty  $24\times24$  maze, GPS achieves almost perfect performance of  $0.97\pm0.06$  at  $1.5\mathrm{M}$  steps, whereas DQN and TempoRL only reach  $0.26\pm0.03$  and  $0.54\pm0.06$  respectively after  $1.5\mathrm{M}$  steps. In complex  $24\times24$  environments with 15% obstacles, our approach achieves an ASR=0.82 after  $1.5\mathrm{M}$  steps, almost eight times its closest competitor. Even in the most difficult environments ( $24\times24$  with 25% obstacles), GPS maintains a significant relative advantage (ASR=0.36 vs. 0.14 for TempoRL).

Performance on medium-sized and structured environments. In structured "rooms" environments, GPS achieves an ASR of  $0.89 \pm 0.07$  at  $1.5 \mathrm{M}$  steps, significantly outperforming DQN  $(0.66 \pm 0.01)$  and TempoRL  $(0.68 \pm 0.06)$ , while DAR struggles considerably  $(0.14 \pm 0.05)$ . This demonstrates GPS's effectiveness in navigating complex structured layouts that require coordinated multi-step strategies. In the challenging " $16 \times 16$  mazes, GPS achieves strong performance with an ASR of  $0.95 \pm 0.03$ , showing moderate improvements over DQN  $(0.85 \pm 0.03)$ , TempoRL  $(0.8 \pm 0.01)$ , and DAR  $(0.84 \pm 0.06)$ . In the most challenging " $16 \times 16$  obstacles 25 %" setup, GPS achieves an ASR of  $0.64 \pm 0.17$ , which is comparable with TempoRL  $(0.74 \pm 0.09)$  and DQN  $(0.73 \pm 0.08)$ , showing that extremely dense obstacle distributions can limit the advantages of sequence-level decision making when precise navigation is required. The performance gap becomes most pronounced in larger  $24 \times 24$  environments, where GPS achieves excellent performance (ASR of  $0.97 \pm 0.06$  for empty  $24 \times 24$  and  $0.82 \pm 0.01$  for  $24 \times 24$  obstacles  $_{15}\%$ ) compared to much lower success rates for all baselines.

The baseline methods demonstrate varying strengths across different environments. DQN achieves competitive performance in several cases, particularly in medium-complexity environments like 16x16 "obstacles\_25%" (ASR 0.73) and structured rooms (ASR 0.66). TempoRL shows strong performance in specific configurations, excelling in corridors (ASR 0.90) and achieving solid results in obstacle environments (ASR 0.74-0.80). DAR exhibits more variable performance, performing well in some obstacle configurations (ASR 0.84 in 16x16 "obstacles\_15%") but struggling significantly in structured environments like rooms (ASR 0.14). For both TempoRL and DAR, which might require more extensive training to converge optimally, we observed improved performance with larger training datasets (see Table 4 in the appendix), though computational constraints limited further exploration. While DAR also showed improvements under these conditions, it still faced significant challenges in complex structured environments, suggesting that simple action repetition may have inherent limitations for certain maze types. GPS's strong performance stems from operating in the space of action sequences rather than individual actions, enabling more strategic exploration and the discovery of long-horizon rewards that would be difficult to find using single actions or simple repetition methods.

Table 3: Comparative Performance Analysis: Convergence Speed and Efficiency Metrics. Values represent mean + standard deviation over three seeds

		A CD	Converge>	0.0 Stop				PER					SGF		
Environment		ASIL	Converge	о.э элер				1 1.11					our		
	DQN	GPS	GPS-D	TempoRL	DAR	DQN	GPS	GPS-D	TempoRL	DAR	DQN	GPS	GPS-D	TempoRL	DAR
8x8	>1M	100k	100k	200k	>1M	$0.98 \pm 0.03$	$0.89 \pm 0.01$	$\textbf{0.99}\pm\textbf{0.01}$	$0.80 \pm 0.01$	$0.61 \pm 0.02$	-	$3.36 \pm 0.05$	$2.87\pm0.07$	$5.28 \pm 0.26$	$3.99 \pm 0.20$
16x16	>1M	200k	400k	> 1.5 M	> 1.5 M	$0.91 \pm 0.15$	$0.84 \pm 0.01$	$\textbf{0.99}\pm\textbf{0.01}$	$0.94 \pm 0.01$	$0.79 \pm 0.04$	-	$8.79 \pm 0.66$	$7.69\pm0.86$	$12.50 \pm 0.97$	$8.10 \pm 0.36$
16x16_obstacles_15%	> 1.5 M	500k	> 1.5 M	> 1.5 M	> 1.5 M	$\textbf{0.96}\pm\textbf{0.07}$	$0.74 \pm 0.03$	$\textbf{0.96}\pm\textbf{0.01}$	$\textbf{0.96}\pm\textbf{0.01}$	$0.65 \pm 0.03$	-	$14.25 \pm 1.08$	$11.04\pm0.43$	$15.56 \pm 1.35$	$11.23 \pm 0.79$
16x16_obstacles_25%	> 1.5 M	> 1.5 M	> 1.5 M	> 1.5 M	> 1.5 M	$0.95 \pm 0.09$	$0.67 \pm 0.01$	$0.94 \pm 0.02$	$\textbf{0.98}\pm\textbf{0.02}$	N/R	-	$17.63 \pm 0.75$	$12.36\pm0.91$	$20.71 \pm 1.54$	N/R
16x16_rooms	> 1.5 M	900k	> 1.5 M	> 1.5 M	> 1.5 M	$0.92 \pm 0.14$	$0.68 \pm 0.01$	$0.93 \pm 0.01$	$\textbf{0.96}\pm\textbf{0.02}$	N/R	-	$15.82 \pm 0.81$	$13.18\pm2.8$	$17.56 \pm 4.35$	N/R
16x16_corridors	>1M	100k	200k	500k	>1M	$0.95 \pm 0.09$	$0.80 \pm 0.01$	$\textbf{0.97}\pm\textbf{0.01}$	$0.90 \pm 0.01$	$0.71 \pm 0.01$	-	$8.54 \pm 0.50$	$7.16 \pm 0.05$	$9.40 \pm 0.17$	$5.71\pm2.09$
24x24	> 1.5 M	1M	> 1.5 M	>1.5M	>1.5M	N/R	$0.80 \pm 0.07$	$\textbf{0.99}\pm\textbf{0.01}$	$0.96 \pm 0.01$	N/R	-	$12.40\pm1.80$	$12.78 \pm 5.98$	$17.10 \pm 0.41$	N/R
24x24_obstacles_15%	> 1.5 M	> 1.5 M	> 1.5 M	-	> 1.5 M	N/R	$0.58 \pm 0.01$	N/R	-	N/R	-	$14.55\pm0.34$	N/R	-	N/R
24x24 obstacles 25%		-	-	_	-	_	_	_	_	_	-	_	_	_	_

Note: We present four key performance metrics. The first column shows training steps required to achieve a 90% success rate (lower is better), with highlighted values indicating the fastest convergence. Path Efficiency Ratio (PER) measures trajectory optimality (higher is better, max=1.0), with **bold values** showing best performance. Sequence Generation Frequency (SGF) indicates the average number of decision points needed per episode (lower generally indicates better temporal abstraction). '-' indicates N/A and N/R' for ASR < 0.5, indicating insufficient success rate for meaningful PER/SGF evaluation.

# 4.4.2 Evaluating the Path Efficiency Ratio (PER).

The results of our evaluation are presented in Table 3. PER is calculated at the final training checkpoint using the total time steps per environment detailed in Table 1. We report PER for GPS, GPS-D, and two baselines introduced in Section 4.2, allowing for direct comparison across methods. A key GPS characteristic is Self-Correction Through Sequential Decision Points. GPS can adjust its course at subsequent decision points without requiring an initially perfect action sequence. This sequence-level closed-loop control enables course corrections while retaining the benefits of temporal abstraction (Empirical validation of this self-correction capability is provided in Appendix E, demonstrating GPS's ability to generate correction sequences with high success rates (62.4-99.6% immediate correction) across diverse maze environments). Leveraging this capability, GPS adopts a strategy of Trading Path Efficiency for Robust Navigation, prioritizing directional correctness over strict path optimality. This approach develops more transferable navigation skills—particularly evident in larger or more obstacle-dense mazes—explaining cases where PER is lower despite higher ASR and faster convergence (see Tables 2 and 3).

GPS-D consistently yields higher PER than GPS in all environments. For example, in the  $16 \times 16$  empty maze, GPS-D's PER is 0.99 versus GPS's 0.84; in the  $24 \times 24$  empty maze, PER is 0.99 for GPS-D and 0.8 for GPS. GPS's Gumbel-Softmax sampling introduces stochasticity that enables broader exploration but can cause path deviations. GPS-D's deterministic argmax decoder produces more consistent trajectories, trading exploration advantages for improved exploitation. Among the baselines, DQN and TempoRL often show high PER, frequently achieving near-optimal paths. For instance, in  $16 \times 16$  "obstacle\_15%" maze, DQN reached PER 0.96 and TempoRL 0.96. However, GPS often surpasses their ASR in complex environments. DAR generally shows a lower PER.

These findings reveal a conceptual trade-off in our approach: GPS's stochasticity boosts exploration and rapid convergence to high ASR, while GPS-D's determinism excels in path efficiency once a good policy is learned. This sequence-level exploration-exploitation trade-off offers practitioners a choice between prioritizing solution discovery (GPS) or execution efficiency (GPS-D) based on their specific requirements.

# 4.4.3 Evaluating Sequence Generation Frequency (SGF).

The results are presented in Table 3, where lower values generally indicate superior temporal abstraction due to fewer policy invocations per episode. For methods reporting SGF, our approaches GPS and GPS-D demonstrate competitive performance across environments. In  $16 \times 16$  empty maze, GPS-D achieves an SGF of 7.69, outperforming DAR's 8.10 and TempoRL's 12.5, while GPS records 8.79. However, in the  $16 \times 16$  corridors environment, DAR (SGF 5.71) outperform GPS-D (7.16) and GPS (8.54). Despite this environment-dependent variation, our methods often operate with limited interventions-such as GPS-D's 12.78 SGF in  $24 \times 24$  empty mazes versus TempoRL's 17.10-demonstrating effective generation of extended proto-sequences.

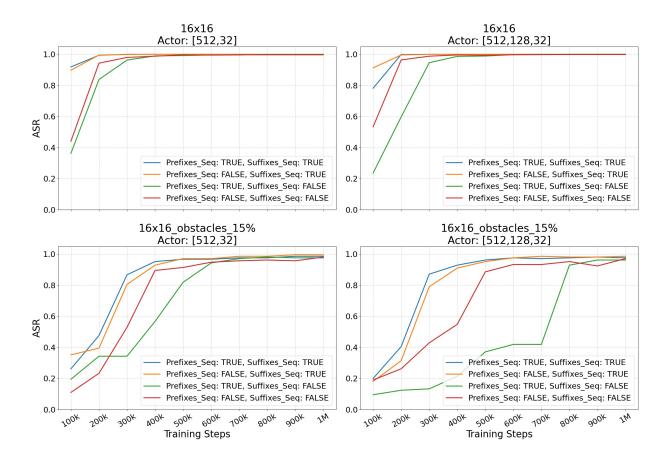


Figure 3: Impact of Subsequence Buffering Strategy on Average Success Rate.

GPS-D's generally low SGF combined with its high PER indicates capability for efficient, strategic trajectory generation through robust behavioral patterns, making it ideal for scenarios requiring predictable execution or constrained resources. GPS offers a compelling trade-off with competitive PER and favorable SGF compared to TempoRL (e.g., 8.79 vs. 12.5 in  $16 \times 16$  empty; 15.82 vs. 17.56 in  $16 \times 16$  rooms), alongside faster Average Success Rate convergence as discussed in Section 4.4.1. It balances path efficiency, sequence compactness, and learning speed effectively.

In conclusion, SGF analysis confirms our sequence-generation paradigm's effectiveness for temporal abstraction. GPS-D provides efficient, long-term utility with fewer, optimal decisions, while GPS balances competitive SGF, good PER, and rapid ASR. The choice between them depends on application priorities: efficiency vs. predictability or adaptation vs. broader performance.

# 5 Analysis and Discussion

# 5.1 Analyzing GPS's Ability to Generate Novel Sequences

While we provided our Decoder with a diverse training set, the latter did not include all possible trajectories. Our reasons were twofold. First, while including all possible action combinations was feasible in our (relatively small) action space, doing the same for larger, more complex action spaces would be infeasible or very costly. secondly, we wanted to evaluate GPS's ability to generalize and produce trajectories that were not in the training set. We consider the ability to generalize important, because the lack of it may limit the usefulness of our approach in large action spaces.

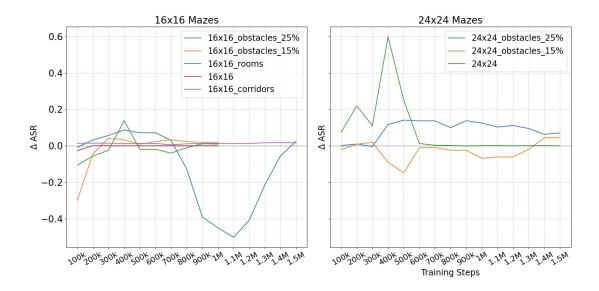


Figure 4: Impact of Actor Network Scaling on Average Success Rate ( $\Delta$  ASR) in Mazes.

As described in Appendix M.1, the PSD was pre-trained on a set of 400 synthetic sequences generated according to simple, common-sense heuristics for navigation tasks: a) each sequence contained at most two distinct action types. b) Actions of the same type appeared in contiguous blocks (e.g., "up, up, left" allowed; "up, left, up" disallowed). c) No immediately contrasting actions were allowed (e.g., "up, down" prohibited). d) Maximum sequence length was capped at  $L_{\text{max}}$  (shorter sequences permitted). e) Avoidance of loops. After full training, we gathered 15 action sequences by sampling states from the GPS replay buffer and generating the corresponding action sequences through the actor and PSD. Eleven of these did not appear in the PSD's training set and were not fully aligned with at least one of the navigation patterns described above. Using the encoding up $\rightarrow$ 0, down $\rightarrow$ 1, left $\rightarrow$ 2, right $\rightarrow$ 3, the novel sequences were:

$$\begin{split} &[\,[1,2,1,1],\,[1,1,0],\,[3,1,0,1],\,[0,0,3,1],\\ &[3,3,2],\,[3,0,3],\,[1,1,2,3],\,[2,0,1],\\ &[0,0,3,3,0,3],\,[3,3,1,2],\,[2,1,1,3]] \end{split}$$

These results show that GPS can create new action sequences not seen during training because it works in a structured embedding space. In this space, sequences with similar structures are grouped together, making it possible to blend known patterns and generate new ones, as shown in Figure 6 in the appendix.

# 5.2 Sequence Subsets Augmentation

We investigate our subsequence buffering approach (Section 3.4), implemented through *prefixes* (fixed start, varying end point) and *suffixes* (fixed goal state, varying starting point). Figure 3 shows that the baseline without subsequence buffering (red) consistently learns most slowly and often converges sub-optimally, while all subsequence buffering variants substantially improve learning efficiency. Using prefixes and suffixes simultaneously (blue) generally produces the most rapid learning, though the suffix-only configuration (orange) performs nearly as well, suggesting backward sampling provides particularly valuable learning signals. The prefix-only approach (green) typically shows slower convergence than other subsequence methods. These performance patterns remain consistent across different maze structures and actor networks.

## 5.3 Performance Under Stochastic Dynamics

Up to this point, we conducted our evaluation in a deterministic setting. We now evaluate GPS's robustness in two stochastic settings, which are more challenging to our approach because of its "commitment" to fully implementing its generated trajectories:

- "Sticky" actions. We used the well established "sticky actions" mechanism Dabney et al. (2020), with 25% probability of repeating the previous action instead of executing the planned one.
- Random actions. In this setup, which is also common in the literature Liu et al. (2024), each action has a 25% probability of being replaces with a random action.

These two setups introduce temporal correlations and execution uncertainty that challenge sequence-based methods, as errors can potentially cascade throughout multi-step sequences.

Results for the sticky actions setup show that GPS maintains strong performance under the sticky actions setup, achieving near-perfect or perfect success rates across all tested environments while preserving competitive convergence times and reasonable path efficiency. In the random actions setup, which proved to be more challenging to all evaluated algorithms (see Appendix C for details), GPS is again the top performer. Moreover, while all evaluated algorithms suffer from a degradation in their performance in this setup, GPS's relative degradation is the smallest. This result is noteworthy, since we would expect DQN, with its single actions, to be the most robust.

The robust performance under stochastic conditions provides evidence that GPS's sequence generation approach may confer resilience beyond deterministic settings. Operating at the sequence level appears to offer some natural buffering against action execution uncertainties, though comprehensive evaluation across diverse stochastic environments would be needed to fully establish this robustness. Detailed experimental details and quantitative results supporting this analysis are presented in Appendix C.

### 5.4 Performance Under Partial Observability

Beyond stochastic dynamics, we evaluated GPS's robustness under partial observability constraints, where agents only perceive a  $7 \times 7$  local view window centered on their current position rather than the full maze structure. This setting is particularly challenging for sequence-based methods as it requires generating multi-step action sequences without complete environmental information, testing whether our approach can maintain effective navigation strategies when operating with limited sensory input.

Results reveal environment-dependent performance patterns under partial observability. In empty and corridor environments, GPS demonstrates strong performance, maintaining clear advantages over DQN in empty mazes throughout training (0.54 vs 0.15 at 100K steps, 0.71 vs 0.65 at convergence) and achieving competitive results in corridors (0.78 vs 0.78 ASR at 1M steps), as action sequences help maintain consistent exploration patterns despite limited visibility. However, the moderate path efficiency ratios (PER 0.50) in these environments indicate suboptimal navigation compared to DQN's perfect efficiency, suggesting that while GPS reaches goals successfully, the generated sequences include unnecessary detours when operating with incomplete information.

In contrast, GPS faces increased challenges in obstacle-dense and rooms environments. While maintaining early advantages in obstacle mazes (0.30 vs 0.05 at 100K steps), DQN eventually surpasses GPS (0.78 vs 0.34 at 1.5M steps), and rooms environments show reduced performance (0.48 vs 0.65 ASR). We hypothesize these difficulties stem from GPS generating multi-step sequences without seeing upcoming barriers or doorways, while DQN can immediately incorporate newly revealed environmental features. The increased Sequence Generation Frequency in complex environments (16.63 in obstacles, 11.71 in empty mazes) reflects less effective navigation with incomplete spatial knowledge. No configuration reaches 0.9 ASR within 1.5M steps, indicating substantial sample complexity increases under partial observability compared to full observability settings where convergence typically occurs within 500K-1M steps. Complete experimental analysis is provided in Appendix D.

