
Games and Adversarial Search

Artificial Intelligence

Slides are mostly adapted from AIMA, MIT Open Courseware and
Svetlana Lazebnik (UIUC)



World Champion chess player Garry Kasparov is defeated by IBM's Deep Blue chess-playing computer in a six-game match in May, 1997

[\(link\)](#)



© Telegraph Group Unlimited 1997

Why study games?

- Games are a traditional hallmark of intelligence
 - Games are easy to formalize
 - Games can be a good model of real-world competitive or cooperative activities
 - Military confrontations, negotiation, auctions, etc.
-

Games – history of chess playing

- 1949 – Shannon paper – originated the ideas
 - 1951 – Turing paper – hand simulation
 - 1958 – Bernstein program
 - 1955-1960 – Simon-Newell program
 - 1961 – Soviet program
 - 1966 – 1967 – MacHack 6 – defeated a good player
 - 1970s – NW chess 4.5
 - 1980s – Cray Bitz
 - 1990s – Belle, Hitech, Deep Thought,
 - 1997 - Deep Blue - defeated Garry Kasparov
-

Games

- Multi agent environments : any given agent will need to consider the actions of other agents and how they affect its own welfare.
 - The unpredictability of these other agents can introduce many possible contingencies
 - There could be competitive or cooperative environments
 - Competitive environments, in which the agent's goals are in conflict require adversarial search – these problems are called as games
-

Games vs. single-agent search

- We don't know how the opponent will act
 - The solution is not a fixed sequence of actions from start state to goal state, but a *strategy* or *policy* (a mapping from state to best move in that state)
 - Efficiency is critical to playing well
 - The time to make a move is limited
 - The branching factor, search depth, and number of terminal configurations are huge
 - In chess, *branching factor* ≈ 35 and *depth* ≈ 100 , giving a search tree of 10^{154} nodes
 - Number of atoms in the observable universe $\approx 10^{80}$
 - This rules out searching all the way to the end of the game
-

Types of game environments

	Deterministic	Stochastic
Perfect information (fully observable)	Chess, checkers, go	Backgammon, monopoly
Imperfect information (partially observable)	Battleships	Scrabble, poker, bridge

Games

- In game theory (economics), any multiagent environment (either cooperative or competitive) is a game provided that the impact of each agent on the other is significant
 - AI games are a specialized kind - deterministic, turn taking, two-player, zero sum games of perfect information
 - In our terminology – deterministic, fully observable environments with two agents whose actions alternate and the utility values at the end of the game are always equal and opposite (+1 and -1)
-

Alternating two-player zero-sum games

- Players take turns
- Each game outcome or **terminal state** has a **utility** for each player (e.g., 1 for win, -1 for loss, 0 for draw)
- The sum of both players' utilities is a constant



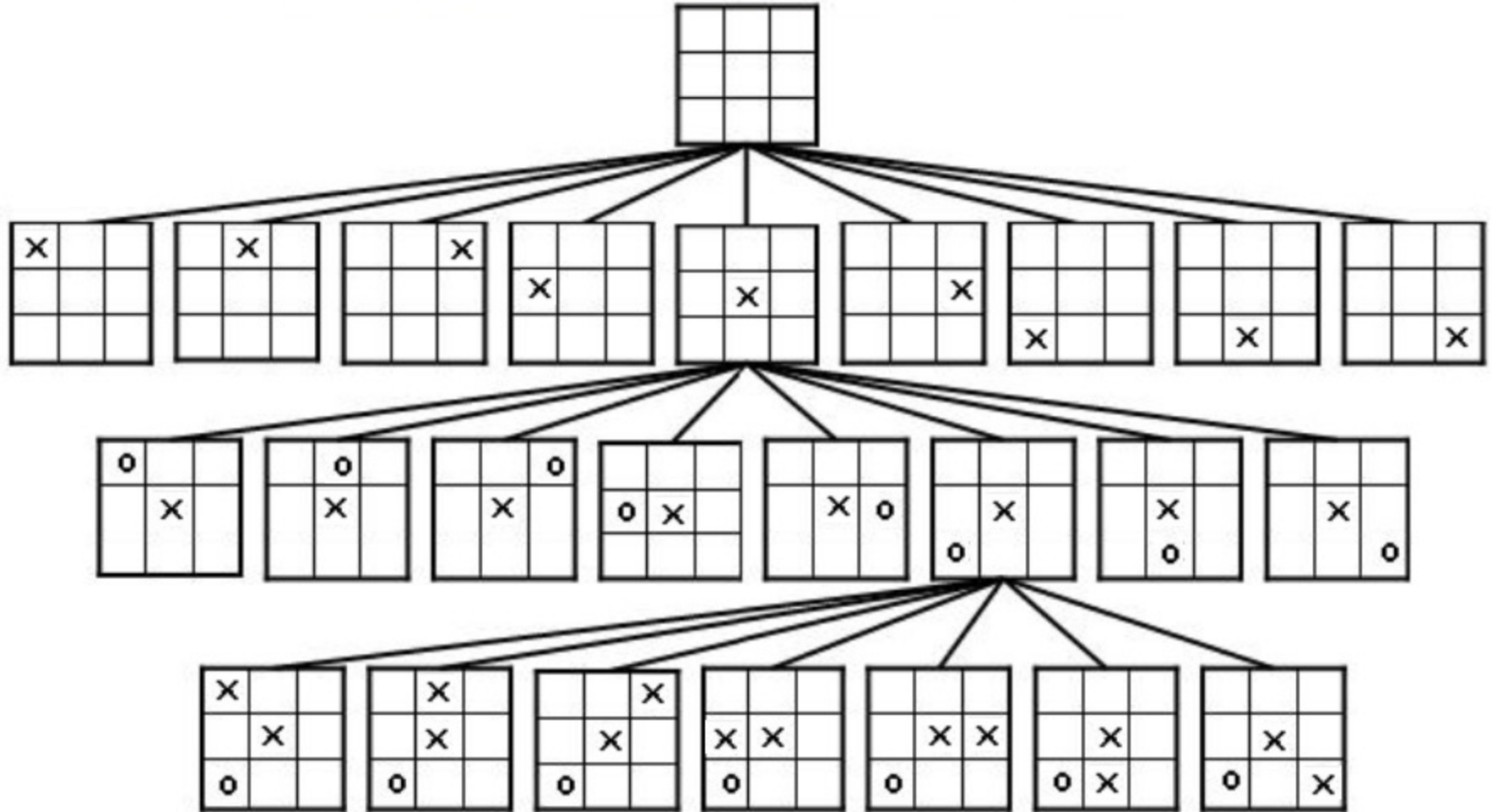
Game Tree search

- **Initial state: initial board position and player**
 - **Operators: one for each legal move**
 - **Goal states: winning board positions**
 - **Scoring function: assigns numeric value to states**
 - **Game tree: encodes all possible games**
-
- **We are not looking for a path, only the next move to make (that hopefully leads to a winning position)**
 - **Our best move depends on what the other player does**
-

Optimal strategies

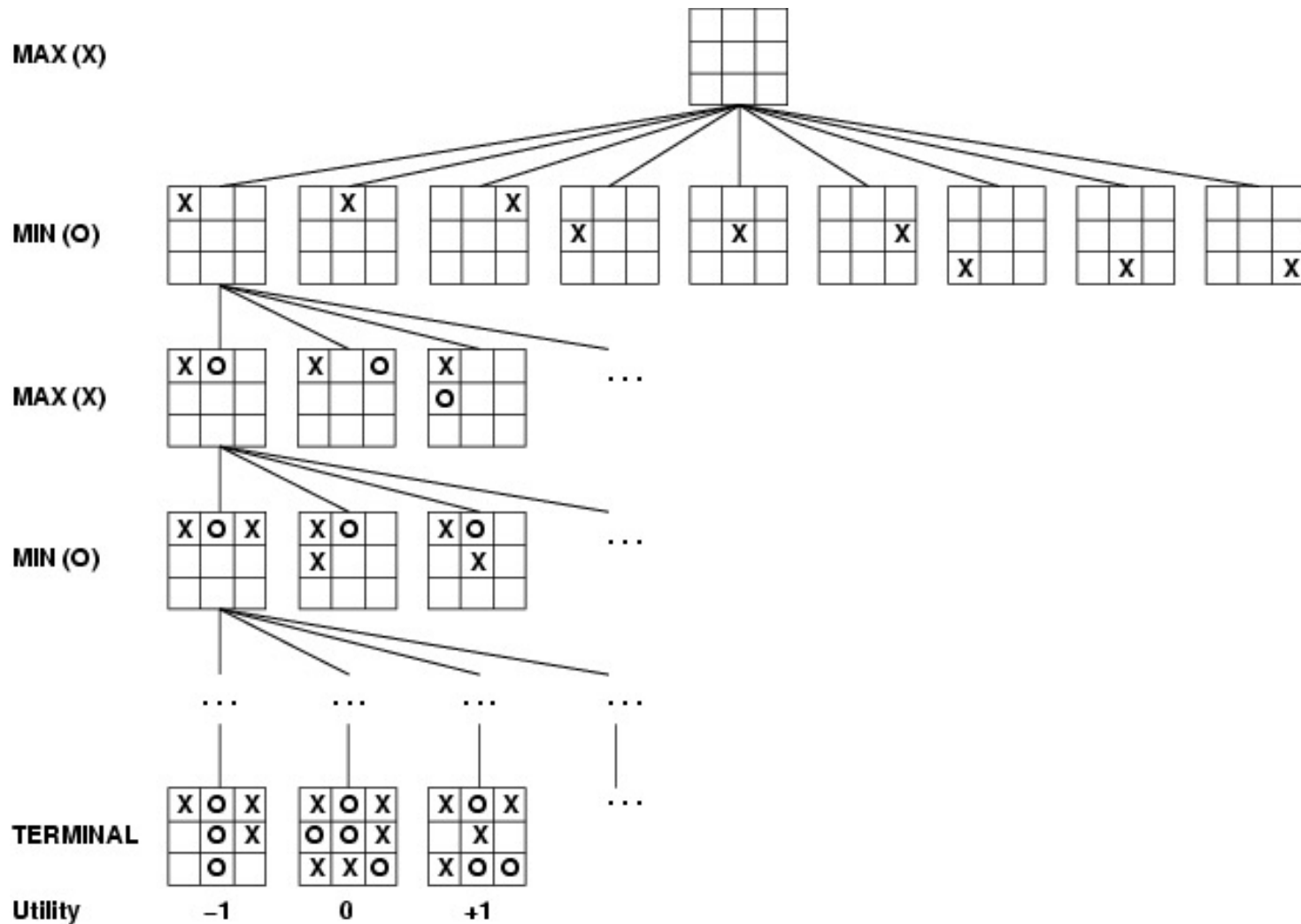
- In a normal search problem, the optimal solution would be a sequence of moves leading to a goal state - a terminal state that is a win
 - In a game, MIN has something to say about it and therefore MAX must find a contingent strategy, which specifies
 - MAX's move in the initial state,
 - then MAX's moves in the states resulting from every possible response by MIN,
 - then MAX's moves in the states resulting from every possible response by MIN to those moves
 - ...
 - An optimal strategy leads to outcomes at least as good as any other strategy when one is playing an infallible opponent
-

Partial Game Tree for Tic-Tac-Toe



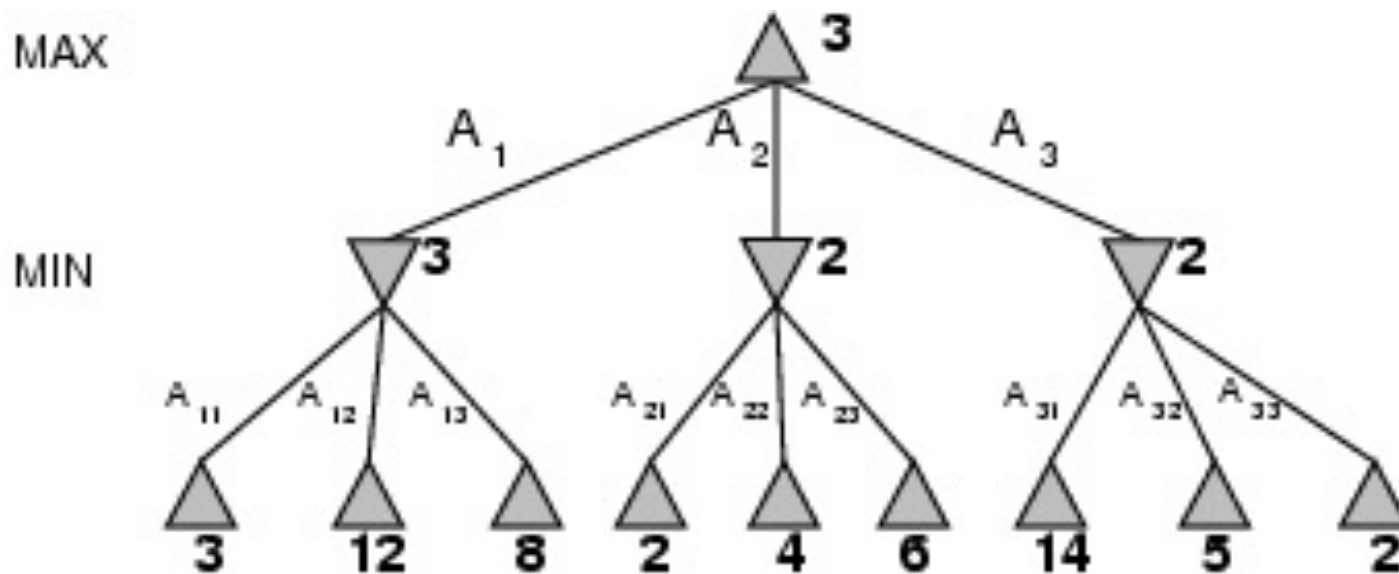
Game tree

- A game of tic-tac-toe between two players, “max” and “min”



Minimax

- Perfect play for deterministic games
- Idea: choose move to position with highest **minimax value**
= best achievable payoff against best play
- E.g., 2-ply game:



Minimax value

- Given a game tree, the optimal strategy can be determined by examining the minimax value of each node (MINIMAX-VALUE(n))
 - The minimax value of a node is the utility of being in the corresponding state, assuming that both players play optimally from there to the end of the game
 - Given a choice, MAX prefer to move to a state of maximum value, whereas MIN prefers a state of minimum value
-

Minimax algorithm

function MINIMAX-DECISION(*state*) *returns an action*

$v \leftarrow \text{MAX-VALUE}(state)$

return the *action* in SUCCESSORS(*state*) with value v

function MAX-VALUE(*state*) *returns a utility value*

if TERMINAL-TEST(*state*) **then return** UTILITY(*state*)

$v \leftarrow -\infty$

for a, s in SUCCESSORS(*state*) **do**

$v \leftarrow \text{MAX}(v, \text{MIN-VALUE}(s))$

return v

function MIN-VALUE(*state*) *returns a utility value*

if TERMINAL-TEST(*state*) **then return** UTILITY(*state*)

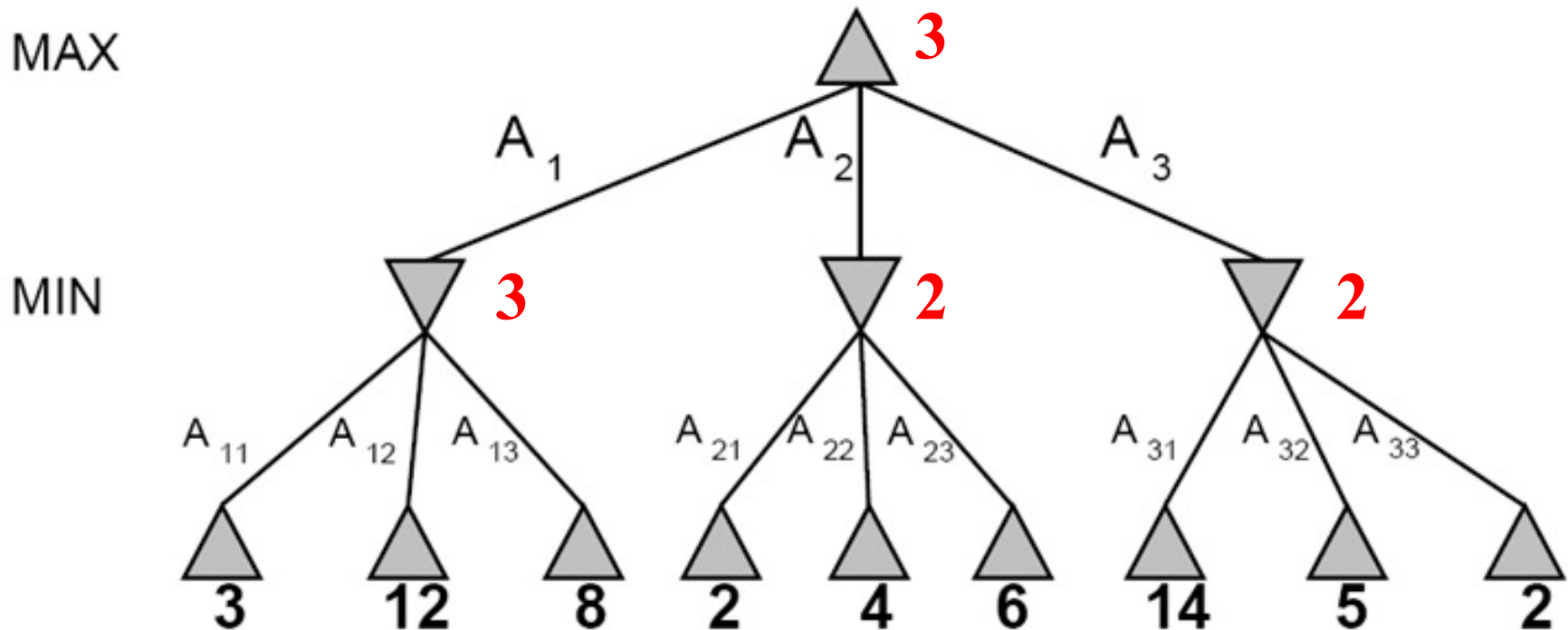
$v \leftarrow \infty$

for a, s in SUCCESSORS(*state*) **do**

$v \leftarrow \text{MIN}(v, \text{MAX-VALUE}(s))$

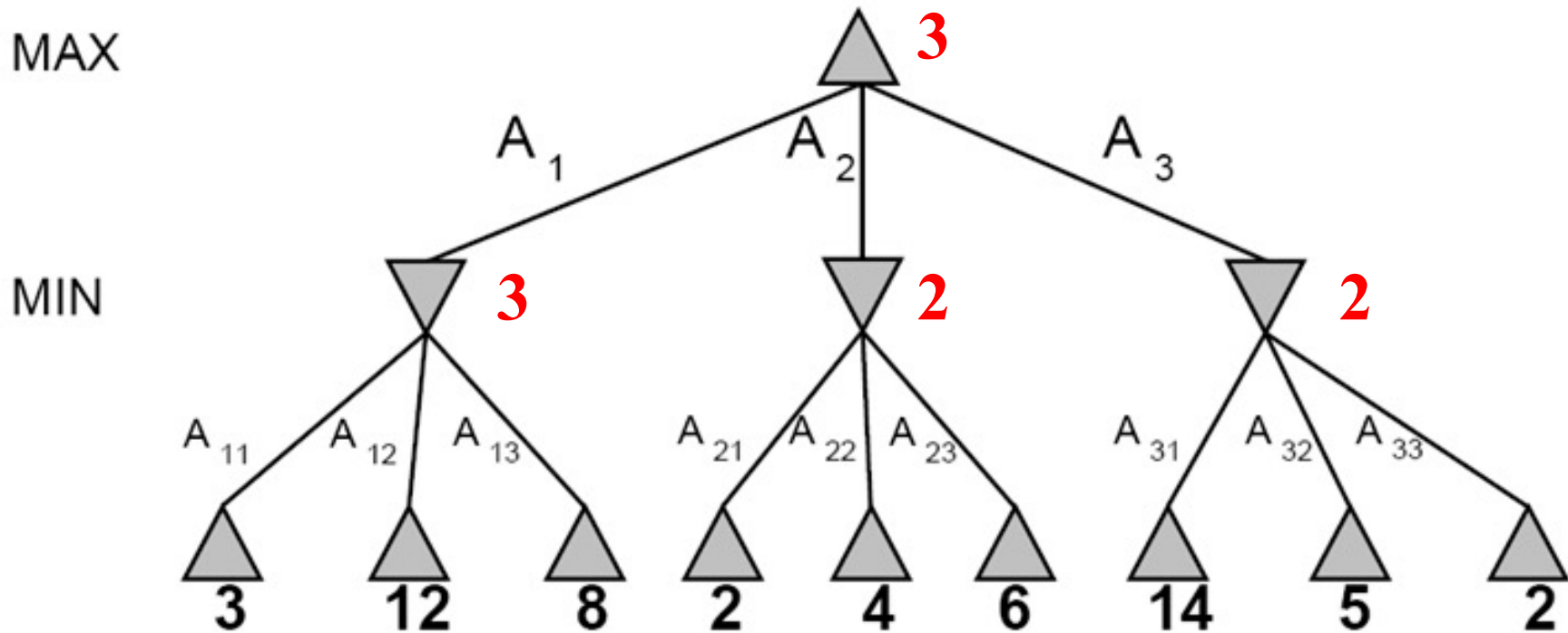
return v

Game tree search



- **Minimax value of a node:** the utility (for MAX) of being in the corresponding state, assuming perfect play on both sides
- **Minimax strategy:** Choose the move that gives the best worst-case payoff

