# Uncertainty-Guided Exploration for Efficient AlphaZero Training

Scott Cheng[1], Meng-Yu Tsai[2], Ding-Yong Hong[3], Mahmut Taylan Kandemir[1]

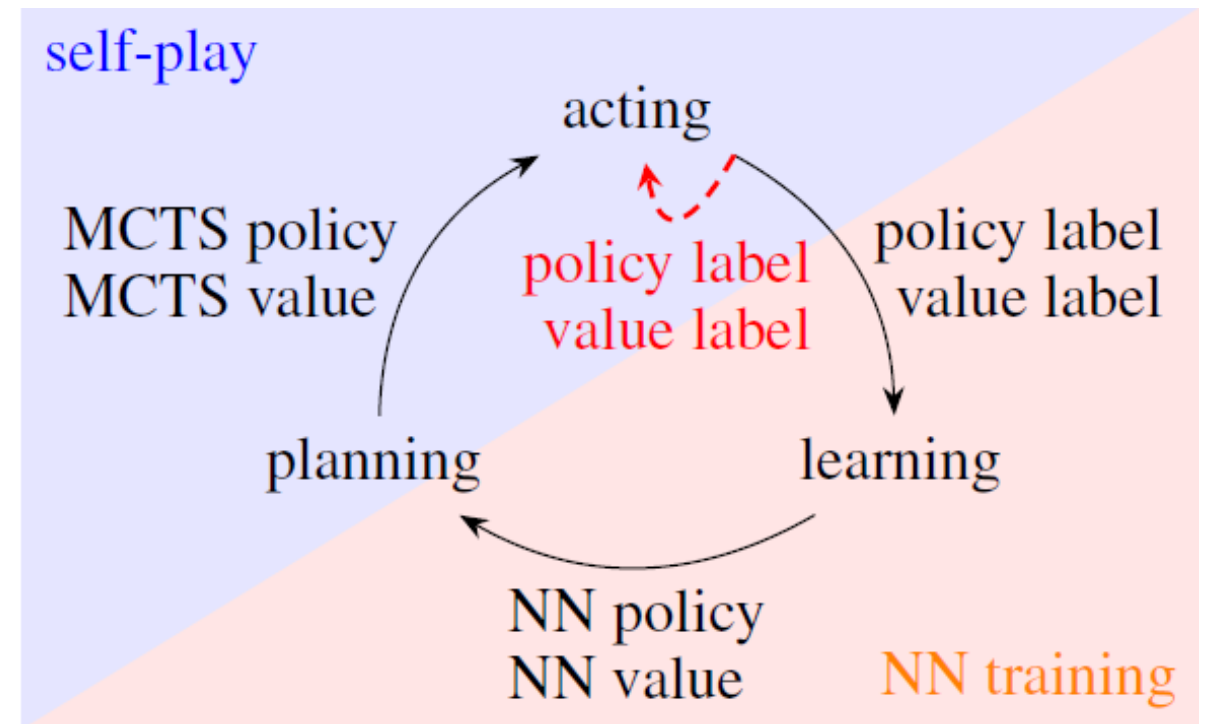