# Complete Chess Games Enable LLM Become A Chess Master

**Yinqi Zhang, Xintian Han, Haolong Li, Kedi Chen, Shaohui Lin***

## Abstract

Large language models (LLM) have shown remarkable abilities in text generation, question answering, language translation, reasoning and many other tasks. It continues to advance rapidly and is becoming increasingly influential in various fields, from technology and business to education and entertainment. Despite LLM's success in multiple areas, its ability to play abstract games, such as chess, is underexplored. Chess-playing requires the language models to output legal and reasonable moves from textual inputs. Here, we propose the Large language model ChessLLM to play full chess games. We transform the game into a textual format with the best move represented in the Forsyth-Edwards Notation. We show that by simply supervised fine-tuning, our model has achieved a professional-level Elo rating of 1788 in matches against the standard Elo-rated Stockfish when permitted to sample 10 times. We further show that data quality is important. Long-round data supervision enjoys a 350 Elo rating improvement over short-round data.

## 1 Introduction

Recently, Large Language Models (LLMs) based on transformer architectures (Vaswani et al., 2017) have demonstrated capabilities well beyond language modeling. A key milestone was the advent of ChatGPT (Ouyang et al., 2022). Extensive research has focused on developing efficient LLM base models (Du et al., 2021; Biderman et al., 2023; Black et al., 2022; Computer, 2023; Touvron et al., 2023a), including supervised models (Taori et al., 2023a; Chiang et al., 2023; Anand et al., 2023; Köpf et al., 2023) and models using Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2017; Ouyang et al., 2022; Rando and Tramèr, 2023; Bai et al., 2023). Recent research (Wei et al., 2022; Li et al., 2024) shows that as models scale, their capabilities increase. This raises questions about LLMs' intelligence and

learning structures. Chess, an ancient game, has dialogue-like characteristics in its notational structures such as Forsyth-Edwards Notation (FEN), Standard Algebraic Notation (SAN), and Universal Chess Interface (UCI). Machine learning in chess has evolved to include reinforcement learning and neural networks based on supervised learning from human gameplay. Developments include AI-based engines like Leela Chess Zero (LC0)[1] and Stockfish NNUE[2], which refine their algorithms through new learning. Deep learning has shown the potential of AI in strategic games. The ChessGPT (Feng et al., 2023) model demonstrated the ability to choose optimal moves by learning from human language and chess data. However, models like ChessGPT cannot generate the best move based on the current game state and complete an entire match. Our focus is on match completeness and quality of gameplay.

Our contributions can be listed as follows:

- **Dataset.** We collected a large dataset of chess games with over 20B tokens from open-source platforms. Data quality matters; long round data supervision outperforms short-round data by 350 Elo points.

- **Model.** Our ChessLLM is designed to play entire chess games through dialogues. After fine-tuning, it achieved an Elo rating of 1788, winning 61% of games against Stockfish at skill level 0, 56% at skill level 1, and 30% at skill level 2.

- **Eval Method.** We propose evaluation methods based on full games against Stockfish, including move validity, Elo rating, and win rate. We are the only ones using a large language model for chess that can complete full games.

---

*Corresponding Author

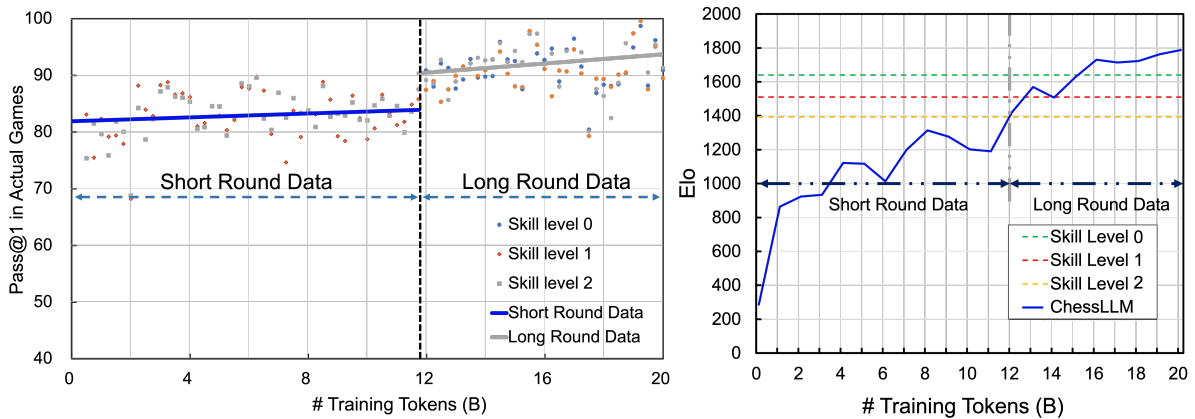[1] https://lczero.org
[2] https://stockfishchess.org

