



# REGRET-GUIDED SEARCH CONTROL FOR EFFICIENT LEARNING IN ALPHAZERO

Yun-Jui Tsai<sup>1,2</sup>, Wei-Yu Chen<sup>2,3</sup>, Yan-Ru Ju<sup>2</sup>, Yu-Hung Chang<sup>2</sup>, Ti-Rong Wu<sup>2</sup>

<sup>1</sup> National Yang Ming Chiao Tung University Academia Sinica

<sup>2</sup> Academia Sinica

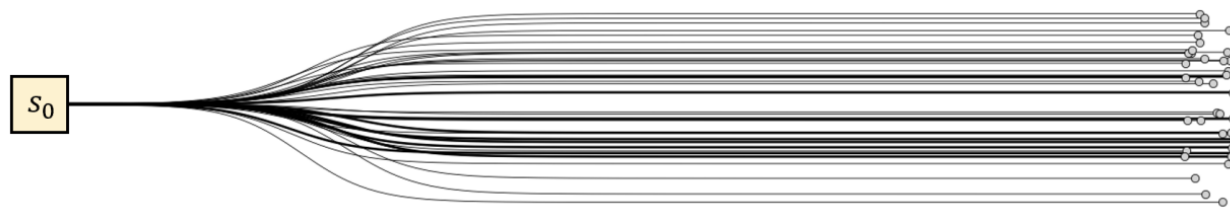
<sup>3</sup> College of Computing, Georgia Institute of Technology



# AlphaZero and Human's Learning Efficiency

## AlphaZero Learning

- Always restarts from an empty board
- Required large training resources

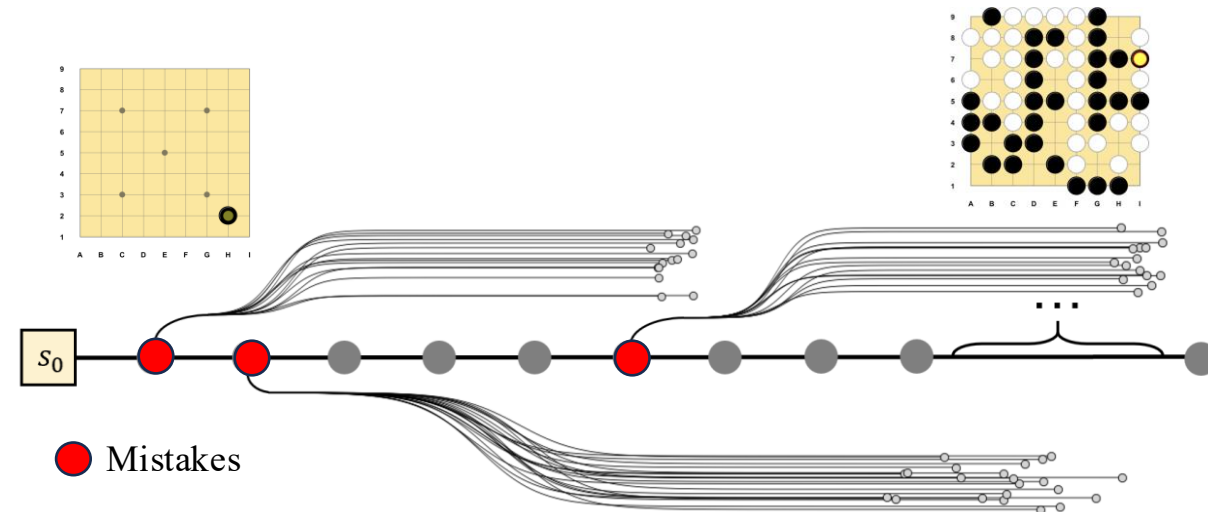


AlphaZero Zero-knowledge Learning

$s_0$  represents the initial state

## Human Learning

- Review the critical positions repeatedly
- Learn faster under restricted training resources
- A kind of self-play search control



Human Learning

# Search Control in AlphaZero

- Search Control: “*The process that selects the starting states and actions for the simulated experiences generated by the model*”<sup>[1]</sup>”
- Previous works for search control in AlphaZero:
  - Starting from random states in the past trajectories
    - Accelerating Self-Play Learning in Go (KataGo)<sup>[2]</sup>
    - Targeted Search Control in AlphaZero for Effective Policy Improvement (Go-Exploit)<sup>[3]</sup>
- We proposed RGSC, which extends AlphaZero
  - **Automatically identify unfamiliar (high-regret) states** for the current model

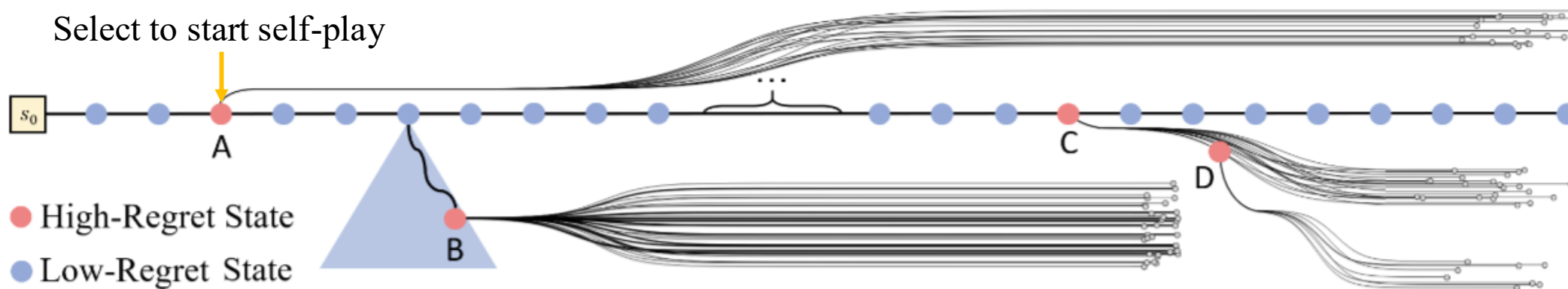
[1] Richard S. Sutton and Andrew G. Barto. Reinforcement Learning: An Introduction. Adaptive Computation and Machine Learning Series. MIT Press, Cambridge, MA, USA, 2 edition, November 2018.

[2] David J. Wu. Accelerating Self-Play Learning in Go. In Proceedings of the AAAI Workshop on Reinforcement Learning in Games, November 2020.

[3] Alexandre Trudeau and Michael Bowling. Targeted Search Control in AlphaZero for Effective Policy Improvement. In Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems, AAMAS '23, pp. 842–850, Richland, SC, May 2023. International Foundation for Autonomous Agents and Multiagent Systems.