## 5.5 Impact of Actor Network Scaling

We examined actor network size impact (small: two-layer (512, 32); large: three-layer (512, 128, 32)) on maze navigation performance, measured by  $\Delta ASR$  (Large - Small) (see Figure 4). In simpler 16x16 mazes (empty or corridor), both architectures performed similarly. With 15% obstacles, the smaller network initially outperformed ( $\Delta ASR \approx -0.3$  at 100K steps) before convergence at 300K steps. In denser 25% obstacles, the smaller network significantly outperformed from 800K steps, peaking at  $\Delta ASR \approx -0.5$  at 1.1M steps. In larger 24x24 mazes with 15% obstacles, the small network generally led, despite the large network's brief advantage (500K steps). However, in the most complex 25% obstacles maze, the large network consistently outperformed, maintaining  $\Delta ASR$  between 0.1-0.15. The 24x24 empty maze showed fluctuating performance with occasional spikes for the larger network around 200K and 500K steps.

These results suggest a trade-off: smaller networks suffice or excel in smaller or moderately complex environments (possibly due to better regularization or more stable sequence generation learning), while larger networks demonstrate clear benefits in more complex environments.

## 5.6 Computational Complexity and Runtime

GPS introduces additional runtime overhead compared to baselines, as its Actor-Decoder pipeline makes each decision 1.6–3.3x slower in wall-clock time. However, because GPS operates at a sequence level and makes decisions less frequently, the effective runtime per episode is comparable to baselines in moderately complex settings (e.g., 16x16 mazes). In larger, more difficult environments GPS shows higher normalized costs, but this is partly because it continues to solve tasks where baselines fail. Despite the per-decision slowdown, GPS converges to high success rates much faster overall (90% in 3.3 hours vs. DQN's 5.83 hours), demonstrating that its superior sample efficiency compensates for runtime overhead in practice. This makes GPS more useful in challenging tasks where training effectiveness outweighs raw inference speed. A comprehensive analysis is presented in Section F in the Appendix.

# 6 Conclusions, Limitations, and Future Work

GPS is a novel actor-critic method that generates variable-length action sequences in a single step. GPS maps state observations to proto-sequences, which are decoded into discrete action sequences. This approach enhances credit assignment and exploration in long-horizon tasks by moving beyond sequential single-action selection. Our evaluation shows GPS consistently surpasses leading action repetition and temporal methods in complex maze environments, achieving higher success rates and faster convergence.

Although our approach shows benefits, particularly in complex environments, several limitations should be acknowledged. First, a new PSD needs to be trained for each unique action space, which adds to the complexity of our approach. GPS has not been evaluated on large action spaces, so adaptations to the decoder component may be needed. Furthermore, our analysis suggests that a partial information setting (e.g., limited visibility) sometimes reduces the effectiveness or our proposed approach. Finally, we have not yet adapted GPS to continuous action spaces.

Another important aspect of our approach is its suitability for domains with varying requirements. While GPS significantly outperforms the baselines in large, long horizon setups and is able to generate solutions rapidly, the resulting solutions are not always strictly optimal. This trade-off between solution quality and computational efficiency is a well-documented phenomenon in deep reinforcement learning, where methods that prioritize rapid inference may forgo perfect optimality in favor of operational practicality and responsiveness. Similar observations have been reported in recent DRL research Sohaib et al. (2025); Wu et al. (2023), highlighting how faster-converging methods frequently produce solutions that, while feasible and effective in real-world applications, may not match the optimality of simpler, slower approaches such as DQN.

Future work will focus on extending to more complex domains, such as those with larger action spaces or continuous control. As part of this research direction, we plan to explore advanced initialization strategies for our Decoder, so that our approach can more efficiently explore large action spaces. Furthermore, we will

aim to adapt GPS to function more effectively in partial information and stochastic settings. To this end, we are considering the creation of a re-planning component, and mechanism that will adapt the maximal length of the generated sequences based on the available information. Finally, we will explore modifications to our approach that will enable us to define desired trade-offs between path optimality and success. By doing so, our aim is to automatically adapt GPS's strategy to match the requirements of various domains.

# References

- Anurag Ajay, Seungwook Han, Yilun Du, Shuang Li, Abhi Gupta, Tommi Jaakkola, Josh Tenenbaum, Leslie Kaelbling, Akash Srivastava, and Pulkit Agrawal. Compositional foundation models for hierarchical planning. Advances in Neural Information Processing Systems, 36:22304–22325, 2023.
- Jose A Arjona-Medina, Michael Gillhofer, Michael Widrich, Thomas Unterthiner, Johannes Brandstetter, and Sepp Hochreiter. Rudder: Return decomposition for delayed rewards. *Advances in Neural Information Processing Systems*, 32, 2019.
- André Biedenkapp, Raghu Rajan, Frank Hutter, and Marius Lindauer. Temporl: Learning when to act. In *International Conference on Machine Learning*, pp. 914–924. PMLR, 2021.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Modular & customizable reinforcement learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023.
- Josiah D Coad, James Ault, Jeff Hykin, and Guni Sharon. A framework for predictable actor-critic control. In *Deep Reinforcement Learning Workshop NeurIPS 2022*, 2022.
- Will Dabney, Georg Ostrovski, and André Barreto. Temporally-extended {\epsilon}-greedy exploration. arXiv preprint arXiv:2006.01782, 2020.
- Zihang Dai, Qizhe Xie, and Eduard Hovy. From credit assignment to entropy regularization: Two new algorithms for neural sequence prediction. arXiv preprint arXiv:1804.10974, 2018.
- Alan Dao and Dinh Bach Vu. Alphamaze: Enhancing large language models' spatial intelligence via grpo. arXiv preprint arXiv:2502.14669, 2025.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt, Peter Sunehag, Timothy Lillicrap, Jonathan Hunt, Timothy Mann, Theophane Weber, Thomas Degris, and Ben Coppin. Deep reinforcement learning in large discrete action spaces. arXiv preprint arXiv:1512.07679, 2015.
- Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O Stanley, and Jeff Clune. First return, then explore. Nature, 590(7847):580–586, 2021.
- Liad Giladi and Gilad Katz. Feedback decision transformer: Offline reinforcement learning with feedback. In *IEEE International Conference on Data Mining (ICDM)*. IEEE, 2023.
- Chia-Chun Hung, Timothy Lillicrap, Josh Abramson, Yan Wu, Mehdi Mirza, Federico Carnevale, Arun Ahuja, and Greg Wayne. Optimizing agent behavior over long time scales by transporting value. *Nature communications*, 10(1):5223, 2019.
- Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence modeling problem. Advances in neural information processing systems, 34:1273–1286, 2021.

- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
- Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016.
- Nishanth Kumar, Tom Silver, Willie McClinton, Linfeng Zhao, Stephen Proulx, Tomás Lozano-Pérez, Leslie Pack Kaelbling, and Jennifer Barry. Practice makes perfect: Planning to learn skill parameter policies. arXiv preprint arXiv:2402.15025, 2024.
- Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(39):1–40, 2016.
- Jinghan Li, Zhicheng Sun, and Yadong Mu. Closed-loop long-horizon robotic planning via equilibrium sequence modeling. arXiv preprint arXiv:2410.01440, 2024.
- Wenhao Li. Efficient planning with latent diffusion. arXiv preprint arXiv:2310.00311, 2023.
- Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In 4th International Conference on Learning Representations (ICLR), 2016.
- Guanlin Liu, Zhihan Zhou, Han Liu, and Lifeng Lai. Efficient action robust reinforcement learning with probabilistic policy execution uncertainty. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=9sZsjfZV3q.
- Jinxin Liu, Donglin Wang, Qiangxing Tian, and Zhengyu Chen. Learn goal-conditioned policy with intrinsic motivation for deep reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pp. 7558–7566, 2022.
- Fan-Ming Luo, Tian Xu, Hang Lai, Xiong-Hui Chen, Weinan Zhang, and Yang Yu. A survey on model-based reinforcement learning. *Science China Information Sciences*, 67(2):121101, 2024.
- Marlos C Machado, Marc G Bellemare, Erik Talvitie, Joel Veness, Matthew Hausknecht, and Michael Bowling. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018.
- Thomas Mesnard, Théophane Weber, Fabio Viola, Shantanu Thakoor, Alaa Saade, Anna Harutyunyan, Will Dabney, Tom Stepleton, Nicolas Heess, Arthur Guez, et al. Counterfactual credit assignment in model-free reinforcement learning. arXiv preprint arXiv:2011.09464, 2020.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602, 2013.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, and Georg Ostrovski. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Jurgis Pašukonis, Timothy P Lillicrap, and Danijar Hafner. Evaluating long-term memory in 3d mazes. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=yHLvIlE9RGN.
- Devdhar Patel and Hava Siegelmann. Overcoming slow decision frequencies in continuous control: Model-based sequence reinforcement learning for model-free control. arXiv preprint arXiv:2410.08979, 2024.
- Roberta Raileanu and Tim Rocktäschel. Ride: Rewarding impact-driven exploration for procedurally-generated environments. arXiv preprint arXiv:2002.12292, 2020.

- Erick Rosete-Beas, Oier Mees, Gabriel Kalweit, Joschka Boedecker, and Wolfram Burgard. Latent plans for task-agnostic offline reinforcement learning. In *Conference on Robot Learning*, pp. 1838–1849. PMLR, 2023.
- Tankred Saanum, Noémi Éltető, Peter Dayan, Marcel Binz, and Eric Schulz. Reinforcement learning with simple sequence priors. Advances in Neural Information Processing Systems, 36:61985–62005, 2023.
- Younggyo Seo and Pieter Abbeel. Coarse-to-fine q-network with action sequence for data-efficient robot learning. arXiv preprint arXiv:2411.12155, 2024a.
- Younggyo Seo and Pieter Abbeel. Coarse-to-fine q-network with action sequence for data-efficient robot learning. arXiv preprint arXiv:2411.12155, 2024b.
- Sahil Sharma, Aravind Srinivas, and Balaraman Ravindran. Learning to repeat: Fine grained action repetition for deep reinforcement learning. arXiv preprint arXiv:1702.06054, 2017.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Rana Muhammad Sohaib, Syed Tariq Shah, Muhammad Ali Jamshed, Oluwakayode Onireti, and Poonam Yadav. Optimizing urllc in open ran: A deep reinforcement learning-based trade-off analysis. *IEEE Communications Standards Magazine*, 2025.
- Aravind Srinivas, Sahil Sharma, and Balaraman Ravindran. Dynamic action repetition for deep reinforcement learning. In *Proc. AAAI*, 2017.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- Giulia Vezzani, Dhruva Tirumala, Markus Wulfmeier, Dushyant Rao, Abbas Abdolmaleki, Ben Moran, Tuomas Haarnoja, Jan Humplik, Roland Hafner, Michael Neunert, et al. Skills: Adaptive skill sequencing for efficient temporally-extended exploration. arXiv preprint arXiv:2211.13743, 2022.
- Shuo Wang, Zhihao Wu, Xiaobo Hu, Youfang Lin, and Kai Lv. Skill-based hierarchical reinforcement learning for target visual navigation. *IEEE Transactions on Multimedia*, 25:8920–8932, 2023.
- Pengfei Wu, Chen Chen, Dexiang Lai, and Jian Zhong. A safe drl method for fast solution of real-time optimal power flow. arXiv preprint arXiv:2308.03420, 2023.
- Mengda Xu, Manuela Veloso, and Shuran Song. Aspire: Adaptive skill priors for reinforcement learning. Advances in Neural Information Processing Systems, 35:38600–38613, 2022.
- Dongkun Zhang, Jiaming Liang, Ke Guo, Sha Lu, Qi Wang, Rong Xiong, Zhenwei Miao, and Yue Wang. Carplanner: Consistent auto-regressive trajectory planning for large-scale reinforcement learning in autonomous driving. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 17239–17248, 2025.
- Haichao Zhang, Wei Xu, and Haonan Yu. Generative planning for temporally coordinated exploration in reinforcement learning. In *International Conference on Learning Representations*, 2022a.
- Haichao Zhang, Wei Xu, and Haonan Yu. Generative planning for temporally coordinated exploration in reinforcement learning. arXiv preprint arXiv:2201.09765, 2022b.

# A Analysis of Reward Strategy Impact on ASR

This appendix details the comparison of two step penalty strategies, illustrated in Figure 5. The strategies are the default Max-Level-based penalty  $(-1/l_{\text{max}})$ , where  $l_{\text{max}}$  is a normalization factor related to task depth, e.g.,  $l_{\text{max}} \approx 30$ ) and an alternative Map-Size-based penalty (-1/maze size), e.g., -1/256 for a  $16 \times 16$  maze). Their relative efficacy is measured by  $\Delta \text{ASR} = \text{ASR}_{\text{Max-Level}} - \text{ASR}_{\text{Map-Size}}$ , where positive values indicate superior performance for the Max-Level strategy.

Figure 5 reveals distinct performance patterns across the tested environments.

- In simple 16 × 16 mazes, the Max-Level strategy provides a significant initial learning speedup
  (ΔASR ≈ +0.8 at 200k steps), although its final Average Success Rate (ASR) is matched by the
  Map-Size strategy after approximately 400k training steps.
- When 15% obstacles are introduced in the 16 × 16 maze, increasing its complexity, the Max-Level strategy maintains a consistent performance advantage throughout the training. ΔASR peaks at approximately +0.44 and remains positive (settling around +0.1).
- In larger 24 × 24 mazes, the Max-Level strategy's superiority becomes more pronounced. ΔASR dramatically increases after 500k steps, reaching and sustaining a value of approximately +0.9. This highlights the diminishing effectiveness of the Map-Size penalty (e.g., -1/576 for 24 × 24) as it becomes increasingly diluted in larger state spaces.

The consistently superior performance of the Max-Level strategy, particularly in more complex or larger environments, can be attributed to several factors. Firstly, it provides a more **impactful and relevant penalty signal**. The  $l_{\text{max}}$ -normalized penalty (e.g.,  $\approx -1/30$ ) offers a substantially stronger and more consistent learning feedback compared to the Map-Size penalty, which diminishes significantly with increasing maze size. Secondly,  $l_{\text{max}}$  serves as a normalization factor that likely **correlates better with the intrinsic task difficulty** and typical solution length than the raw cell count of the maze, which does not inherently capture navigational complexity. Consequently, the Max-Level penalty structure appears to offer more **effective exploration guidance** and promotes **greater learning efficiency**.

In summary, normalizing step penalties by  $l_{\text{max}}$  (Max-Level-based strategy) leads to a more robust and effective reward scheme for the navigation tasks studied. This approach fosters more efficient learning and achieves higher success rates by aligning the penalty signal more accurately with the inherent challenges of the environment, proving especially advantageous as task complexity and scale increase.

# B Evaluating the Average Success Rate (ASR) With Larger Train dataset

To assess performance on a larger dataset, we trained the agent on 2000 mazes. Table 4 presents the Average Success Rates (ASR) across various maze configurations.

Table 4. Ave	Table 4. Average buces states (Abri) with 2000 mazes in train dataset across maze types.											
Maze Type	Max	Train	Dist. from	DQN	GPS	GPS-D	TempoRL	DAR				
Maze Type	Steps	Size	Start to Goal	DQN	GIS	GI S-D	Temport	DAIL				
8x8	1M	2000	[1 - 14]	1.00	1.00	1.00	1.00	0.81				
16x16	1M	2000	[16 - 26]	1.00	1.00	1.00	1.00	0.89				
16x16_obst_15%	1.5M	2000	[20 - 30]	0.94	0.99	0.94	0.91	0.78				
16x16_obst_25%	1.5M	2000	[20 - 30]	0.98	0.98	0.94	0.39	0.16				
16x16_rooms	1.5M	2000	[20 - 30]	1.00	0.95	0.82	0.98	0.13				
$16x16\_corr$	1M	2000	[10 - 30]	1.00	1.00	1.00	1.00	0.67				
24x24	1.5M	2000	[20 - 30]	0.98	1.00	0.98	1.00	0.75				
24x24 obst $15%$	1.5M	2000	[10 - 30]	0.05	0.09	0.05	0.03	0.13				

Table 4: Average Success Rates (ASR) with 2000 mazes in train dataset across maze types.

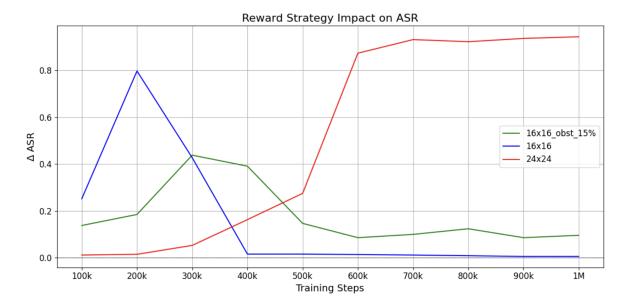


Figure 5: Comparison of average success rate differences ( $\Delta$ ASR) between two reward strategies: **Max-Level** (the default strategy) and **Map-Size**, which is identical except that it replaces  $l_{\text{max}}$  with the total number of cells in the maze (e.g., 256 for a  $16\times16$  maze). The values are evaluated across training steps. A positive  $\Delta$ ASR indicates that the Max-Level reward strategy yields better performance.

# C Evaluation on Stochastic Environments

Real-world environments feature inherent stochasticity through action execution noise, sensor uncertainty, and environmental dynamics. To evaluate GPS's robustness under such conditions, we extended our evaluation framework to include two types of stochastic dynamics:

- "Sticky" actions. We implemented a setup where there is a 25% probability that the agent executes the previous action in the sequence instead of the current planned action. This "sticky actions" mechanism follows established protocols in reinforcement learning evaluation Machado et al. (2018); Dabney et al. (2020), introducing temporal correlations and systematic execution biases that create realistic uncertainty and test an agent's ability to adapt to imperfect action execution.
- Random actions. In this setup each action has a 25% probability of being replaces with a random action. These setups are also common in the literature, where they are used to simulate adversarial environments Liu et al. (2024), or environments with noisy signals Sutton et al. (1999).

