# HEX playing AI

By

**Arapsih Güngör, Daniel Pötscher, Philipp Jonas, Thomas Gruber**

# Contents

# 1. Introduction:

The goal of this project is to implement and train an AI to play the game HEX [1] using the AlphaZero reinforcement learning algorithm (RL) [2], [3] [4] created by DeepMind [5]. The institutes SLURM [6] cluster was used to offload training data generation to up to 20 machines in parallel. There are various publicly available solutions that have implemented that algorithm for different games and in different languages. This was helpful to get an overview and a better understanding for the actual implementation.

## 1.1. HEX Game

The Hex game was invented in 1942 by Piet Hein and re-invented in 1948 by John Nash. It is a two-player game where each player tries to connect the horizontal or vertical boundaries of an n x n Hex board. The Size of the board can vary but common sizes are 11x11, 13x13, 14x14, 17x17 and 19x19 [7].



*Figure 1 - The Hex Game [1]*

The players place their stones one after the other, trying to reach the opposite side while preventing their opponent from doing so.

Since there can never be a draw in this game, white will always win in a perfect game.

To compensate for the advantage of the first move, a swap rule is usually used: after white has made the first move, black may choose whether or not to swap pieces with white, which would mean that white would have to play against its own opening move. [1]

## 1.2. Related Work

There are several projects on GitHub and other platforms where the AlphaZero algorithm has been implemented for various games. The most forked repository is alpha-zero-general [8] which implements the core algorithm of AlphaZero according to the original paper and contains example applications for games like Othello [9], TicTacToe [10], and others. One of the forks of this repository is called alpha-zero-hex [11] which implements the HEX game using the AlphaZero implementation of alpha-zero-general. This implementation doesn't use any parallelization techniques and therefore is slow to train. Further first tests do use this repo did not result in a converging model that played better than random.

## 1.3.   AlphaZero

AlphaZero is a reinforcement learning algorithm developed by DeepMind that has gained lots of attention because its application in games like GO or Chess has led to superhuman level of play and convincingly defeated world-champion programs in each case [3]. It is astoundingly simple and yet very powerful as it combines techniques as the Monte-Carlo-Tree-Search (MCTS) [12] with a deep-learning approach [13] to overcome shortcomings in both approaches.

The MCTS is a heuristic search algorithm that simulates several game steps given an initial state and returns the best action based on the information it gathered during the simulations. In principle the search algorithm uses Upper-Confidence-Bound based exploration which focuses on expanding nodes that have been marked as leading to a win in previous simulations or randomly choose new nodes to expand.

The integration of deep learning into MCTS happens at the sampling level. Where the classical MCTS chooses randomly to expand on new nodes AlphaZero uses the current state of the incrementally trained neural network to predict probabilities for the next actions to choose.

While a pure MCTS algorithm is limited by the available memory to store simulation data and a pure neural network has only limited knowledge about a games development the combination of both make it decrease their shortcomings. In each training iteration a new instance of MCTS is created starting without any prior knowledge and new nodes are expanded using a trained network which can make more informed decisions than a random number generator. After several simulations these results are used to further train the network which in the next iteration should improve the predictive behaviour of the network and hence the quality of the MCTS search.

## 1.4.   SLURM

SLURM is an open source, fault-tolerant cluster management and job scheduling system. The FHTW has such a cluster active which students can use to run computationally expensive tasks distributed on several machines. It consists of 19 machines which each have a NVIDIA RTX 2080 TI installed of which two were down at the time of training.

The cluster can be controlled with view simple commands which are listed below, and their usage will be explained later.

- sinfo: shows the state of the cluster nodes.
- squeue: shows the current job queue.
- sacct: lists completed and running jobs.
- sview: graphical monitoring tool for jobs and nodes.
- salloc: allocates some resources on the cluster without starting a task.
- srun: runs a single job interactively.
- sbatch: runs a job in the background. Also runs job arrays.
- scancel: cancels a running job.

## 2. The Implementation

As a starting point the code from the alpha-zero-hex repository was analysed and tested but found to not converge to a better model after several days of training. In course to understand more in depth what was going on and to overcome some flaws of the existing code, we decided to reinvent the wheel and code everything again from scratch while using the existing code as guidance but with the focus on parallelizing it to be able to run on the SLURM cluster.