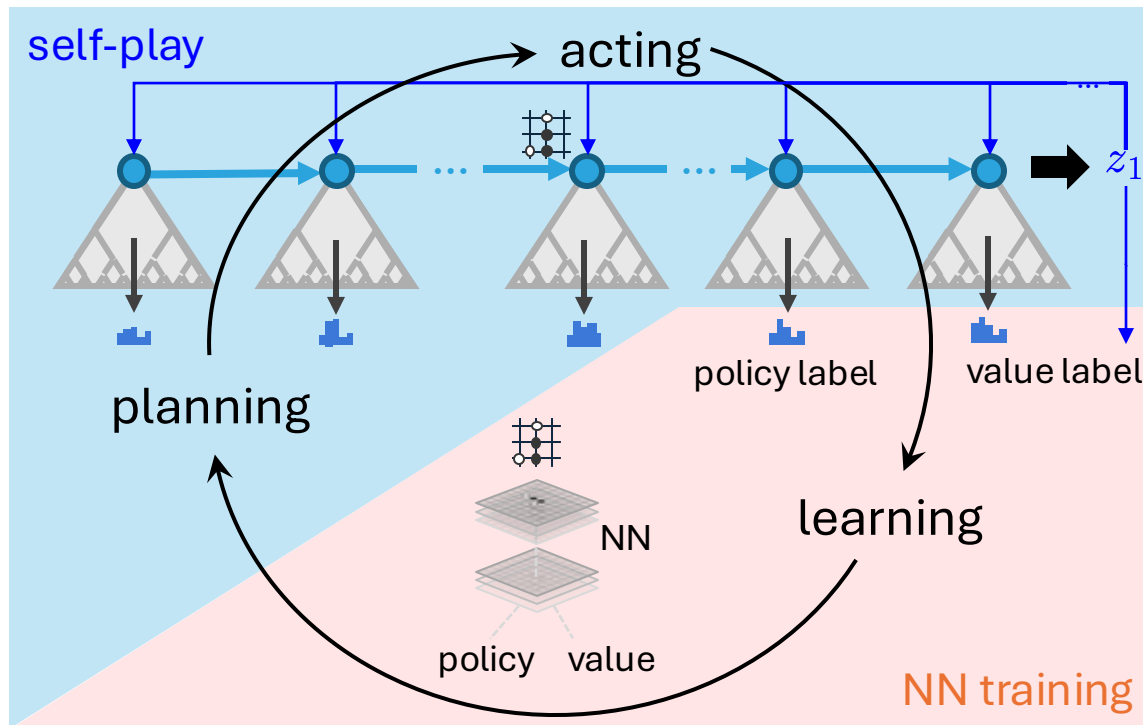[1] The Pennsylvania State University, USA
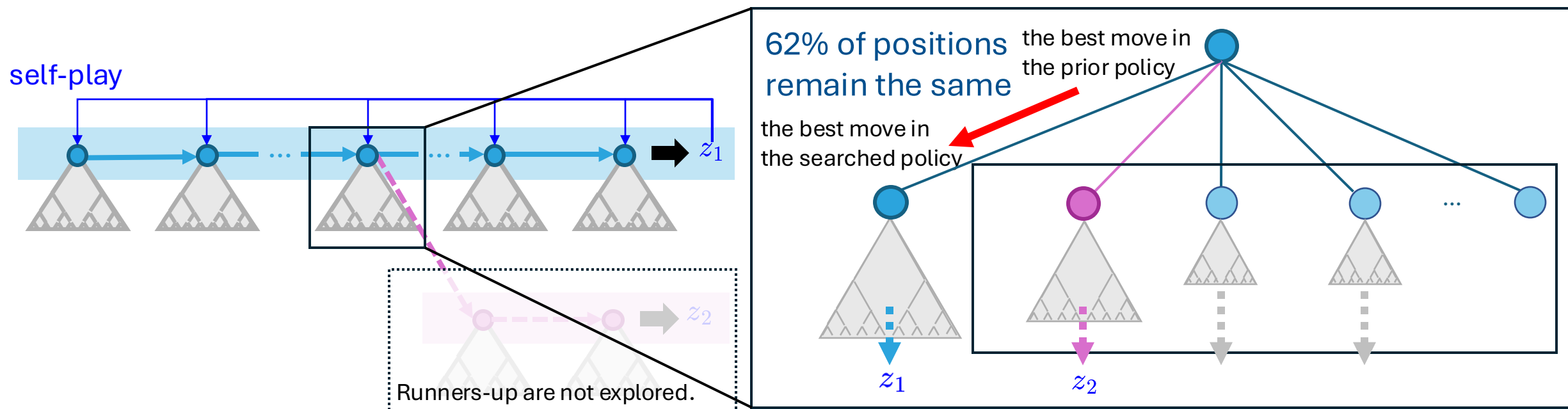[2] Independent
[3] Academia Sinica, Taiwan

