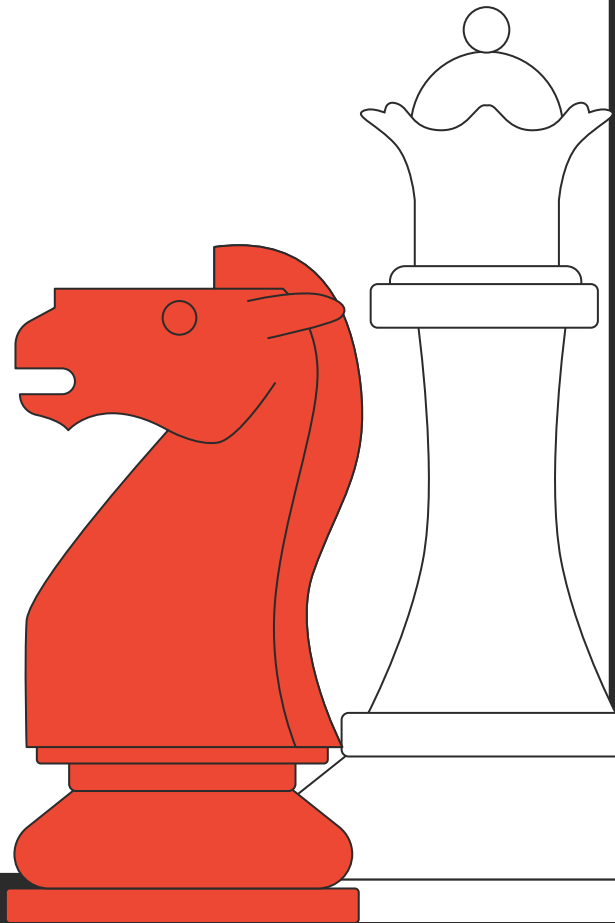




Probing a ChessLLM

ANLP - Project Presentation

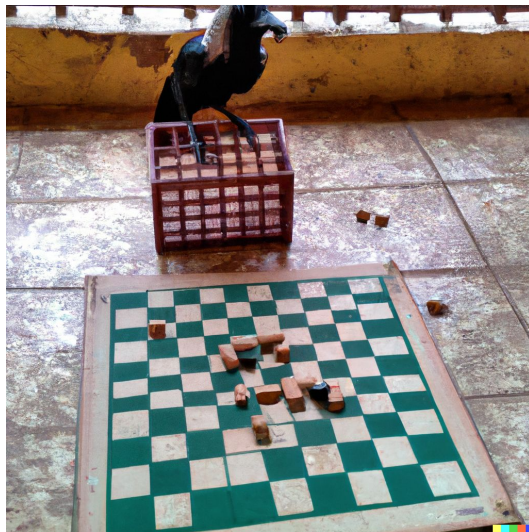
Silvano Vento Maddonni



Introduction

Idea:

- Take a LLM and explore how it represents the world state underlying the text in a task not (directly) related to NLP.
- Chess
 - Chess provides a simple, constrained, and deterministic domain where the exact world state is known.
 - Chess games can also be transcribed exactly and unambiguously using chess notations
 - The form of chess notations allows us to probe our language models for aspects of the board.



Chess as a Testbed for Language Model State Tracking

Shubham Toshniwal¹, Sam Wiseman², Karen Livescu¹, Kevin Gimpel¹

¹Toyota Technological Institute at Chicago

²Duke University

{shtoshni, klivescu, kgimpel}@ttic.edu, swiseman@cs.duke.edu

Watching a Language Model Learning Chess

Andreas Stöckl

University of Applied Sciences Upper Austria
Digital Media Department
Hagenberg, Austria

andreas.stoeckl@fh-hagenberg.at

Abstract

We analyse how a transformer-based language model learns the rules of chess from text data of recorded games. We show how it is possible to investigate how the model capacity and the number of training data influence the success of a language model with the chess-specific metrics. With these metrics, we show that more games used for training offer significantly better results than for the same training time. However, size does not show such a clear correlation. It is also interesting to observe that accuracy and perplexity, given by the model, reveal how they store information about the board state in the activations of the groups, and how the overfitting moves influence the model's performance.

model capacity is able to learn the rules of arithmetic to a certain degree by training it with the data of crawled websites. In the process, elementary operations were learned in a certain number of space, but not beyond. Is this limitation due to the lack of capacity of the model, insufficient training time or training data that did not contain sufficient information?

