# Exploring the Limits of OthelloGPT's Emergent Representations

Julian Baldwin

# OVERVIEW

➤ Motivations
➤ Background
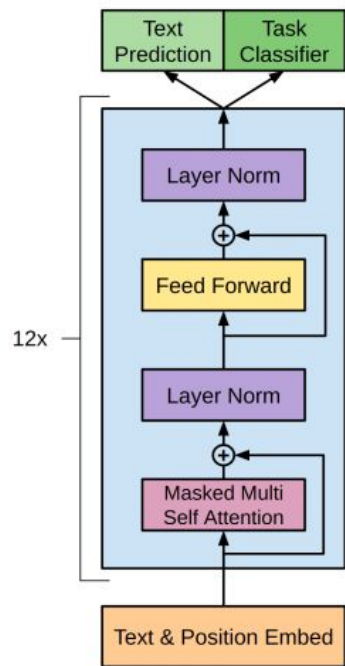  ○ Othello-GPT
➤ Results
➤ High-Level Takeaways

# Broad Goals / Motivations

**Mechanistic Interpretability:** the study of attempting to reverse engineer neural networks to human interpretable algorithms

Current techniques for steering LLM outputs are limited, difficult to make any rigorous guarantees about model behavior

Sub-question: how do transformer models represent concepts?
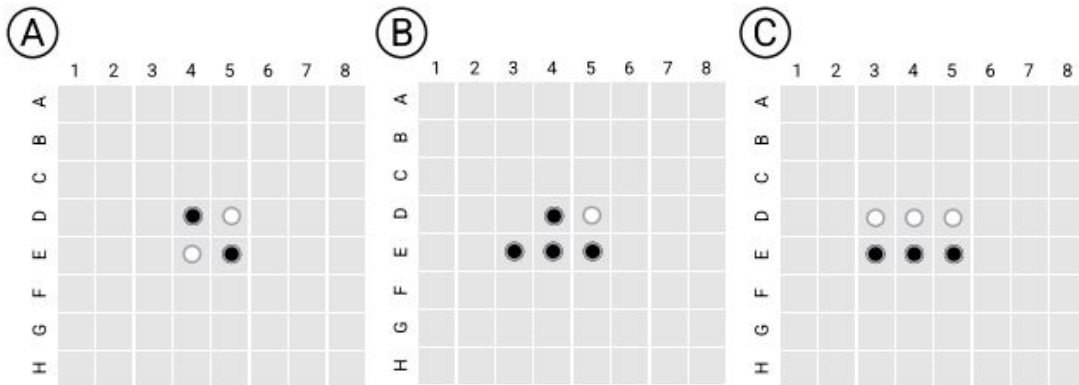
# OTHELLO-GPT



Used the toy task of Othello as a testbed; model is trained on randomly generated games to predict legal moves

Even with no a priori knowledge of the rules or game board, model learns an interpretable representation of the board state
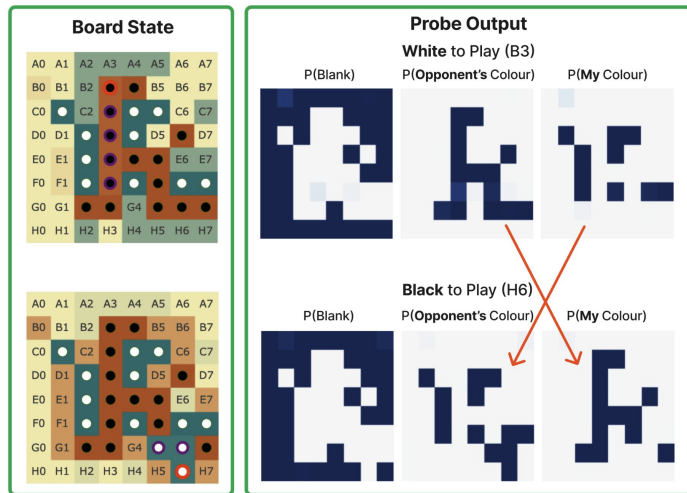
We can extract this representation from the transformer's residual stream using a probe—a classifier trained on model activations