# Background: AlphaZero

Many recent breakthroughs in artificial intelligence have been driven by the AlphaZero algorithm, which consists of self-play and NN training:
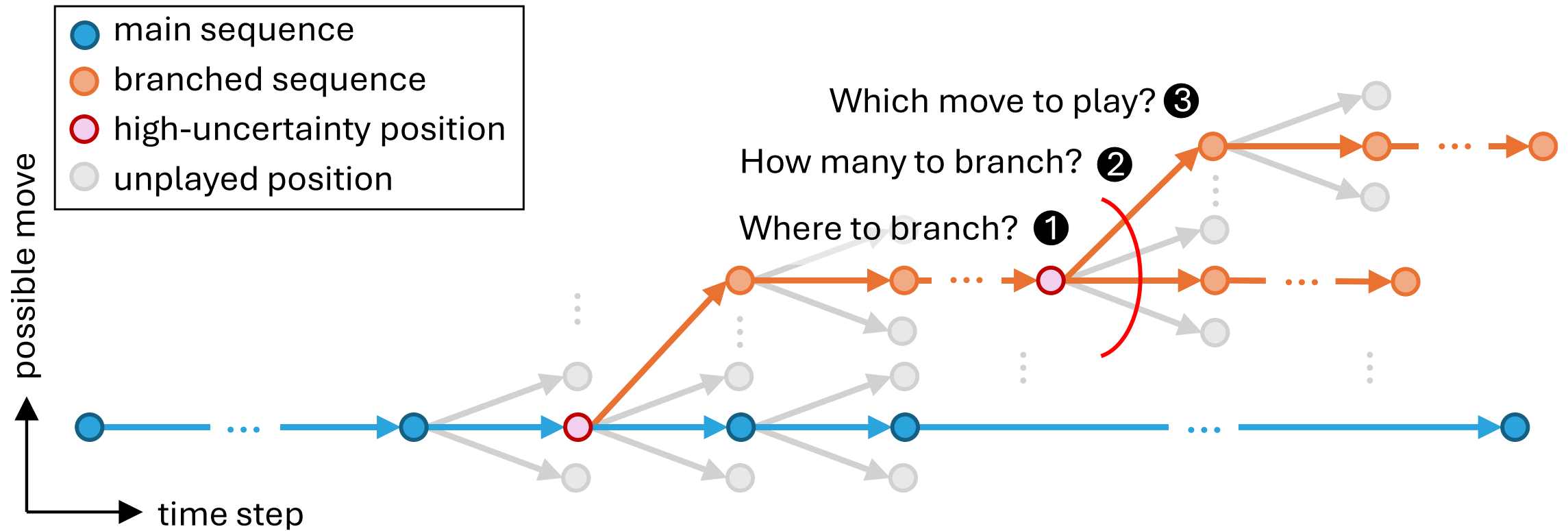
# AlphaZero's self-play process remains inefficient

(1) Value labels have high variance, as they are derived from a single result.

(2) A few positions exhibit high-uncertainty, but these positions are not further explored.



self-play

$z_1$

Runners-up are not explored.

