# LEARNING TO SEARCH FROM DEMONSTRATION SEQUENCES

Dixant Mittal, Liwei Kang, Wee Sun Lee
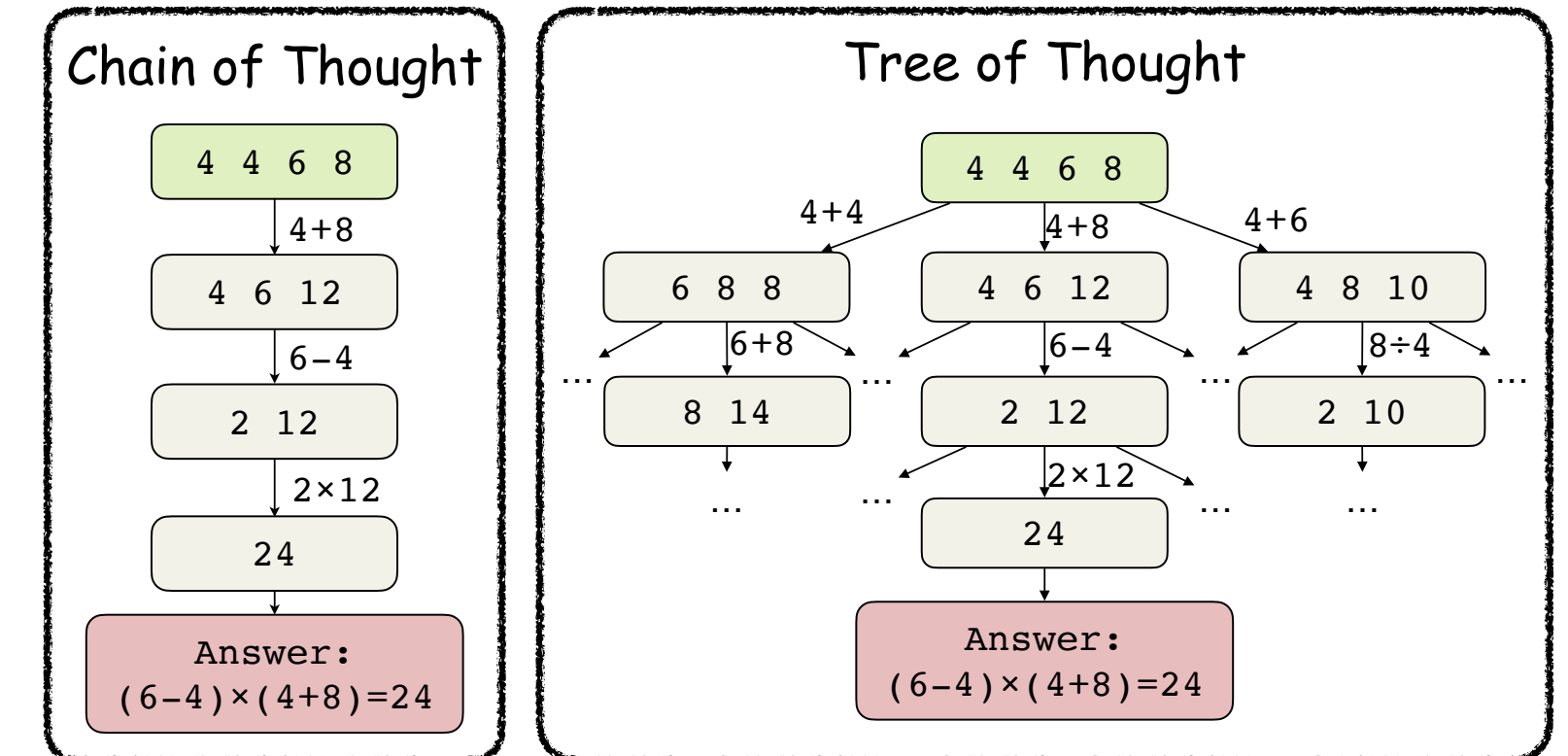
**Search and Planning**

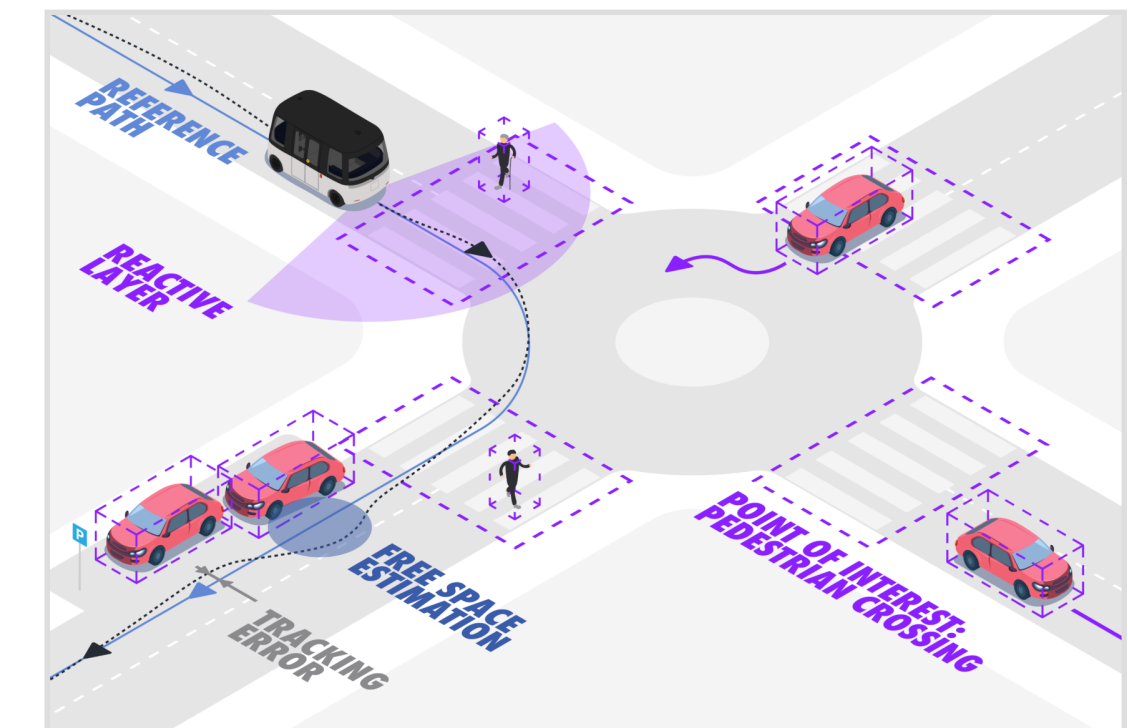are fundamental for complex reasoning tasks

**Data is Often Limited**

in real-world and all we may have is a collection of

**Demonstration Sequences**

Learning to Search from Demonstrations is tricky because:

- Limited **state space coverage**

- **No exploration**

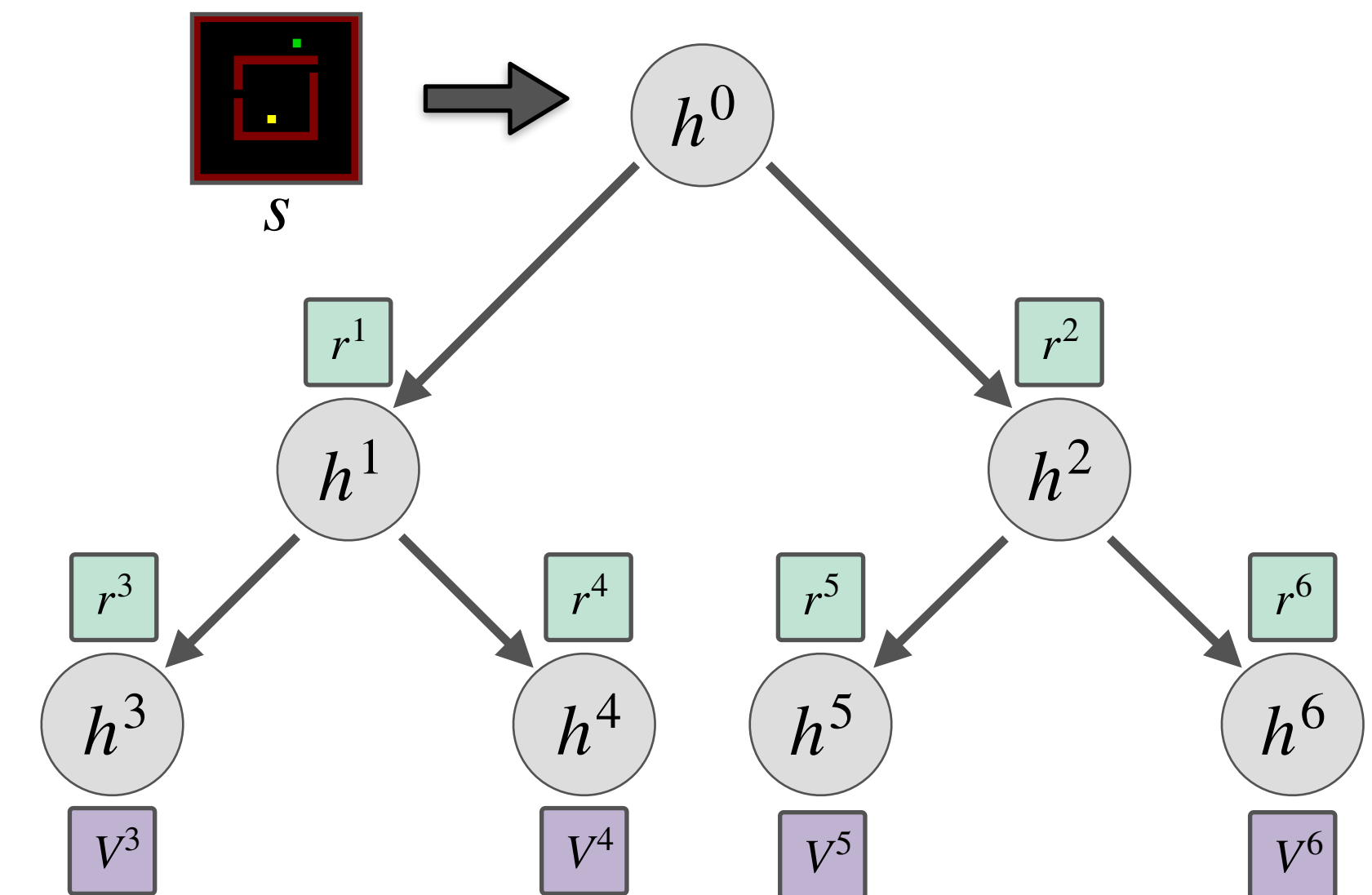- **Compounding Errors** for planning



**Reasoning in LLMs**



**Real World
Autonomous Driving**

**Search Tree as Parametric Policy**

and directly learn the mapping from state to action

**TreeQN**

expands the full search tree upto a fixed depth
and backups the value to the root node



**TreeQN**

# TREEQN: **STRENGTHS & LIMITATIONS**

## STRENGTHS

**Planning** inductive bias

**No dependency** on Simulator

## LIMITATIONS

TreeQN is **exponential** in depth

**Infeasible** to perform **deeper** search

Performance suffers in **complex** problems

**Modular Neural Network Architecture**

that comprises of several *learnable submodules*

**Algorithmic Inductive Bias**

of a flexible and scalable *best-first tree search* algorithm
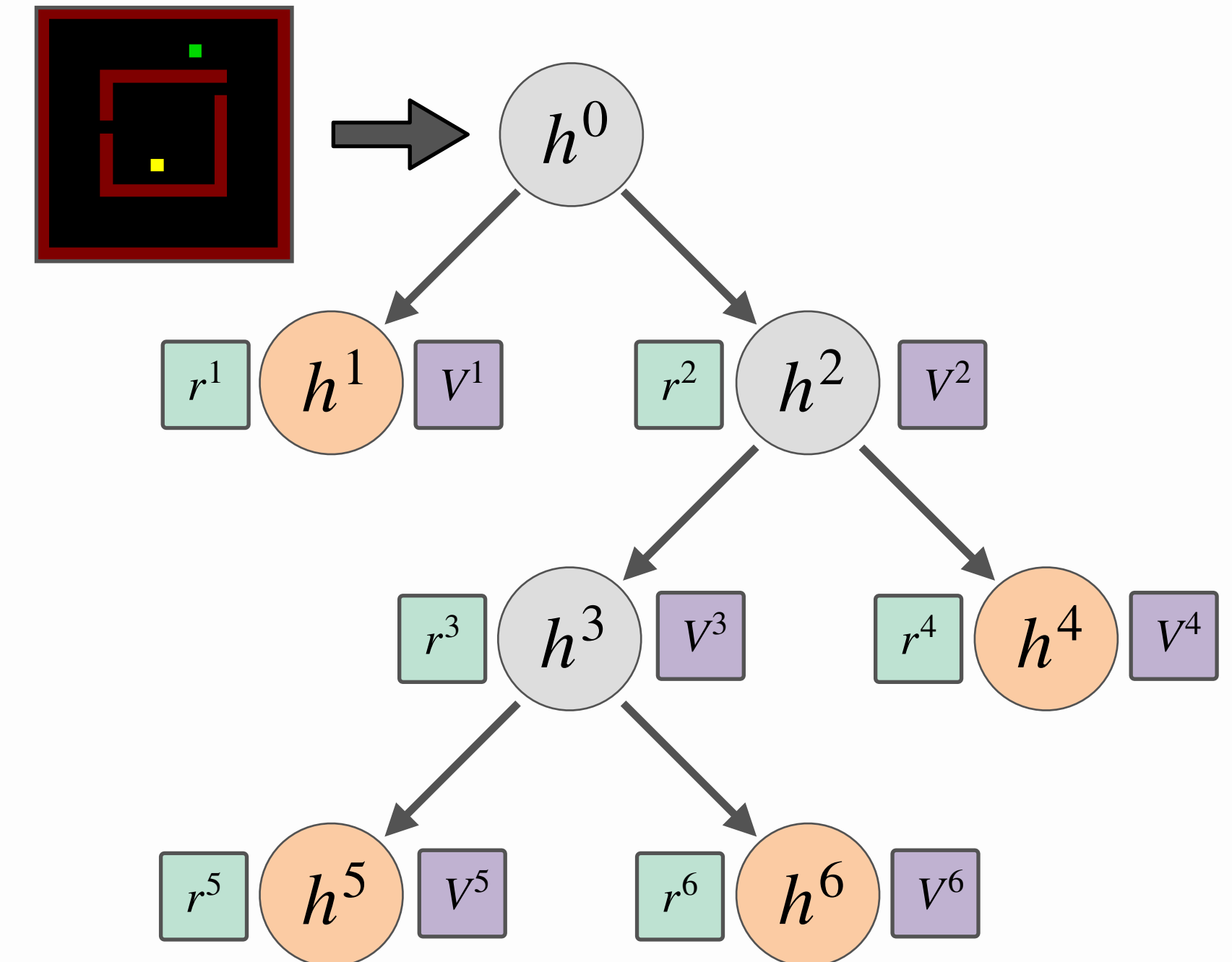
**Learnt World Model**

that is trained to be *useful* for the online search, even if inaccurate
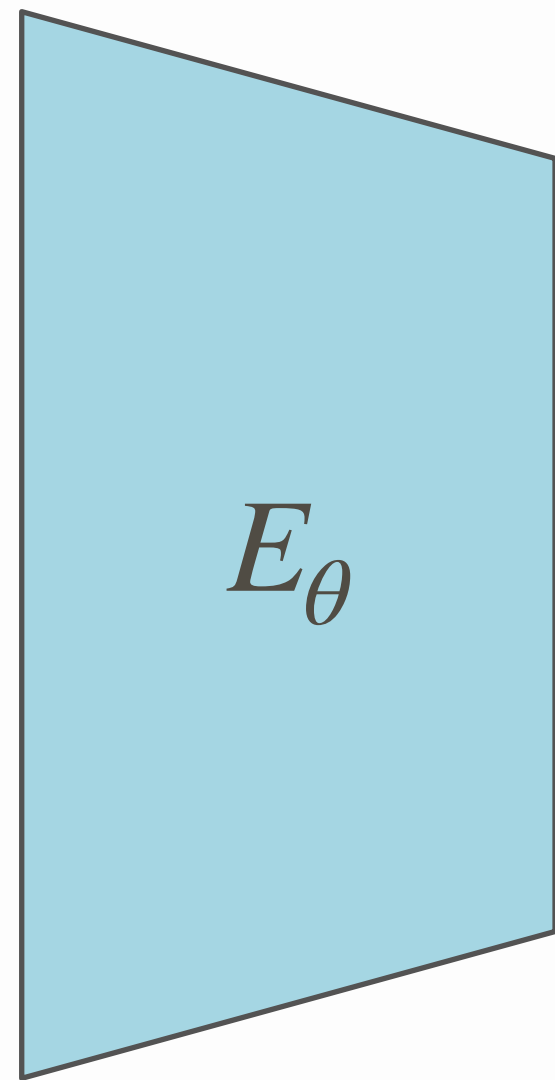
**Joint Optimisation**