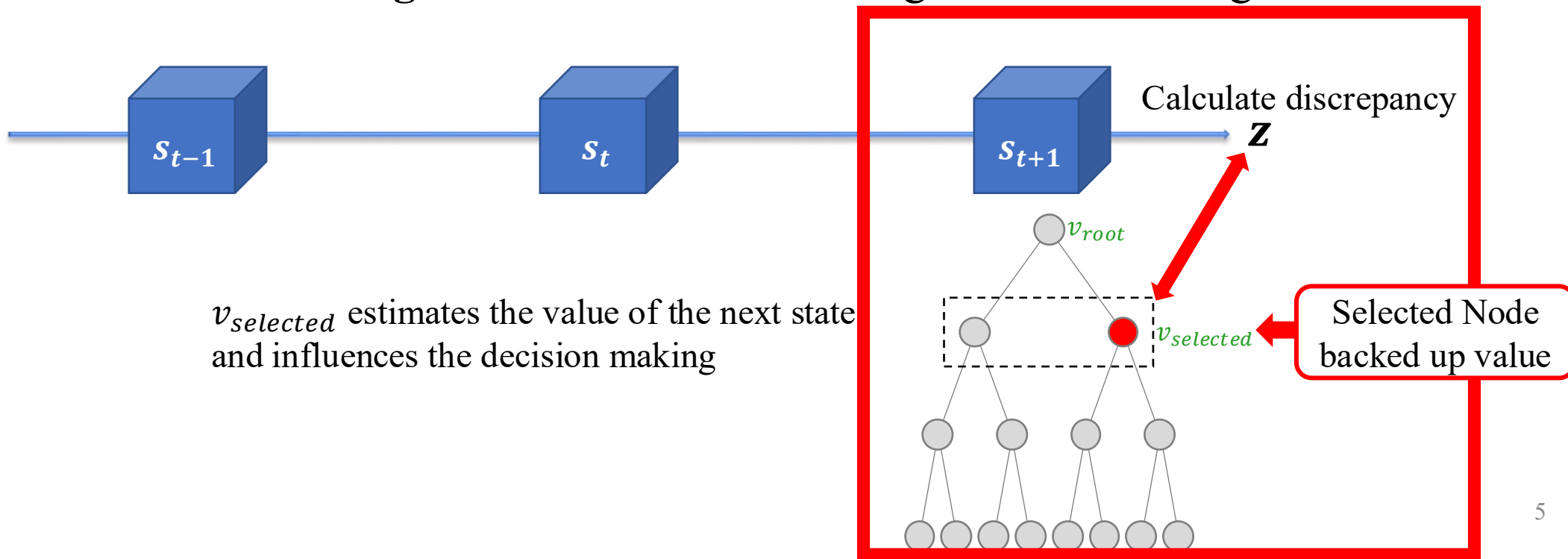
# Regret in AlphaZero

- Motivation:
  - Aim to **select critical states** as starting states, like humans reviewing their mistakes
- Regret is used to determine the critical state
  - Regret approximates the distance between the current policy and the optimal policy



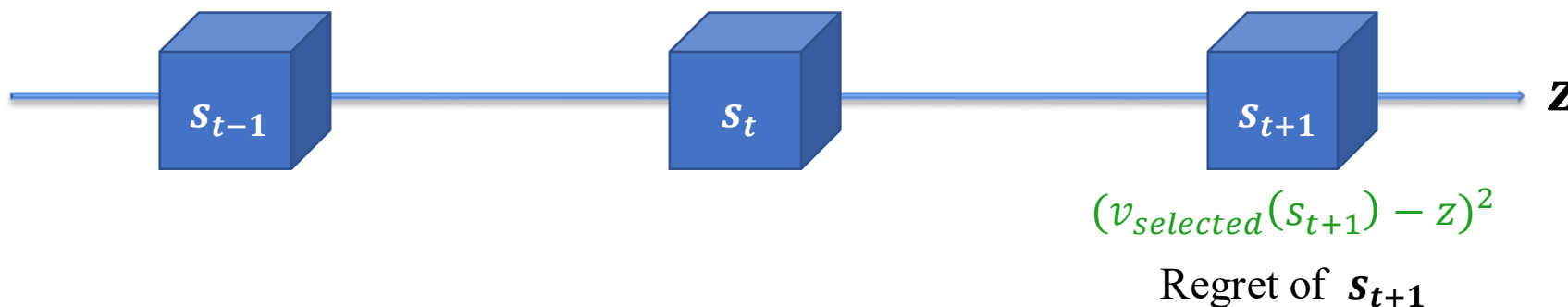
# Regret in AlphaZero

- For every single self-play trajectory  $s_0, s_1, s_2 \dots s_T$ , we aim to find a state with the largest regret
- In RGSC, we estimate regret by measuring the average discrepancy between the **agent's evaluation** and **the game outcome** when the agent start from a given state



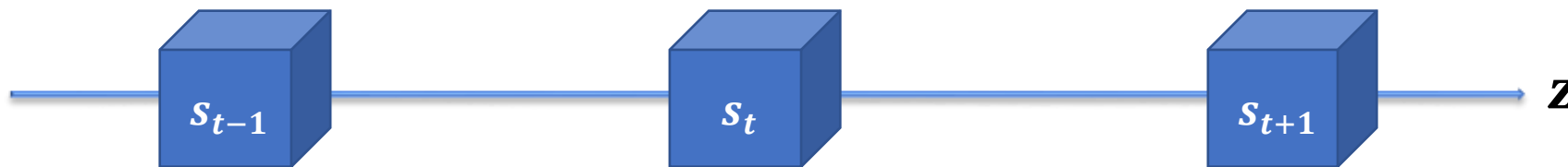
# Regret in AlphaZero

- For every single self-play trajectory  $s_0, s_1, s_2 \dots s_T$ , we aim to find a state with the largest regret
- In RGSC, we estimate regret by measuring the average discrepancy between the **agent's evaluation** and **the game outcome** when the agent start from a given state



# Regret in AlphaZero

- For every single self-play trajectory  $s_0, s_1, s_2 \cdots s_T$ , we aim to find a state with the largest regret
- In RGSC, we estimate regret by measuring the average discrepancy between the **agent's evaluation** and **the game outcome** when the agent start from a given state