$z_2$

62% of positions remain the same

the best move in the prior policy

the best move in the searched policy

$z_1$

$z_2$

(3) When deeper search changes the best move, 79% of the new best moves come from the top-3 runners-up.
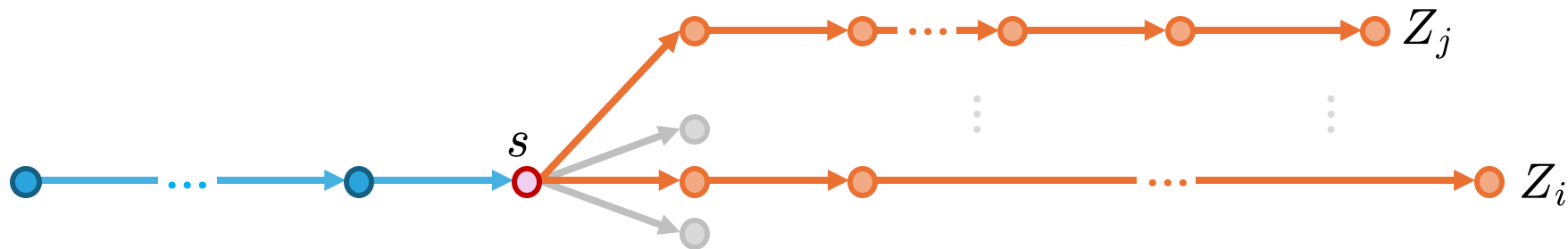
# Uncertainty-Guided Exploration during Self-play

# Label Change Rate (LCR)

Game result: $Z \sim \mathrm{Bernoulli}(w)$    Game plays: $Z_1, \dots, Z_n \overset{\text{i.i.d.}}{\sim} Z$

Uncertainty metric: $\theta = \mathrm{LCR}(s)$