These two setups present unique challenges for sequence-based methods like GPS. Since our approach commits to entire action sequences, execution errors can potentially cascade throughout the sequence, leading to significant deviations from intended trajectories. This makes stochastic environments a critical test for the robustness of our sequence generation paradigm.

## C.1 Sticky Actions Results

Table 5 presents GPS performance results under the "sticky" stochastic conditions across representative maze configurations, comparing GPS under deterministic conditions versus GPS with sticky actions (GPS-SA) under stochastic conditions. The results demonstrate that GPS exhibits remarkable resilience to stochastic dynamics. Despite the 25% action execution uncertainty, GPS-SA achieves near-perfect or perfect success rates in all environments, reaching ASR = 1.00 in empty mazes (8×8, 16×16, and 24×24) and maintaining strong performance (ASR = 0.94) even in the challenging  $16\times16$  maze with 15% obstacles. Comparing GPS

Table 5: GPS Performance on Stochastic Maze Environments (25% Sticky Actions). GPS-SA denotes GPS with sticky actions. Both GPS and GPS-SA results are based on single-seed experiment.

				ASR at Tra	aining S	teps				Converge .9 Step	]	PER	SGF	
Environment	1	100K		500K		1M	1	1.5M						
	GPS	GPS-SA	GPS	GPS-SA	GPS	GPS-SA	GPS	GPS-SA	GPS	$\operatorname{GPS-SA}$	$\operatorname{GPS}$	$\operatorname{GPS-SA}$	GPS	$\operatorname{GPS-SA}$
8x8	1.00	0.99	1.00	1.00	1.00	1.00	-	-	100k	100k	0.90	0.71	2.9	5.45
16x16	0.96	0.77	1.00	1.00	1.00	1.00	1.00	1.00	100k	300k	0.84	0.80	6.9	8.83
$16x16\_obstacles\_15\%$	0.22	0.33	0.96	0.89	0.99	0.94	0.96	0.94	300k	600k	0.72	0.61	10.8	16.97
24x24	0.14	0.13	1.00	1.00	1.00	1.00	1.00	1.00	500k	300k	0.78	0.78	9.6	12.41

Table 6: GPS Performance on Stochastic Maze Environments (25% Uniform Random Actions). Both GPS and DQN results are based on single-seed experiment.

			AS	R at Tra	aining S	teps				onverge Step	P	ER	SG	F
Environment	10	0K	50	00K	1	M	1	5M						
	GPS	DQN	GPS	DQN	GPS	DQN	GPS	DQN	GPS	$_{\mathrm{DQN}}$	GPS	DQN	GPS	$_{\mathrm{DQN}}$
8x8	0.98	0.74	0.99	0.86	1.0	0.86	-	-	100k	>1M	0.5	0.63	6.93	
16x16	0.39	0.1	0.79	0.53	0.9	0.59	0.91	0.59	900k	> 1.5 M	0.52	0.7	14.13	-
$16x16\_obstacles\_15\%$	0.11	0.05	0.28	0.58	0.61	0.63	0.75	0.7	> 1.5 M	> 1.5 M	0.49	0.62	20.26	-
24x24	0.08	0.07	0.33	0.09	0.72	0.19	0.72	0.23	> 1.5 M	> 1.5 M	0.54	0.58	19.33	-

performance across conditions reveals execution uncertainty's impact: while GPS achieves rapid convergence in deterministic settings, sticky actions extend training requirements. The  $16\times16$  empty maze exemplifies this trade-off: convergence time triples from 100K to 300K steps under sticky actions, yet final ASR remains perfect (1.00) in both conditions. Path Efficiency Ratio shows moderate degradation from deterministic to stochastic conditions (0.84 to 0.80), indicating reduced optimality while maintaining effective navigation. Obstacle-dense environments face more pronounced challenges, with the  $16\times16$  maze requiring doubled convergence time (300K to 600K steps) and greater PER reduction (0.72 to 0.61), as execution errors in constrained spaces increase collision likelihood.

The Path Efficiency Ratio (PER) values of 0.61-0.80 indicate that GPS maintains reasonably efficient navigation despite the stochastic perturbations. While slightly lower than deterministic performance, these values indicate that the agent successfully adapts to execution uncertainty while maintaining reasonable navigation efficiency. The Sequence Generation Frequency (SGF) increases under sticky actions, ranging from 5.45 to 16.97 sequence generations per episode, reflecting the need for more frequent replanning when execution deviates from intended sequences.

Notably, GPS under sticky actions still outperforms the baseline methods (DQN, TempoRL, DAR) operating in deterministic environments. For instance, GPS-SA achieves perfect ASR in the  $24 \times 24$  maze at 500K steps, while DQN only reaches 0.24 ASR even at 1M steps in the deterministic setting. This superior performance eliminates the need to evaluate baselines under sticky actions, as they already struggle in the simpler deterministic case.

# C.2 Random Actions Results

The stochastic setup of random actions presents a potential larger challenge to our proposed approach. Although this setup has the same probability as the previous one of injecting an unintended action into the sequence, these actions may be more diverse (and potentially disruptive) than those of the previous setup. For example, an optimal trajectory for a given maze may include only "up" and "right" actions; in the sticky setup, the injected actions may only include these two actions. In the random setup, every action – including those not included in the optimal trajectory – is equally possible.

The results of our evaluation are presented in Table 6. We present a comparison to DQN, since this baseline, which performs a single action at a time, is likely to be the most robust against the injection of random actions. The results clearly show that both algorithms achieve reduced performance in this stochastic

setup: in all evaluated setups except 8x8 (the smallest and easiest), the evaluated algorithms underperform compared to their non-stochastic results. However, GPS maintains a large lead over DQN: the former's final ASR performance in all evaluated scenarios is significantly higher.

The results of the PER metric reveals that uniform random action injection causes uniform degradation across all environments. PER values under uniform random conditions cluster tightly between 0.49-0.54, representing drops of 0.17-0.29 from sticky action performance. Under sticky actions, PER varied substantially by environment (0.61-0.80 range), but uniform random perturbations compress all environments to approximately 0.50 PER. This uniform degradation reflects how unpredictable disruptions systematically undermine sequence coherence, reducing the path optimality benefits of temporal abstraction compared to more predictable sticky perturbations. Sequence Generation Frequency (SGF) increases substantially under uniform random versus sticky conditions: +27% (8×8: 6.93 vs 5.45), +60% (16×16: 14.13 vs 8.83), +19% (16×16\_obstacles\_15%: 20.26 vs 16.97), and +56% (24×24: 19.33 vs 12.41). The obstacle-dense environment shows highest absolute SGF (20.26) but smallest relative increase, having already required frequent replanning under sticky conditions. Open environments show larger relative increases (56-60%), indicating unpredictable perturbations necessitate more frequent sequence regeneration than correlated sticky disruptions. These elevated SGF values, particularly in complex environments, demonstrate that uniform random perturbations substantially erode temporal abstraction benefits, though GPS's sustained ASR advantage over DQN indicates that multi-step sequence generation retains value over purely reactive single-action decisions.

## C.3 Summary

These findings suggest that GPS's sequence-level decision making provides inherent robustness against action execution noise. By generating coherent multi-step action sequences, the approach creates a natural buffer against individual action failures. The sequence generation paradigm appears to offer resilience to stochastic perturbations, as the method can adapt its behavioral patterns over time while maintaining overall navigational effectiveness despite local execution uncertainties.

The robust performance under stochastic conditions provides initial evidence that temporal abstraction through sequence generation may offer advantages beyond sample efficiency. The ability to operate at the sequence level appears to confer resilience to environmental uncertainties, suggesting potential applicability to scenarios where perfect action execution cannot be guaranteed, though further evaluation across diverse stochastic environments would be needed to establish broader generalizability of these findings.

# D Evaluation on Partial Observability Environments

Real-world tasks often involve partial observability, where agents must make decisions based on limited local information rather than complete environmental knowledge. To evaluate GPS's robustness under such constraints, we tested our approach with agents restricted to a  $7 \times 7$  observation window centered on their current position, obscuring the majority of the  $16 \times 16$  maze structure.

We implemented partial observability by modifying the observation space to show only a  $7 \times 7$  window around the agent's current position. Cells outside this window are marked as unobserved, preventing the agent from accessing global maze structure or distant goal locations. This constraint fundamentally challenges sequence-based methods like GPS, which must generate entire action sequences without visibility of the complete path to the goal. The partial observability setting introduces limited planning horizons (visibility only 3 cells in any direction), implicit memory requirements for previously observed areas, and increased exploration complexity without global visibility.

Table 7 presents GPS's performance under partial observability constraints compared to the DQN baseline across representative maze configurations. The results reveal distinct performance patterns across environment types. In empty and corridor environments, GPS demonstrates strong performance under partial observability. In empty  $16\times16$  mazes, GPS maintains clear advantages throughout training (0.54 vs 0.15 at 100K steps, 0.71 vs 0.65 at convergence), while in corridor environments, it achieves competitive results with an ASR of 0.78 at 1M steps compared to DQN's 0.78. This suggests that generating action sequences provides inherent advantages under partial observability: multi-step sequences naturally maintain movement

Table 7. GI	Dane	1 DQ1	<b>1</b> 1 C1.	loi ilia	ncc c	nuci .	1 61 616	ii Obs		v (	\ 1 VI	OVV VVI	ildow j	
			AS	R at Tra	aining S	teps				onverge Step	P	ER	SG	F
Environment	10	0K	50	0K	1	M	1.	5M						
	GPS	DQN	GPS	DQN	GPS	DQN	GPS	DQN	GPS	$_{\mathrm{DQN}}$	$\operatorname{GPS}$	DQN	GPS	$_{\mathrm{DQN}}$
16x16	0.54	0.15	0.67	0.63	0.71	0.65	0.71	0.65	>1.5M	>1.5M	0.50	1.00	11.71	
$16x16\_obstacles\_15\%$	0.30	0.05	0.34	0.74	0.34	0.78	0.34	0.78	> 1.5 M	> 1.5 M	0.53	1.00	16.63	-
16x16_corridors	0.48	0.49	0.76	0.78	0.78	0.78	-	-	> 1M	> 1M	0.50	0.99	13.68	-
16x16_rooms	0.11	0.06	0.39	0.62	0.48	0.63	0.48	0.65	> 1.5 M	> 1.5 M	0.66	0.98	14.08	-

Table 7: GPS and DQN Performance Under Partial Observability (7 × 7 view window)

direction and exploration consistency, helping the agent navigate effectively even when it cannot see its destination. The narrow pathways of corridor mazes may provide implicit guidance that reduces the disadvantage of partial observability—once an agent commits to moving down a corridor, the constrained environment naturally channels movement toward valid destinations.

GPS faces increased challenges in obstacle-dense and rooms environments. While maintaining early advantages in obstacle mazes (0.30 vs 0.05 at 100K steps), DQN eventually surpasses GPS in 15% obstacle environments (0.78 vs 0.34 at 1.5M steps), and rooms environments present significant challenges with GPS achieving 0.48 ASR at 1.5M steps. We hypothesize these difficulties stem from the compounded uncertainty of partial observability and complex spatial navigation: in obstacle environments, GPS must generate multistep sequences without seeing upcoming barriers, leading to frequent collisions with unseen walls, while rooms environments require precise doorway navigation and chamber transitions that are particularly challenging when global layout information is unavailable. DQN's step-by-step approach allows it to immediately incorporate each newly revealed environmental feature into its next decision, while GPS remains committed to sequences planned from incomplete information, accumulating penalties from suboptimal actions in complex spaces. In addition, none of the tested configurations achieved the 0.9 ASR convergence threshold within 1.5M training steps, indicating significantly increased sample complexity under limited visibility.

The efficiency metrics further illuminate these trade-offs. GPS maintains moderate Path Efficiency Ratios (0.50-0.66) indicating consistent but suboptimal navigation, while DQN achieves perfect efficiency when successful. The elevated Sequence Generation Frequency under partial observability (11.71-16.63 vs 8.79-15.8 with full visibility) reflects the increased challenges of operating with limited environmental information—more frequent sequence generation becomes necessary to maintain navigation progress. The higher SGF in obstacle environments (16.63 vs 11.71) indicates that GPS must generate new sequences more frequently when navigating uncertain terrain, likely due to sequences that prove less effective when executed with incomplete knowledge of obstacle locations, requiring more decision points to navigate around newly discovered barriers and supporting our hypothesis about the challenges of sequence planning in complex partially-observed spaces.

These findings highlight both the potential and limitations of sequence-based decision making under partial observability. GPS's early-stage advantages and moderate path efficiency (PER 0.5) demonstrate that multi-step action sequences can provide useful structure even with limited sensory input, though performance degrades significantly compared to full observability. The increased generation frequency shows GPS attempting to adapt to limited visibility, though obstacle-dense environments remain challenging. Future work might investigate confidence-conditioned sequence horizons, where the decoder learns to generate shorter sequences when observation uncertainty is high, balancing temporal abstraction benefits with the need for frequent replanning.

## **E** Self-Correction Analysis

To substantiate our claim that GPS can self-correct through sequential sequences (Section 4.4.2), we conducted a comprehensive analysis of GPS's ability to recover from deviations in its trajectory.

## E.1 Methodology

We developed a quantitative framework to measure GPS's self-correction behavior by tracking progress toward the goal using Breadth-First Search (BFS) distance. Our analysis focuses on identifying when GPS deviates from optimal progress and measuring its ability to generate correction sequences in subsequent decisions.

#### E.1.1 Deviation Detection

For each action sequence executed during an episode, we calculate:

- 1. Initial BFS distance ( $d_{before}$ ): The shortest path distance to the goal before sequence execution
- 2. Final BFS distance ( $d_{after}$ ): The shortest path distance to the goal after sequence execution

A deviation is detected when the sequence does not bring the agent closer to the goal:

$$d_{\text{after}} \ge d_{\text{before}}$$
 (1)

This indicates the agent either maintained the same distance or moved further from the goal, representing a sequence that requires correction.

## E.1.2 Correction Sequence Classification

When a deviation is identified, we track whether GPS generates a correction sequence in subsequent decisions. We measure correction at multiple time horizons:

- Immediate Correction: The very next sequence reduces the optimal path length from its own starting position to its ending position, demonstrating GPS's ability to generate locally optimal moves that recover from the previous suboptimal decision
- Near-term Correction (within 2 sequences): Recovery occurs within the next two decision points
- Medium-term Correction (within 3-4 sequences): Recovery within 3-4 subsequent decisions
- Long-term Correction (5+ sequences): Recovery before episode termination

A correction sequence is considered successful when the agent's BFS distance falls below what would have been expected from optimal execution of the original deviating sequence.

## E.2 Results

Table 8 presents the self-correction analysis results across representative maze environments. The analysis was performed only on successful episodes to focus on GPS's recovery capabilities when it ultimately reaches the goal.

## E.3 Analysis and Discussion

The results provide strong empirical evidence for GPS's self-correction capability:

High Immediate Correction Rates: Across all environments, GPS demonstrates substantial immediate correction rates, ranging from 62.4% in complex room environments to 99.6% in open spaces. This indicates that GPS frequently generates correction sequences at the very next decision point after detecting a deviation.

## **Environment-Specific Patterns:**

• Open spaces (16×16, 24×24): Nearly perfect immediate correction (>89%), reflecting GPS's ability to quickly identify and generate correction sequences for directional errors

Table 8: CP9	S Self-Correction	Success Rates	Across Mazz	Environments
Table of Great	э эен-Сопесион	Ouccess nates	ACTOSS MAZE	e raivironnems

Environment	Immediate Correction	$egin{aligned} \mathbf{Near-term} \ (\leq 2 \ \mathbf{seq}) \end{aligned}$	$egin{array}{l} { m Medium ext{-}term} \ (\leq 4 { m  seq}) \end{array}$	$   \begin{array}{c}     \text{Long-term} \\     (5 + \text{seq})   \end{array} $
8×8	$0.885 \pm 0.095$	$0.889 \pm 0.092$	$0.896 \pm 0.086$	$1.000 \pm 0.000$
$16\times16$	$0.996 \pm 0.006$	$1.000 \pm 0.000$	$1.000 \pm 0.000$	$1.000 \pm 0.000$
$16 \times 16$ _obstacles_ $15\%$	$0.753 \pm 0.311$	$0.779 \pm 0.313$	$0.794 \pm 0.292$	$1.000 \pm 0.000$
$16 \times 16$ _obstacles_ $25\%$	$0.765 \pm 0.324$	$0.789 \pm 0.298$	$0.808 \pm 0.272$	$1.000 \pm 0.000$
$16 \times 16$ _rooms	$0.624 \pm 0.114$	$0.773 \pm 0.088$	$0.805 \pm 0.079$	$1.000 \pm 0.000$
$16 \times 16$ _corridors	$0.965 \pm 0.046$	$1.000 \pm 0.000$	$1.000 \pm 0.000$	$1.000 \pm 0.000$
$24 \times 24$	$0.899 \pm 0.143$	$1.000 \pm 0.000$	$1.000 \pm 0.000$	$1.000 \pm 0.000$

Note: Values represent mean  $\pm$  standard deviation across representative models from three different training seeds per maze type. Analysis conducted only on successful episodes.

- Structured environments (rooms, obstacles): Lower immediate correction rates (62-77%) but strong recovery within 2-4 sequences, suggesting GPS adapts its correction strategy to environmental complexity
- Corridors: High immediate correction (96.5%), likely due to the constrained nature limiting deviation possibilities

Guaranteed Long-term Recovery: All environments show 100% long-term correction rates in successful episodes, confirming that GPS consistently generates effective correction sequences even after significant trajectory deviations.

Variance Patterns: Higher standard deviations in obstacle-rich environments ( $\pm 0.3$ ) compared to open spaces ( $\pm 0.006$ ) indicate that correction difficulty varies with local maze structure, yet GPS maintains robust recovery capabilities.

The evidence presented here directly supports our claim that GPS employs sequence-level closed-loop control, generating correction sequences at subsequent decision points without requiring initially perfect action sequences. This self-correction mechanism helps explain GPS's performance characteristics, where the ability to recover from deviations enables reliable goal-reaching (high ASR) even when individual sequences may not follow strictly optimal paths (resulting in lower PER).

# F Computational Cost Analysis

While GPS demonstrates superior sample efficiency, we analyze its computational trade-offs compared to baseline methods to provide a complete picture of the method's practicality.

**Pre-training Costs.** The VAE decoder requires approximately 10-15 minutes of pre-training on an Apple M1 Max with 64 GB RAM. This is a one-time cost per action space, and the trained decoder can be reused across different environments with the same action space, amortizing this cost over multiple experiments.