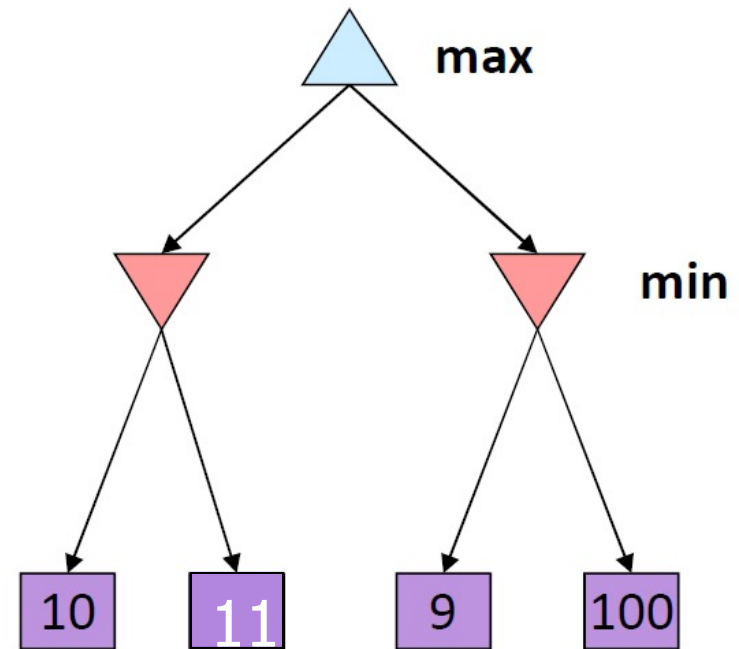
$$\begin{aligned} \text{MINIMAX-VALUE}(\text{root}) &= \max(\min(3,12,8), \min(2,4,6), \min(14,5,2)) \\ &= \max(3,2,2) \\ &= 3 \end{aligned}$$



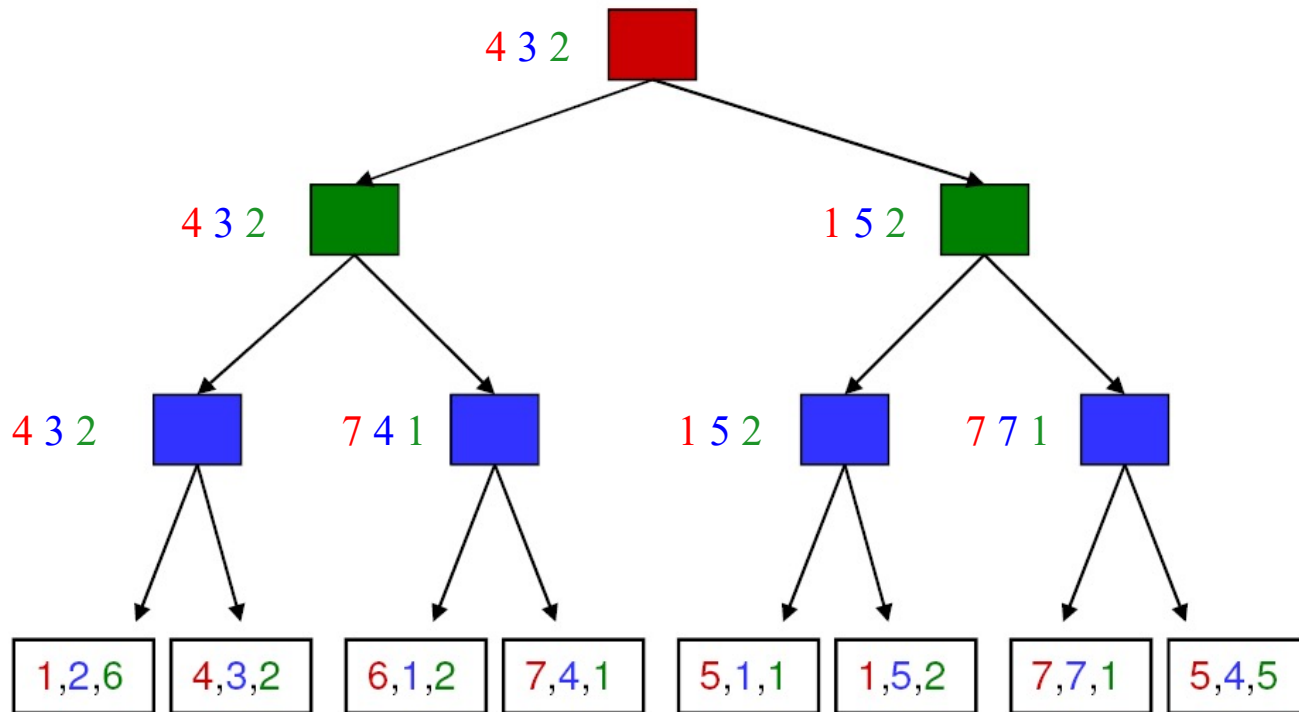
- **Minimax**(*node*) =
 - Utility(*node*) if *node* is terminal
 - $\max_{\text{action}} \text{Minimax}(\text{Succ}(\text{node}, \text{action}))$ if *player* = MAX
 - $\min_{\text{action}} \text{Minimax}(\text{Succ}(\text{node}, \text{action}))$ if *player* = MIN

Optimality of minimax

- The minimax strategy is optimal against an optimal opponent
- What if your opponent is suboptimal?
 - Your utility can only be higher than if you were playing an optimal opponent!
 - A different strategy may work better for a sub-optimal opponent, but it will necessarily be worse against an optimal opponent

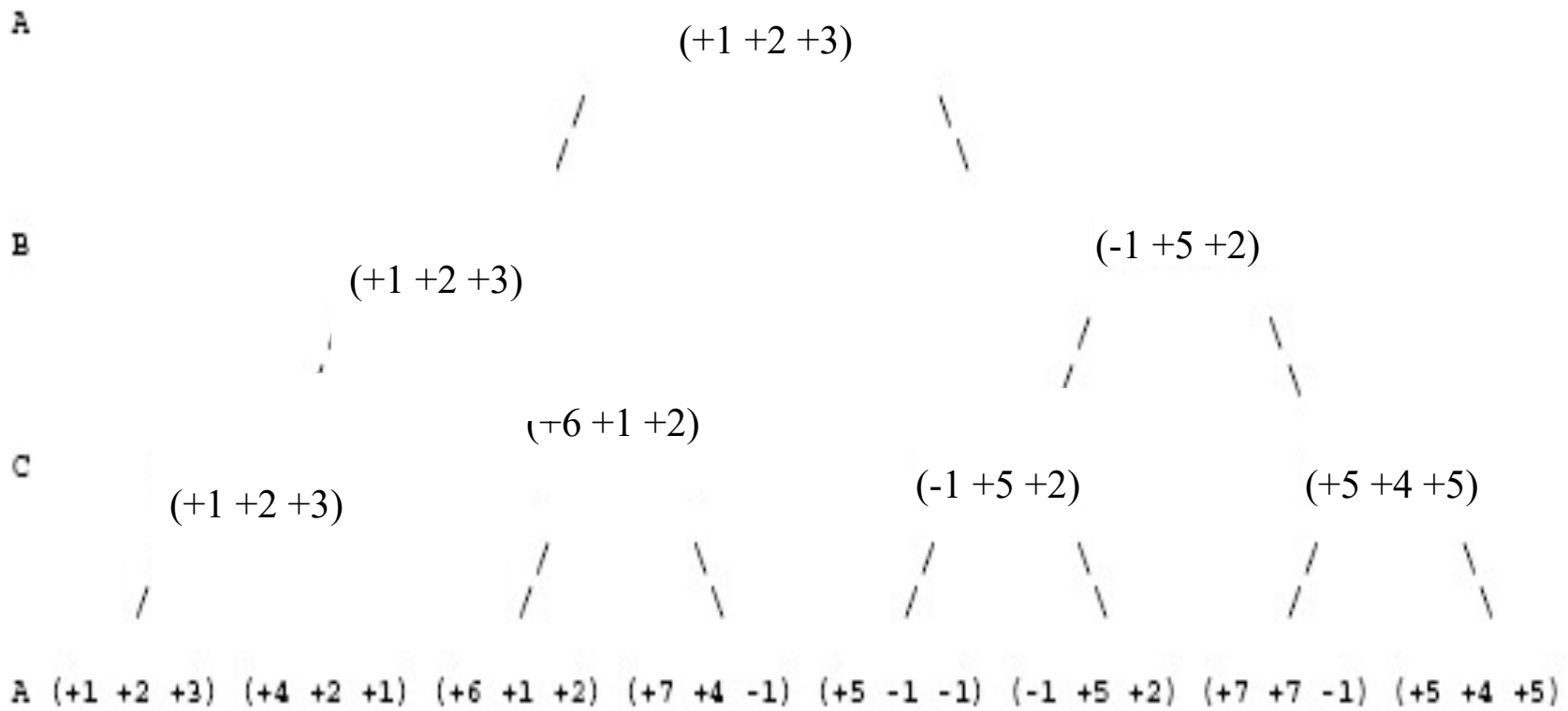


More general games



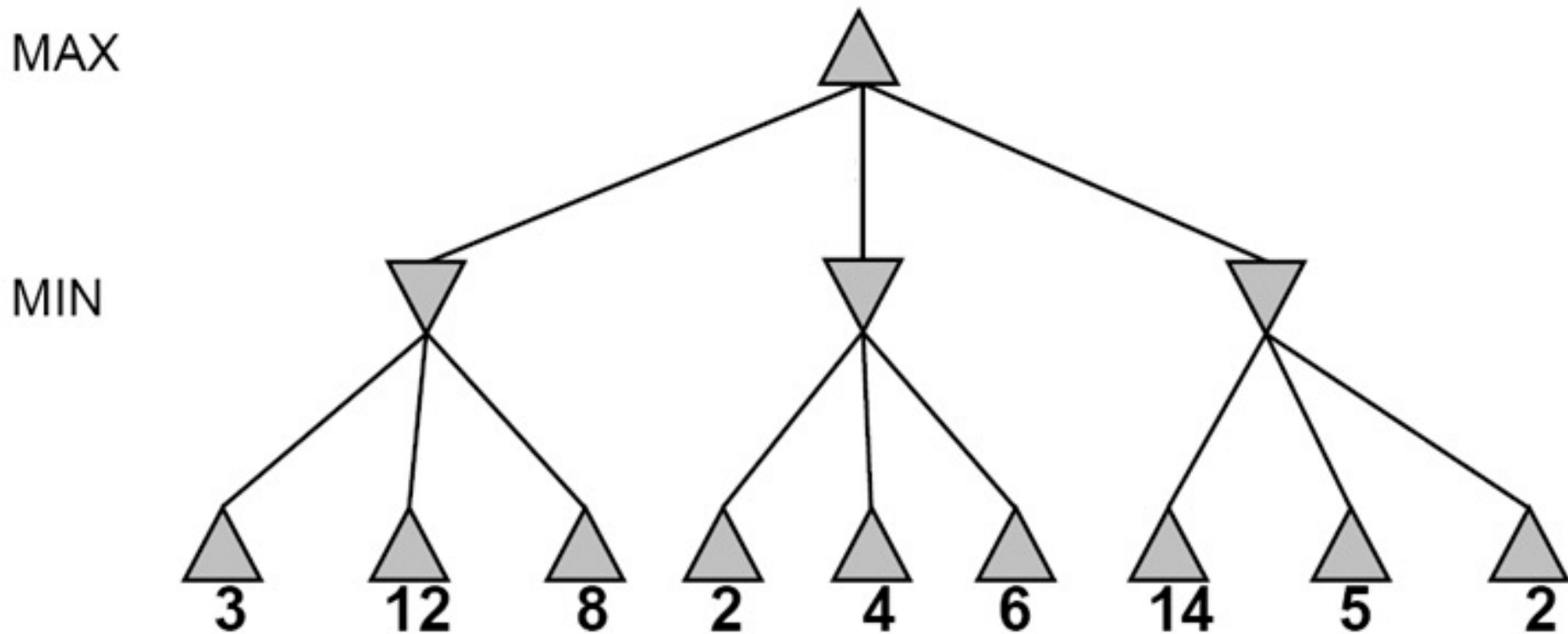
- More than two players, non-zero-sum
- Utilities are now tuples
- Each player maximizes their own utility at their node
- Utilities get propagated (*backed up*) from children to parents

Tree Player and Non-zero sum games



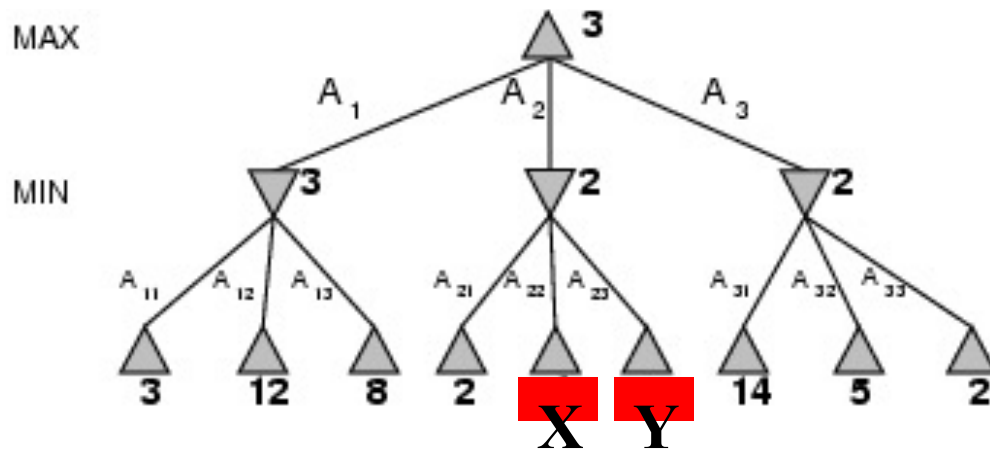
Alpha-beta pruning

- It is possible to compute the exact minimax decision without expanding every node in the game tree

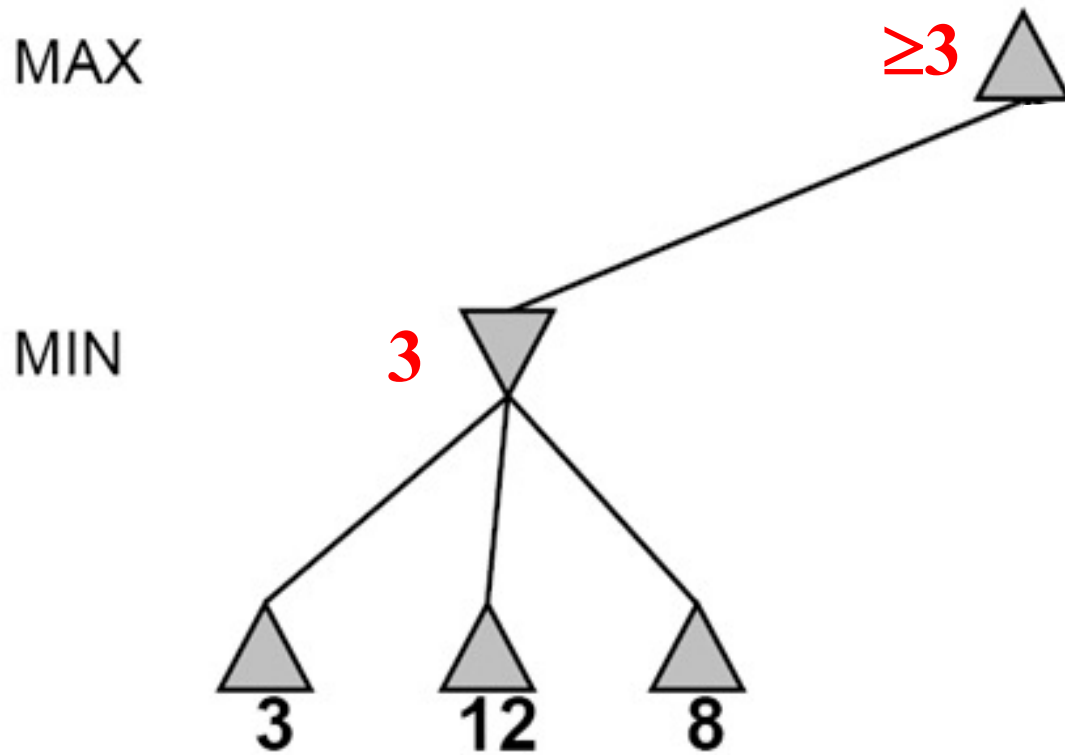


α - β pruning

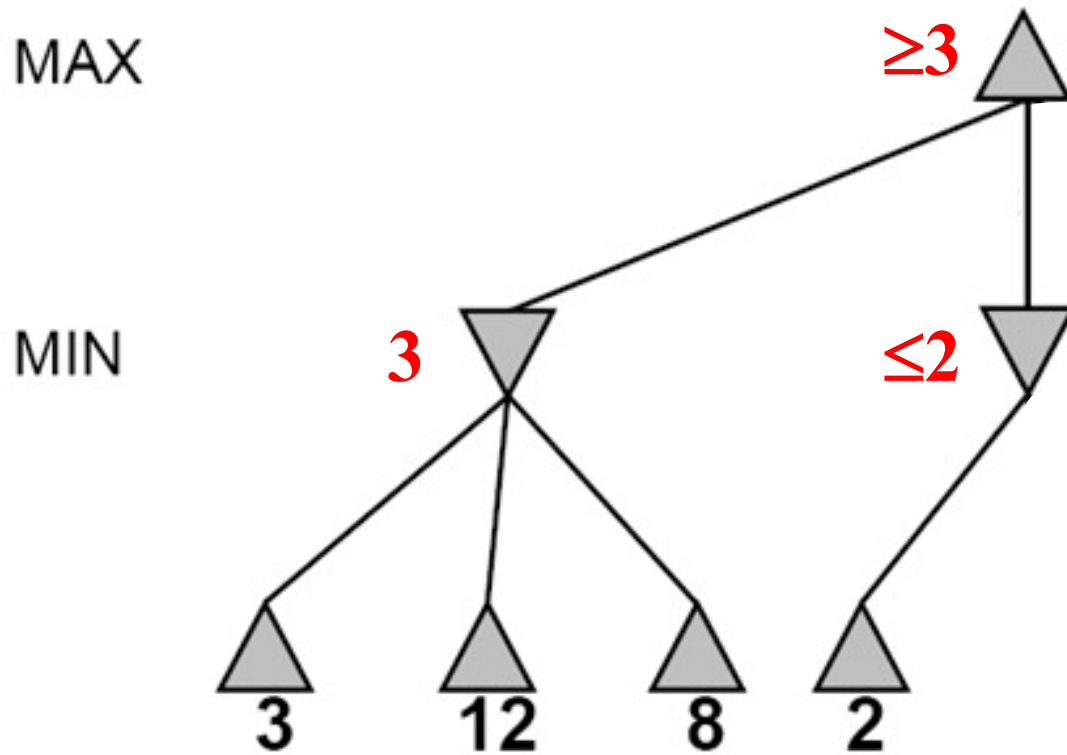
$$\begin{aligned}
 \text{MINIMAX-VALUE}(\text{root}) &= \max(\min(3,12,8), \min(2,x,y), \min(14,5,2)) \\
 &= \max(3, \min(2,x,y), 2) \\
 &= \max(3, z, 2) \quad \text{where } z \leq 2 \\
 &= 3
 \end{aligned}$$



