# Highway Value Iteration Networks

Yuhui Wang*[1], Weida Li*[2], Francesco Faccio[1,3], Qingyuan Wu[4], Jürgen Schmidhuber[1,3]

[1]AI Initiative, King Abudullah University of Science and Technology (KAUST), [2]National University of Singapore, [3]The Swiss AI Lab IDSIA/USI/SUPSI, [4]The University of Liverpool

## Motivation

**Maze Navigation**

**Performance on tasks with various shortest path lengths**

## Background

**VIN/Highway VIN**

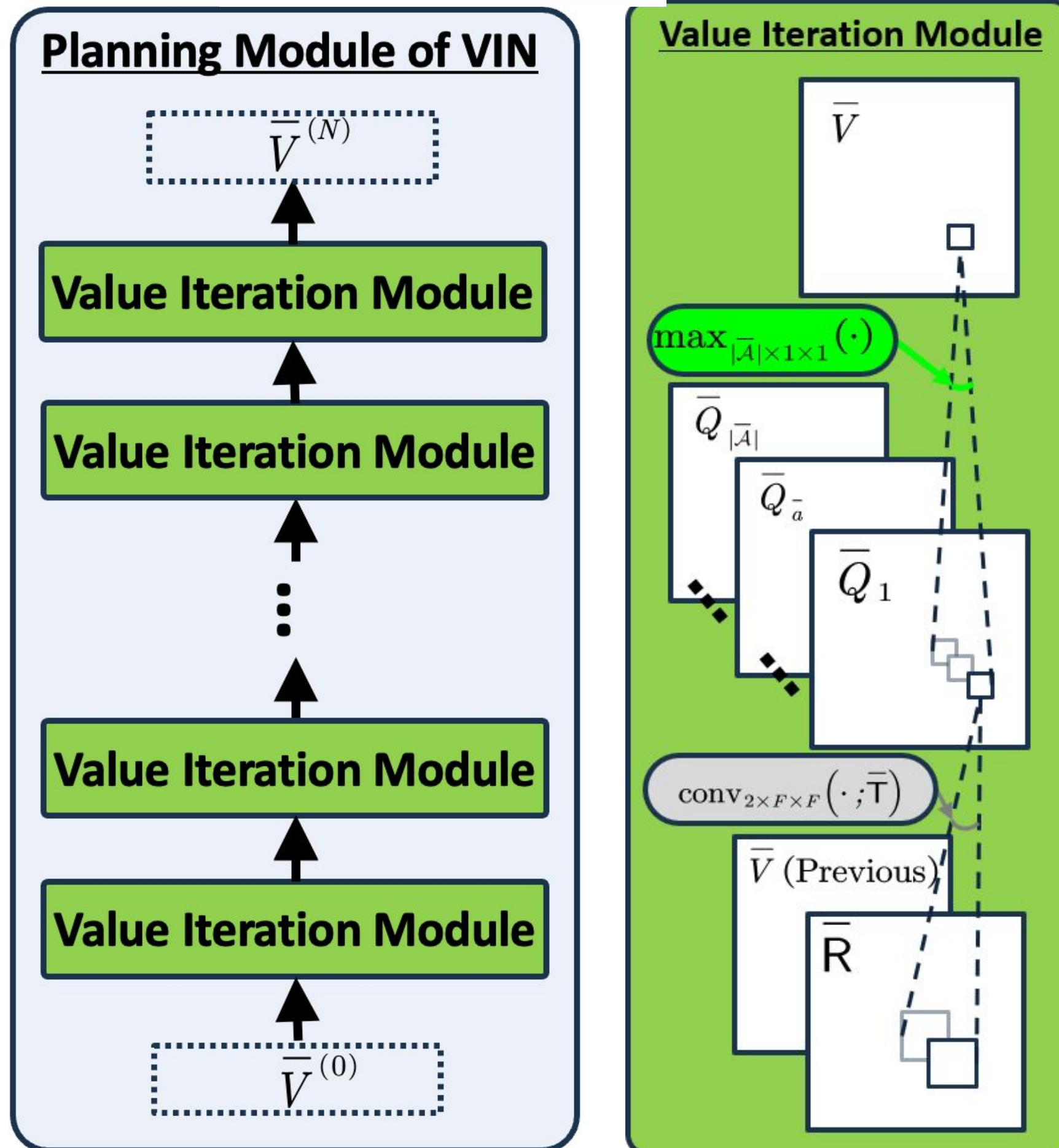**Architecture of Value Iteration Network (VIN) [1]**

## Method

**Bellman Optimality/Expectation Operator**

$$(\mathcal{B}V)(s) \triangleq \max_a \sum_{s'} \mathcal{T}(s'|s,a)[\mathcal{R}(s,a,s') + \gamma V(s')], \quad (\mathcal{B}^{\pi}V)(s) \triangleq \sum_a \pi(a|s) \sum_{s'} \mathcal{T}(s'|s,a)[\mathcal{R}(s,a,s') + \gamma V(s')]$$