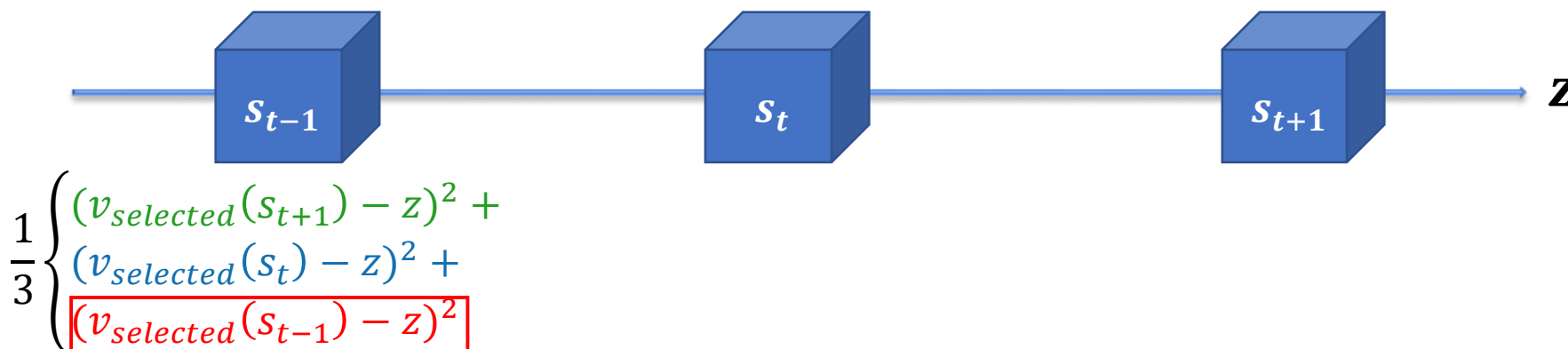


$$\frac{1}{2} \left\{ \begin{array}{l} (v_{selected}(s_{t+1}) - z)^2 + \\ (v_{selected}(s_t) - z)^2 \end{array} \right.$$

Regret of  $s_t$  represents the average discrepancy when an agent starts from  $s_t$

# Regret in AlphaZero

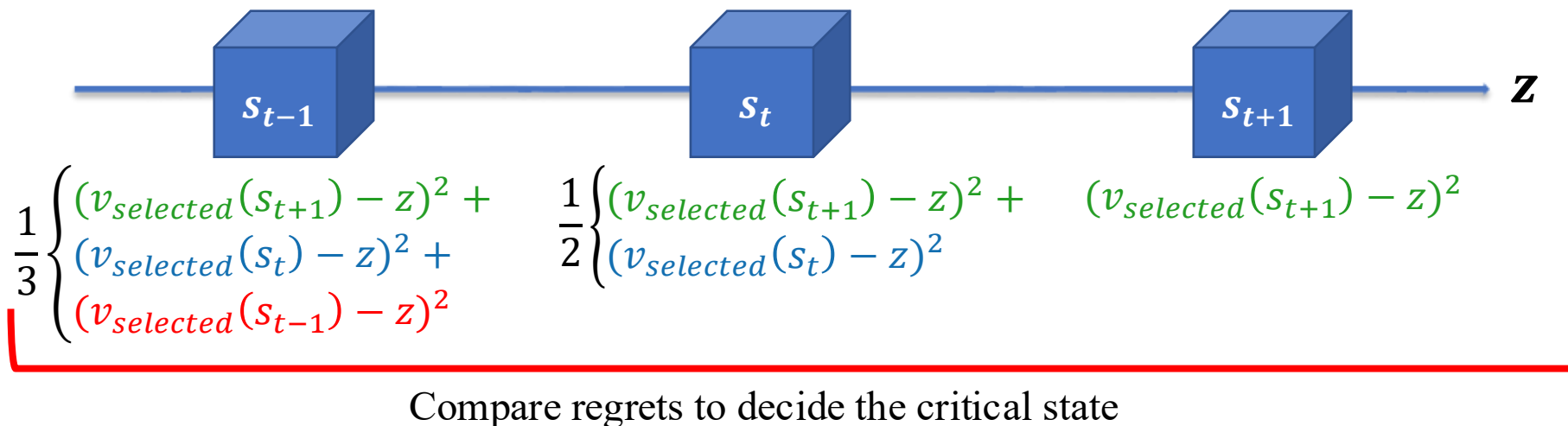
- For every single self-play trajectory  $s_0, s_1, s_2 \dots s_T$ , we aim to find a state with the largest regret
- In RGSC, we estimate regret by measuring the average discrepancy between the **agent's evaluation** and **the game outcome** when the agent start from a given state



Regret of  $s_{t-1}$  represents the average discrepancy when an agent starts from  $s_{t-1}$

# Regret in AlphaZero

- For every single self-play trajectory  $s_0, s_1, s_2 \dots s_T$ , we aim to find a state with the largest regret
- In RGSC, we estimate regret by measuring the average discrepancy between the **agent's evaluation** and **the game outcome** when the agent start from a given state

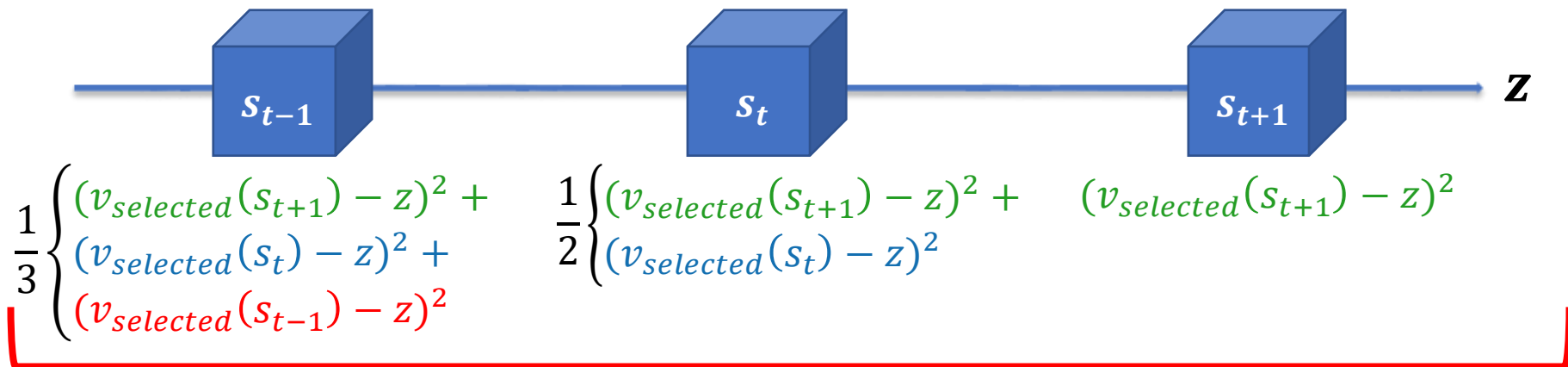


# Regret in AlphaZero

- We define regret as :

