



Artificial Intelligence

Chapter 5: Adversarial Search

Andreas Zell

After the Textbook: Artificial Intelligence,
A Modern Approach
by Stuart Russel and Peter Norvig (3rd Edition)

5.1 Games

5.2 Optimal Decisions in Games

5.3 Alpha-Beta Pruning

5.4 Imperfect Real-Time Decisions

5.5 Stochastic Games

5.6 Partially Observable Games

5.7 State-of-the-Art Game Programs

5.8 Alternative Approaches

5.9 Summary



- Chapter 2 introduced **multiagent environments**.
- Each agent must consider the **actions** of the **other agents**.
- Unpredictability of these actions can introduce **contingencies** into the problem-solving process.
- **Games** are **competitive environments**, in which agents' goals are in conflict, giving rise to **adversarial search** problems.

- Mathematical **game theory** views any multiagent environment as a game, in which the impact of each agent on the others is significant, be it **cooperative** or **competitive**.
- In AI, the most common games are:
 - deterministic
 - turn-taking
 - two-player games
 - zero-sum games (result is draw, or a win and a loss)
 - have perfect information (i.e., fully observable)

- Games have many advantages for AI research
- The state is **easy to represent**, there is usually a **small number of possible actions**, whose outcomes are defined by **precise rules**.
- Games are often too hard to solve, making them more interesting
- “easy to learn, but hard to master”

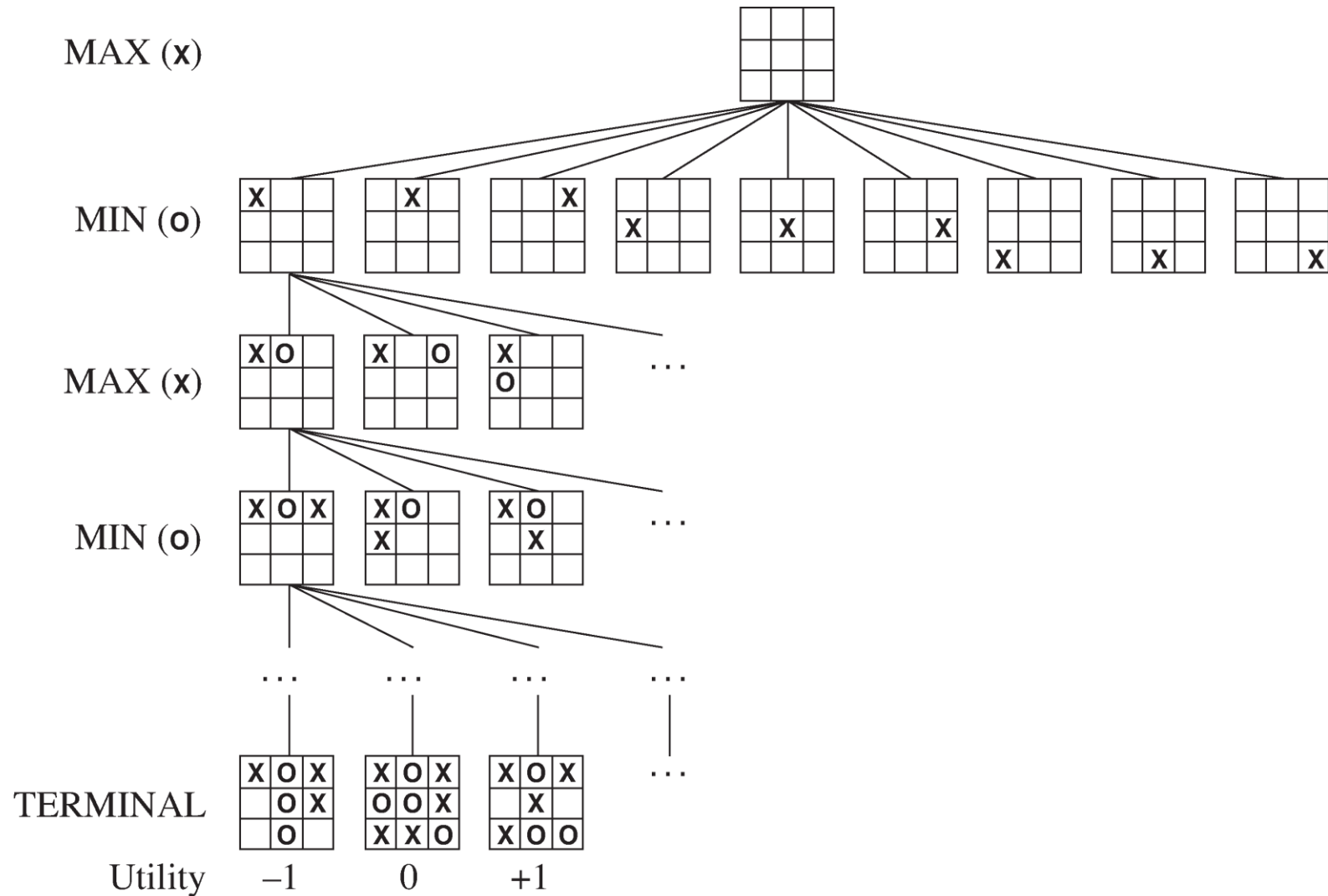
- Example: Chess
- Average branching factor of 35
- Games often go to 50 moves by each player
- Search tree has about 35^{100} or 10^{154} nodes
- Search graph has about 10^{40} distinct nodes
- Calculating the optimal decision is infeasible
- Games often have time constraints, so efficiency plays an important role

- Techniques for finding an optimal move when time is limited:
- **Pruning**: Ignore portions of the search tree that make no difference to the final choice
- **Evaluation functions**: Heuristics to **approximate** the true utility of a state without doing a complete search
- **Imperfect information**: How to find an optimal move if not all information are available?

- Definition of a game as a **search problem**
 - S_0 : **initial state**, game setup
 - $\text{PLAYER}(s)$: player that has the move in state s
 - $\text{ACTIONS}(s)$: legal moves in state s
 - $\text{RESULT}(s,a)$: **transition model**, result of move a
 - $\text{TERMINAL-TEST}(s)$: **terminal test** or goal test. States where the game has ended are called **terminal states**
 - $\text{UTILITY}(s,p)$: **utility function**, defines the numeric value in a terminal state s for player p
 - Example: Loss is -1, draw is 0, and win is +1
 - zero-sum: Sum of utility of both players is zero

- Initial state S_0 , $ACTIONS(s)$ and $RESULT(s,a)$ define the **game tree**
- Nodes are game states
- Edges are moves
- Two players, let's name them MIN and MAX
- Alternating turns, MAX begins
- Example: TicTacToe
 - only $9! = 362,880$ terminal nodes in search tree
 - remember, chess has over 10^{40} nodes

5.1 Games



- In a normal search problem, the optimal solution is a sequence of actions leading to a goal state.
- In adversarial search, every other action is chosen by MIN.
- MAX must find a **strategy** that specifies his move after every possible response to his move by MIN, then his moves after every possible response by MIN to those moves, and so on.
- An **optimal strategy** must lead to outcomes at least as good as any other strategy when playing against an **infallible opponent**.

5.2 Optimal Decisions in Games

- Similar to AND-OR search algorithm
- MAX is OR, MIN is AND
- Given a game tree, the optimal strategy can be determined from the **minmax value** of each node
- It is assumed that both players play **optimally**