Trains the *search* and *world model* submodules jointly
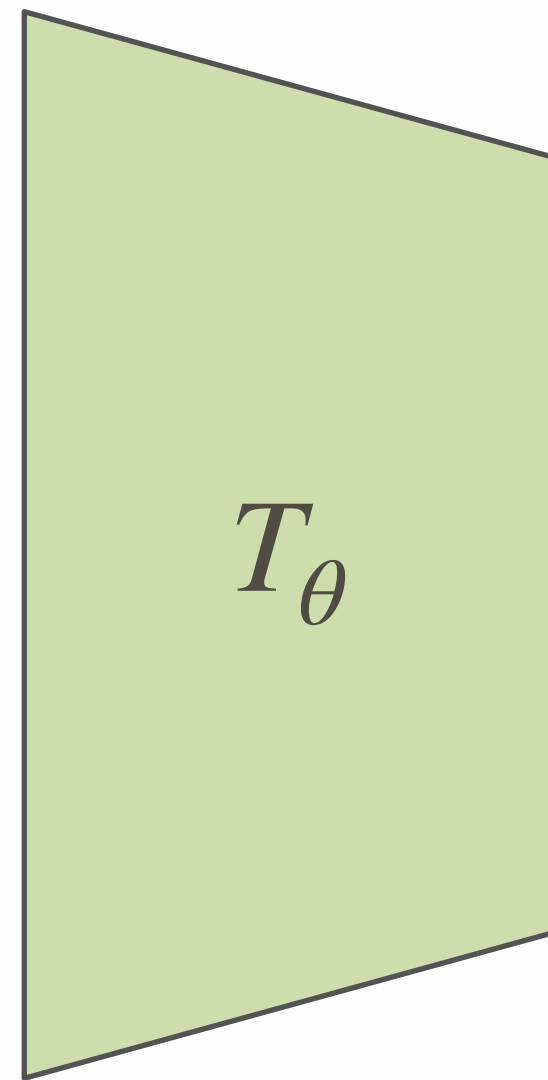
**Additional Technical Details**

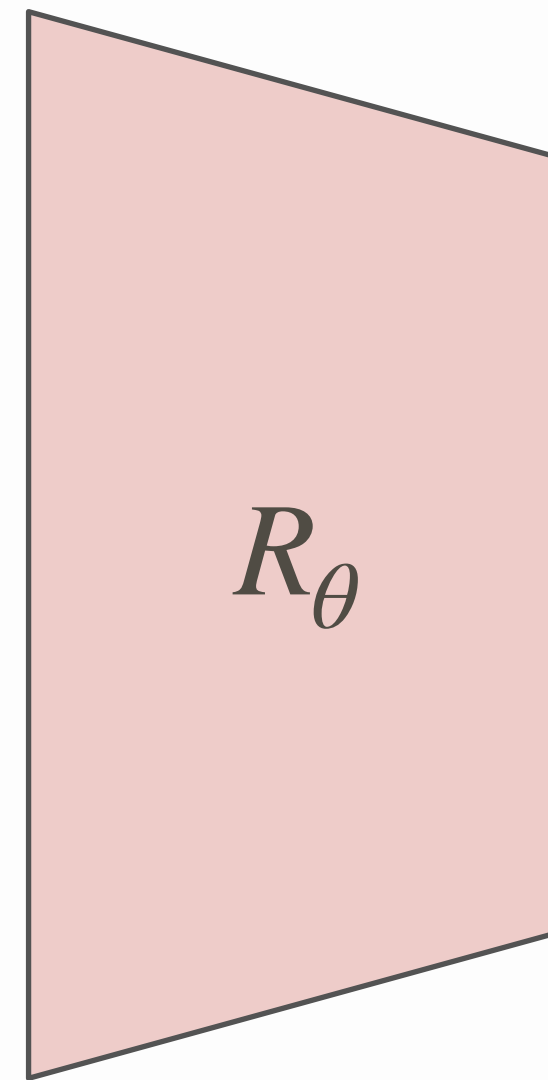like *Tree Expansion Policy* and *Telescopic Sum for Variance Reduction*
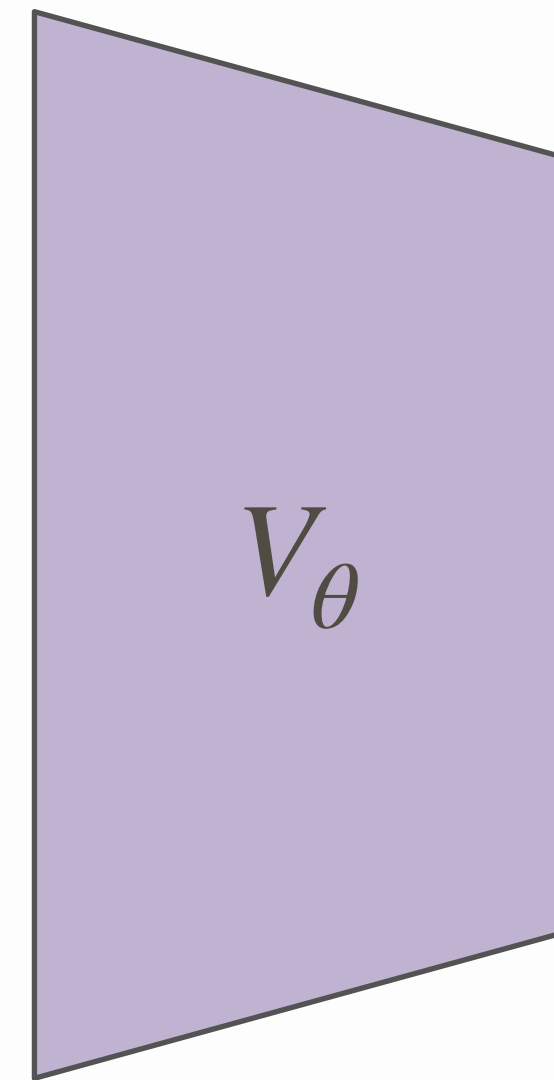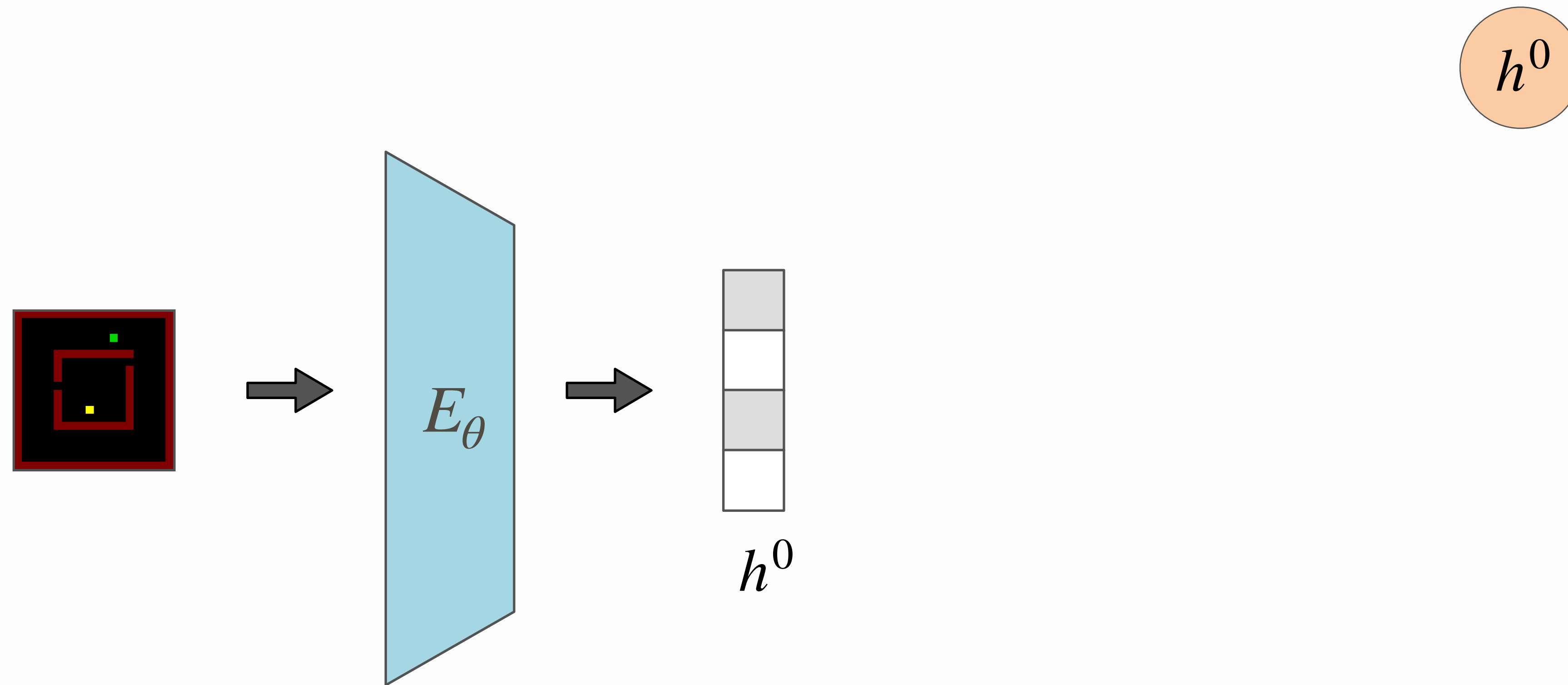
$E_\theta$

