

Playing Chess with RL

Michał Kurek, Joel Runevic, Vladimir Petrov
Harvard College

1 Introduction

1.1 Chess as an Algorithmic Challenge

Invented in the 6th century [3], various strategies and algorithms for the game of chess have since been developed. The 1950s and 1960s saw the first fully developed chess programs that leveraged alpha-beta search algorithms and various heuristics, being able to beat humans in tournament play [6]. Most recently, neural networks and Monte-Carlo methods (such as MCTS) have been used to develop and/or learn evaluation functions that have enabled chess-engines to attain superhuman levels of performance, such as in the case of Stockfish and AlphaZero.

1.2 Our Contribution

Our work explores the performance of the baseline MCTS with Upper Confidence Bounds (UCB) algorithm for the purposes of playing chess, as well as its variations where we calibrate exploration vs. exploitation trade-off with self-play. We also compare the algorithms to baseline models, such as a pure Alpha-Beta search to analyze which particular variation is empirically the most effective at playing the game of chess. In our analysis, we also consider various implementation tradeoffs (e.g. time and memory efficiency). The core infrastructure of our chess engine (move generation, maintaining board state et cetera) was written by Michał in C++ over the summer. Thus, our work focused on developing various models for the evaluation and search functions for the engine. Interestingly, our MCTS-based approach was able to obtain a good grasp of the chess fundamentals (e.g. material values of individual pieces and how they change over time as the game progresses).

2 Related Methods

Heuristic-based methods of tree search for the game of chess include performing an exhaustive breadth-first search of the game tree for the first few several layers before utilizing heuristic thresholds at deeper levels of the game tree [2] in order to prune unpromising branches, thereby allowing the algorithm to search deeper levels of the game tree, whilst circumnavigating the problem of combinatorial explosion. Fogel et al [7] utilized an evolutionary algorithm, allowing the model to learn to play chess by simply playing games against itself. Moreover, alpha-beta pruning was used to search the game tree itself and was done so in combination with a neural network trained to assess specific sections of the chessboard (i.e. the output of a neural network was used to determine the evaluation score assigned to a particular board state s). Bandit-based algorithms have been applied to the problem of tree search, particularly in the context of games such as Chess and Go; Coquelin et al [5] discuss UCT - a UCB-based Bandit Algorithm applied to game trees - which allows for a tree-based algorithm that statistically balances exploration and exploitation. Takeuchi et al [10] analyze various methods in which evaluation functions (in the context of game programming) can be tuned in order to strongly correlate with win probabilities (i.e. a low evaluation score corresponds to a low probability of winning and vice versa).

2.1 Evaluation Functions

An evaluation function is a function $f_{\text{eval}} : \mathcal{S} \rightarrow \mathbb{R}$ that assigns a real value $r \in \mathbb{R}$ to each state $s \in \mathcal{S}$. For any arbitrary state $s \in \mathcal{S}$, we note that $f_{\text{eval}}(s) \in \mathbb{R}$ corresponds to the *centipawn* score of state s . This is simply a unit of measure in chess that quantifies the current advantage of a given board state with respect to the side that is to move. If $f_{\text{eval}}(s) > 0$, then the side-to-move is winning and $f_{\text{eval}}(s) < 0$ is thus a losing state with $f_{\text{eval}}(s) = 0$ corresponding to a draw.

43 Suppose that it is Player A’s turn to move, where the chess board is in state s . Ideally, the evaluation
 44 function should be close to the *true* reward function r of the chess game itself, so that Player A would
 45 simply select the action a that maximizes expected reward:

$$s^* = \operatorname{argmax}_{s' \in \Delta(s)} f_{\text{eval}}(s')$$

46 where $\Delta(s)$ denote the set of all of the possible states $s' \in \mathcal{S}$ that can be transitioned to from state s
 47 in one legal move. That is to say, Player A will simply take the unique move a that transitions the
 48 game state from the current state s to the state s^* that is most favorable according to the evaluation
 49 function.

50 2.2 Tree Search + Evaluation Functions

51 In reality, learning such an "ideal" $f_{\text{eval}}(s)$ just as a function of the present state of the board s
 52 is practically impossible, as the evaluation itself cannot (usually) predict things like traps or some
 53 nuanced lines of play, especially in a game like chess. Therefore, it is essential to "look ahead" and
 54 attempt to simulate how the game could play out from next states s' . Ideally, we would like to explore
 55 all potential future outcomes (i.e. sequences of moves) with the "look-ahead". However, as we know
 56 the game of chess suffers from *combinatorial explosion*: looking ahead by just a few moves in a
 57 brute-force manner is computationally infeasible.