Inference Overhead. During inference, GPS requires 1.6-3.3x the wall-clock time per decision compared to baselines due to the Actor-Decoder pipeline. In the 16x16 maze with 15% obstacles, GPS takes 2.72ms  $\pm$  0.25ms per decision versus DAR's 1.39ms  $\pm$  0.31ms (1.96x overhead) and DQN's 1.67ms  $\pm$  0.05ms (1.63x overhead), while being 1.17x faster than TempoRL's 3.18ms  $\pm$  0.76ms. In the larger 24x24 environment, the overhead pattern persists: GPS requires 3.82ms  $\pm$  0.81ms per decision compared to DAR's 1.16ms  $\pm$  0.09ms (3.29x overhead), DQN's 1.68ms  $\pm$  0.05ms (2.27x overhead), and TempoRL's 2.89ms  $\pm$  0.77ms (1.32x overhead). However, since GPS makes decisions less frequently due to its sequence-level abstraction, we also measure normalized episode efficiency—the total episode wall-clock time divided by the optimal path length to normalize for task difficulty and enable fair comparison across algorithms with different success rates. In the 16x16 with 15% obstacles environment, GPS (4.6ms  $\pm$  0.19ms per optimal step) performs

comparably to baselines (DQN: 4.03ms, TempoRL: 4.55ms, DAR: 4.65ms), indicating that the per-decision overhead is largely offset by reduced decision frequency. However, in the 24x24 environment, GPS shows higher normalized costs (7.24ms  $\pm$  1.51ms per optimal step) compared to baselines (DQN: 5.83ms, TempoRL: 5.83ms, DAR: 4.92ms). Note that this metric is computed only on successful episodes, which may bias results toward baselines that solve only easier test instances (GPS: 90% ASR vs DQN: 24%, TempoRL: 46%, DAR: 23%), suggesting the true computational cost difference for solving challenging instances may be smaller than indicated.

Training Time to Convergence. Despite higher per-step costs, GPS achieves faster wall-clock convergence in medium complexity environments like 16x16 mazes with 15% obstacles due to its superior sample efficiency. In 16x16 mazes with 15% obstacles, GPS reaches 90% ASR in 3.3 hours (500K steps), while other methods fail to reach 90% ASR within the same training budget: DQN achieves only 85% ASR after 5.83 hours (1.5M steps), TempoRL reaches 82% ASR after 9.2 hours (1.5M steps), and DAR achieves 64% ASR after 4.08 hours (1.5M steps). For the larger 24x24 maze, GPS's computational overhead becomes more apparent, with training times becoming comparable to the fastest baseline despite superior convergence and final performance. GPS achieves 90% ASR in 11.5 hours (1M steps), while all baselines fail to reach 90% ASR: DQN reaches only 24% ASR after 10.5 hours (1.5M steps), TempoRL achieves 46% ASR after 13.02 hours (1.5M steps), and DAR reaches 23% ASR after 6.26 hours (1.5M steps), illustrating the trade-off between computational efficiency and learning effectiveness.

Computational Trade-offs and Practical Considerations. While GPS incurs additional computational overhead per decision compared to most baselines, its reduced decision frequency often compensates for this cost. The computational trade-offs of GPS vary with environment complexity: in medium complexity environments like 16x16 with 15% obstacles maze, GPS provides clear training time advantages, reaching target performance levels significantly faster than baselines. In larger environments such as 24x24 maze, training times become comparable to the fastest baselines, though GPS achieves substantially superior final performance (90% vs 24% ASR for DQN). This pattern suggests that GPS becomes most practical when learning effectiveness is prioritized over raw computational speed, particularly in scenarios where baseline methods struggle to reach acceptable performance levels rather than purely on training efficiency grounds.

Table 9: Computational cost analysis across environments and methods. Values represent mean  $\pm$  standard deviation over 5 mag.

Environment	Metric	GPS	DQN	TempoRL	DAR
	Inference time per decision (ms)	$2.72 \pm 0.25$	$1.67 \pm 0.05$	$3.18 \pm 0.76$	$1.39 \pm 0.31$
16x16_obst_15%	Normalized episode efficiency (ms/opt. step)	$4.6 \pm 0.19$	$4.03 \pm 0.08$	$4.55 \pm 0.75$	$4.65 \pm 0.56$
10X10_008t_1570	Training wall-clock to 90% ASR	3.3h	5.83h	9.2h	4.08h
	Total training wall-clock (h)	$9.97 \pm 0.31$	$5.83 \pm 0.35$	$9.2 \pm 1.72$	$4.08 \pm 0.65$
	Inference time per decision (ms)	$3.82 \pm 0.81$	$1.68 \pm 0.05$	$2.89 \pm 0.77$	$1.16 \pm 0.09$
24x24	Normalized episode efficiency (ms/opt. step)	$7.24 \pm 1.51$	$5.83 \pm 0.17$	$5.83 \pm 1.18$	$4.92 \pm 0.54$
Z4XZ4	Training wall-clock to 90% ASR	11.5h	10.5h	13.02h	6.26h
	Total training wall-clock (h)	$17.35 \pm 2.11$	$10.5 \pm 0.39$	$13.02 \pm 1.67$	$6.26 \pm 0.12$
VAE Pre-trainin	10	-15 minutes or	Apple M1 Ma	X	

In summary, GPS represents a favorable computational trade-off for challenging navigation tasks where baseline methods fail to achieve acceptable performance, offering superior learning effectiveness at competitive training costs.

# G Comparison: Pre-trained Decoder vs. End-to-End Training

When training GPS, one needs to choose whether to pre-train the Proto-Sequence Decoder (PSD) or train it jointly end-to-end with the actor and critic networks. To assess the impact of this choice, we conducted a comprehensive comparison across multiple maze environments. This analysis was motivated by the observation that decoder pre-training could potentially encode implicit biases about feasible action sequences, raising questions about whether GPS's advantages stem from the sequence-level decision-making paradigm itself or from initialization artifacts.

Table 10: Performance comparison between GPS with pre-trained decoder and end-to-end training (GPS-E2E). All results use seed=123.

				ASR at Tra	ining S	teps				Converge ).9 Step		PER	5	SGF
Environment		100K		500K		1M		1.5M						
	GPS	GPS-E2E	GPS	GPS-E2E	GPS	GPS-E2E	GPS	GPS-E2E	GPS	$\operatorname{GPS-E2E}$	GPS	$\operatorname{GPS-E2E}$	GPS	$\operatorname{GPS-E2E}$
16x16	0.71	0.62	1.00	1.00	1.00	1.00	1.00	1.00	300k	200k	0.82	0.88	9.19	11.54
$16\mathrm{x}16\_\mathrm{obs}\_15\%$	0.43	0.31	0.91	0.92	0.92	0.92	0.92	0.92	300k	300k	0.74	0.78	13.30	19.70
$16x16\_rooms$	0.21	0.13	0.72	0.89	0.91	0.90	0.91	0.90	700k	1M	0.68	0.74	16.27	28.48
$16x16\_corr$	0.93	0.96	0.99	0.99	0.99	0.99	-	=	100k	100k	0.75	0.82	9.75	15.96
24x24	0.28	0.11	1.00	1.00	1.00	1.00	1.00	1.00	200k	300k	0.83	0.86	11.42	14.78

## G.1 Experimental Setup

We evaluated two **GPS** variants:

- GPS: Our default implementation with a pre-trained VAE decoder (frozen during actor-critic training)
- GPS-E2E: End-to-end training where all three components (actor, decoder, critic) are jointly optimized from random initialization

Both variants were evaluated on the same set of maze environments using identical hyperparameters for the actor and critic networks. The decoder architecture remained unchanged; only the training procedure differed.

# G.2 Results and Analysis

Table 10 presents the comparative results across representative environments. The findings reveal several important insights:

# **G.2.1** Performance Parity in Success Rate and Convergence

The most significant finding is that **GPS-E2E** achieves comparable final performance to GPS with pre-trained decoder across all tested environments. Both variants converge to similarly high success rates, with convergence speeds being environment-dependent (e.g.,  $16\times16$ : GPS-E2E reaches ASR > 0.9 at 200k steps vs GPS at 300k steps;  $24\times24$ : GPS converges at 200k steps vs GPS-E2E at 300k steps). The final ASR values are nearly identical across environments (e.g.,  $16\times16$ : both reach 1.00;  $16\times16$ \_obs\_15%: 0.92 vs 0.92;  $16\times16$  rooms: 0.91 vs 0.90). This demonstrates that:

- 1. The performance advantage of GPS stems primarily from the **sequence-level decision-making paradigm** rather than from implicit biases encoded in decoder pre-training.
- 2. GPS can successfully learn effective action sequence representations from scratch during joint training.
- 3. Decoder pre-training is not a fundamental requirement. However, as we show below, there are trade-offs with regard to path optimality and the number of sequences that need to be generated per trajectory.

#### G.2.2 Path Efficiency vs. Temporal Abstraction Trade-off

An interesting pattern emerges when examining the Path Efficiency Ratio (PER) and Sequence Generation Frequency (SGF) metrics:

• **Higher PER for GPS-E2E**: End-to-end training consistently produces slightly higher path efficiency (e.g.,  $16 \times 16$ : 0.88 vs 0.82;  $16 \times 16$ \_obs\_15%: 0.78 vs 0.74). This indicates that GPS-E2E learns to generate action sequences that more closely follow optimal paths.

• Lower SGF for GPS: The pre-trained decoder variant requires fewer sequence generations per episode (e.g., 16×16: 9.19 vs 11.54; 16×16\_rooms: 16.27 vs 28.48), suggesting it produces longer, more temporal abstraction action sequences.

This trade-off reveals an important behavioral difference: GPS with pre-trained decoder appears to develop more aggressive temporal abstraction strategies—generating longer sequences that reduce decision frequency at the cost of some path optimality. In contrast, GPS-E2E learns to balance temporal abstraction with path efficiency, generating somewhat shorter sequences that more closely track optimal trajectories.

We hypothesize this difference arises because:

- 1. The pre-trained decoder's latent space structure, learned from diverse synthetic sequences, encourages exploration of longer, more varied action patterns
- 2. End-to-end training jointly optimizes sequence generation with the specific task objectives, potentially leading to more task-specific (and thus more efficient but less exploratory) sequence patterns

This suggests that decoder pre-training may provide a form of **implicit exploration regularization**, encouraging the agent to commit to longer temporal abstractions even when shorter sequences might be locally optimal.

# G.3 Decoder Transferability Across Tasks

An important practical advantage of the pre-trained decoder approach is its reusability. In our experiments, the same pre-trained decoder was successfully used across:

- All maze configurations (empty, obstacles with varying densities, rooms, corridors)
- All maze sizes  $(8\times8, 16\times16, 24\times24)$
- Multiple experimental conditions (deterministic dynamics, stochastic environments and partial observability)

This transferability demonstrates that the decoder learns general-purpose sequence generation capabilities that apply across diverse navigation scenarios within the same action space. For applications involving multiple related tasks with shared action spaces, pre-training the decoder once and reusing it can provide practical benefits. However, our end-to-end results confirm that decoder reusability is an added benefit rather than a necessity for GPS's effectiveness. The choice between variants ultimately reflects the PER-SGF trade-off: pre-training encourages more aggressive temporal abstraction (lower SGF) while end-to-end training optimizes for path efficiency (higher PER).

## **G.4** Conclusion

Our analysis demonstrates that GPS achieves strong performance with both pre-trained and end-to-end training, confirming that its advantages stem from sequence-level decision-making rather than initialization artifacts. The comparable performance of GPS-E2E validates that GPS's superior performance compated to the baselines (Section 4.4) is derived by its architectural approach to temporal abstraction.

The pre-trained variant offers decoder reusability and more aggressive temporal abstraction (lower SGF), while the end-to-end variant provides superior path efficiency (higher PER) with comparable success rates. This PER-SGF trade-off suggests decoder initialization influences the exploration-exploitation balance. Practitioners can select either variant based on specific requirements without sacrificing fundamental performance.

Future work could explore combining both approaches' strengths through adaptive decoder training schedules (starting pre-trained for exploration, transitioning to task-specific optimization) or curriculum learning that leverages the complementary benefits of both paradigms.

# **H** ASR Statistical Significance Testing

To assess the statistical significance of the differences in Average Success Rates (ASR) between our proposed method (GPS) and the baseline (DQN), we employed McNemar's test. This section details the methodology and presents the results of these tests.

## Methodology

McNemar's test is a non-parametric test suitable for paired nominal data. It is used to determine whether there is a significant difference in the proportions of two related samples, such as when two algorithms are evaluated on the same set of test instances. In our context, each maze evaluation episode serves as a paired instance, and the outcome for each algorithm (GPS or DQN) is categorized as either a success or a failure.

An episode was deemed a **success** if the agent reached the goal in an episodic length of less than 75 steps. Otherwise, it was considered a **failure**.

For each pair of algorithms (GPS vs. DQN) on a given maze type, we constructed a  $2 \times 2$  contingency table based on the outcomes of common evaluation episodes:

		Algorith	m B (DQN)
		Success	Failure
Algorithm A (GPS)	Success	a	b
	Failure	c	d

Where:

- a: Number of episodes where both GPS and DQN succeeded.
- b: Number of episodes where GPS succeeded and DQN failed.
- c: Number of episodes where GPS failed and DQN succeeded.
- d: Number of episodes where both GPS and DQN failed.

McNemar's test focuses on the discordant pairs (b and c). The null hypothesis ( $H_0$ ) is that the two algorithms have the same ASR. The test statistic is calculated as:

$$\chi^2 = \frac{(b-c)^2}{b+c}$$

This statistic follows a chi-squared distribution with 1 degree of freedom. We used the version of the test without continuity correction, as implemented in 'statsmodels.stats.contingency\_tables.mcnemar'.

The significance level was set at  $\alpha = 0.05$ . If the calculated p-value was less than 0.05, we rejected the null hypothesis and concluded that there is a statistically significant difference in the ASR performance of the two algorithms.

## Results: GPS vs. DQN

The results of McNemar's test comparing GPS (Algorithm A) to DQN (Algorithm B) across various maze configurations are summarized in Table 11. The Average Success Rate (ASR) reported in the table for each algorithm is based on Table 2:

- ASR (GPS) = (a+b)/(a+b+c+d)
- ASR (DQN) = (a+c)/(a+b+c+d)

Table 11: McNemar's Test Results for GPS vs. DQN. All p-values < 0.05 indicate a statistically significant difference in performance, favoring GPS in all listed cases. P-values reported as 0.0000 by the script are presented as < 0.0001.

Maze Type	ASR (GPS)	ASR (DQN)	McNemar Stat.	p-value
8x8	1.00	0.95	54.0000	< 0.0001
16x16	1.00	0.69	377.0000	< 0.0001
$16x16\_obs\_15$	0.96	0.85	22.1538	< 0.0001
$16x16\_obs\_25$	0.90	0.8	28.4462	< 0.0001
$16x16\_rooms$	0.92	0.65	136.5329	< 0.0001
16x16_corridors	1.00	0.80	107.0000	< 0.0001
24x24	1.00	0.24	757.0000	< 0.0001
24x24_obs_15	0.91	0.15	760.0208	< 0.0001
$24x24\_obs\_25$	0.36	0.11	192.9627	< 0.0001

All comparisons in Table 11 yield p-values substantially less than 0.05, demonstrating statistically significant improvements of GPS over DQN across all tested maze environments. The consistent outcomes and significant p-values robustly support the conclusion that the GPS method offers superior performance compared to the DQN baseline under the specified experimental conditions.

The following sections provide the detailed per-run summaries logged and the specific contingency tables used for McNemar's test for each maze configuration.

## Maze: 8x8

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 1000, Failures: 0, Errors: 0
- DQN (Algorithm B) Summary: Total episodes: 1000, Successes: 946, Failures: 54, Errors: 0

Contingency Table (GPS vs. DQN):

		DQN (Algorithm B	
		Success	Failure
GPS (Alg. A)	Success	946 (a)	54 (b)
	Failure	0 (c)	0 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 54.0000, p-value: < 0.0001

## Maze: 16x16

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 1000, Failures: 0, Errors: 0
- DQN (Algorithm B) Summary: Total episodes: 1000, Successes: 623, Failures: 377, Errors: 0

 $Contingency\ Table\ (GPS\ vs.\ DQN):$ 

		DQN (Algorithm B)	
		Success Failure	
GPS (Alg. A)	Success $Failure$	623 (a) 0 (c)	377 (b) 0 (d)

Common episodes for comparison: 999

McNemar's Statistic: 342.0000, p-value: <  $\bf 0.0001$ 

# Maze: 16x16\_obs\_15

- GPS (Algorithm A) Summary: Total episodes: 210, Successes: 202, Failures: 8, Errors: 0
- DQN (Algorithm B) Summary: Total episodes: 210, Successes: 178, Failures: 32, Errors: 0

Contingency Table (GPS vs. DQN):

		DQN (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	177 (a)	25 (b)
	Failure	1 (c)	7 (d)

Common episodes for comparison: 210

McNemar's Statistic: 22.1538, p-value: < 0.0001

## Maze: 16x16\_obs\_25

- GPS (Algorithm A) Summary: Total episodes: 399, Successes: 360, Failures: 39, Errors: 0
- DQN (Algorithm B) Summary: Total episodes: 402, Successes: 320, Failures: 82, Errors: 0

Contingency Table (GPS vs. DQN):

		DQN (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	306 (a)	54 (b)
	Failure	11 (c)	28 (d)

Common episodes for comparison: 399

McNemar's Statistic: 28.4462, p-value: < 0.0001

# Maze: 16x16\_rooms

- GPS (Algorithm A) Summary: Total episodes: 586, Successes: 530, Failures: 56, Errors: 0
- DQN (Algorithm B) Summary: Total episodes: 586, Successes: 379, Failures: 207, Errors: 0

Contingency Table (GPS vs. DQN):

DQN (Algorithm B)	
Failure	
159 (b)	
48 (d)	
_	

Common episodes for comparison: 583

McNemar's Statistic: 132.7872, p-value: < 0.0001

# Maze: 16x16\_corridors

- GPS (Algorithm A) Summary: Total episodes: 545, Successes: 545, Failures: 0, Errors: 0
- DQN (Algorithm B) Summary: Total episodes: 545, Successes: 438, Failures: 107, Errors: 0

Contingency Table (GPS vs. DQN):

		DQN (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success $Failure$	438 (a) 0 (c)	107 (b) 0 (d)

Common episodes for comparison: 545

McNemar's Statistic: 107.0000, p-value: < 0.0001

# Maze: 24x24

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 999, Failures: 1, Errors: 0
- DQN (Algorithm B) Summary: Total episodes: 1000, Successes: 242, Failures: 758, Errors: 0

Contingency Table (GPS vs. DQN):

		DQN (Algorithm B)	
		Success Failure	
GPS (Alg. A)	Success $Failure$	242 (a) 0 (c)	757 (b) 1 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 757.0000, p-value: < 0.0001

# Maze: 24x24\_obs\_15

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 914, Failures: 86, Errors: 0
- DQN (Algorithm B) Summary: Total episodes: 1000, Successes: 150, Failures: 850, Errors: 0

Contingency Table (GPS vs. DQN):

		DQN (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success $Failure$	148 (a) 2 (c)	766 (b) 84 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 760.0208, p-value: <  $\bf 0.0001$ 

# Maze: 24x24\_obs\_25

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 368, Failures: 632, Errors: 0
- DQN (Algorithm B) Summary: Total episodes: 1000, Successes: 99, Failures: 901, Errors: 0

Contingency Table (GPS vs. DQN):

		DQN (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	46 (a)	322 (b)
	Failure	53 (c)	579 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 192.9627, p-value: < **0.0001** 

This detailed breakdown for each environment shows the specific data underlying the McNemar's tests.