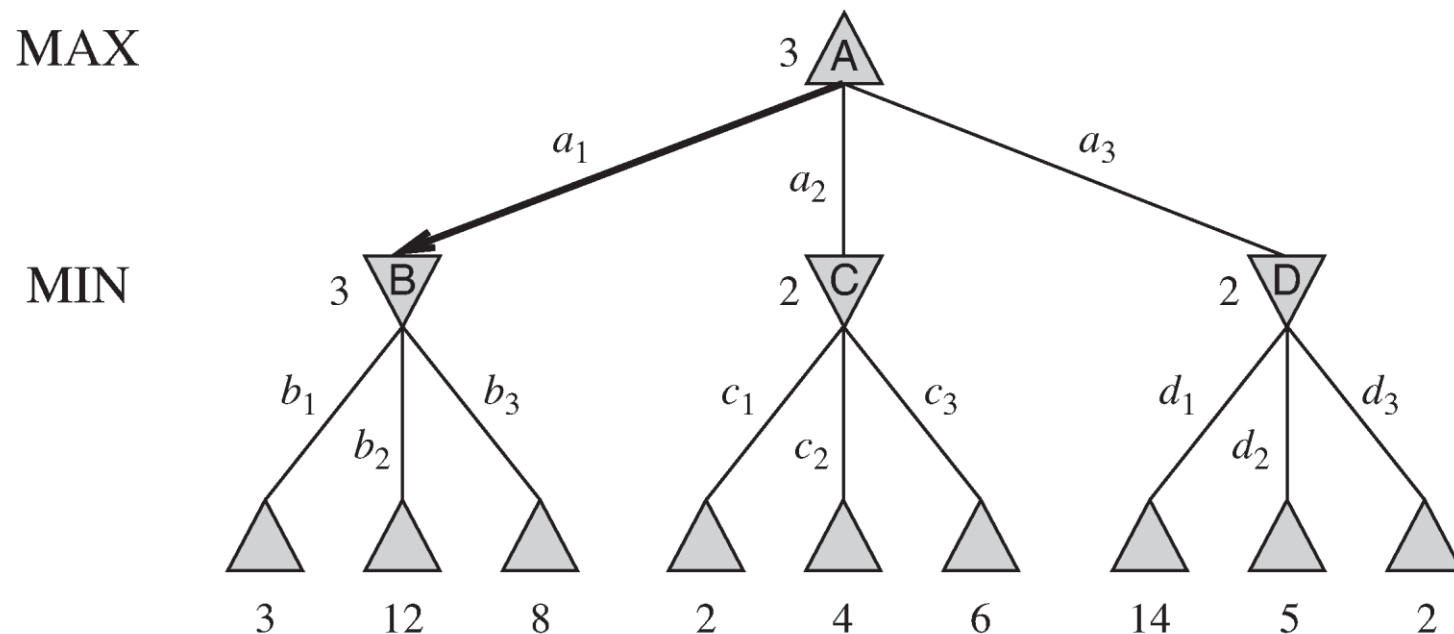
$MINIMAX(s) =$

$$\begin{cases} UTILITY(s) & \text{if } TERMINAL-TEST(s) \\ \max_{a \in Actions(s)} MINIMAX(RESULT(s, a)) & \text{if } PLAYER(s) = MAX \\ \min_{a \in Actions(s)} MINIMAX(RESULT(s, a)) & \text{if } PLAYER(s) = MIN \end{cases}$$

5.2 Optimal Decisions in Games



- Example game tree: MAX always chooses the move that maximizes the minimax value (a_1), while MIN selects the one that minimizes it (b_1)



5.2 Optimal Decisions in Games



- In **game parlance**, the previous game tree is considered one move deep
- a move consists of two **half-moves**
- a half-move is also called a **ply**

- The optimal strategy of MAX **assumes** that MIN also **plays optimally**
- If MIN does not play optimally, MAX will do better, **at least equally good**
- There might be other strategies that are better if MIN plays suboptimally, but those will always be worse against optimal opponents
- The **minimax algorithm** computes the optimal decision using the minimax values by **traversing the game tree** all the way **down to the leaves**

5.2 Optimal Decisions in Games

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

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**
 $v \leftarrow \text{MAX}(v, \text{MIN-VALUE}(\text{RESULT}(s, a)))$
 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**
 $v \leftarrow \text{MIN}(v, \text{MAX-VALUE}(\text{RESULT}(s, a)))$
 return *v*

- The minimax algorithm performs a **depth-first search**
- For a tree with maximum depth m and b possible actions at each state
 - time complexity is $O(b^m)$
 - space complexity is $O(bm)$ if all actions are generated at once, $O(m)$ if actions are generated one at a time
- For most real games, this **time complexity** is **totally impractical**

- Minimax can be extended for games that have **more than two players**
- Instead of a utility value, there is now a **utility vector**, for example $\langle v_a, v_b, v_c \rangle$ for three players
- Instead of MAX choosing the action maximizing the minimax value and MIN minimizing it, each player chooses the action maximizing his utility vector component
- Instead of a value being backed up the tree, a vector is backed up

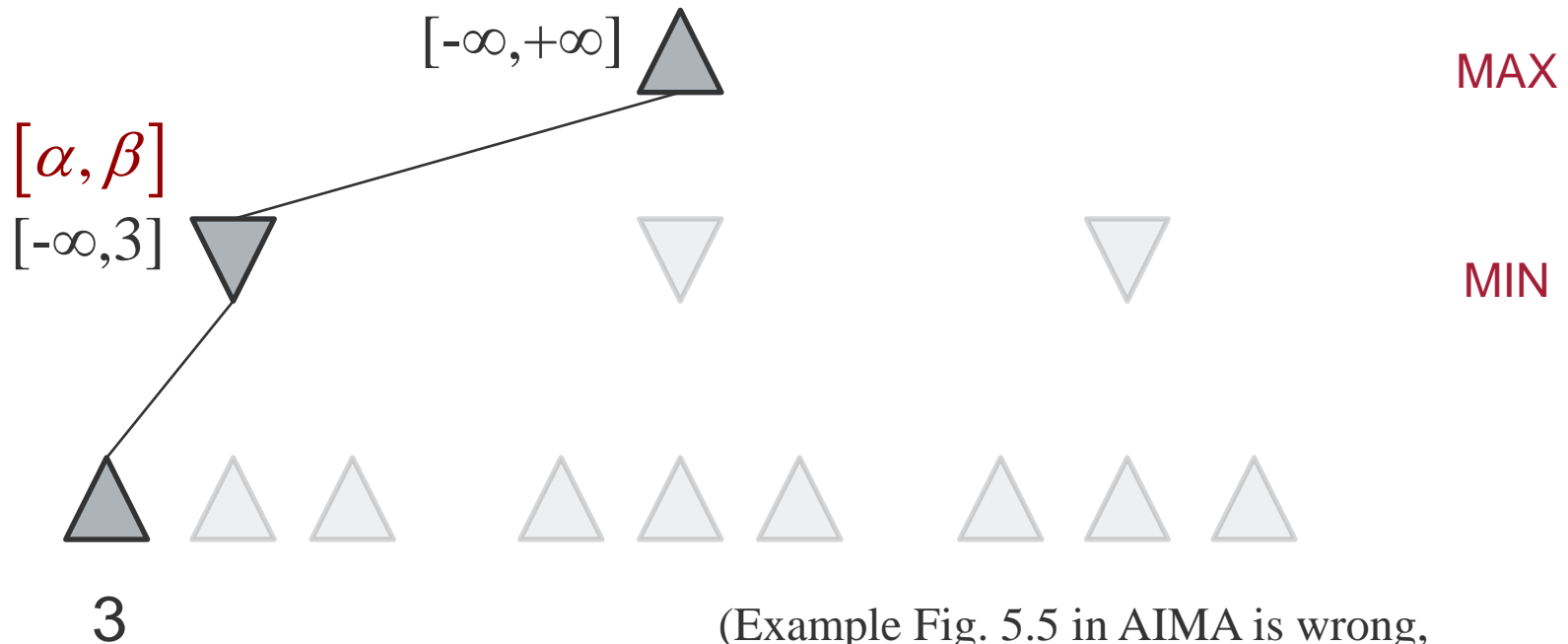
5.3 Alpha-Beta Pruning

- Obviously, the exponential time complexity of the minimax algorithm is a problem
- Approach: **Prune** away those parts of the tree that we don't need to examine, because they **wouldn't change the optimal decision**
- Examine the previous game tree again, this time applying **alpha-beta pruning**

5.3 Alpha-Beta Pruning

- α : value of the best (highest-value) choice found so far any choice point along the path for MAX
- β : value of the best (lowest-value) choice found so far any choice point along the path for MIN
- If an action of MAX leads to a state where MIN has at least one action that leads to a worse result than an already evaluated action of MAX, we don't need to examine it further
- Same holds vice versa