Figure 1: **Left:** $pass@1$ increases with the number of tokens. After introducing long-round data, $pass@1$ further increases. **Right:** The Elo Rating of ChessLLM with the number of training tokens. Skill level indicates the level of Stockfish.

## 2 Related work

### 2.1 Large Language Model

The emergence of Large language models (LLMs) GPT-4 (Achiam et al., 2023), stands as a noteworthy testament to the significant advancements in natural language understanding and generation. Unlike commercial models, open-source models, such as Alpaca (Taori et al., 2023b), Vicuna (Zheng et al., 2023), and Llama2 (Touvron et al., 2023b), have recently become more accessible. Due to their proficiency in text reasoning, LLMs are increasingly being utilized in everyday applications(Chen et al., 2024). Comprehensive benchmarks, such as MMLU (Hendrycks et al., 2021) and HELM (Lee et al., 2023), have been devised for thorough assessments of the LLMs' overall capabilities. Our work takes this evaluation process one step further, particularly highlighting and investigating the capacity of LLMs' ability to play abstract games.

### 2.2 Supervised Fine-tuning

Supervised Fine-tuning has emerged as a revolutionary technique within the field of machine learning and has been the subject of a multitude of studies. Owing to the continuous advancements in the domain of transfer learning, pre-trained models, fine-tuned in a supervised manner, have demonstrated superior performance in numerous tasks. Notably, in the context of natural language processing (NLP), the work by Howard and Ruder became a pioneering model of this technique. Their method (Howard and Ruder, 2018) leverages the power of transfer learning for comprehensive language modeling tasks, thus effectively surpassing previous benchmarks. Manipulating the same concept, BERT (Devlin et al., 2019), an innovative model fine-tuned in a supervised manner for a wide

array of NLP tasks. BERT demonstrated remarkable success within various NLP tasks, setting new performance standards.

In this work, we trained ChessLLM with supervised fine-tuning.

### 2.3 Chess

The quest to develop artificial intelligence capable of playing chess can be traced back to the inception of computer science (Turing, 1953; Campbell et al., 2002). The application of machine learning, particularly deep learning, in the domain of chess has been explored extensively in recent years (Silver et al., 2018; McGrath et al., 2022). One of the pivotal works in this field is the study by DeepChess (David et al., 2016), which presented an end-to-end learning method for chess based solely on deep neural networks, demonstrating the powerful capabilities of machine learning in comprehending and mastering strategic games without a priori knowledge.

In this work, we applied LLMs to chess and evaluated them with Elo rating.

## 3 A Large Scale Dataset of Chess

We introduce a large-scale dataset by collecting chess games online and generating the best moves based on Stockfish's evaluations. Previous research relied on Portable Game Notation (PGN) for strategy learning, interpreting moves as actions in a Markov Decision Process. ChessGPT sees additional value in PGN data, such as Elo ratings indicating player strength and annotated moves providing computer-generated evaluations. These annotations aid in value function learning, thus ChessGPT retains all this information for easier strategy learning. We argue that the core of chess is making the best decision for a given Forsyth-Edwards Notation

(FEN) position. Human players focus on the current position rather than past moves. While Chess-GPT uses historical moves, formats like PGN can be inefficient for large language models (LLMs) due to their expanding token length. The FEN format remains constant, making it more suitable for LLMs. Therefore, we constructed our dataset as FEN-Best move pairs.

**Best Move Construction** Our Best Move dataset was created through a search method using Stockfish. It consists of two parts: the short round dataset from Chessdb[3] and the long round dataset from self-play endgames based on Stockfish evaluations. Stockfish evaluates positions using heuristic functions and an alpha-beta game tree search. We searched for valid moves from current positions, with search depths of 12-50 for short rounds and 50-200 for long rounds, limiting each search to two seconds. The highest win-rate moves were selected as the best moves.