58 **Ultimately, evaluation & search are complementary:** We first develop a reasonably good heuristic
 59 for how good the states are without any look-ahead – our actual f_{eval} . Then, we attempt to improve
 60 the look-ahead exhaustive tree search by not exploring branches where we already have sufficiently-
 61 high confidence about the state’s outcome due to f_{eval} . Luckily, RL approaches, namely self-play,
 62 can help us both to tune f_{eval} and to efficiently incorporate it in the MCTS search.

63 The subsequent sections of the paper are organized as follows. In the remains of this section, we
 64 describe current State-of-the-Art (SOTA) chess engines and how they relate to our work. Section
 65 3 goes into technical details of our methods for both f_{eval} and the MCTS search strategy. Section
 66 4 outlines implementation, while Section 5 presents results and our analysis. Finally, Section 6
 67 discusses future directions and concludes.

68 2.3 Chess SOTA

69 **Stockfish:** utilizes Sequential Probability Ratio Testing (SPRT) in its tuning, determining through
 70 self-play whether one version of the engine is stronger than another. In SPRT applied to chess engine
 71 testing, we set up a null hypothesis (e.g. the new version of the engine is at least 0 Elo stronger
 72 than the old one) and an alternative hypothesis (e.g. the new engine is stronger by 5 Elo), with a
 73 certain confidence level (typically 95%). The number of self-play games needed to get statistically
 74 significant results often turns out to be very large (in the order of hundreds of thousands of games for
 75 each little change to the engine’s parameters), especially if the performance difference between the
 76 two versions of Stockfish tested is subtle. As such, Stockfish uses a distributed framework (Fishtest)
 77 where volunteers donate CPU time to run self-play games between new and old versions of Stockfish.
 78 As of December 2023, it’s reported that Fishtest has used over 11,600 years of CPU time to play
 79 more than 6.8 billion chess games [1].

80 **AlphaGo Zero:** though not a chess engine, AlphaGo Zero takes an approach much closer to ours.
 81 Unlike its predecessor AlphaGo, it doesn’t rely on human game data at all, instead learning Go from
 82 scratch through self-play. These games are generated by the most effective player (version of the
 83 engine) from all previous iterations. The performance of a new player iteration was measured against
 84 the best player to date; if the new player won by a margin of 55% or more, it would replace the best
 85 player. Subsequent self-play games would then be generated by this new, superior player.

86 **AlphaZero:** a significant departure from the way traditional (Alpha-Beta based) chess engines were
 87 tuned in the past, the DeepMind developed engine utilizes self-play without reliance on historical
 88 game data or opening books to train a combination of two deep neural networks (separate policy
 89 and value networks) and MCTS hyperparameters. Unlike AlphaGo Zero, which created distinct
 90 generations of players, AlphaZero maintains a single neural network during training that gets updated
 91 continually with each iteration. AlphaZero’s training involved 5000 first-gen Google TPUs (for
 92 self-play) and 64 second-gen TPUs (for training the neural networks) [9].

93 3 Our Methodology & RL Components

94 Drawing inspiration from some of the most successful chess and Go engines described in the previous
95 section, we decided to incorporate self-play to derive both a good f_{eval} and tune the MCTS search
96 hyperparameters. This section first goes into detail of fitting f_{eval} with self-play using SGD and
97 self-play tournaments between different instances of our policies (as elaborated upon in Section 3.1).
98 Then, we introduce our initial MCTS and strategies that helped to improve our MCTS implementation
99 further.

100 3.1 Self-play Setup

101 Lishex (Ours)

102 We adopt a DAGGER-like (Dataset Aggregation) approach for iterative improvement of our engine.
103 Self-play games are generated by the current best iteration of our player (engine), and from every
104 game a set of chess positions is randomly selected and annotated with a label (we elaborate on the
105 specifics of the labels later on in this section). After each training iteration, we evaluate the player’s
106 performance based on the mean squared error (MSE) against our dataset of self-play games, rather
107 than direct win-rate comparison as in the case of AlphaZero. This helps reduce the computational
108 needs of training; namely, we do not require potentially hundreds of thousands of games to get
109 statistically significant results. One theoretical disadvantage of our approach, however, is that the
110 respective method can potentially lead to less accurate results (i.e. there is no statistical guarantee),
111 but we found this approach very effective in practice.