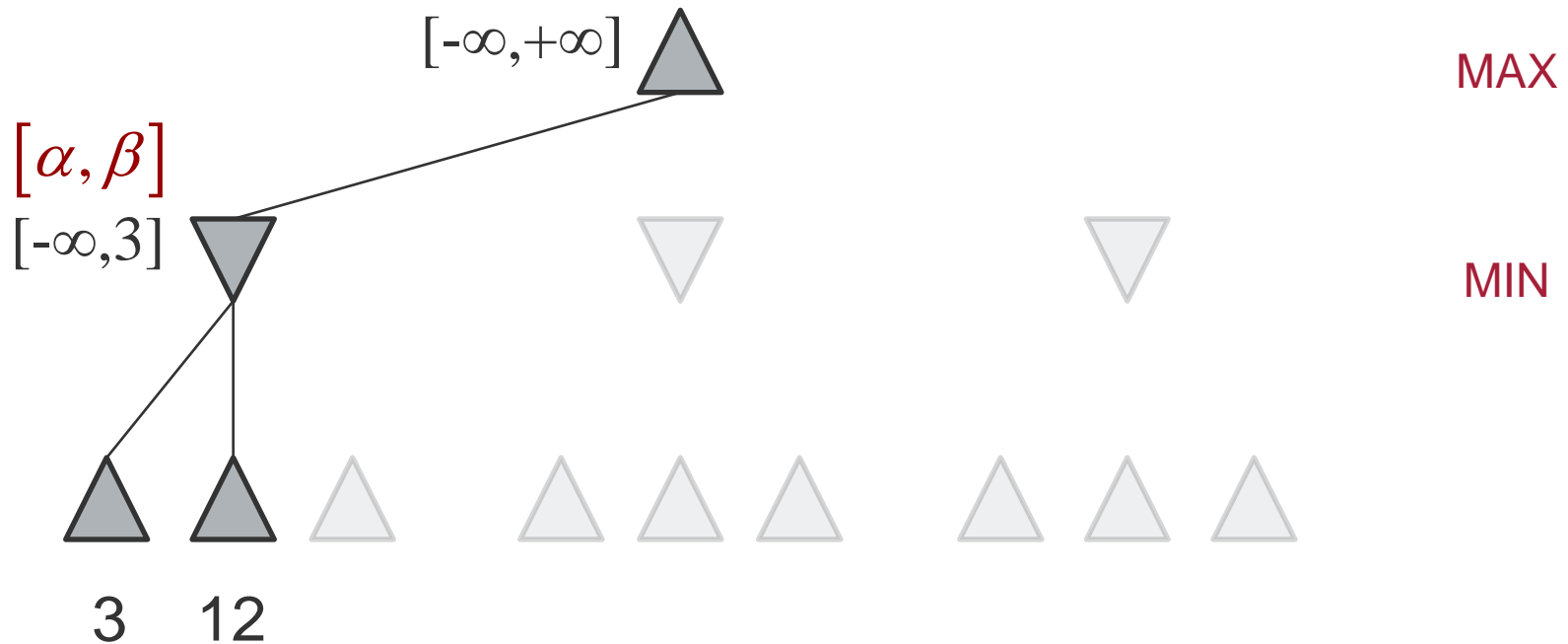
5.3 Alpha-Beta Pruning



(Example Fig. 5.5 in AIMA is wrong,
is not the algorithm in Fig. 5.7)

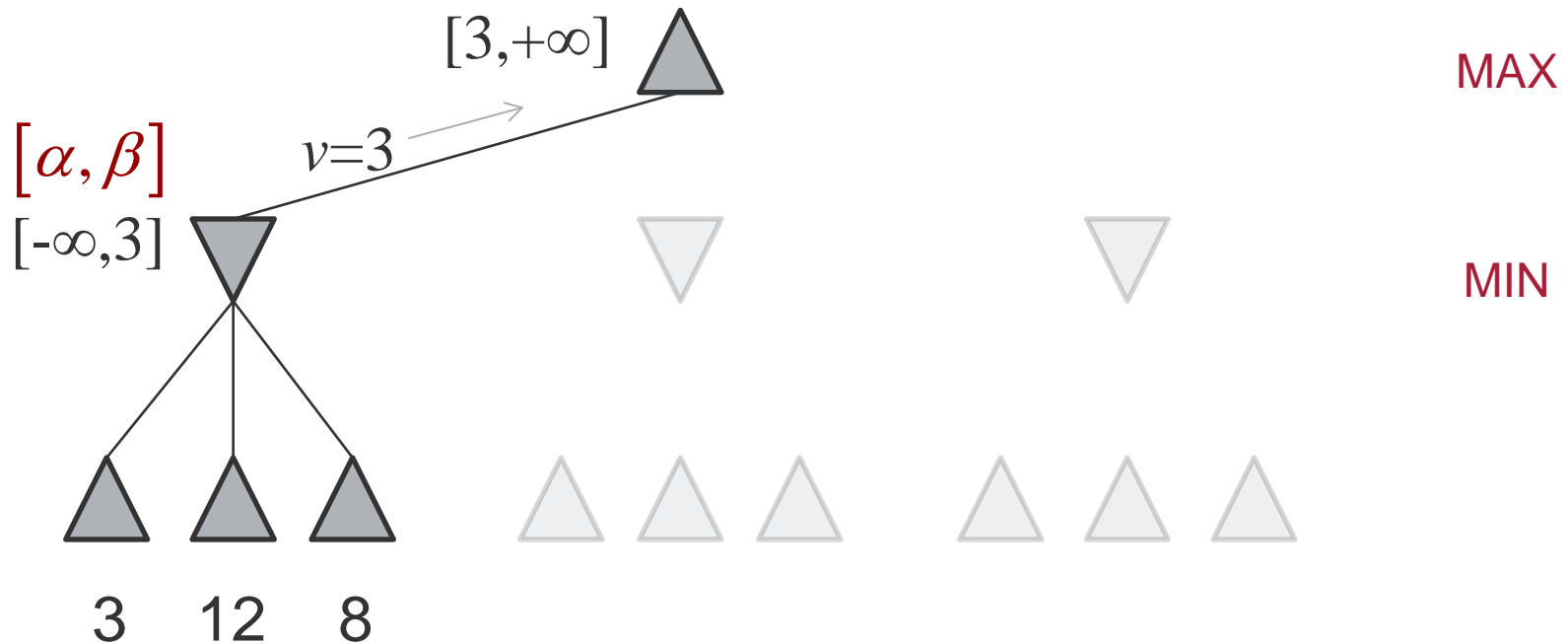
- Best choice for MIN so far, set $\beta = 3$
- Best choice for MAX can't be known yet

5.3 Alpha-Beta Pruning



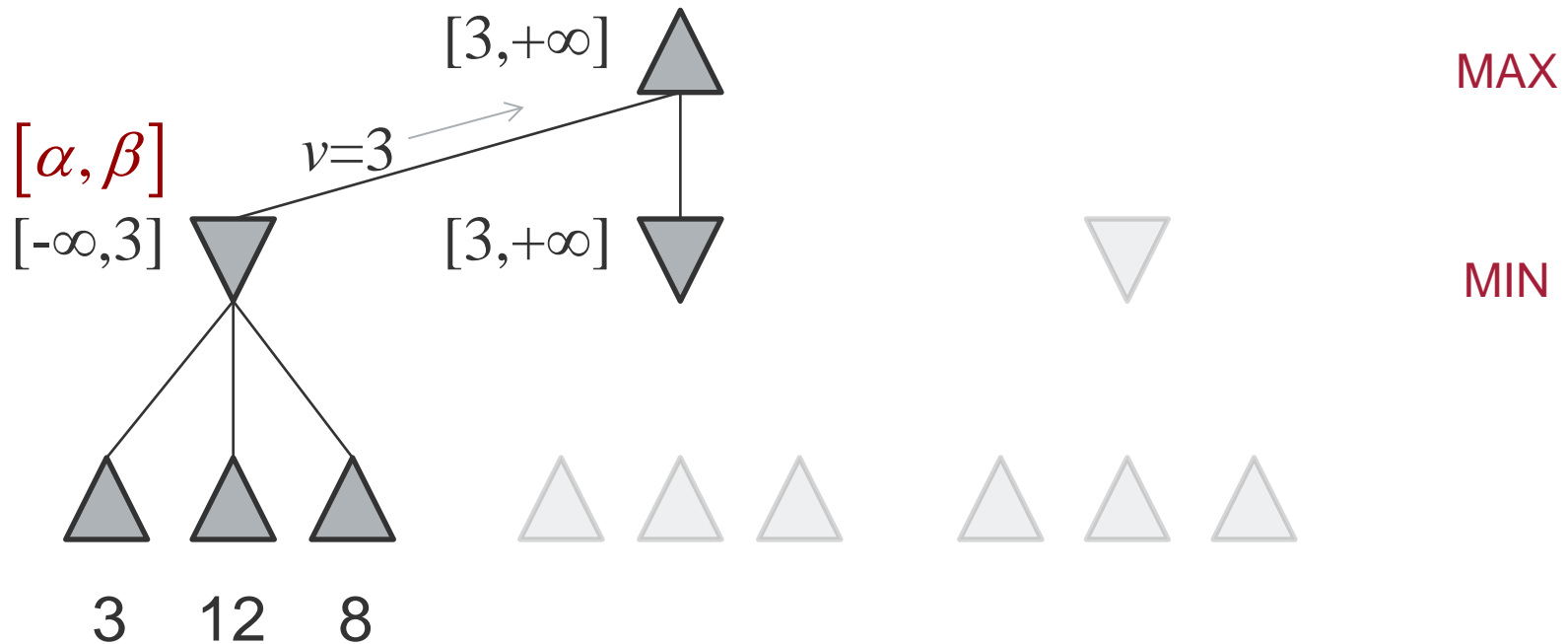
- Next action has higher value than 3, so MIN won't choose it