For any independent plays $i \neq j$ ,

$$\mathrm{LCR}(s) := \mathrm{Pr}(Z_i \neq Z_j) = 2w(1-w)$$



$$\hat{\theta} = 2v(1-v) \quad \text{for MCTS value } v \in [0,1]$$

# Bayesian Inference

Bayesian view: model the LCR $\theta$ as a random variable.

Label change: $X = \mathbf{1}_{Z_i \neq Z_j}$

Observed data: $D = \{X_1, X_2, \ldots, X_m\}$ where $X_i \sim \mathrm{Bernoulli}(\theta)$

Prior LCR:

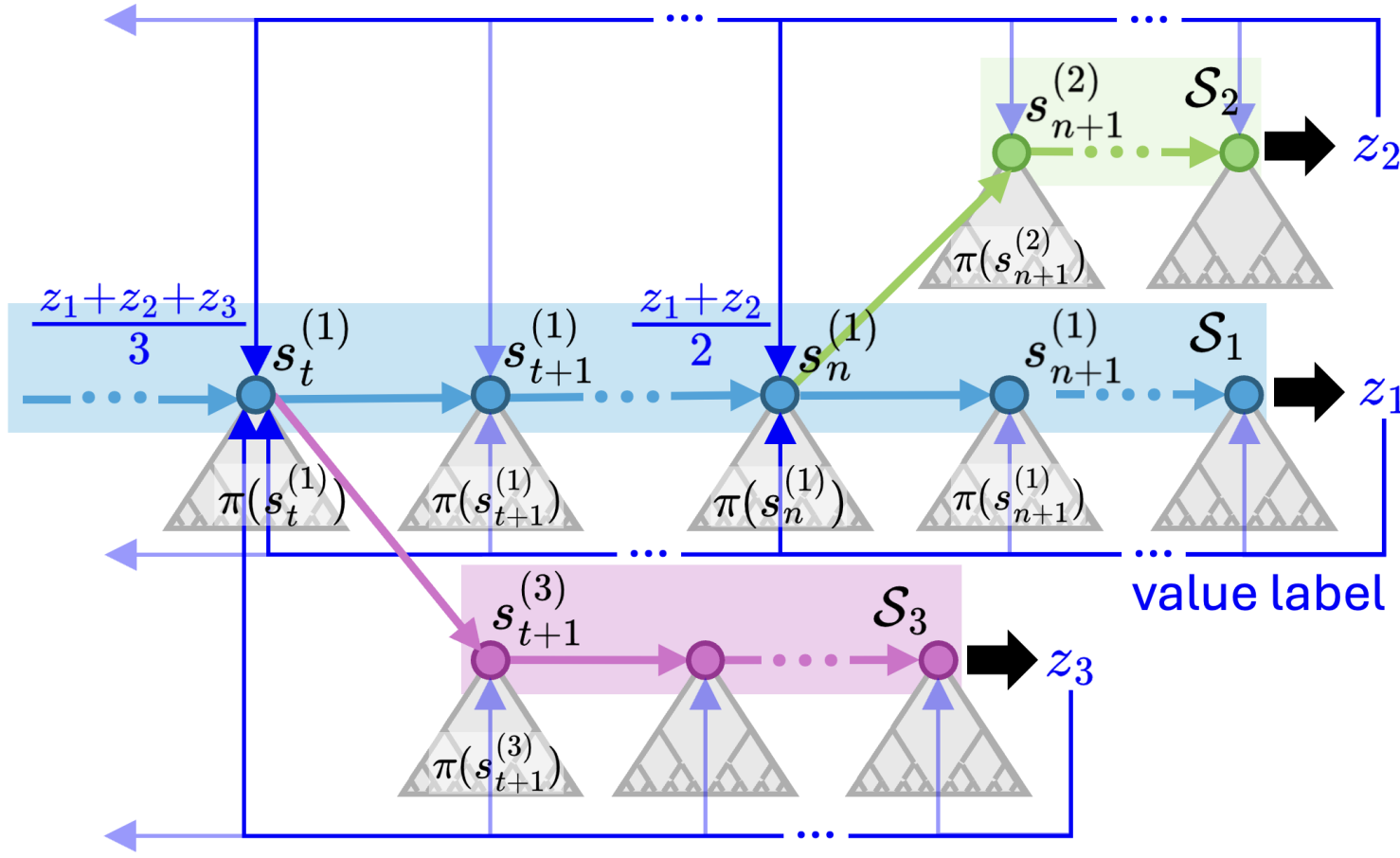$$\mathbb{E}[\theta] = \frac{\alpha}{\alpha + \beta}$$

