# Scaling Value Iteration Networks to 5000 Layers for Extreme Long-Term Planning

Yuhui Wang*[1], Qingyuan Wu*[2], Weida Li[3], Dylan R. Ashley[4], Francesco Faccio[4], Chao Huang[2], Jürgen Schmidhuber[1,4,5]

[1] Center of Excellence in GenAI, King Abdullah University of Science and Technology (KAUST), [2] The University of Southampton, [3] National University of Singapore, [4] The Swiss AI Lab IDSIA/USI/SUPSI, [5] NNAISENSE
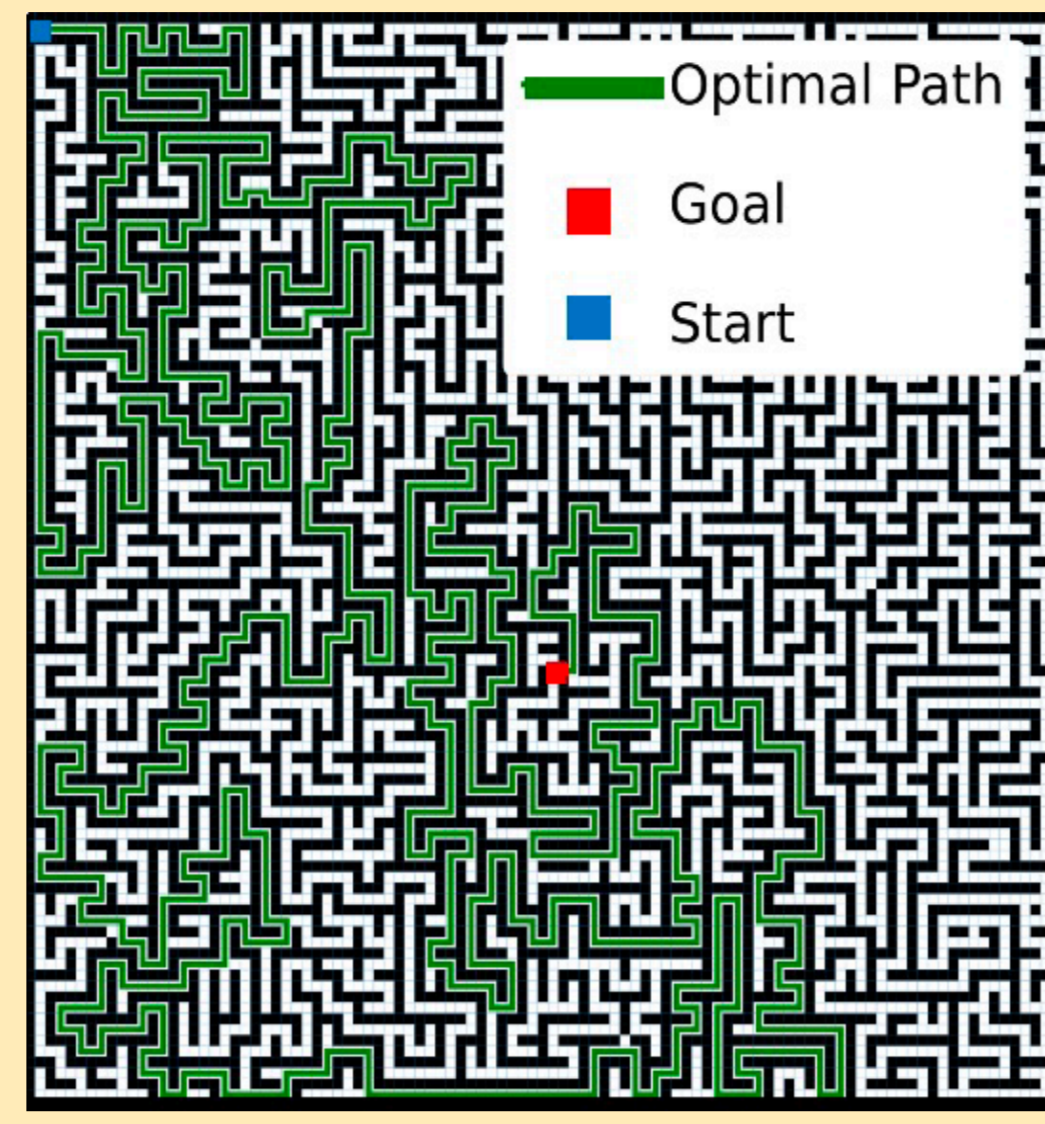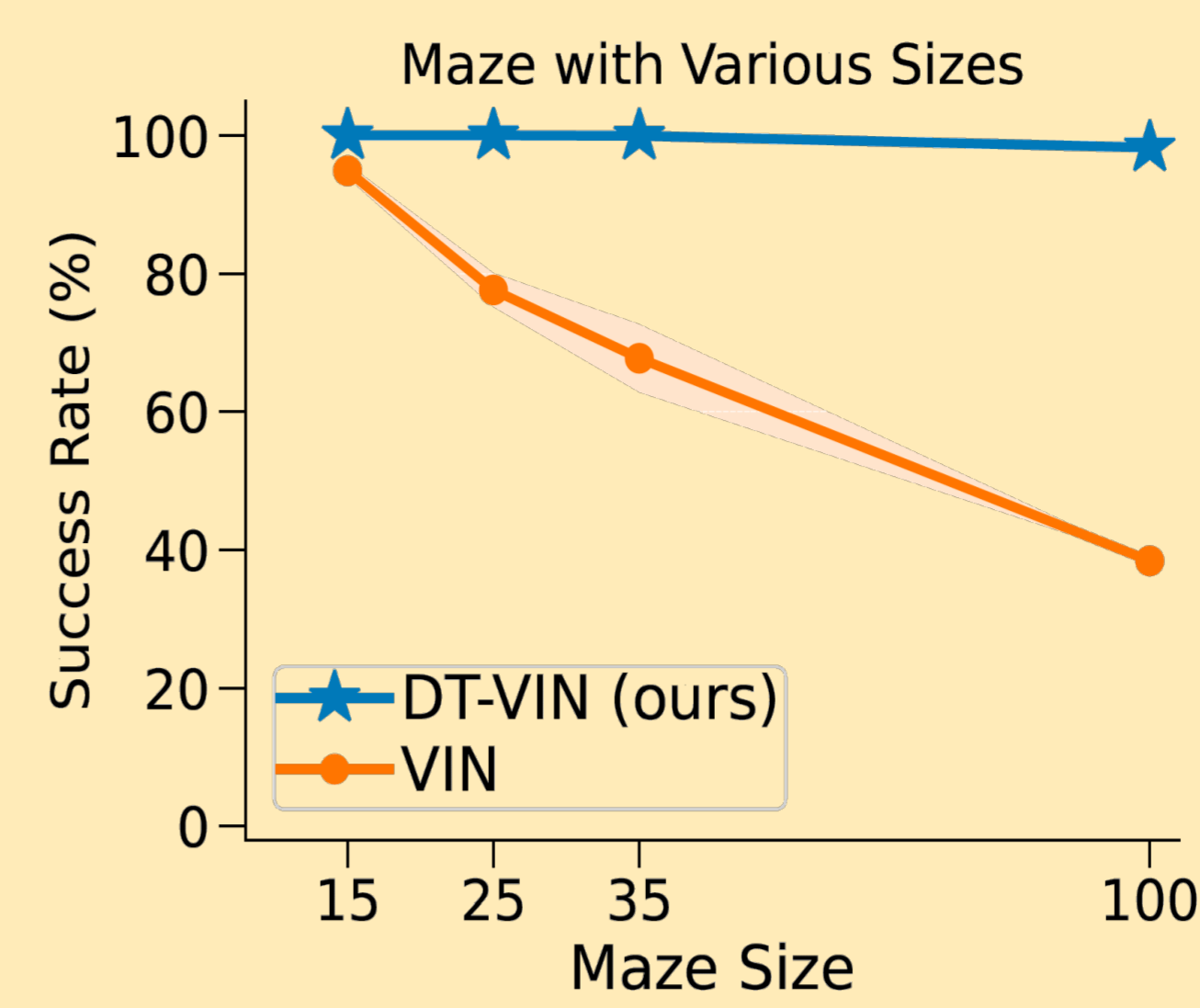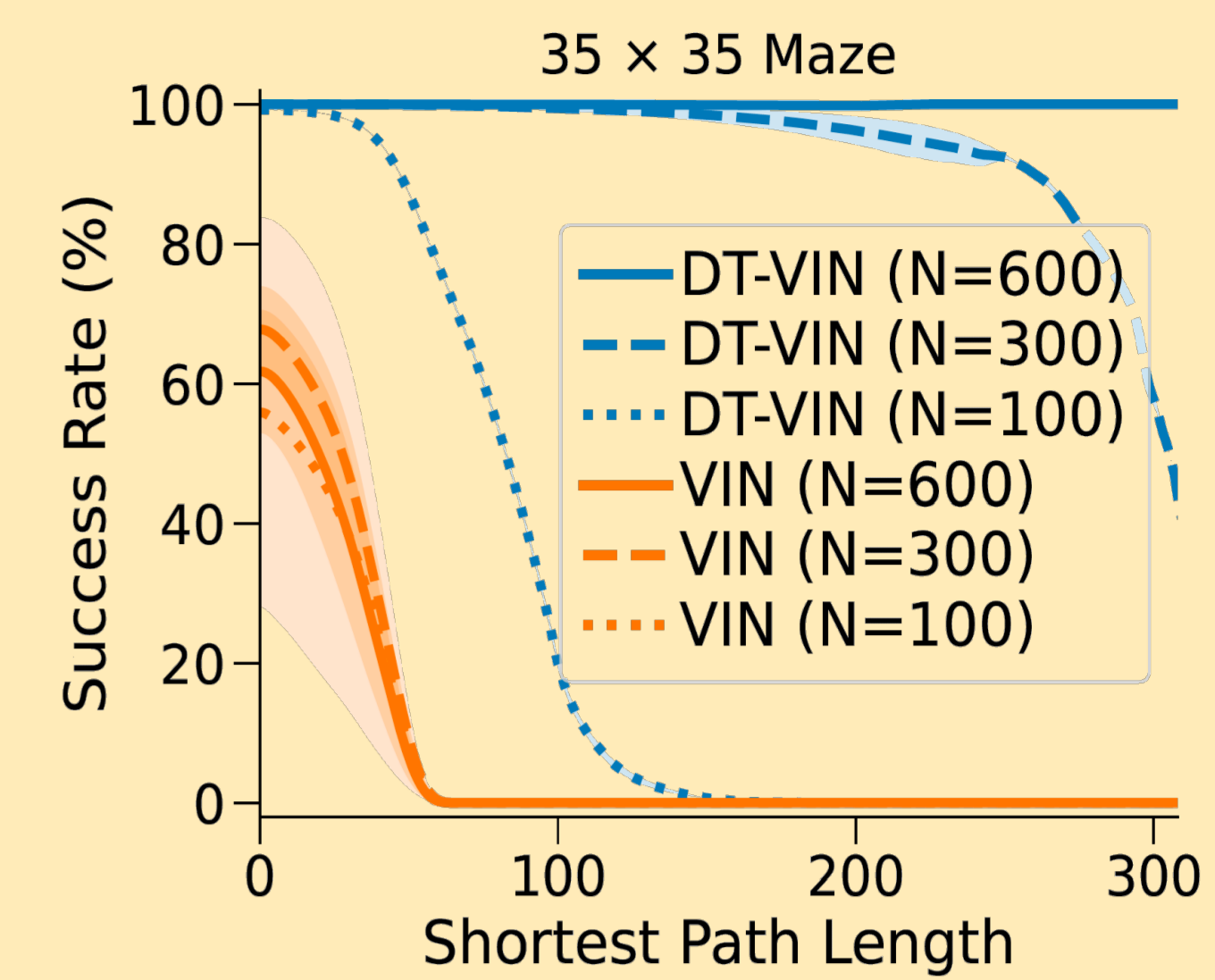
## Abstract