$$\mathcal{R}(s_t) = \frac{1}{T-t} \sum_{i=t}^T (v_{selected}(s_i) - z)^2$$

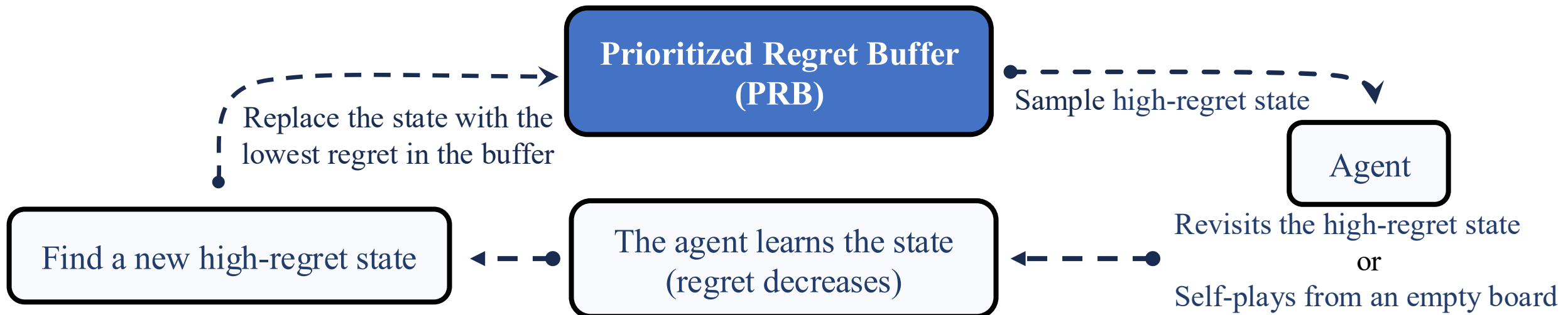
where,  $V_{selected}(s_i)$  represents MCTS value of the selected action at state  $s_i$



Compare regrets to decide the critical state

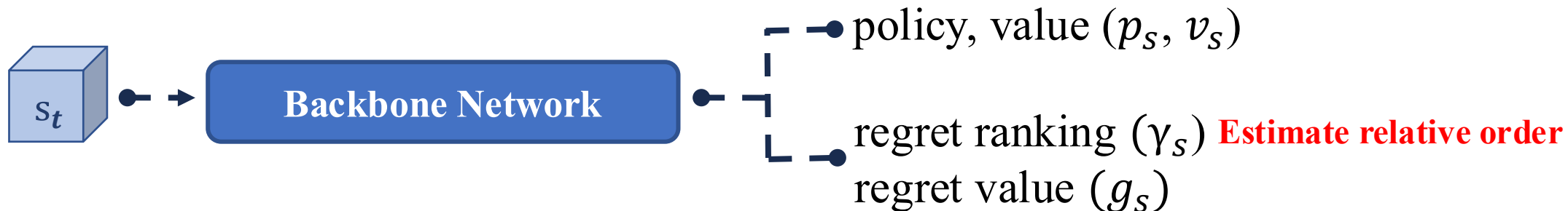
# Prioritized Regret Buffer (PRB)

- After getting high-regret states, **Prioritized Regret Buffer (PRB)** is used to mimic how humans learn



# Approximate Regret for MCTS Node

- Motivation:
  - Measure and select critical state among MCTS nodes
  - Make the model focus on learning the **relative order** among all nodes
- Regret ranking head predicts the ranking value  $\gamma_s$  for **selecting the opening**
- Regret value head predicts the regret value  $g_s$  for **initializing the opening**



# Regret Ranking Head Training

- According to Exploring Restart Distribution<sup>[1]</sup>, the training objective with search control is

$$L = \sum_{s \in S} \rho(s) \sum_{a \in A} \pi(a|s) (q(s, a) - \hat{q}(s, a))^2$$

where  $\rho(s)$  is the probability for sampling state  $s$

- Similarly, ranking training objective aims to maximize the regret of sampled states,

$$\mathcal{J}_{rank} = \sum_{s \in S} \rho(s) \mathcal{R}(s)$$

- Regret ranking head aims to approximate  $\rho(s)$  to get the rank of states

$$\mathcal{J}_{rank} = \sum_{s \in S} \text{softmax}(\gamma_s) \mathcal{R}(s)$$

where  $\gamma_s$  is the ranking head output, higher  $\mathcal{R}(s)$  implies higher  $\gamma_s$

# Regret Ranking Head Training

- Optimized by using a surrogate objective

$$\mathcal{J}_{rank} = \sum_{s \in S} \text{softmax}(\gamma_s) \mathcal{R}(s); \tilde{\mathcal{J}}_{rank} = \sum_{s \in S} \text{softmax}(\gamma_s) \exp\{\mathcal{R}(s)\}$$

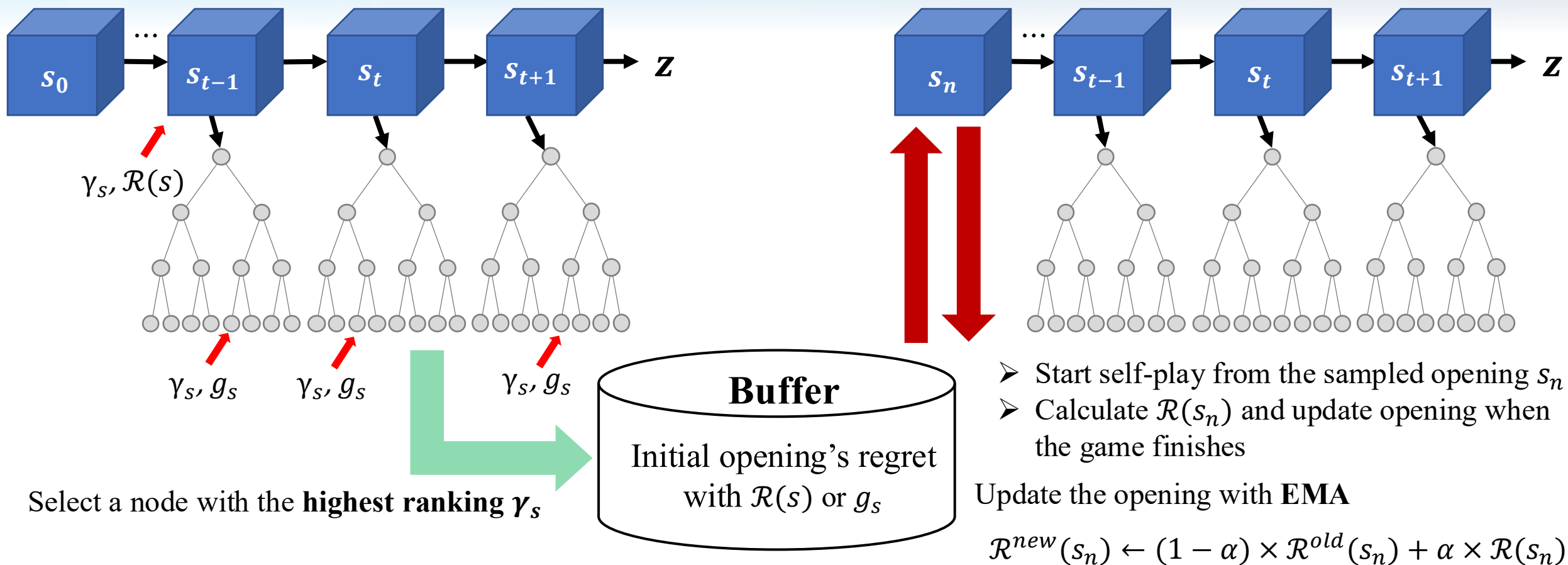
- Since  $\exp\{\mathcal{R}(s)\}$  is a strictly increasing function, the order of  $\gamma_s$  to optimize  $\mathcal{J}_{rank}$  equals the order to optimize  $\tilde{\mathcal{J}}_{rank}$
- The loss function of regret ranking head is