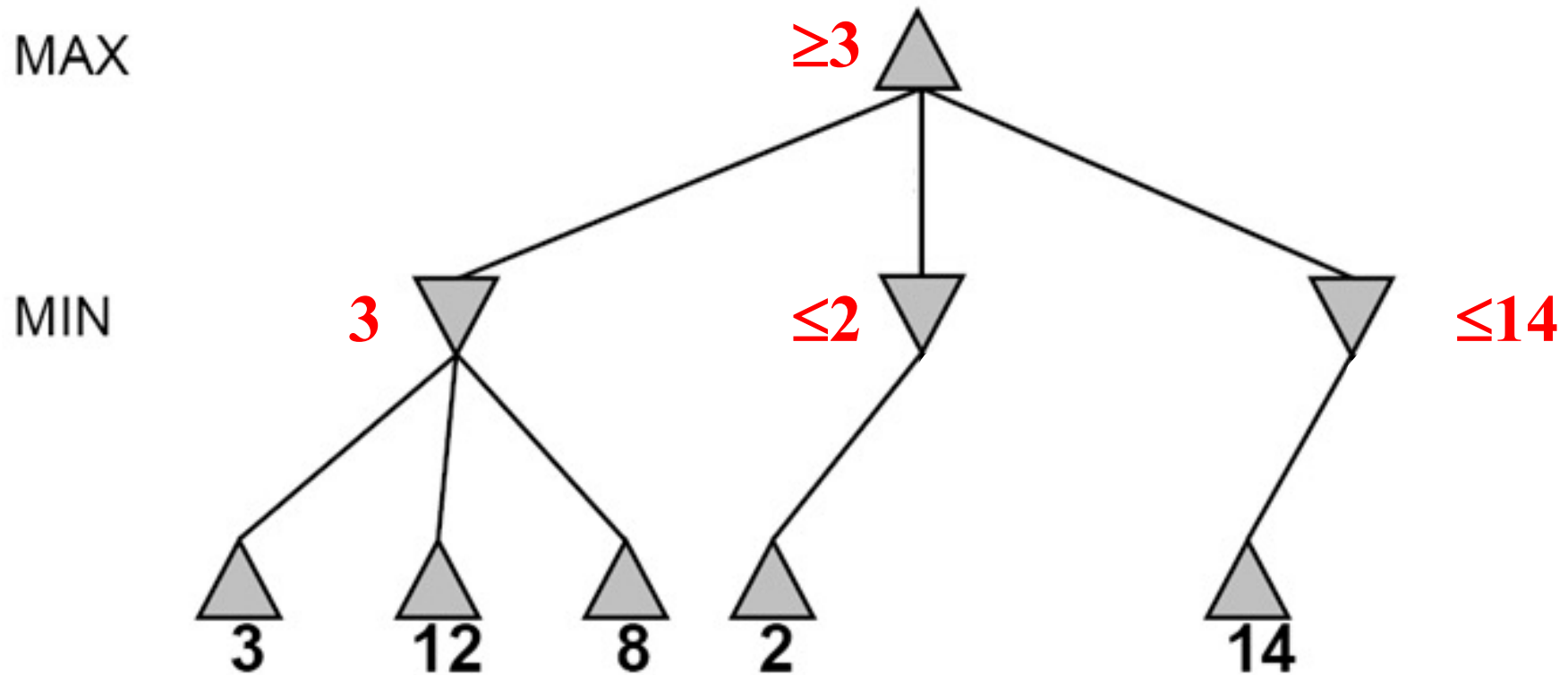
Alpha-beta pruning



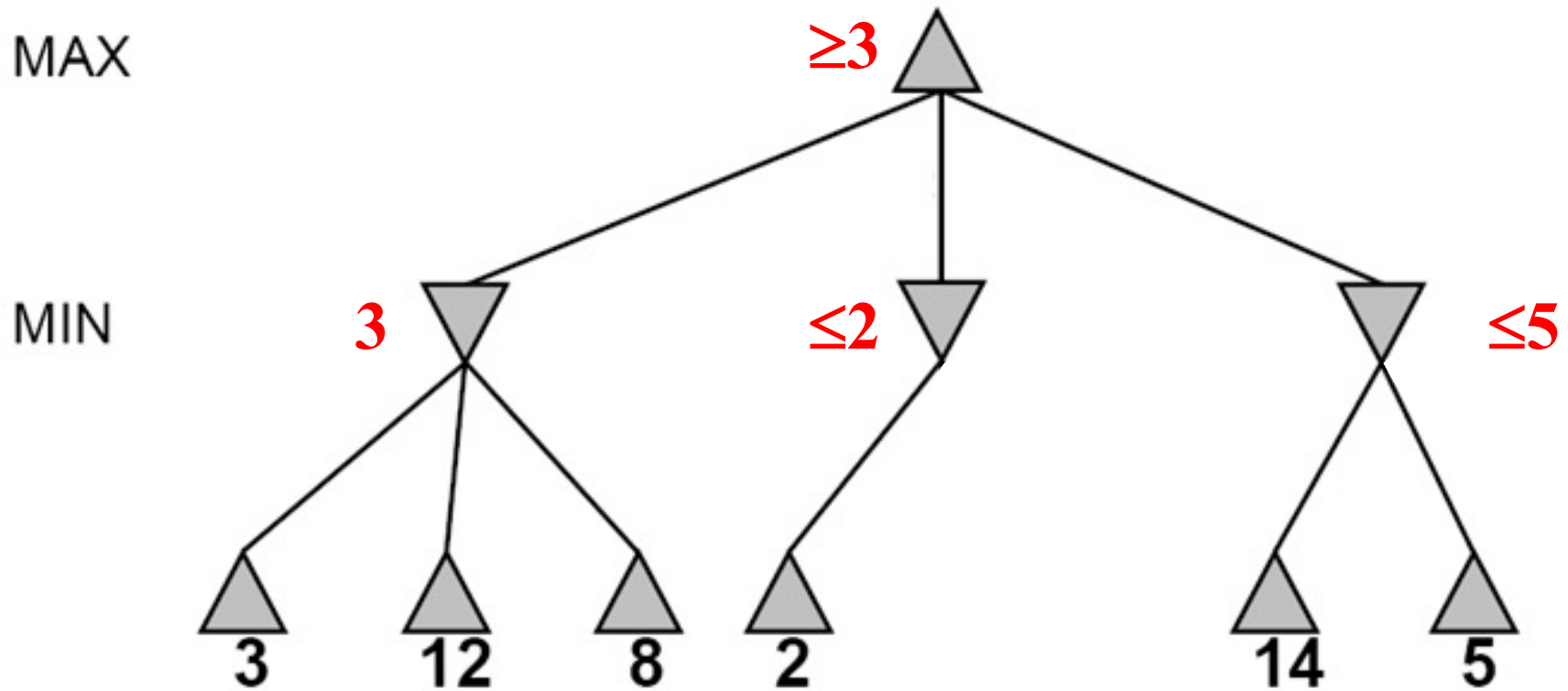
Alpha-beta pruning



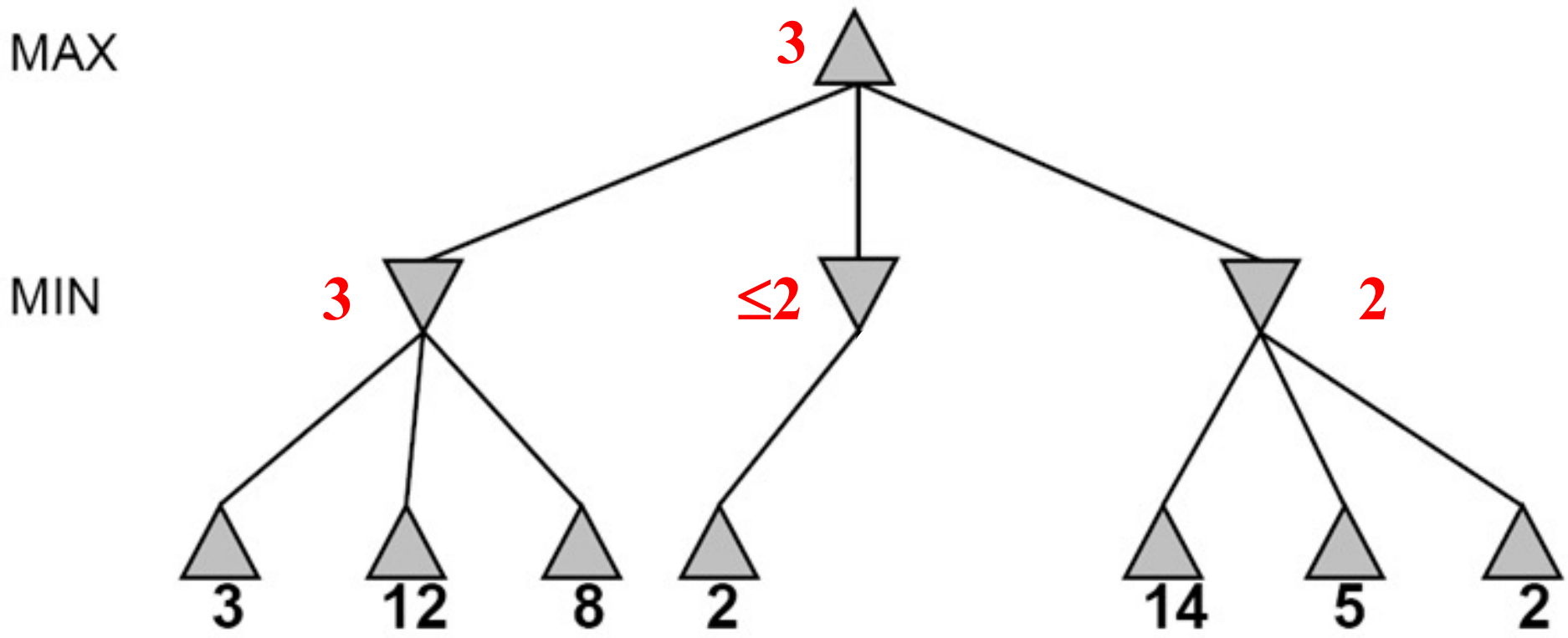
Alpha-beta pruning



Alpha-beta pruning

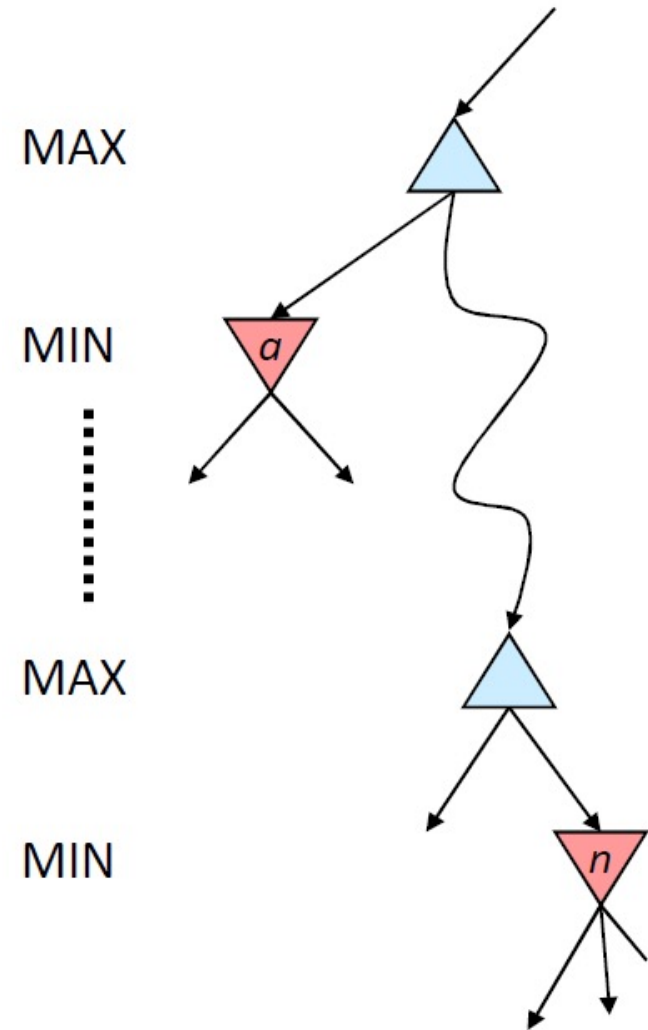


Alpha-beta pruning



Alpha-beta pruning

- α is the value of the best choice for the MAX player found so far at any choice point above node n
- We want to compute the MIN-value at n
- As we loop over n 's children, the MIN-value decreases
- If it drops below α , MAX will never choose n , so we can ignore n 's remaining children
- Analogously, β is the value of the lowest-utility choice found so far for the MIN player



The α - β algorithm

function ALPHA-BETA-SEARCH(*state*) **returns** an action
 $v \leftarrow \text{MAX-VALUE}(\text{state}, -\infty, +\infty)$
return the *action* in ACTIONS(*state*) with value v

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

function MIN-VALUE(*state*, α , β) **returns** a utility value
if TERMINAL-TEST(*state*) **then return** UTILITY(*state*)
 $v \leftarrow +\infty$
for each a **in** ACTIONS(*state*) **do**
 $v \leftarrow \text{MIN}(v, \text{MAX-VALUE}(\text{RESULT}(s, a), \alpha, \beta))$
if $v \leq \alpha$ **then return** v
 $\beta \leftarrow \text{MIN}(\beta, v)$
return v

Alpha-beta pruning

Function $action = \text{Alpha-Beta-Search}(node)$

$v = \text{Max-Value}(node, -\infty, \infty)$

return the *action* from *node* with value v

α : best alternative available to the Max player

β : best alternative available to the Min player

Function $v = \text{Max-Value}(node, \alpha, \beta)$

if Terminal(*node*) return Utility(*node*)

$v = -\infty$

for each *action* from *node*

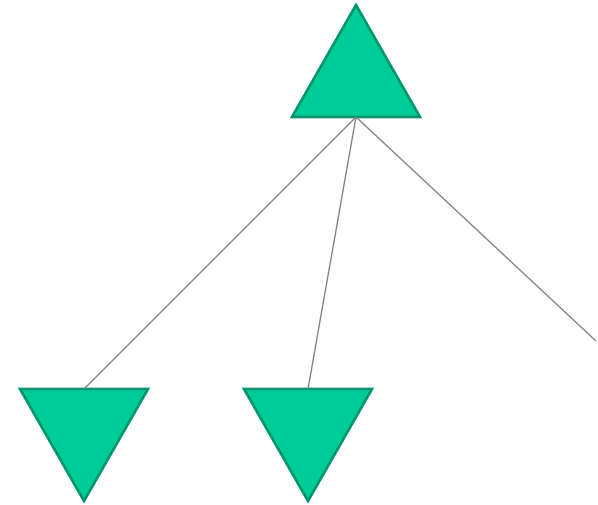
$v = \text{Max}(v, \text{Min-Value}(\text{Succ}(node, action), \alpha, \beta))$

if $v \geq \beta$ return v

$\alpha = \text{Max}(\alpha, v)$

end for

return v



Function $action = \text{Alpha-Beta-Search}(node)$

$v = \text{Min-Value}(node, -\infty, \infty)$

return the $action$ from $node$ with value v

α : best alternative available to the Max player

β : best alternative available to the Min player

Function $v = \text{Min-Value}(node, \alpha, \beta)$

if $\text{Terminal}(node)$ return $\text{Utility}(node)$

$v = +\infty$

for each $action$ from $node$

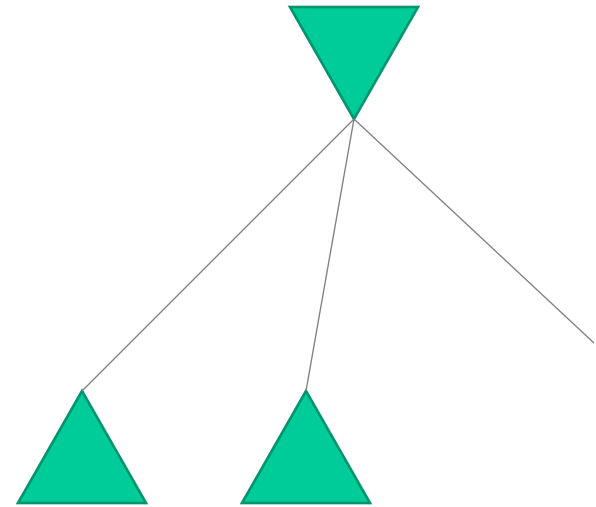
$v = \text{Min}(v, \text{Max-Value}(\text{Succ}(node, action), \alpha, \beta))$