The new implementation consists of two main applications which are the DataGenerator and the ModelTrainer. Both are standalone applications that only communicate via file exchange with each other. Internally they use the MCTS implementation to either generate new training data with the latest neural network or evaluate the quality of a new trained model.

### 2.1. The Neural Network

The model is a convolutional neural network that consists of 11 hidden layers as can be seen in Figure 2. The first six layers consists of 3 convolutions and batch normalizations and the next four are fully connected layers, also with batch normalizations to narrow the model down the final two parallel output layers that should predict the probabilities for the next best action and the expected outcome of the game given a game state as an input.

```python
class HexNetModel(nn.Module):
    def __init__(self, game: Game) -> None:
        super(HexNetModel, self).__init__()
        self.game = game

        self.num_layers = 64

        self.conv1 = nn.Conv2d(1, self.num_layers, 3, padding=1)
        self.conv1_bn = nn.BatchNorm2d(self.num_layers)
        self.conv2 = nn.Conv2d(self.num_layers, self.num_layers, 3, padding=1)
        self.conv2_bn = nn.BatchNorm2d(self.num_layers)
        self.conv3 = nn.Conv2d(self.num_layers, self.num_layers, 3, padding=1)
        self.conv3_bn = nn.BatchNorm2d(self.num_layers)
        self.fc1 = nn.Linear(self.num_layers * self.game.size * self.game.size, 512)
        self.fc1_bn = nn.BatchNorm1d(512)
        self.fc2 = nn.Linear(512, 256)
        self.fc2_bn = nn.BatchNorm1d(256)
        self.fc_pi = nn.Linear(256, game.action_size)
        self.fc_v = nn.Linear(256, 1)

    def forward(self, x):
        x = F.relu(self.conv1_bn(self.conv1(x)))
        x = F.relu(self.conv2_bn(self.conv2(x)))
        x = F.relu(self.conv3_bn(self.conv3(x)))
        x = x.view(-1, self.num_layers * self.game.size * self.game.size)
        x = F.relu(self.fc1_bn(self.fc1(x)))
        x = F.relu(self.fc2_bn(self.fc2(x)))
        pi = self.fc_pi(x)
        v = self.fc_v(x)
        return F.log_softmax(pi, dim=1), torch.tanh(v)
```

*Figure 2 - The Convolutional Network Model*

## 2.2.    The Monte Carlo Tree Search

The implementation of the MCTS has one main entry point which is the predict function. Given any game state as a parameter it will perform a number of searches to unravel new leaf nodes that have not been seen before and while doing so gathers heuristic information about the potential outcomes of the game from this state. After these searches have been performed a simple probability distribution is returned depicting the best next actions to choose.

In a search step the algorithm expands already seen nodes or chooses to expand to a new node depending on the UCT score each possible action achieves. This UCT score consists of two major parts – one term for exploration and one exploitation term. While this exploration term would use a random distribution in a standalone MCTS implementation the AlphaZero one uses the neural networks prediction at this stage. Also, another output of the network is the possible outcome of the game. If a new node is expanded then this value is taken as if the search has progressed to the end of the game and returns the "winner" to the parent nodes to update their respective state values. At first these predictions are random themselves but with increasing convergence of the model they become more and more accurate and therefore improve the quality of the predictions.

The state of any node is defined by five key values which get updated at each search step and are as follows:

- Has the game ended
- PI – the networks action prediction
- N – The number of times a state has been visited
- Na – The number of times an action has been chosen for a state
- Qa – The quality parameter used in the UCT calculations

For each node that has been visited the respective counters get increased and the Q values updated.

Finally, after the search has finished the predict function returns simply returns Na / N – the probability of how often an action for a state was visited.

## 2.3.    The Data Generator Application

As mentioned above, training the network takes a long time and the most time-consuming part is the generation of training data. The data gets generated using the MCTS with the current best model. Unfortunately, the prediction of the model is relatively slow and must be executed many times for each MCTS prediction. To improve training data generation speeds the application has been designed to independently work by reading a current model from a shared path and writing results immediately back to disk where the model-trainer picks it up and incorporates it in the training of the next model version. Since this application is so loosely coupled to the trainer, we can spawn several instances of it to speed up data generation. In our training session up to 100 instances were running in parallel.