$T_\theta$

$R_\theta$

$V_\theta$

**ENCODER MODULE**

**TRANSITION MODULE**

**REWARD MODULE**

**VALUE MODULE**

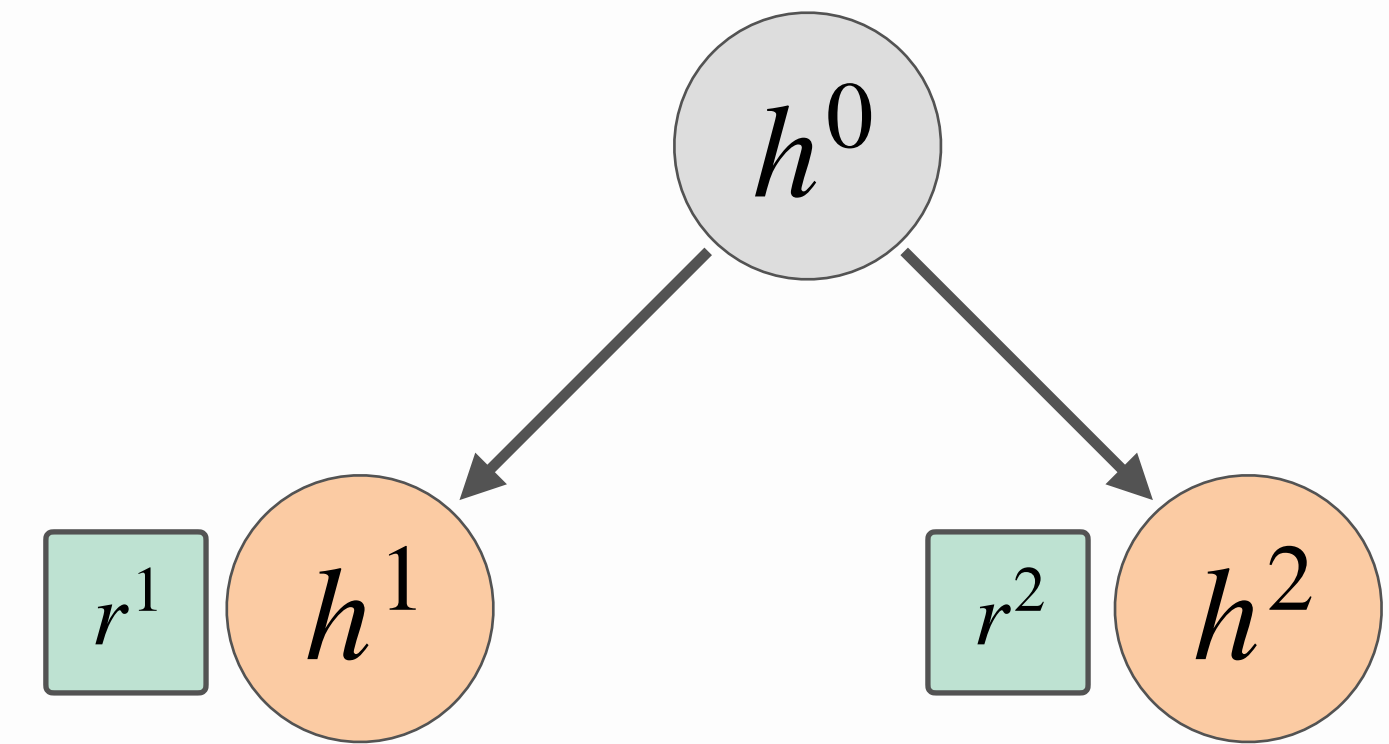$$E_\theta$$

$$h^0$$

$$h^0$$

**Stochastic Tree Expansion Policy**

$$n \sim \pi(\tau)$$



$$\pi(\tau_1) = \mathbf{softmax}\Big( h^1 \quad h^2 \Big)$$

$$= \mathbf{softmax}\Big( \begin{array}{c} r^1 + V^1 \\ r^2 + V^2 \end{array} \Big)$$

**Stochastic Tree Expansion Policy**

$$n \sim \pi(\tau)$$

$$\pi(\tau_2) = \textbf{softmax}\left( h^1 \quad h^3 \quad h^4 \right)$$

$$= \textbf{softmax}\left( \begin{array}{c} r^1 + V^1 \\ r^2 + r^3 + V^3 \\ r^2 + r^4 + V^4 \end{array} \right)$$



$$\tau_2$$

NUS
National University
of Singapore

**BELLMAN EQUATION**

$$Q(h, a) = r(h, a) + V(h')$$

$$V(h) = \max_{a} \left[ Q(h, a) \right]$$

Model's Output: $Q_\theta(s, a)$

Path returning the highest value

$$Q_\theta\left( \; h^0 \; , \; \nearrow \; \right) = \boxed{r^1} + \boxed{V^1}$$

$$Q_\theta\left( \; h^0 \; , \; \searrow \; \right) = \boxed{r^2} + \boxed{r^3} + \boxed{r^6} + \boxed{V^6}$$

TreeQN

D-TSN

**Exponential** in Depth

Can go **Deeper**

**TreeQN**

$$Q(s, Left) = \boxed{r^1} + max \begin{bmatrix} \boxed{r^3} + \boxed{V^3} \\ \boxed{r^4} + \boxed{V^4} \end{bmatrix}$$

$$Q(s, Right) = \boxed{r^2} + max \begin{bmatrix} \boxed{r^5} + \boxed{V^5} \\ \boxed{r^6} + \boxed{V^6} \end{bmatrix}$$



$Q_\theta(s, a)$ from a **Full Tree Search** are **continuous** in parameter space

$L\left(Q_\theta(s, a)\right)$ is **continuous** in parameter space $\theta$

**D-TSN**

A **best-first search** solution is an **approximation** to the **Full Tree Search** solution

$Q_\theta(s, a)$ depends upon the sampled tree $\tau$

$L\Big(Q_\theta(s, a \mid \tau)\Big)$ can be **discontinuous** in parameter space $\theta$

$$\textbf{Loss} = \mathbb{E}_\tau \left[ L\Big( Q_\theta(s,a \,|\, \tau) \Big) \right] = \boxed{\sum_\tau \pi_\theta(\tau) \; L\Big( Q_\theta(s,a \,|\, \tau) \Big)}$$

**Continuous**

in the parameter space

$$\textbf{Loss} = \mathbb{E}_\tau \left[ L\Big( Q_\theta(s, a \mid \tau) \Big) \right]$$

$$\nabla_\theta(\textbf{Loss}) = \mathbb{E}_\tau \left[ \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(n_t \mid \tau_t) \, L\Big( Q_\theta(s, a \mid \tau_T) \Big) + \nabla_\theta L\Big( Q_\theta(s, a \mid \tau_T) \Big) \right]$$

**Policy Gradient**

(to improve tree expansion policy)

**Regular Loss Gradient**

(on output Q-values)

$$\nabla_\theta(\mathbf{Loss}) = \mathbb{E}_\tau \left[ \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(n_t \,|\, \tau_t) \, L\Big( Q_\theta(s, a \,|\, \tau_T) \Big) \;+\; \nabla_\theta L\Big( Q_\theta(s, a \,|\, \tau_T) \Big) \right]$$

**High Variance!**

$$\nabla_\theta(\textbf{Loss}) = \mathbb{E}_\tau \left[ \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(n_t \,|\, \tau_t) \boxed{L\Big(Q_\theta(s, a \,|\, \tau_T)\Big)} + \nabla_\theta L\Big(Q_\theta(s, a \,|\, \tau_T)\Big) \right]$$

**Let** $L_T = L\Big(Q_\theta(s, a \,|\, \tau_T)\Big)$ **and** $L_0 = 0$

$$L_T = L_T - L_0$$

$$= L_T - L_{T-1} + L_{T-1} - L_0$$

$$= \sum_{t=1}^{T} L_t - L_{t-1}$$

Telescopic Sum

Treat **Tree Expansion** as another decision making problem with goal of reducing Loss value after $t^{th}$ expansion.