if $v \leq \alpha$ return v

$\beta = \text{Min}(\beta, v)$

end for

return v



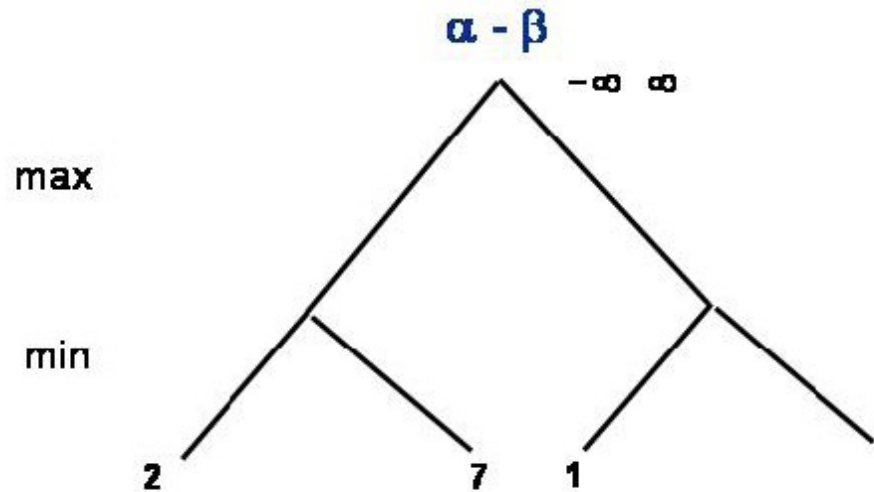
α - β pruning example

α - β

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



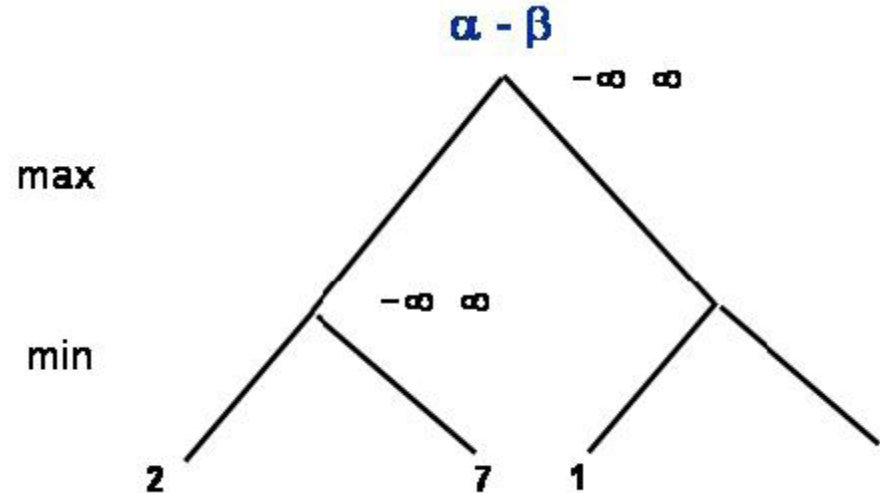
α - β pruning example

α - β

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



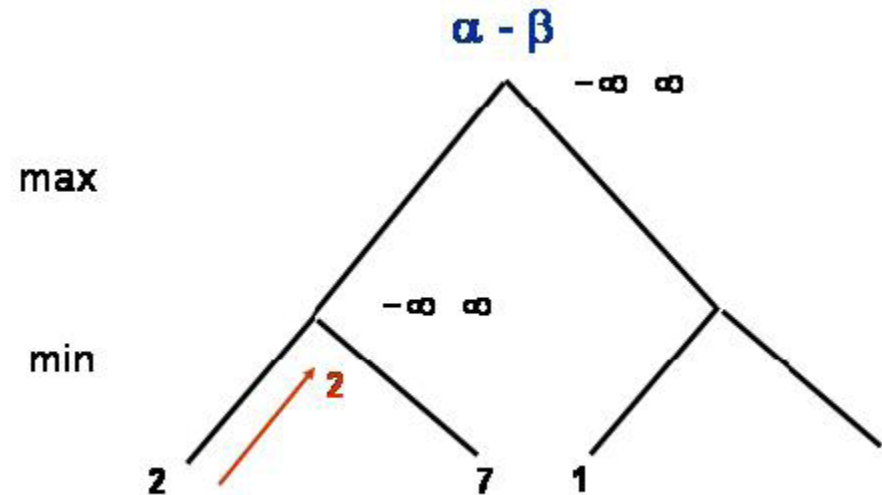
α - β pruning example

$\alpha - \beta$

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



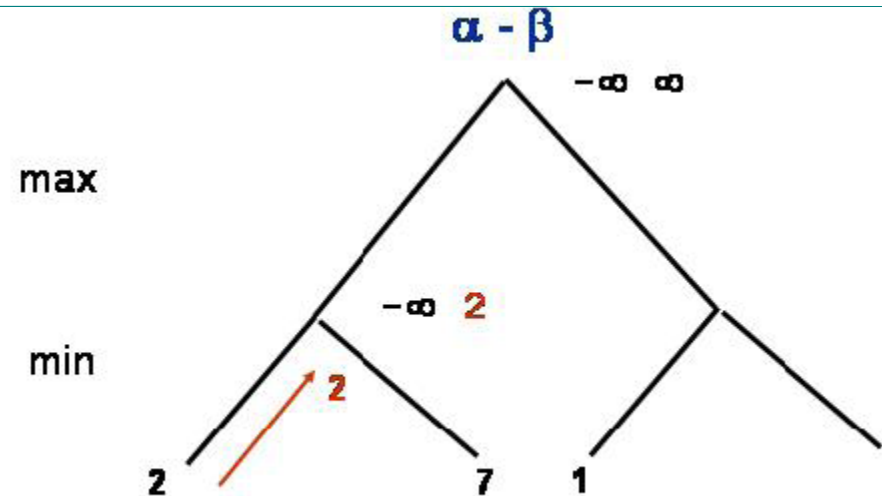
α - β pruning example

$\alpha - \beta$

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



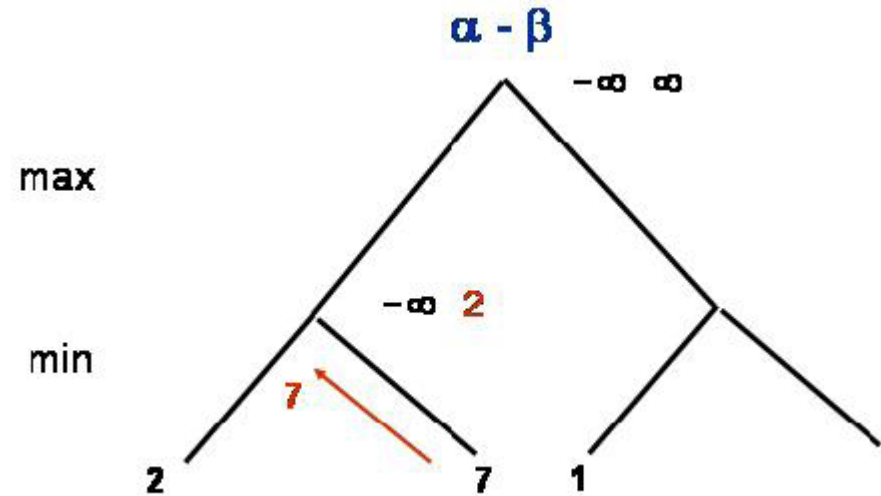
α - β pruning example

$\alpha - \beta$

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



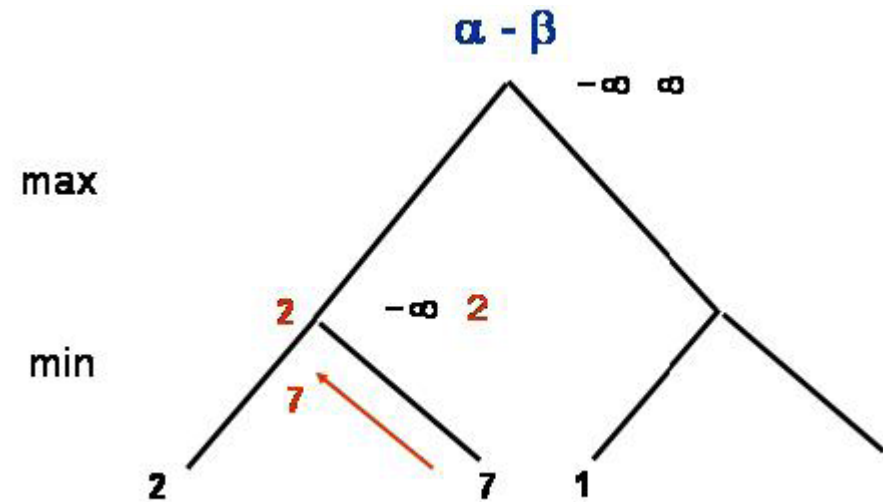
α - β pruning example

$\alpha - \beta$

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



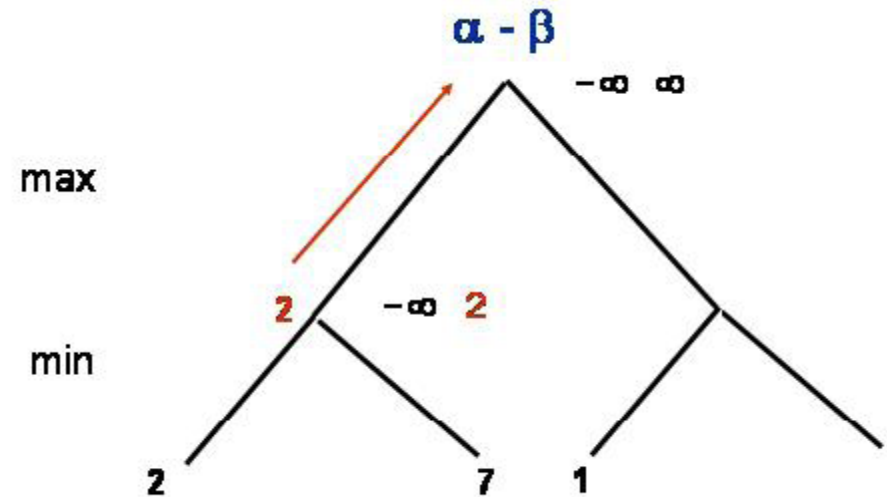
α - β pruning example

α - β

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



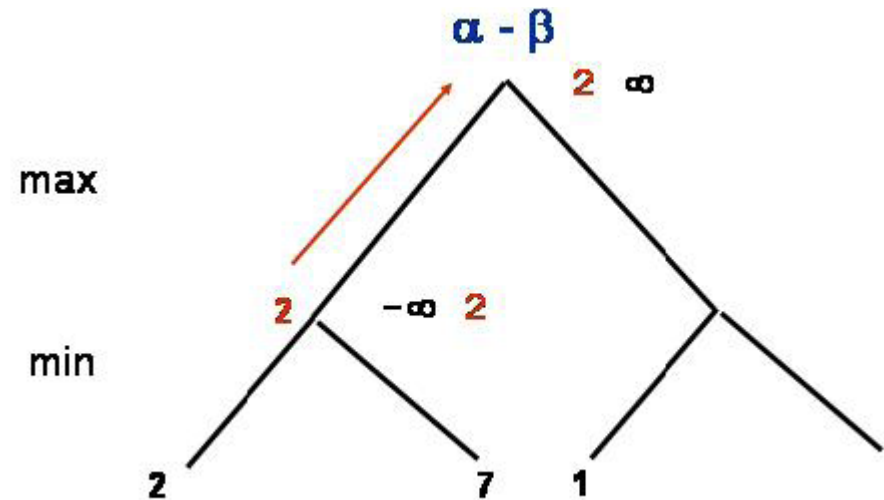
α - β pruning example

α - β

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



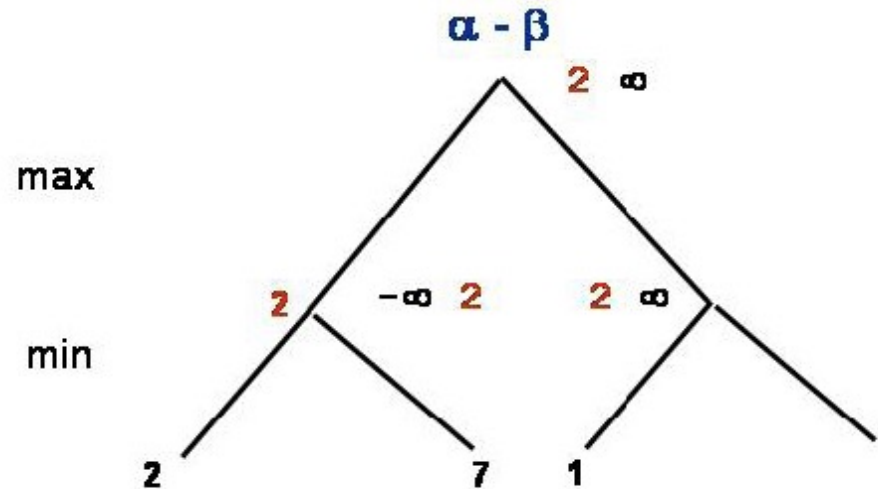
α - β pruning example

$\alpha - \beta$

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



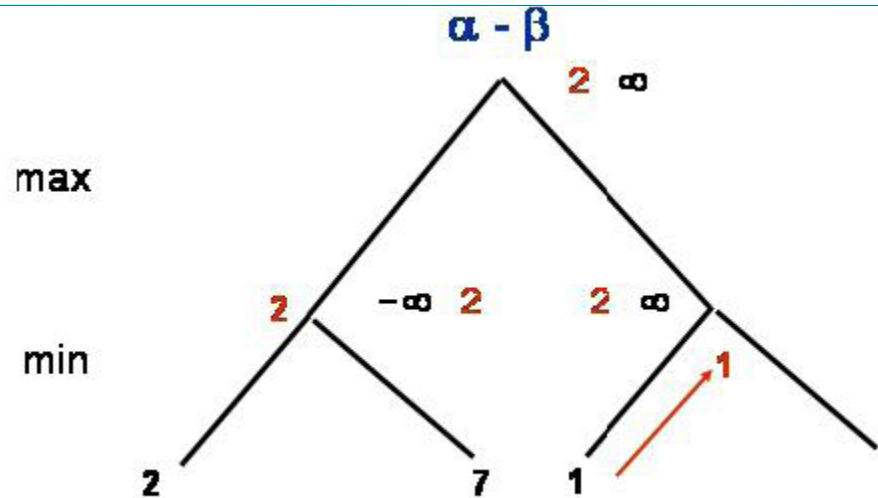
α - β pruning example

$\alpha - \beta$

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



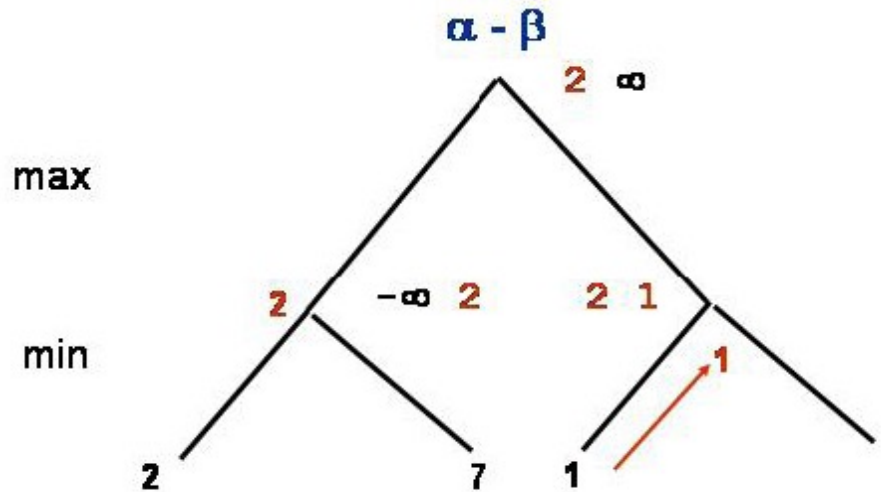
α - β pruning example

α - β

// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



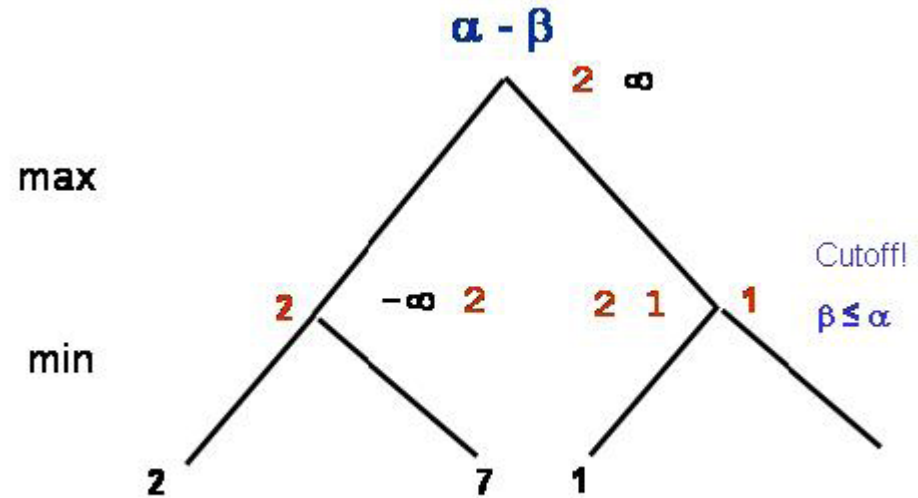
α - β pruning example

α - β

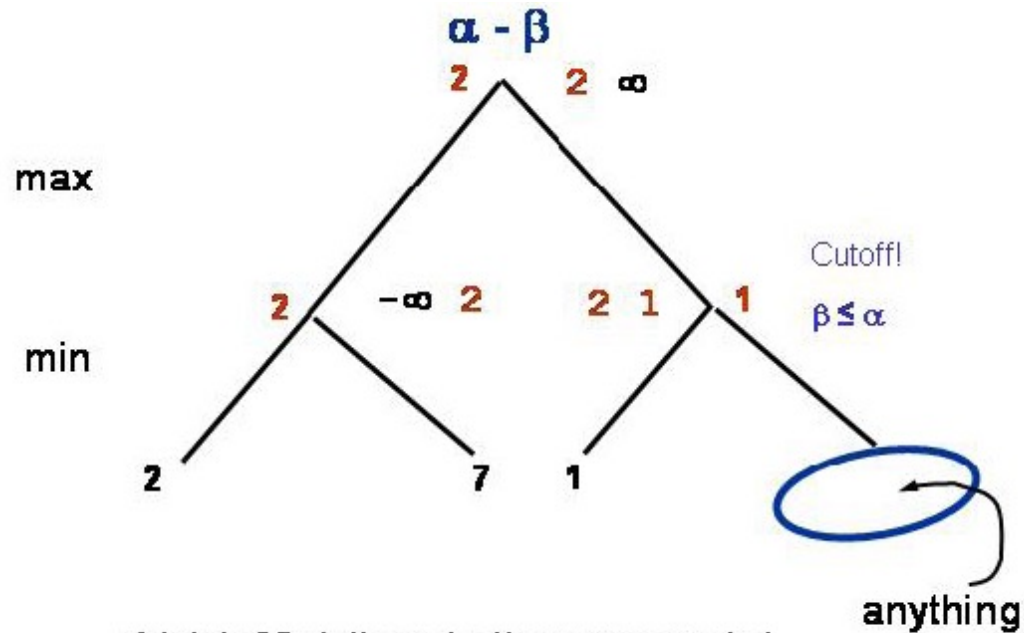
// α = best score for MAX, β = best score for MIN
 // initial call is MAX-VALUE(state, $-\infty$, ∞ , MAX-DEPTH)

```
function MAX-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\alpha$  = MAX ( $\alpha$ , MIN-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\alpha \geq \beta$  then return  $\alpha$  // cutoff
  end
  return  $\alpha$ 
```

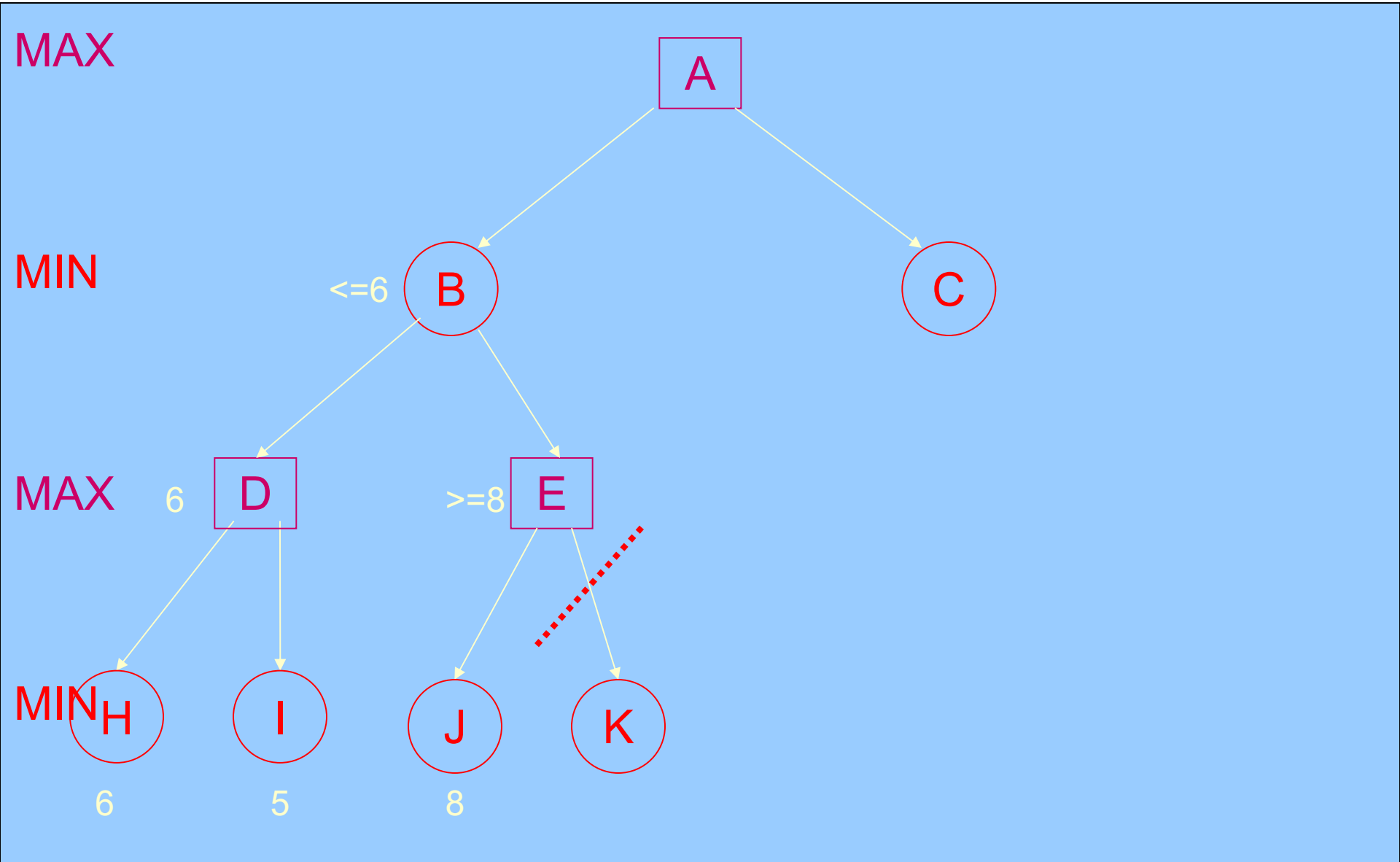
```
function MIN-VALUE (state,  $\alpha$ ,  $\beta$ , depth)
  if (depth == 0) then return EVAL (state)
  for each s in SUCCESSORS (state) do
     $\beta$  = MIN ( $\beta$ , MAX-VALUE (s,  $\alpha$ ,  $\beta$ , depth-1))
    if  $\beta \leq \alpha$  then return  $\beta$  // cutoff
  end
  return  $\beta$ 
```