# Results: GPS vs. TempoRL

Table 12 summarizes the Average Success Rates (ASR) for GPS and TempoRL, along with the McNemar test statistics and p-values derived from common paired evaluation episodes.

All comparisons in Table 12 yield p-values substantially less than 0.05, demonstrating statistically significant improvements of GPS over the TempoRL baseline across all tested maze environments. The consistent outcomes and significant p-values robustly support the conclusion that the GPS method offers superior performance compared to the TempoRL baseline under these experimental conditions.

Table 12: Summary of McNemar's Test Results for GPS vs. TempoRL. All p-values < 0.05 indicate a statistically significant difference in performance, favoring GPS.

Maze Type	ASR (GPS)	ASR (TempoRL)	McNemar Stat.	p-value
8x8	1.00	0.97	26.0000	< 0.0001
16x16	1.00	0.84	178.0000	< 0.0001
$16x16\_obs\_15$	0.96	0.82	19.5652	< 0.0001
$16x16\_obs\_25$	0.90	0.79	18.9804	< 0.0001
$16x16\_rooms$	0.92	0.63	119.0088	< 0.0001
$16x16$ _corridors	1.00	0.9	56.0000	< 0.0001
24x24	1.00	0.46	511.0078	< 0.0001
$24x24\_obs\_15$	0.91	0.20	651.0968	< 0.0001
$24x24\_obs\_25$	0.36	0.09	171.6100	< 0.0001

The following sections provide the detailed per-run summaries and the specific contingency tables used for McNemar's test for each maze configuration when comparing GPS with TempoRL.

## Maze: 8x8 (GPS vs. TempoRL)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 1000, Failures: 0, Errors: 0
- TemporL (Algorithm B) Summary: Total episodes: 1000, Successes: 974, Failures: 26, Errors: 0

Contingency Table (GPS vs. TemporL):

		TemporL (Algorithm B	
		Success	Failure
GPS (Alg. A)	Success	974 (a)	26 (b)
	Failure	0 (c)	0 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 26.0000, p-value: < 0.0001

# Maze: 16x16 (GPS vs. TempoRL)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 1000, Failures: 0, Errors: 0
- TemporL (Algorithm B) Summary: Total episodes: 1000, Successes: 822, Failures: 178, Errors: 0

Contingency Table (GPS vs. TemporL):

		TemporL Success	(Algorithm B) Failure
GPS (Alg. A)	Success $Failure$	822 (a) 0 (c)	178 (b) 0 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 178.0000, p-value: < 0.0001

# Maze: 16x16\_obs\_15 (GPS vs. TempoRL)

- GPS (Algorithm A) Summary: Total episodes: 210, Successes: 202, Failures: 8, Errors: 0
- TemporL (Algorithm B) Summary: Total episodes: 210, Successes: 172, Failures: 38, Errors: 0

Contingency Table (GPS vs. TemporL):

		TemporL (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	164 (a)	38 (b)
	Failure	8 (c)	0 (d)

Common episodes for comparison: 210

McNemar's Statistic: 19.5652, p-value: < 0.0001

# Maze: 16x16\_obs\_25 (GPS vs. TempoRL)

- GPS (Algorithm A) Summary: Total episodes: 399, Successes: 360, Failures: 39, Errors: 0
- TemporL (Algorithm B) Summary: Total episodes: 402, Successes: 318, Failures: 84, Errors: 0

Contingency Table (GPS vs. TemporL):

		TemporL (Algorithm B)	
		Success Failure	
GPS (Alg. A)	Success	287 (a)	73 (b)
	Failure	29 (c)	$10 \; (d)$

Common episodes for comparison: 399

McNemar's Statistic: 18.9804, p-value: < 0.0001

## Maze: 16x16\_rooms (GPS vs. TempoRL)

- GPS (Algorithm A) Summary: Total episodes: 586, Successes: 530, Failures: 56, Errors: 0
- TemporL (Algorithm B) Summary: Total episodes: 586, Successes: 366, Failures: 220, Errors: 0

Contingency Table (GPS vs. TemporL):

		TemporL (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	335 (a)	195 (b)
	Failure	31 (c)	25 (d)

Common episodes for comparison: 586

McNemar's Statistic: 119.0088, p-value: <  $\bf 0.0001$ 

# Maze: 16x16\_corridors (GPS vs. TempoRL)

- GPS (Algorithm A) Summary: Total episodes: 545, Successes: 545, Failures: 0, Errors: 0
- TemporL (Algorithm B) Summary: Total episodes: 545, Successes: 489, Failures: 56, Errors: 0

Contingency Table (GPS vs. TemporL):

		TemporL (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	489 (a)	56 (b)
	Failure	0 (c)	0 (d)

Common episodes for comparison: 545

McNemar's Statistic: 56.0000, p-value: < 0.0001

# Maze: 24x24 (GPS vs. TempoRL)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 999, Failures: 1, Errors: 0
- TemporL (Algorithm B) Summary: Total episodes: 1000, Successes: 486, Failures: 514, Errors: 0

Contingency Table (GPS vs. TemporL):

		TemporL (Algorithm B)	
		Success Failure	
GPS (Alg. A)	Success	485 (a)	514 (b)
	Failure	1 (c)	0 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 511.0078, p-value: < **0.0001** 

# Maze: 24x24\_obs\_15 (GPS vs. TempoRL)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 914, Failures: 86, Errors: 0
- TemporL (Algorithm B) Summary: Total episodes: 1000, Successes: 218, Failures: 782, Errors: 0

Contingency Table (GPS vs. TemporL):

		TemporL (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	194 (a)	720 (b)
	Failure	24 (c)	$62 \; (d)$

Common episodes for comparison: 1000

McNemar's Statistic: 651.0968, p-value: < 0.0001

# Maze: 24x24\_obs\_25 (GPS vs. TempoRL)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 368, Failures: 632, Errors: 0
- TemporL (Algorithm B) Summary: Total episodes: 1000, Successes: 106, Failures: 894, Errors: 0

Contingency Table (GPS vs. TemporL):

		TemporL Success	(Algorithm B) Failure
GPS (Alg. A)	Success $Failure$	37 (a) 69 (c)	331 (b) 563 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 171.6100, p-value: < 0.0001

This detailed breakdown for each environment when comparing GPS to TempoRL shows the specific data underlying McNemar's tests.

# Results: GPS vs. DAR

Table 13 summarizes the Average Success Rates (ASR) for GPS and DAR, along with the McNemar test statistics and p-values derived from common evaluation episodes.

Table 13: Summary of McNemar's Test Results for GPS vs. DAR. All p-values < 0.05 indicate a statistically significant difference in performance, favoring GPS.

Maze Type	ASR (GPS)	ASR (DAR)	McNemar Stat.	p-value
8x8	1.00	0.76	244.0000	< 0.0001
16x16	1.00	0.61	394.0000	< 0.0001
$16x16\_obs\_15$	0.96	0.64	57.8000	< 0.0001
$16x16\_obs\_25$	0.90	0.14	298.2038	< 0.0001
$16x16\_rooms$	0.92	0.15	438.0800	< 0.0001
16x16_corridors	1.00	0.61	213.0000	< 0.0001
24x24	1.00	0.23	711.0000	< 0.0001
$24x24\_obs\_15$	0.91	0.12	614.2257	< 0.0001
$24x24\_obs\_25$	0.36	0.07	236.2798	< 0.0001

All comparisons in Table 13 yield p-values substantially less than 0.05, demonstrating statistically significant improvements of GPS over the DAR baseline across all tested maze environments. The consistent outcomes and significant p-values robustly support the conclusion that GPS offers superior performance compared to the DAR baseline under these experimental conditions.

The following sections provide the detailed per-run summaries logged by the script and the specific contingency tables used for McNemar's test for each maze configuration when comparing GPS with DAR.

## Maze: 8x8 (GPS vs. DAR)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 1000, Failures: 0, Errors: 0
- DAR (Algorithm B) Summary: Total episodes: 1000, Successes: 756, Failures: 244, Errors: 0

Contingency Table (GPS vs. DAR):

		DAR (Algorithm B)	
		Success Failure	
GPS (Alg. A)	Success $Failure$	756 (a) 0 (c)	244 (b) 0 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 244.0000, p-value: < 0.0001

# Maze: 16x16 (GPS vs. DAR)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 1000, Failures: 0, Errors: 0
- DAR (Algorithm B) Summary: Total episodes: 1000, Successes: 606, Failures: 394, Errors: 0

Contingency Table (GPS vs. DAR):

		DAR (Algorithm B)	
		Success Failure	
GPS (Alg. A)	Success	606 (a)	394 (b)
	Failure	0 (c)	0 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 394.0000, p-value: < 0.0001

# Maze: 16x16\_obs\_15 (GPS vs. DAR)

- GPS (Algorithm A) Summary: Total episodes: 210, Successes: 202, Failures: 8, Errors: 0
- DAR (Algorithm B) Summary: Total episodes: 210, Successes: 134, Failures: 76, Errors: 0

Contingency Table (GPS vs. DAR):

		DAR (Algorithm B)	
		Success Failure	
GPS (Alg. A)	Success	128 (a)	74 (b)
	Failure	6 (c)	2 (d)

Common episodes for comparison: 210

McNemar's Statistic: 57.8000, p-value: < 0.0001

# Maze: 16x16\_obs\_25 (GPS vs. DAR)

- GPS (Algorithm A) Summary: Total episodes: 399, Successes: 360, Failures: 39, Errors: 0
- DAR (Algorithm B) Summary: Total episodes: 402, Successes: 54, Failures: 348, Errors: 0

Contingency Table (GPS vs. DAR):

		DAR (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	50 (a)	310 (b)
	Failure	4 (c)	35 (d)

Common episodes for comparison: 399

McNemar's Statistic: 298.2038, p-value: < 0.0001

# Maze: 16x16\_rooms (GPS vs. DAR)

- GPS (Algorithm A) Summary: Total episodes: 586, Successes: 530, Failures: 56, Errors: 0
- DAR (Algorithm B) Summary: Total episodes: 586, Successes: 86, Failures: 500, Errors: 0

Contingency Table (GPS vs. DAR):

		DAR (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success $Failure$	83 (a) 3 (c)	447 (b) 53 (d)

Common episodes for comparison: 586

McNemar's Statistic: 438.0800, p-value: < **0.0001** 

## Maze: 16x16\_corridors (GPS vs. DAR)

- GPS (Algorithm A) Summary: Total episodes: 545, Successes: 545, Failures: 0, Errors: 0
- DAR (Algorithm B) Summary: Total episodes: 544, Successes: 331, Failures: 213, Errors: 0

Contingency Table (GPS vs. DAR):

		DAR (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	331 (a)	213 (b)
	Failure	0 (c)	0 (d)

Common episodes for comparison: 544

McNemar's Statistic: 213.0000, p-value: < 0.0001

## Maze: 24x24 (GPS vs. DAR)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 999, Failures: 1, Errors: 0
- DAR (Algorithm B) Summary: Total episodes: 1000, Successes: 288, Failures: 712, Errors: 0

Contingency Table (GPS vs. DAR):

		DAR (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	288 (a)	711 (b)
	Failure	0 (c)	1 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 711.0000, p-value: < 0.0001

## Maze: 24x24\_obs\_15 (GPS vs. DAR)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 914, Failures: 86, Errors: 0
- DAR (Algorithm B) Summary: Total episodes: 4833, Successes: 524, Failures: 4309, Errors: 0

Contingency Table (GPS vs. DAR):

		DAR (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success $Failure$	38 (a) 6 (c)	632 (b) 55 (d)

Common episodes for comparison: 731

McNemar's Statistic: 614.2257, p-value: <  $\bf 0.0001$ 

## Maze: 24x24\_obs\_25 (GPS vs. DAR)

- GPS (Algorithm A) Summary: Total episodes: 1000, Successes: 368, Failures: 632, Errors: 0
- DAR (Algorithm B) Summary: Total episodes: 1000, Successes: 66, Failures: 934, Errors: 0

Contingency Table (GPS vs. DAR):

		DAR (Algorithm B)	
		Success	Failure
GPS (Alg. A)	Success	24 (a)	344 (b)
	Failure	42 (c)	590 (d)

Common episodes for comparison: 1000

McNemar's Statistic: 236.2798, p-value: < 0.0001

This detailed breakdown for each environment when comparing GPS to DAR shows the specific data underlying McNemar's tests.

#### I Baselines and Architecture

Each baseline is evaluated using a grid search over multiple hyperparameter configurations; Tables 14, 15, 16, 17, 18, 19 and 22 detail the specific value ranges for these parameters.

All baseline models employ the same CNN feature extractor architecture followed by similarly sized linear layers, differing only in the final output layer size. For example, DQN outputs 4 Q-values (one per action), while DAR outputs 12 (4 actions × 3 repetition heads). TempoRL requires an additional network head to implement the skip policy, adding architectural complexity but gaining flexibility in temporal decision-making. In our GPS method, the actor and critic networks each have their own separate CNN state feature extractors. In future work, we plan to explore a shared CNN feature extraction architecture as implemented in TempoRL, which could potentially improve computational efficiency and state representation learning.

For TempoRL, we configured the model with a maximum skip length between 1..10 to allow variable sequence lengths of action repetition. For DAR, we evaluated possible coarse control values of 1,5,10 to allow the same maximum sequence length and mid-sequence capability, with the fine control value fixed at 1 to allow for actions at every time step. We based our implementations on the publicly available code at https://github.com/automl/TempoRL but reimplemented from scratch to enrich with more metrics and employ our evaluation methodology. Detailed architectures, hyperparameter configurations, and implementation specifics can be found in Appendix J, K, L, M and N.

While GPS, like hierarchical methods, targets long-horizon tasks, its approach is fundamentally different. Hierarchical RL learns temporal abstractions via skill discovery and multi-level policies, whereas GPS generates full action sequences directly from state observations without skills or sub-goal decomposition. The proto-sequence decoder is pre-trained and fixed, providing a structured output space rather than a learned controller. Because of these differences and GPS 's focus on discrete action spaces, we compare against DAR and TempoRL, which similarly extend temporal horizons through action repetition and sequence commitment.

# J DQN Baseline Implementation Details

This section outlines the architecture and configuration of the Deep Q-Network (DQN) agent used as a baseline. It details the neural network structure, hyperparameter settings, exploration strategy, optimization method, and other relevant training aspects.

## J.1 Model Architecture (QNetwork)

The Q-Network is a neural network designed to approximate the action-value function Q(s, a). It consists of a convolutional part for feature extraction from the input observation and a linear part for producing Q-values for each action.

In our maze environments, as depicted in Subsection 4.3, the input observation has a shape (C, H, W), where the number of input channels C is 3. The number of output channels,  $n\_output\_channels$ , corresponds to the number of available actions, which is 4 (right, left, up, down).

## J.1.1 Convolutional Neural Network (CNN) Part

The CNN component processes the input observation through a sequence of convolutional layers:

## 1. Conv2D Layer 1:

Input channels: 3Output channels: 16

• Kernel size: 2 • Stride: 1

2. Activation: ReLU

## 3. Conv2D Layer 2:

Input channels: 16Output channels: 32

Kernel size: 2Stride: 1

4. Activation: ReLU

#### 5. Conv2D Layer 3:

Input channels: 32Output channels: 64

Kernel size: 2Stride: 1

6. Activation: ReLU

7. **Flatten Layer**: The output of the convolutional layers is flattened into a 1D vector. The size of this vector,  $n\_flatten$ , is computed automatically.

#### J.1.2 Linear Part

The flattened output from the CNN (n\_flatten) is fed into a sequence of fully connected linear layers:

- The hidden layer sizes are configurable via grid search (see Table 14 for details). The activation function for these hidden layers is Leaky ReLU (negative slope 0.1).
- The final linear layer maps the last hidden layer's output to n\_output\_channels (4 actions).

## J.2 Hyperparameters

The agent's behavior and training process are governed by a set of hyperparameters, detailed in Tables 14, 15, and 16.

Table 14: General Experiment Hyperparameters for DQN Baseline

Parameter	Default Value
seed	123
torch_deterministic	True
save_model_strategy	SUCCESS_RATE
val_eval_freq	5000
train_eval_freq	5000
eval_test_dataset_training_freq	100000

Table 15: Environment-Specific Hyperparameters for DQN Baseline

Parameter	Default Value
max_episode_steps	75
reward_strategy	NEGATIVE_BASED_ON_MAX_LEVEL_WITH_PENALTIES
observation_encoding_strategy	DEFAULT
Max Path Length (max_level)	Varies (see Env. Def. in Table 1)
Min Path Length (start_level)	Varies (see Env. Def. in Table 1)

Table 16: Algorithm Specific Hyperparameters for DQN Baseline

Parameter	Default Value / Options
total_timesteps	Environment specific
learning_rate	$[1 \times 10^{-3}, 1 \times 10^{-4}]$
buffer_size	[10000, 50000]
$\gamma$ (discount factor)	0.99
$\tau$ (target update rate)	[0.01, 0.005]
target_network_frequency	[10, 100] (soft-target update freq.)
batch_size	256
start_e	1.0 (initial $\varepsilon$ )
end_e	$0.1 \text{ (final } \varepsilon)$
exploration_fraction	[0.1, 0.3, 0.5]
learning_starts	1000 (timestep to begin learning)
train_frequency	2 (Q-network update freq.)
linear_layers	["512,128,32", "512,32"]
activation_function	Leaky ReLU (slope 0.1)

# J.3 Epsilon-Greedy Exploration

The agent uses an epsilon-greedy strategy for action selection. The value of epsilon  $(\epsilon)$  is linearly annealed from start\_e (1.0) to end\_e (0.1) over a duration. This duration is calculated as  $\lfloor \exp frac \times total\_timesteps \rfloor$ , where  $\exp frac$  is the selected exploration\_fraction (from options in Table 16) and t is the current global timestep. The epsilon at timestep t is:

$$\epsilon_t = \max((((\text{end}_e - \text{start}_e)/\text{duration}) \times t + \text{start}_e), \text{end}_e)$$

With probability  $\epsilon_t$ , a random action is chosen; otherwise, the action with the highest Q-value is selected.