5.3 Alpha-Beta Pruning



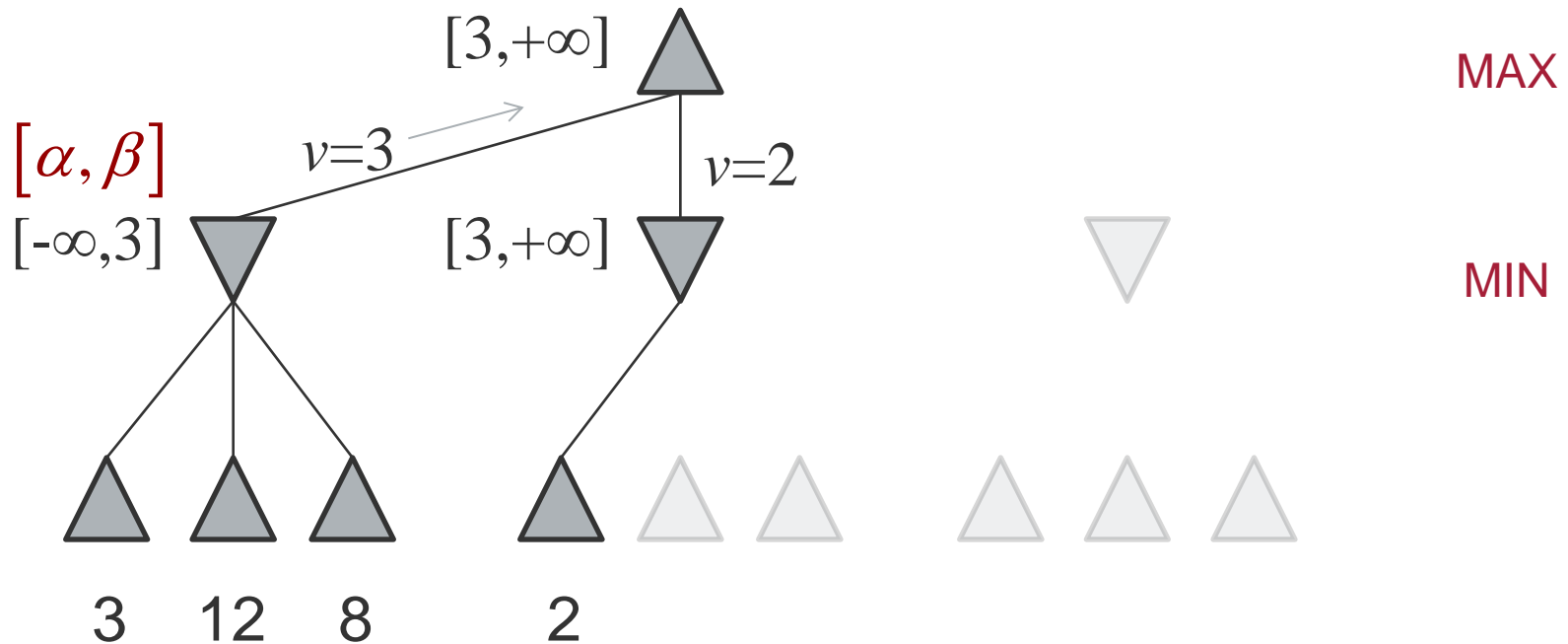
- Still no better choice for MIN, so back up minimal value and set as $\alpha = 3$ on MAX

5.3 Alpha-Beta Pruning



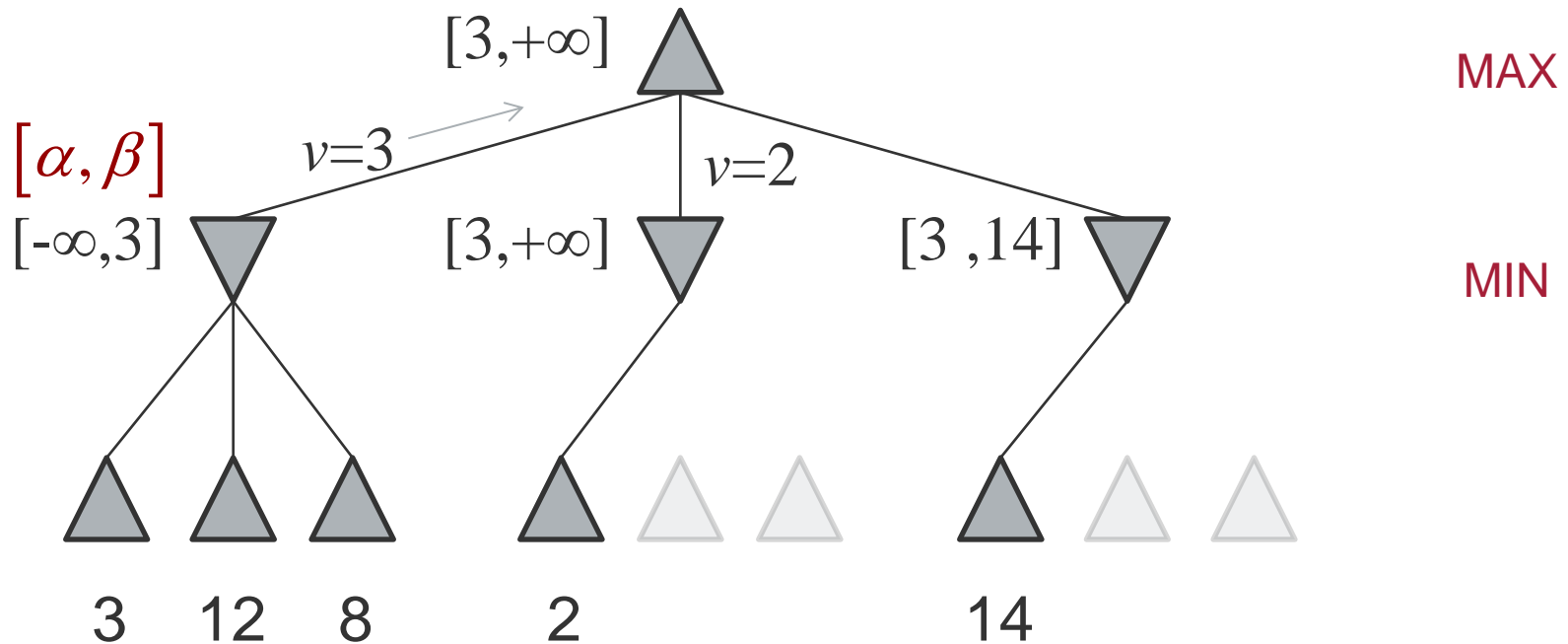
- Backed up α value, best value for MAX so far
- Go down to next node, pass α and β along

5.3 Alpha-Beta Pruning



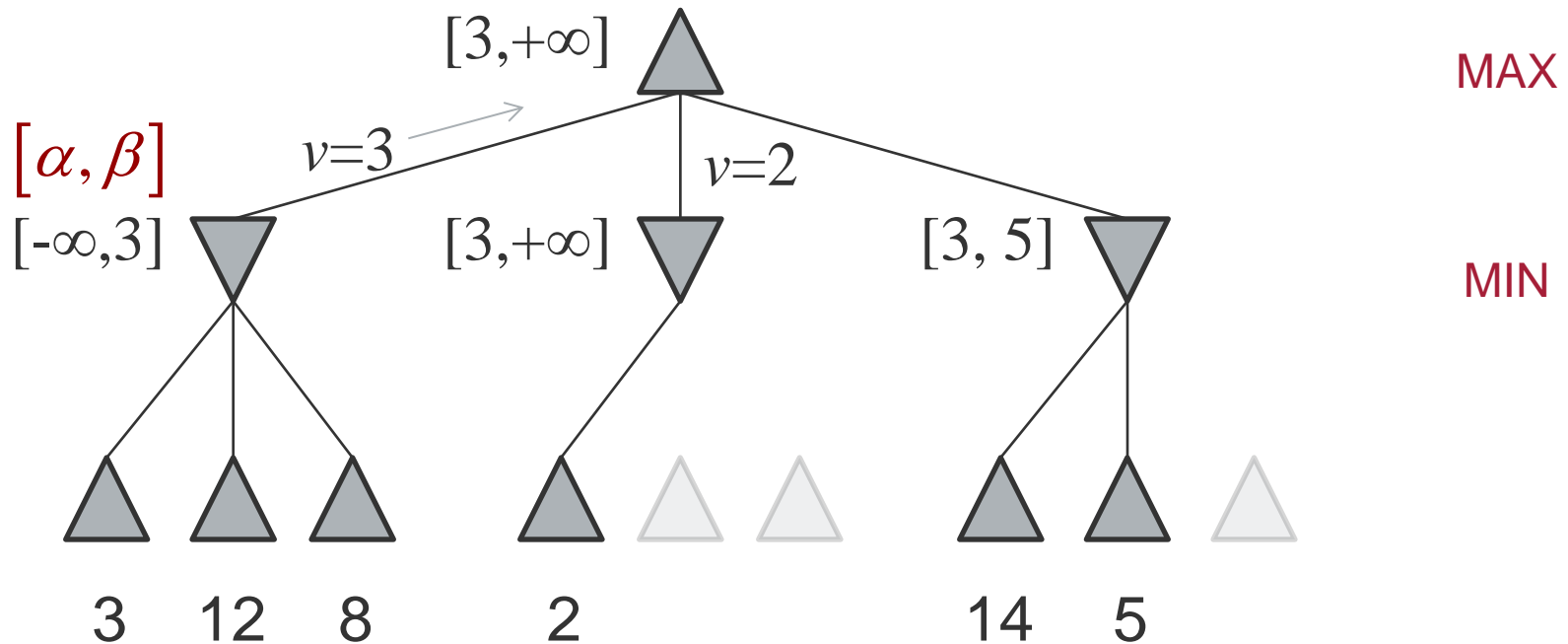
- MIN will choose at most a value of 2, but MAX already knows a better move ($2 < \alpha$). Prune here