α - β pruning example



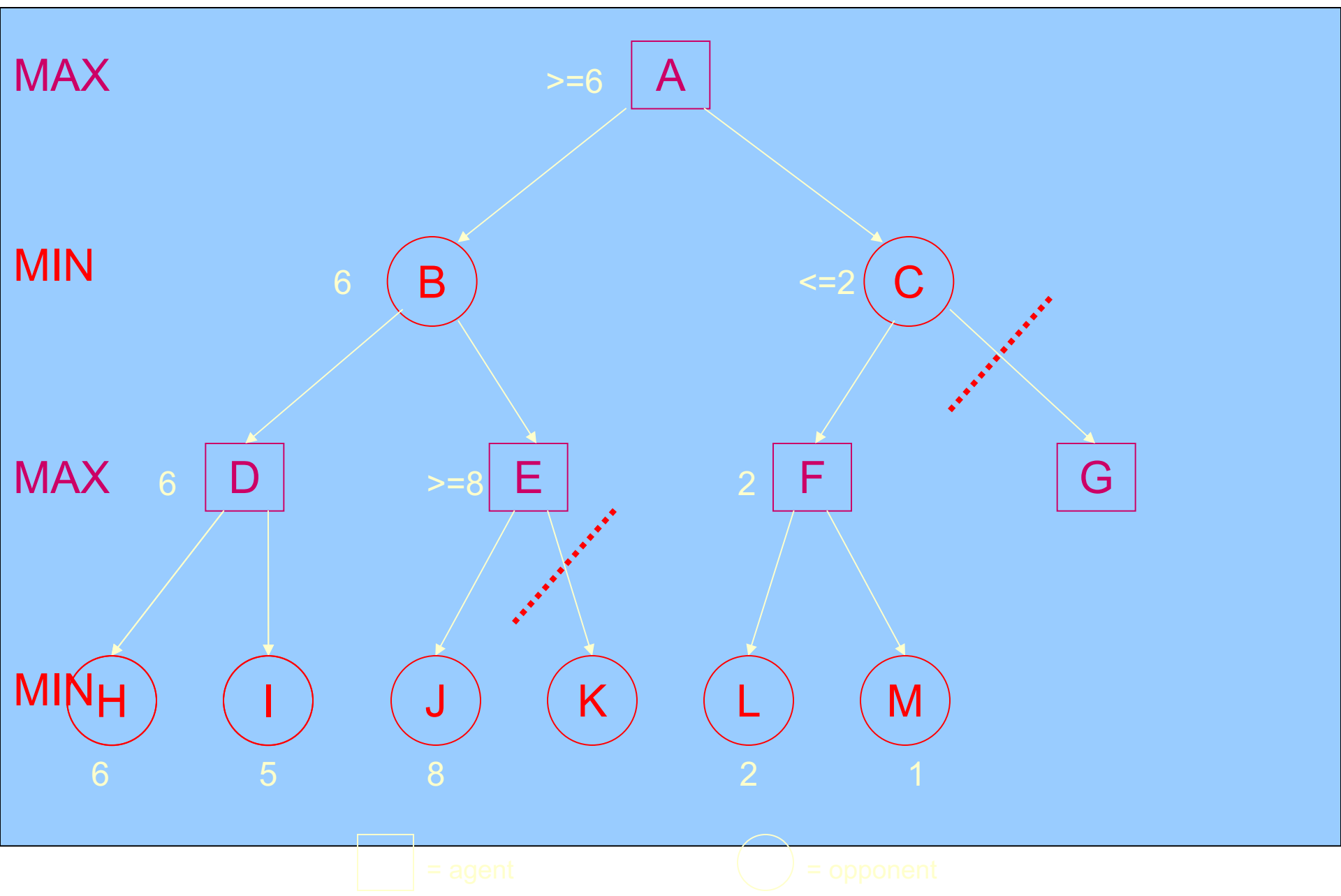
A total of 3 static evaluations were needed to obtain the value for the tree.

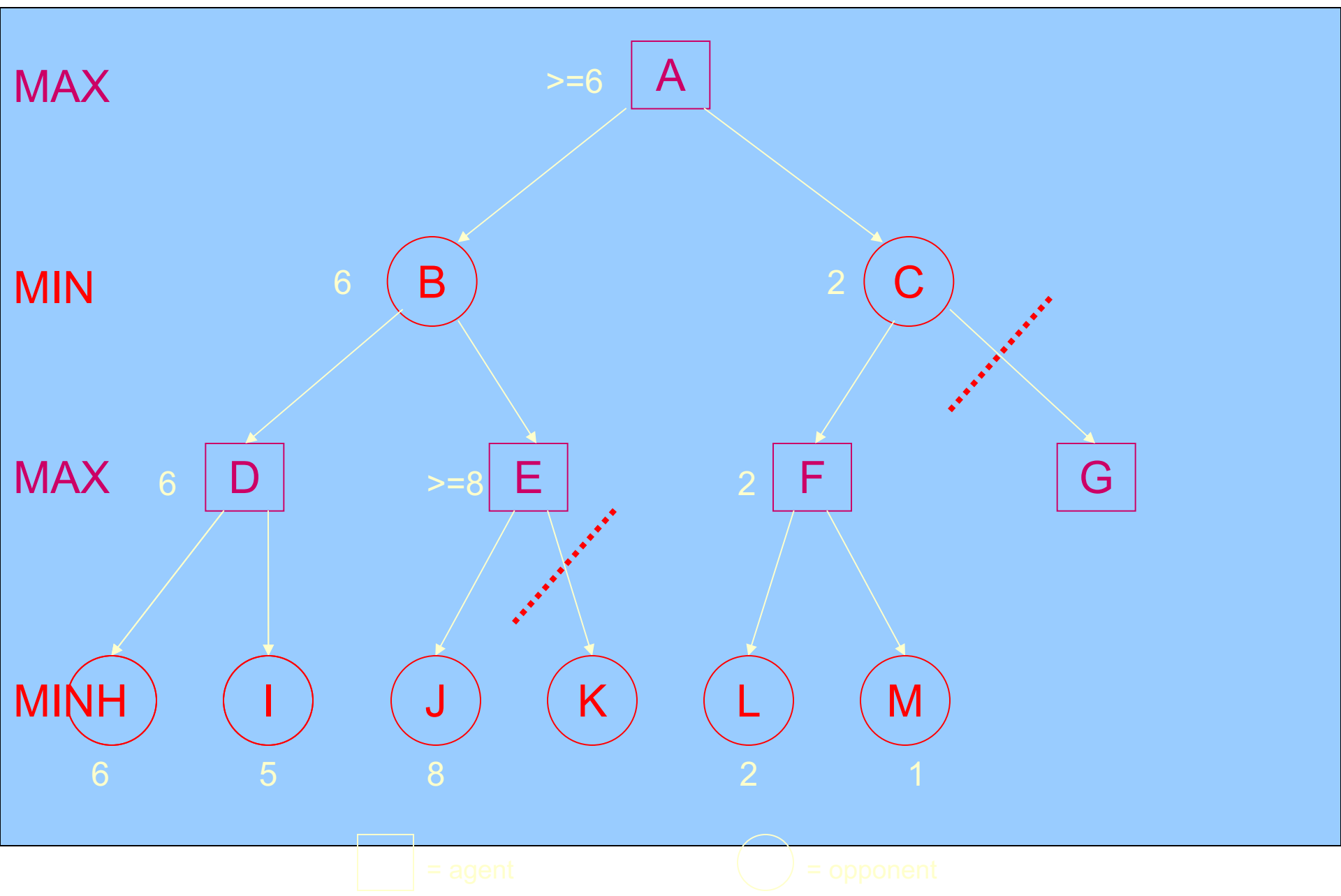


= agent

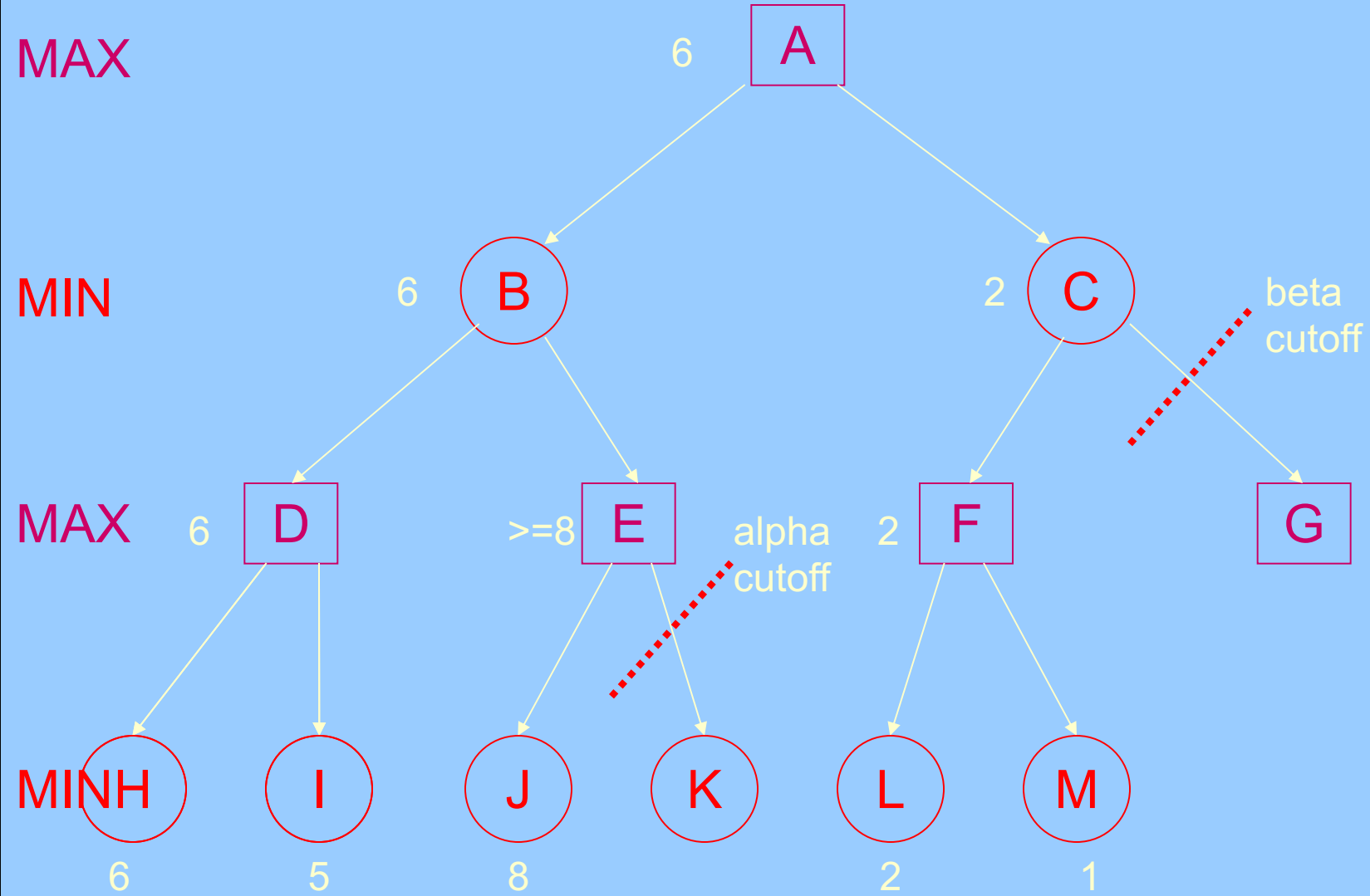


= opponent





Alpha-beta Pruning

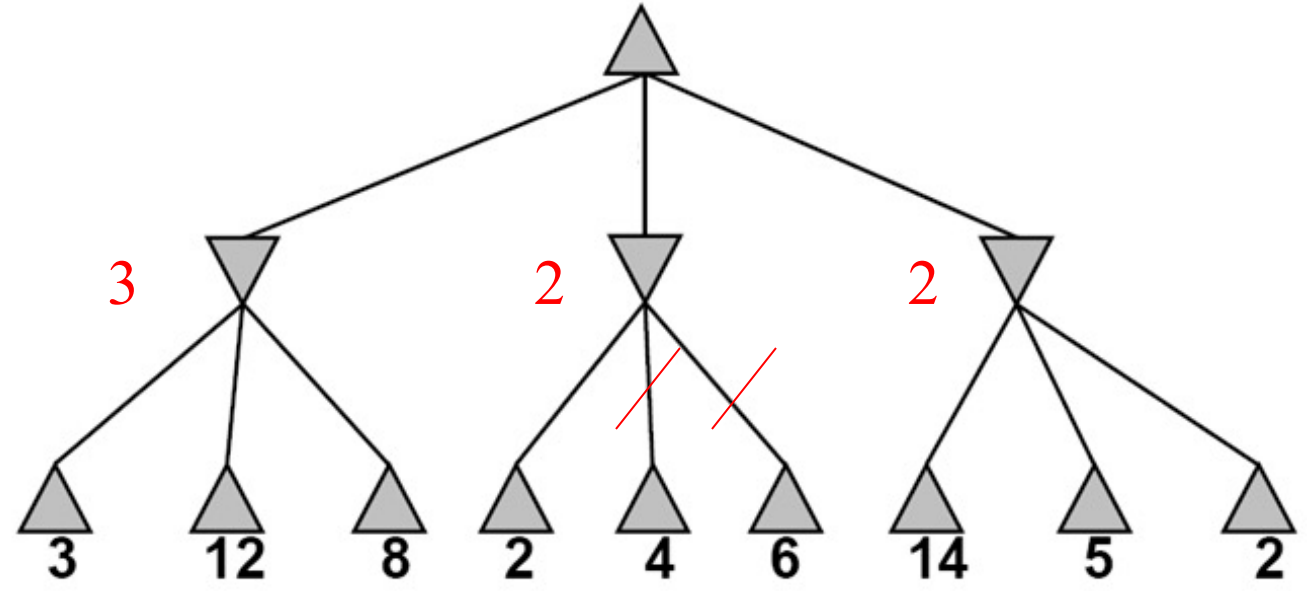


= agent

= opponent

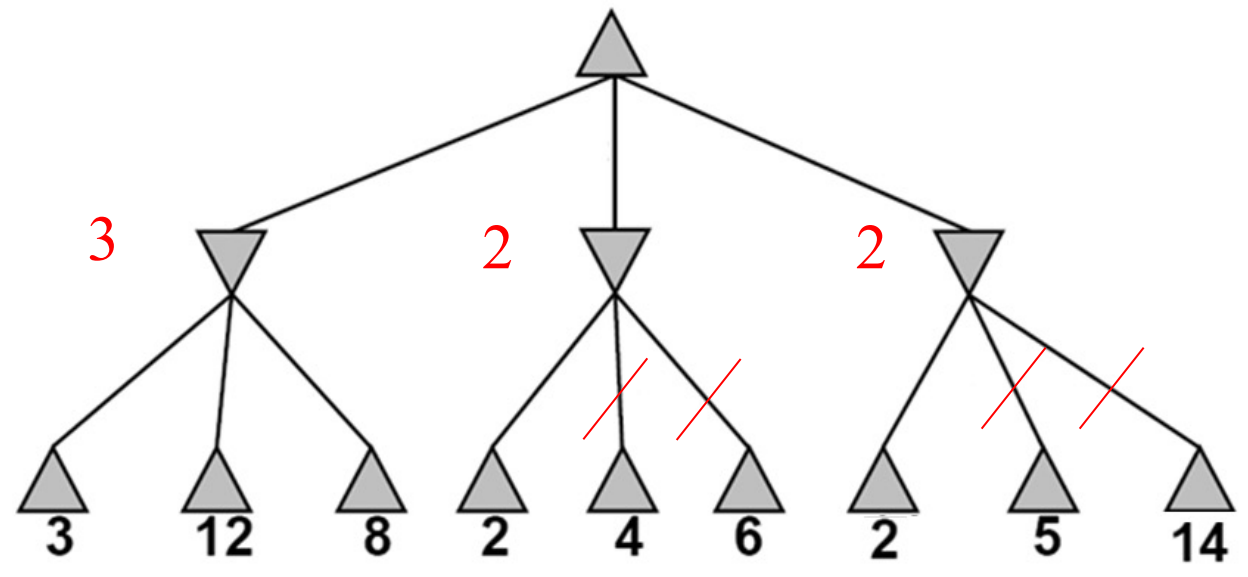
MAX

MIN



MAX

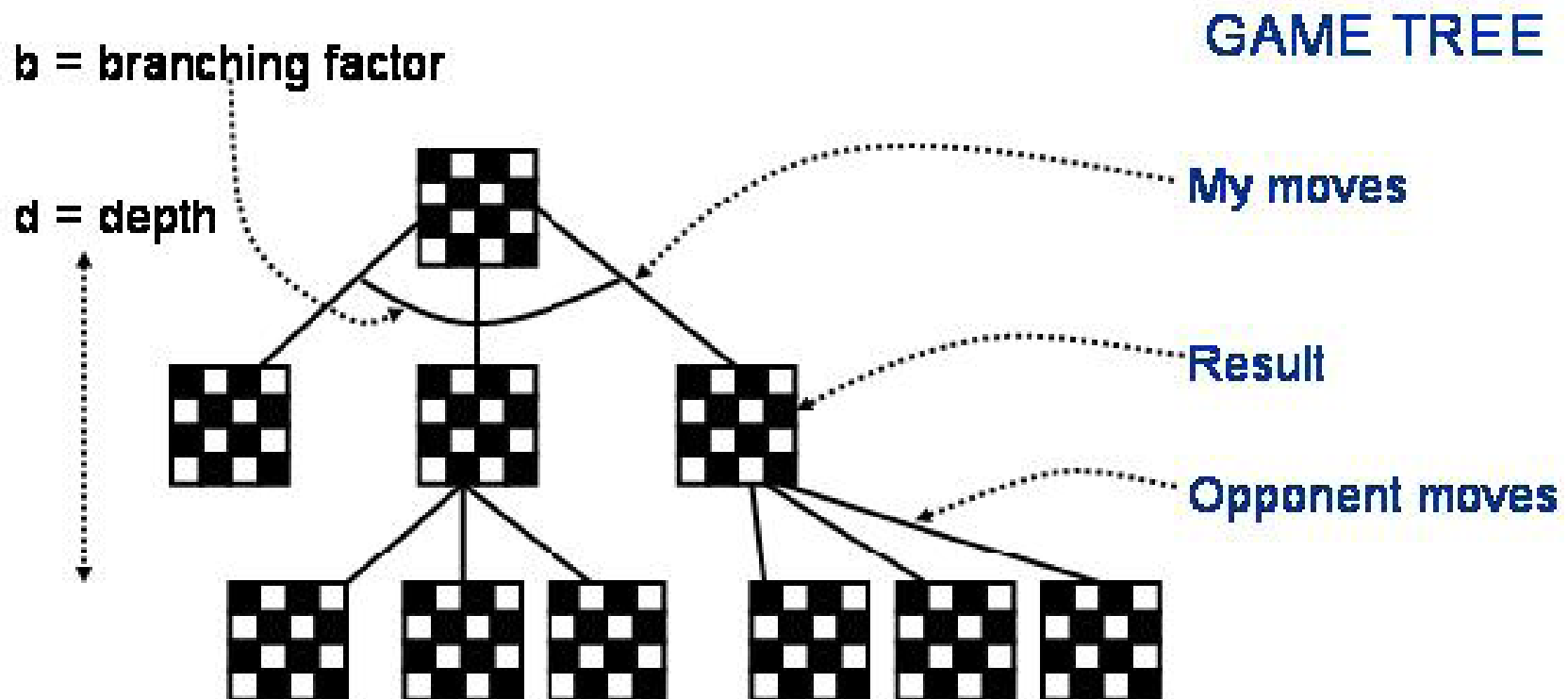
MIN



Alpha-beta pruning

- Pruning does not affect final result
 - Amount of pruning depends on move ordering
 - Should start with the “best” moves (highest-value for MAX or lowest-value for MIN)
 - For chess, can try captures first, then threats, then forward moves, then backward moves
 - Can also try to remember “killer moves” from other branches of the tree
 - With perfect ordering, the time to find the best move is reduced to $O(b^{m/2})$ from $O(b^m)$
 - Depth of search is effectively doubled
-