## 2.4.    The Model Trainer Application

The model-trainer is the counter part of the data generator and reads the generated data from disk and starts incremental trainings on a current model. After each training step an evaluation phase of a number of self-plays is held to determine if the new trained model performs better than the previous version. If so, this new model gets stored to disk where the data-generators can pick it up again to generate even better data. If the new model performs worse than training is continued with new incoming data.

## 3. Results

The training was done over a period of two days with said 100 data generators which generated in total 32.286.084 records and were further used in 1244 model trainings of which 100 were actually accepted as better models. Not all these models were kept but only occasionally samples of them were downloaded to test their performance.

From the models that are available a final self-play has been conducted where every model played against each other without the help of MCTS to see their raw performance and a clear progression can be seen in Figure 3.



*Figure 3 - Evolutionary Model Performance*

Further the most promising models have been compared using also MCTS in various configurations to see how they perform then as can be seen in the Figures 4-8 in the appendix. While all perform subjectively well it is hard to make out a best model as their performance varies depending on the configurations, they are used in. Since the results of the self-plays of the top 11 models is not very conclusive, we will submit the latest model to compete against our class-mates results.

# 4. References

## 4.1. Bibliography

[1] '*Hex* (board game)', *Wikipedia*. Dec. 13, 2022. Accessed: Dec. 18, 2022. [Online]. Available: https://en.wikipedia.org/w/index.php?title=Hex_(board_game)&oldid=1127241285

[2] D. Silver *et al.*, 'Mastering the game of Go without human knowledge', *Nature*, vol. 550, no. 7676, pp. 354–359, Oct. 2017, doi: 10.1038/nature24270.

[3] D. Silver *et al.*, 'Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm'. arXiv, Dec. 05, 2017. Accessed: Dec. 18, 2022. [Online]. Available: http://arxiv.org/abs/1712.01815

[4] R. S. Sutton and A. G. Barto, *Reinforcement learning: an introduction*, Second edition. Cambridge, Massachusetts: The MIT Press, 2018.

[5] 'DeepMind'. https://www.deepmind.com/about (accessed Dec. 18, 2022).

[6] 'Slurm Workload Manager - Documentation'. https://slurm.schedmd.com/documentation.html (accessed Dec. 18, 2022).

[7] 'suragnair/alpha-zero-general: A clean implementation based on AlphaZero for any game in any framework + tutorial + Othello/Gobang/TicTacToe/Connect4 and more'. https://github.com/suragnair/alpha-zero-general (accessed Dec. 18, 2022).

[8] 'Reversi', *Wikipedia*. Dec. 16, 2022. Accessed: Dec. 18, 2022. [Online]. Available: https://en.wikipedia.org/w/index.php?title=Reversi&oldid=1127715743#Othello

[9] 'Tic-tac-toe', *Wikipedia*. Dec. 14, 2022. Accessed: Dec. 18, 2022. [Online]. Available: https://en.wikipedia.org/w/index.php?title=Tic-tac-toe&oldid=1127337707

[10] A. Agha, 'alpha-zero-hex'. May 27, 2019. Accessed: Dec. 18, 2022. [Online]. Available: https://github.com/aebrahimian/alpha-zero-hex

[11] C. B. Browne *et al.*, 'A Survey of Monte Carlo Tree Search Methods', *IEEE Trans. Comput. Intell. AI Games*, vol. 4, no. 1, pp. 1–43, Mar. 2012, doi: 10.1109/TCIAIG.2012.2186810.

[12] 'DeepLearning'. https://www.deeplearningbook.org/contents/intro.html (accessed Dec. 18, 2022).