where $\theta \sim \mathrm{Beta}(\alpha, \beta)$

Posterior LCR:

$$\mathbb{E}[\theta \mid D] = \frac{\alpha + s}{\alpha + \beta + m}$$

where $\theta \mid D \sim \mathrm{Beta}(\alpha + s, \beta + m - s)$

with $s$ observed result changes

# Uncertainty-Guided Branching



$$\mathrm{Var}[\bar{Z}] = w(1-w)/n$$

**Algorithm 1:** Uncertainty-Guided Branching

**Input:** $G$ number of states to collect
**Output:** $\mathcal{G} = \{(\text{state, policy label, value label})\}$
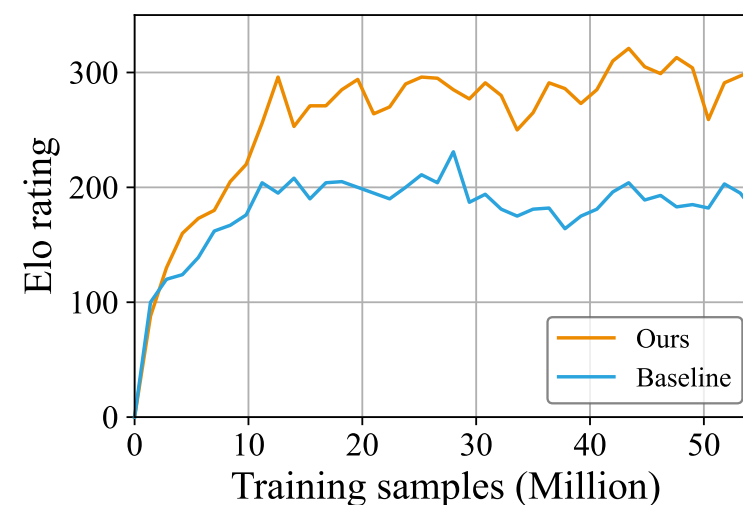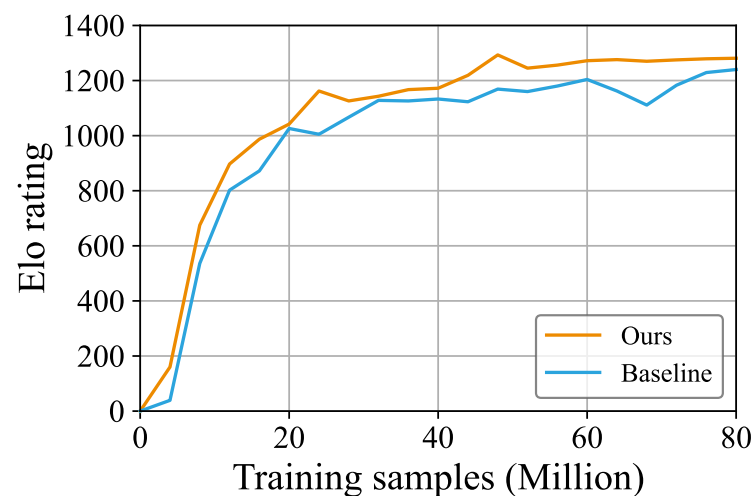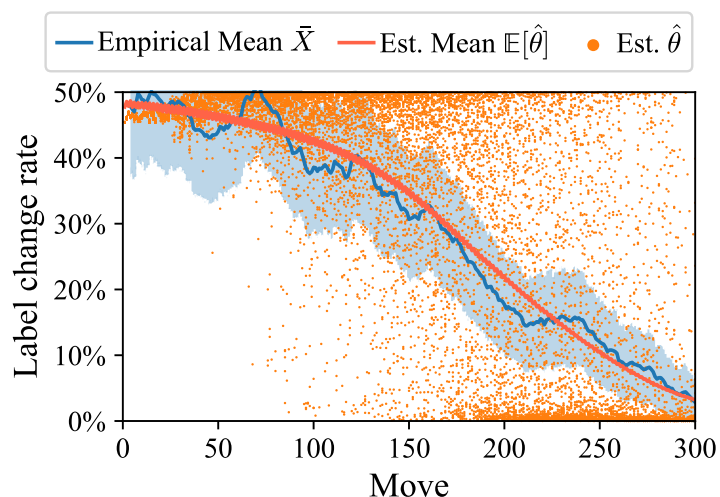
1   $\mathcal{G} \leftarrow \emptyset$
2   **while** $|\mathcal{G}| < G$ **do**
3     $(s, a) \leftarrow$ initial state and action