$$\begin{aligned} \mathcal{L}_{rank} &= -\log \tilde{\mathcal{J}}_{rank} = -\log \sum_{s \in S} \text{softmax}(\gamma_s) \exp\{\mathcal{R}(s)\} \\ &= -\log \sum_{s \in S} \exp(\log(\text{softmax}(\gamma_s)) + \mathcal{R}(s)) \\ &= -\log \sum_{s \in S} \exp(\log(\text{softmax}(\gamma_s)) + \mathcal{R}(s)) \end{aligned}$$

# Regret-Guided Search Control (RGSC)

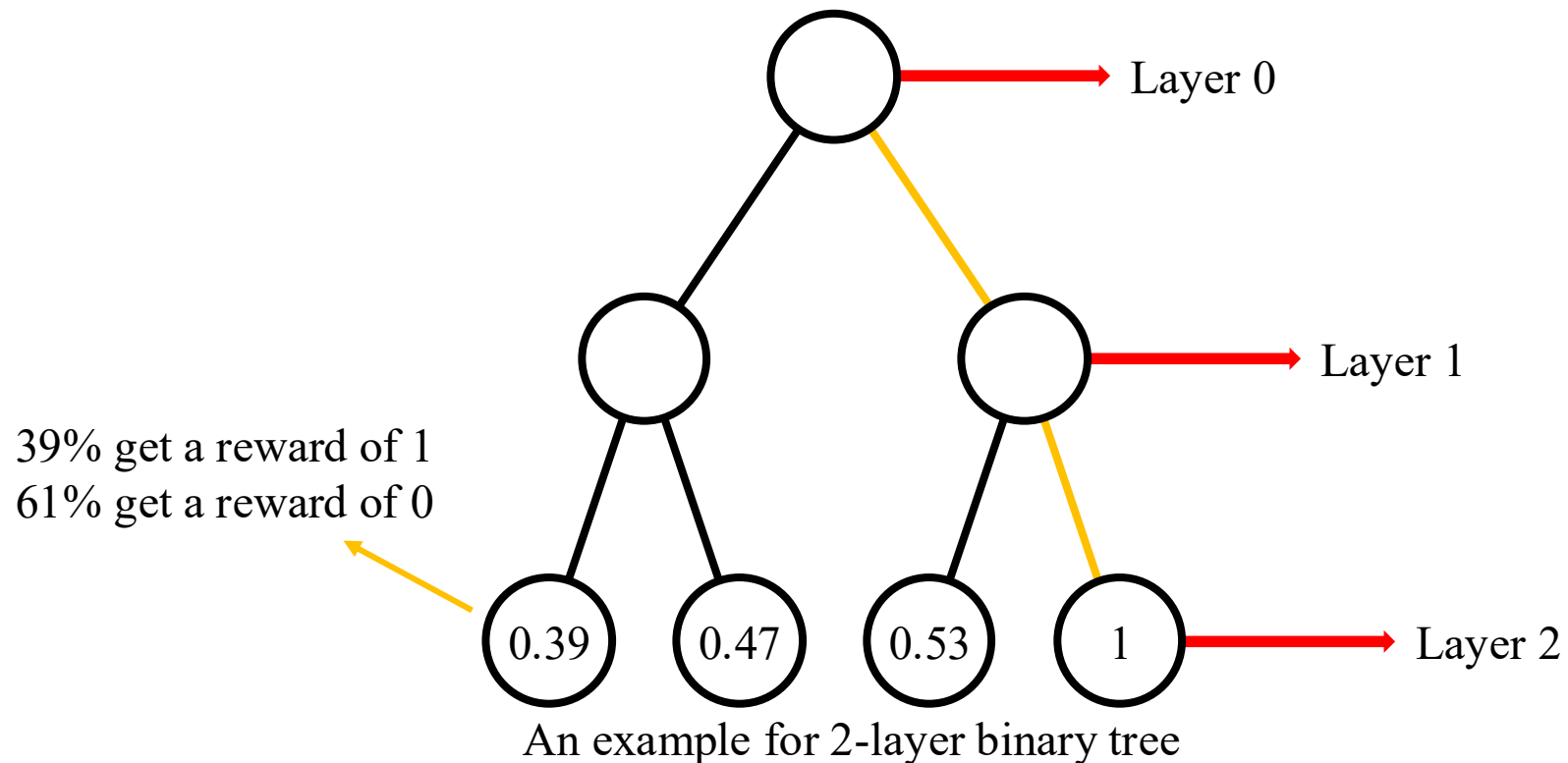
- RGSC consists of two parts:
  - $\lambda$  probability to start self-play from the opening sampled from PRB
    - Update the opening's regret in PRB
  - $1 - \lambda$  probability to start self-play from the initial state
    - Selecting a new opening into PRB

