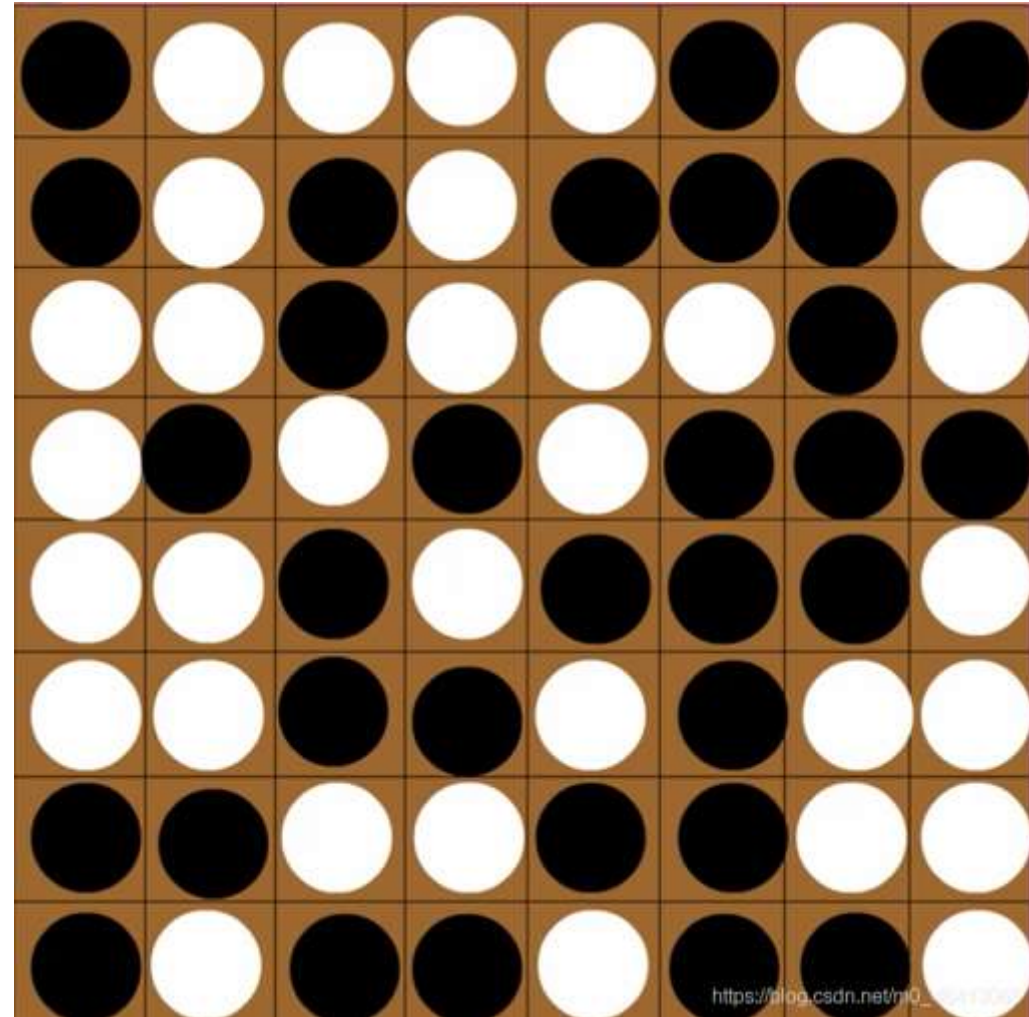# Artificial Intelligence Project

鞠屹昂 杨在洲 李盛忻 黄冠杰 刁斫

# Topic: Reversi(黑白棋) Agent

Rules:

1. Start from four pieces
2. when placing a new piece, opponent's pieces in the middle will be reversed

3. piece must be placed at locations that can reverse opponent's pieces.

4. Winner is the player with most pieces

# Motivation

Reason why we choose this topic:

Latest news:

Reversi(or Othello) was finally solved by Hiroki Takizawa at 2023.10 using the latest computer cluster and improved algorithm.