## 4 Model

The Generative Pre-trained Transformer (GPT-3) is an autoregressive language model that generates human-like text through deep learning. It trains on casual language modeling, predicting the next word based on previous words. We trained a GPT-like model using open-llama-3B (Geng and Liu, 2023) and the chess resources from Section 3. Unlike policy behavior data in robotics or gaming, chess state and move data can be expressed textually. This allows chess to be rendered as a text-based game, enabling imitation learning for policy through casual language modeling of the game dataset (Figure 2). This innovative approach of applying language modeling to chess signifies a novel shift in policy learning, leveraging the game's unique aspects to develop superior gameplay tactics.
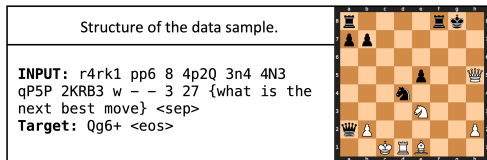


Figure 2: One example of training data.

## 5 Evaluation Methods

Chess requires a dynamic evaluation method beyond a fixed set typical of NLP tasks. We propose

---

[3]http://chessdb.sourceforge.net

supplementing the evaluation set with actual games to better assess the model's capabilities.

### 5.1 Actual Games

Playing against Stockfish, a top chess engine, offers a strategic challenge. Stockfish uses advanced algorithms to determine optimal moves. Players can choose time controls (blitz, quick, or traditional) to set the gameplay tempo. The engine analyzes moves and positions to find the best move using its evaluative function. In our experiments, we analyzed metrics such as $pass@1$ and win rate. We believe using Stockfish against our model more authentically simulates real-world human-model interactions and offers greater robustness than a static evaluation set.

**Pass@1 in Actual Games.** We evaluated our model's performance across different data scales, focusing on its ability to generate legal moves successfully.

**Win Rating.** The win rating refers to victories, draws, and losses out of 100 rounds when the model competes against Stockfish or other engines.

**Elo Rating.** We ran a series of matches between our model and Stockfish, recording strategies and moves. The Elo rating is calculated using the formula

$$Elo_N = Elo_O + (R_A - R_E)K, \qquad (1)$$

$$R_E = \frac{1}{1 + 10^{\frac{Elo_S - Elo_M}{400}}}. \qquad (2)$$

where $Elo_N$ is the updated Elo rating after the game. $Elo_O$ is the previous Elo rating before the game. $K$ is the weight of the tournament. In professional chess, $K$ is often set to 10 for high-ranked players and 20 for low-ranked players. $R_A$ is the actual result of the game (1 for win, 0.5 for draw, 0 for loss). $R_E$ is the expected result of the game. $Elo_S$ is the old Elo rating of Stockfish. $Elo_M$ is the old Elo rating of the model. Moreover, we refer to the method introduced by Stockfish to convert between its skill level and Elo rating. The specific calculation method is shown as follows.

$$SK = 37.247e^3 - 40.852e^2 + 22.294e - 0.311, \qquad (3)$$

$$e = \frac{Elo - 1320}{1870}, \qquad (4)$$

where $SK$ represents skill level $SK = 0, 1, 2, ..., 20$, and $Elo$ represents Stockfish's Elo rating.
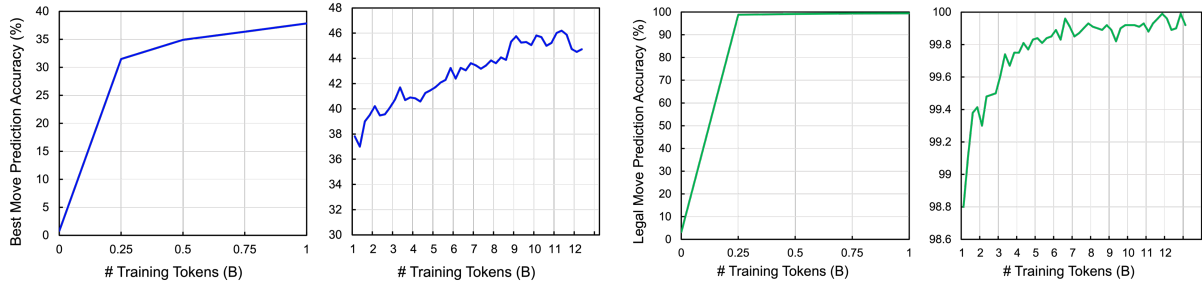
Figure 3: **Left:** Best Move Accuracy of ChessLLM training with short round data. The accuracy of the best move increases with the number of training tokens. **Right:** Legal Move Accuracy of ChessLLM training with short round data. The accuracy of the legal move increases with the number of training tokens.
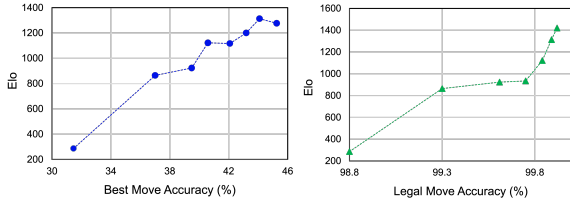


