

Emergent world representations:

Exploring a sequence model trained on a synthetic task.

Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg.



 Harvard John A. Paulson
School of Engineering
and Applied Sciences



 Northeastern University
Khoury College
of Computer
Sciences

Are models learning something meaningful, causal, about the world...

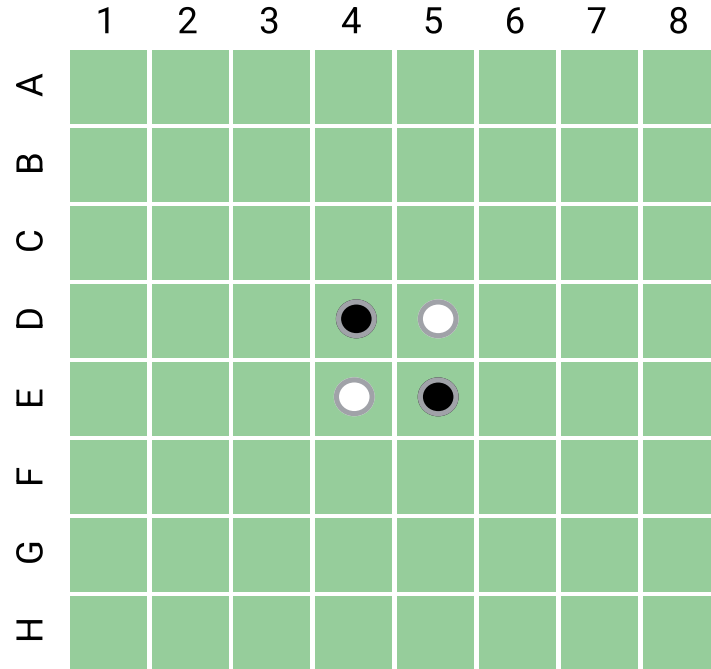
**Are models learning something meaningful, causal, about the world...
Or are they just memorizing data (then regurgitating it)?**



Looking for meaningful representations in large models is like looking for a needle in a haystack

Let's look at tiny pile of hay

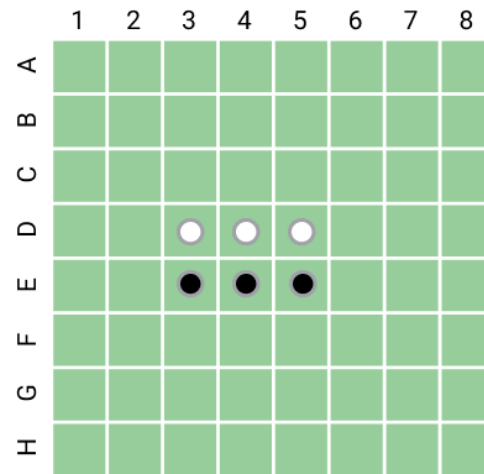
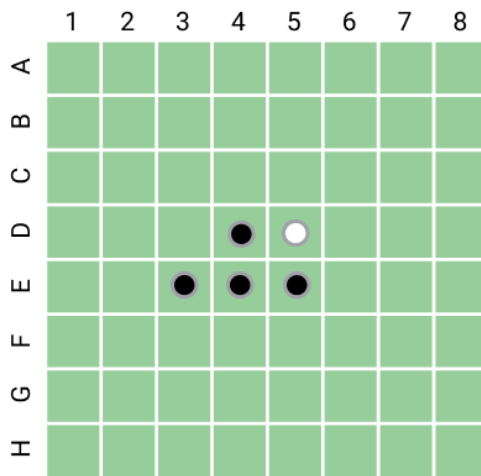
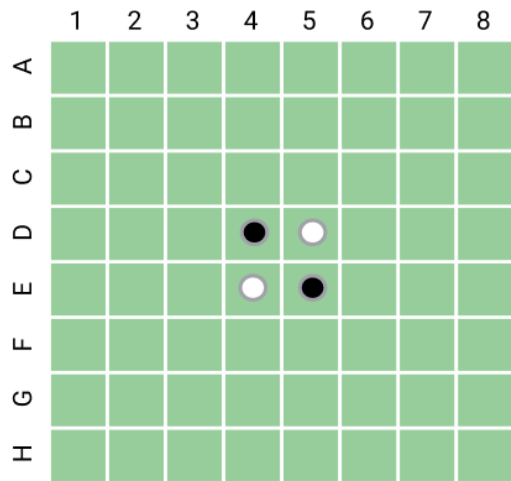
Game of Othello: A Toy Model



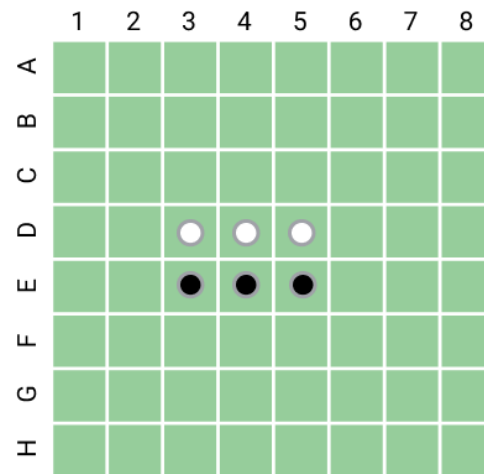
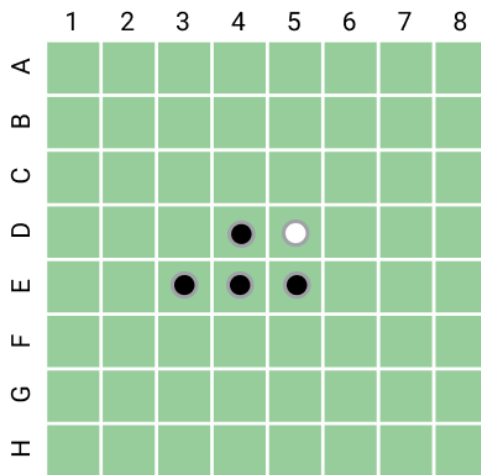
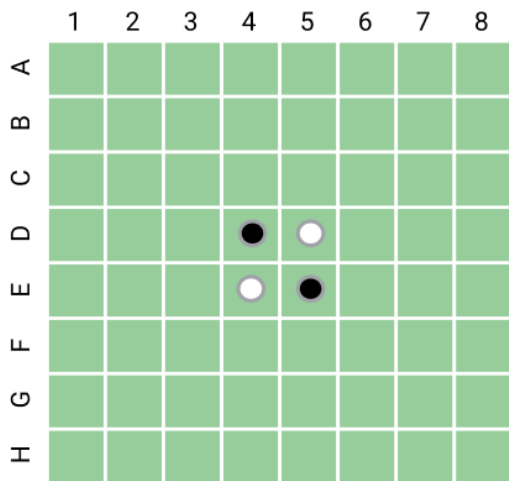
Othello

	1	2	3	4	5	6	7	8
A								
B								
C								
D				●	○			
E			●	●	●			
F								
G								
H								

Game of Othello: A Toy Model



Game of Othello: A Toy Model

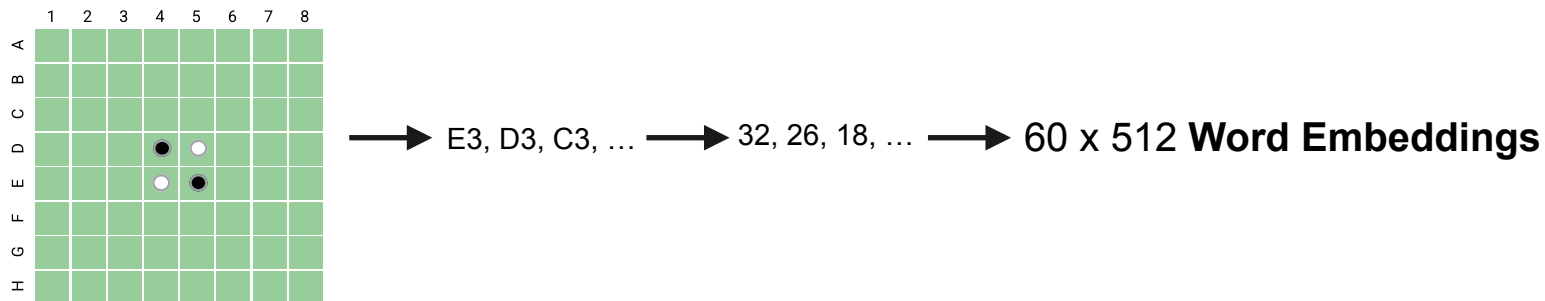


Toy model captures aspects of the general case:

- Legality is nonlinear function of board state and
- The board state is a nonlinear function of the moves

Trained “GPT-Othello” to predict tokens in transcripts of Othello games

Example task: C4 C3 D3 C5 D6 F4 B4 C6 B5 B3 B6 E3 C2 __?