## Define

– Reward as $$r_t = L_t - L_{t-1}$$

– Sum of Rewards as $$R_t = \sum_t^T r_t = L_T - L_{t-1}$$

$$\nabla_\theta(\textbf{Loss}) = \mathbb{E}_\tau \left[ \sum_{t=1}^{T} \left[ \nabla_\theta \log \pi_\theta(n_t \,|\, \tau_t) \right] L_T \quad + \quad \nabla_\theta L_T \right]$$

$$= \mathbb{E}_\tau \left[ \sum_{t=1}^{T} \left[ \nabla_\theta \log \pi_\theta(n_t \,|\, \tau_t) \right] \left[ L_T - L_{t-1} \right] \quad + \quad \nabla_\theta L_T \right]$$

**Lower** variance!

NUS
National University
of Singapore

**Differentiable Tree Search Network**

Differentiable Search + Joint Optimisation

**vs**

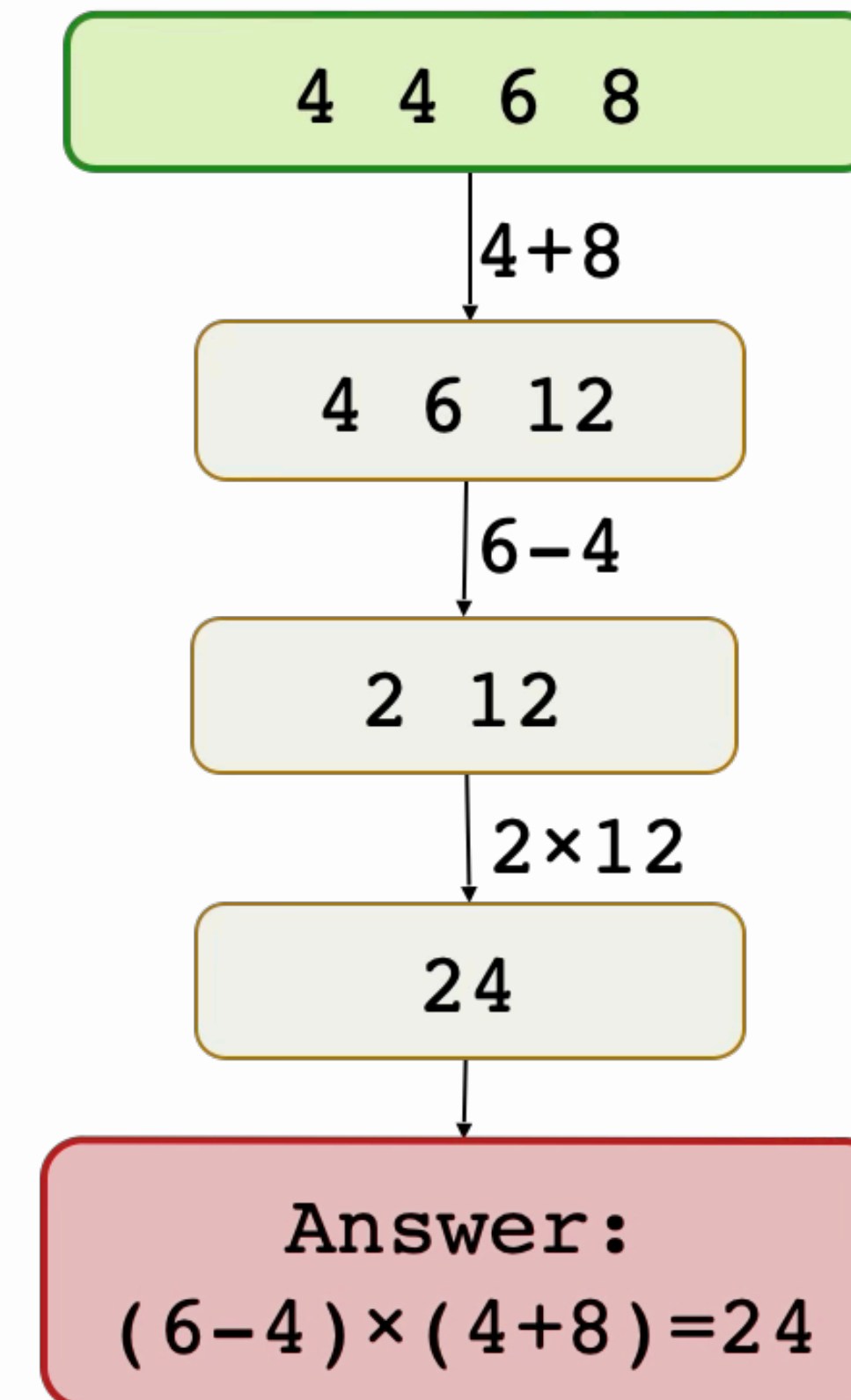**Model-free Q-network**

A generic Neural Network based Q-function

**Model-based Search**

Best-first Tree Search + Learnt World Model
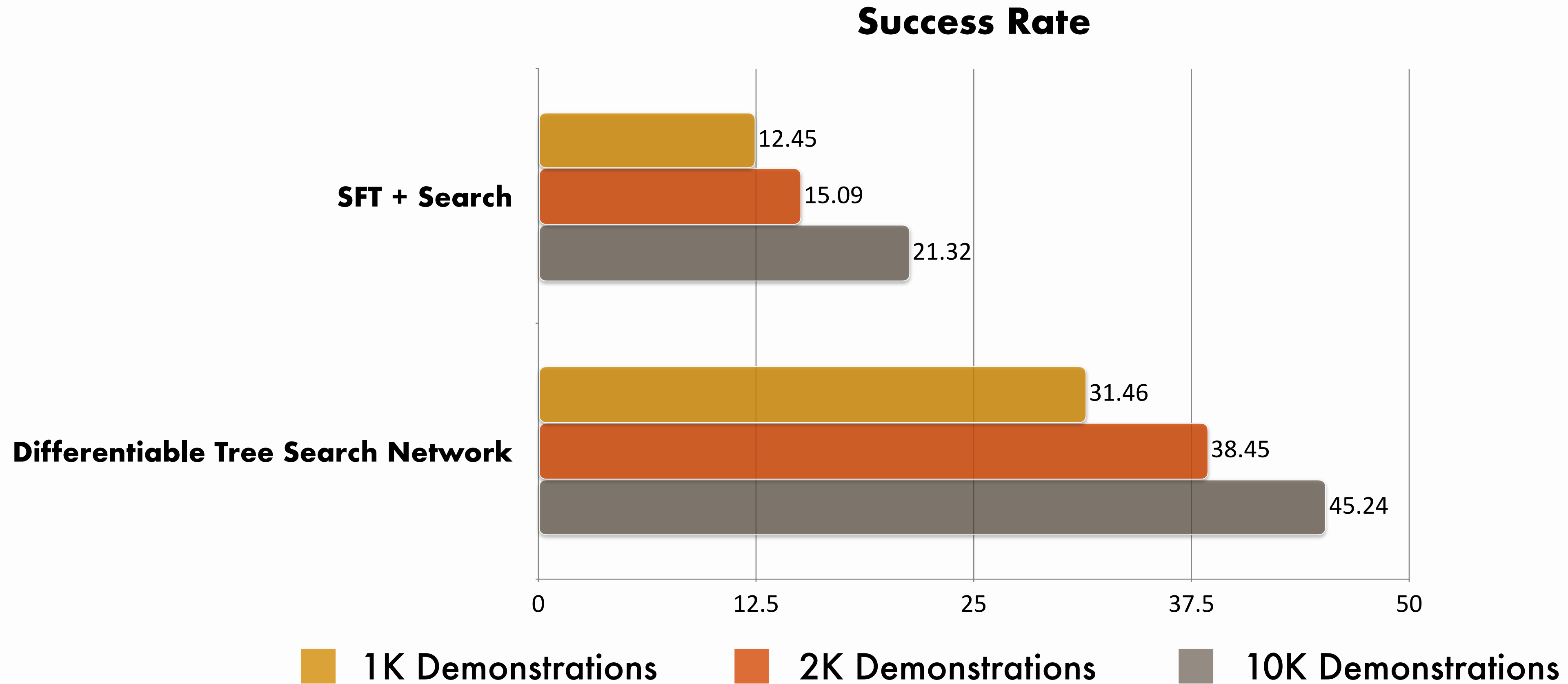(No joint optimisation)