Figure 4: **Left:** Correlation between ChessLLM's best move accuracy and its Elo rating. **Right:** Correlation between ChessLLM's legal move accuracy and its Elo rating.

## 5.2 Evaluation Set

While games against Stockfish provide a robust performance assessment, their length introduces substantial evaluation costs. Thus, we also use an evaluation set to measure the model's prowess. Data distribution in the evaluation set focuses on games spanning 10-20 rounds (30%) and 20-40 rounds (50%), emphasizing the model's middle-game capabilities. This approach manages the inherent uncertainty in chess match lengths, ensuring the model does not exhibit forgetting phenomena after exposure to long rounds.

**Distribution of Training set and Evaluation set** Our training data was generated with $depth = 1$ and $timelimited = 0.1$, while the data used in the game process was generated with $timelimited = 10$ and without depth limited. The eval set is produced by $depth = 1$ and $timelimited = 0.1$, the same as the train set. These two datasets are from different domains, so our method is effective not only on in-domain data.

**Legal Move Accuracy.** We used Stockfish to generate legal move responses for 10,000 unique board positions not in the training set, evaluating our model's proposed moves for legality to ensure proper convergence.

**Best Move Accuracy.** Stockfish generated best move responses, allowing us to compare its outcomes with our model to calculate the accuracy

Table 1: Exhibition of Match Results and Computed Elo Scores of ChessLLM *vs.* Stockfish at Different Skill Levels. The table enumerates the number of wins, losses, and draws, along with the calculated Elo scores of ChessLLM when competing against Stockfish at varying skill levels.

| | Stockfish | | ChessLLM | | | |
|---|---|---|---|---|---|---|
| | Skill level | Elo | Win | Lose | Draw | Elo |
| ChessLLM | 0 | 1350-1440 | 61 | 29 | 10 | $1632 \pm 45$ |
| *vs.* | 1 | 1450-1560 | 56 | 37 | 7 | $1753 \pm 55$ |
| Stockfish | 2 | 1570-1720 | 30 | 69 | 1 | $1788 \pm 75$ |

Table 2: General policy evaluation in Black. Note LLAMA denotes the LLAMA-7B

| Elo Rating | Move Scores (%) | | | | |
|---|---|---|---|---|---|
| | LLAMA | RedPajama | ChessGPT-Base | ChessGPT-Chat | ChessLLM |
| 700-1000 | $52.9 \pm 0.9$ | $46.2 \pm 1.0$ | $51.9 \pm 0.1$ | $52.1 \pm 0.9$ | $\mathbf{90.96 \pm 1.4}$ |
| 1200-1500 | $53.2 \pm 0.9$ | $46.9 \pm 0.9$ | $53.0 \pm 1.0$ | $52.4 \pm 1.0$ | $\mathbf{95.11 \pm 0.8}$ |
| 1700-2000 | $52.1 \pm 0.8$ | $46.6 \pm 1.0$ | $52.0 \pm 1.0$ | $52.0 \pm 1.0$ | $\mathbf{96.88 \pm 0.9}$ |
| 2700-3000 | $53.6 \pm 0.9$ | $47.3 \pm 1.0$ | $52.2 \pm 0.9$ | $52.1 \pm 1.1$ | $\mathbf{97.14 \pm 0.6}$ |

rate for best move predictions.

## 6 Experiment Analysis

### 6.1 Evaluation Set

We evaluated in-distribution data to analyze our model's performance on the evaluation set under varying computing power. From Fig. 3, we observed that on in-distribution data, model performance improves with an increase in training tokens, but at a diminishing rate. This relationship is crucial for understanding model scalability and resource allocation during training. Note that "same distribution" refers to the FEN board state distribution and its corresponding best move.