No prior knowledge of game, rules, board.
Just sequences of tokens.

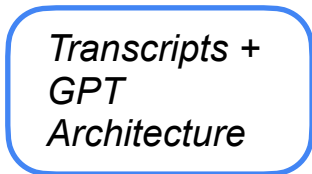
Game Play Transformer



Finetuning



OR

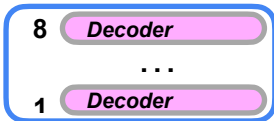
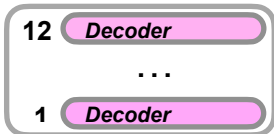


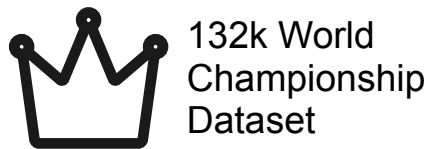
Training



GPT-2

GPT-Othello

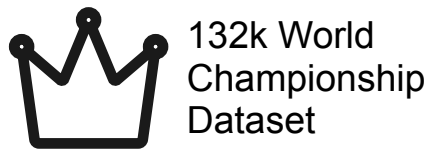




GPT-Othello
Championship



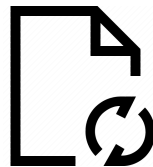
GPT-Othello
Synthetic



132k World
Championship
Dataset



GPT-Othello
Championship



20m Synthetic
Dataset



GPT-Othello
Synthetic

Can GPT-Othello predict a (legal) next move?

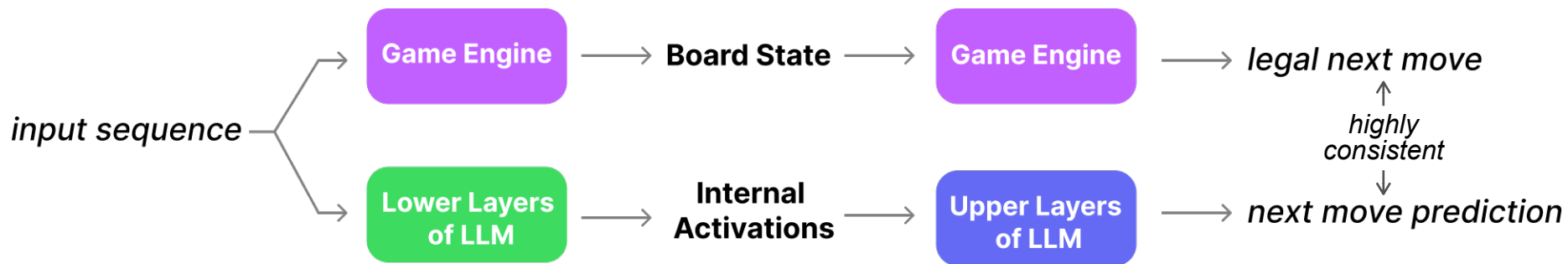
Percentage of illegal top predictions (error) in validation

Synthetic 0.01%

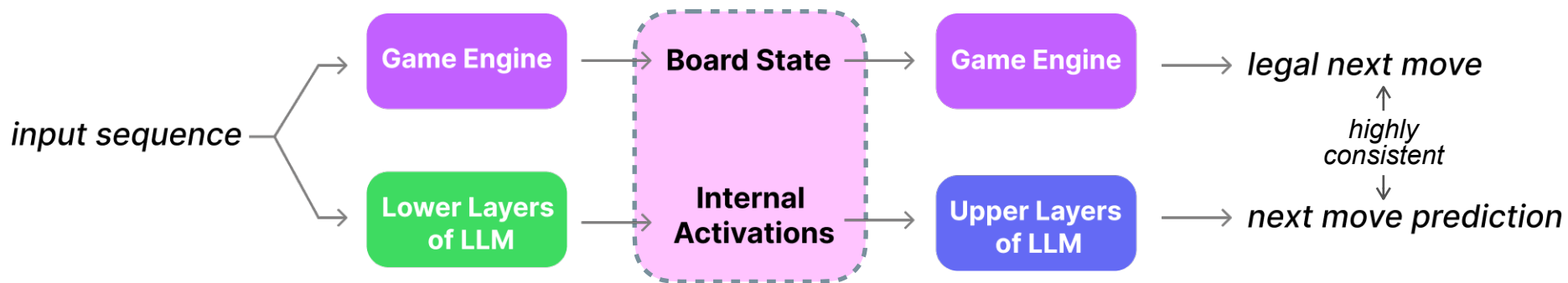
Championship 5.17%

Randomly Initialized 93.29%

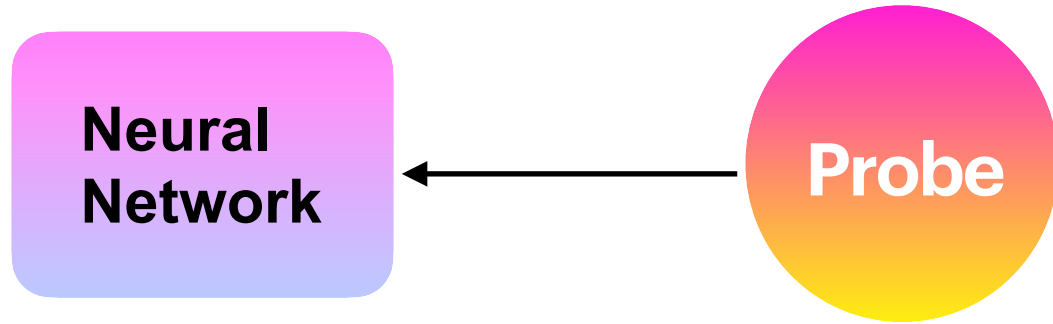
Comparison of Game Engine to GPT



Comparison of Game Engine to GPT

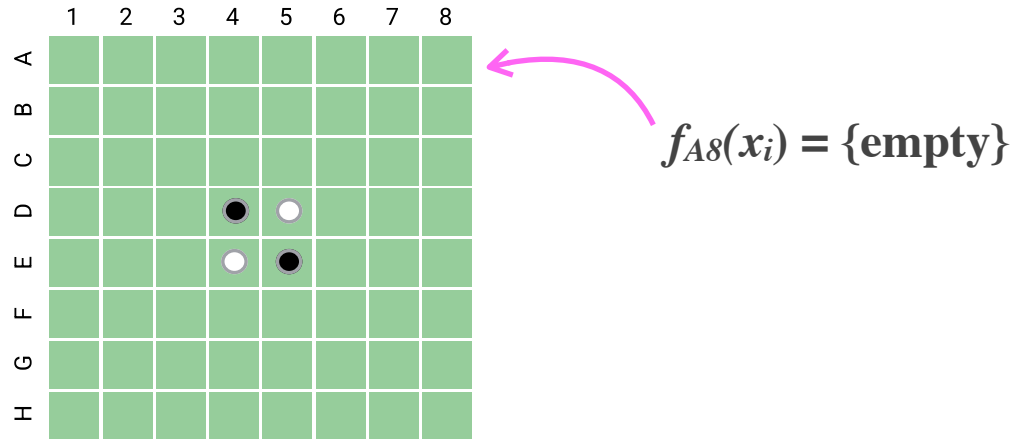


**In language models, we do this by training probes—
classifiers—on the language model’s layer activations.**



Defining Board State Probe for Othello GPT

For each square \mathbf{s} on the board, can we train a simple classifier f_s , such that $f_s(x_i) = \{\mathbf{white}, \mathbf{black}, \mathbf{empty}\}$ (where x_i represents the value of concept \mathbf{C} in the input) reflecting whether \mathbf{s} is white, black, or empty?



