# Othello Mech Int playground [1]
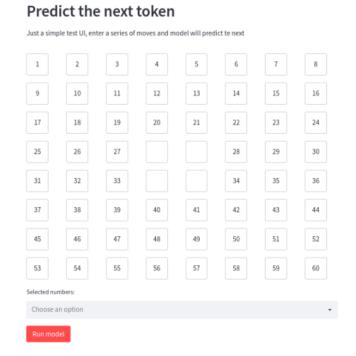
**Victor Levoso**        **Kunvar Thaman**

**Edoardo Pona**        **Abhay Sheshadri**

Made an interactive app to run basic interpretability experiments on othello.

*Keywords: Mechanistic Interpretability, Othello GPT,streamlit app.*

This is a modification of the "Trafo Mech Int playground" project (by Stefan Heimersheim and Jonathan Ng) to work on Othello-GPT instead of LLM.
Maybe available in streamlit but might crash at some point due to memory limitations.
Also available in a github repository to run locally.

## Screenshots



---

# Attention Pattern Visualization

Enter a prompt, show attention patterns

Prompt:

`20 ×` `5 ×` `8 ×` `6 ×` `9 ×`    ⚙ ▾

**Run model**

Attention patterns Layer 0:

**Attention Patterns**   **Head selector** (hover to focus, click to lock)

| | |
|---|---|
|  | Head 0   Head 1   Head 2   Head 3 <br> Head 4   Head 5   Head 6   Head 7 |

**Tokens** (click to focus)   `Source ← Destination ▾`

`20 5 8 6 9`

# Residual stream patching

Enter a clean prompt, correct answer, corrupt prompt and corrupt answer, the model will compute the patching effect

Clean Prompt:

`20 ×` `21 ×` `14 ×` `11 ×` `34 ×`    ⚙ ▾

Correct Answer:

`16 ×`    ⚙ ▾

Corrupt Prompt:

`20 ×` `21 ×` `14 ×` `19 ×` `34 ×`    ⚙ ▾

Corrupt Answer:

`43 ×`    ⚙ ▾

**Run model**

**Patching residual stream at specific layer and position**