# LINEAR REPRESENTATION



Othello-GPT's Linear Model of Board State

Input: F4 F3 D2 F5 G2 F2 G3 C4 E5 F6 D6 E2 B4 C5 G7 C1 G6 F7 G5 C3 **B3 H6**
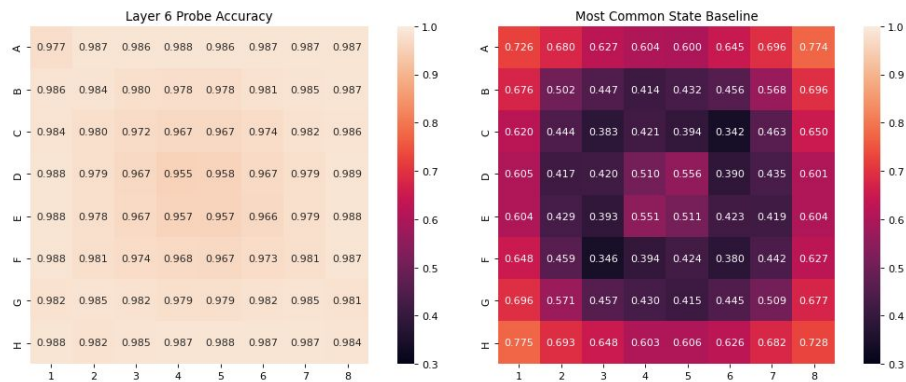
Published by Neel Nanda in March 2023

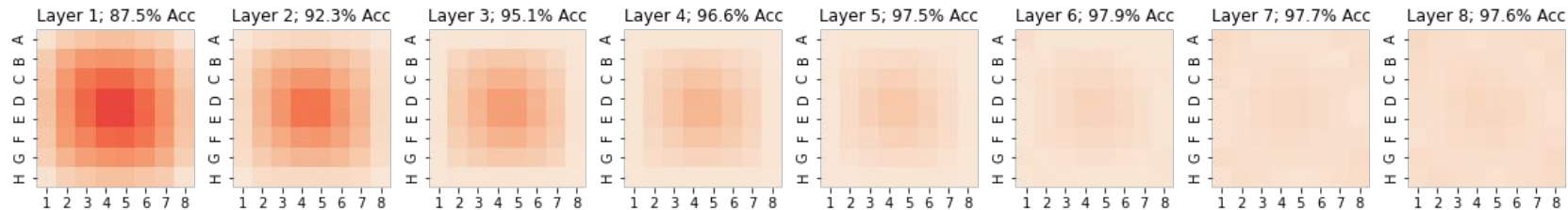Demonstrates that board state can be recovered with linear probes