Probe Experiment

Modeling: 3-way classification (black/blank/white) for each square (64 in total)

Features: Layer activations between Transformer blocks

Linear Probe:

$$\text{softmax}(Wx)$$

$$W \in \mathbb{R}^{F \times 3}$$

Two-layer Probe:

$$\text{softmax}(W_1 \text{ReLU}(W_2 x))$$

$$W_1 \in \mathbb{R}^{H \times 3}$$

$$W_2 \in \mathbb{R}^{F \times H}$$

Probe Performance

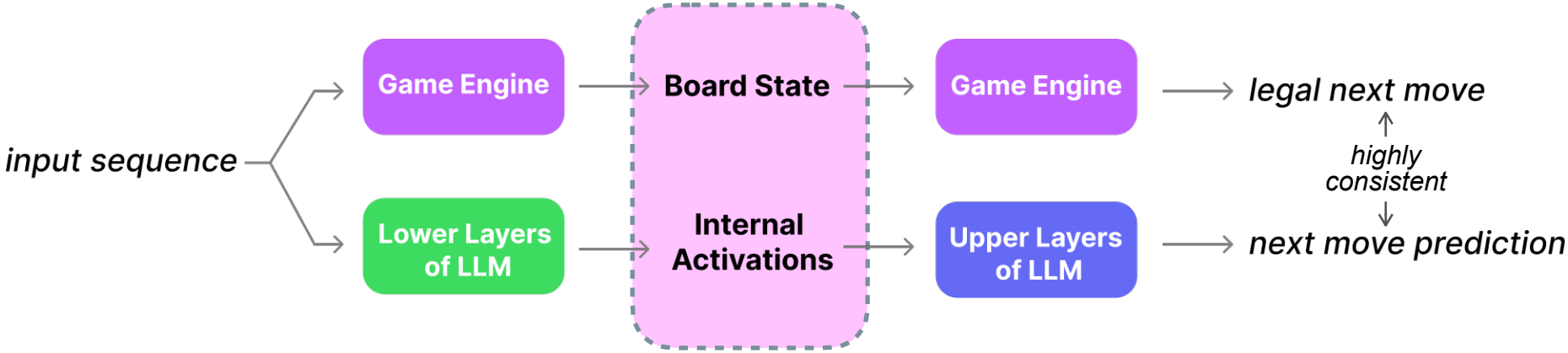
Error rates of linear probe across different layers

	x^1	x^2	x^3	x^4	x^5	x^6	x^7	x^8
Randomized	26.7	27.1	27.6	28.0	28.3	28.5	28.7	28.9
Championship	24.2	23.8	23.7	23.6	23.6	23.7	23.8	24.3
Synthetic	21.9	20.5	20.4	20.6	21.1	21.6	22.2	23.1

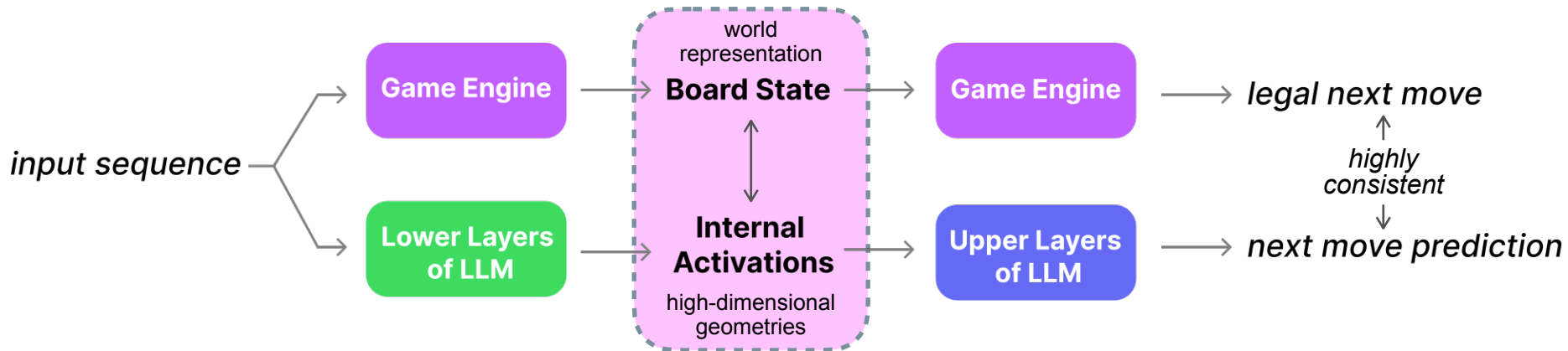
Error rates of nonlinear probe across different layers

	x^1	x^2	x^3	x^4	x^5	x^6	x^7	x^8
Randomized	25.5	25.4	25.5	25.8	26.0	26.2	26.2	26.4
Championship	12.8	10.3	9.5	9.4	9.8	10.5	11.4	12.4
Synthetic	11.3	7.5	4.8	3.4	2.4	1.8	1.7	4.6