## J.4 Optimizer

The Q-Network is trained using the Adam optimizer (torch.optim.Adam). The learning rate is controlled by the learning\_rate hyperparameter (see Table 16).

## J.5 Replay Buffer

A replay buffer (stable\_baselines3.common.buffers.ReplayBuffer) stores experiences  $(s_t, a_t, r_t, s_{t+1}, d_t)$ .

The buffer size is specified in Table 16. Key configurations include optimize\_memory\_usage = False and handle\_timeout\_termination = False.

## J.6 Training Details

**Loss Function.** The Q-Network parameters ( $\theta$ ) are updated by minimizing the Mean Squared Error (MSE) loss:

$$L(\theta) = \mathbb{E}_{(s,a,r,s',d)\sim\mathcal{B}} \left[ (y_t - Q(s,a;\theta))^2 \right]$$

where the TD target:

$$y_t = r_t + \gamma \max_{a'} Q_{target}(s_{t+1}, a'; \theta^-)(1 - d_t)$$

Here,  $r_t$  is the reward,  $\gamma$  is the discount factor,  $Q_{target}$  is the target network with parameters  $\theta^-$ , and  $d_t$  indicates if  $s_{t+1}$  is terminal. This is implemented via torch.nn.functional.mse\_loss.

**Target Network.** A separate target network  $Q_{target}$  with parameters  $\theta^-$  stabilizes training. Its weights are updated using Polyak averaging:  $\theta^- \leftarrow \tau\theta + (1-\tau)\theta^-$ . The soft update rate  $\tau$  and update frequency target\_network\_frequency are specified in Table 16.

## Training Procedure.

- Learning Starts: Training begins after learning starts timesteps (see Table 16).
- Training Frequency: The Q-network is updated every train\_frequency global steps (see Table 16).
- Batch Size: Number of experiences sampled per training step is batch\_size (see Table 16).

#### J.7 Evaluation

The agent's performance is evaluated periodically on validation and test datasets.

- Evaluation on the validation dataset occurs every val\_eval\_freq steps.
- Evaluation on the test dataset can occur during training every eval\_test\_dataset\_training\_freq steps.
- During evaluation, actions are chosen greedily (or with a small fixed epsilon, e.g., 0.05 or 0.0).
- Metrics logged include mean episodic return, success rate, and agent step ratio.
- Model saving is based on performance metrics (e.g., highest success rate or reward on validation) as per save\_model\_strategy.

# K DAR Baseline Implementation Details

This section outlines the architecture and configuration of the Dyanmic Action Repetition (DAR) agent used as a baseline. The DAR agent builds upon the Deep Q-Network (DQN) architecture and training methodology. Therefore, for aspects not explicitly mentioned here, such as the general experiment configuration (Table 14), environment-specific arguments (Table 15), epsilon-greedy exploration strategy (Section J.3), optimizer (Section J.4), replay buffer (Section J.5), general training procedure (Section J.6), and evaluation methodology (Section J.7), please refer to the corresponding descriptions in the DQN baseline implementation details (Section J).

The primary distinctions of the DAR baseline are its modified network architecture to support an expanded action space and an additional algorithm-specific hyperparameter, dar\_r\_1, related to action repetition.

#### K.1 Model Architecture

The DAR network for the DAR agent, similar to DQN, approximates the action-value function Q(s, a). It comprises a convolutional part for feature extraction and a linear part for producing Q-values.

The input observation from the maze environments has a shape (C, H, W), where C = 3, identical to the DQN baseline (Section J.1).

## K.1.1 Convolutional Neural Network (CNN) Part

The CNN component is identical to the one used in the DQN baseline. For details on the architecture (number of layers, channels, kernel sizes, strides, and activations), please refer to Section J.1. The output of this part is a flattened 1D vector of size  $n\_flatten$ .

#### K.1.2 Linear Part

The flattened output (n\_flatten) from the CNN is processed by a sequence of fully connected linear layers:

- The hidden layer sizes are configurable via grid search, with the same options as the DQN baseline (see Table 17 for linear\_layers). The activation function for these hidden layers is Leaky ReLU (negative slope 0.1).
- The final linear layer maps the last hidden layer's output to n\_output\_channels. For the DAR agent, n\_output\_channels = 12, corresponding to 4 base actions (right, left, up, down) each associated with 3 repetition heads/levels.

#### K.2 Hyperparameters

The general experimental configuration and environment-specific hyperparameters for the DAR baseline are the same as those for the DQN baseline, as detailed in Table 14 and Table 15, respectively.

## K.2.1 Algorithm Specific Arguments

The algorithm-specific hyperparameters for the DAR baseline, including the newly introduced dar\_r\_l parameter, are listed in Table 17. These parameters are subject to grid search to find the optimal configuration for each environment.

Table 17: Algorithm Specific Hyperparameters for DAR Baseline

Parameter	Default Value / Options
total_timesteps	Environment specific
learning_rate	$[1 \times 10^{-3}, 1 \times 10^{-4}]$
buffer_size	[10000, 50000]
$\gamma$ (discount factor)	0.99
au (target update rate)	[0.01, 0.005]
target_network_frequency	[10, 100] (soft-target update freq.)
batch_size	256
start_e	1.0 (initial $\varepsilon$ )
end_e	$0.1 \text{ (final } \varepsilon)$
exploration_fraction	[0.1, 0.3, 0.5]
learning_starts	1000 (timestep to begin learning)
train_frequency	2 (Q-network update freq.)
linear_layers	["512,128,32", "512,32"]
activation_function	Leaky ReLU (slope 0.1)
dar_r_l	[1, 5, 10] (repetition level parameter)

# K.3 Training Details

**Loss Function.** For DAR, the Q-Network parameters  $(\theta)$  are updated by minimizing Hubber loss. This is implemented via torch.nn.SmoothL1Loss.

## L TempoRL Baseline Implementation Details

This section describes the architecture and configuration of the TempoRL agent, a baseline designed for temporal abstraction by learning how long to repeat actions. TempoRL shares several components and procedures with the DQN baseline. For details on the general experiment configuration (Table 14), environment-

specific arguments (Table 15), replay buffer (Section J.5), and evaluation methodology (Section J.7), please refer to the corresponding descriptions in the DQN baseline implementation details (Section J).

Key distinctions of the TempoRL agent include its specialized network architecture with separate heads for action selection and skip duration, unique hyperparameters related to these mechanisms (skip\_dim, weight sharing), and the use of Huber loss for training.

#### L.1 Model Architecture

The TempoRL network processes input observations to produce Q-values for primitive actions and Q-values for skip durations. The input observation from the maze environments has a shape (C, H, W), where C = 3, identical to the DQN baseline (Section J.1).

## L.1.1 Convolutional Neural Network (CNN) Part

The CNN component used for initial feature extraction is identical to the one in the DQN baseline. For details on its architecture (number of layers, channels, kernel sizes, strides, and activations), please refer to Section J.1. The output of this CNN part is a flattened 1D vector of size n flatten.

## L.1.2 Linear Heads for Action and Skip Policies

Following the CNN, the network processes features through a structure that leads to two distinct output heads: one for action selection and one for determining the skip duration. The MLP for each pathway (from CNN output to pre-output layer) consists of layers with output units [512, 128, 32].

## • Feature Processing and Weight Sharing:

- If weight\_sharing = True (default configuration): The n\_flatten vector is first processed by a shared linear layer producing 512 output units, followed by a Leaky ReLU activation (negative slope 0.1). This 512-unit feature vector serves as the common input to the subsequent differing layers of the action and skip heads.
- If weight\_sharing = False: The  $n\_flatten$  vector is independently fed into the first linear layer (512 output units, Leaky ReLU) of both the action and skip processing streams. Each stream then continues with its own [128, 32] layers.

## • Action Head:

- Starting from the 512-unit feature vector (either shared or head-specific), it is processed through two subsequent linear layers with 128 and 32 output units, respectively. Each of these hidden layers uses a Leaky ReLU activation (negative slope 0.1).
- The final linear layer of the action head maps the 32-unit feature vector to  $n\_output\_actions$  Q-values, where  $n\_output\_actions = 4$  (corresponding to right, left, up, down).

## • Skip Head:

- Similarly starting from the 512-unit feature vector, it is processed through two subsequent linear layers with 128 and 32 output units, each followed by a Leaky ReLU activation (negative slope 0.1).
- The final linear layer of the skip head maps the 32-unit feature vector to skip\_dim. Each corresponds to the utility of repeating the chosen primitive action for a specific number of steps, from 1 up to skip\_dim.

## L.2 Hyperparameters

General experimental configuration (Table 14) and environment-specific arguments (Table 15) are consistent with the DQN baseline. Algorithm-specific hyperparameters for TempoRL, including those unique to its architecture, are detailed in Table 18.

Table 18: Algorithm Specific Hyperparameters for TempoRL Baseline

Parameter	Default Value / Options
total_timesteps	Environment specific
learning_rate	$[1 \times 10^{-3}, 1 \times 10^{-4}]$
buffer_size	[10000, 50000]
gamma $(\gamma)$	0.99 (discount factor)
$tau(\tau)$	[0.01, 0.005] (target network update rate)
target_network_frequency	[10, 100] (frequency of applying soft target network update)
batch_size	256
start_e	1.0 (starting epsilon for exploration)
end_e	0.1 (ending epsilon for exploration)
exploration_fraction	[0.1, 0.3, 0.5]
learning_starts	1000 (timestep to start learning)
train_frequency	2 (frequency of training the Q-network)
activation_function	Leaky ReLU (negative slope 0.1 for hidden layers)
skip_dim	10 (maximum skip size)
weight_sharing	True (whether to share the first 512-unit layer)

#### L.3 Action Selection and Exploration

TempoRL employs a two-step  $\epsilon$ -greedy strategy for exploration and action selection:

- 1. **Primitive Action Selection**: Given the current state  $s_t$ , a primitive action  $a_t$  (e.g., right, left, up, down) is chosen. With probability  $\epsilon$ ,  $a_t$  is selected randomly from the set of  $n\_output\_actions$ . Otherwise (with probability  $1 \epsilon$ ),  $a_t = \operatorname{argmax}_{a'}Q(s_t, a'; \theta)$ , where  $Q(s_t, \cdot; \theta)$  are the Q-values produced by the action head of the online network.
- 2. **Skip Duration Selection**: Conditioned on the current state  $s_t$  and the chosen primitive action  $a_t$ , a skip duration  $k_t$  (number of times to repeat  $a_t$ , from 1 to  $\mathtt{skip\_dim}$ ) is selected. With probability  $\epsilon$ ,  $k_t$  is chosen randomly from  $\{1,\ldots,\mathtt{skip\_dim}\}$ . Otherwise,  $k_t = \mathrm{argmax}_{k'}Q_{skip}(s_t,a_t,k';\theta_{skip})$ , where  $Q_{skip}(s_t,a_t,\cdot;\theta_{skip})$  are the Q-values for different skip durations produced by the skip head (which might use shared parameters if weight\_sharing = True).

The selected primitive action  $a_t$  is then executed in the environment for  $k_t$  consecutive timesteps. The value of  $\epsilon$  is typically linearly annealed from start\_e to end\_e over exploration\_fraction of total timesteps, as detailed for the DQN baseline (see Section J.3 and Table 18).

#### L.4 Optimizer

Separate Adam optimizers (torch.optim.Adam) are used for the action Q-network parameters and the skip Q-network parameters. The learning rate for both optimizers is controlled by the learning\_rate hyperparameter (see Table 18). Gradients for both networks are clipped (grad\_clip\_val = 40.0).

### L.5 Replay Buffers

TempoRL utilizes two distinct replay buffers with capacity buffer\_size (see Table 18) to store experiences for training its action and skip policies:

- Action Replay Buffer: This is a standard replay buffer (ReplayBuffer from Stable Baselines3) that stores transitions corresponding to individual primitive actions. Each experience tuple is of the form  $(s_t, a_t, r_t, s_{t+1}, d_t)$ , where:
  - $-s_t$ : The state at time t.

- $-a_t$ : The primitive action taken at time t.
- $r_t$ : The reward received at time t+1.
- $-s_{t+1}$ : The state at time t+1.
- $d_t$ : A boolean flag indicating if  $s_{t+1}$  is a terminal state.

Experiences sampled from this buffer are used to train the action Q-network (the action head).

- Skip Replay Buffer: This is a custom replay buffer (referred to as NoneConcatSkipReplayBuffer in the implementation) specifically designed to store experiences related to the execution of multi-step skip actions. Each experience tuple is of the form  $(s_j, k_j, s_{j+k_j}, R_j, d_{j+k_j}, k_j^{len}, a_j^{behav})$ , representing:
  - $s_j$ : The state from which the skip action (repeating  $a_i^{behav}$ ) commenced.
  - $-k_i$ : The selected skip duration (i.e., the 'action' taken by the skip policy).
  - $-s_{j+k_i}$ : The state reached after the primitive action  $a_j^{behav}$  was executed  $k_j^{len}$  times.
  - $R_j$ : The accumulated (and potentially discounted, depending on exact calculation before storage) reward received over the course of the  $k_j^{len}$  steps of the skip.
  - $-d_{j+k_j}$ : A boolean flag indicating if  $s_{j+k_j}$  (the state after the skip) is a terminal state.
  - $-k_j^{len}$  (length): The actual number of steps the primitive action  $a_j^{behav}$  was repeated (this is equivalent to  $k_j$ ).
  - $-a_i^{behav}$ : The underlying primitive action that was chosen to be repeated for  $k_i^{len}$  steps.

Experiences sampled from this buffer are used to train the skip Q-network (the skip head).

## L.6 Training Details

TempoRL involves separate training updates for the action Q-network and the skip Q-network, both utilizing the Huber loss function.

#### L.6.1 Loss Function and Updates

The network parameters are updated by minimizing the Huber loss (This is implemented via torch.nn.SmoothL1Loss) for both action and skip predictions.

• Action Q-Network Update: Experiences  $(s_j, a_j, r_j, s_{j+1}, d_j)$  are sampled from a standard replay buffer. The target value  $y_i^{action}$  is computed using a Double DQN-style approach:

$$y_j^{action} = r_j + \gamma (1 - d_j) Q_{target}(s_{j+1}, \arg \max_{a'} Q(s_{j+1}, a'; \theta); \theta^-)$$

where Q is the online action Q-network with parameters  $\theta$ , and  $Q_{target}$  is its target network with parameters  $\theta^-$ . The loss is then:

$$L_{action}(\theta) = \mathbb{E}_{(s_j, a_j, r_i, s_{j+1}, d_i) \sim \mathcal{B}} \left[ \text{HuberLoss}(y_i^{action} - Q(s_j, a_j; \theta)) \right]$$

• Skip Q-Network Update: Experiences  $(s_j, a_j^{behav}, k_j, R_j, s_{j+k_j}, d_{j+k_j})$  are sampled from a separate replay buffer for skips. Here,  $a_j^{behav}$  is the primitive action executed,  $k_j$  is the skip duration (number of times  $a_j^{behav}$  was repeated),  $R_j$  is the accumulated discounted reward during these  $k_j$  steps, and  $s_{j+k_j}$  is the state after  $k_j$  steps. The target value  $y_j^{skip}$  is calculated as:

$$y_j^{skip} = R_j + \gamma^{k_j} (1 - d_{j+k_j}) Q_{target}(s_{j+k_j}, \arg \max_{a'} Q(s_{j+k_j}, a'; \theta); \theta^-)$$

Note that the future value component  $Q_{target}(s_{j+k_j},\dots)$  uses the main action Q-network and its target, reflecting the value of the optimal next primitive action after the skip concludes. The current prediction is  $Q_{skip}(s_j, a_j^{behav}, k_j; \theta_{skip})$ , where  $Q_{skip}$  has parameters  $\theta_{skip}$  (which may share some parameters with  $\theta$  if weight\_sharing = True). The loss is:

$$L_{skip}(\theta_{skip}) = \mathbb{E}_{(\dots) \sim \mathcal{B}_{skip}} \left[ \text{HuberLoss}(y_j^{skip} - Q_{skip}(s_j, a_j^{behav}, k_j; \theta_{skip})) \right]$$

Both loss functions are optimized using their respective Adam optimizers, and gradients are clipped to prevent large updates.

## L.6.2 Target Network Updates

To stabilize training, target networks are employed.

- A target network  $Q_{target}$  (with parameters  $\theta^-$ ) is maintained for the primary action Q-network Q (parameters  $\theta$ ). This  $Q_{target}$  is always updated using Polyak averaging:  $\theta^- \leftarrow \tau\theta + (1-\tau)\theta^-$ .
- If weight\_sharing is 'False', the distinct skip Q-network  $Q_{skip}$  (parameters  $\theta_{skip}$ ) has its own separate target network,  $Q_{skip\_target}$  (parameters  $\theta_{skip}^-$ ). This  $Q_{skip\_target}$  is similarly updated using Polyak averaging with  $\theta_{skip}$  and  $\theta_{skip}^-$ .
- If weight\_sharing is 'True', the parameters of the skip mechanism are part of the overall network structure whose online parameters are  $\theta$  (which includes the shared trunk and potentially specific skip head layers not part of the action head). In this scenario, a separate Polyak update for a distinct  $Q_{skip\_target}$  is not performed; target values for the skip component's loss are derived using  $Q_{target}$  for estimating future state-action values, as shown in the skip target formula.

The soft update rate  $\tau$  and the update frequency target\_network\_frequency are specified in Table 18.

