Google DeepMind
Challenge Match

ALPHAGO
00:45:23

LEE SEDOL
00:43:04

QUT

$(p, v) = f_\theta(s)$ and $l = (z - v)^2 - \pi^T \log p + c\|\theta\|^2$

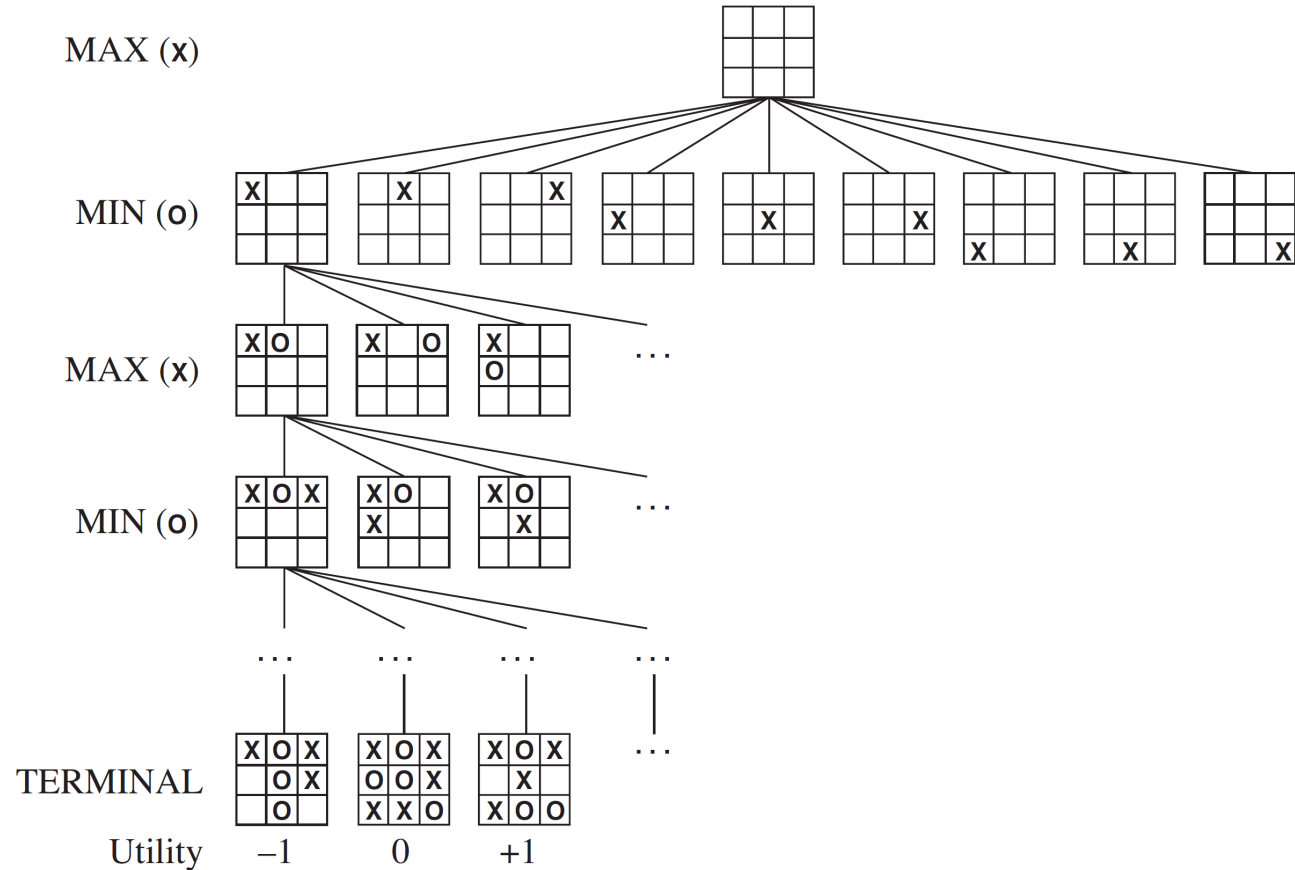# AlphaZero

# Outline

- Game trees

- MinMax

- TD leaf

- Monte Carlo Tree Search

- AlphaZero

# Learning to Optimize Rewards

- In Reinforcement Learning, a software agent makes **observations** and takes **actions** within an **environment**, and in return it receives rewards.