5.3 Alpha-Beta Pruning



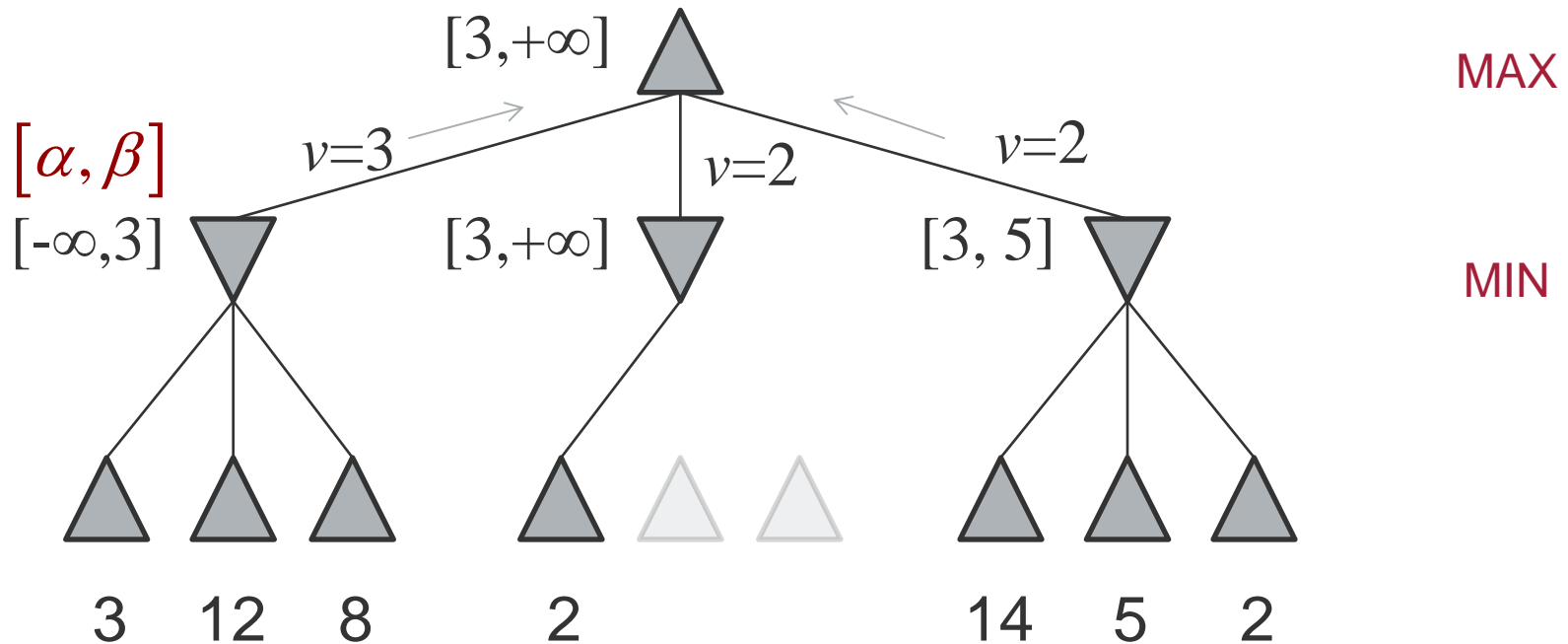
- Go down to next node, pass α and β along
- Leaf $14 \geq \alpha$, update β

5.3 Alpha-Beta Pruning



- Leaf 5 $\geq \alpha$, update β

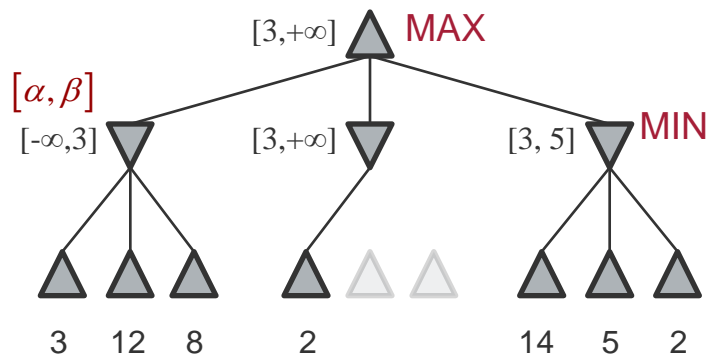
5.3 Alpha-Beta Pruning



- MIN would choose action leading to 2
- $2 < \alpha$, so no β update, pass up 2 to MAX
- Best move for MAX in root known now

5.3 Alpha-Beta Pruning

- Algorithm for minimax search with alpha-beta pruning



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

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

- Last subtree in example could have been skipped earlier, if action with result 2 would have been the first examined action
- Is it possible to **order** the examined **moves** in a way to allow for **earlier alpha-beta pruning**?
- If **best moves** are examined **first**, alpha-beta only needs to examine $O(b^{m/2})$ nodes, instead of $O(b^m)$ for minimax
- Branching factor b becomes \sqrt{b}
- If successors are examined in **random order**, alpha-beta needs to examine $O(b^{3m/4})$ nodes

- If **perfect ordering** was possible, we would already know the best move without searching
- But **heuristics** can be used to select moves first that are probably better (e.g. attacking moves)
- Best moves are called **killer moves** and are selected by the killer move heuristic
- States may occur multiple times, so it is possible to create a **transposition table**, storing game states that already have been evaluated
- But: Storing all states is not practical, because there are too many

5.4 Imperfect Real-Time Decisions

- Alpha-beta prunes large portions of game tree away, but still has to go down to terminal nodes
- Approach: Apply a **heuristic evaluation function** on nonterminal nodes to limit search depth
- A **cutoff-test** decides when to apply heuristic

$$H\text{-}MINIMAX(s, d) = \begin{cases} EVAL(s) & \text{if } CUTOFF\text{-}TEST(s, d) \\ \max_{a \in Actions(s)} H\text{-}MINIMAX(RESULT(s, a), d + 1) & \text{if } PLAYER(s) = MAX \\ \min_{a \in Actions(s)} H\text{-}MINIMAX(RESULT(s, a), d + 1) & \text{if } PLAYER(s) = MIN \end{cases}$$

- Heuristic evaluation functions in Chapter 3 return the **estimated distance** to the goal
- Heuristic evaluation functions here return the **estimated utility** of a **game state**
- Performance of a game-playing program highly depends on the quality of its evaluation function
- **Incorrect evaluation** of a game state might direct the search into a **wrong direction**

Requirements

1. The evaluation function should order terminal nodes in the same way as the utility function
 - Win > Draw > Loss
2. The computation of this function must be fast
3. The result of the evaluation function should strongly correlate with the chances of winning for nonterminal states
 - Here, “chance” means the uncertainty due to computational limitations, even if the game itself is deterministic

- One way to define an evaluation function is to group states into categories by their features
- Features can be for example (in chess):
 - number of pawns, queens, knights for black and white
 - number of threatened pieces for black and white
- A large knowledge-base could be used to calculate expected values, e.g.
 - In 80% of previous games where white has three pawns and black has two pawns, white won the game, in 10% it was a draw and in 10% white lost
$$\text{expected value} = (0.8 \times 1) + (0.1 \times 0.5) + (0.1 \times 0)$$



- In practice, this requires the analysis of too many categories and too much experience
- The evaluation function doesn't need to return expected values, as long as it orders the states correctly
- Chess often uses approximate material values
 - a pawn is worth 1, a knight or bishop is worth 3, etc.
 - features like “king safety” might be worth 0.5
- Values are added for each player

- Here, the evaluation function becomes a weighted linear function

$$\text{EVAL}(s) = w_1 f_1(s) + w_2 f_2(s) + \cdots + w_n f_n(s) = \sum_{i=1}^n w_i f_i$$

where w_i is a weight and f_i is a feature.