**TreeQN**

Shallow Tree Search + Joint Optimisation

NUS
National University
of Singapore

- **Known** world model

- Actions are: Select **operands** or **operators**

- **527** problems for training

**Success Rate**



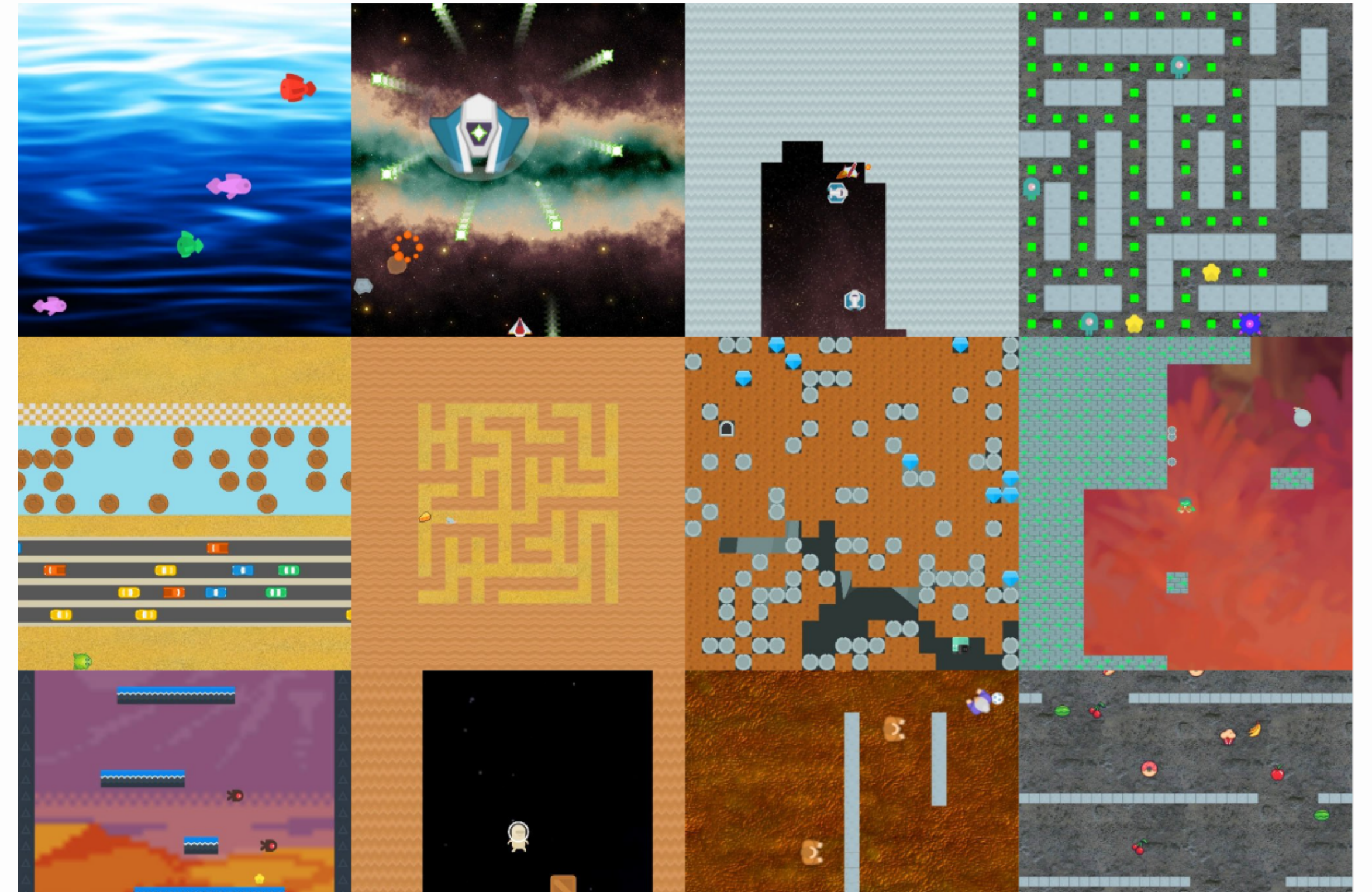| | Success Rate |
|---|---|
| **SFT + Search** | 12.45 (1K), 15.09 (2K), 21.32 (10K) |
| **Differentiable Tree Search Network** | 31.46 (1K), 38.45 (2K), 45.24 (10K) |

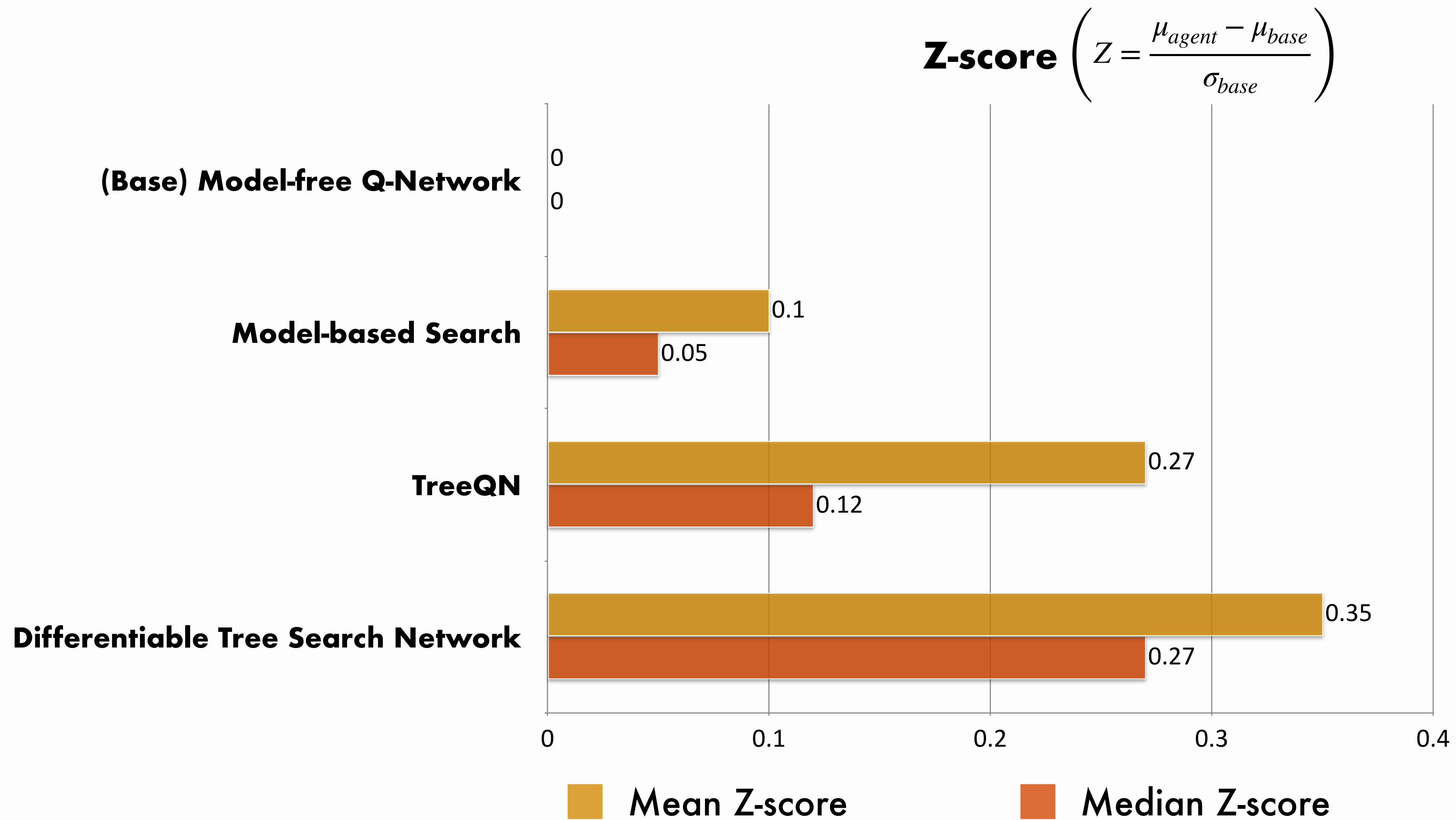■ 1K Demonstrations ■ 2K Demonstrations ■ 10K Demonstrations