View board as empty/mine/theirs rather than empty/black/white

# RECREATING LINEAR PROBES



Layer 6 Probe Accuracy
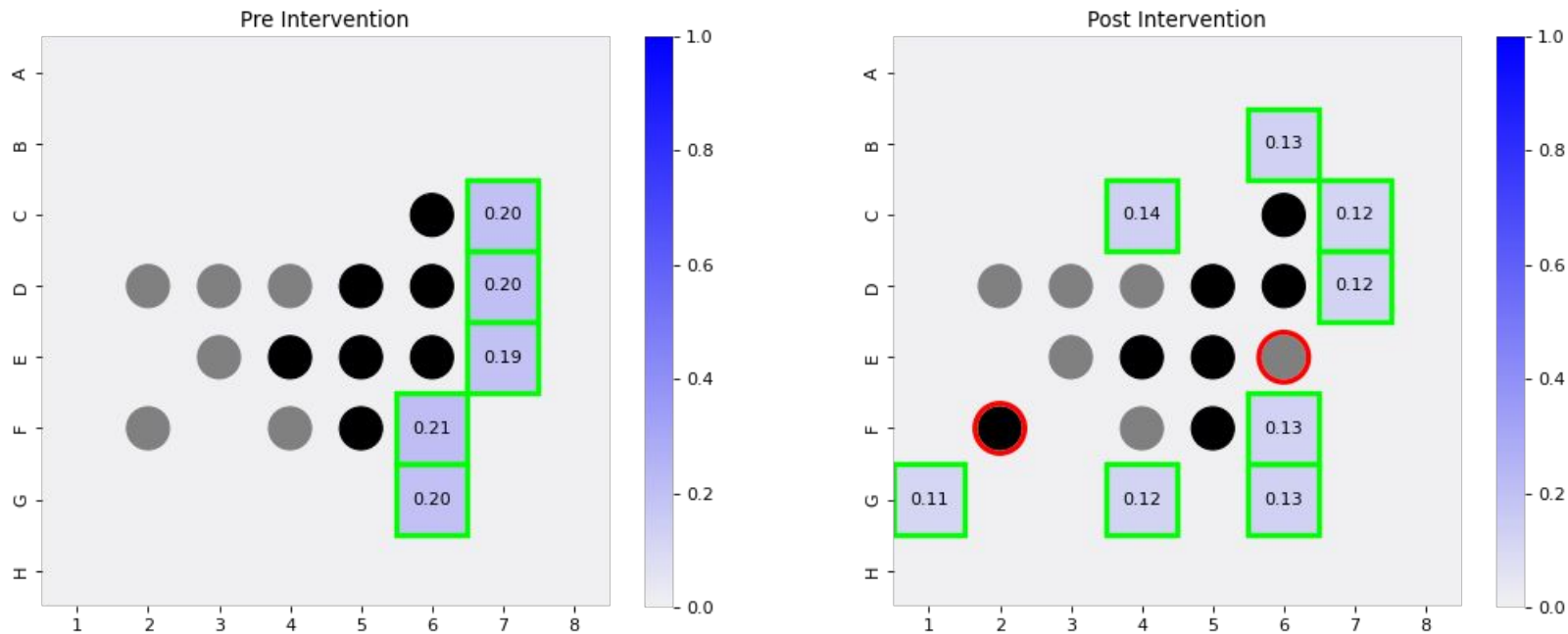
Most Common State Baseline

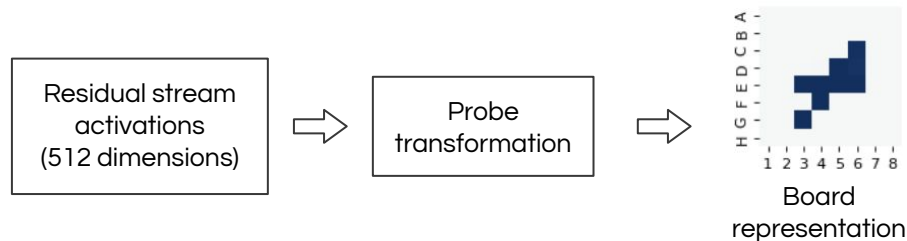Starting from the original OthelloGPT codebase, I trained a new model on randomly generated othello sequences

Recreated Neel Nanda's work, achieved high accuracy of board prediction with linear probes



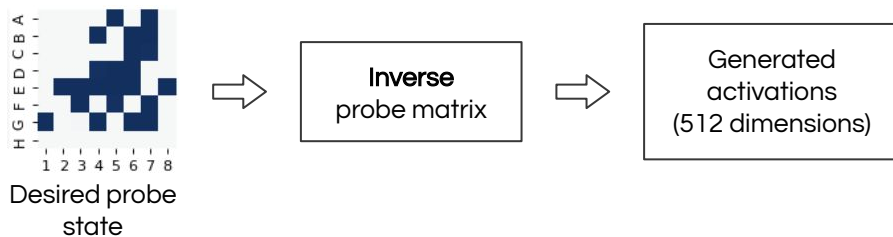Layer 1; 87.5% Acc    Layer 2; 92.3% Acc    Layer 3; 95.1% Acc    Layer 4; 96.6% Acc    Layer 5; 97.5% Acc    Layer 6; 97.9% Acc    Layer 7; 97.7% Acc    Layer 8; 97.6% Acc