The Value Iteration Network (VIN) is an end-to-end differentiable architecture that performs value iteration on a latent MDP for planning in reinforcement learning (RL). However, VINs struggle to scale to long-term and large-scale planning tasks, such as navigating a 100x100 maze—a task which typically requires thousands of planning steps to solve. We observe that this deficiency is due to two issues: the representation capacity of the latent MDP and the planning module's depth. **We address these by augmenting the latent MDP with a dynamic transition kernel, dramatically improving its representational capacity, and, to mitigate the vanishing gradient problem, introduce an "adaptive highway loss" that constructs skip connections to improve gradient flow.** We evaluate our method on both 2D maze navigation environments and the ViZDoom 3D navigation benchmark. We find that our new method, named *Dynamic Transition VIN (DT-VIN)*, easily scales to 5000 layers and casually solves challenging versions of the above tasks. Altogether, we believe that DT-VIN represents a concrete step forward in performing long-term large-scale planning in RL environments.
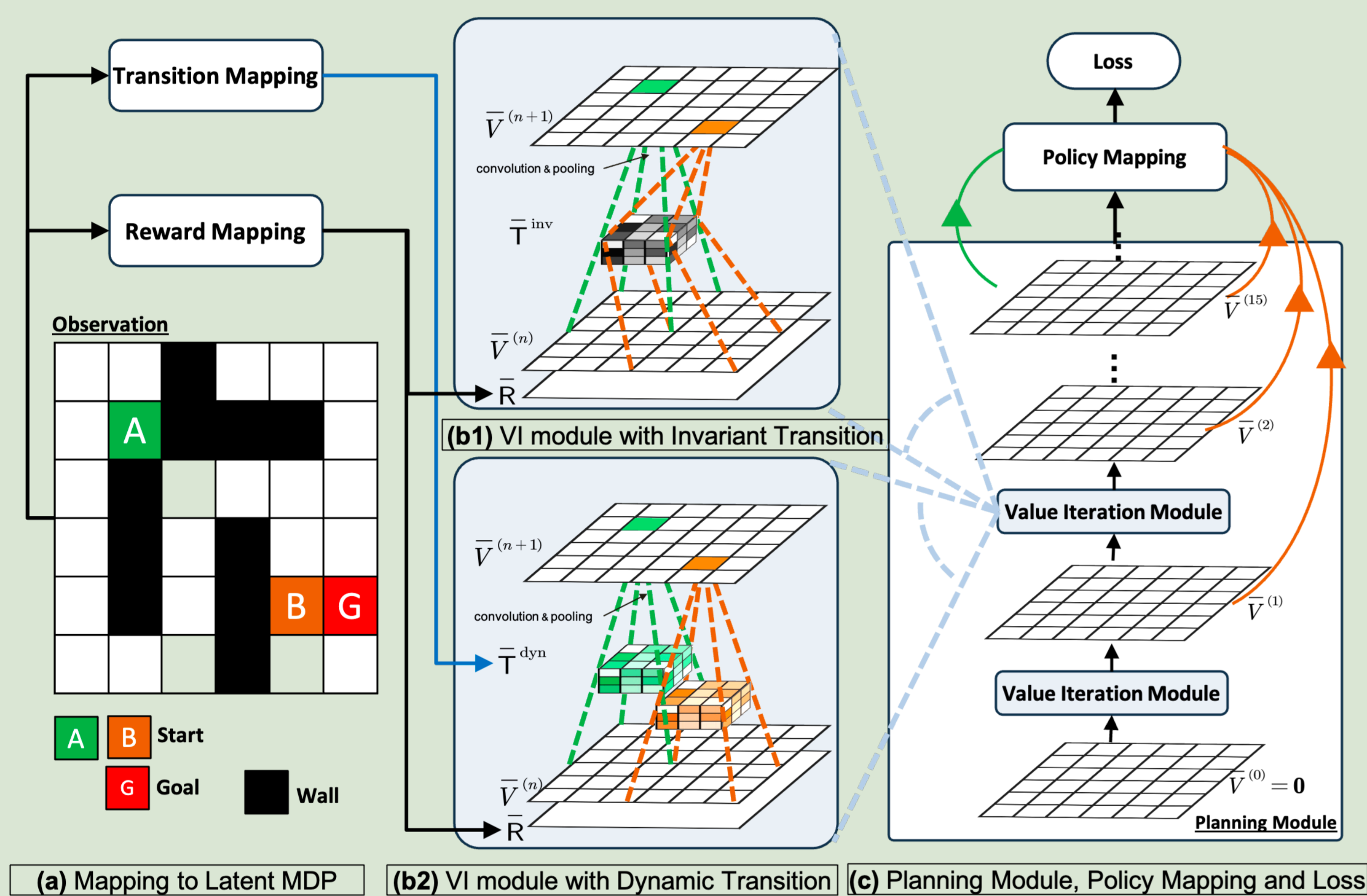
**100 × 100 Maze Navigation**