Probe provides evidence of board state model

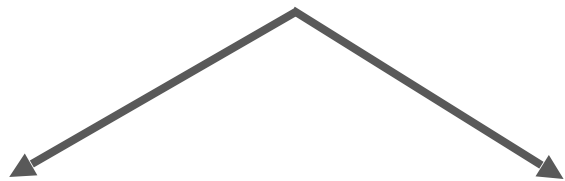


Probe provides evidence of board state model

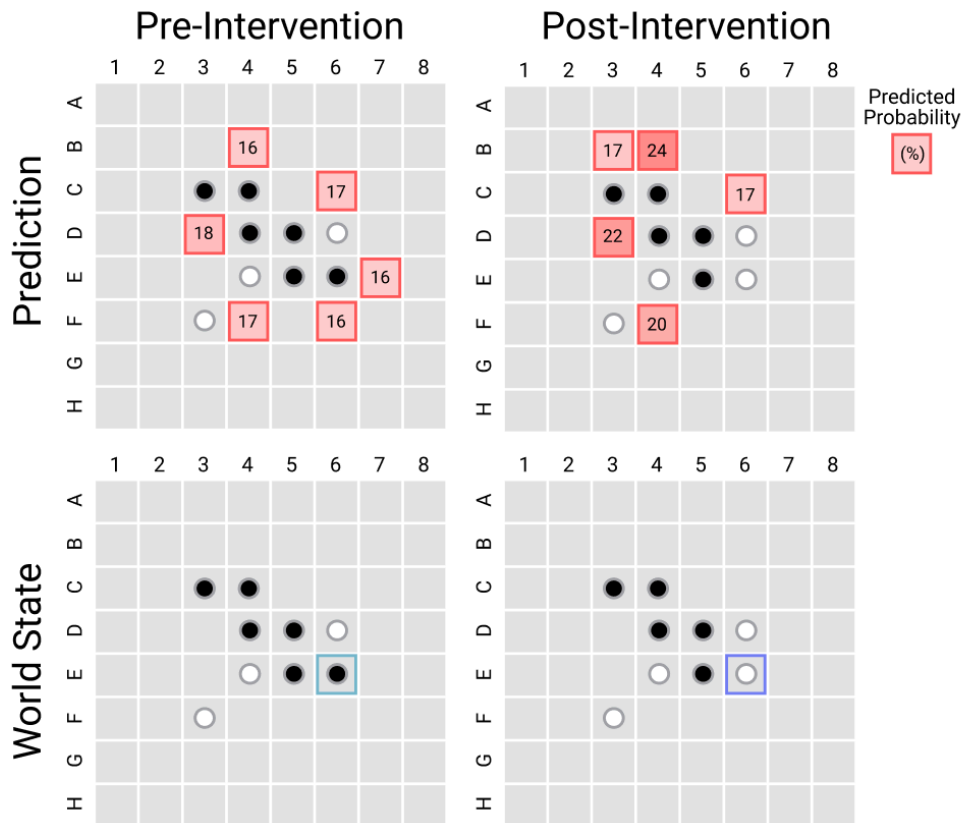


Controllable World Model

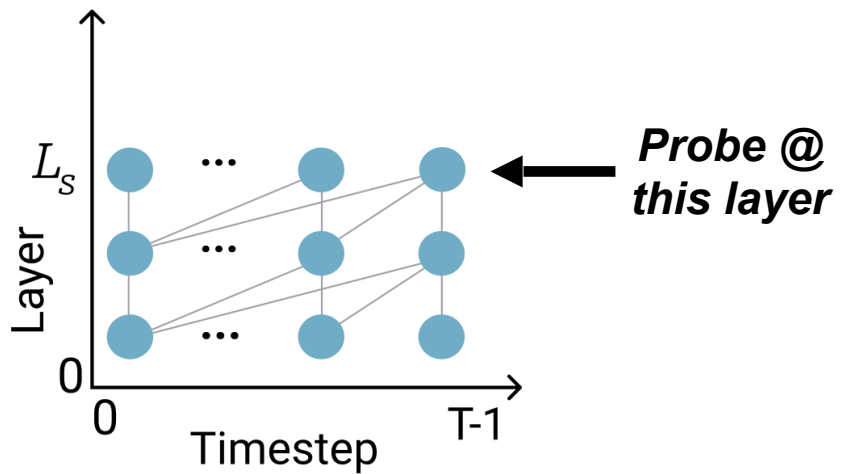
**Internal
Representation**



World Models - - - - - **Prediction**

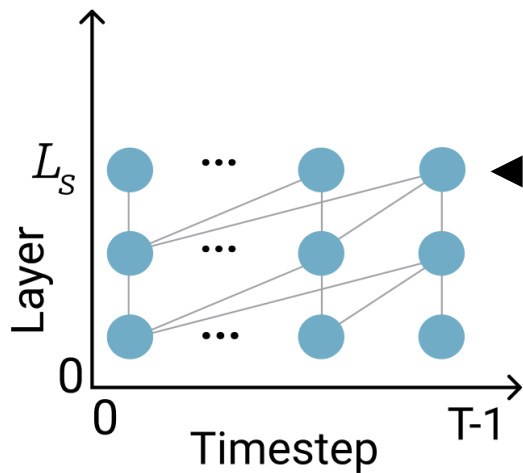


Intervention

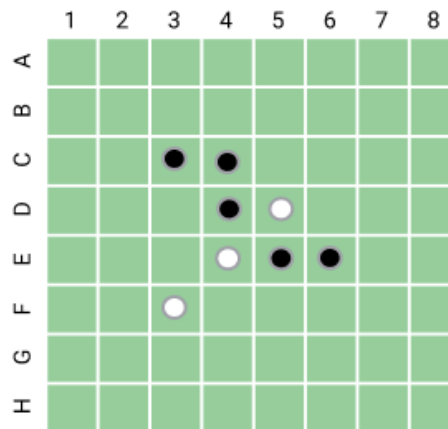


● : Factual feature

Intervention

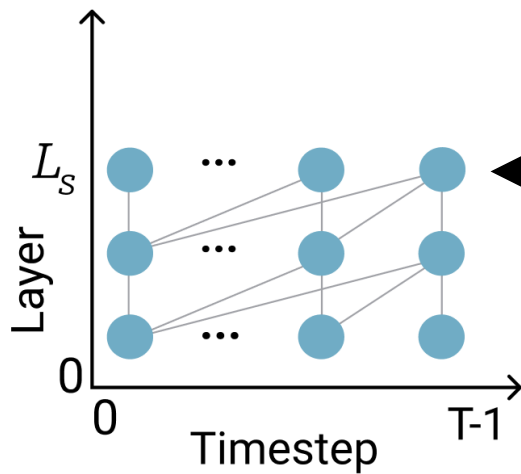


*Probe @
this layer*

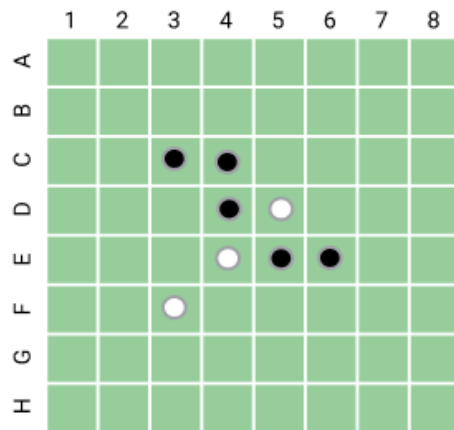


● : Factual feature

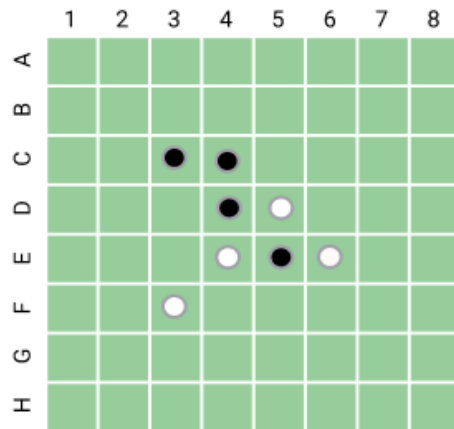
Intervention



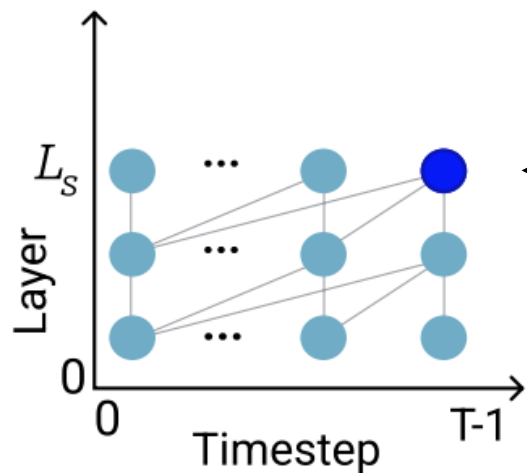
● : Factual feature



Intervene with gradient descent

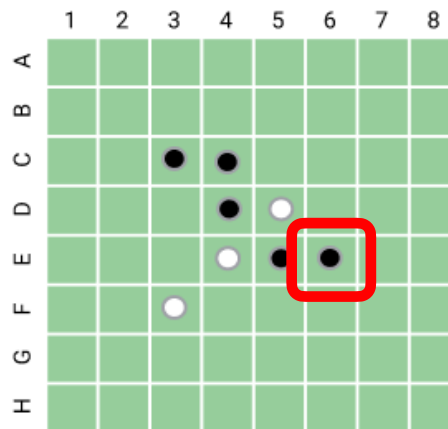


Intervention

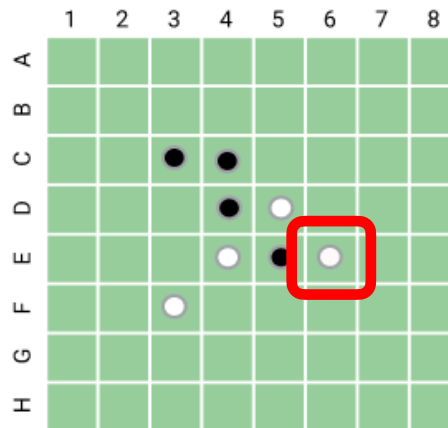


*Copy
back*

- : Factual feature
- : Intervened feature

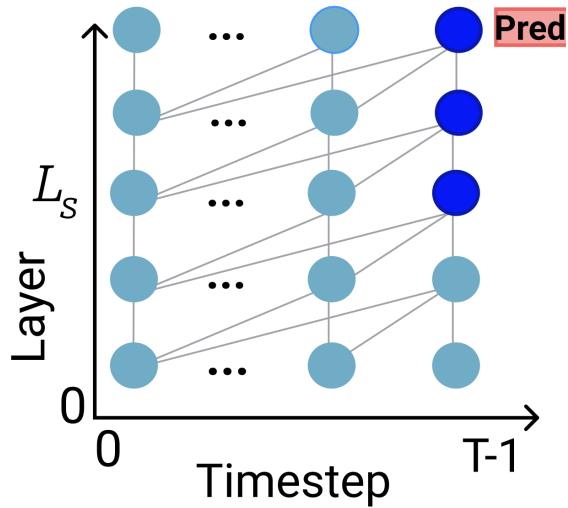


*Intervene with
gradient descent*



pause, hack, put-back (continue)

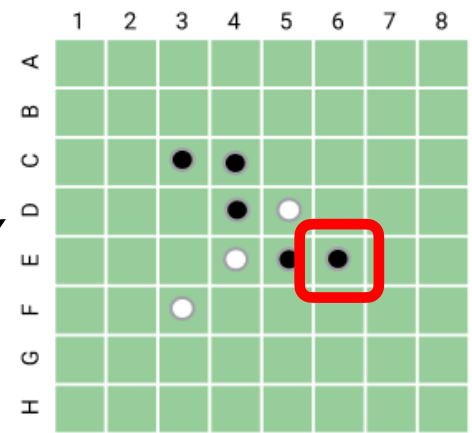
Intervention



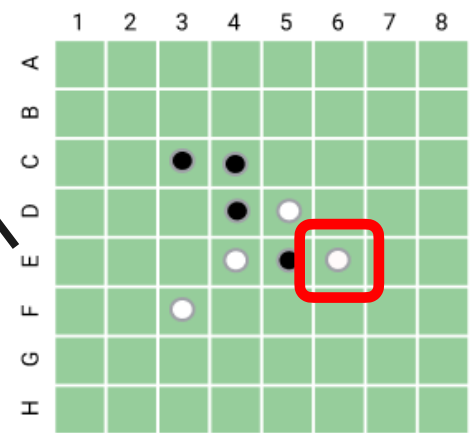
- : Factual feature
- : Intervened feature
- : Self-corrected