# M Proto Sequence Decoder (PSD) Implementation Details

This section outlines the architecture and configuration of the Proto-Sequence Decoder (PSD) module used in our model. It describes the decoder network structure, training objectives, latent space regularization, sequence reconstruction process, and other relevant design choices that enable effective decoding of proto-sequence embeddings into action-sequence.

## M.1 Training Dataset Setup

To ensure that the model learns from meaningful and structured data rather than arbitrary noise, we constructed the training dataset according to the following constraints:

- 1. Sequence Length Constraint: All action sequences have lengths between 1 and  $L_{max} = 10$  steps.
- 2. Action Diversity Constraint: Each sequence includes at most two distinct action types. For example, valid sequences include [up, up] or [up, left], whereas a sequence like [up, left, down] is considered invalid.
- 3. Switch Constraint: Each sequence may contain at most one switch between action types. For instance, [up,left] is allowed, but [up,left,up] is not.
- 4. **Directional Conflict Constraint**: Sequences cannot include **opposite directions**, such as both up and down, or both left and right.
- 5. Avoidance of loops.

Following these criteria, we generated a total of **400 unique action sequences**. These sequences were one-hot encoded and padded with EOS tokens to a fixed length of  $L_{max}$  to suit VAE input requirements. These sequences form the basis of the training data for the Proto-Sequence Decoder (PSD).

### M.2 Model Architecture

The Proto-Sequence Decoder (PSD) is implemented as a Variational Autoencoder (VAE) designed to map proto-sequence embeddings into action sequences of varying length. It is trained using reconstruction loss combined with a Kullback-Leibler divergence (KLD) regularization toward a standard Gaussian prior. The decoder operates on flattened one-hot sequence representations of actions and outputs reconstructed sequences over a predefined action vocabulary.

## M.2.1 Input Representation

Each action in the sequence is represented as a one-hot vector over a vocabulary of size  $n_{\text{words}} = 5$  corresponding to {up, down, right, left, eos\_token}. The decoder models sequences of up to input\_length = 10 actions, resulting in an input vector of dimension  $10 \times 5 = 50$ .

#### M.2.2 Encoder Network

The encoder receives a flattened 50-dimensional input vector and passes it through a series of fully connected layers:

## 1. Linear Layer 1:

Input size: 50Output size: 32

Normalization: InstanceNorm1dActivation: LeakyReLU (slope 0.2)

## 2. Linear Layer 2:

• Output size: 16

• Normalization: InstanceNorm1d

• Activation: LeakyReLU

## 3. Linear Layer 3:

• Output size: 16

• Normalization: InstanceNorm1d

• Activation: Tanh

The output is then projected into two parallel linear layers to produce the latent mean  $\mu \in \mathbb{R}^{16}$  and log-variance  $\log \sigma^2 \in \mathbb{R}^{16}$ . A latent sample z is drawn using the reparameterization trick:  $z = \mu + \sigma \cdot \epsilon$ , where  $\epsilon \sim \mathcal{N}(0,1)$ .

#### M.2.3 Decoder Network

The sampled latent vector  $z \in \mathbb{R}^{16}$  is decoded through a symmetric feedforward network:

## 1. Linear Layer 1:

• Output size: 16

• Normalization: InstanceNorm1d

• Activation: LeakyReLU

## 2. Linear Layer 2:

• Output size: 32

• Normalization: InstanceNorm1d

Activation: LeakyReLU

## 3. Linear Layer 3:

x§

• Output size: 50 (reconstructed sequence)

• Normalization: InstanceNorm1d

• Final Activation: Sigmoid (applied element-wise)

# M.3 Training Objective

The PSD is optimized using a combination of:

- Reconstruction Loss: Binary cross-entropy loss between the input sequence and its reconstruction, normalized by sequence length.
- KL Divergence Loss: Encourages the latent distribution to match a unit Gaussian prior.

All input sequences are EOS-padded to the maximum length of 10 to ensure uniform input dimensionality across batches.

## M.4 Hyperparameters

The training of the Proto-Sequence Decoder (PSD) is governed by a set of fixed hyperparameters, detailed in Table 19. These parameters control aspects such as optimization, batch processing, and reproducibility.

Table 19: Proto-Sequence Decoder (PSD) Training Hyperparameters

Parameter	Value
train_on_entire_dataset	True
seed	42
optimizer	Adam
optimizer_learning_rate	$1 \times 10^{-4}$
optimizer_weight_decay	$1 \times 10^{-3}$
batch_size	32

## M.5 Optimizer

The PSD is trained using the Adam optimizer (torch.optim.Adam). The learning rate is controlled by the learning\_rate hyperparameter (see Table 19).

## M.6 Training Details

#### M.6.1 Loss Function

The Proto-Sequence Decoder (PSD) parameters are updated by minimizing a combined loss:

$$L = L_{\rm rec} + L_{\rm KL}$$

where:

•  $L_{\rm rec}$  is the label-smoothed binary cross-entropy over the reconstructed sequence:

$$L_{\text{rec}} = -\frac{1}{T} \sum_{t=1}^{T} \left[ y_t \log \hat{y}_t + (1 - y_t) \log(1 - \hat{y}_t) \right],$$

with  $y_t$  replaced by  $\tilde{y}_t = y_t(1-\epsilon) + \frac{\epsilon}{2}$ ,  $\epsilon = 0.1$ , and T = 10 is the sequence length.

 $\bullet$   $L_{\mathrm{KL}}$  is the Kullback–Leibler divergence between the approximate posterior and a unit Gaussian:

$$L_{\mathrm{KL}} = D_{\mathrm{KL}} (\mathcal{N}(\mu, \sigma^2) \| \mathcal{N}(0, I)).$$

## M.6.2 Training Procedure

- Maximum Steps: Train for up to 20,000 epochs.
- Batch Composition: Split sequences into
  - Short ( $\leq 5$  actions) and
  - Long (> 5 actions),

and sample each batch with a 50/50 ratio of short and long sequences.

#### M.7 Evaluation

The Proto-Sequence Decoder's (PSD) performance was assessed on the entire training set using two key metrics. Evaluations were conducted every 50 epochs, and the checkpoint yielding the highest exact match accuracy was retained. After 20,000 epochs, the following results were achieved:

- Exact Match Accuracy: This metric measures the proportion of sequences reconstructed with zero errors. The PSD achieved an Exact Match Accuracy of 0.978.
- Per-Step Accuracy: This metric calculates the fraction of correctly reconstructed actions across all positions within the sequences. The PSD achieved a Per-Step Accuracy of 0.99.

## M.8 Visualization of the Learned Embedding Space

Figure 6 illustrates the two-dimensional t-SNE projection of the learned proto-action-sequence embeddings. Each point corresponds to one sequence from the dataset, with colors representing the effective sequence length. As can be seen, sequences with similar structural properties tend to form dense clusters, indicating that the embedding space preserves meaningful relationships between sequences.

A particularly noteworthy observation is the position of the red star, which represents a previously unseen sequence not included during training. This sequence is located within the cluster of its closest structural neighbors, suggesting that the learned representation generalizes effectively to new data. In other words, the embedding model is able to position novel sequences near the most similar examples from the training set, supporting its potential for robust retrieval, similarity search, and downstream predictive tasks.

## N GPS (Generative Proto Sequence) Implementation Details

This section outlines the architecture and configuration of our GPS (Generative Proto Sequence) method. GPS is an actor-critic based algorithm where the actor network generates a latent representation, termed "proto-sequence." This proto-sequence is then processed by a pre-trained generative decoder model (PSD) to produce a sequence of discrete actions. A critic network evaluates this action sequence to guide learning. The subsequent subsections detail neural network structures, hyperparameter settings, exploration strategy, optimization methods, and other relevant training aspects.

#### N.1 Model Architecture and State Representation

The GPS agent consists of three primary neural network components: an Actor, a Critic, and a pre-trained Decoder. It processes observations from the environment. For details on the specific state representation, input shape (e.g., height, width, channels), and preprocessing, please refer to the "State Representation" within Section 4.3.

### N.1.1 Convolutional Neural Network (CNN) Part

Both the Actor and Critic networks use separate but identical Convolutional Neural Network (CNN) architecture to extract features from the input observation. The CNN component used for initial feature

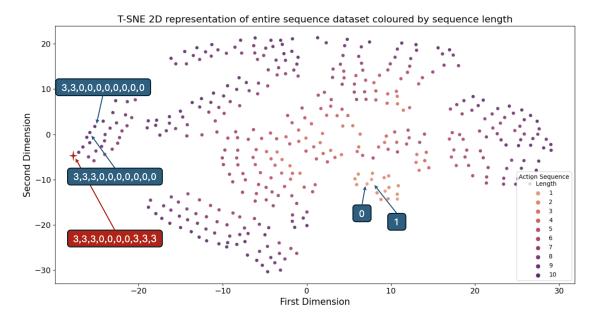


Figure 6: T-SNE 2-D projection of the proto-action-sequence embedding space. The map contains embeddings for the 400 original sequences in the dataset together with one previously unseen sequence (red star). Points are colored by their effective sequence length, and sequences with similar structure form tight neighborhoods. The unseen sequence falls naturally inside the cluster of its closest structural neighbors, showing that the learned representation places new, out-of-training sequences adjacent to the most similar known examples.

extraction is identical to the one in the DQN baseline. For details on its architecture (number of layers, channels, kernel sizes, strides, and activations), please refer to Appendix J.1. The output of this CNN part is a flattened 1D vector of size  $n\_flatten$ .

Note: in our current implementation, the actor and critic networks each have their own separate CNN state feature extractors. In future work, we plan to explore a shared CNN feature extraction architecture as implemented in TempoRL, which could potentially improve computational efficiency and state representation learning.

#### N.1.2 Actor Network

The Actor network takes the extracted features from the CNN and produces a proto-sequence embedding.

- Input: The flattened feature vector  $n_{\text{features}}$  from the CNN is concatenated with a positional encoding. The embedding dimension of this positional encoding is specified by pe\_embedding\_dim.
- Architecture: The combined features are processed through a series of fully connected linear layers, defined by the actor\_linear\_layers parameter (e.g., [512, 128, 32]). The activation function for these hidden layers is specified by actor\_linear\_layers\_activation\_function which is "leaky\_relu" with negative slope 0.1.
- Output: The actor generates a single proto-sequence embedding. This embedding is a vector of size actor\_n\_output\_channels and serves as input to the Decoder PSD network.

### N.1.3 Position Encoding in Proto-Sequence Generation

To enhance the expressiveness of our action sequence generation, we incorporate positional encoding within the Actor network. This technique, inspired by transformer architectures (Dosovitskiy et al., 2020), helps the

Actor generate more contextually aware proto-sequence embeddings by providing explicit spatial information about the agent and goal positions.

**Implementation Details.** Our positional encoding implementation combines both row and column information for each grid cell in the observation space:

1. We create sinusoidal encodings for both dimensions (height and width) separately:

```
pe_row[:, 0::2] = torch.sin(position_row * div_term)
pe_row[:, 1::2] = torch.cos(position_row * div_term)
```

- 2. These encodings are combined into a unified representation where the first half of each cell's embedding encodes its row position and the second half encodes its column position.
- 3. During forward passes, we extract the agent and goal positions from the observation and retrieve their respective positional encodings:

```
agent_pe = self.position_encoding[agent_row, agent_col]
goal_pe = self.position_encoding[goal_row, goal_col]
```

4. These position-specific features are concatenated with the CNN-extracted features before being processed by the linear layers of the Actor.

Motivation and Benefits. Integrating positional encoding within the Actor network provides several advantages:

- 1. **Enhanced Spatial Reasoning**: By explicitly encoding agent and goal positions, the Actor can better understand spatial relationships, which is crucial for navigation tasks.
- 2. **Improved Exploration Early in Training**: The position encodings enable the generation of more diverse proto-sequence embeddings in the initial training phases, facilitating better exploration before the CNN features become well-trained.
- 3. **Direction-Aware Sequence Generation**: The relative positions of agent and goal inform the Actor about the general direction of movement required, allowing it to generate more purposeful action sequences even with limited experience.
- 4. **Invariance to Visual Feature Quality**: Especially early in training when the CNN features may be unreliable, position encodings provide a stable signal that enables meaningful action sequence generation.

Our manual investigations and targeted experiments suggest that incorporating position encodings enhances the model's capabilities in several ways. We observed that the position-enriched Actor generates protosequence embeddings with greater contextual awareness of spatial relationships, which in turn produces more diverse and situation-appropriate action sequences. Without position encoding, the proto-sequence embeddings tended to cluster more closely in the latent space, resulting in less differentiated action patterns. This difference was particularly evident in larger and more complex maze environments, where the position-encoded model demonstrated an improved ability to generate directionally coherent sequences that efficiently navigated toward goals. The positional information appears to provide a structural prior that helps the Actor formulate meaningful navigation strategies even before the CNN features are fully refined through training.

Quantitative Ablation Study. To provide empirical evidence for the benefits of positional encoding, we conducted a controlled ablation study comparing GPS with and without positional encoding on two challenging environments. Table 20 presents the comparative results.

The ablation results provide strong quantitative evidence for positional encoding's critical importance. GPS with positional encoding achieves convergence (ASR > 0.9) in 300K-1.5M steps, while the variant without positional encoding fails to converge after 1.5M steps. Performance gaps are substantial: 44 percentage

Table 20: GPS with vs. without Positional Encoderate	ding
--	------

Environment	Positional Encoding	ASR at Training Steps				ASR Converge		
		100K	500K	1M	1.5M	>0.9 Step	PER	$\operatorname{SGF}$
16x16_obstacles_15%	With Without	$0.22 \\ 0.15$	0.96 0.38	$0.99 \\ 0.51$	$0.96 \\ 0.52$	<b>300k</b> >1.5M	<b>0.72</b> 0.62	<b>10.8</b> 17.67
24x24_obstacles_15%	With Without	$0.12 \\ 0.07$	$0.45 \\ 0.12$	$0.80 \\ 0.13$	0.91	<b>1.5M</b> >1.5M	0.48 <b>0.54</b>	<b>6.34</b> 16.91

in the  $16\times16$  with 15% obstacles environment (0.96 vs 0.52 final ASR) and 78 percentage in the  $24\times24$  with 15% obstacles environment (0.91 vs 0.13 final ASR). Path efficiency results show mixed patterns, with positional encoding improving PER by 0.10 in the  $16\times16$  with 15% obstacles environment (0.72 vs 0.62) but showing lower PER in the  $24\times24$  with 15% obstacles environment (0.48 vs 0.54). However, this apparent PER advantage for the variant without positional encoding is misleading, as its extremely low ASR (0.13) indicates fundamental task learning failure—it likely only succeeds on the easiest instances, artificially inflating the path efficiency metric. Both environments show substantial reductions in sequence generation frequency with positional encoding: 39% reduction in  $16\times16$  (10.8 vs 17.67 SGF) and 63% reduction in  $24\times24$  (6.34 vs 16.91 SGF), indicating more coherent action sequences requiring fewer decision points.

## N.1.4 Critic Network

The Critic network estimates the Q-value of a state and a decoded action sequence.

## • Input:

- The flattened feature vector  $n_{\text{features}}$  from the CNN, representing the current state.
- The action sequence generated by the Decoder from the actor's proto-sequence. The representation used for this action sequence is ACTION\_SEQ\_AS\_ONE\_HOT. Namely, the input is a tensor where each action in the sequence of length n\_actions\_in\_seq is one-hot encoded. Shorter sequences are padded to this length using an End-of-Sequence (EOS) token. Each one-hot vector has a dimension equal to the action\_space\_size plus one (for the EOS token). Consequently, the total input dimension for the action sequence part, action\_seq\_dim, is n\_actions\_in\_seq × (action\_space\_size + 1).
- Architecture: The inputs are processed through a series of fully connected linear layers, defined by critic\_linear\_layers (e.g., [512, 128, 32]). The activation function is specified by critic linear layers activation function which is "leaky relu" with negative slope 0.1.
- Output: A single scalar Q-value. The output Q-value can be optionally clipped between min\_qf\_value and max\_qf\_value.

## N.1.5 Decoder Network

A pre-trained generative model, specifically a Variational Autoencoder, acts as the Decoder (PSD).

- Loading: The Decoder is loaded from a pre-trained model specified by decoder\_model\_path.
- Input: The proto-sequence embedding (size actor\_n\_output\_channels) generated by the Actor.
- Output: A sequence of n\_actions\_in\_seq discrete actions. Each action is selected from a vocabulary of action\_space\_size primitive actions plus an end\_of\_sequence\_token token.

• Generation: The Decoder can use Gumbel-Softmax for differentiable sampling if use\_gumble\_in\_decoder is true, or a deterministic argmax with a Straight-Through Estimator otherwise. See explanation in Section 4.2.

## N.2 Hyperparameters

The GPS method is configured by a wide range of hyperparameters. General experiment settings and environment-specific configurations are typically managed as detailed for the DQN baseline (see Tables 14 and 15). Key algorithm-specific hyperparameters for GPS are listed in Tables 21, 22.

Table 21: Algorithm Specific Hyperparameters for GPS

Parameter	Value / Options / Description			
total_timesteps	Env. specific			
actor_learning_rate	Actor LR (e.g., $1 \times 10^{-4}$ ).			
critic_learning_rate	Critic LR (e.g., $1 \times 10^{-4}$ ).			
buffer_size	[10000, 50000]			
gamma $(\gamma)$	Discount factor.			
$tau\left( au ight)$	Target net. soft update rate (e.g., 0.005).			
batch_size	Experiences per train step.			
learning_starts	Timestep train begins.			
actor_policy_frequency	Actor net. update freq. rel. to Critic (e.g., 2).			
actor_target_network_frequency	Target Actor net. update freq. (e.g., 10 steps).			
<pre>critic_target_network_frequency</pre>	Target Critic net. update freq. (e.g., 10 steps).			
start_e	Initial $\epsilon$ for $\epsilon$ -greedy (e.g., 1.0).			
end_e	Final $\epsilon$ value (e.g., 0.1).			
total_steps_e	Timesteps for $\epsilon$ annealing (e.g., 15000).			
sub_sequences_move_start_point	Boolean; varying sub-seq. start points.			
sub_sequences_move_end_point	Boolean; varying sub-seq. end points.			
sub_sequences_min_jump_move_start_point	Min. jump for start-moved sub-seq. gen.			
sub_sequences_min_jump_move_end_point	Min. jump for end-moved sub-seq. gen.			
every_one_step_transition_to_buffer	Boolean; all single-step trans. stored.			
actor_n_output_channels	Proto-action-seq. embed dim. (e.g., 16).			
actor_linear_layers	Actor MLP hidden layer sizes (e.g., [512, 128, 32]).			
actor_linear_layers_activation_function	Actor MLP activation (e.g., "leaky_relu").			
actor_weight_decay	Actor L2 reg. strength (e.g., $1 \times 10^{-2}$ ).			
pe_embedding_dim	Positional encoding dim. in Actor (e.g., 16).			
critic_linear_layers	Critic MLP hidden layer sizes (e.g., [512, 128, 32]).			
critic_linear_layers_activation_function	Critic MLP activation (e.g., "leaky_relu").			
critic_linear_layers_activation_lunction critic_weight_decay	Critic L2 reg. strength (e.g., $1 \times 10^{-3}$ ).			
decoder_model_path	Path to pre-trained VAE Decoder model.			
n_actions_in_seq	Decoder action seq. length (e.g., 10).			
<b>-</b>	Unique primitive actions in env. (e.g., 4).			
action_space_size end_of_sequence_token	Decoder integer token for end of seq. (e.g., 4).			
use_gumble_in_decoder	Boolean; Decoder uses Gumbel-Softmax.			
re= Romone Tulacoder	boolean; Decoder uses Gumber-Solumax.			