Move generation



Chess

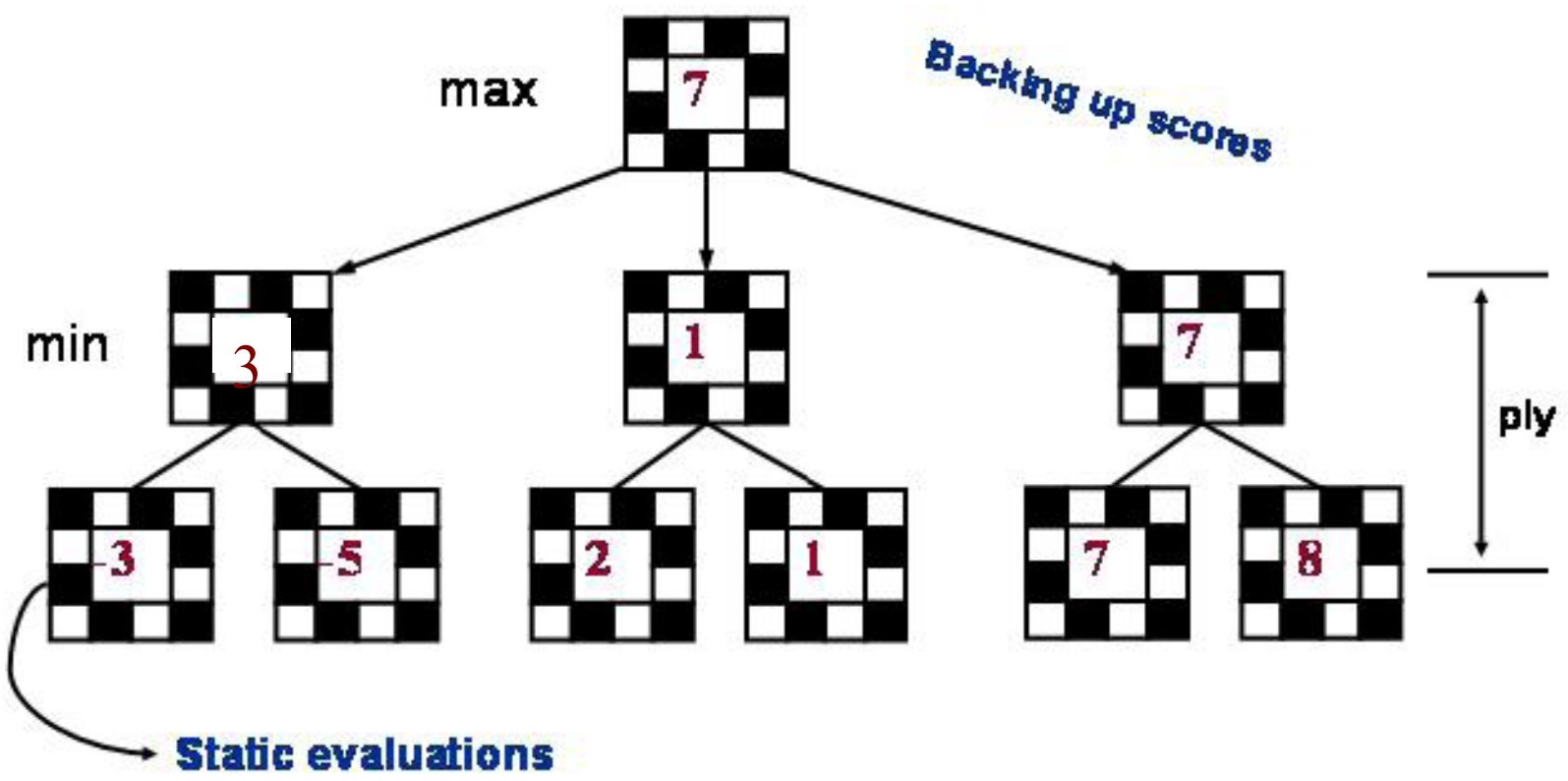
b = 36

d > 40

36⁴⁰

is big!

Min-Max



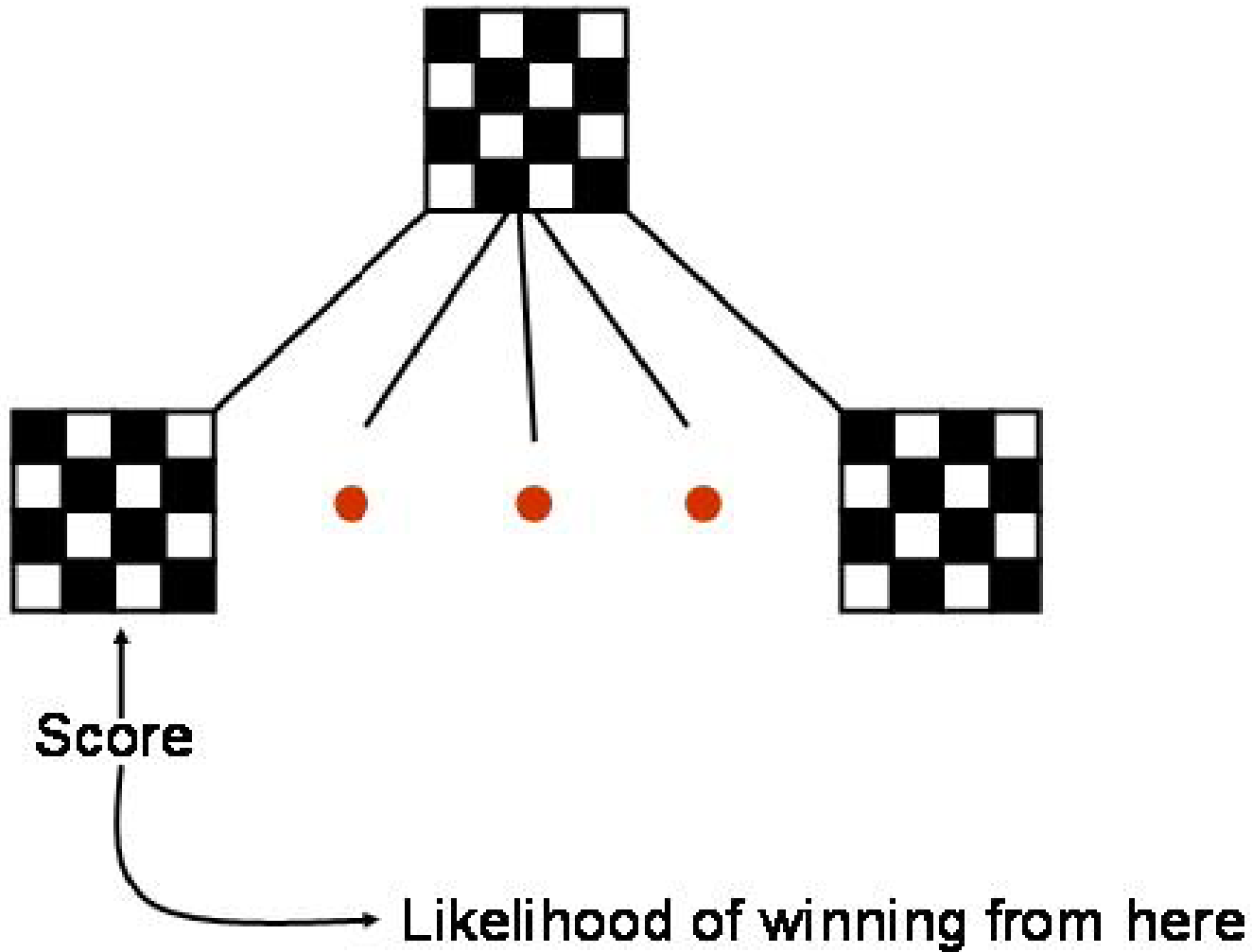
Resource limits

Suppose we have 100 secs, explore 10^4 nodes/sec
→ 10^6 nodes per move

Standard approach:

- **cutoff test:**
e.g., depth limit (perhaps add **quiescence search**)
 - **evaluation function**
= estimated desirability of position
-

Evaluation function



Evaluation function

S = **c**₁ **x** **material**
 + **c**₂ **x** **pawn structure**
 + **c**₃ **x** **mobility**
 + **c**₄ **x** **king safety**
 + **c**₅ **x** **center control**
 + ...

P	1
K	3
B	3.5
R	5
Q	9

- "material", : some measure of which pieces one has on the board.
- A typical weighting for each type of chess piece is shown
- Other types of features try to encode something about the distribution of the pieces on the board.

Evaluation functions

- A typical evaluation function is a linear function in which some set of coefficients is used to weight a number of "features" of the board position.
- weighted sum of *features*:

$$Eval(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

- For chess, w_k may be the **material value** of a piece (pawn = 1, knight = 3, rook = 5, queen = 9) and $f_k(s)$ may be the advantage in terms of that piece
 - Eg. $w_1 = 9$ with
 $f_1(s) = (\text{number of white queens}) - (\text{number of black queens})$
-

Evaluation function

- Cut off search at a certain depth and compute the value of an **evaluation function** for a state instead of its minimax value
 - The evaluation function may be thought of as the probability of winning from a given state or the *expected value* of that state
 - If a position A has a 100% chance of winning it should have the evaluation 1
 - If position B have a 50% chance of winning and 25% os loosing and 25% of being a draw, the evaluation value would be $+1 \times 0.50 + -1 \times 0.25 + 0 \times 0.25 = 0.25$
 - Evaluation functions may be *learned* from game databases or by having the program play many games against itself
-

Cutting off search

MinimaxCutoff is identical to *MinimaxValue* except

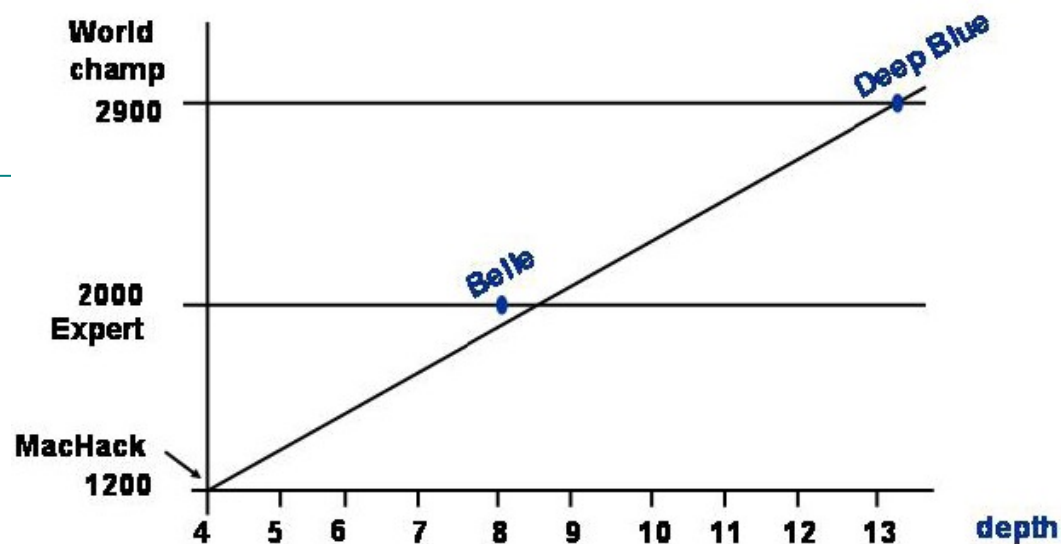
1. *Terminal?* is replaced by *Cutoff?*
2. *Utility* is replaced by *Eval*

Does it work in practice?

$$b^m = 10^6, b=35 \quad m=4$$

4-ply lookahead is a hopeless chess player!

- 4-ply \approx human novice
- 8-ply \approx typical PC, human master
- 12-ply \approx Deep Blue, Kasparov



Chess playing systems

- Baseline system: 200 million node evaluations per move (3 min), minimax with a decent evaluation function and quiescence search
 - 5-ply \approx human novice
 - Add alpha-beta pruning
 - 10-ply \approx typical PC, experienced player
 - Deep Blue: 30 billion evaluations per move, singular extensions, evaluation function with 8000 features, large databases of opening and endgame moves
 - 14-ply \approx Garry Kasparov
 - More recent state of the art ([Hydra](#), ca. 2006): 36 billion evaluations per second, advanced pruning techniques
 - 18-ply \approx better than any human alive?
-

Practical issues

Variable branching



Iterative deepening

- └ order best move from last search first
- └ use previous backed up value to initialize $[\alpha, \beta]$
- └ keep track of repeated positions (transposition tables)

Horizon effect

- └ quiescence
- └ Pushing the inevitable over search horizon

Parallelization

Cutting off search

- **Horizon effect:** you may incorrectly estimate the value of a state by overlooking an event that is just beyond the depth limit
 - For example, a damaging move by the opponent that can be delayed but not avoided
 - Possible remedies
 - **Quiescence search:** do not cut off search at positions that are unstable – for example, are you about to lose an important piece?
 - **Singular extension:** a strong move that should be tried when the normal depth limit is reached
-

Types of game environments

	Deterministic	Stochastic
Perfect information (fully observable)	Chess, checkers, go	Backgammon, monopoly
Imperfect information (partially observable)	Battleships	Scrabble, poker, bridge

Stochastic games

- How to incorporate dice throwing into the game tree?



Stochastic games

MAX

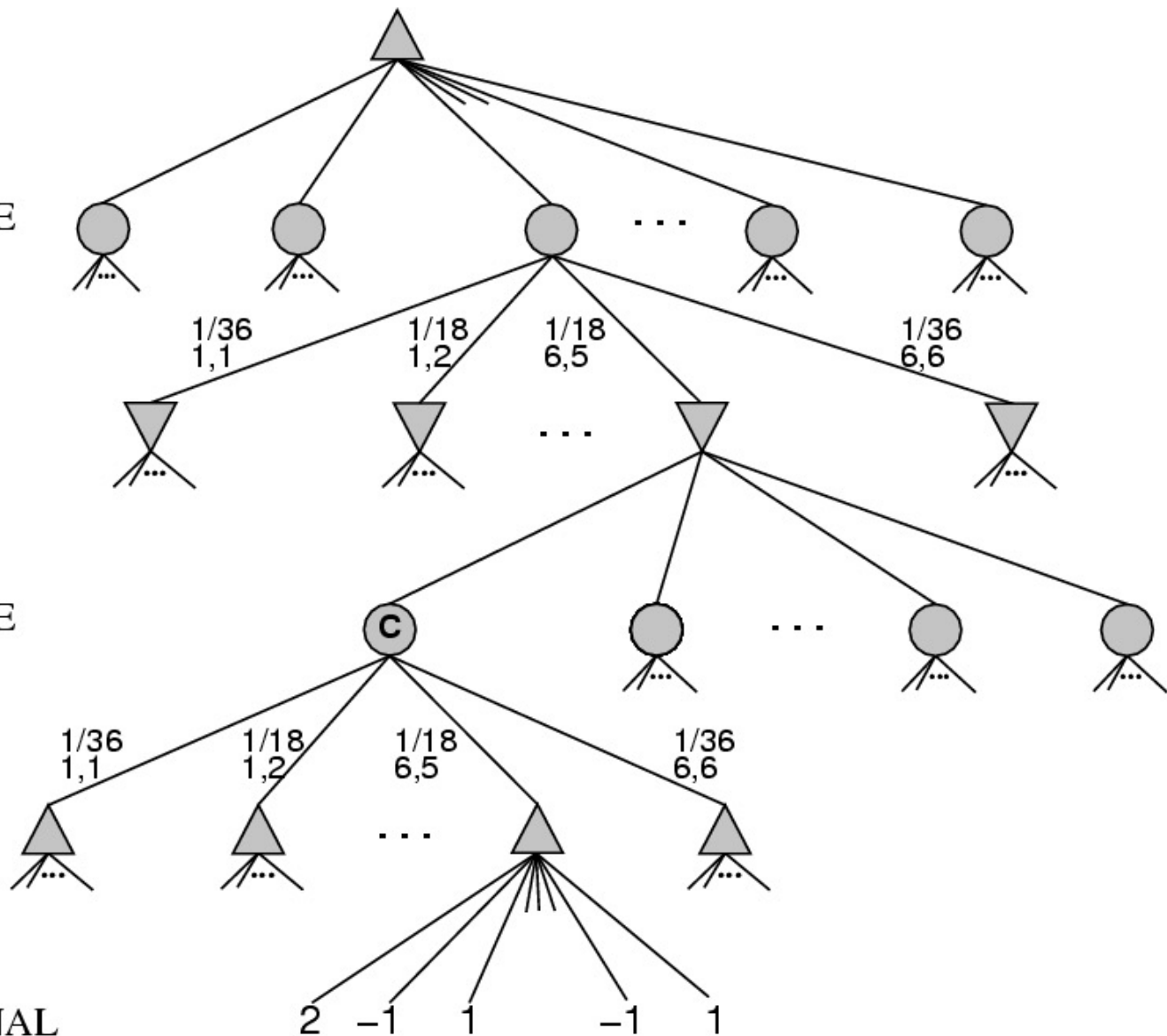
CHANCE

MIN

CHANCE

MAX

TERMINAL



Minimax vs. Expectiminimax

- **Minimax:**

- **Maximize** (over all possible moves I can make) the
- **Minimum** (over all possible moves Min can make) of the
- Reward

$$Value(node) = \max_{my\ moves} \left(\min_{Min's\ moves} (Reward) \right)$$

- **Expectiminimax:**

- **Maximize** (over all possible moves I can make) the
- **Minimum** (over all possible moves Min can make) of the
- **Expected** reward

$$Value(node) = \max_{my\ moves} \left(\min_{Min's\ moves} (\mathbb{E}[Reward]) \right)$$

$$\mathbb{E}[Reward] = \sum_{outcomes} Probability(outcome) \times Reward(outcome)$$

Stochastic games

- **Expectiminimax:** for chance nodes, sum values of successor states weighted by the probability of each successor

 - **Value**(*node*) =
 - $Utility(node)$ if *node* is terminal
 - $\max_{action} \mathbf{Value}(Succ(node, action))$ if *type* = MAX
 - $\min_{action} \mathbf{Value}(Succ(node, action))$ if *type* = MIN
 - $\sum_{action} P(Succ(node, action)) * \mathbf{Value}(Succ(node, action))$ if *type* = CHANCE
-

Expectiminimax summary

- All of the same methods are useful:
 - Alpha-Beta pruning
 - Evaluation function
 - Quiescence search, Singular move
 - Computational complexity is pretty bad
 - Branching factor of the random choice can be high
 - Twice as many “levels” in the tree
-

