

Game Solving with Online Fine-Tuning

Ti-Rong Wu,¹ Hung Guei,¹ Ting Han Wei,² Chung-Chin Shih,^{1,3} Jui-Te Chin,³ I-Chen Wu^{1,3}

¹ Academia Sinica
² University of Alberta
³ National Yang Ming Chiao Tung University



Game Solving with AlphaZero

- AlphaZero not only demonstrates super-human levels in game playing, but also **serves as heuristics in game solving**
- To solve a game, a winning response must be found for all possible moves by the losing player, which **includes very poor lines of play**

Self-Play

Out-Of-Distribution

⇒ For game solving, the fixed, pre-trained AlphaZero heuristics can be highly inaccurate



Game Solving with Online Fine-Tuning

- We investigate online fine-tuning to learn tailor-designed heuristics
- The online fine-tuning solver comprises three components:
 - Manager
 - Workers
 - Online Fine-Tuning Trainer



Online Fine-Tuning Trainer



Manager

- Maintain the proof search tree
 - Perform Monte Carlo tree search
- Employ a heuristic to assign jobs to workers to solve
 - Use *Proof Cost Network* to predict the cost for solving the position
 - Determine whether to assign to workers by a *cost threshold* v_{thr}
- Forward **solved/critical positions** to online fine-tuning trainer





中央研究院資訊科學研究所 Institute of Information Science, Academia Sinica

Workers

- Attempt to solve the jobs in parallel
 - Within given computing constraints
 - Employ the same heuristic as manager
- Return the solutions to manager





Online Fine-Tuning Trainer

- Fine-tune the heuristic using
 - Solved positions: where theoretic outcomes are found
 - ⇒ Guide the model to learn their theoretic outcomes
 - **Critical positions**: where manager is currently focused on
 - ⇒ Use them as initial positions to perform self-play
- Update the heuristic employed by manager and workers



Experiments

• We evaluate online fine-tuning solvers on 16 challenging 7x7 Killall-Go opening problems

JA

- SP: w/ solved positions
- CP: w/ critical positions
- In general, online fine-tuning solvers significantly reduce the solving time by using only 23.54% of computation time



Behavior Analysis

- We observe several positions in which the winning moves differed between the two solvers
- For a crucial sub-position of JA, the baseline and the online finetuning solver chose moves A and B to search, respectively
 - Both solvers have a higher chance to explore A at the beginning
 - The online fine-tuning solver **quickly realizes that solving B is faster**



	# Nodes	% Nodes on A	% Nodes on B
BASELINE	1,628,054,640	91.00%	2.18%
ONLINE-CP	136,470,552	16.21%	69.74%

Summary

- Pre-trained AlphaZero-based models provide less accurate heuristics
 ⇒ Not optimal for solving problems
- Online fine-tuning solvers learn tailor-designed heuristics
 - Dynamically during solving
 - According to the manager's attention
 - ⇒ Find faster solutions
- We expect it has the potential to extend to
 - Single-player games such as Rubik's Cube
 - Even other non-game fields



Thank You for Your Attention

Our code and data are available at https://rlg.iis.sinica.edu.tw/papers/neurips2023-online-fine-tuning-solver