112
113 We periodically play a large number of games with the new player, using positions from
114 these games to augment our dataset. This approach ensures our model is always trained on fresh,
115 up-to-date data. To improve the rate at which our model learns (as well as to combat memory
116 requirements), as we aggregate new data into the dataset we discard all training data, generated by
117 past players. We utilize chess tournament software, ‘cutechess’ to perform and record the results of
118 N self-play games of our engine (with N sufficiently large, e.g. $N = 10000$, all played in parallel on
119 an 8 core machine). To manage time constraints, we set short time controls for these games, such as
120 10 second time limit with 0.05s per move.

121
122 The resulting dataset consists of N game results, stored in the Portable Game Notation
123 (PGN) format: a human-readable text format keeping track of each move performed in the game,
124 as well as the final result, which can be a win for white (1-0), draw (1/2-1/2), or loss (0-1). From
125 each game, we randomly select a small subset of positions, say $n = 10$, for further analysis. These
126 positions form the input of our training data. Correspondingly, the actual game outcomes form the
127 output ‘label.’ Each row of our dataset thus represents a unique position paired with the game’s final
128 result. For the labels, we encode the label with respect to the side-to-move in a given position. For
129 example, if in a given position X_i white is to move, and white wins the game, the corresponding
130 label y_i is $y_i = 1$, whereas if black was to move the label would have been $y_i = 0$ (we encode draws
131 as 0.5).

132 This paradigm of self-play was primarily used to tune evaluation function f_{eval} , as described in the
133 following section.

134 3.1.1 Tuning (Training)

135 To take a step-back, currently the evaluation function consists of numerous heuristics as the earlier
136 mentioned "centipawn scores" and others. The question is how to combine them all together into
137 a working, effective f_{eval} , which is just a tuning process of choosing the "weights of importance"
138 for those different heuristics. During this tuning process, the parameters of our evaluation function
139 are updated through stochastic gradient descent. The training is carried out in a mini-batch fashion,
140 where gradients are estimated numerically.

141 Each position is fed into our static evaluation function, resulting in some centipawn score P . We then
142 transform this score into a perceived probability W of winning for the player whose turn it is using
143 the following formula (for an in-depth explanation as to why this approach works, see [10]):

$$W = \frac{1}{1 + 10^{-P/4}}$$

Our loss function then becomes the Mean Squared Error of our evaluation function’s predictions and the actual outcomes of the games.

Care has to be taken to maintain a balanced distribution of classes in our dataset. Chess is considered a *drawish* game, where a vast majority of games between top players results in draws. The high frequency of draws can skew our evaluation function towards predominantly predicting draws (to minimize the loss). To counter this, we trim our dataset to maintain a balanced ratio of 30% wins, 40% draws, and 30% losses (W-D-L) in our training dataset.

3.2 MCTS + UCB: baseline

As discussed earlier, the game breaks down to the effective decision of which move to make at each timestep from the current position r ("root"). In our tree, each node p represents a position in the game reachable from the root r with a indicator of whose turn it is. Our goal is to understand how promising are the children **from the perspective of the player** whose turn it currently is in the node p . We accomplish this by estimating $AveValue(s)$, which importantly stands for the probability of **the parent** of s to win if they were to make a move into s :
 $\pi(\text{win for } p \mid \text{move } p \rightarrow s) \approx AveValue(s)$.

In this way, our version of MCTS is closer to the spirit of AlphaZero rather than original MCTS. Below we provide the pseudocode for this search, while the actual code can be found in the "MCTS_Search" function in our GitHub repo. Figure 1 also shows the schematic presentation of the algorithm.

Repeat $t = 1$ to N ($\approx 100k$):

(1) **Selection.** Start in the root, go down until current node s has all children explored. Every time choose the child with max UCB score of winning probability for the parent s : Define $s' = \text{nextState}(s, a)$ and define

$$UCB(s, a) = AveValue(s') + C \cdot \sqrt{\left(\frac{\log(\#visits(s))}{\#visits(s')}\right)}$$

choose the action $\hat{a} = \text{argmax}_a(UCBscore(s, a))$
 set $s = \text{nextState}(s, \hat{a})$

(2) **Expansion + Simulation.** Simulate the game to get estimated *reward* at s . This requires obtaining a rollout deeper in the tree until the terminal node (reward = ± 1 or 0 when ties) or until we ran out of budget¹, in which case we estimate the reward with f_{eval} .