## N.3 Reward Function and Empty Sequence Handling

The GPS agent utilizes a structured reward function designed to encourage efficient navigation while penalizing inefficient or invalid behaviors. The reward function is defined as:

$$R = r_{\rm goal} - \frac{1}{l_{\rm max}} \times n_{\rm valid} - \frac{3}{l_{\rm max}} \times n_{\rm invalid}$$

Table 22: GPS Hyperparameters Settings

Parameter	Meters Settings Value		
total_timesteps	Environment specific		
buffer_size	[10000, 50000]		
gamma	0.99		
tau	[0.01, 0.005]		
batch_size	256		
start_e	1		
end_e	0.1		
learning_starts	1000		
actor_learning_rate	[1e-3, 1e-4, 1e-5]		
critic_learning_rate	1e-04		
actor_policy_frequency	2		
sub_sequences_move_start_point	TRUE		
sub_sequences_move_end_point	TRUE		
sub_sequences_min_jump_move_start_point	1		
sub_sequences_min_jump_move_end_point	1		
every_one_step_transition_to_buffer	TRUE		
actor_target_network_frequency	[10, 100]		
<pre>critic_target_network_frequency</pre>	10		
total_steps_e	15000		
actor_n_output_channels	16		
actor_linear_layers	["512, 32", "512, 128, 32"]		
actor_linear_layers_activation_function	leaky_relu (negative slope 0.1)		
actor_weight_decay	1e-04		
<pre>pe_embedding_dim</pre>	128		
critic_linear_layers	[512, 128, 32]		
critic_linear_layers_activation_function	leaky_relu (negative slope 0.1)		
critic_weight_decay	1e-04		
max_level	Environment specific		
start_level	Environment specific		
use_gumble_in_decoder	TRUE		

where  $r_{\text{goal}}$  is 1 if the agent reached the goal (0 otherwise),  $l_{\text{max}}$  is the maximal start-goal distance acting as a regularizer,  $n_{\text{valid}}$  is the number of valid actions taken, and  $n_{\text{invalid}}$  is the number of invalid actions (e.g., bumping into a wall).

To handle cases during training where the actor generates proto-sequences that are decoded as empty sequences (i.e., where the PSD outputs the EOS token as the first action), we apply a harsh penalty of -20. For these instances, we also hard-code the action to be 1 (DOWN) to ensure the agent always takes some action. This direct negative reinforcement was found to be highly effective in guiding the actor to produce valid proto-sequence embeddings that decode into meaningful action sequences. Without this penalty, the actor might frequently produce embeddings that map to empty, significantly hampering exploration and learning progress. This approach provides a clear signal to the actor network about the importance of generating proto-sequences that translate to substantive action sequences, accelerating the learning process and improving the overall stability of training. Empirically, we observed that this simple yet effective mechanism substantially reduced the occurrence of empty sequences.

## N.4 Action Selection and Exploration

At each decision step, the Actor network generates a proto-sequence embedding. The Decoder then translates this embedding into a corresponding action-sequence of primitive actions. The agent employs an  $\epsilon$ -greedy exploration strategy:

- With probability  $\epsilon_t$  (where  $\epsilon_t$  anneals from start\_e to end\_e over total\_steps\_e steps): An exploratory action sequence is selected. This sequence is typically chosen randomly from a pre-defined set of valid action sequences see Section 4.3 for more details.
- With probability  $1 \epsilon_t$  (exploitation): The action sequence generated by the Actor-Decoder pipeline is used for execution.

The value of  $\epsilon_t$  is linearly annealed:

$$\epsilon_t = \max \Big(end\_e, start\_e - (start\_e - end\_e) \cdot (current\_step/total\_steps\_e) \Big)$$

The chosen action sequence is subsequently trimmed using the end\_of\_sequence\_token before execution.

## N.5 Optimizer

The Actor and Critic networks are trained using separate Adam optimizers (torch.optim.Adam).

- The Actor's optimizer is configured with a learning rate of actor\_learning\_rate and applies L2 weight decay with a coefficient of actor\_weight\_decay.
- The Critic's optimizer uses a learning rate of critic\_learning\_rate and L2 weight decay with a coefficient of critic\_weight\_decay.

#### N.6 Replay Buffer

A replay buffer (ReplayMemory) with a capacity of buffer\_size stores past experiences. Each stored transition typically includes: the current state observation  $(s_t)$ , the next state observation  $(s_{t+1})$ , the selected action sequence (act\_seq<sub>t</sub>), the received reward  $(r_t)$ , a terminal flag  $(d_t)$ , and the actor's proto-action-sequence embedding that generated act\_seq<sub>t</sub> (emb<sub>t</sub>). The system may also store sub-sequences derived from executed plans if parameters such as push\_every\_one\_step\_transition\_to\_buffer, push\_sub\_sequences\_to\_buffer\_move\_start\_point, and push\_sub\_sequences\_to\_buffer\_move\_end\_point are enabled, potentially enriching the diversity of experiences in the buffer.

## N.7 Training Details

Network training commences after learning\_starts timesteps have been collected. Updates are performed using batches of batch\_size experiences sampled from the replay buffer.

## N.7.1 Critic Network Update

The Critic network parameters ( $\theta_C$ ) are updated by minimizing the Mean Squared Error (MSE) loss:

$$L(\theta_C) = \mathbb{E}_{(s, \text{act\_seq}, r, s', d, \text{emb}) \sim \mathcal{B}} \left[ \left( Q(s, \text{emb}, \text{act\_seq}; \theta_C) - y_t \right)^2 \right]$$

The target Q-value  $y_t$  is computed using the target Actor (target\_actor\_network) and target Critic (target\_critic\_network) networks to ensure stability:

$$y_t = r + \gamma (1 - d) Q_{\text{target}}(s', \text{Actor}_{\text{target}}(s'), \text{Decoder}(\text{Actor}_{\text{target}}(s')); \theta_C^-)$$

where  $Actor_{target}(s')$  is the proto-sequence embedding from the target actor for state s',  $Decoder(\cdot)$  converts it to an action sequence, and  $\theta_C^-$  are the parameters of the target critic. The Q-values from the target critic can be clipped using min\_qf\_value and max\_qf\_value.

## N.7.2 Actor Network Update

The Actor network parameters  $(\theta_A)$  are updated with a frequency of actor\_policy\_frequency (delayed policy update). The goal is to adjust the actor's parameters to produce a proto-sequence that leads to a higher Q-value as estimated by the current Critic. For a sampled batch of states, the Actor generates a proto-sequence embedding. This is decoded into an action sequence, which is then evaluated by the online Critic network  $Q(\cdot; \theta_C)$ . The actor loss is designed to maximize this Q-value:

$$L(\theta_A) = -\mathbb{E}_{s \sim \mathcal{B}} [Q(s, \text{Actor}(s), \text{Decoder}(\text{Actor}(s)); \theta_C)]$$

#### N.7.3 Target Network Updates

Separate target networks are maintained for both the Actor (Actor<sub>target</sub> with parameters  $\theta_A^-$ ) and the Critic ( $Q_{\text{target}}$  with parameters  $\theta_C^-$ ). Their parameters are updated using Polyak averaging with the parameters of their corresponding online networks ( $\theta_A, \theta_C$ ):

$$\theta^- \leftarrow \tau\theta + (1-\tau)\theta^-$$

The soft update rate is  $\tau$ . Target network updates for the Actor occur every actor\_target\_network\_frequency steps, and for the Critic every critic\_target\_network\_frequency steps.

## N.7.4 Training Procedure Summary

- Initialization: Networks and target networks are initialized. Replay buffer is empty.
- Data Collection: Agent interacts with the environment using the action selection strategy (Section N.4), storing experiences  $(s_t, s_{t+1}, \text{act\_seq}_t, r_t, d_t, \text{emb}_t)$  in the replay buffer.
- Learning Phase (after learning\_starts steps):
  - 1. Sample a batch\_size of experiences from the replay buffer.
  - 2. Update Critic network parameters by minimizing the MSE loss with the computed TD targets.
  - 3. Periodically (every actor\_policy\_frequency steps), update Actor network parameters to maximize the Q-value of the generated sequence as estimated by the Critic.
  - 4. Periodically (every actor\_target\_network\_frequency and critic\_target\_network\_frequency steps respectively), update target Actor and target Critic networks using Polyak averaging.

#### N.8 Evaluation

The performance of the GPS agent is assessed periodically during training and/or at the end of the training process.

- Frequency: Evaluations on validation datasets typically occur every val\_eval\_freq steps, and on subsets of the training data every train\_eval\_freq steps. Less frequent evaluations may occur on a dedicated test dataset (e.g., every eval\_test\_dataset\_training\_freq steps) or at the end of training.
- Method: During evaluation, the actor generates a proto-sequence, the decoder converts it to an action sequence, and this sequence is executed. The Decoder may operate in a deterministic mode (deterministic\_inference = True).
- Metrics: Standard reinforcement learning metrics are logged, such as mean episodic return and success rate. Additional metrics might include the average number of decoder generations per episode or properties of the generated action sequences.
- Model Saving: If save\_model is enabled, the best performing models (actor and critic) are saved based on criteria defined by save\_model\_strategy (e.g., best success rate or mean reward on the validation set).

## O Maze Evaluation Environments Benchmark

In our research, our evaluation environments consist of procedurally generated mazes with varying structures and complexity. We utilized synthetic maze environments created using Large Language Models (LLMs) to ensure unbiased benchmark construction. The following details our approach to maze generation for the different environment types used in our experiments.

## 0.1 Synthetic Maze Generation Process

We generated our maze environments using LLM. For each maze type (rooms, obstacles, and corridors), we provided specific prompts instructing the LLM to generate Python code that would create the maze environments according to our requirements. Importantly, all maze generation was performed programmatically without manual intervention, ensuring reproducibility and eliminating human bias. We specifically used the OpenAI o3-mini model for code generation, which was instructed to create five different variants for each maze type. To avoid experimenter bias in seed selection, we also employed an LLM to generate code for choosing random seeds:

```
import hashlib
def get_consistent_seed():
"""
Selects a seed number from a list consistently across multiple runs.
Returns:
int: The selected seed number.
"""
seed_list = [42, 1234, 9999, 2024, 2025]

Create a hash of the function name to ensure consistency
hash_object = hashlib.sha256(b'get_consistent_seed')
hash_value = int(hash_object.hexdigest(), 16)

Use the hash value to select a seed from the list
seed_index = hash_value % len(seed_list)
return seed_list[seed_index]

Get the consistent seed
seed = get_consistent_seed()
print(f"Selected seed: {seed}")
```

The code generated by the LLM for each maze type is available in our code repository. This approach ensured that the maze generation process was fully automated and free from experimenter bias, providing a consistent and fair benchmark for evaluating our GPS algorithm against the baselines.

## 0.2 Maze Type Generation Prompts

For each maze type, we provided detailed prompts to the LLM to guide the generation process:

# 0.2.1 Obstacles Maze Prompt

#### Maze Generation with 15% Obstacles

Write a Python script that programmatically generates five distinct  $16 \times 16$  mazes with randomly placed obstacles. The grid follows these rules:

## **Grid & Output Format:**

• The maze is a  $16 \times 16$  grid.

- Each cell is either open space (0) or a wall/obstacle (1).
- The outermost border (first and last rows and columns) must remain walls (1).
- 15% of the inner cells (excluding the border) should be randomly chosen as obstacles (1), while the remaining are open spaces (0).
- The final maze should be output as a Python dictionary:

Replace the dots with appropriate 0s (open spaces) and 1s (walls/obstacles).

## Maze Generation Method: Initialize the Maze:

- Create a  $16 \times 16$  grid where every cell is an **open space** (0).
- The outer border (first and last rows/columns) must always remain walls (1).
- The inner 14×14 area (excluding the border) will contain open spaces (0) and obstacles (1).

#### Randomly Place Obstacles:

- 15% of the inner  $14 \times 14$  cells should be converted into obstacles (1).
- The placement of these obstacles should be **random**.
- Ensure that at least one path remains between any two open spaces for potential connectivity.

# Generate Five Distinct Mazes:

- Use different random seeds to create five unique mazes.
- Ensure that each maze has exactly 15% obstacles inside the inner area.

## 0.2.2 Rooms Maze Prompt

You are a maze designer responsible for enhancing a predefined base maze structure for a navigation simulation. Your task is to decide where entrances should be located while ensuring the maze meets the following requirements:

# Maze Design Requirements Base Maze Structure:

- The maze is predefined and consists of a  $16 \times 16$  grid.
- The base structure must remain intact, but you will determine the placement of the entrances and ensure connectivity between all rooms and open spaces.

### **Entrance and Room Connectivity:**

• The maze is divided into four equal-sized quadrants (rooms) separated by walls.

- Each room must have only **two entrances of size 1**. Use a seeded random choice for entrance placement.
- Entrances should be placed strategically to ensure the maze is **fully connected**, meaning an agent can navigate between any two open cells (0) using **up**, **down**, **left**, **or right** movements.
- Passageways between rooms must be narrow and preserve the integrity of the maze's challenge.

## **Obstacle Coverage:**

- Obstacles (1) must make up 5% of the total grid ( $\approx 13$  cells).
- Don't consider the obstacles that are part of the maze's frame.
- Obstacles may be added or removed **within constraints** to maintain connectivity and alignment with the entrance placement.
- Don't place obstacles in nearby squares close to any entrance. Make sure that an obstacle doesn't block any entrance.
- Use a seeded random choice for obstacles placement.

## Reproducibility:

• Use a **specific random seed** to ensure the design is reproducible.

#### **Output Format:**

- Generate code for creating the maze.
- Generate the maze as a Python dictionary with a key (e.g., 'maze') and represent each row as a binary string.
- Output 5 mazes using different seeds and make sure that obstacles don't block entrances.

## Base Maze Layout: The base structure is as follows:

```
{
   'base maze': [
        "1111111111111111",
        "1000000100000001",
        "1000000100000001",
        "1000000100000001",
        "1000000100000001",
        "1000000100000001",
        "1000000100000001"
        "1111111111111111",
        "1000000100000001",
        "1000000100000001",
        "1000000100000001",
        "1000000100000001",
        "1000000100000001",
        "1000000100000001",
        "1000000100000001",
        "1111111111111111"
     ]
}
```

**Design Task:** Modify the maze by:

- Placing **entrances** in the walls separating the quadrants.
- The logic should be based on seed for determining the entrances position and obstacles positions.
- Ensuring the maze is **fully connected**.
- Making any minor adjustments to obstacles (1) to meet connectivity and percentage requirements.
- Verify by code that each room must have only **two entrances of size 1**.
- Verify by code that each obstacle doesn't block any entrance.

Output the modified maze as a Python dictionary with the format below:

```
{
    'maze': [
          "updated_row_1",
          "updated_row_2",
          ...,
          "updated_row_16"
]
```

## General Steps:

- Choose two entrances for each room by selecting a random square on each wall.
- Randomly select a room and place obstacles within it until the obstacle budget is reached.

## 0.2.3 Corridors Maze Prompt

Write a Python script that programmatically generates five distinct  $16 \times 16$  mazes. In each maze, start with a grid completely filled with wall cells (represented by 1), then carve out corridors by selecting one or more vertical lines and one or more horizontal lines to convert wall cells to open cells (represented by 0). The corridors will be 1-cell-wide, and they must intersect so that every open cell is reachable from any other via up, down, left, and right moves. The outer border of the maze should always remain as walls.

### **Grid & Output Format:**

- The maze is a  $16 \times 16$  grid.
- Each cell is either open (0) or a wall (1).
- The outermost border (first and last rows and columns) must remain walls.
- The final maze should be output as a Python dictionary:

]

}

Replace the dots with the appropriate 0s and 1s as per the carved corridors.

## Corridor Carving Method: Initialize the Maze:

• Create a  $16 \times 16$  grid where every cell is a wall (1), with the outer border fixed as walls.

#### **Select Corridor Lines:**

- Vertical Corridors: Choose 2 up to 4 vertical columns (not including the outer borders) that will serve as corridors.
- Restriction: Ensure that no consecutive vertical columns are selected—there must be at least one wall column between any two chosen corridor columns.
- Horizontal Corridors: Choose 2 up to 4 horizontal rows (again, not including the outer borders) that will serve as corridors.
- Restriction: Ensure that no consecutive horizontal rows are selected—there must be at least one wall row between any two chosen corridor rows.
- These lines will form a network of corridors that cross each other.

## Carve the Corridors:

- For each selected vertical column, change all cells in that column (except the outer border) from 1 (wall) to 0 (open space).
- Similarly, for each selected horizontal row, change all cells in that row (except the outer border) from 1 to 0.
- The intersections of these corridors (where a selected vertical column crosses a selected horizontal row) will naturally be open, ensuring connectivity.

#### **Ensure Full Connectivity:**

- The chosen vertical and horizontal corridors should intersect, guaranteeing that every open cell (in the corridors) is reachable from any other open cell.
- Optionally, you can add additional corridor "branches" (by clearing cells adjacent to the main corridors) to create a more interesting maze layout, as long as all open cells remain interconnected.

#### Randomness:

• Generate five distinct mazes by using different random seeds and varying the selected vertical and horizontal corridor positions.