**Legal Move and Best Move Accuracy.** Fig. 3 Left shows that with only 0.5B tokens, our model achieves a legal move accuracy of 99.11% on in-distribution boards, indicating its impressive preliminary chess playing ability. As data volume increases, performance improves, demonstrating the model's scalability and potential for further enhancement. The high accuracy with just 0.5B tokens underscores the model's efficiency and ef-

Table 3: The win rates of various LMs when competing in Chess. Note LLAMA denotes the LLAMA-7B.

| | LLAMA | RedPajama | ChessGPT-Base | ChessGPT-Chat |
|---|---|---|---|---|
| LLAMA | - | - | - | - |
| RedPajama | 22.2 ± 4.2 | - | - | - |
| ChessGPT-Base | 61.3 ± 2.4 | 73.6 ± 1.1 | - | - |
| ChessGPT-Chat | 59.8 ± 1.5 | 70.8 ± 0.7 | 48.8 ± 2.7 | - |
| ChessLLM(Ours) | **89.8 ± 0.8** | **95.5 ± 0.1** | **91.7 ± 0.3** | **92.3 ± 0.1** |

fectiveness. Fig. 3 Right shows the Best Move accuracy under the same distribution. With 2.75B tokens, the model achieved a Best Move accuracy of 40.11%. Although the logic is similar, the generation steps differ, highlighting our model's ability to accurately predict the best moves in most cases, proving its practical utility.

## 6.2 Actual Games

**Pass@1 in Actual Games.** The $temperature$ and $top_p$ parameters were both set at 1.0, and $top_k$ was set at 50 we generated once to calculate Pass@1. Matches against Stockfish, using only one sampling iteration per match, evaluated the legality of our model's moves. Figure 1 shows our model's results. Despite fluctuations from incorporating more endgame strategies, the model consistently achieves over 90% move legality. The legality remains stable against opponents of varying strengths.

**Elo rating.** Table 1 shows our model's performance in 100 rounds each against Stockfish at skill levels 0, 1, 2etc., computing Elo ratings. With $temperature$ and $top_p$ parameters were both set at 0.7, and $top_k$ was set at 50. we used up to 10 sampling iterations, performing the move upon obtaining a legal one. Our model achieves an Elo score of about 1788, positioning it at the top of amateur chess performance.

## 6.3 Eval Set Accuracy and Actual Games

Figure 4 shows that within the evaluation set, an increase in Best Move accuracy correlates with Elo rating gains. A significant Elo rating jump occurs when the model's Legal Move accuracy reaches 99.8%. This increase is due to the reduction in errors after the model learns to generate legal moves, reinforcing that continuous error correction and learning the correct moves significantly improve Elo ratings.

## 6.4 Compare with Other LMs

**General Policy.** General Policy is proposed by ChessGPT (Feng et al., 2023). Table 2 showcases

the results, delineating the effectiveness of various models in identifying the most fitting move for the black chess piece.

**Win Rating.** We conduct matches between ChessLLM and other Language Models (LMs) such as LLAMA (Touvron et al., 2023a), RedPajama (Computer, 2023), ChessGPT-Base (Feng et al., 2023), and ChessGPT-Chat (Feng et al., 2023), calculating their respective win rating. As other models cannot guarantee the legality of the moves they generate, we bring in Stockfish to aid in this process. Should the model fail to produce a valid move even after 50 sampling efforts, a mechanism is employed wherein there's a 50% chance of favoring either the best move identified by Stockfish or a randomly picked move from the list of all possible legal moves. Similarly, as ChessGPT is unable to generate the best move for the next step, we generate all legal moves through Stockfish and utilize their proposed general policy for selection, picking the most optimal move as recognized by the model.

## 6.5 Impact of Token Quantity and Quality

We have investigated the impact of data quantity and quality on the generation of legal moves. Figure 1 Left presents the $Pass@1$ indicators for two groups of data. It can be observed that the model performance significantly improves with the addition of more high-quality data, supplementing the data beyond the original distribution. Figure 1 Right presents an augmentation in the number of tokens, it is observed that the model's Elo rating experiences an enhancement. Concurrently, the enrichment of the model with data not within the distribution can expedite the elevation of the model's Elo rating.