(3) **Back-propagation.** For all board states "visited" on our path to s , update the stats: let $V(s') = \text{visits}(s')$, then

$$AveValue(s') = \frac{V(s')}{V(s') + 1} \cdot AveValue(s') + \frac{1}{V(s') + 1} \cdot (\pm 1) \cdot \text{reward}$$

$$V(s') = V(s') + 1$$

where the sign \pm depends on whose turn it was in the s' node.

(4) **Output.** Finally, in the root we return the action that maximizes the "empirical" score across root’s children $\text{argmax}_s AveValue(s)$ for all s children of r – without the exploration component.

The overall idea is very much of the RL spirit: balance between exploration of the moves that give high empirical probability of winning (based of what we have seen so far) and the moves that we have not visited much yet. This exact trade-off is controlled by the const C in the UCB score (which we tune), while the "learning" of empirical probabilities happens via rolling-out policies again and again.

3.3 MCTS limitations

There are several limitations of the initial MCTS + UCB approach described in the previous section.

(1) First, as mentioned, in the game of chess the tree "explodes", i.e. already on the 4th level (4 moves

¹Time and memory budget constraints are discussed in the implementation section.

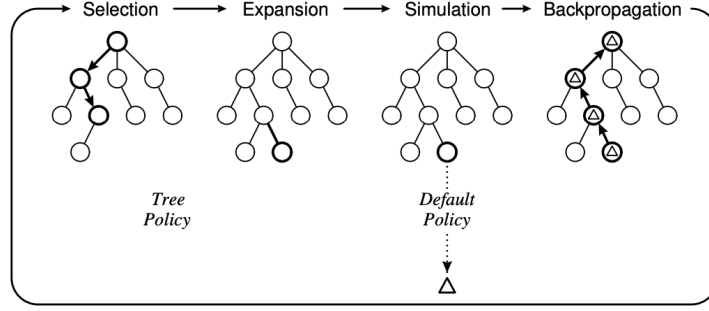


Figure 1: Diagram of a single MCTS search loop. Tree represents chess board positions reachable from the root, default policy can be random or fully deterministic determined by f_{eval} .

191 ahead) we deal with $\approx 1\text{M}$ of potential nodes. Initial algorithm’s **selection** stage always prioritizes
 192 unexplored moves: we traverse the tree until the first moment we find the node s who has at least one
 193 unvisited child. Thus, most probably, we spend lots of time and memory on moves that are clearly
 194 unpromising.

195 (2) Second, the **simulation** stage as it’s proposed originally is very *slow* and very *noisy*. It’s slow
 196 because we choose to traverse the tree until the very end (terminal node), which might require up to
 197 100 moves. It’s noisy because our default policy is completely random (all children are unvisited
 198 and so we choose actions by random), and thus the game’s output *reward* after 100 moves is very
 199 uninformative of initial’s winning probability for s .

200 To address concern (1), we tried two approaches:

201 **Epsilon sampling**: during the selection stage, with some probability ϵ we decide to exploit current
 202 children, and with probability $1 - \epsilon$ we stick to the previous approach of exploration (i.e. if there are
 203 unseen children – explore them). ϵ is a hyperparameter that we tune.

204 **Evaluation based sampling**: when exploring unseen children s of current node p , instead of picking
 205 them by random, we sample from the distribution *inverse*-proportional to $f_{\text{eval}}(s)$. In this way,
 206 children with high values of $f_{\text{eval}}(s)$, i.e. promising positions for the opponent (since at the position s
 207 it’s opponent’s turn!), would get decreased chances of being explored – because they are less lucrative
 208 for the player at parent node p .

209 To address concern (2), we tried the following ideas:

210 **Budget limits for depth**: our simulation rollouts do not go deeper than some *BUDGET* hyperpa-
 211 rameter. Then, we output $\text{reward} = f_{\text{eval}}(s_{\text{end}})$ for the node at which a given rollout stopped. This
 212 speeds up the simulation part (no need to wait until the end of the game), but also puts more trust on
 213 the f_{eval} correctness.