Probe @ this layer

Copy back

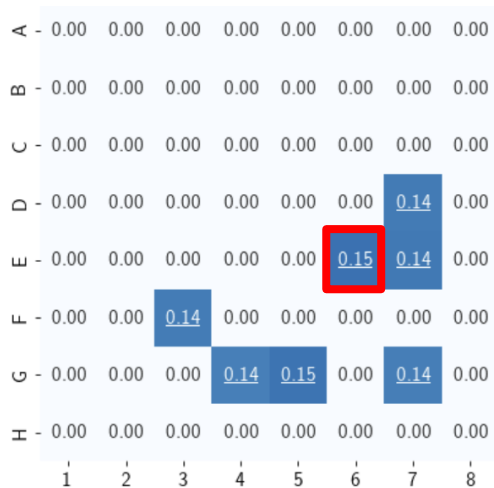


Intervene with gradient descent

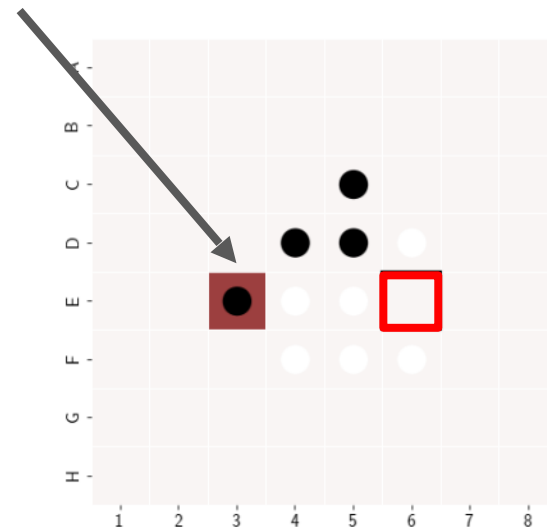
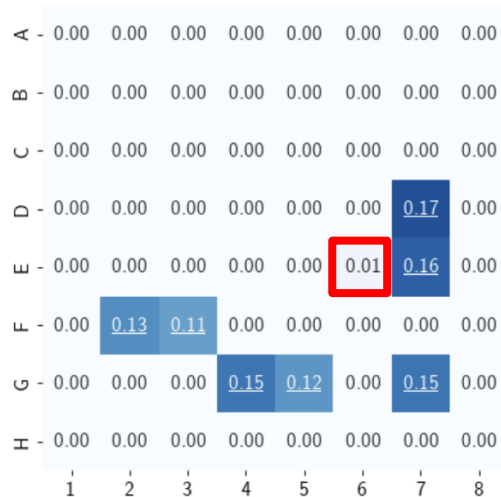


Attribution via intervention

Logit map *before*
flipping E3 world representation



Logit map *after*
flipping E3 world representation




0.15

-

0.01

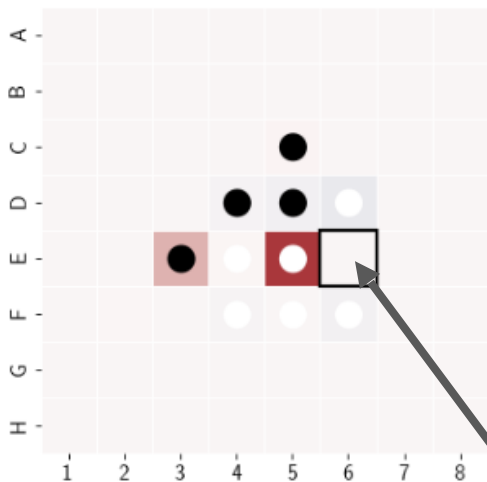
=

0.14

 : E6, the square we attribute

Latent saliency maps

By holding the world model of the colored tile as it is,
Othello-GPT thinks it is...



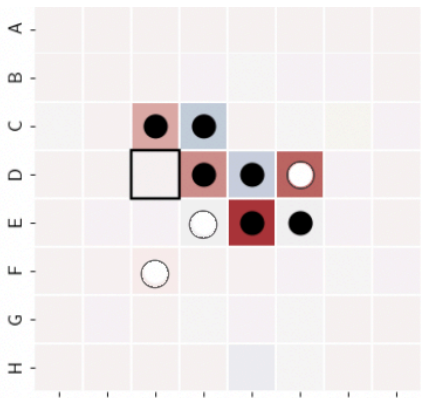
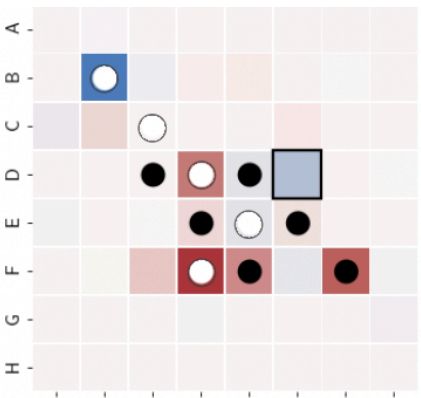
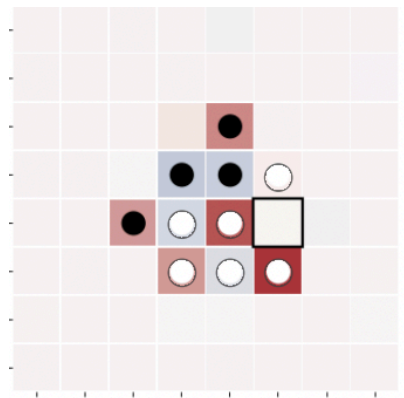
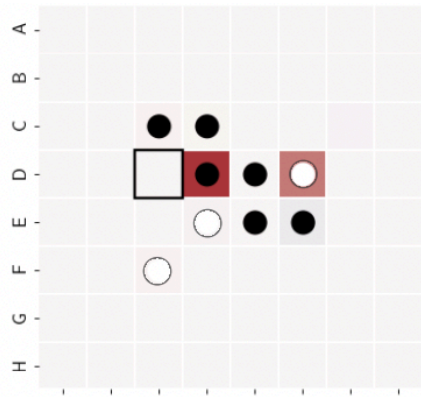
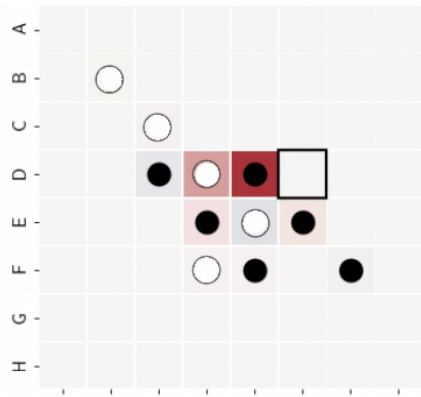
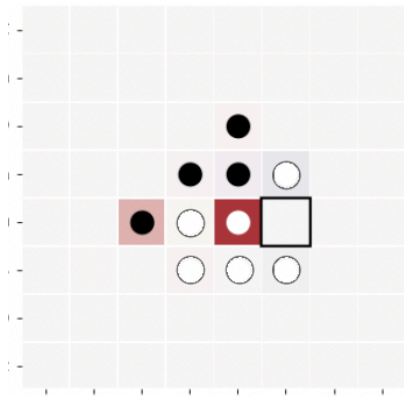
Negatively

Positively

...contributing to the prediction of the
next step.

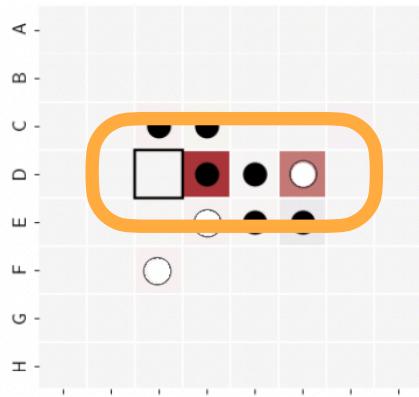
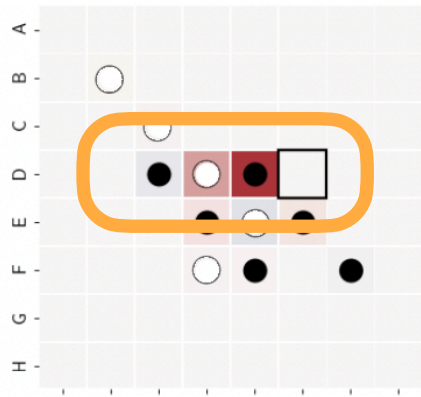
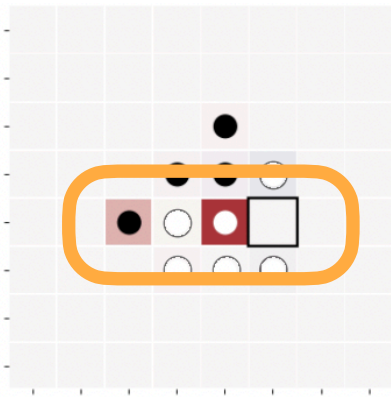
Enclosed E6: a legal next step we attribute

Which row is strategic? Which is just “legal”