- Strong assumption here: The value of each piece is independent of the others
- In chess, bishops are more powerful in the endgame (move number is high or number of remaining pieces is low)
- So, in an endgame, bishops are worth more

- Deciding when to cut off the search is important
- Simple approach: Stop at a fixed depth
- Better: Use iterative deepening
- The deepest completed search at the end of the time limit specifies the selected move
- Bonus advantage: Previous depth searches can give additional information for move ordering
- But, there still might be situations, where an important move is just one move deeper than calculated

- **Quiescence search:** The evaluation function is only applied to quiescent positions, i.e. that are unlikely to dramatically change in the near future
- Positions with favorable captures are expanded until quiescent positions are reached
- There still is the horizon effect that some important moves are “pushed over the horizon” and are not included in the evaluation
- There are approaches to mitigate this effect, but we won’t discuss them here

- A more human-like approach is forward pruning
- In chess, humans usually consider only a few moves in each position, pruning the others without further consideration
- Beam search only considers the n best moves in each position (according to the evaluation)
- This might prune away the subtree with the best overall move

- For (phases of) games with only a limited number of states, a lookup table can be feasible
- Opening tables in chess are very successful (every game has the same starting position, so there is a large data base for statistics)
- In endgames, the lookup table can be made into a graph, where for each position the best move is already known
- E.g. King-bishop-and-knight-versus-king (KBNK) has $462 \times 62 \times 61 \times 2 = 3,494,568$ positions

- There are also many stochastic games with random elements, like throwing a dice
- Example: Backgammon
- Each player knows his own legal moves after throwing the dice, but he doesn't know what will be the legal moves of the opponent
- Standard game tree can't be constructed
- Add chance nodes to the game tree

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



- How to compute the best move here?
- Calculate the **expected value** of each move
- That is, the **weighted average** over all possible outcomes of chance nodes
- Generalization of the **minimax value** to the **expected minimax value**

- additional player CHANCE
- r represents a chance event (e.g. dice roll)
- $RESULT(s, r)$ is the same state as s , with the additional fact that the chance event was r

$$EXPECTIMINIMAX(s) =$$

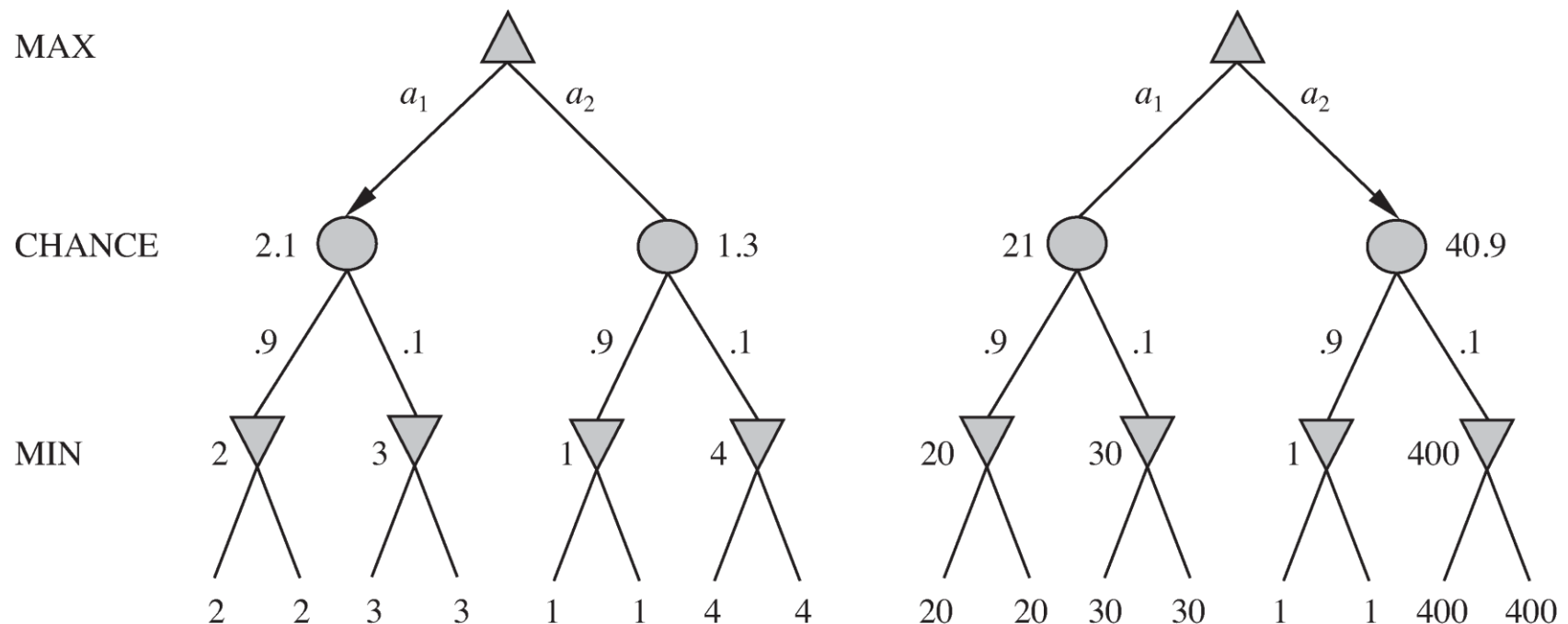
$$\begin{cases} UTILITY(s) & \text{if } TERMINAL-TEST(s) \\ \max_a EXPECTIMINIMAX(RESULT(s, a)) & \text{if } PLAYER(s) = MAX \\ \min_a EXPECTIMINIMAX(RESULT(s, a)) & \text{if } PLAYER(s) = MIN \\ \sum_r P(r) EXPECTIMINIMAX(RESULT(s, r)) & \text{if } PLAYER(s) = CHANCE \end{cases}$$

- Obvious approximation is to cut the search off at a certain point and apply an evaluation function
- For deterministic games, the evaluation function only needed to order the leaves correctly
- For stochastic games, this can lead to problems
- Solution is a stronger requirement:
- Evaluation function must be a positive linear transformation of the probability of winning

5.5 Stochastic Games



- both evaluation functions order leaves the same
- but they lead to different expected values at the chance nodes



- Time complexity for normal minimax is $O(b^m)$
 - b is the branching factor
 - m is the maximum depth of the game tree
- With addition of chance events, it is $O(b^m n^m)$
 - n is the number of distinct chance events
- In backgammon, the maximal feasible search depth was three plies
- Alpha-beta pruning ignores future developments that won't happen given best play
- With chance events, this becomes harder

- For alpha-beta pruning, we need find upper and lower bounds for nodes
- Is it possible to find these bounds for chance nodes without examining all children?
- Yes, if there are bounds on the possible values of the utility function, the expected values are also bounded
- Alternate approach: Monte Carlo simulation
- Simulate thousands of games from a position with random dice rolls and use the resulting win percentage as an evaluation function