# DOMAIN: **PROCGEN**

- Procedurally generated environments

- **16** different games

- Designed to test **generalisation capability** of an agent's policy

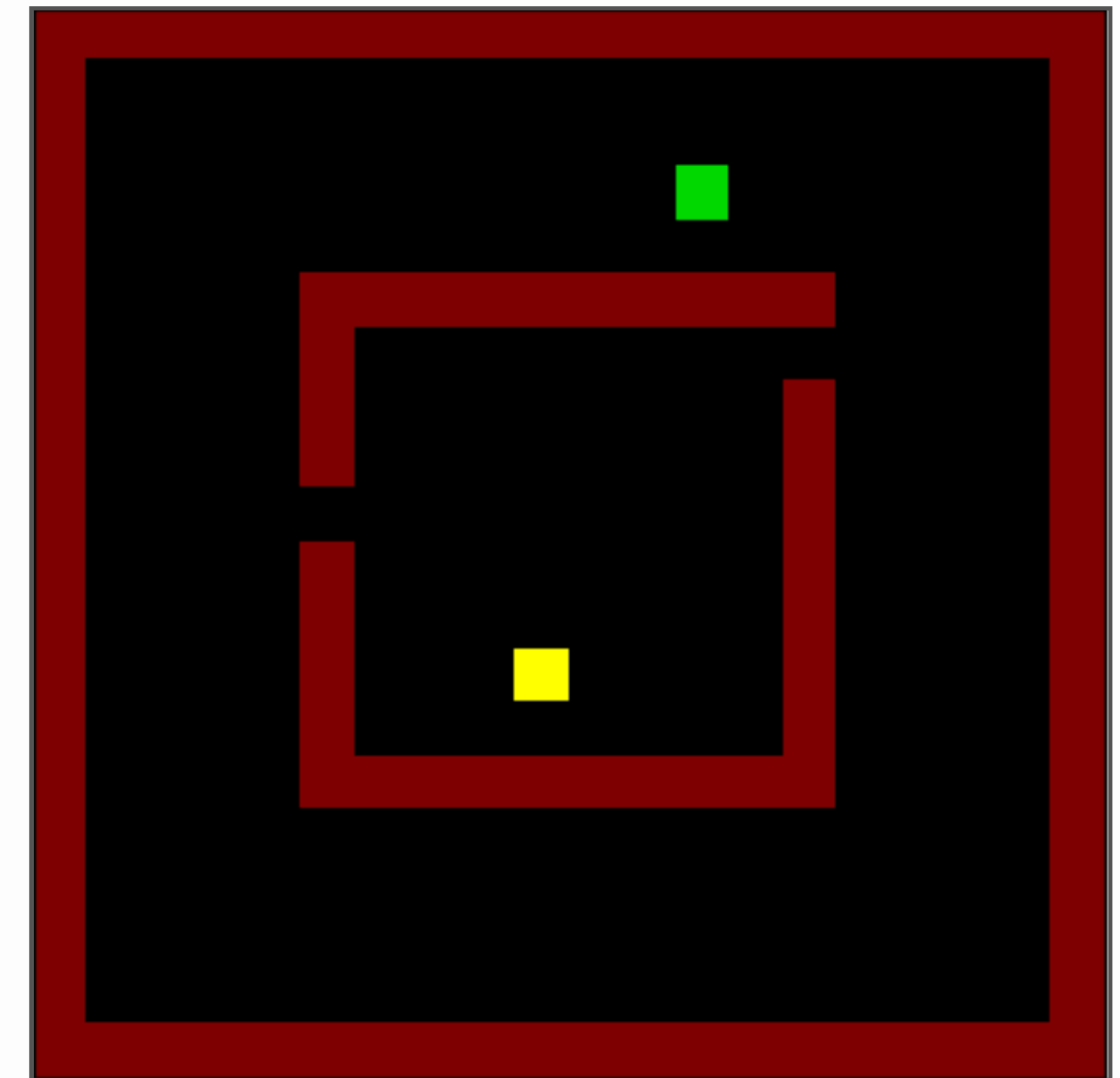- Collected **1000 Trajectories** for training

$$\text{Z-score} \left( Z = \frac{\mu_{agent} - \mu_{base}}{\sigma_{base}} \right)$$

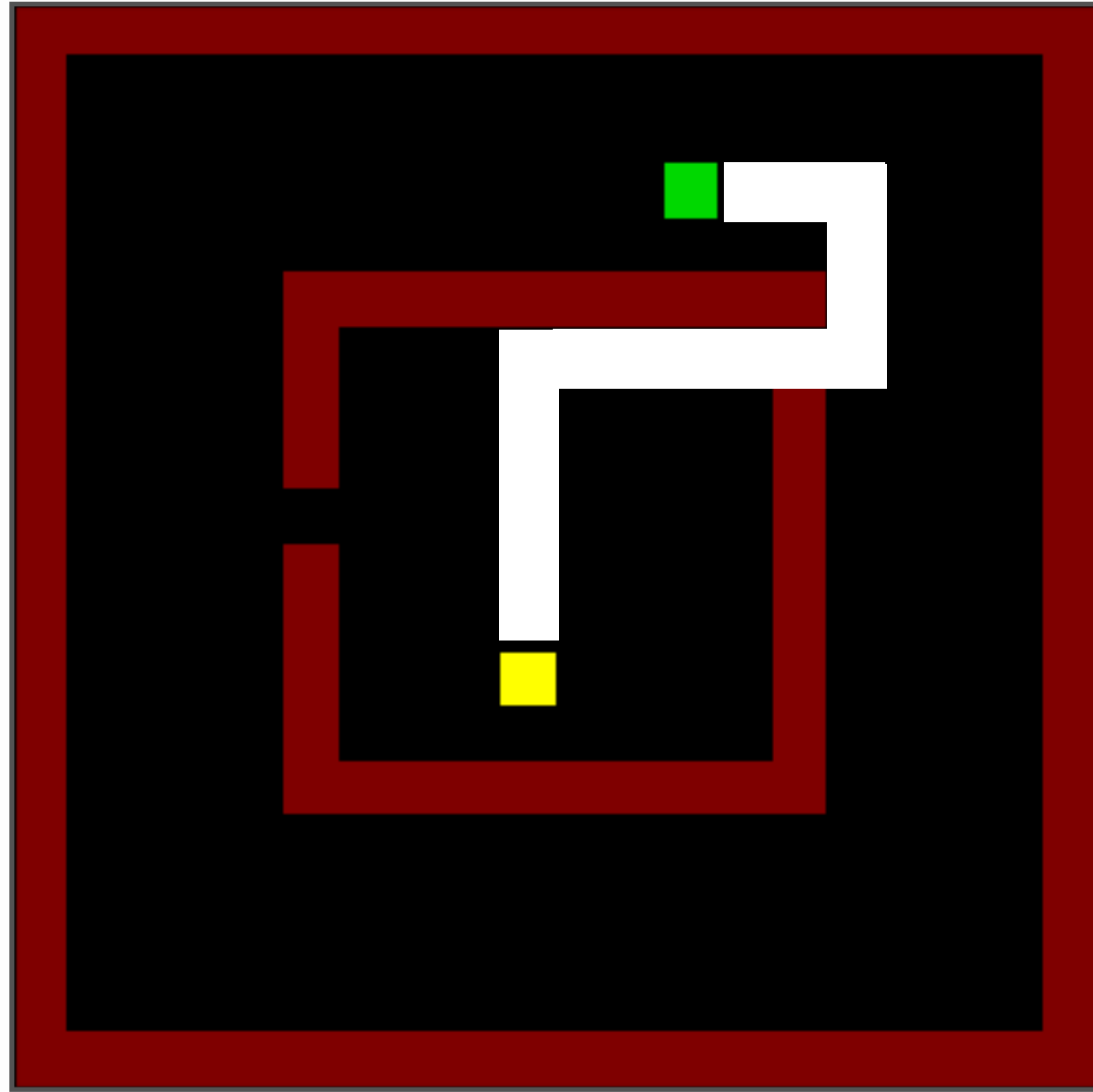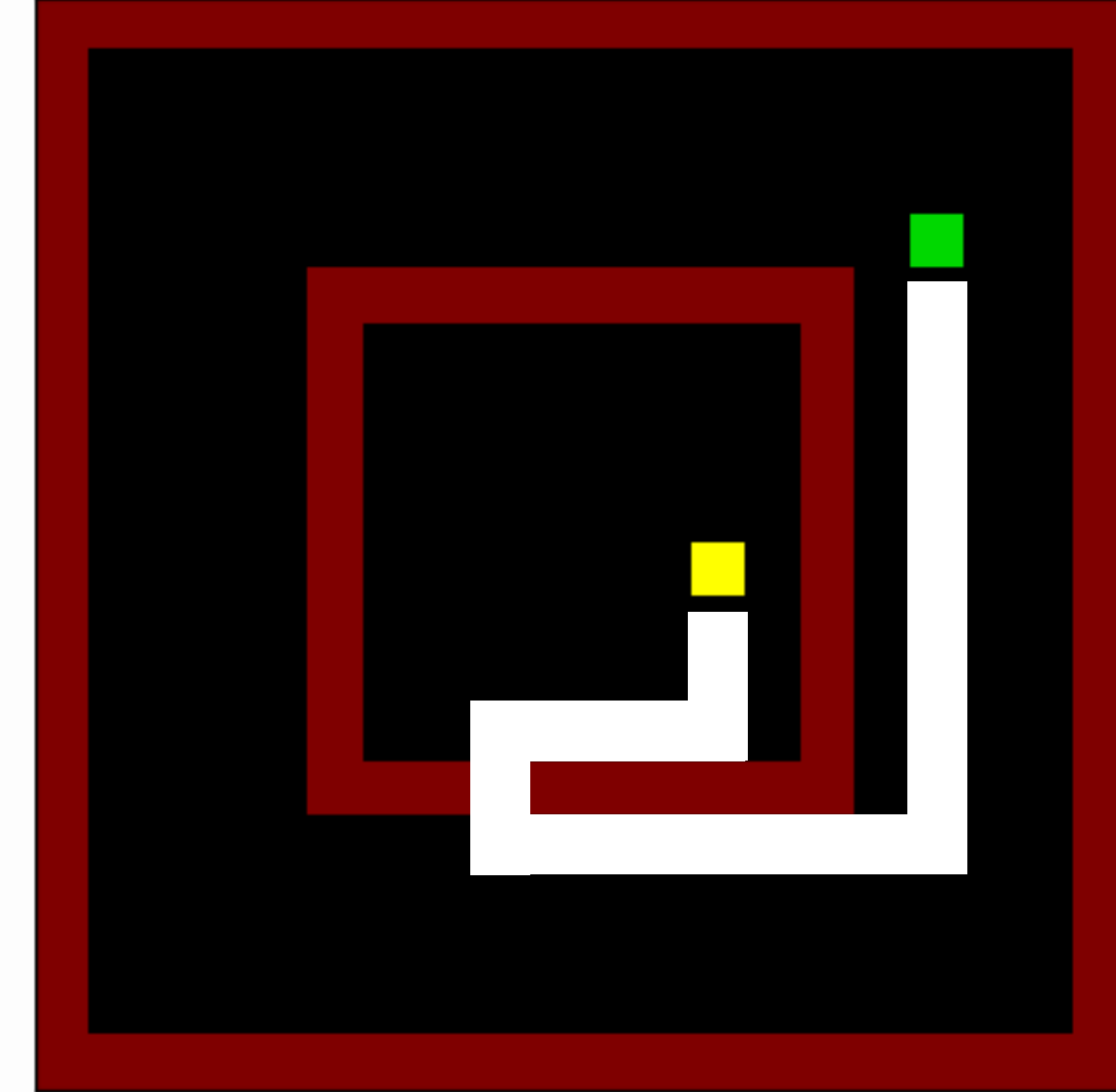| Model | Mean Z-score | Median Z-score |
|---|---|---|
| (Base) Model-free Q-Network | 0 | 0 |
| Model-based Search | 0.1 | 0.05 |
| TreeQN | 0.27 | 0.12 |
| Differentiable Tree Search Network | 0.35 | 0.27 |