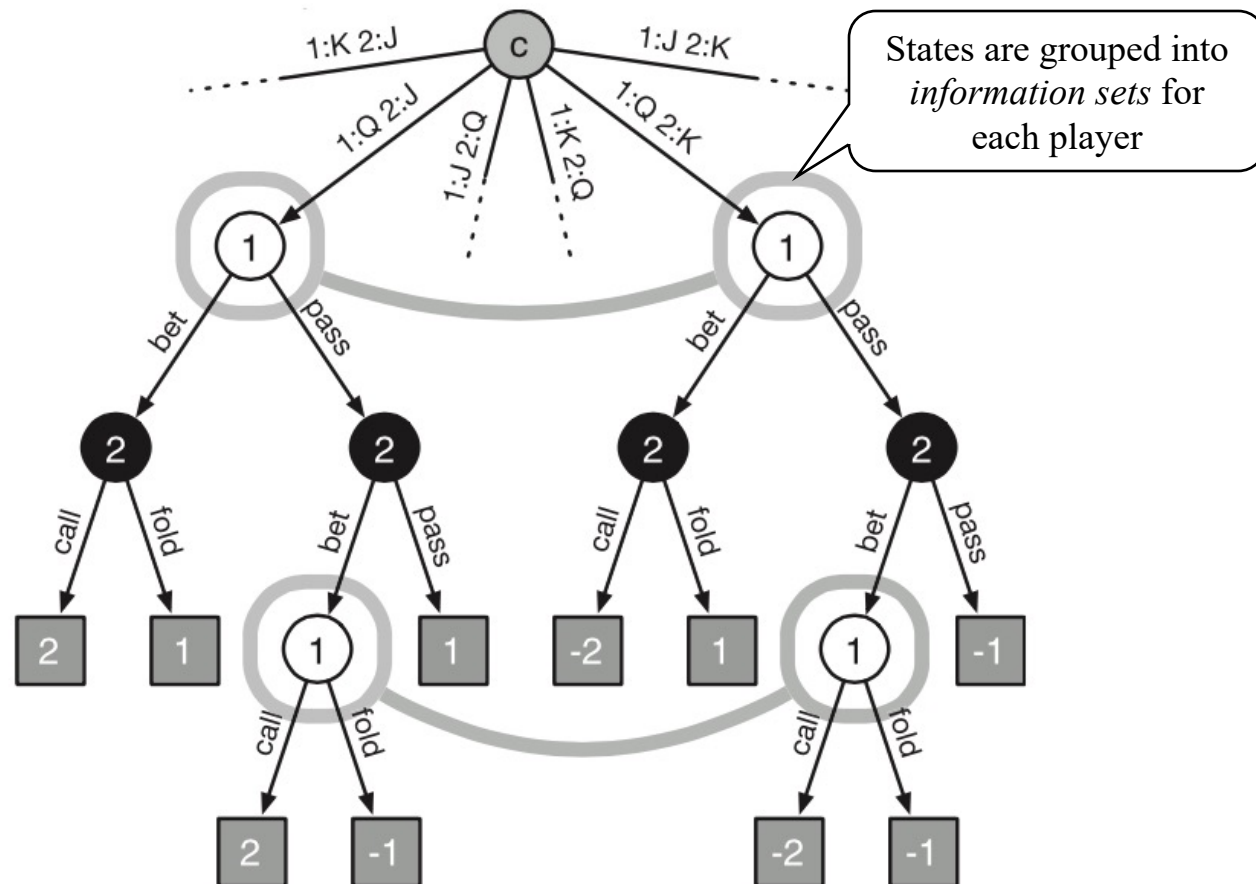
Stochastic games

- **Expectiminimax:** for chance nodes, sum values of successor states weighted by the probability of each successor
 - Nasty branching factor, defining evaluation functions and pruning algorithms more difficult
 - **Monte Carlo simulation:** when you get to a chance node, simulate a large number of games with random dice rolls and use win percentage as evaluation function
 - Can work well for games like Backgammon
-

Stochastic games of imperfect information

Fig. 1. Portion of the extensive-form game representation of three-card Kuhn poker (16).

Player 1 is dealt a queen (Q), and the opponent is given either the jack (J) or king (K). Game states are circles labeled by the player acting at each state ("c" refers to chance, which randomly chooses the initial deal). The arrows show the events the acting player can choose from, labeled with their in-game meaning. The leaves are square vertices labeled with the associated utility for player 1 (player 2's utility is the negation of player 1's). The states connected by thick gray lines are part of the same information set; that is, player 1 cannot distinguish between the states in each pair because they each represent a different unobserved card being dealt to the opponent. Player 2's states are also in information sets, containing other states not pictured in this diagram.



Stochastic games of imperfect information

- Simple Monte Carlo approach: run multiple simulations with random cards pretending the game is fully observable
 - “Averaging over clairvoyance”
 - Problem: this strategy does not account for bluffing, information gathering, etc.

Miniminimax with imperfect information

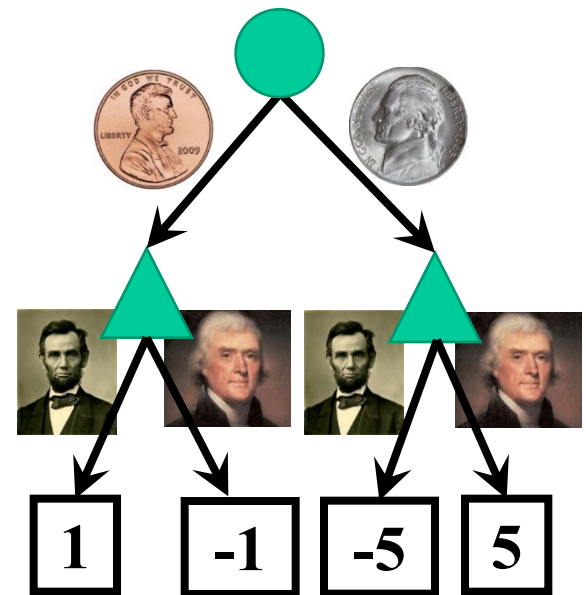
- **Minimax:**

- **Maximize** (over all possible moves I can make) the
- **Minimum**
 - (over all possible states of the information I don't know,
 - ... over all possible moves Min can make) the
- Reward.

$$Value(node) = \max_{my\ moves} \left(\min_{\substack{missing\ info, \\ Min's\ moves}} (Reward) \right)$$

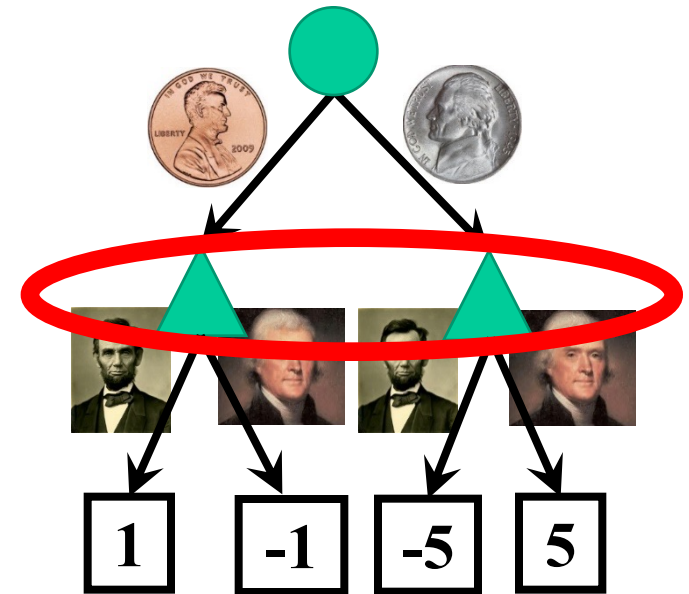
Imperfect information example

- Min chooses a coin.
- I say the name of a U.S. President.
 - If I guessed right, she gives me the coin.
 - If I guessed wrong, I have to give her a coin to match the one she has.





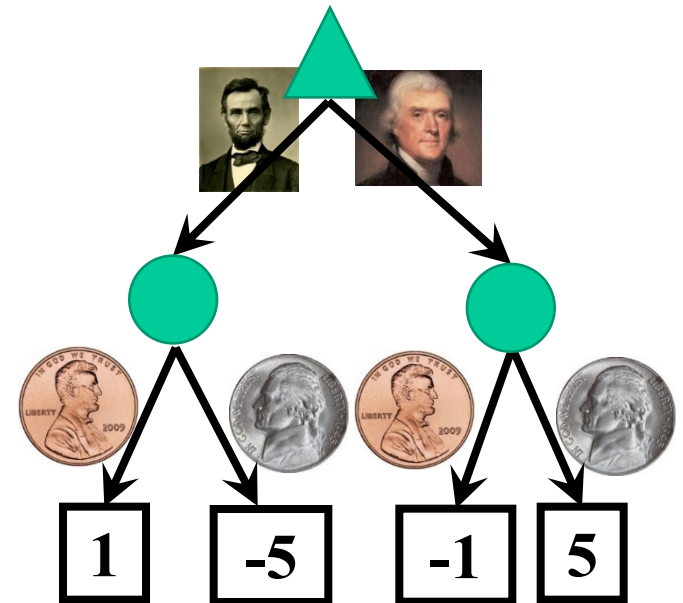
Method #1: Treat “unknown” as “unknown”

- The problem: I don't know which state I'm in. I only know it's one of these two.
- The solution: choose the policy that maximizes my minimum reward.
 - “Lincoln”: minimum reward is -5.
 - “Jefferson”: minimum reward is -1.
- Minimax policy: say “Jefferson”.



Method #2: Treat “unknown” as “random”

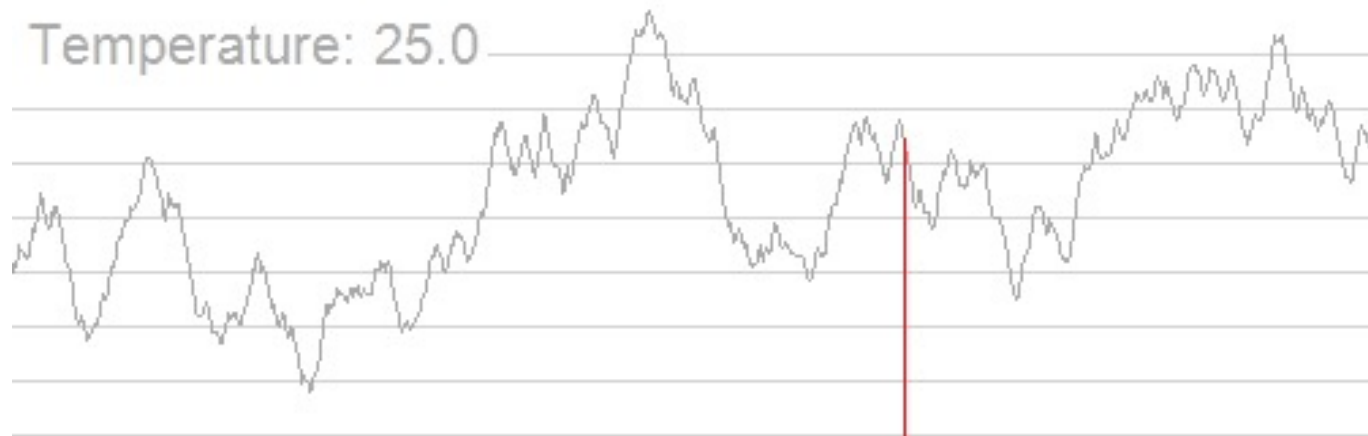
- Expectiminimax: treat the unknown information as random.
- Choose the policy that maximizes my expected reward.
 - “Lincoln”: $\frac{1}{2} \times 1 + \frac{1}{2} \times (-5) = -2$
 - “Jefferson”: $\frac{1}{2} \times (-1) + \frac{1}{2} \times 5 = 2$
- Expectiminimax policy: say “Jefferson”.
- BUT WHAT IF:  and  are not equally likely?



How to deal with imperfect information

- If you think you know the probabilities of different settings, and if you want to maximize your average winnings (for example, you can afford to play the game many times):
expectiminimax
 - If you have no idea of the probabilities of different settings; or, if you can only afford to play once, and you can't afford to lose: **miniminimax**
 - If the unknown information has been selected intentionally by your opponent: use **game theory**
-

Stochastic search



Stochastic search for stochastic games

- The problem with expectiminimax: huge branching factor (many possible outcomes)

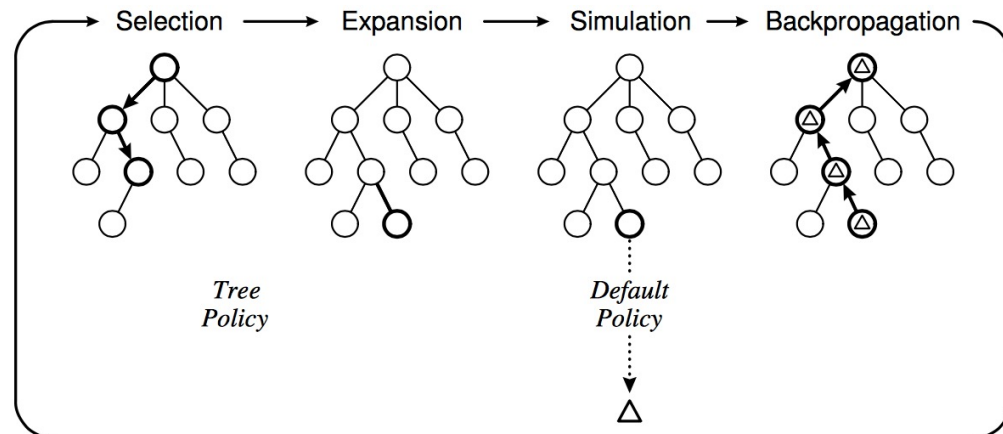
$$\mathbb{E}[\textit{Reward}] = \sum_{\textit{outcomes}} \textit{Probability}(\textit{outcome}) \times \textit{Reward}(\textit{outcome})$$

- An approximate solution: Monte Carlo search

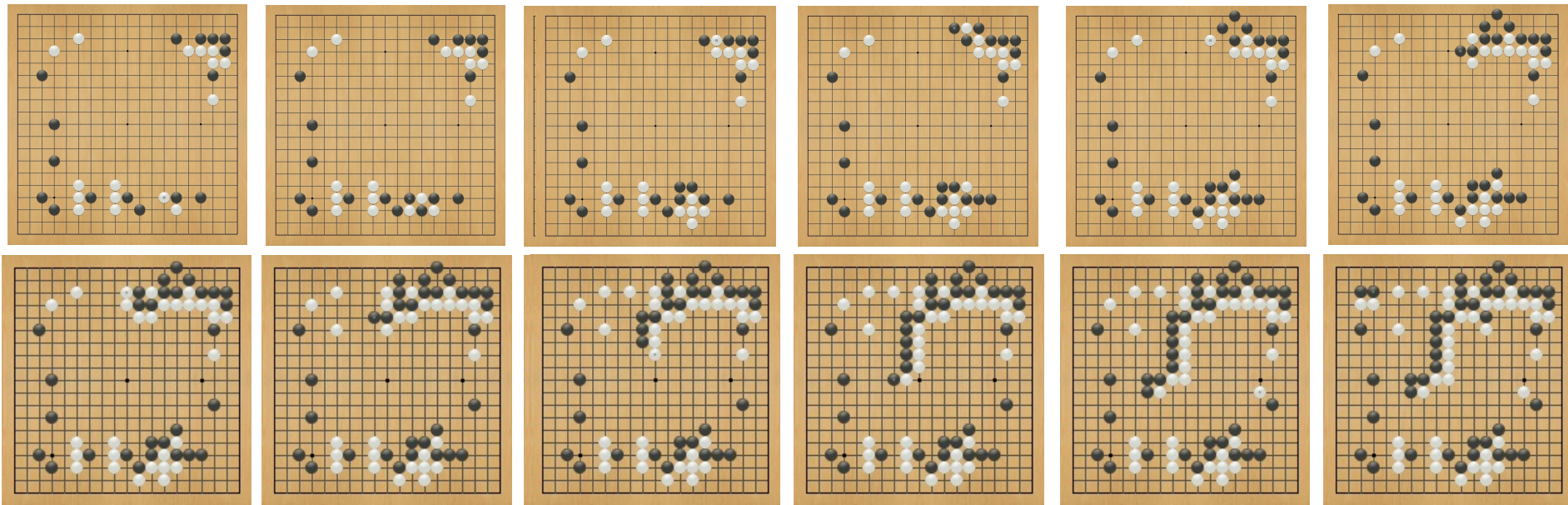
$$\mathbb{E}[\textit{Reward}] \approx \frac{1}{n} \sum_{i=1}^n \textit{Reward}(i' \textit{th random game})$$

Monte Carlo Tree Search

- What about *deterministic* games with deep trees, large branching factor, and no good heuristics – like Go?
- Instead of depth-limited search with an evaluation function, use randomized simulations
- Starting at the current state (root of search tree), iterate:
 - Select a leaf node for expansion using a *tree policy* (trading off *exploration* and *exploitation*)
 - Run a simulation using a *default policy* (e.g., random moves) until a terminal state is reached
 - Back-propagate the outcome to update the value estimates of internal tree nodes



Learned evaluation functions



Stochastic search off-line

Training phase:

- Spend a few weeks allowing your computer to play billions of random games from every possible starting state
- Value of the starting state = average value of the ending states achieved during those billion random games

Testing phase:

- During the alpha-beta search, search until you reach a state whose value you have stored in your value lookup table
 - Oops.... Why doesn't this work?
-

Evaluation as a pattern recognition problem

Training phase:

- Spend a few weeks allowing your computer to play billions of random games from billions of possible starting states.
- Value of the starting state = average value of the ending states achieved during those billion random games

Generalization:

- Featurize (e.g., x_1 = number of  patterns, x_2 = number of  patterns, etc.)