- Its objective is to learn to act in a way that will **maximize** its **expected rewards** over time.

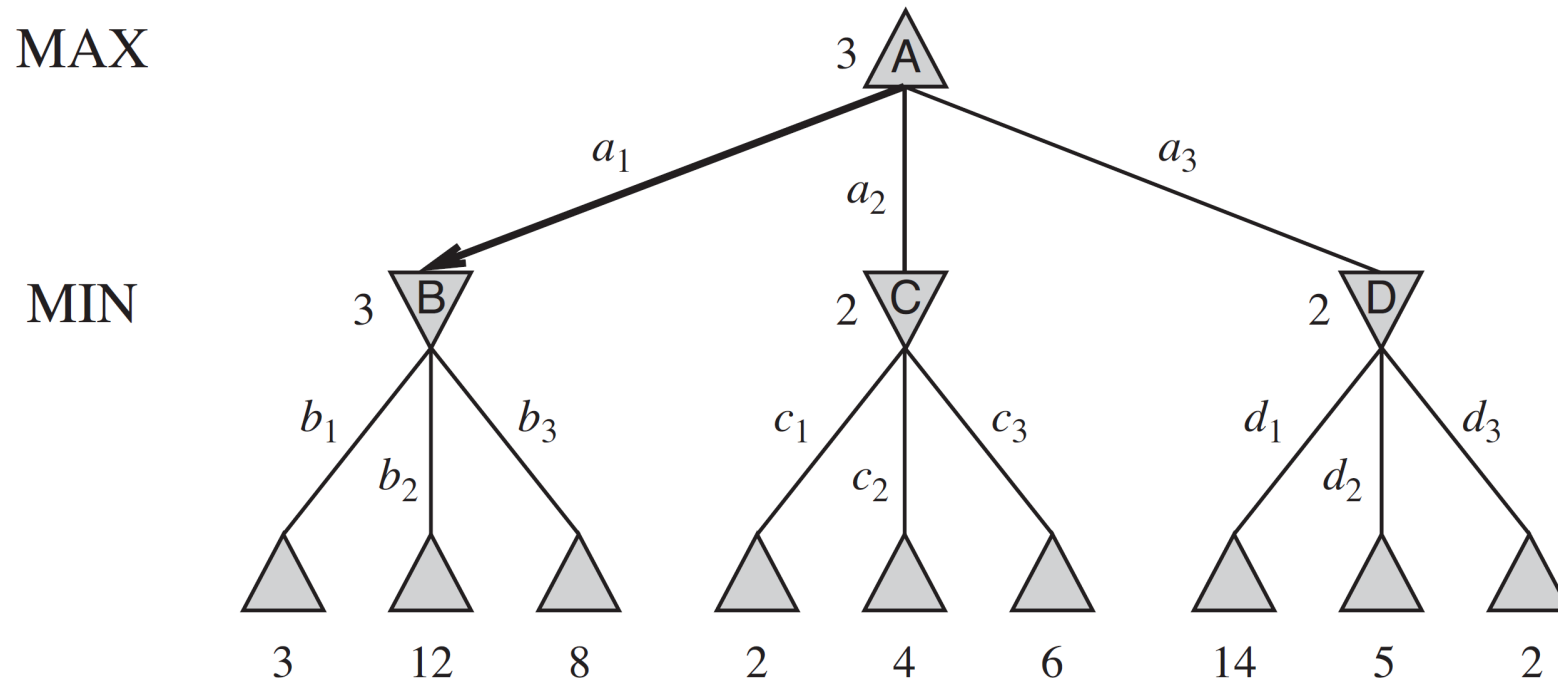Refresher from last lecture

# A (partial) game tree for Tic-Tac-Toe

# Minimax Algorithm

- How to find the optimal action for a given possible game state?
  - Minimax value of a state: the value of the best outcome possible for the player who needs to move at the state assuming both players will take the best action in the corresponding states after the current move
  - Minimax Algorithm
    - Labeling the minimax value of each state given the minimax values of its successors
    - For player , take the max value of the values of its successors
    - For player , take the min value of the values of its successors
    - Essentially a backward induction algorithm (von Neumann & Morgenstern 1944)
    - Often implemented in a depth-first manner (recursive algorithm)