- Optimal Path
- Goal
- Start

**Performance for Different Domain Scales**

Maze with Various Sizes

- DT-VIN (ours)
- VIN

**Performance for Different Planning Steps**

35 × 35 Maze

- DT-VIN (N=600)
- DT-VIN (N=300)
- DT-VIN (N=100)
- VIN (N=600)
- VIN (N=300)
- VIN (N=100)

## Method



(a) Mapping to Latent MDP
(b1) VI module with Invariant Transition
(b2) VI module with Dynamic Transition
(c) Planning Module, Policy Mapping and Loss

A, B: Start | G: Goal | Wall

### 1) VI Module

**VIN: invariant transition kernel**

$$\overline{V}_{i,j}^{(n)} = \max_{\overline{a}} \sum_{i',j'} \overline{T}_{\overline{a},i',j'}^{\mathrm{inv}}\left(\overline{R}_{i-i',j-j'} + \overline{V}_{i-i',j-j'}^{(n-1)}\right)$$

$$\overline{T}^{\mathrm{inv}} \in \mathbb{R}^{|\overline{\mathcal{A}}| \times F \times F}$$

**Ours: dynamic transition kernel**

$$\overline{V}_{i,j}^{(n)} = \max_{\overline{a}} \sum_{i',j'} \overline{T}_{i,j,\overline{a},i',j'}^{\mathrm{dyn}}\left(\overline{R}_{i-i',j-j'} + \overline{V}_{i-i',j-j'}^{(n-1)}\right)$$

$$\overline{T}^{\mathrm{dyn}} = softmax\left(f^{\overline{T}}(x)\right) \in \mathbb{R}^{m \times m \times |\overline{\mathcal{A}}| \times F \times F}$$

where $f^{\overline{T}}$ is the transition mapping module

## 2) Loss function

**VIN: Normal Loss**

$$\mathcal{L}(\theta) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \ell\left(f^{\pi}\left(\overline{V}^{(N)}(x)\right)\right)$$

where

$f^{\pi}$: policy mapping module
$l$: loss function, e.g., cross entropy loss
$x$: observation
$y$: label