## 4.2. Images:

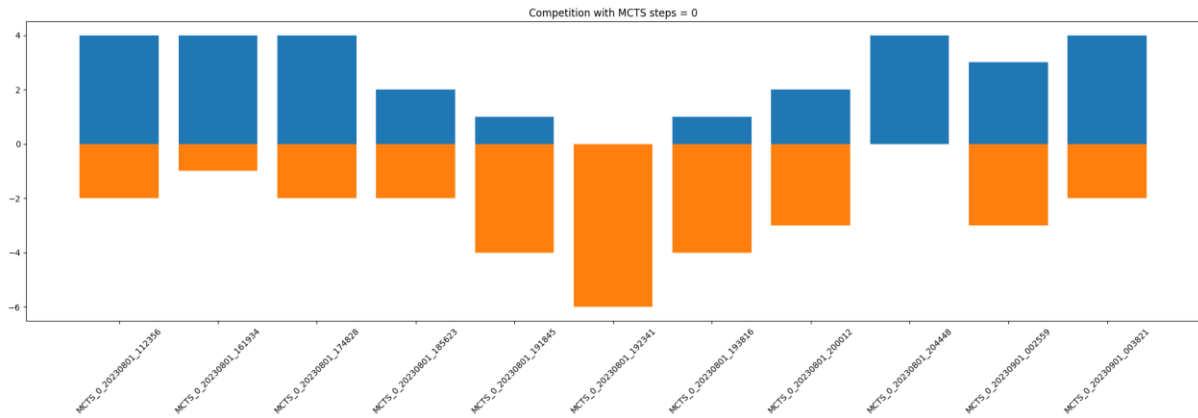# 5. Appendix

## 5.1.  Further Model Evaluations



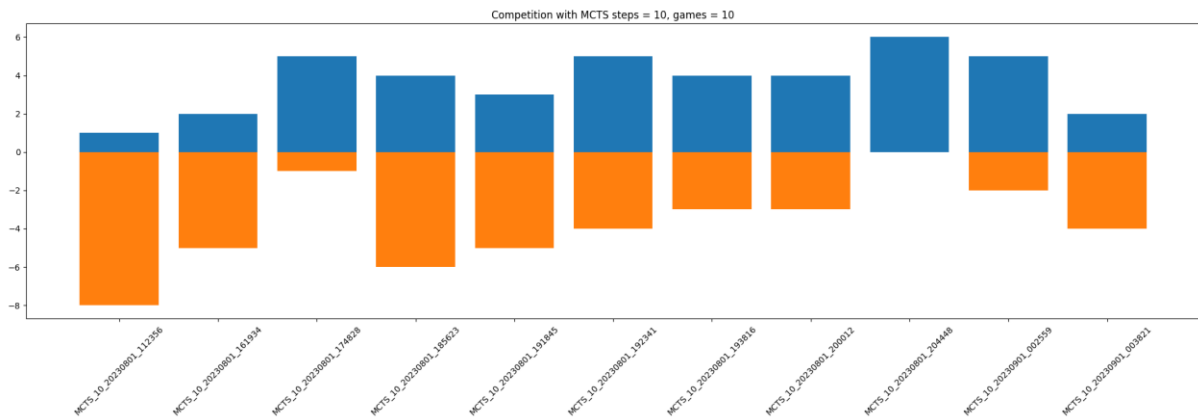*Figure 4 – Self-Play of the Top 11 Models without MCTS*



*Figure 5 - 10 Self-Play games of the Top 11 Models with MCTS 10 iterations*
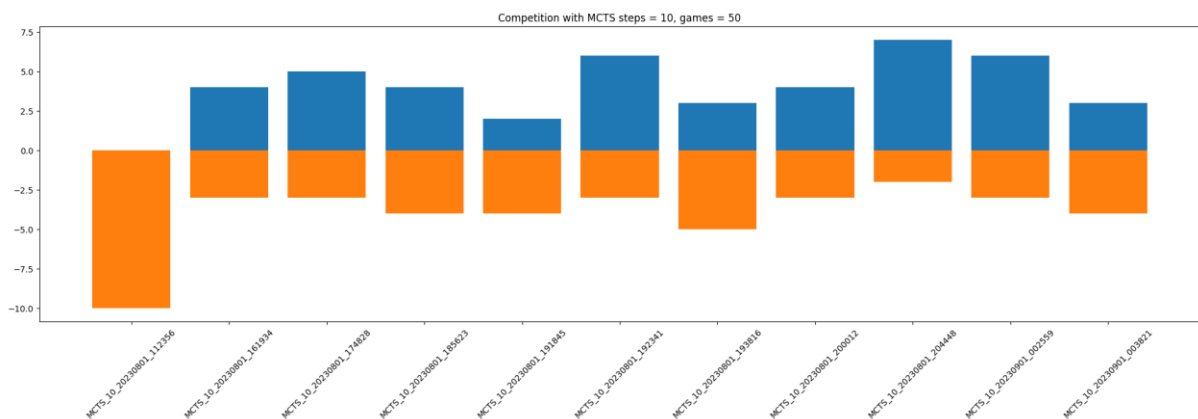