## 7 Conclusion

In this paper, we convert chess to a text game and introduce a large-scale Fen-Best Move pair dataset. With the dataset, we propose the Large language model ChessLLM that can play a complete chess game. Considering the limitation of the evaluation set in out-of-distribution data, we propose the need to evaluate model capabilities in actual games. ChessLLM finally achieves an Elo rating of 1788 through the SFT method. In subsequent work, we will discuss how to improve ChessLLM by improving the data quality.

## 8 Limitations

In this study, we explored the problem of LLM playing chess games and found that with high-quality synthetic data of complete games, LLM can have the extrapolation and combat capabilities of chess games. In the future, we will continue to explore this capability by improving the data quality, RLHF, and self-play + MCTS so that LLM can become better at chess games. Our ultimate goal is to enable LLM to excel in various games through high-quality game data.

## 9 Ethics Statement

In this research, we adhere to strict ethical guidelines and principles. The study has been designed and implemented with respect for the rights, privacy, and well-being of all individuals involved. Our findings and conclusions are reported accurately and objectively, avoiding any misrepresentation or manipulation of data. The entire process and outcomes are free from intellectual property and ethical legal disputes.

## Acknowledgments

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Yuvanesh Anand, Zach Nussbaum, Brandon Duderstadt, Benjamin Schmidt, and Andriy Mulyar. 2023. Gpt4all: Training an assistant-style chatbot with large scale data distillation from gpt-3.5-turbo. https://github.com/nomic-ai/gpt4all.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*.

Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit,

USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.

Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, et al. 2022. Gpt-neox-20b: An open-source autoregressive language model. *arXiv preprint arXiv:2204.06745*.

Murray Campbell, A Joseph Hoane Jr, and Feng-hsiung Hsu. 2002. Deep blue. *Artificial intelligence*, 134(1-2):57–83.

Kedi Chen, Qin Chen, Jie Zhou, Yishen He, and Liang He. 2024. Diahalu: A dialogue-level hallucination evaluation benchmark for large language models. *Preprint*, arXiv:2403.00896.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.

Together Computer. 2023. Redpajama: an open dataset for training large language models.

Omid E. David, Nathan S. Netanyahu, and Lior Wolf. 2016. *DeepChess: End-to-End Deep Neural Network for Automatic Learning in Chess*, page 88–96. Springer International Publishing.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *Preprint*, arXiv:1810.04805.

Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2021. Glm: General language model pretraining with autoregressive blank infilling. *arXiv preprint arXiv:2103.10360*.

Xidong Feng, Yicheng Luo, Ziyan Wang, Hongrui Tang, Mengyue Yang, Kun Shao, David Mguni, Yali Du, and Jun Wang. 2023. Chessgpt: Bridging policy learning and language modeling. *arXiv preprint arXiv:2306.09200*.

Xinyang Geng and Hao Liu. 2023. Openllama: An open reproduction of llama.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.

Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. *Preprint*, arXiv:1801.06146.

Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, et al. 2023. Openassistant conversations–democratizing large language model alignment. *arXiv preprint arXiv:2304.07327*.

Tony Lee, Michihiro Yasunaga, Chenlin Meng, Yifan Mai, Joon Sung Park, Agrim Gupta, Yunzhi Zhang, Deepak Narayanan, Hannah Benita Teufel, Marco Bellagente, Minguk Kang, Taesung Park, Jure Leskovec, Jun-Yan Zhu, Li Fei-Fei, Jiajun Wu, Stefano Ermon, and Percy Liang. 2023. Holistic evaluation of text-to-image models. *Preprint*, arXiv:2311.04287.

Haolong Li, Yu Ma, Yinqi Zhang, Chen Ye, and Jie Chen. 2024. Exploring mathematical extrapolation of large language models with synthetic data. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 936–946, Bangkok, Thailand. Association for Computational Linguistics.

Thomas McGrath, Andrei Kapishnikov, Nenad Tomašev, Adam Pearce, Martin Wattenberg, Demis Hassabis, Been Kim, Ulrich Paquet, and Vladimir Kramnik. 2022. Acquisition of chess knowledge in alphazero. *Proceedings of the National Academy of Sciences*, 119(47):e2206625119.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

Javier Rando and Florian Tramèr. 2023. Universal jailbreak backdoors from poisoned human feedback. *arXiv preprint arXiv:2311.14455*.

David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. 2018. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023a. Stanford alpaca: An instruction-following llama model.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023b. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

Alan M Turing. 1953. Digital computers applied to games. *Faster than thought*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Preprint*, arXiv:2306.05685.