**Ours: Adaptive Highway Loss**

$$\mathcal{L}(\theta) = \frac{1}{K} \sum_{(x,y,l) \in \mathcal{D}} \sum_{1 \leq n \leq N} \mathbb{1}_{\{n \geq l\}} \mathbb{1}_{\{n \bmod l_j = 0\}} \ell\left(f^{\pi}\left(\overline{V}^{(n)}(x)\right), y\right)$$
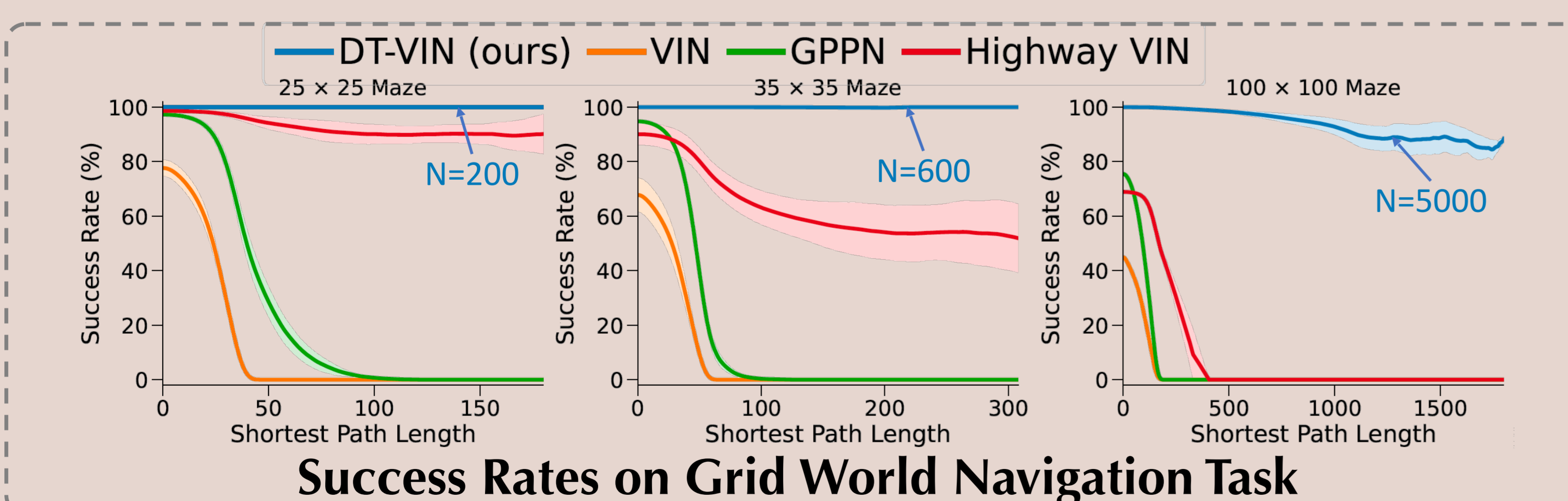
where

$n$: the index of the layer
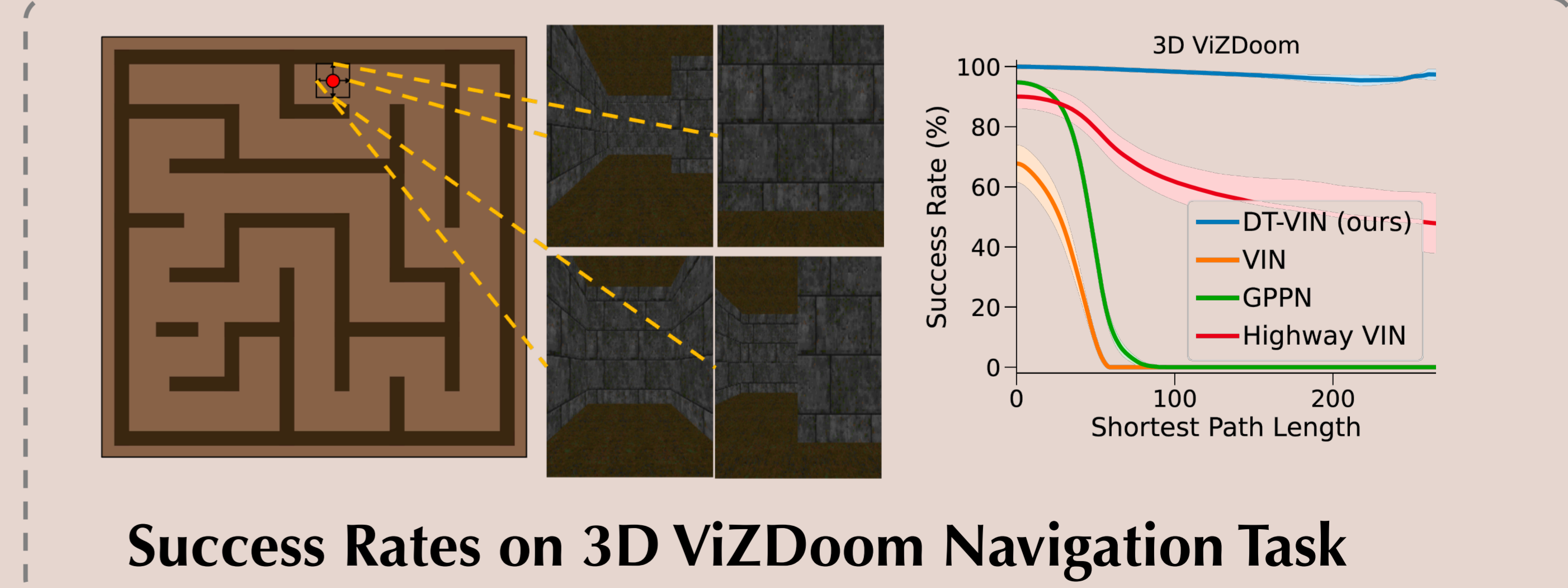$l$: the length of the trajectory
$l_j$: jumping hyperparameter
$\mathbb{1}$: indictor function