214 **Heavy Rollouts**: during the simulation assume that players have preference towards nodes that have
 215 lower f_{eval} scores of the children (i.e. opponents is in disadvantage). We tried both probabilistic
 216 sampling with children’s f_{eval} scores (similar to 1.2 idea), and deterministic $\text{argmin}_s(f_{\text{eval}}(s))$.
 217 This reduces "the noise" of completely random, "light" rollouts, but also slows down the process
 218 since f_{eval} is quite time-expensive to compute.

219 **Averaging out Rollouts**: Finally, another option to reduce the noise is simply to do numerous rollouts
 220 and output the average reward among all of them.

221 3.4 Tuning MCTS hyperparameters: manual self-play

222 In addition to utilizing self-play to tune our evaluation function, we utilize it to tune the hyperparam-
 223 eters of our MCTS algorithm: the rollout budget, the UCB exploration constant C , and the ϵ noise
 224 constant used during the selection phase, as well as the rollout policy used. Notably, the methodology
 225 we used for hyperparameter tuning is akin to a regular hyperparameter search combined with the
 226 methodology employed by AlphaGo Zero: we start with an existing version of our engine, using some
 227 initial set of hyperparameters. We then manually² adjust these parameters and create a new version of
 228 the engine for testing. We then perform a couple hundred self-play games between the two versions

²That’s why we call it "manual" self-play

229 of the engine and pick a winner to move onto the next iteration according to the win-draw-lose ratio.
230 Despite its simplicity (and lack of more sophisticated approaches like SPRT used by Stockfish),
231 considering the initial suboptimality of our hyperparameters and the limited computational resources
232 and time, we were happy with the results.

233 4 Implementation

234 We have implemented a single-threaded pure MCTS search routine, as well as a UCT variant (UCB
235 applied to Trees). The code for our implementation can be found in our GitHub repository. As
236 discussed in section 3.3, we experimented with multiple rollout policies for our MCTS routine:
237 fully random, evaluation based greedy, and evaluation based sampling. We discuss the comparative
238 performance of these strategies in section 5 (results).

239 We implement a *Node* class, as well as multiple helper functions for helping perform the MCTS
240 search. Note that we structured our code in a plug-n-play fashion: allowing us to quickly switch
241 between utilizing AlphaBeta, Pure MCTS, or MCTS w/ UCB for our experiments further down the
242 line.

243 4.1 Time management

244 In the original formulation of MCTS we perform some number K of iterations, during each iteration
245 performing each step of MCTS: selection, simulation, playout, backpropagation. In practice however,
246 a game like chess imposes an additional constraint: time controls. To address this, our implementation
247 periodically checks during the MCTS process whether the allocated time for a given move has been
248 exhausted, and halts the search if necessary. Then, the best move is reported to the user according to
249 the win/visit ratios accumulated at the root during the search. Note that this is safe to do with MCTS,
250 since it's an *anytime* algorithm: it can be halted at any time during execution and still return a valid
251 move (though not necessarily the most optimal one!).

252 4.2 Memory management

253 Due to the anytime nature of MCTS and the way it grows the search tree it's performance improves
254 the longer it runs. As it operates, it allocates more nodes and progressively builds up the search tree. In
255 practice, given constraints of finite memory resources, especially in a game of chess where numerous
256 searches are conducted, it's essential to bound the memory usage of MCTS. In our implementation,
257 we utilize an approach of *stunting*: we impose an upperbound on the memory usage of MCTS
258 and after we run out of this memory, we stop expanding the tree: the selection, simulation, and
259 backpropagation steps continue to be executed until the time limit is reached [8]. To achieve the
260 above we implement a custom arena memory allocator: a big chunk of memory, for instance 2GB, is
261 preallocated for the entire tree during the startup of the engine and each new node allocation during
262 the expansion step utilizes only this part of memory. This not only bounds the memory used by the
263 engine, but also drastically improves our search performance by eliminating the need for dynamic
264 memory allocation at each expansion step: our engine experienced a fivefold improvement in its
265 nodes (positions) searched per second (NPS) rate after switching to our arena allocator.