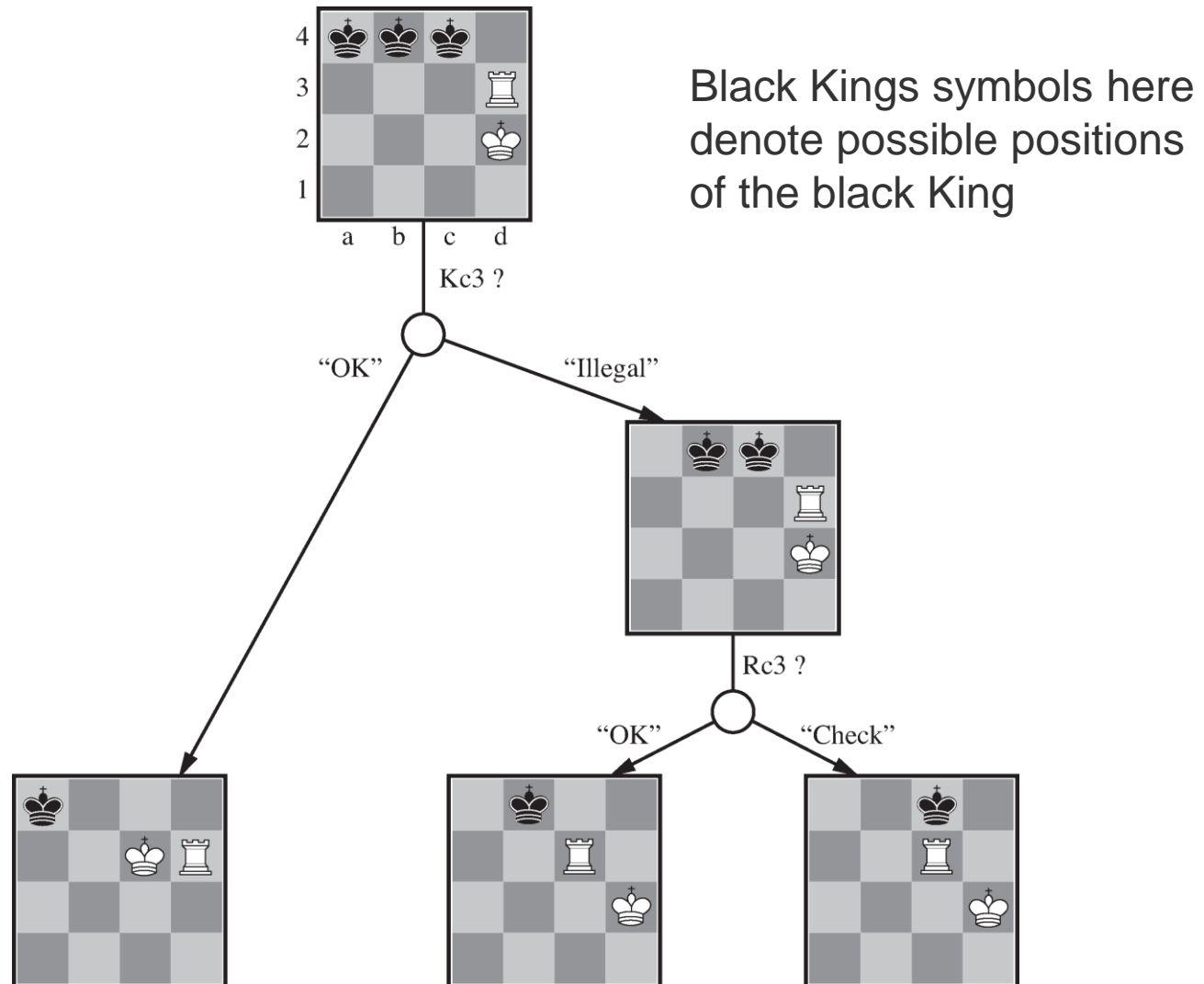


- So far, all games were fully observable
- But some games have the “fog of war”
- Existence and positions of enemy units is unknown until revealed by direct contact
- This adds new aspects to games
- The purpose of some actions can be to gather information only
- Successful strategies may involve bluffing

- “Kriegspiel”:
- A partially observable variant of chess
- Moves and pieces of the opponent are hidden
- Players propose an action to a referee
 - If the action is illegal, the referee announces “illegal”
 - The player keeps proposing moves until a legal one is found, which is then executed
 - The referee then announces things like
 - “Capture on square X”
 - “Check by [direction]”
 - It is the other player’s turn to move

- Players have to act upon **belief states** (see 4.4)
- The initial belief state is a singleton, because the starting positions are known
- With each move, the belief state contains more positions
- A winning strategy, or **guaranteed checkmate**, is one that, for each possible **percept sequence**, leads to an actual checkmate for every possible board state in the current belief state, regardless of how the opponent moves

5.6 Partially Observable Games



- With this definition of a strategy, the belief state of the opponent is **irrelevant**, which greatly simplifies computation
- There is also a **probabilistic checkmate**:
A strategy that contains randomized moves and will still lead to a checkmate, eventually
- A strategy that will only work for some of the positions in a belief state is called an **accidental checkmate**
- For that, it can be helpful to consider whether some positions are more likely than others



- Many card games also provide examples of **stochastic partial observability**
- E.g., cards are dealt **randomly at the beginning** of the game and a player's hand is not visible to the others (like in Poker, Bridge, Skat, ...)
- All random effects occurred at the beginning
- Solve each possible state as if it were a fully observable game, then choose the move with the best outcome **averaged over all deals**

- With an exact minimax

$$\arg \max_a \sum_a P(s) MINIMAX (RESULT(s, a))$$

- If computationally infeasible, run *H-MINIMAX*
- If the number of possible deals is too large, one can use a **Monte Carlo approximation**

$$\arg \max_a \frac{1}{N} \sum_{i=1}^N P(s) MINIMAX (RESULT(s_i, a))$$

- But both approaches assume a fully observable game after the first move

- Those strategies will never select actions that **gather information**
- They never choose moves that **hide information**, or that provide information to partners
- They will also never **bluff**, because they assume that the opponent can see all cards



- Chess: Deep Blue (IBM, 1997)
 - 30 IBM RS/6000 processors for alpha-beta search
 - 480 custom VLSI chess processors for move generation, move ordering and leaf node evaluation
 - up to 30 billion positions per move, search depth 14
 - evaluation function with 8000 features
 - opening book of 4000 positions
 - database of 700,000 grandmaster games
 - large endgame database of solved positions containing all positions with 5 pieces and many with 6

- Chess on standard PCs
 - Algorithmic improvements allow standard PCs to win World Computer Chess Championships
 - Effective pruning heuristics reduce the effective branching factor to less than 3
 - Most important one is the **null move heuristic**, where the opponent moves twice in the search, generating a good lower bound on the current position
 - There is also futility pruning, helping to decide in advance which moves will cause a beta cutoff in successor nodes
- Top chess programs beat all human contenders

- Checkers: Chinook
 - It defeated the human world champion in 1990
 - Since 2007, it can play perfectly by using alpha-beta search and an endgame database with 39 trillion positions
- Backgammon: TD-Gammon
 - The evaluation function was improved using reinforcement learning and neural networks
 - After playing more than 1 trillion training games against itself, it is competitive with top human players
 - Some of its opening moves radically altered the previous general game play knowledge

5.7 State-of-the-Art Game Programs



- Go: AlphaGo
 - AlphaGo, developed by Google DeepMind (London) defeated the best professional Go player Lee Sedol (9-dan, Korea) 4:1 in 5 matches, March 2016



5.7 State-of-the-Art Game Programs

- AlphaGo used 1,920 CPUs and 280 GPUs in the match against Lee Sedol
- It used **Monte Carlo tree search**, guided by a value network and a policy network, both implemented using **deep neural network technology**
- AlphaGo was initially trained on a database of around 30 million moves of expert human players
- Then it was trained further by playing large numbers of games against other instances of itself, using **reinforcement learning** to improve its play
- <https://de.wikipedia.org/wiki/AlphaGo>

5.7 State-of-the-Art Game Programs

- TCEC (Top Chess Engine Championship) 2017
 - 2nd round

Rang	Name	Punkte	Elo
1	 Houdini 6.02	18½/28	3184
2	 Komodo 1959.00	18½/28	3232
3	 Stockfish 051117	18/28	3228
4	 Fire 6.2	15/28	3112
5	 Chiron 251017	11½/28	3013
6	 Ginkgo 2.01	10½/28	3052
8 Spieler		Elo-Durchschnitt: 3127	

Top 100 Players December 2017

Rank	Name	Title	Country	Rating
1	<u>Carlsen, Magnus</u>	g	NOR	2837
2	<u>Aronian, Levon</u>	g	ARM	2805
3	<u>Mamedyarov, Shakhriyar</u>	g	AZE	2799
4	<u>Caruana, Fabiano</u>	g	USA	2799
5	<u>Vachier-Lagrave, Maxime</u>	g	FRA	2789
6	<u>So, Wesley</u>	g	USA	2788
7	<u>Kramnik, Vladimir</u>	g	RUS	2787
8	<u>Anand, Viswanathan</u>	g	IND	2782
9	<u>Nakamura, Hikaru</u>	g	USA	2781
10	<u>Ding, Liren</u>	g	CHN	2777

- **Superfinal:** Houdini 6.02 vs. Komodo 1959.00
 - 100 games played between them (120 min. each)
 - Houdini won 53:47 (15 wins, 76 draws, 9 losses)

Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm

David Silver,^{1*} Thomas Hubert,^{1*} Julian Schrittwieser,^{1*}
Ioannis Antonoglou,¹ Matthew Lai,¹ Arthur Guez,¹ Marc Lanctot,¹
Laurent Sifre,¹ Dhharshan Kumaran,¹ Thore Graepel,¹
Timothy Lillicrap,¹ Karen Simonyan,¹ Demis Hassabis¹

¹DeepMind, 6 Pancras Square, London N1C 4AG.

Abstract: The game of chess is the most widely-studied domain in the history of artificial intelligence. The strongest programs are based on a combination of sophisticated search techniques, domain-specific adaptations, and handcrafted evaluation functions that have been refined by human experts over several decades. In contrast, the **AlphaGo Zero** program recently achieved superhuman performance in the game of Go, by tabula rasa reinforcement learning from games of self-play. In this paper, we generalise this approach into a single **AlphaZero** algorithm that can achieve, tabula rasa, superhuman performance in many challenging domains. Starting from random play, and given no domain knowledge except the game rules, AlphaZero achieved within 24 hours a superhuman level of play in the games of chess and shogi (Japanese chess) as well as Go, and convincingly defeated a world-champion program in each case.

arXiv:1712.01815v1 [cs.AI] 5 Dec 2017

- The **AlphaZero** algorithm is a more generic version of the **AlphaGo Zero** algorithm
- It uses a deep neural network and a tabula rasa reinforcement learning algorithm
- This neural network $(\mathbf{p}, v) = f_{\theta}(s)$ takes the board position s as an input and outputs a vector of move probabilities \mathbf{p} with components $p_a = \Pr(a|s)$ for each action a , and a scalar value v estimating the expected outcome z from position s , so $v \approx E[z|s]$.
- AlphaZero learns these move probabilities and value estimates entirely from self play.
- Instead of an alpha-beta search, AlphaZero uses a **Monte-Carlo tree search (MCTS)** algorithm.