4     **for** $i \leftarrow 1$ **to** $V$ **do**
5       $\{(s', \pi(s'), z_i) \mid s' \in \mathcal{S}_i\} \leftarrow$
6        play episode from $(s, a)$
7       Compute LCR and sampling weights $u$ (defined in Equation 4)
8       $s \sim Y(u)$ (defined in Equation 5)
9       $a \sim \pi(\cdot \mid s)$
10       **while** $(s, a)$ has been played **do**
11        $s \leftarrow$ preceding state of $s$
12        $a \sim \pi(\cdot \mid s)$
13     $\mathcal{S} := \bigcup_{i=1}^{V} \mathcal{S}_i$
14     $\mathcal{G} \leftarrow \mathcal{G} \cup \left\{ \left(s', \pi(s'), \frac{\sum_{i:\, s' \in \mathcal{S}_i} z_i}{|\{i:\, s' \in \mathcal{S}_i\}|}\right) \mid s' \in \mathcal{S} \right\}$
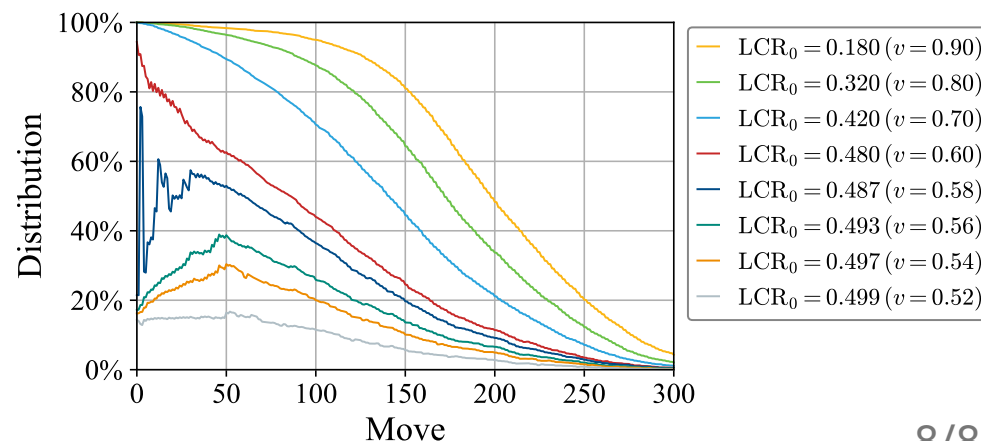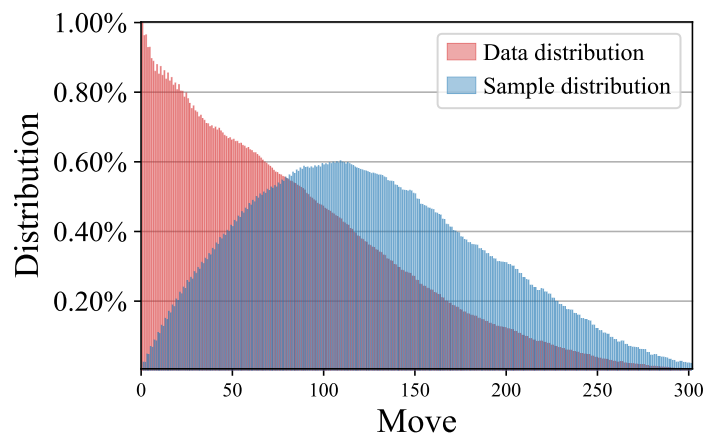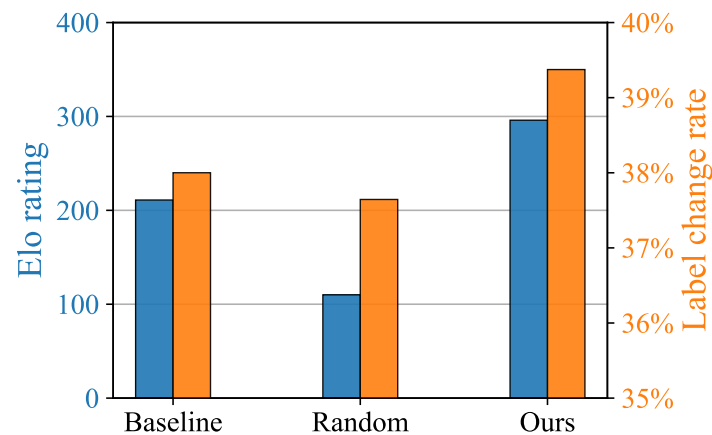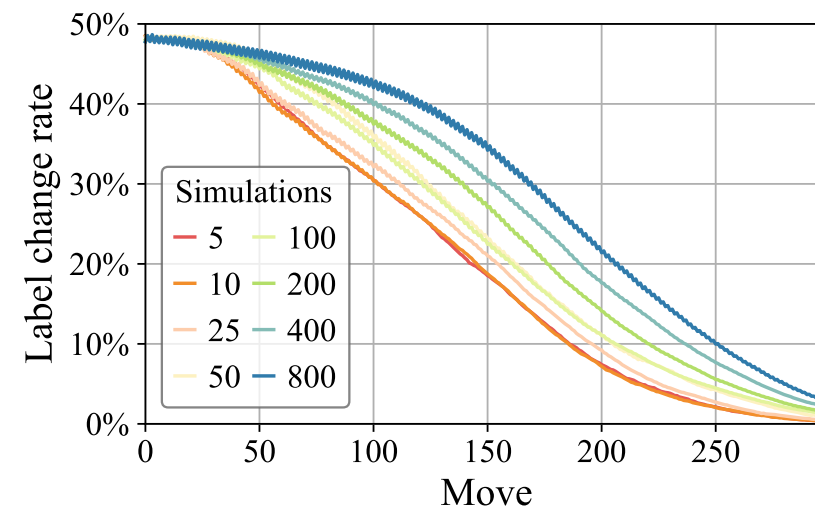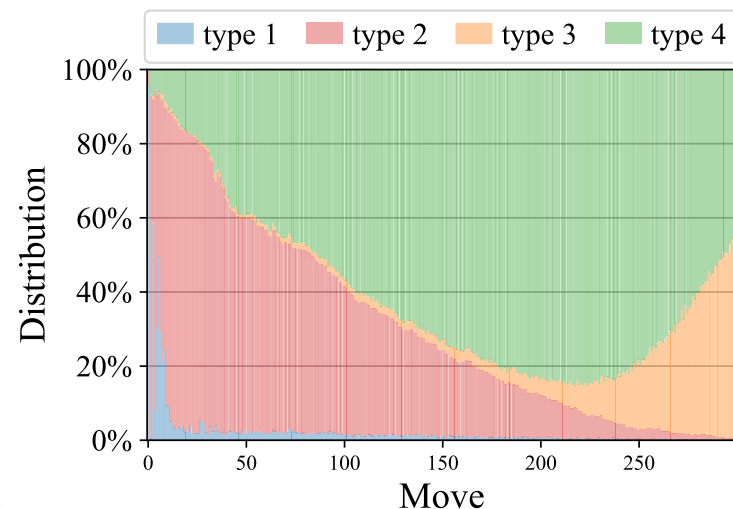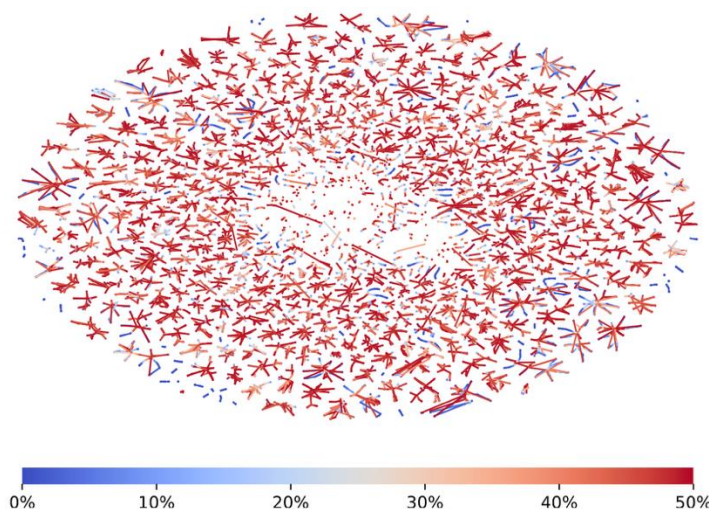15   **return** $\mathcal{G}$

# Evaluation

1. End-to-end training: improves sample efficiency up to 58.5% in 9x9 Go and 47.3% in 19x19 Go.
2. LCR uncertainty estimator closely matches the empirical results with an average RMSE of 0.02.
3. Our design space exploration shows that V=10 achieves the best balance for exploration within and between games.

and more...

# For more details... Please check our paper!

# Thank you

NEURAL INFORMATION PROCESSING SYSTEMS