266 4.3 Evaluation function

267 Lishex employs a traditional handcrafted evaluation function (HCE), as opposed to neural-based
268 approaches of modern engines like Stockfish and AlphaZero. The evaluation at its core is a linear
269 combination of various independent features, each assigned a specific weight. These features range
270 from human-derived chess concepts (such as king safety and mobility of pieces) to direct features of
271 the chessboard like the number and types of pieces and their spatial distribution. We also consider
272 more complex interactions between pieces: a pawn supporting a minor piece or a queen attacking a
273 hanging pawn.

274 In addition we split our evaluation function into two distinct phases: the middlegame and the endgame,
275 each with a separate set of weights. For any given board state s , we calculate two separate scores:
276 one for the middlegame and another for the endgame. As the game progresses, we linearly interpolate
277 between the two to reflect the transition from one phase to the other. We use a material balance based

278 heuristic to determine the current phase of the game: the less pieces there is on the board the closer to
 279 the endgame we expect to be.

280 5 Results

281 We split this section into two parts: sole analysis of learned patterns in the evaluation function and
 282 best strategies of MCTS as identified by our manual self-play.

283 5.1 Evaluation Function analysis

284 We rescale the centipawn scores assigned by Lishex so that one pawn is equivalent to 100 centipawns,
 285 as detailed in Table 1. Note that our results (especially for the endgame) align well with the commonly
 286 held human scale of material values: 1 for a pawn, 3 for knights and bishops, 5 for rooks, and 9 for
 287 queens [4].

288 In figure 2 we observe quite a few interesting trends. Though initially random, over time Lishex
 289 steadily increases its belief in the high value of queens and the relatively lower significance of
 290 pawns. Notably, Lishex learns that pawns gain value in the endgame, likely due to their potential for
 291 promotion (playing a decisive role in the outcome of many games). It also realizes that bishops are
 292 marginally stronger than knights and interestingly suggests that queens and rooks become slightly
 293 less valuable in the endgame. This latter finding is particularly intriguing: for queens we hypothesize
 294 that Lishex learns that their higher mobility becomes less and less useful as there are fewer pieces on
 295 the board to be attacked. For rooks however, the typical human heuristic says that they become more
 296 and more valuable as the game progresses, as their movement and attacks are less obstructed than in
 297 the middlegame. That said, the values of the parameters don't deviate much from the human nor the
 298 AlphaZero values, making us confident in Lishex learning correctly!

299 Arguably our engine didn't learn chess from scratch the way AlphaZero did: except material values,
 300 sensible defaults were used for all other parameters of the evaluation function, giving us a head start.
 301 Despite that, we were happy to see our self-play work out and Lishex learning sensible material
 302 values for the different piece types.

Table 1: Comparison of Material Values (pawn relative)

Piece	Human	AlphaZero	Lishex (middle game)	Lishex (endgame)
Pawn	1	1	1	1
Knight	3	3.05	5.75	2.72
Bishop	3	3.33	6.42	2.84
Rook	5	5.63	9.29	4.95
Queen	9	9.5	17.51	9.56

303 5.2 MCTS analysis

304 Overall, the best choice of strategies was **Epsilon sampling** with $\epsilon \approx 0.5$ and utilizing **light rollouts**
 305 for the simulation part. Below we provide the discussion of the comparative performances between
 306 different policies and why believe the final choice turned out to work the best.

307 5.2.1 Choice of selection

308 As described in the methodology section 3.3, Epsilon sampling approach tends to **exploit** more:
 309 in roughly half of the cases we use the children we have seen already (and thus we are confident
 310 about their scores), rather than unseen children. The problem with aggressive **exploration** approach
 311 proposed initially, i.e. when we prioritize any child we have not seen yet (which also corresponds
 312 to $\epsilon = 0$) is that we spent too much of compute time and memory to explore all the children *in*
 313 *breadth* from the root. As a result, our tree depth practically never exceeded a value of 6 (with the
 314 searching time limit of 10 seconds), which corresponds to only looking 3 steps ahead for both of the
 315 players. This means that for child nodes 3 steps ahead we had just a small number of estimates of
 316 rewards from simulation rollouts, perhaps implying that those estimates were still very noisy and
 317 uninformative.