In Nogueira et al. (2021), it was demonstrated that regardless of the number of parameters and training examples, *Transformer* Vaswani et al. (2017) models are unable to learn addition rules seen during training. To test the ability of the model to learn rules, and assess its capacity to learn the rules of chess, we use a transformer-based language model trained solely on chess-related data. We probe the model's ability to learn the rules of chess by training it with the data of crawled websites. In the process, elementary operations were learned in a certain number of space, but not beyond. Is this limitation due to the lack of capacity of the model, insufficient training time or training data that did not contain sufficient information?

Introduction

Large language models have stretched the simple self-supervised training paradigm, becoming a fixture in state-of-the-art (Vaswani et al. 2017; De-

language modeling objective or introducing any new factors.¹

Due to the simplicity and precision of univariate language model predictions at the state level, rather than merely comparing the model's output to the ground truth move, the model's prediction is legal given the current state of the board.

Under review as a conference paper at COLM 2024

Emergent World Models and Latent Variable Estimation in Chess-Playing Language Models

Adam Karvonen

Independent

adam.karvonen@gmail.com

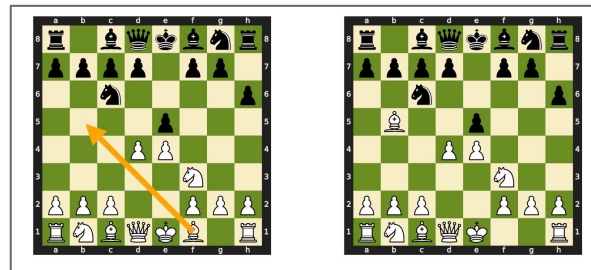
Abstract

Language models have shown unprecedented capabilities, sparking debate over the source of their performance. Is it merely the outcome of learning syntactic patterns and surface level statistics, or do they extract semantics and a world model from the text? Prior work by Li et al. investigated this by training a GPT model on synthetic, randomly generated Othello games and found that the model learned an internal representation of the board state, extending this work into the more complex domain of chess, training a GPT model on synthetic chess games. We investigate our model's internal representations of board state, and contrastive activations. The model is given no a priori information about the board state, and is solely trained on next character prediction. We find that the model's internal representations of board state, and contrastive activations, are able to make interventions on the board state, and predict the next move. Unlike Li et al., we find that the model also learns to predict the next move, and is able to make interventions on the board state, and predict the next move. We find that the model's internal representations of board state, and contrastive activations, are able to make interventions on the board state, and predict the next move. Unlike Li et al., we find that the model also learns to predict the next move, and is able to make interventions on the board state, and predict the next move.

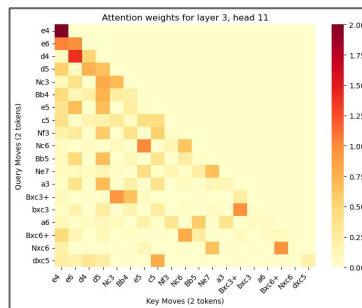
The project

The project is divided in three parts:

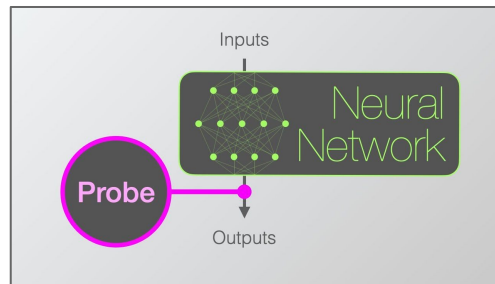
- Part 1: LLM next move prediction + legal moves



- Part 2: Attention Heatmaps

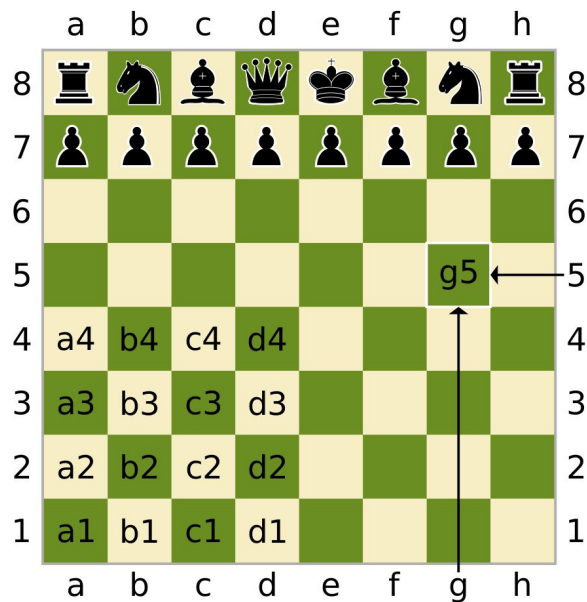


- Part 3: Probing Classifier for board pieces locations



Chess as text

- Chess games and moves can be transcribed in many different ways.
- Algebraic notation is the standard:





PGN vs UCI

- The sequence of moves composing a game are generally encapsulated using the Portable Game Notation (**PGN**).