# RGSC Overview



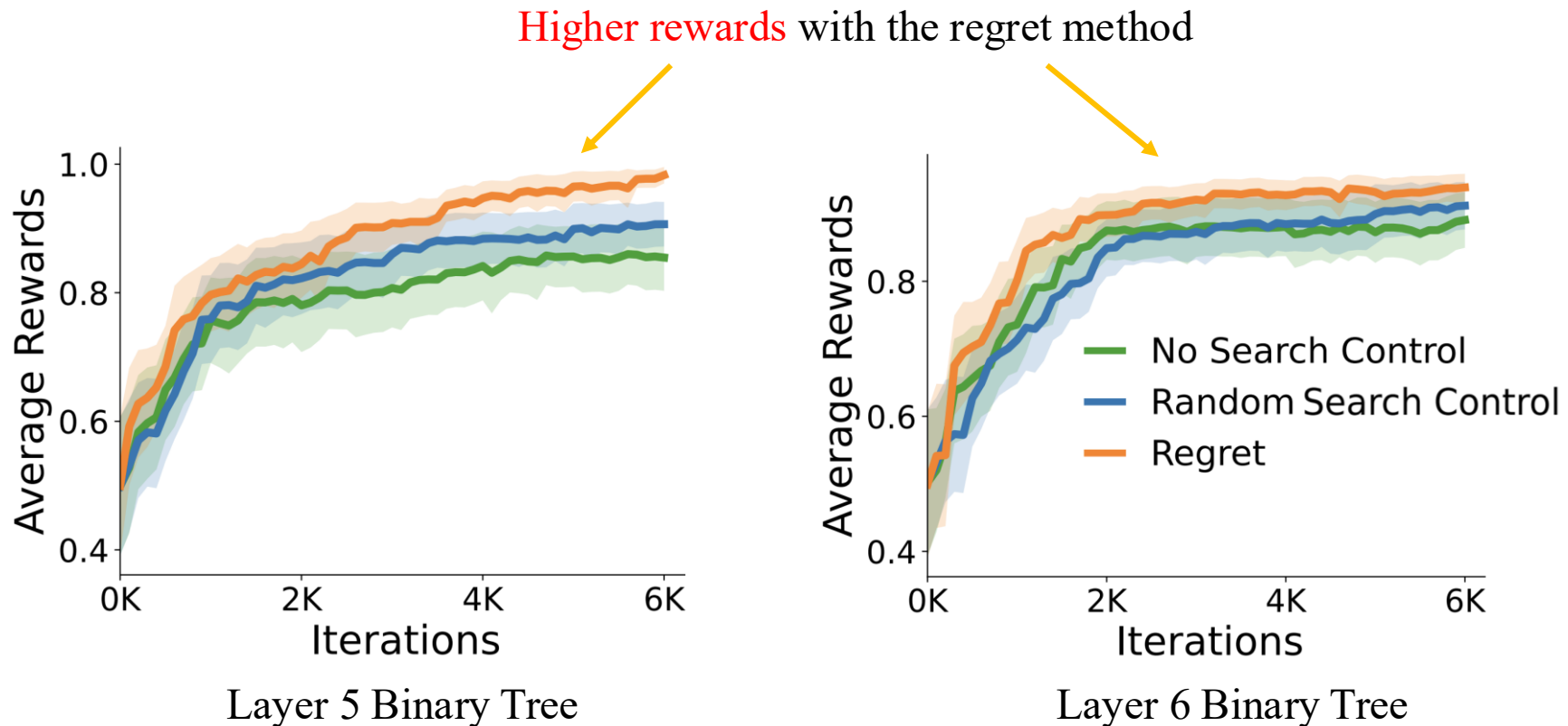
# Toy Example

- Motivation: Prove regret method **accelerates Q-learning** in a small binary tree
- Goal: Find the optimal path (orange path)



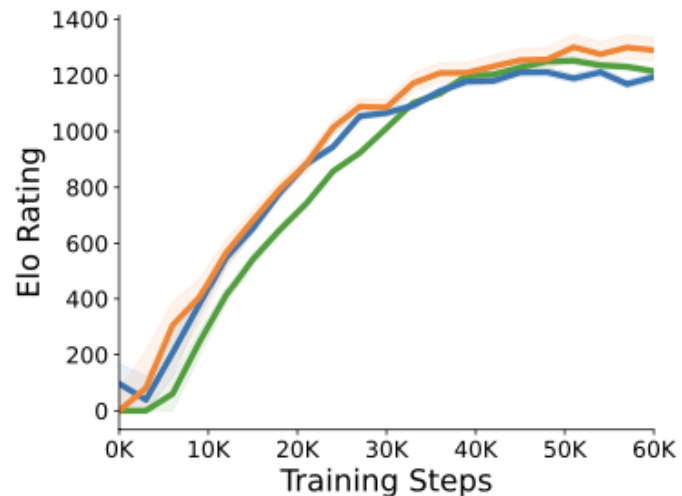
# Toy Example

- Average Reward: **Regret** > Random > No Search Control
- These results show the **effectiveness of the regret-guided search control**

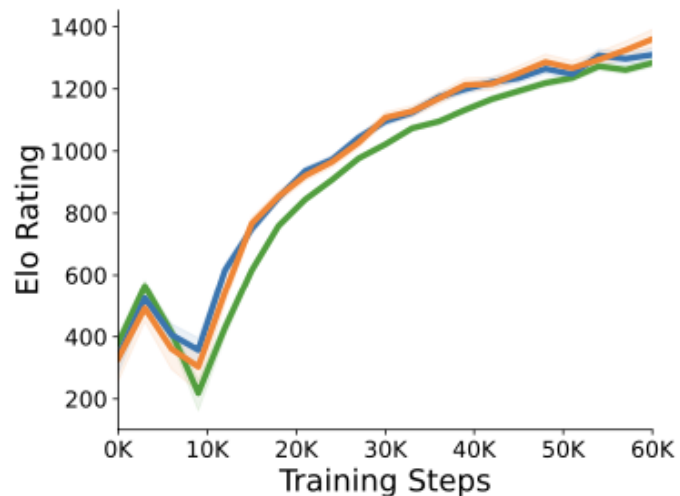


# RGSC in Board Games

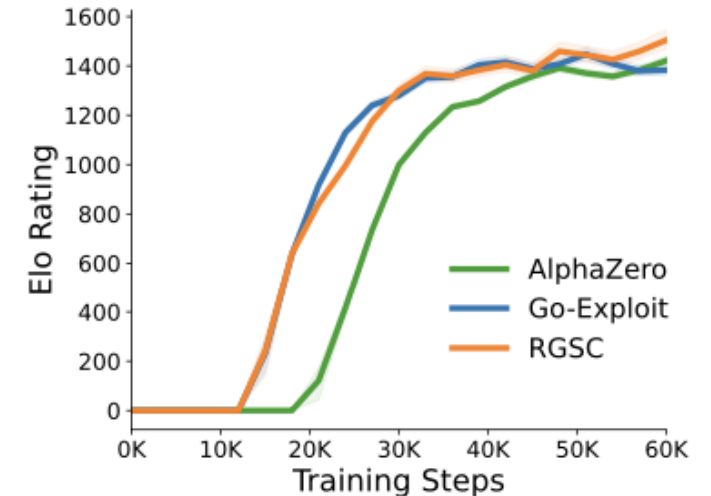
- Although Go-Exploit achieves a higher Elo than AlphaZero in the early stages of training, its advantage diminishes as training converges
- RGSC, however, **maintains its advantage** in the later stages of training



