

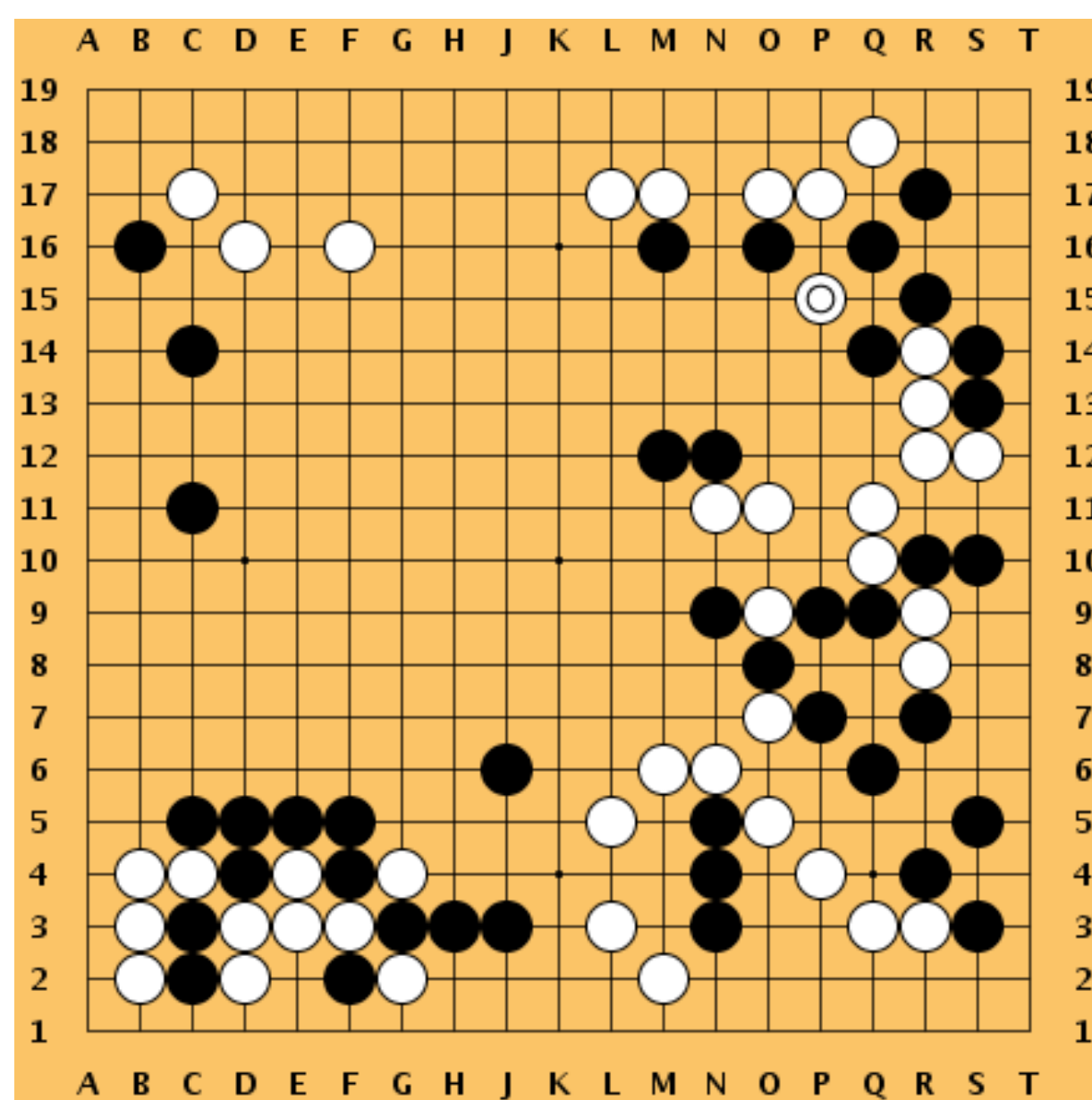


LEARN TO PLAY GO

ALBERT LIU (ALBERTPL@STANFORD.EDU)



MOTIVATION & CHALLENGES



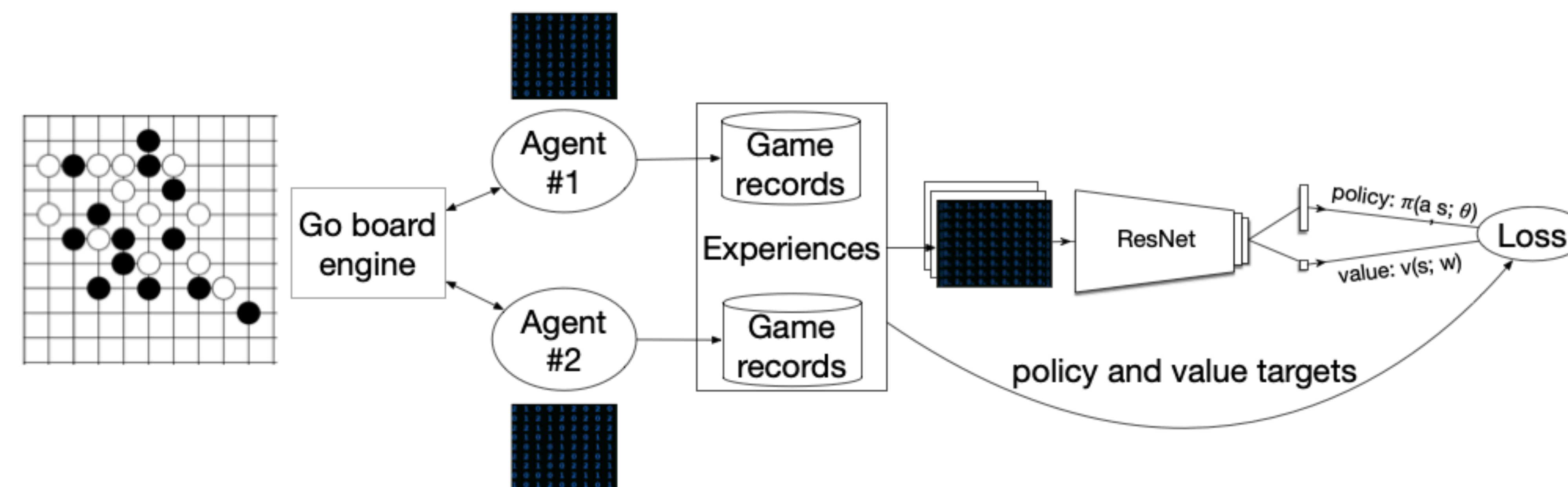
- simple rule but mastering requires many years of study by human
- state based search is intractable due enormous search space ($\sim 10^{170}$), large number of legal move per state (~ 250)
- the difficulty to handcraft a heuristic evaluation function
- AlphaGo and AlphaGoZero have rocked the Go and AI world