# LOCAL LINEAR INTERVENTIONS



Flipping the color of two board squares, layer 4 linear intervention

# Global Intervention via Probe inverse



Residual stream activations (512 dimensions) → Probe transformation → Board representation

Usual procedure for extracting board representation from model

Desired probe state → Inverse probe matrix → Generated activations (512 dimensions)
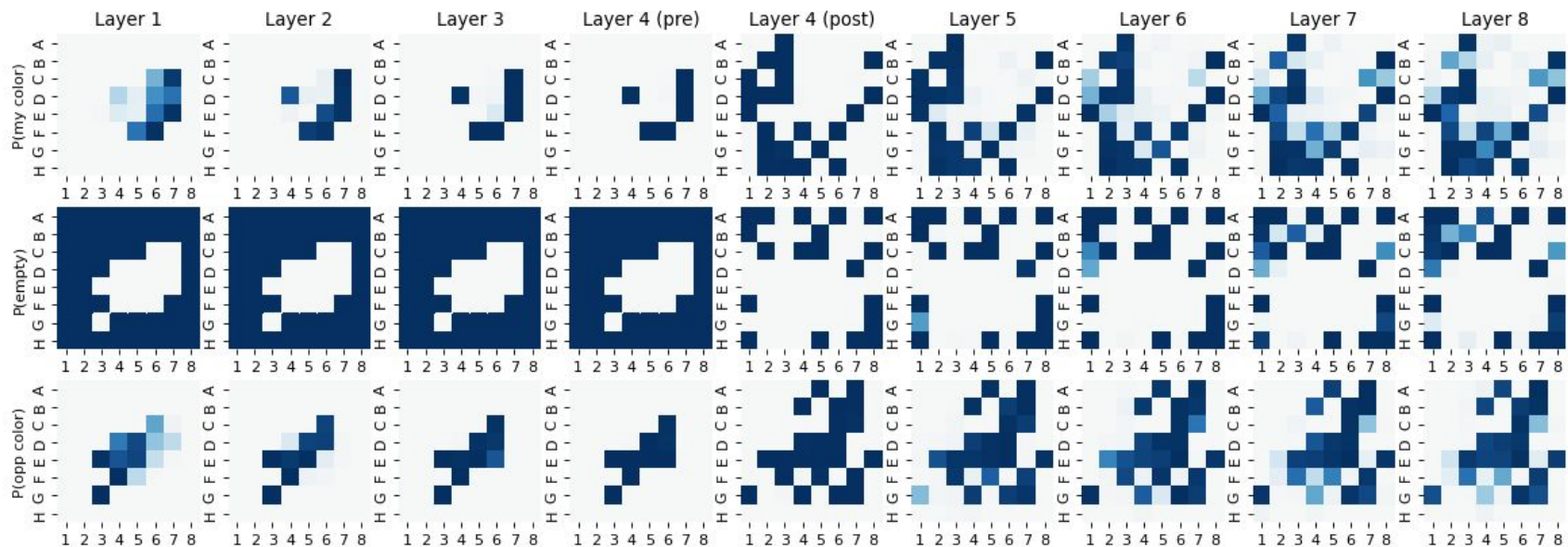
Novel technique for generating activations directly from a desired probe state, which can be inserted into a model during inference
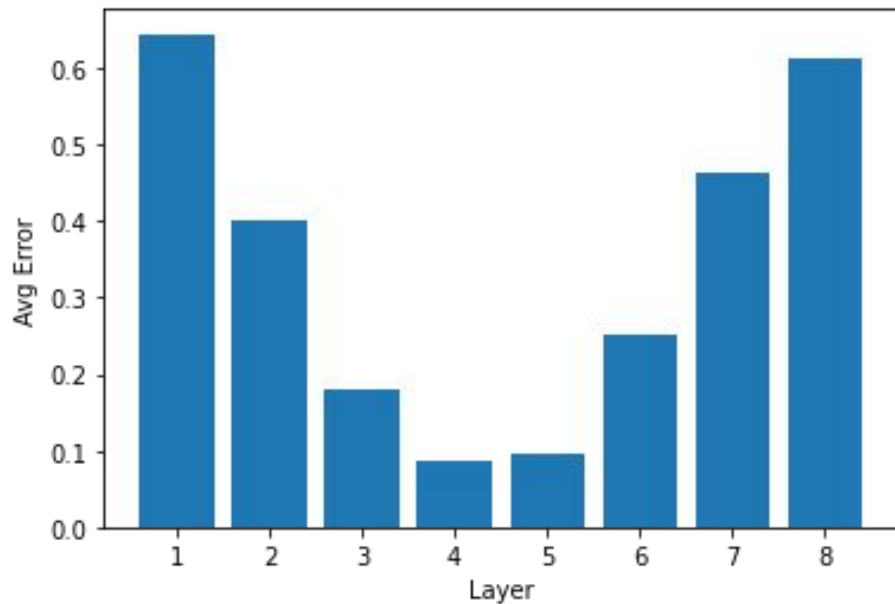
# GLOBAL EDITS