- Linear regression: find a_1, a_2 , etc. so that $\text{Value}(\text{state}) \approx a_1 * x_1 + x_2 * x_2 + \dots$

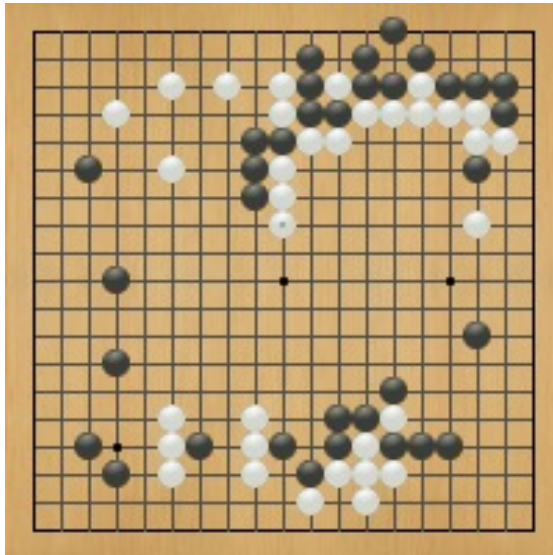
Testing phase:

- During the alpha-beta search, search as deep as you can, then estimate the value of each state at your horizon using $\text{Value}(\text{state}) \approx a_1 * x_1 + x_2 * x_2 + \dots$
-

Pros and Cons

- Learned evaluation function
 - Pro: off-line search permits lots of compute time, therefore lots of training data
 - Con: there's no way you can evaluate every starting state that might be achieved during actual game play. Some starting states will be missed, so generalized evaluation function is necessary
 - On-line stochastic search
 - Con: limited compute time
 - Pro: it's possible to estimate the value of the state you've reached during actual game play
-

Case study: AlphaGo



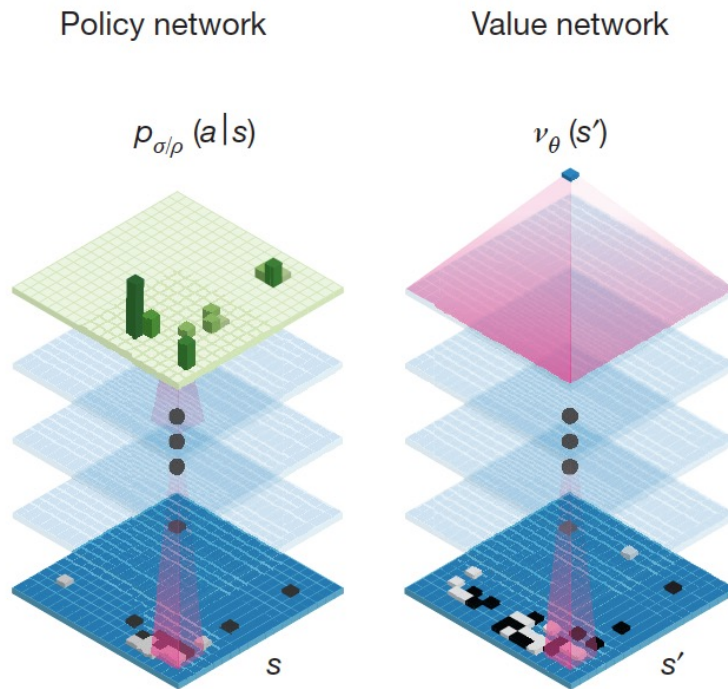
- **“Gentlemen should not waste their time on trivial games -- they should play go.”**
- *-- Confucius,*
- *The Analects*
- *ca. 500 B. C. E.*

Anton Ninno
antonninno@yahoo.com

Roy Laird, Ph.D.
roylaird@gmail.com

special thanks to Kiseido Publications

AlphaGo



- Deep convolutional neural networks
 - Treat the Go board as an image
 - Powerful function approximation machinery
 - Can be trained to predict distribution over possible moves (*policy*) or expected *value* of position

[Mastering the Game of Go with Deep Neural Networks and Tree Search](#)

AlphaGo

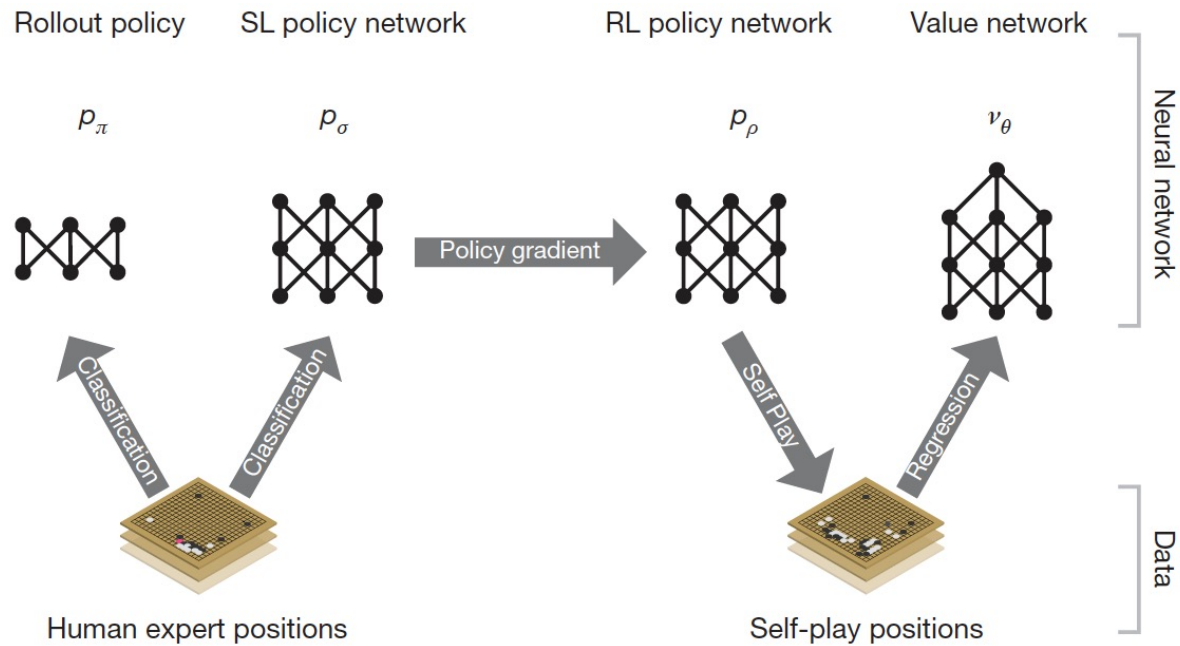
- SL policy network
 - Idea: perform *supervised learning* (SL) to predict human moves
 - Given state s , predict probability distribution over moves a , $P(a|s)$
 - Trained on 30M positions, 57% accuracy on predicting human moves
 - Also train a smaller, faster *rollout policy* network (24% accurate)
- RL policy network
 - Idea: fine-tune policy network using *reinforcement learning* (RL)
 - Initialize RL network to SL network
 - Play two snapshots of the network against each other, update parameters to maximize expected final outcome
 - RL network wins against SL network 80% of the time, wins against open-source Pachi Go program 85% of the time

AlphaGo

- SL policy network
- RL policy network
- Value network
 - Idea: train network for position evaluation
 - Given state s , estimate $v(s)$, expected outcome of play starting with position s and following the learned policy for both players
 - Train network by minimizing mean squared error between actual and predicted outcome
 - Trained on 30M positions sampled from different self-play games

[Mastering the Game of Go with Deep Neural Networks and Tree Search](#)

AlphaGo



[Mastering the Game of Go with Deep Neural Networks and Tree Search](#)

AlphaGo

- Monte Carlo Tree Search
 - Each edge in the search tree maintains *prior probabilities* $P(s,a)$, *counts* $N(s,a)$, *action values* $Q(s,a)$
 - $P(s,a)$ comes from SL policy network
 - Tree traversal policy selects actions that maximize Q value plus exploration bonus (proportional to P but inversely proportional to N)
 - An expanded leaf node gets a value estimate that is a combination of value network estimate and outcome of simulated game using rollout network
 - At the end of each simulation, Q values are updated to the average of values of all simulations passing through that edge

[Mastering the Game of Go with Deep Neural Networks and Tree Search](#)

AlphaGo

- Monte Carlo Tree Search

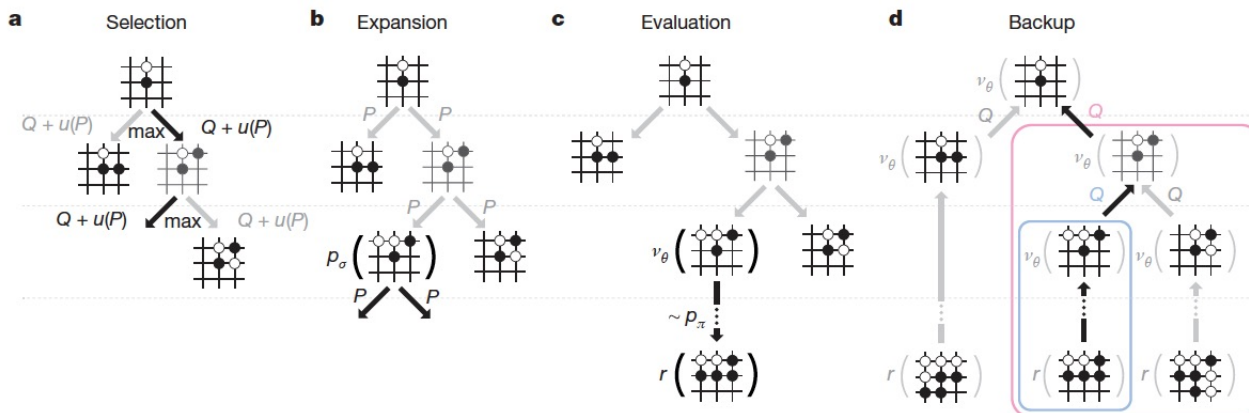


Figure 3 | Monte Carlo tree search in AlphaGo. **a**, Each simulation traverses the tree by selecting the edge with maximum action value Q , plus a bonus $u(P)$ that depends on a stored prior probability P for that edge. **b**, The leaf node may be expanded; the new node is processed once by the policy network p_σ and the output probabilities are stored as prior probabilities P for each action. **c**, At the end of a simulation, the leaf node

is evaluated in two ways: using the value network v_θ ; and by running a rollout to the end of the game with the fast rollout policy p_π , then computing the winner with function r . **d**, Action values Q are updated to track the mean value of all evaluations $r(\cdot)$ and $v_\theta(\cdot)$ in the subtree below that action.

[Mastering the Game of Go with Deep Neural Networks and Tree Search](#)

AlphaGo

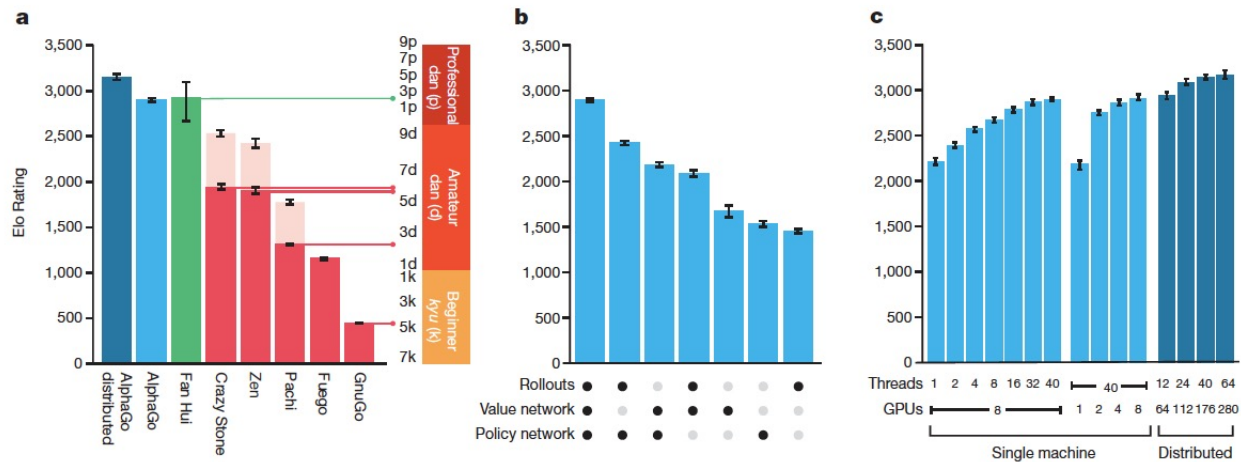
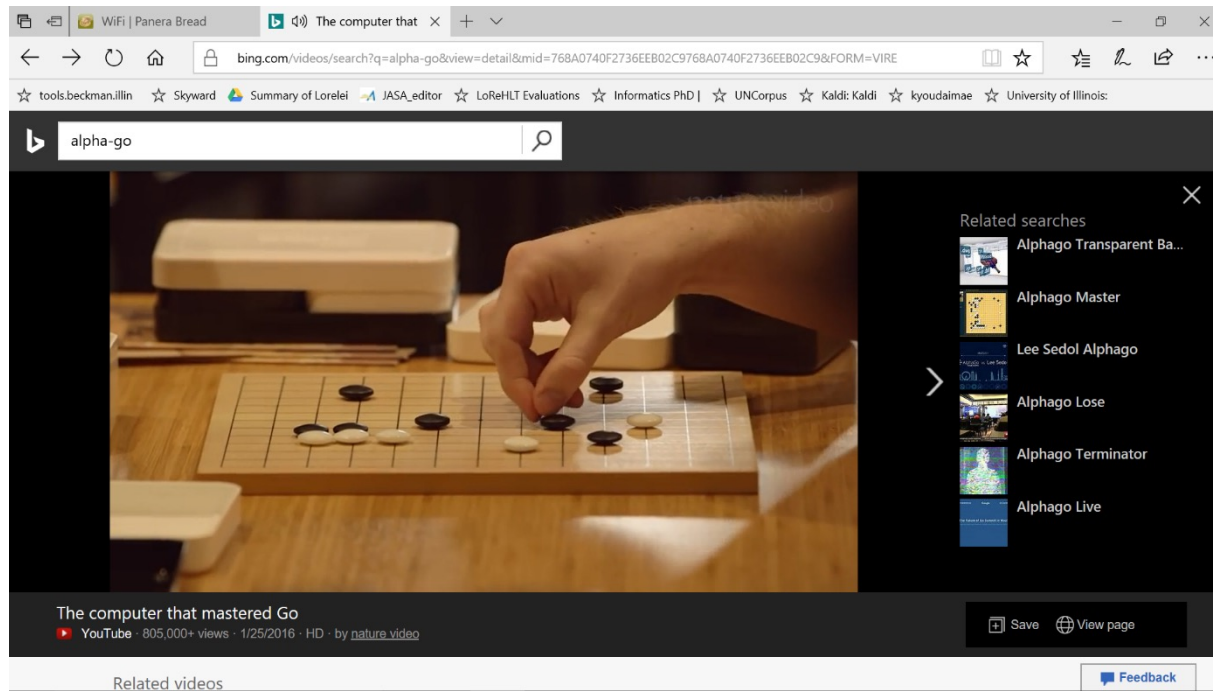


Figure 4 | Tournament evaluation of AlphaGo. a, Results of a tournament between different Go programs (see Extended Data Tables 6–11). Each program used approximately 5 s computation time per move. To provide a greater challenge to AlphaGo, some programs (pale upper bars) were given four handicap stones (that is, free moves at the start of every game) against all opponents. Programs were evaluated on an Elo scale³⁷: a 230 point gap corresponds to a 79% probability of winning, which roughly corresponds to one amateur *dan* rank advantage on KGS³⁸; an approximate correspondence to human ranks is also shown,

horizontal lines show KGS ranks achieved online by that program. Games against the human European champion Fan Hui were also included; these games used longer time controls. 95% confidence intervals are shown. b, Performance of AlphaGo, on a single machine, for different combinations of components. The version solely using the policy network does not perform any search. c, Scalability study of MCTS in AlphaGo with search threads and GPUs, using asynchronous search (light blue) or distributed search (dark blue), for 2 s per move.

[Mastering the Game of Go with Deep Neural Networks and Tree Search](#)

Alpha-Go video



The screenshot shows a web browser window with a Bing search results page for the query "alpha-go". The browser's address bar displays the URL: `bing.com/videos/search?q=alpha-go&view=detail&mid=768A0740F2736EEB02C9768A0740F2736EEB02C9&FORM=VIRE`. The search bar contains the text "alpha-go".

The main content area features a video player showing a close-up of a hand placing a white Go stone on a wooden board. Below the video, the title "The computer that mastered Go" is displayed, along with a YouTube logo, "805,000+ views · 1/25/2016 · HD · by nature video".

On the right side, a "Related searches" panel is visible, listing several related terms with small thumbnail images:

- Alphago Transparent Ba...
- Alphago Master
- Lee Sedol Alphago
- Alphago Lose
- Alphago Terminator
- Alphago Live

At the bottom of the page, there are links for "Related videos" and "Feedback".

Game AI: Origins

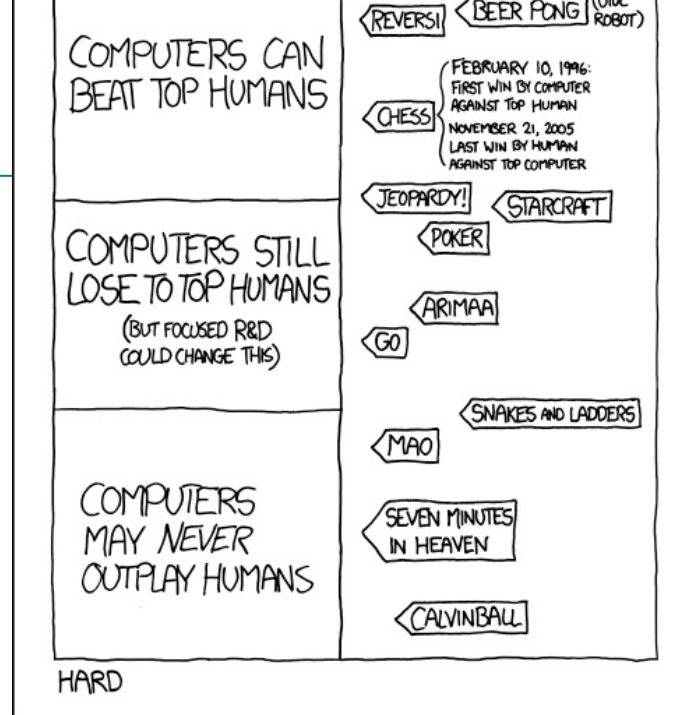
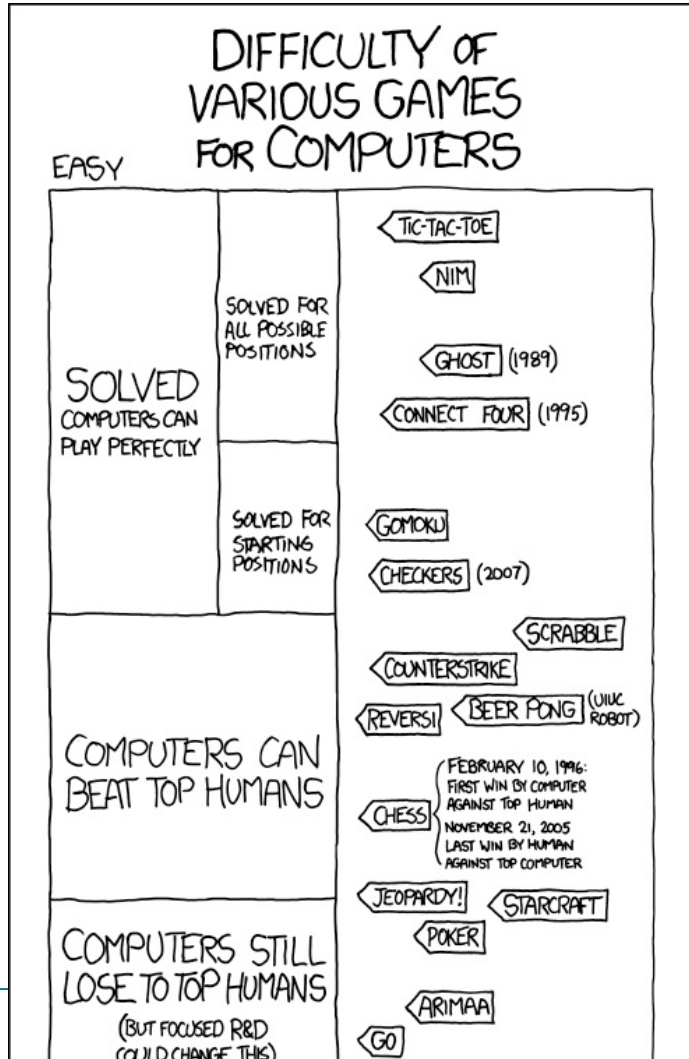
- Minimax algorithm: Ernst Zermelo, 1912
 - Chess playing with evaluation function, quiescence search, selective search:
Claude Shannon, 1949 ([paper](#))
 - Alpha-beta search: John McCarthy, 1956
 - Checkers program that learns its own evaluation function by playing against itself: Arthur Samuel, 1956 ([Rodney Brooks blog post](#))
-

Game AI: State of the art

- Computers are better than humans:
 - **Checkers:** [solved in 2007](#)
 - **Chess:**
 - State-of-the-art search-based systems now better than humans
 - [Deep learning machine teaches itself chess in 72 hours, plays at International Master Level](#) (arXiv, September 2015)
 - Computers are competitive with top human players:
 - **Backgammon:** [TD-Gammon system](#) (1992) used *reinforcement learning* to learn a good evaluation function
 - **Bridge:** top systems use Monte Carlo simulation and alpha-beta search
 - **Go:** computers were not considered competitive until AlphaGo in 2016
-

Game AI: State of the art

- Computers are ~~not~~ competitive with top human players:
 - **Poker**
 - [Heads-up limit hold'em poker is solved](#) (2015)
 - Simplest variant played competitively by humans
 - Smaller number of states than checkers, but partial observability makes it difficult
 - *Essentially weakly solved* = cannot be beaten with statistical significance in a lifetime of playing
 - [CMU's Libratus system beats four of the best human players at no-limit Texas Hold'em poker](#) (2017)
-



<http://xkcd.com/1002/>

<http://xkcd.com/1263/>



calvinball



Calvinball:

- [Play it online](#)
- [Watch an instructional video](#)