# DOMAIN: **NAVIGATION**

- A simple 10 x 10 grid with 2 areas, i.e. central room and hallway.

- Starts inside the central room and Goal is in the hallway

- Random **Exits x2**

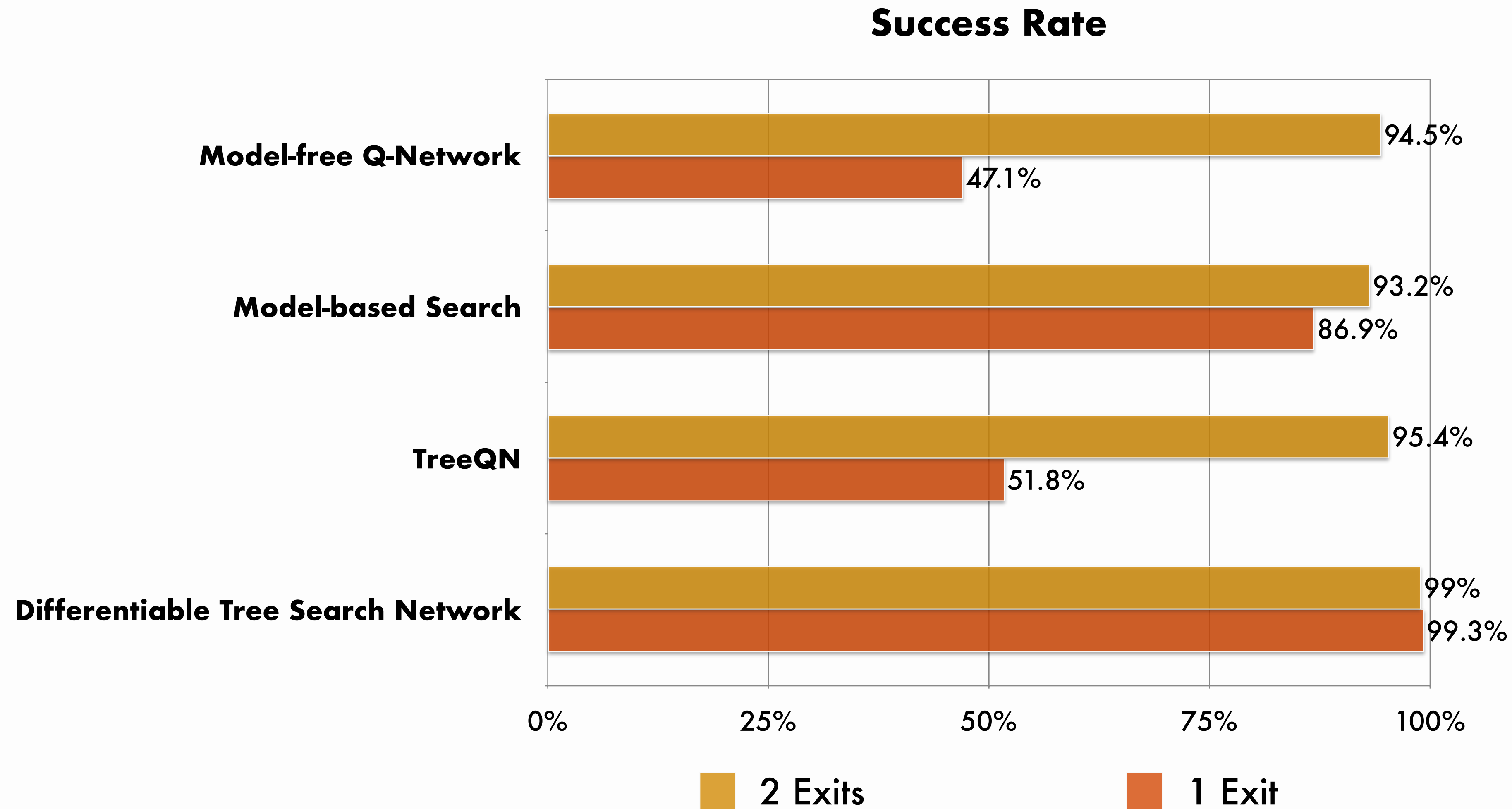- Collected **1000 Trajectories** for training

Trained on 2 Exits



Tested on 1 Exit

Success Rate

# SUMMARY OF CONTRIBUTIONS

## Differentiable Tree Search Network (D-TSN)

**Modular Neural Network Architecture**
that comprises of several *learnable submodules*

**Algorithmic Inductive Bias**
of a flexible and scalable *best-first tree search* algorithm

**Learnt World Model**
that is trained to be *useful* for the online search, even if inaccurate

**Joint Optimisation**
Trains the *search* and *world model* submodules jointly

**Additional Technical Details**
like *Tree Expansion Policy* and *Telescopic Sum for Variance Reduction*

# Thank You!