Figure 2:

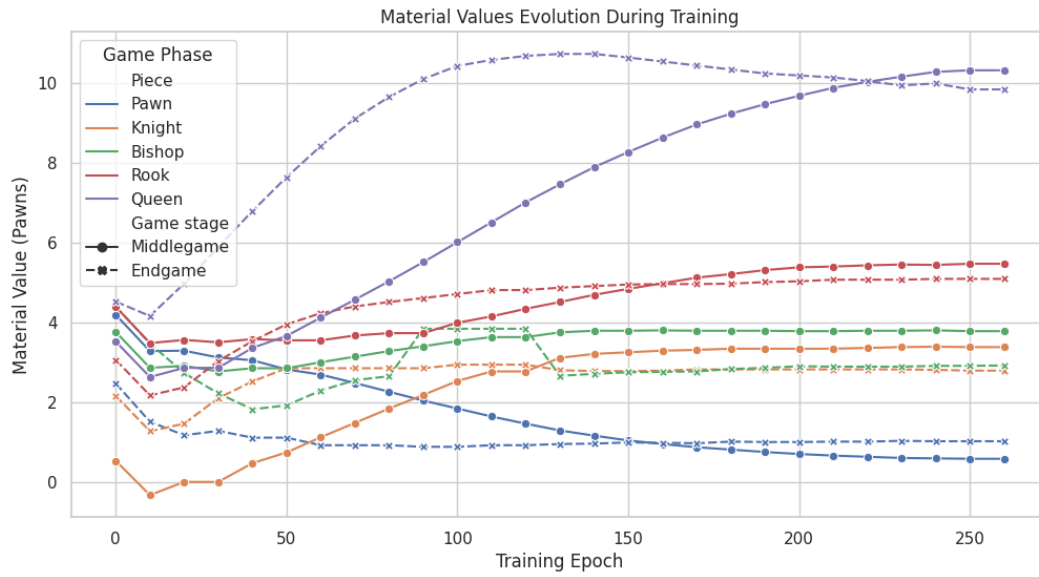
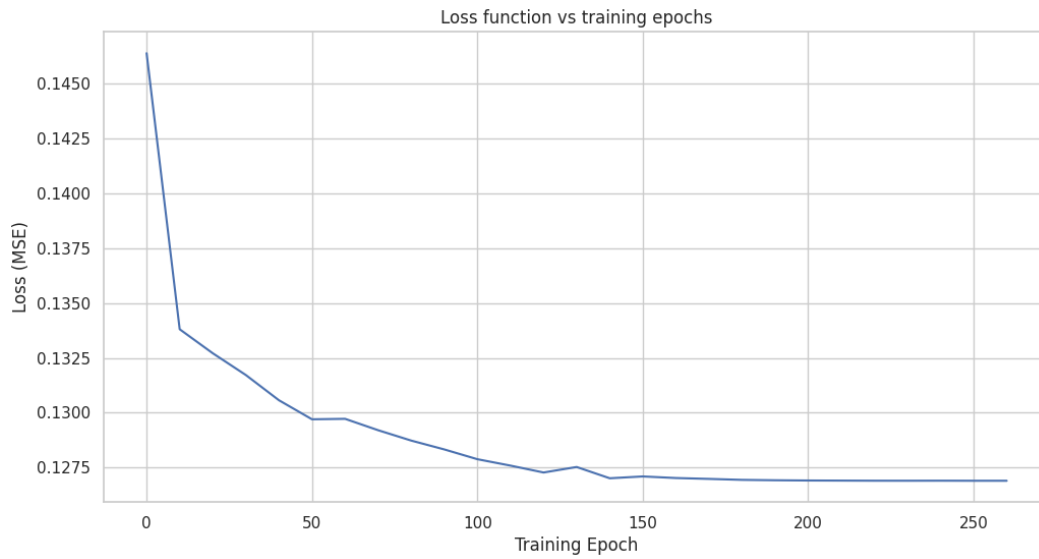


Figure 3:



On the contrary, with $\epsilon \approx 0.5$ our tree depth was always at least 12 (and it increased up to 18 as we approached the endgame), because we were exploring less. In such a scenario, child nodes 3 steps ahead have a much larger number of simulation rollouts, and thus estimated rewards / winning probabilities are more informative, explaining the improved performance.