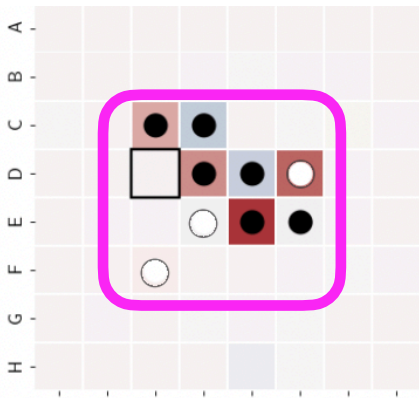
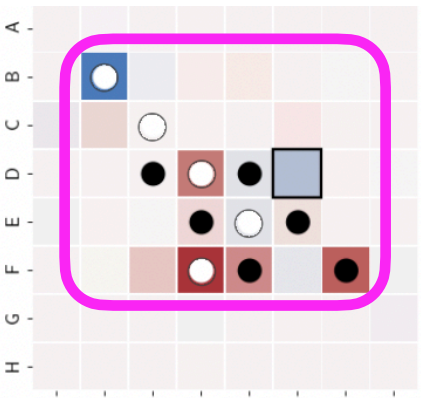
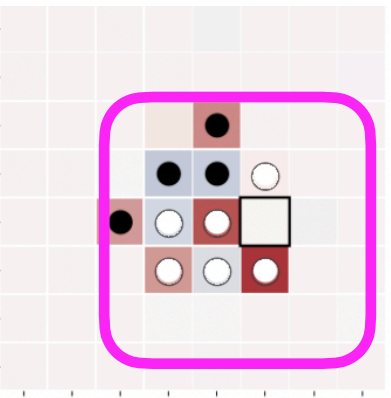


Which row is strategic? Which is just “legal”