$$K = \sum_{(x,y,l) \in \mathcal{D}} \sum_{1 \leq n \leq N} \mathbb{1}_{\{n \geq l\}} \mathbb{1}_{\{n \bmod l_j = 0\}}$$
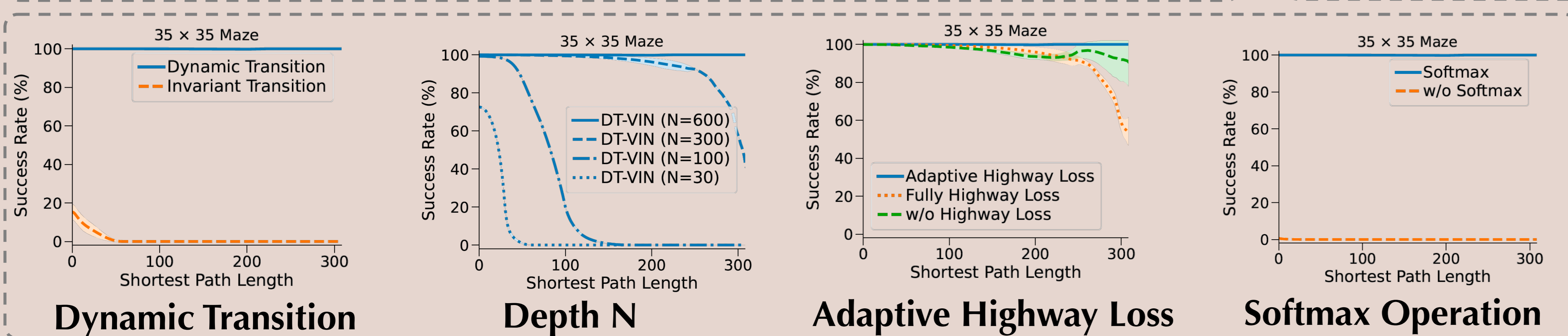
## Results



**Success Rates on Grid World Navigation Task**

- DT-VIN (ours)
- VIN
- GPPN
- Highway VIN

25 × 25 Maze — N=200
35 × 35 Maze — N=600
100 × 100 Maze — N=5000



**Success Rates on 3D ViZDoom Navigation Task**

3D ViZDoom
- DT-VIN (ours)
- VIN
- GPPN
- Highway VIN



**Dynamic Transition**
35 × 35 Maze
- Dynamic Transition
- Invariant Transition

**Depth N**
35 × 35 Maze
- DT-VIN (N=600)
- DT-VIN (N=300)
- DT-VIN (N=100)
- DT-VIN (N=30)

**Adaptive Highway Loss**
35 × 35 Maze
- Adaptive Highway Loss
- Fully Highway Loss
- w/o Highway Loss

**Softmax Operation**
35 × 35 Maze
- Softmax
- w/o Softmax

Correspondence to:
Yuhui Wang, yuhui.wang@kaust.edu.sa,
Dylan R. Ashley, mail@dylanashley.io