# Minimax Algorithm

# Minimax Algorithm

$$\text{M\scriptsize INIMAX}(s) =$$

$$\begin{cases} \text{U\scriptsize TILITY}(s) & \text{if T\scriptsize ERMINAL}\text{-T\scriptsize EST}(s) \\ \max_{a \in Actions(s)} \text{M\scriptsize INIMAX}(\text{R\scriptsize ESULT}(s, a)) & \text{if P\scriptsize LAYER}(s) = \text{M\scriptsize AX} \\ \min_{a \in Actions(s)} \text{M\scriptsize INIMAX}(\text{R\scriptsize ESULT}(s, a)) & \text{if P\scriptsize LAYER}(s) = \text{M\scriptsize IN} \end{cases}$$

# Minimax Algorithm

**function** MINIMAX-DECISION(*state*) **returns** *an action*
 **return** $\arg\max_{a \,\in\, \text{ACTIONS}(s)}$ MIN-VALUE(RESULT(*state*, *a*))

---

**function** MAX-VALUE(*state*) **returns** *a utility value*
 **if** TERMINAL-TEST(*state*) **then return** UTILITY(*state*)
 $v \leftarrow -\infty$
 **for each** *a* **in** ACTIONS(*state*) **do**
  $v \leftarrow$ MAX($v$, MIN-VALUE(RESULT(*s*, *a*)))
 **return** $v$

---

**function** MIN-VALUE(*state*) **returns** *a utility value*
 **if** TERMINAL-TEST(*state*) **then return** UTILITY(*state*)
 $v \leftarrow \infty$
 **for each** *a* **in** ACTIONS(*state*) **do**
  $v \leftarrow$ MIN($v$, MAX-VALUE(RESULT(*s*, *a*)))
 **return** $v$

# Minimax Algorithm

- Intractable for large games
  Chess( nodes), Go ( nodes)
  Cannot even represent the optimal strategy profile in space

- For many problems, you do not need a full contingent plan at the very beginning of the game
  Can solve the problem in a more "online" fashion, just like how human players play Chess/Go: My opponent takes this action and what should I do now?

# Minimax Algorithm

- If we only care about the game value, or the optimal action at a specific game state, can we do better?