Legal

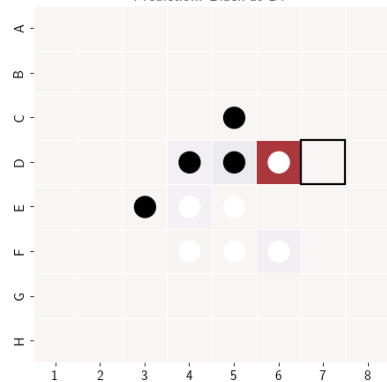


Strategic

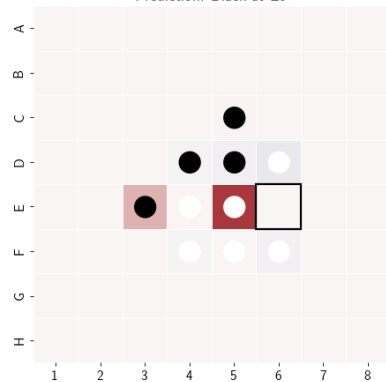


Attributing other legal moves

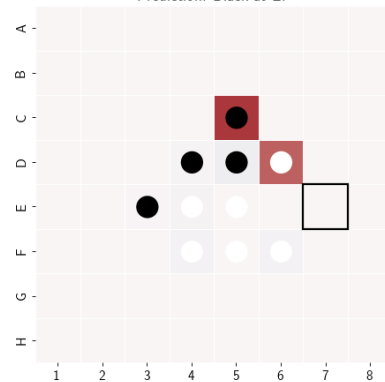
Prediction: Black at D7



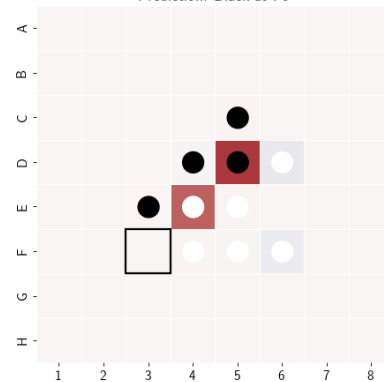
Prediction: Black at E6



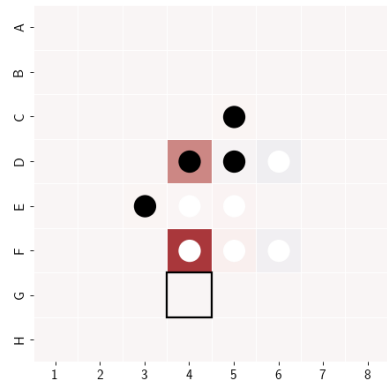
Prediction: Black at E7



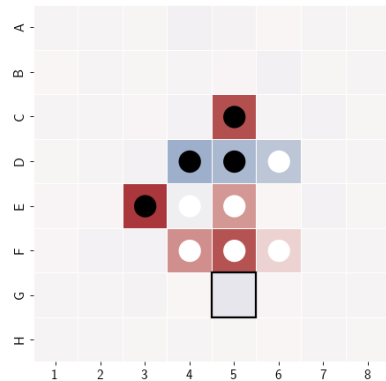
Prediction: Black at F3



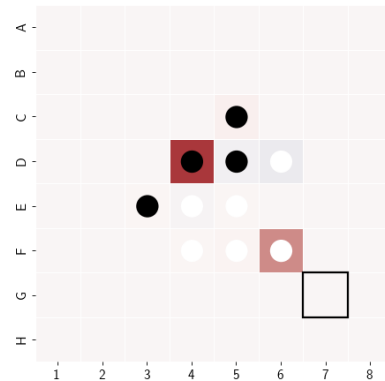
Prediction: Black at G4



Prediction: Black at G5

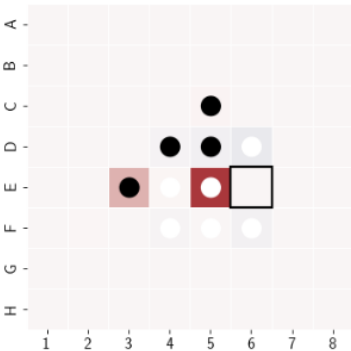


Prediction: Black at G7

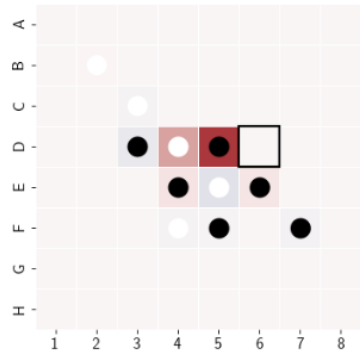


More cases on synthetic Othello-GPT

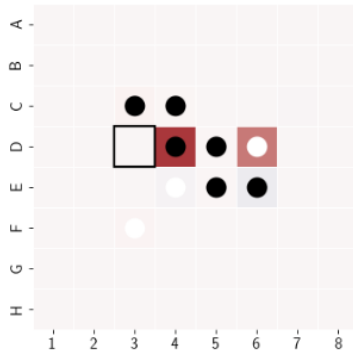
Prediction: Black at E6



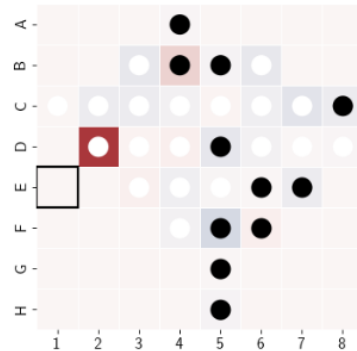
Prediction: White at D6



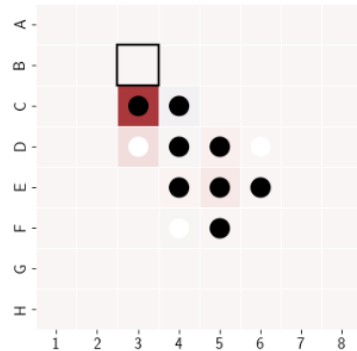
Prediction: White at D3



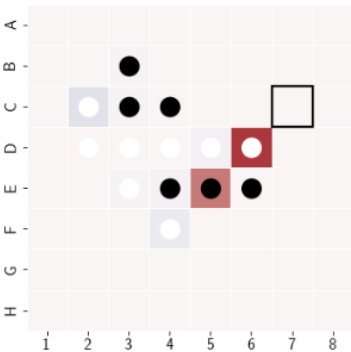
Prediction: Black at E1



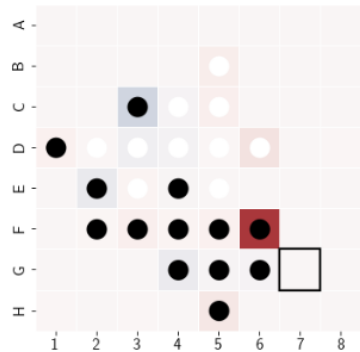
Prediction: White at B3



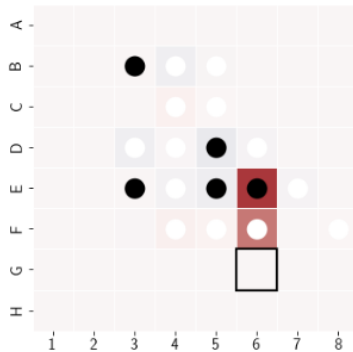
Prediction: Black at C7



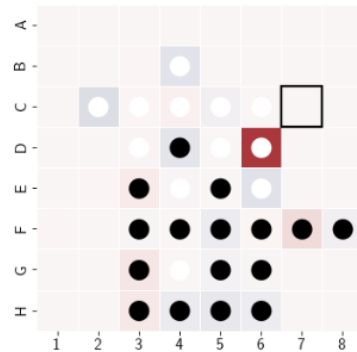
Prediction: White at G7



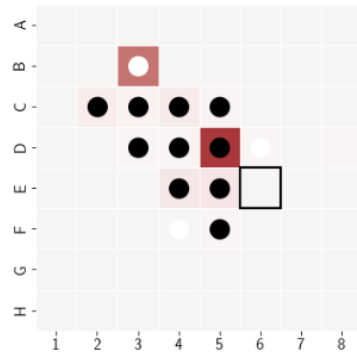
Prediction: Black at G6



Prediction: Black at C7

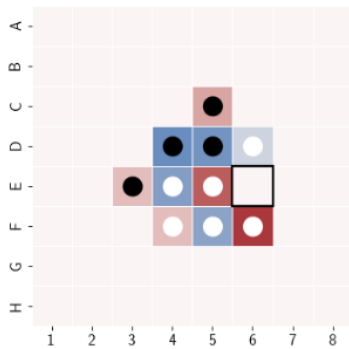


Prediction: White at E6

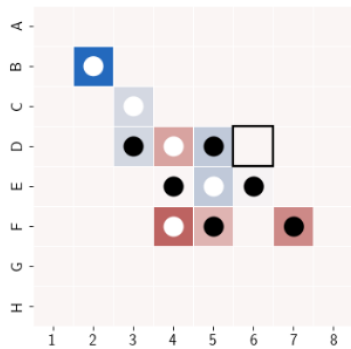


More cases on championship Othello-GPT

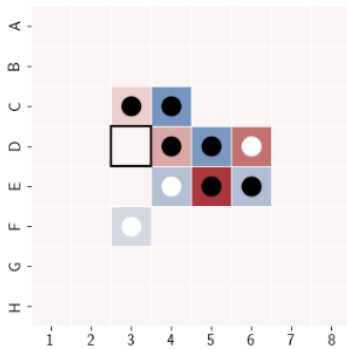
Prediction: Black at E6



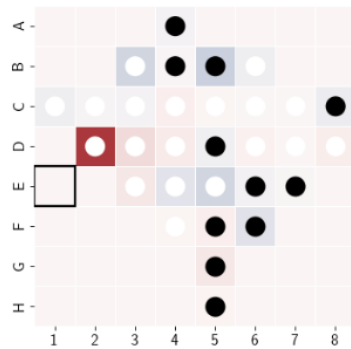
Prediction: White at D6



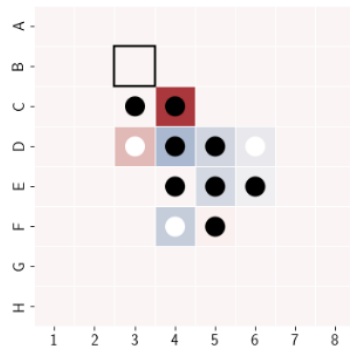
Prediction: White at D3



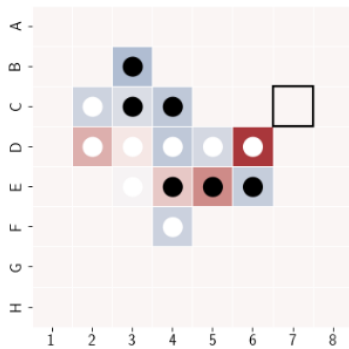
Prediction: Black at E1



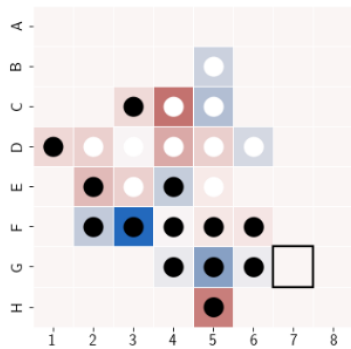
Prediction: White at B3



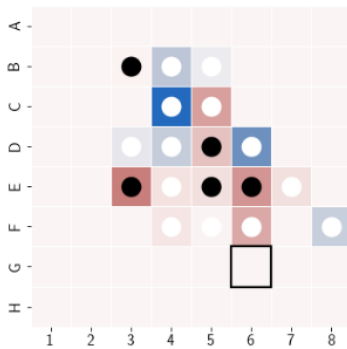
Prediction: Black at C7



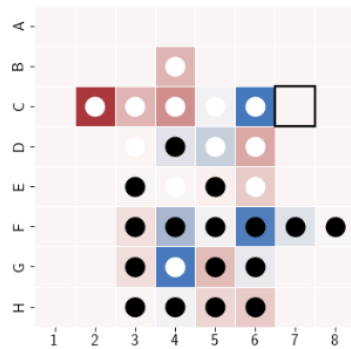
Prediction: White at G7



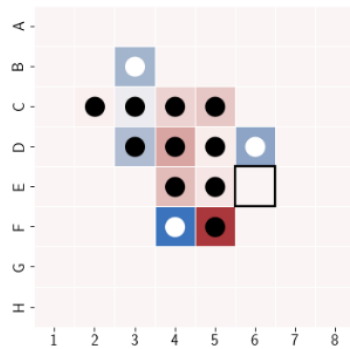
Prediction: Black at G6



Prediction: Black at C7



Prediction: White at E6



1 *Development and evaluation of GPT model, **GPT-Othello***



2 *Comparison of **linear** and **non-linear** probing*



3 *Novel **intervention** technique*



4 *Novel **latent saliency maps** built using intervention*



< INTERPRETING OTHELLO-GPT >

Actually, Othello-GPT Has A Linear Emergent World Representation

by **Neel Nanda** 23 min read 29th Mar 2023 17 comments

Interpretability (ML & AI)