(a) 9x9 Go



(b) 10x10 Othello



(c) 11x11 Hex

Playing performance of AlphaZero, Go-Exploit, and RGSC

# RGSC in Board Games

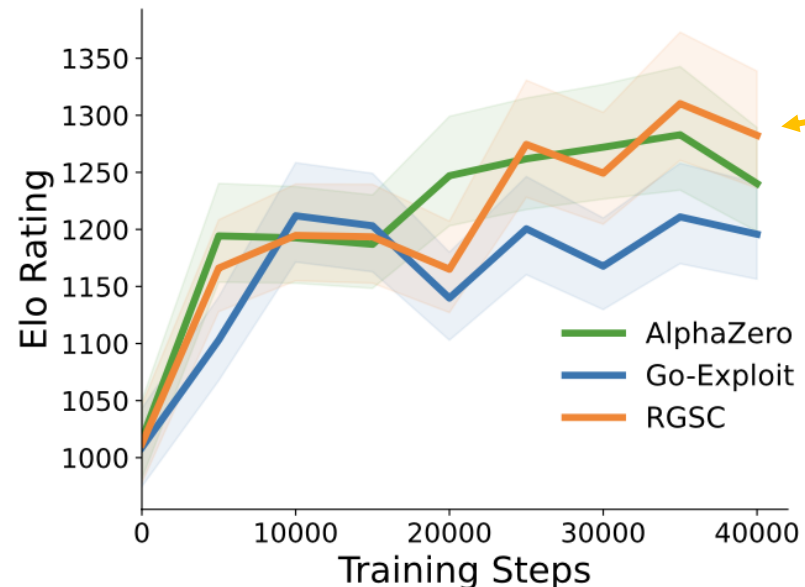
- Although Go-Exploit achieves a higher Elo than AlphaZero in the early stages of training, its advantage diminishes as training converges
- RGSC, however, **maintains its advantage** in the later stages of training

	AlphaZero	Go-Exploit	RGSC
9x9 Go	45.5% $\pm$ 1.5%	49.5% $\pm$ 2.0%	<b>53.6% <math>\pm</math> 2.4%</b>
10x10 Othello	51.7% $\pm$ 2.5%	52.9% $\pm$ 3.3%	<b>57.8% <math>\pm</math> 3.2%</b>
11x11 Hex	83.6% $\pm$ 1.6%	89.2% $\pm$ 1.8%	<b>91.1% <math>\pm</math> 2.0%</b>

Winrate at **60000** steps against established open-source programs  
(KataGo, Ludii alpha-beta search and Mohex)

# RGSC on Well-trained Model

- Motivation: Prove RGSC surpasses other methods when training on a well-trained model
- Training on a well-trained model for 40,000 steps
- RGSC achieves **a substantially higher win rate**, outperforming all baselines

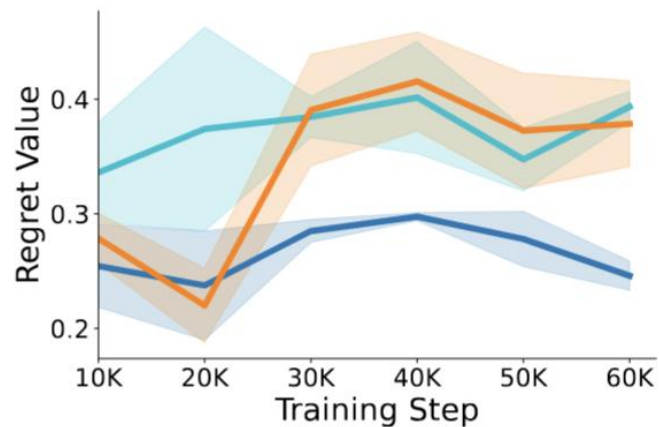


Baseline(0 step)	AlphaZero	Go-Exploit	RGSC
69.3% ± 2.6%	70.2% ± 2.7%	69.2% ± 2.7%	<b>78.2% ± 2.5%</b>

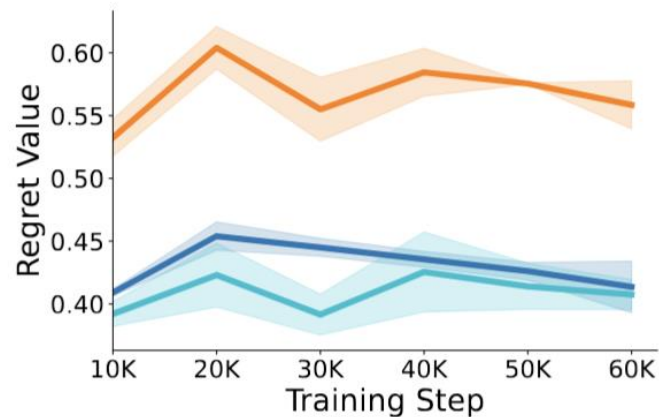
Winrate at 40000 steps against KataGo

# Comparison between Ranking and Regret

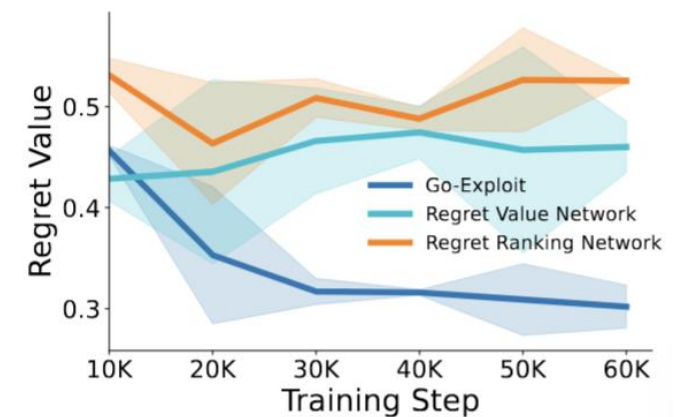
- The regret ranking network consistently **selects states with higher true regret** than the regret value network



(a) 9x9 Go



(b) 10x10 Othello

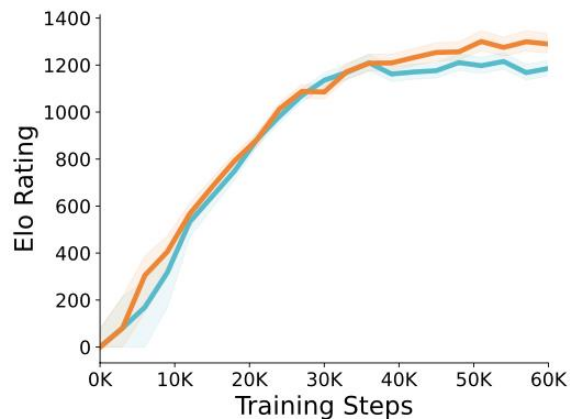


(c) 11x11 Hex

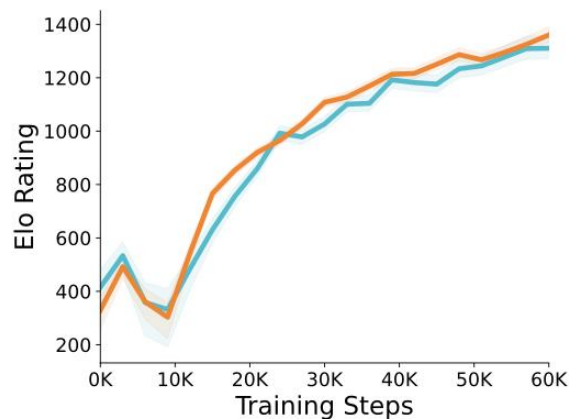
Average regret value of nodes selected by different methods