*Figure 6 - 50 Self-Play games of the Top 11 Models with MCTS 10 iterations*
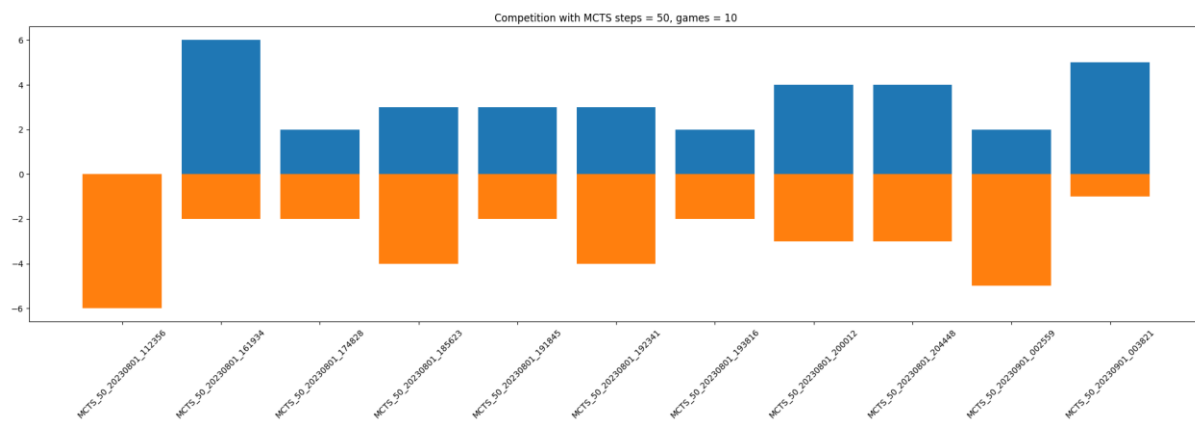
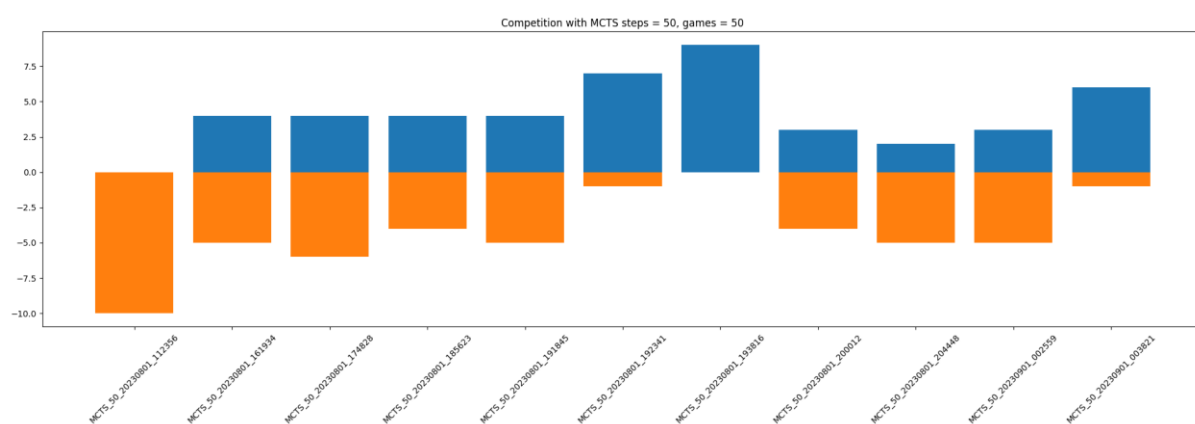*Figure 7 - 10 Self-Play games of the Top 11 Models with MCTS 50 iterations*



*Figure 8 - 50 Self-Play games of the Top 11 Models with MCTS 50 iterations*