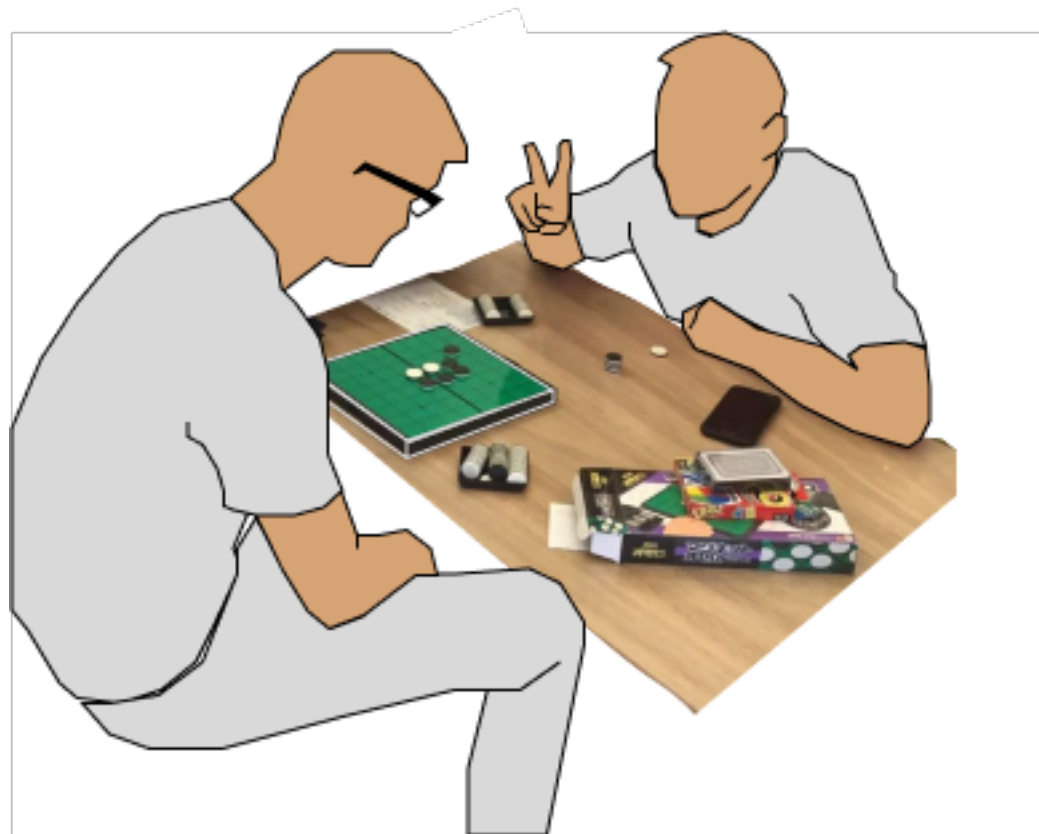
AI

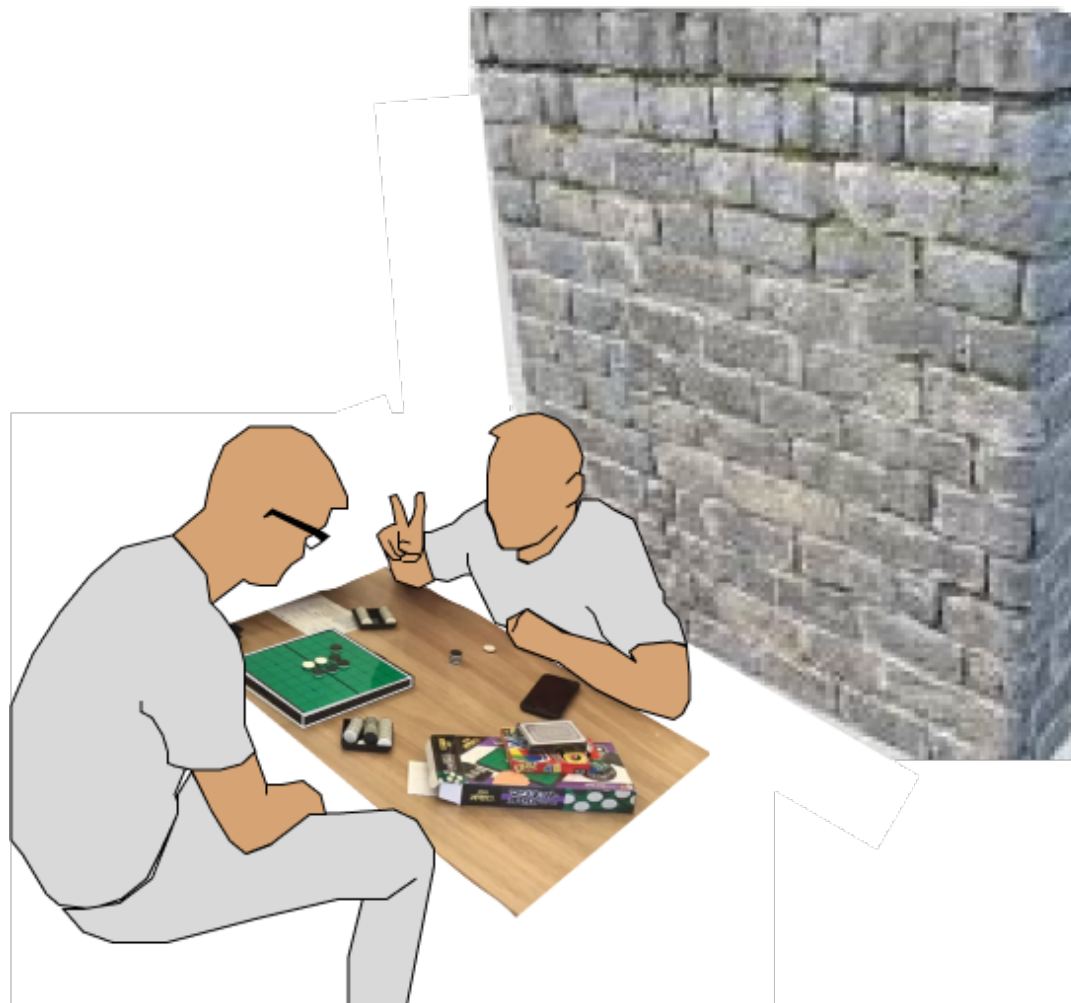
Frontpage

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

This is a linkpost for <https://neelnanda.io/mechanistic-interpretability/othello>

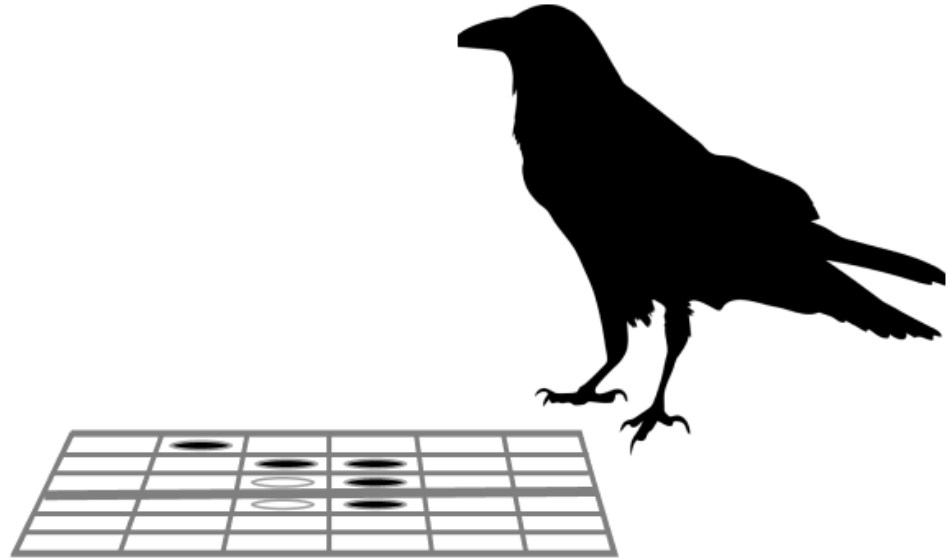




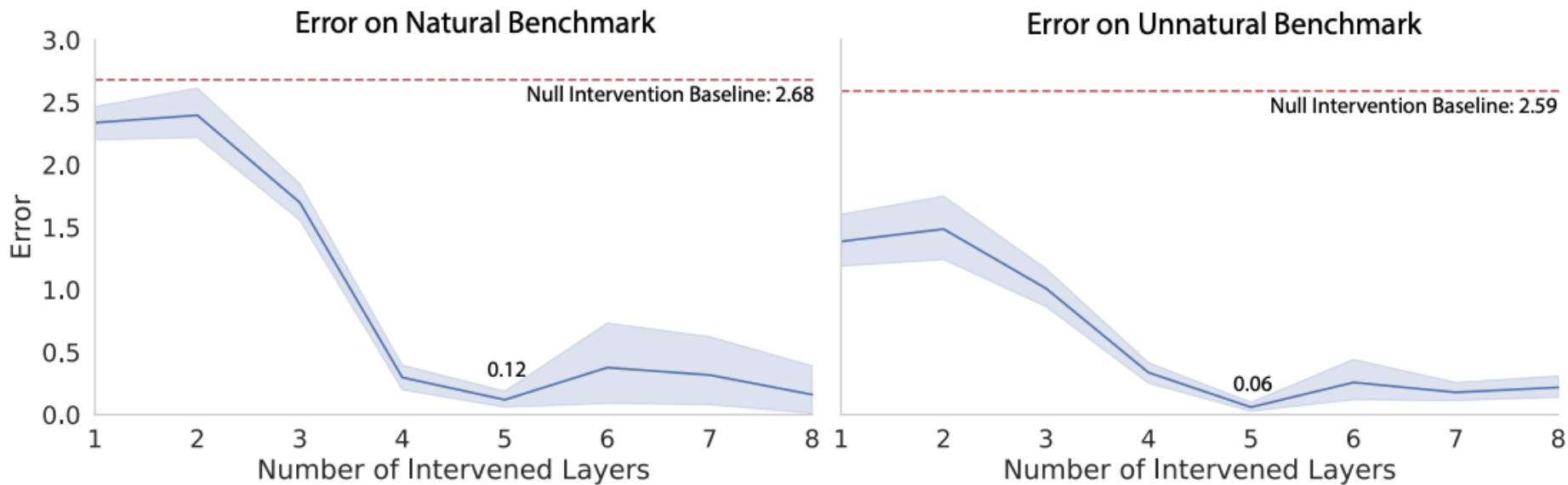


DELETE BOARD





Comparing top-N prediction to the N ground truths



Superficial Statistics vs ***Meaningful Representations***
(World Models)

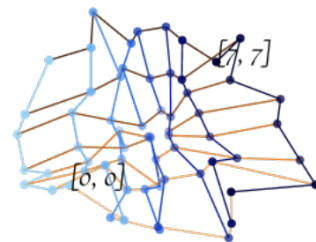
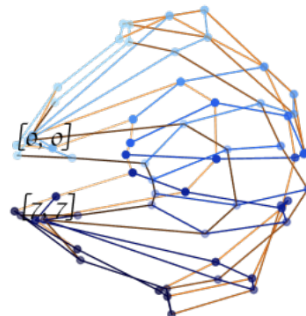
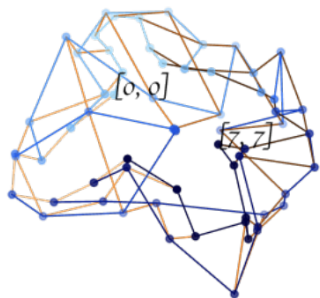
Geometry of probe weights

Randomized

Championship

Synthetic

Linear



Nonlinear