1. e4 e5 2. Nf3 Nc6 3. Bb5 a6 4. Ba4 Nf6 5. 0-0
Be7 6. Re1 b5 7. Bb3 d6 8. c3 0-0 9. h3 Nb8
10. d4 Nbd7 11. c4 c6 12. cxb5 axb5 13. Nc3

- We represent moves using Universal Chess Interface (**UCI**) notation, which combines the starting square and the destination square to represent a move.

e2e4 e7e5 g1f3 b8c6 f1b5 a7a6 b5a4 g8f6 e1g1
f8e7 f1e1 b7b5 a4b3 d7d6 c2c3 e8g8 h2h3 c6b8
d2d4 b8d7 c3c4 c7c6 c4b5 a6b5 b1c3





Tokenizer

- The games represented in UCI notation are tokenized using a simple regular expression based tokenizer, which considers a board square symbol such as b1 as a single token.
- This produces a vocabulary of 77 token types, which includes the 64 squares, piece type symbols, and other special symbols.

Type	Examples	Count
Square names	e4, d1	64
Piece type	P, K, Q, R, B, N	6
Promoted	q, r, b, n	4
Special symbols	BOS, EOS, PAD	3
Total		77

Model Vocabulary



Next Move Prediction + Legal Moves

- The starting model used is a GPT2-small. 12-layers, 12 heads, embedding 768.
- The model is pre-trained on 250K chess games taken from a full dataset of 2.9 million quality chess games. The whole training was done using UCI notation.

PGN seq:

1.Nf3 f5 2.b4 Nf6 3.Bb2 e6 4.b5 d5 5.g3 Bd7 6.a4
a6 7.c4 dxc4 8.Nc3 axb5 9.axb5

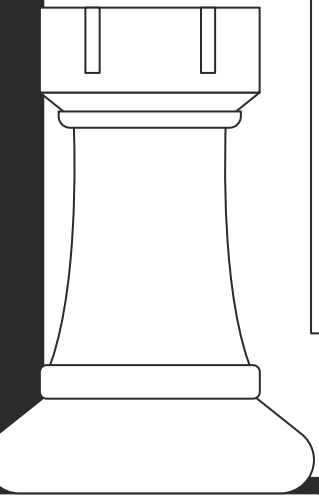
UCI seq:

g1f3 f7f5 b2b4 g8f6 c1b2 e7e6 b4b5 d7d5 g2g3
c8d7 a2a4 a7a6 c2c4 d5c4 b1c3 a6b5 a4b5