Evaluation based sampling suggests sampling based on f_{eval} scores to make decisions about moves during the selection stage. However, it turns out that in most cases these scores tend to be fairly close to each other numerically. That’s because two children of any particular node would differ just a little bit in terms of centipawn score (say, some player loses a pawn), which results in final winning probabilities difference on the magnitudes of 1 – 5%. In other words, such approach would be very similar to that of just random sampling, which we also observed empirically. To address this issue one could do transformations of obtained scores (say, square the winning probability estimates) to let the underlying distribution deviate more from the uniform (i.e. random choice). Further investigation is required for to calibrate the proper approach and see whether it could outperform the **epsilon sampling**.

5.2.2 Choice of simulation

First, **heavy rollouts** (i.e. assuming that players make their moves in the simulation based on f_{eval} values) tended to underperform. Using f_{eval} -based sampling suffered from the same issue to that of described in the previous paragraph: the distribution is very close to uniform, and, to add to this, search is slowed down because computing f_{eval} is not cheap. On the other hand, fully deterministic $\text{argmax}_s f_{\text{eval}}(s)$ (among all children s) strategy seems to be too greedy and over-reliant on the f_{eval} . In other words, there was no exploration in the simulations at all, so the search was missing out some non-obvious moves two-three steps ahead.

Finally, **averaging out** several rollouts during the simulation (instead of using just one of them) didn’t produce neither positive nor negative improvements. We found a good discussion here of why this is the case. In short, even though one play-out per selection step might be a very noisy estimate of winning probability, subsequent selection steps can go through the same node multiple times. This is especially true when using the **epsilon based** sampling as we have higher probabilities to visit a node and perform more rollouts from it.

Engine	NPS (k)
Stockfish	932
AlphaZero	80
Lishex AlphaBeta	524
Lishex MCTS	22

Table 2: Nodes (positions) per Second (NPS) Comparison. Note: AlphaZero ran on DeepMind hardware, whereas we benchmark the other three engines on an Intel Core i7-7700HQ CPU.

6 Conclusion & Future Work

Inevitably, we never expected our results to be close to the state-of-the-art or a carefully tuned alpha beta approach. Our main goal was to see just how far a combination of RL + MCTS could take us within the game of chess given the scope that we had for this project. All in all, we are satisfied with the results; it seems as though our chess engine was able to learn a solid grasp of the chess fundamentals and perform some good openings and moves. As for future work, given our analysis, we would be interested to explore a combination of MCTS and Alpha Beta search. This could be done by deciding which branches to prune during the MCTS search, using α and β values, a proposed method we denote $\alpha\beta\text{MCTS}$. We think that this combined approach could pose as an interesting avenue for future work.

References

- [1] Stockfish testing framework.
- [2] G.M. Adelson-Velskiy, V.L. Arlazarov, and M.V. Donskoy. Some methods of controlling the tree search in chess programs. *Artificial Intelligence*, 6(4):361–371, 1975.

- 360 [3] Yuri Averbakh. *A history of chess: From Chaturanga to the present day*. Russell Enterprises,
361 Inc, 2012.
- 362 [4] J.R. Capablanca and N. de Firmian. *Chess Fundamentals: Completely Revised and Updated for*
363 *the 21st Century*. McKay chess library. Random House Puzzles and Games, 2006.
- 364 [5] Pierre-Arnaud Coquelin and Rémi Munos. Bandit algorithms for tree search. *arXiv preprint*
365 *cs/0703062*, 2007.
- 366 [6] Carl Felstiner. Alpha-beta pruning, 2019.
- 367 [7] D.B. Fogel, T.J. Hays, S.L. Hahn, and J. Quon. A self-learning evolutionary chess program.
368 *Proceedings of the IEEE*, 92(12):1947–1954, 2004.
- 369 [8] Edward Powley, Peter Cowling, and Daniel Whitehouse. Memory bounded monte carlo tree
370 search. *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital*
371 *Entertainment*, 13(1):94–100, Jun. 2021.
- 372 [9] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur
373 Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy P. Lillicrap,
374 Karen Simonyan, and Demis Hassabis. Mastering chess and shogi by self-play with a general
375 reinforcement learning algorithm. *CoRR*, abs/1712.01815, 2017.
- 376 [10] Shogo Takeuchi, Tomoyuki Kaneko, Kazunori Yamaguchi, and Satoru Kawai. Visualization
377 and adjustment of evaluation functions based on evaluation values and win probability. In
378 *PROCEEDINGS OF THE NATIONAL CONFERENCE ON ARTIFICIAL INTELLIGENCE*,
379 volume 22, page 858. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press;
380 1999, 2007.