A global intervention: completely overwriting the residual stream after layer 4

# Global edits



Board representation visualized across layers during global intervention
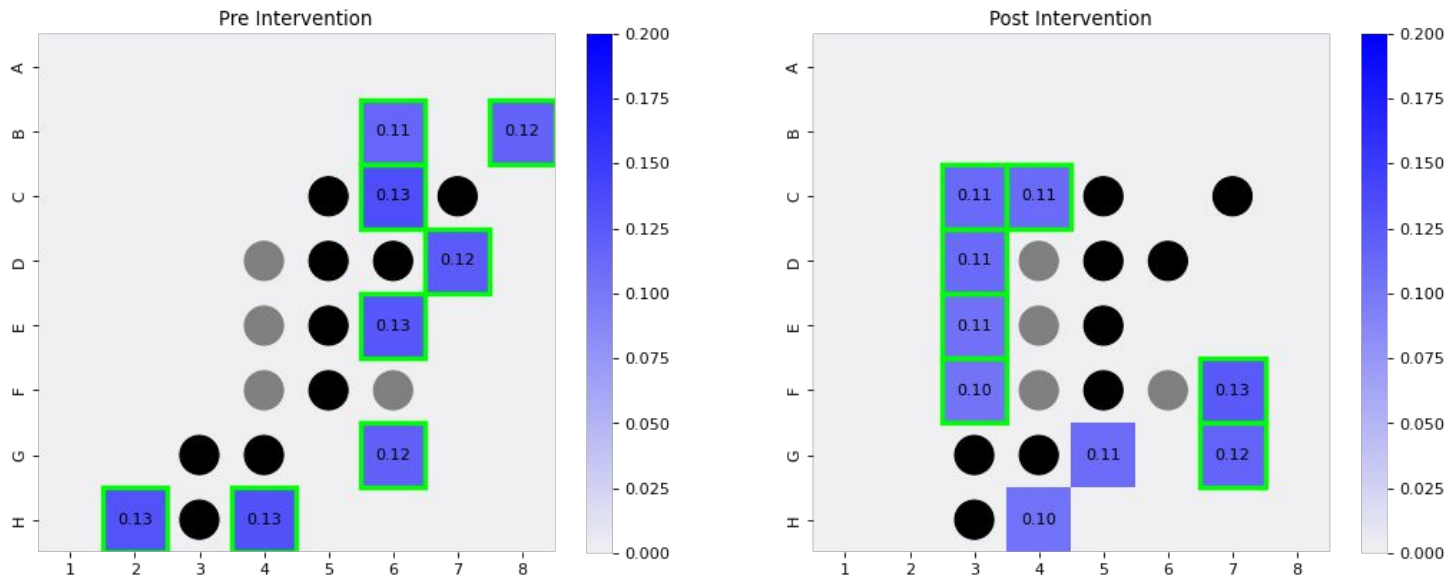
# Evidence of Board locality