**Planning Module of VIN**

$$V^{(n+1)} = \mathcal{B}V^{(n)}$$

**Planning Module of VIN**

**Value Iteration Module**

**Value Iteration Module:**

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

(2) $\overline{V}_{i,j}^{(n)} = \max_{\overline{a}} \overline{Q}_{\overline{a},i,j}^{(n)}$

**Planning Module of Highway VIN**

$$V^{(n+1)} = \mathop{smax}_{\pi \in \Pi}^{\widetilde{\alpha}} \mathop{smax}_{n \in \mathcal{N}}^{\alpha} \max \left\{ (\mathcal{B}^{\pi})^{\circ(N_b - 1)} \mathcal{B}V^{(n)}, \mathcal{B}V^{(n)} \right\}$$

*smax* is the softmax function

**Planning Module of Highway VIN**

**Highway Block** — Aggregate Gate

**Filter and Aggregate Gate** — Value Exploration Module

**Value Iteration Module**

**Value Exploration Module**

**Value Exploration Module:**

(1) $\overline{\mathcal{Q}}_{\pi,\overline{a},i,j}^{(n+n_b)} = \sum_{i',j'} \left( \overline{\mathsf{T}}_{\overline{a},i',j'} \overline{\mathsf{R}}_{i-i',j-j'} + \overline{\mathsf{T}}_{\overline{a},i',j'} \overline{\mathcal{V}}_{\pi,i-i',j-j'}^{(n+n_b-1)} \right)$

(2) $\overline{\mathcal{V}}_{n_p,i,j}^{(n+n_b)} = \sum_{\overline{a}} \overline{\pi}_{n_p,\overline{a},i,j}^{(n+n_b)} \overline{\mathcal{Q}}_{n_p,\overline{a},i,j}^{(n+n_b)}$

where

$$\overline{\pi}_{n_p,\overline{a},i,j}^{(n+n_b)} = \begin{cases} 1, & \overline{a} = \widehat{\overline{a}} \sim P\left(\cdot; \overline{\mathcal{Q}}_{n_p,\cdot,i,j}^{(n+n_b)}, \epsilon\right) \\ 0, & \text{otherwise}, \end{cases}$$

$$P\left(\overline{a}; \overline{\mathcal{Q}}_{n_p,\cdot,i,j}^{(n+n_b)}, \epsilon\right) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}|}, & \overline{a} = \arg\max_{\overline{a}'} \overline{\mathcal{Q}}_{n_p,\overline{a}',i,j}^{(n+n_b)} \\ \frac{\epsilon}{|\mathcal{A}|}, & \text{otherwise}. \end{cases}$$

**Aggregate Gate:**     **Filter Gate:**

$$\overline{V}_{i,j}^{(n+N_b)} = \sum_{n_p=1}^{N_p} \widetilde{\mathsf{A}}_{n_p,i,j}^{(n+N_b)} \underbrace{\sum_{n_b=1}^{N_b} \mathsf{A}_{n_p,i,j}^{(n+n_b)} \max\left\{ \overline{\mathcal{V}}_{n_p,i,j}^{(n+n_b)}, \overline{V}_{i,j}^{(n+1)} \right\}}_{\overline{V}''^{(n+N_b)}_{n_p,i,j}}$$

where

$$\widetilde{\mathsf{A}}_{n_p,i,j}^{(n+N_b)} = \frac{\exp\left(\alpha_{\widetilde{\mathsf{A}}} \overline{V}''^{(n+N_b)}_{n_p,i,j}\right)}{\sum_{n_p'} \exp\left(\alpha_{\widetilde{\mathsf{A}}} \overline{V}''^{(n+N_b)}_{n_p',i,j}\right)} \quad \mathsf{A}_{n_p,i,j}^{(n+n_b)} = \frac{\exp\left(\alpha_{\mathsf{A}} \overline{V}'^{(n+N_b)}_{n_p,i,j}\right)}{\sum_{n_b'} \exp\left(\alpha_{\mathsf{A}} \overline{V}'^{(n+N_b)}_{n_p',i,j}\right)}$$

## Experiments

**15 × 15 Maze**

- Highway VIN (ours)
- VIN
- Highway Net
- GPPN

**25 × 25 Maze**

## Ablation Studies

**15 × 15 Maze**
- w/ filter gate
- w/o filter gate

**Filter Gate**

**15 × 15 Maze**
- w/ VE modules
- w/o VE modules

**Value Exploration Module**

**25 × 25 Maze**
- N=60
- N=90
- N=120
- N=150
- N=180
- N=240
- N=300

**Depths of Highway VIN (ours)**

**25 × 25 Maze**
- N=30
- N=60
- N=90
- N=120
- N=150
- N=180
- N=240
- N=300

**Depths of VIN**

[1] Tamar A, Wu Y, Thomas G, Levine S, & Abbeel P. Value iteration networks[J]. Advances in neural information processing systems, 2016, 29.

[2] Wang Y, Strupl M, Faccio F, Wu Q, Liu H, Grudzien M, Tan X, and Schmidhuber J. Highway reinforcement learning[J]. arXiv preprint arXiv:2405.18289, 2024.