- Each search consists of a series of self-play games that traverse a tree from root to a leaf.
- Each simulation proceeds by selecting in each state s a move a with low visit count, high move probability and high value according to the current neural network $f_\theta(s)$.
- The search returns a vector π representing a probability distribution over moves, either proportionally or greedily with respect to the visit counts at the root state.
- The parameters θ of the deep neural network are trained by self-play reinforcement learning, starting from randomly initialized parameters θ .
- Games are played by selecting moves for both players by MCTS, $a_t \sim \pi_t$.

- At the end of the game, the terminal position s_T is scored according to the rules of the game to compute the game outcome z : -1 for a loss, 0 for a draw, and +1 for a win.
- The neural network parameters θ are updated so as to minimize the error between the predicted outcome v_t and the game outcome z , and to maximize the similarity of the policy vector p_t to the search probabilities π_t .
- Specifically, the parameters θ are adjusted by gradient descent on a loss function l that sums over mean-squared error and cross-entropy losses respectively, where c is a parameter controlling the level of L_2 weight regularization.
$$l = (z - v)^2 - \pi^T \log p + c \|\theta\|^2$$
- The updated parameters are used in subsequent games.

5.8 Google AlphaZero

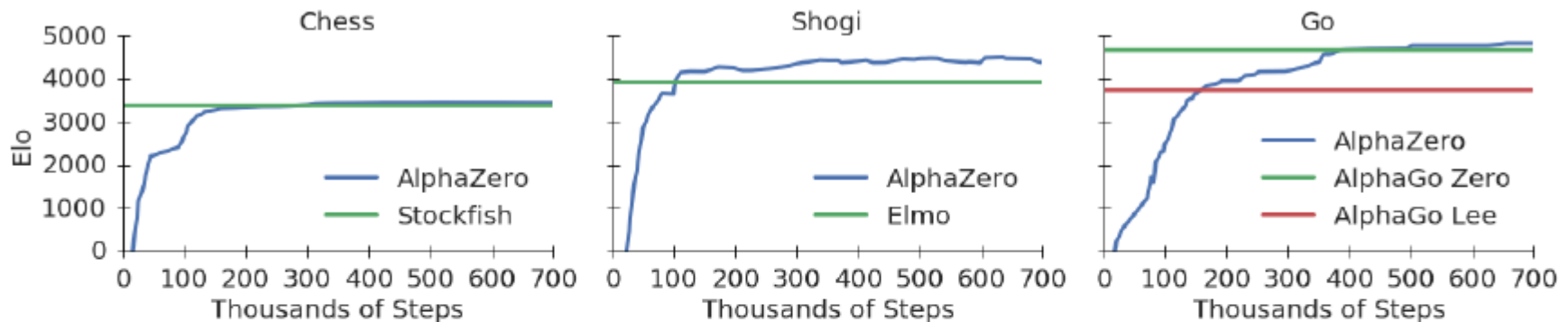


Figure 1: Training *AlphaZero* for 700,000 steps. Elo ratings were computed from evaluation games between different players when given one second per move. **a** Performance of *AlphaZero* in chess, compared to 2016 TCEC world-champion program *Stockfish*. **b** Performance of *AlphaZero* in shogi, compared to 2017 CSA world-champion program *Elmo*. **c** Performance of *AlphaZero* in Go, compared to *AlphaGo Lee* and *AlphaGo Zero* (20 block / 3 day) (29).

They trained a separate instance of AlphaZero for each game. Training used 5,000 first-generation TPUs (Tensor processing units) to generate self-play games and 64 second-generation TPUs (https://de.wikipedia.org/wiki/Tensor_Processing_Unit) to train the neural networks.

5.8 Google AlphaZero

Game	White	Black	Win	Draw	Loss
Chess	<i>AlphaZero</i>	<i>Stockfish</i>	25	25	0
	<i>Stockfish</i>	<i>AlphaZero</i>	3	47	0
Shogi	<i>AlphaZero</i>	<i>Elmo</i>	43	2	5
	<i>Elmo</i>	<i>AlphaZero</i>	47	0	3
Go	<i>AlphaZero</i>	<i>AG0 3-day</i>	31	–	19
	<i>AG0 3-day</i>	<i>AlphaZero</i>	29	–	21

Table 1: Tournament evaluation of *AlphaZero* in chess, shogi, and Go, as games won, drawn or lost from *AlphaZero*'s perspective, in 100 game matches against *Stockfish*, *Elmo*, and the previously published *AlphaGo Zero* after 3 days of training. Each program was given 1 minute of thinking time per move.

5.8 Google AlphaZero

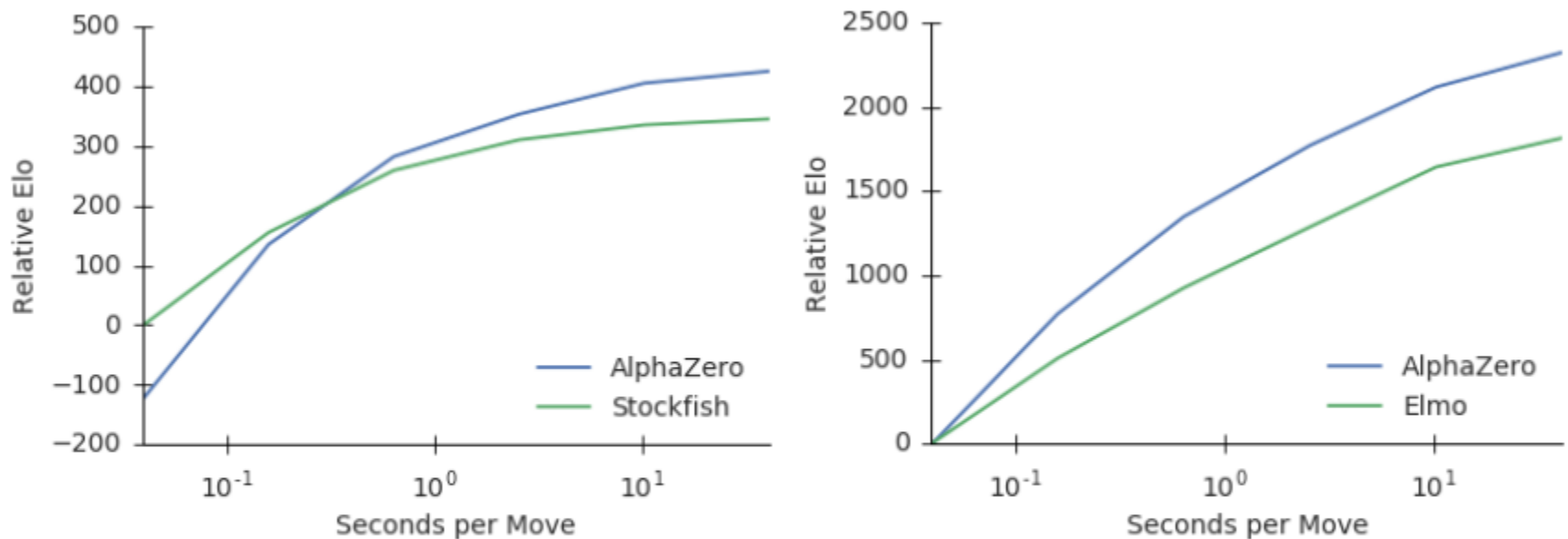


Figure 2: Scalability of *AlphaZero* with thinking time, measured on an Elo scale. **a** Performance of *AlphaZero* and *Stockfish* in chess, plotted against thinking time per move. **b** Performance of *AlphaZero* and *Elmo* in shogi, plotted against thinking time per move.

- A game can be defined by
 - the **initial state**,
 - the **legal actions** in each state,
 - a **terminal test**,
 - and a **utility function** that applies to terminal states.
- In two-player zero-sum games with **perfect information**, minimax can select the **optimal** move via a **depth-first search** of the game tree
- **Alpha-beta** search computes the same optimal moves, but achieves **much greater efficiency**

- Usually, it is **not feasible** to consider the whole game tree (even with alpha-beta pruning)
- We need to **cut off** evaluation and apply a heuristic **evaluation function** to states
- Precomputed **opening** and **endgame tables** are often used in game programs
- **Games of chance** can be handled by extending the minimax algorithm to include **chance nodes**
- Optimal play in games of **imperfect information** requires **reasoning about** current and future **belief states**

- Computer programs have bested champion human players at games such as **chess**, **checkers**, **Othello**, and in **Go** (AlphaGo, 2016).
- **AlphaZero** (2017) learns to play Chess, Shogi and Go only by self-play with deep neural networks, MCTS and reinforcement learning.
- Humans retain the edge in modern real-time computer games with extremely many move options and very strategic play.
- But even this might change soon, as DeepMind is trying to learn to play StarCraft II (Blizzard)