The rules are explicit and the state space is easy to represent

**Hiroki Takizawa**
Preferred Networks, Inc.
Chiyoda-ku, Tokyo, Japan
contact@hiroki-takizawa.name

### ABSTRACT

The game of Othello is one of the world's most complex and popular games that has yet to be computationally solved. Othello has roughly ten octodecillion (10 to the 58th power) possible game records and ten octillion (10 to the 28th power) possible game positions. The challenge of solving Othello, determining the outcome of a game with no mistake made by either player, has long been a grand challenge in computer science. This paper announces a significant milestone: Othello is now solved. It is computationally proved that perfect play by both players lead to a draw. Strong Othello software has long been built using heuristically designed search techniques. Solving a game provides a solution that enables the software to play the game perfectly.

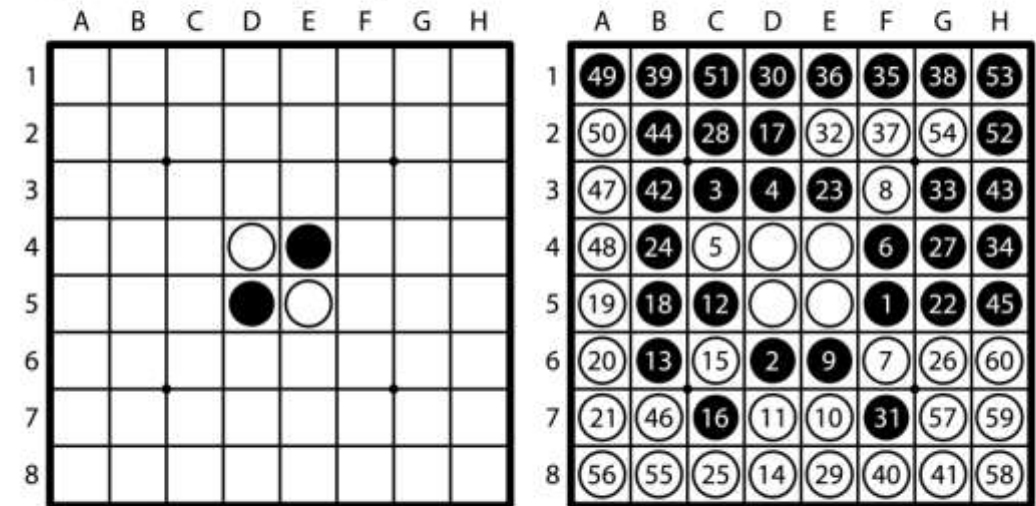**Keywords** Othello · Reversi · Games · alpha-beta search



Figure 1: (Left) The initial board position of 8 × 8 Othello. (Right) A diagram of an optimal game record designated by our study. The game record is "F5D6C3D3 C4F4F6F3 E6E7D7C5 B6D8C6C7 D2B5A5A6 A7G5E3B4 C8G6G4C2 E8D1F7E2 G3H4F1E1 F2G1B1F8 G8B3H3B2 H5B7A3A4 A1A2C1H2 H1G2B8A8 G7H8H7H6". The numbers in the stones indicate the order of moves, and the colors of stones indicate the final result. Our study confirms that if a deviation from this record occurs at any point, our software, playing as the opponent, is guaranteed a draw or win.

# Possible Methods

We plan to implement these methods:

1. Greedy: Simply pursue most reversed pieces at each step
2. Monte Carlo Tree Search
3. Reinforcement Learning based on Q-values

# Monte Carlo Tree Search

Learn from some existing frameworks of Reversi, figure out how to implement MCTS tree search, design an interface for testing.

Roxanne Policy: Combination of some useful policies, could be helpful in guiding our implementation.

# Reinforcement Learning based on Q-values

Considerations:

$2^{28}$ possible states, we cannot store them all, so we fit the Q-value function with a network instead.



Figure 2. Topologies of function approximators. A TD-network (a) tries to approximate the value of the state presented at the input. A Q-learning network (b) tries to approximate the values of all the possible actions in the state presented at the input.

# Evaluation

After implementing these methods,

1. We take the implemented methods to fight Game AI and calculate their win rate

2. We let the three agents fight with each other in order to find an optimal one