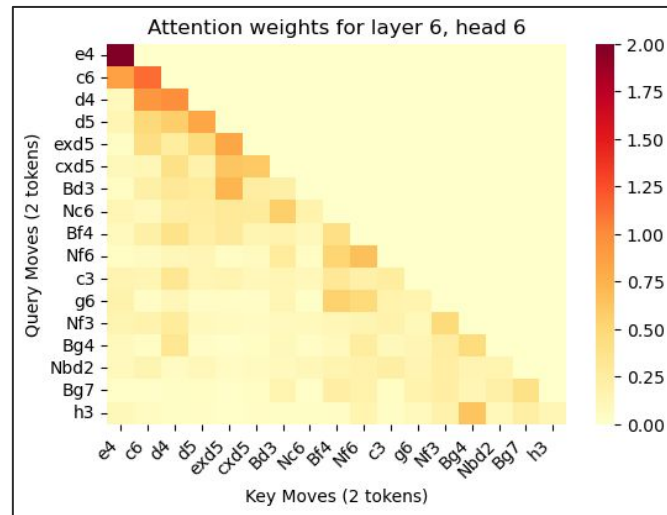
LM plays: a8a1

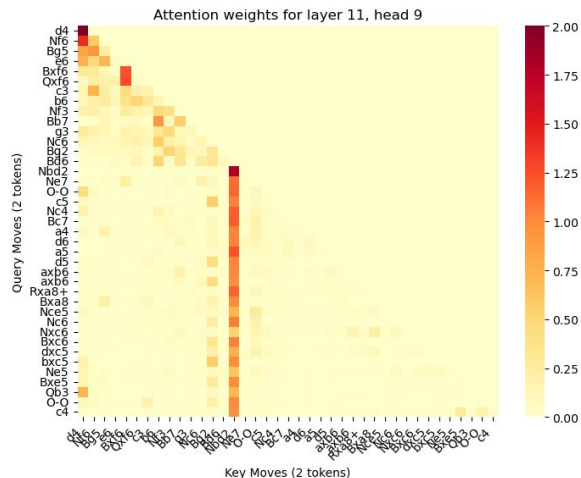
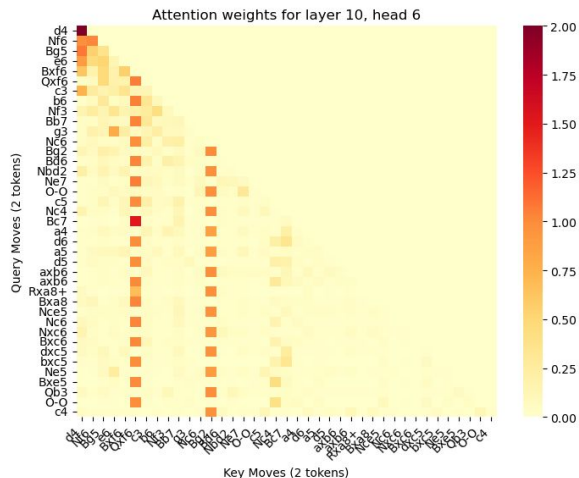
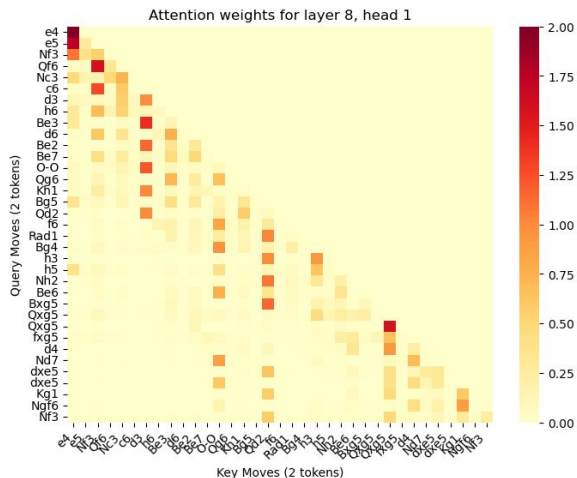
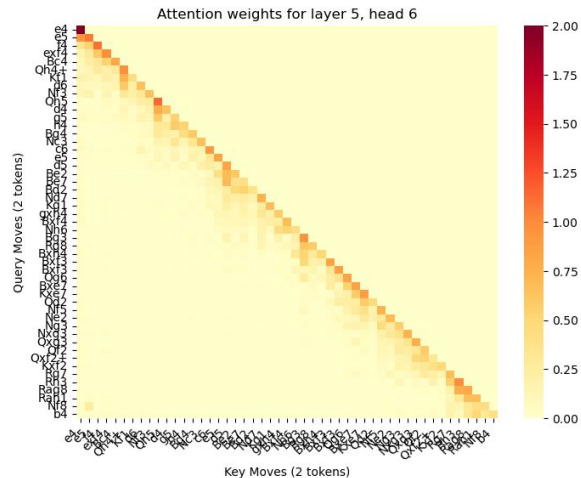
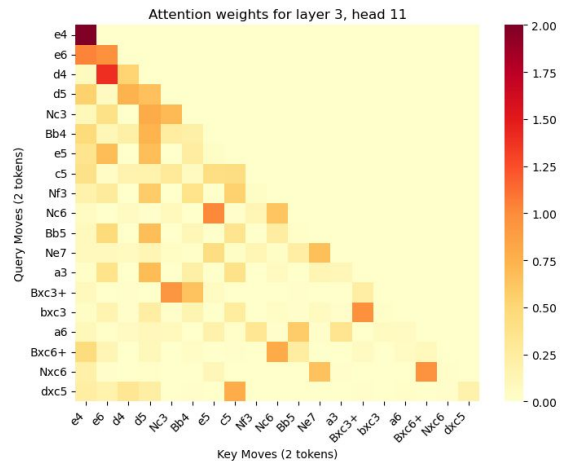
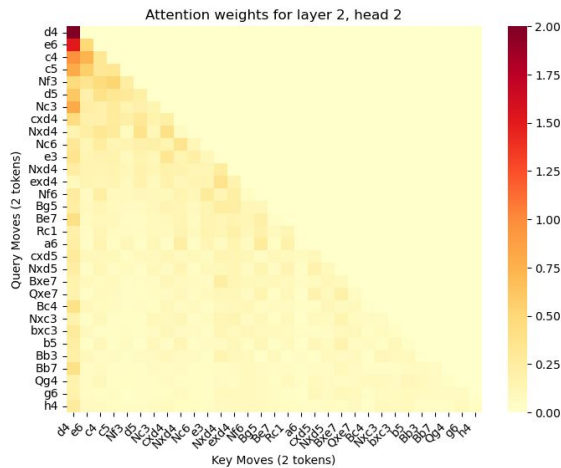
(In the game it was: a8a1)