Aggregated results from systemically attempting thousands of global interventions

Error is minimized in layers 4 and 5 where the board has the cleanest representation within the model
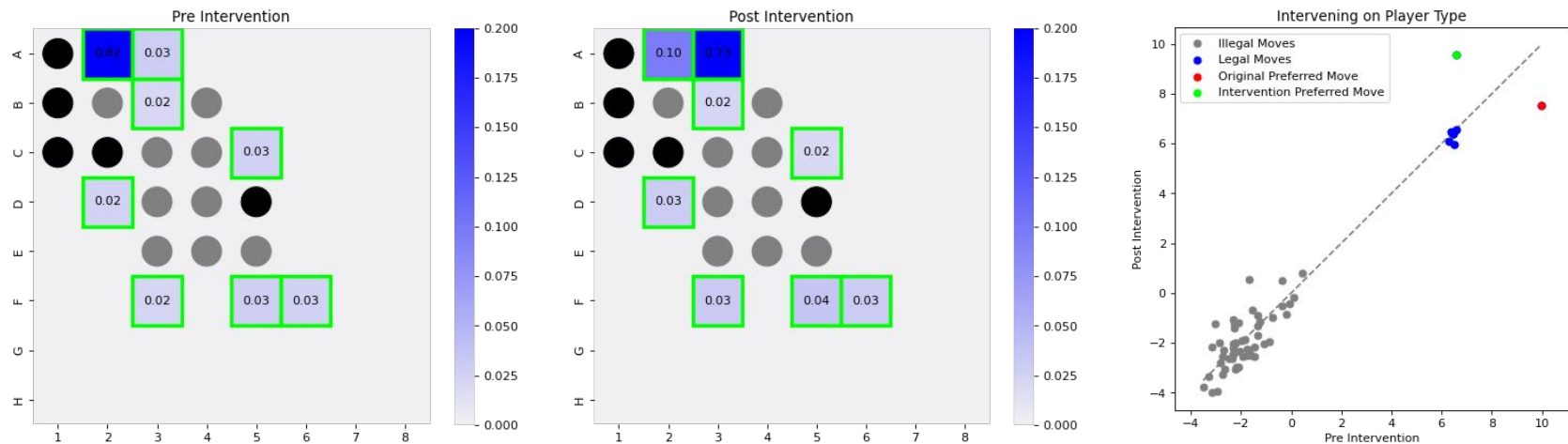
# PROBING OTHER REPRESENTATIONS

**Next color to play:** performed successful interventions to shift model between white/black on the same sequence



Intervention on next color from white to black

# PROBING OTHER REPRESENTATIONS

**Player bias:** created new dataset with artificial feature—each game had a bias to move towards one of the corners. Probes were highly accurate but interventions had mixed success



Intervention on player bias from top left to top right, layers 4-7

# High-level takeaways

Reinforces linear representation hypothesis

Demonstrates impressive resilience of world models

Probes are an effective tool to extract representations, make causal changes to model without needing more fine-grained analysis, but have fundamental weaknesses in identifying fully disentangled representations

Future work should aim to identify if transformers generally tend towards highly interpretable representations, or just a function of data richness (or some domain specific phenomenon)?

# ACKNOWLEDGEMENTS

Huge thanks to **David Reber**, my research mentor for this project, for invaluable advice and guidance throughout the research process

**Victor Veitch**, for detailed feedback on ideas and access to the DSI cluster so I could train models

**Zack**, for overall support and running this program

QUESTIONS?