PROBLEM DEFINITION

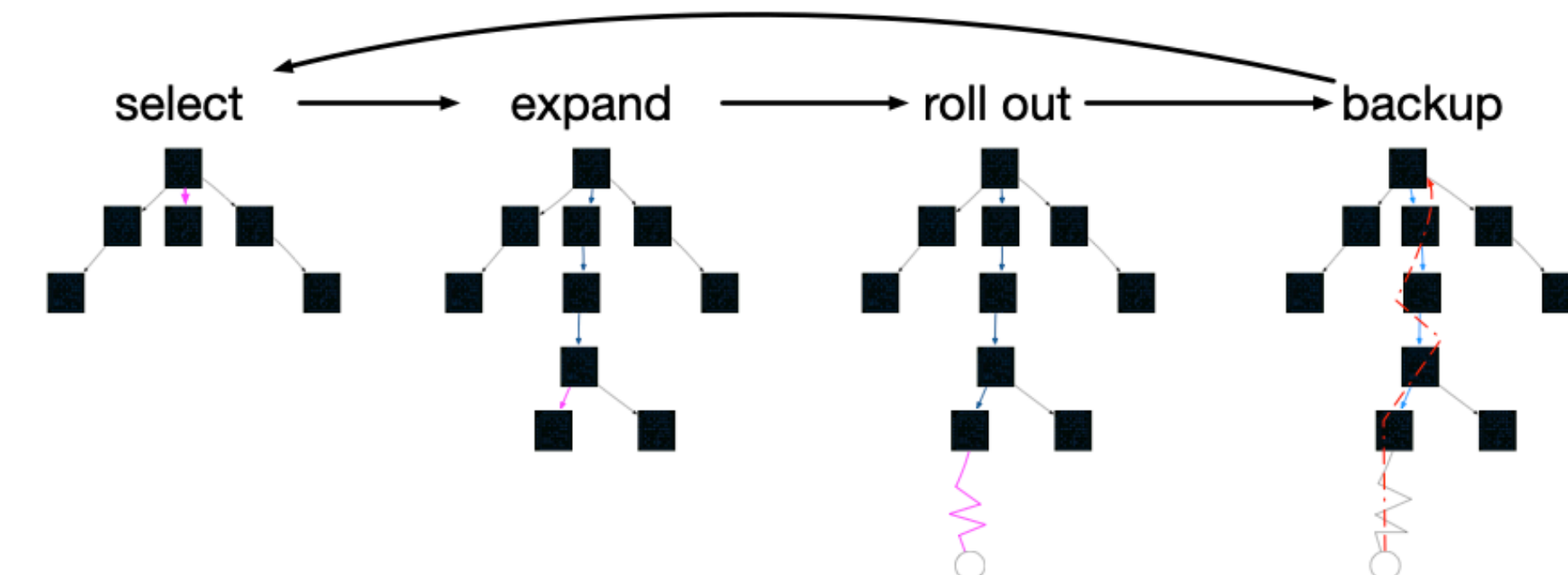
Turn-taking, two player, fully observable zero-sum game.

- The state space is all possible placements of the stones on the board and the player of that turn.
- The action space is any legal position of a stone and a pass action.
- The goal is to maximize territory, by observing the board only.

SYSTEM ARCHITECTURE



MONTE CARLO TREE SEARCH



Input: root node s_0 : current game state

Input: $c > 0$: parameter to control the degree of exploration

```

1 while within computation bound do
2    $s \leftarrow s_0$ 
3    $\Delta \leftarrow \emptyset$ 
4   // selection based on tree policy: Upper Bound Confidence
5   while  $s$  is in the tree do
6      $a \leftarrow \arg \max \{ q(s, a) + c \sqrt{\frac{\log \sum_{a'} n(s, a')}{n(s, a)}} \}$ 
7      $\Delta \leftarrow \Delta \cup \{(s, a)\}$ 
8      $s \leftarrow \text{Succ}(s, a)$ 
9   expand tree with the new node  $s$ 
10  continue the game from  $s$  with random policy and let  $r$  be the reward
11  // update each node on the path
12   $s_0 \rightarrow \dots \rightarrow s$ 
13  for  $s, a \in \Delta$  do
14     $n(s, a) \leftarrow n(s, a) + 1$ 
15     $q(s, a) \leftarrow q(s, a) + \frac{1}{n(s, a)}(r - q(s, a))$ 
16 return  $\arg \max_a q(s_0, a)$ 

```

REINFORCEMENT LEARNING

PRELIMINARY RESULTS

ANALYSIS