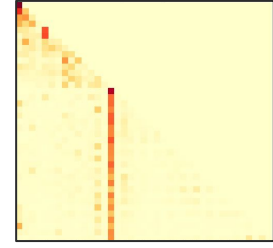
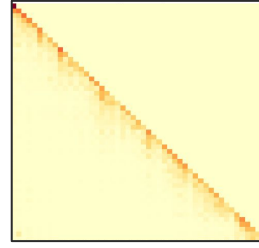
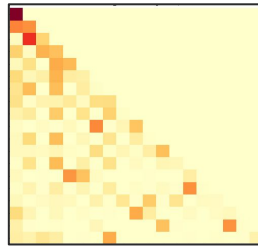
- Legal moves → 96% ratio.



-







Final Prediction = [1 to 4] + [5 to 7] + [8 to 12]

[layers 1-4] = General and Sparse

[layers 5-7] = Reinforced Diagonal

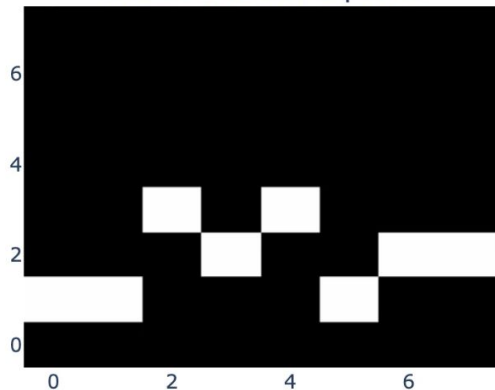
(local connectivity of moves)

[layers 8-12] = Single moves attention

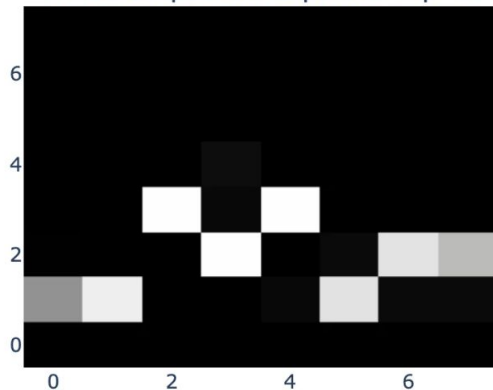
Internal Model World Representation

- Investigate the internal representations of language models in a more constrained setting.
- Train linear probes that recover the internal board state,

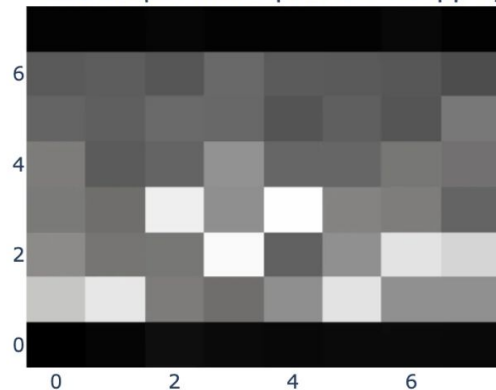
Chess board white pawns



Probe output white pawns clip=5



Probe output white pawns no clipping



Internal Model World Representation

PiecePositionProbe:

- 768x64 linear classifier model with a sigmoid that outputs a 64 elements tensor with probabilities for each square.
- We take the model with 'frozen' weights and extract the features from a particular layer (often the last hidden layer).
- Then, we train the linear classifier on top of these features to see how well it can perform on this task.

Piece	Base	Trained Probe
Pawn	0.54	0.82
Rook	0.55	0.94
Knight	0.51	0.97
King	0.45	0.93

Classifier Probe



Thanks!

Silvano Vento Maddonni