# Comparison between Ranking and Regret

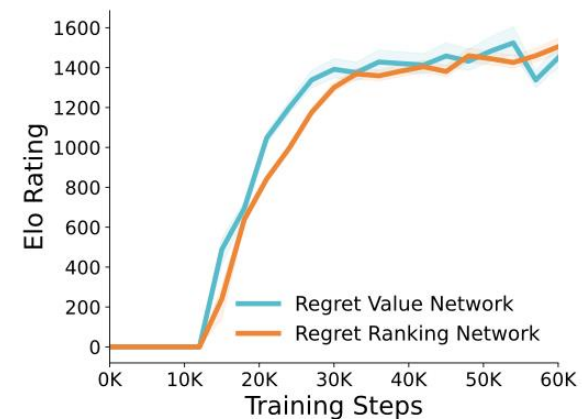
- Directly compare the training performance with regret value network and regret ranking network
- These results demonstrate that the **ranking objective improves the quality** of selected states by prioritizing those with greater learning potential



(a) 9x9 Go



(b) 10x10 Othello

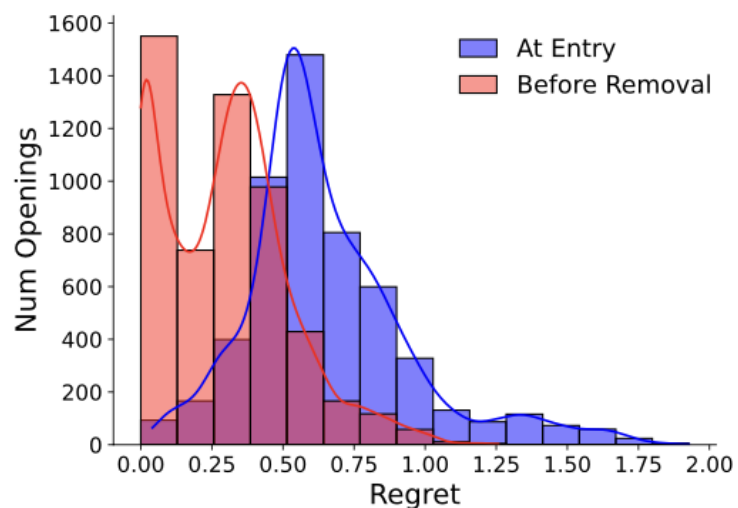


(c) 11x11 Hex

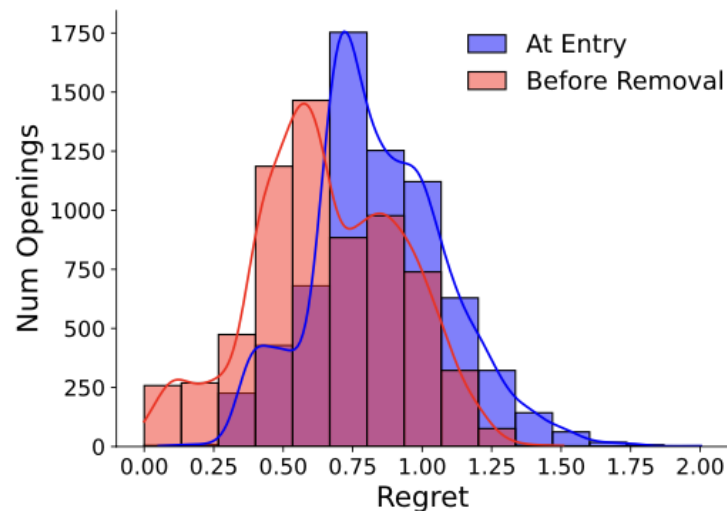
Playing performance of model selecting MCTS node using  
regret ranking network and regret value network

# Regret Decay in Buffer

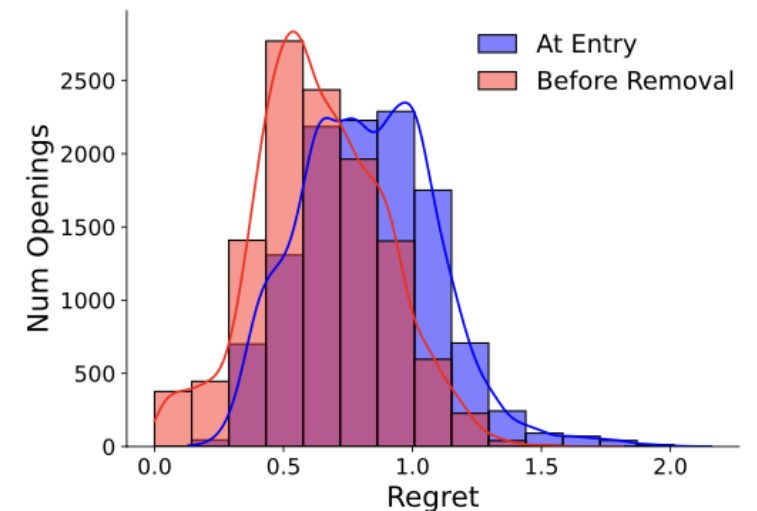
- High-regret states in the buffer gradually decrease during training
  - The model can actually **correct** its mistakes by repeatedly revisiting those states
- This experiment confirms that states initially associated with high regret are eventually corrected through repeated replay, resulting in reduced regret over time



(a) 9x9 Go



(b) 10x10 Othello



(c) 11x11 Hex

# Summary

- RGSC: An extension of AlphaZero that
  - **Identifies high-regret states**
  - Focus on **correcting its most critical mistakes**, mimicking how humans learn
  - Outperform AlphaZero and Go-Exploit when
    - Start training from a random weight
    - Continue training from a well-trained  $9 \times 9$  Go model weight
- Future work:
  - Extend RGSC to MuZero (already a preliminary experiment)
  - Enhance model explanation according to human preference
- Our results highlight regret-guided search control as a promising direction for **improving the efficiency and robustness** of reinforcement learning

# Thank You for Your Attention

Our code and data are available at

<https://rlg.iis.sinica.edu.tw/papers/rgsc/>