    Trick: Depth-limited search (Limit the depth)
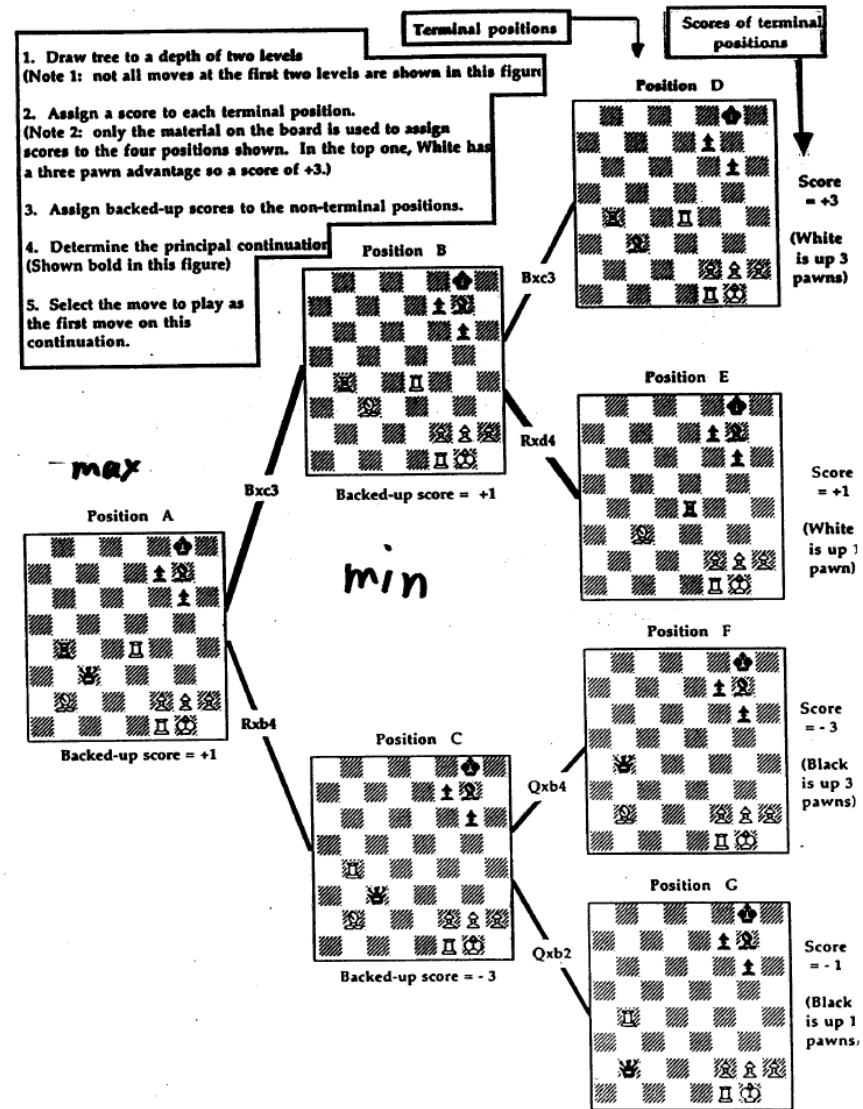
- Minimax algorithm (Shannon, 1950):

    Set  (an integer) as an upper bound on the search depth

    the static evaluation function, returns an estimate of minimax value of state
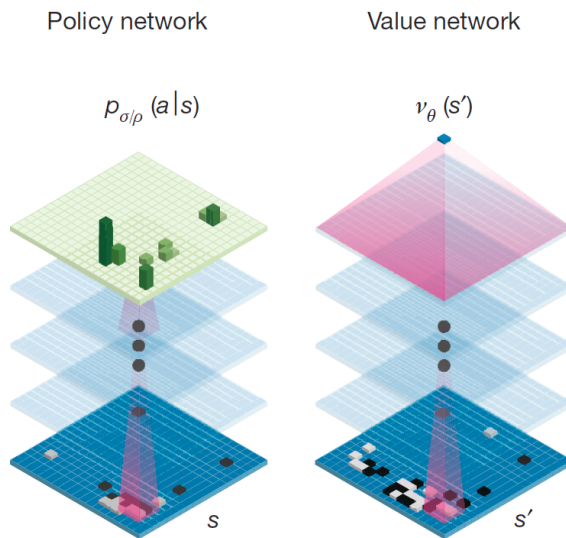
# Minimax Algorithm

Claude Shannon,
Alan Turing:
Minimax search with
scoring function
(1950)

Only show a few
branches here

# Connections with RL

- Learning / Estimating the **state value function** will help the MAX player find the optimal policy against a perfectly rational MIN player



Policy network      Value network

$p_{\sigma/\rho}\,(a|s)$      $v_\theta\,(s')$

$s$      $s'$

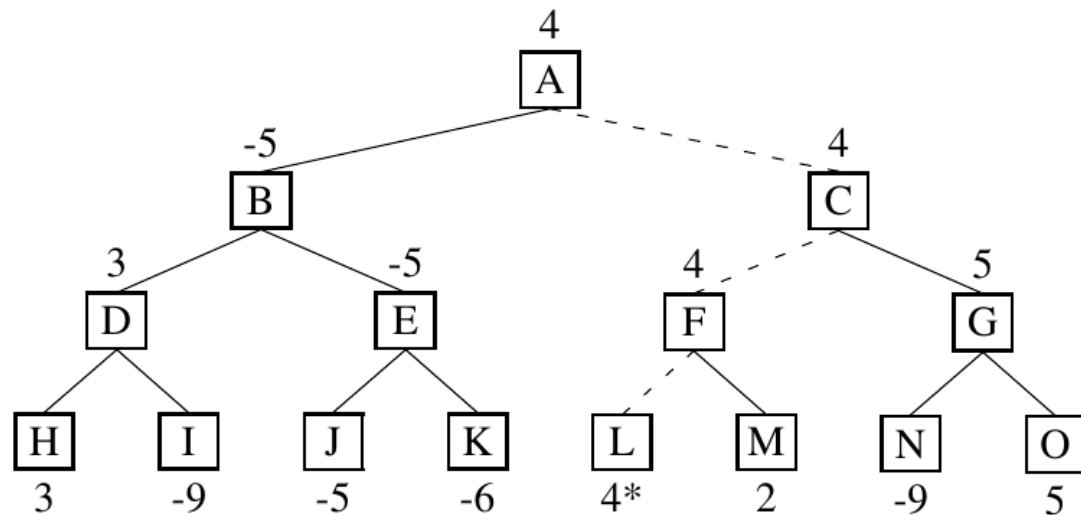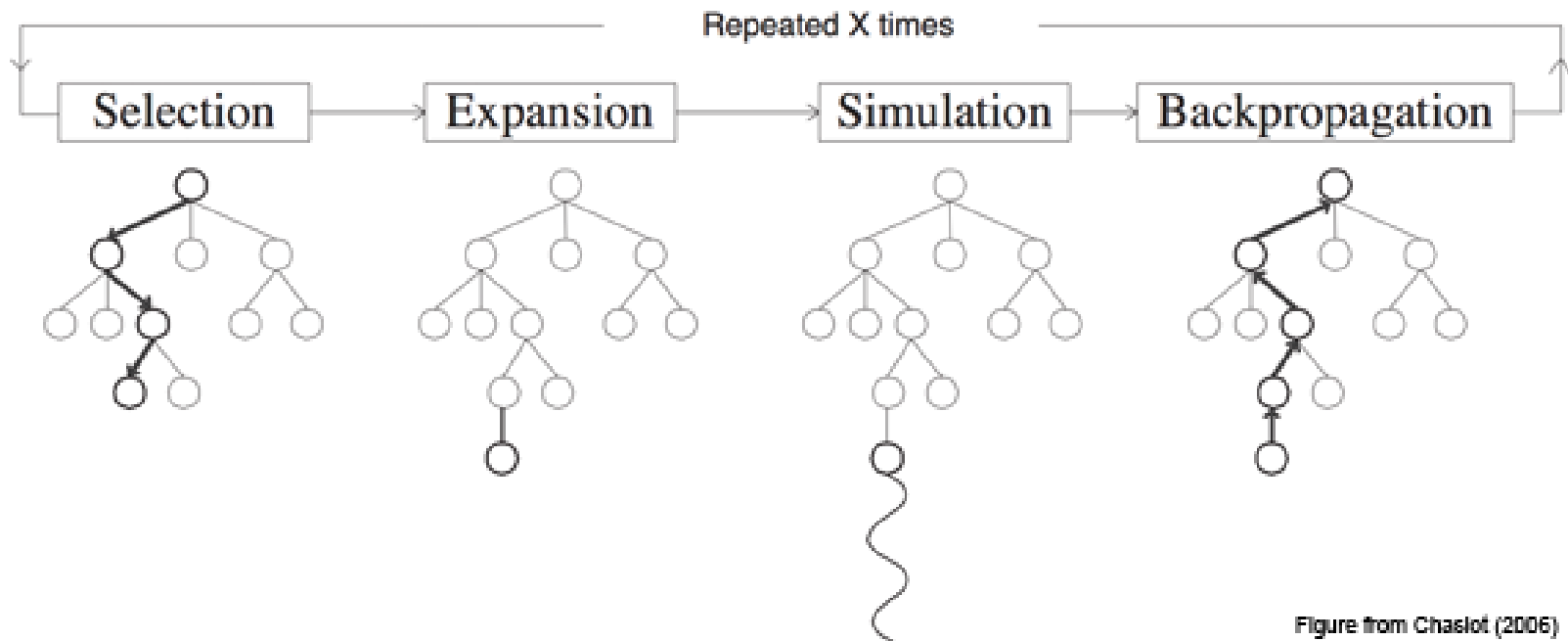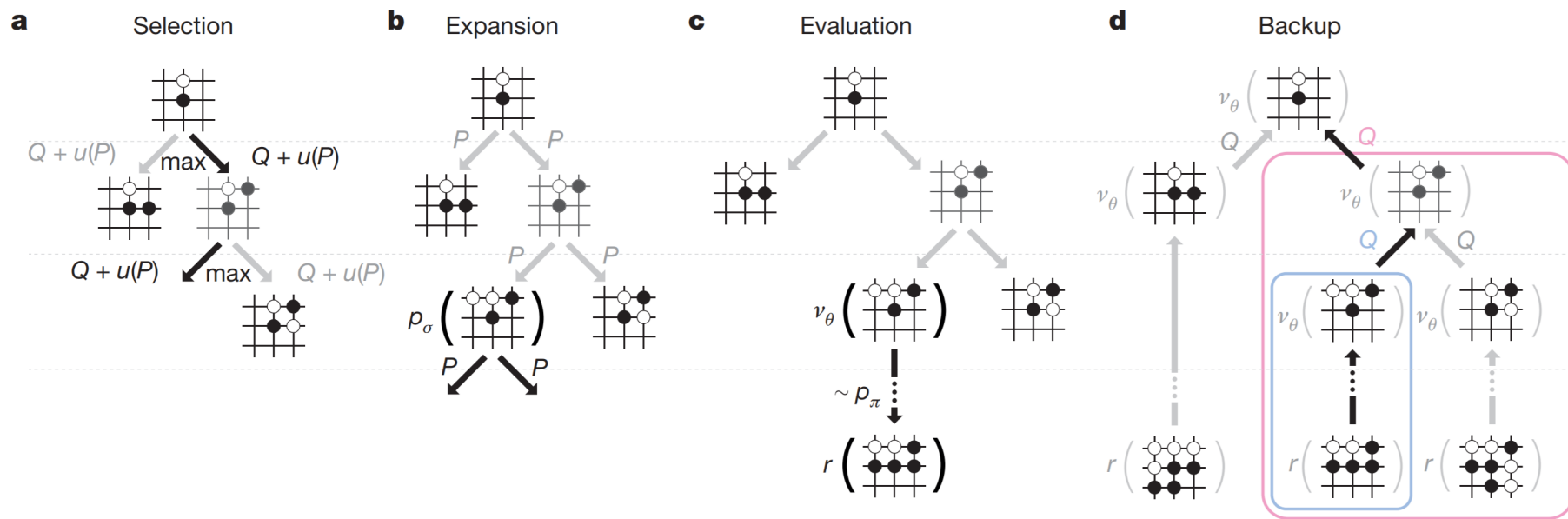# TDLeaf(λ): Combining Temporal Difference Learning with Game-Tree Search



Fig. 1: Full breadth, 3-ply search tree illustrating the minimax rule for propagating values. Each of the leaf nodes (H–O) is given a score by the evaluation function, $\tilde{J}(\cdot, w)$. These scores are then propagated back up the tree by assigning to each opponent's internal node the minimum of its children's values, and to each of our internal nodes the maximum of its children's values. The principle variation is then the sequence of best moves for either side starting from the root node, and this is illustrated by a dashed line in the figure. Note that the score at the root node A is the evaluation of the leaf node (L) of the principal variation. As there are no ties between any siblings, the derivative of A's score with respect to the parameters $w$ is just $\nabla \tilde{J}(L, w)$.

# MCTS tree growth



Repeated X times

Selection → Expansion → Simulation → Backpropagation

Figure from Chaslot (2006)

# How AlphaGo works

- Supervised learning + Reinforcement learning
- Monte-Carlo Tree Search



upper confidence bound $Q(s, a) + U(s, a)$, where $U(s, a) \propto P(s, a) / (1 + N(s, a))$