# Easy Trafo* Mech Int[1]

*Transformer

**Stefan Heimersheim, Jonathan Ng**
SERIMATS scholars

## Abstract

We created a simple web app allowing users to create some standard mechanistic interpretability plots (based on Stefan's explainer) for arbitrary prompts.

The web app computes residual stream patching, attention head output patching, and attention pattern visualizations. The currently-online version allows only 1 & 2 layer models, but in principle the code supports any models supported by TransformerLens.

The code is built on TransformerLens and CircuitsVis for the interpretability tools, with the web page built with Streamlit.

The web app code is available on GitHub. You can run a local version with `streamlit run Home.py`, which also allows you to select arbitrary models (comment them in at the top of the file).

*Keywords: Mechanistic interpretability, ML safety*

## 1. Introduction

The motivation for this project was my surprise how easy and relatively accessible it was to start researching Transformer Interpretability. A couple lines of code with TransformerLens gives you a model, you can look at all the weights, and arbitrarily mess with internal activations!
Then I realized it would be really cool to just have a web page where really anyone could try this out with minimal trivial inconveniences.

A secondary motivation was to finally establish the German abbreviation **Trafo** for **Transformer**.

## 2. Methods

---

[1] Research conducted at the Apart Research Alignment Jam #4 (Mechanistic Interpretability), 2023 (see https://itch.io/jam/mechint)

Really we just wrote a web app that wraps simple TransformerLens functions. The full code is on [GitHub](), the patching functions are based on [Stefan's tutorial text]() and [notebook]().

## 3. Results

Let's put a couple of screenshots here.

### Overview



*The website with the model selector and simple prompt testing tools.*

# Residual stream and Attention head output patching



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

Clean Prompt:

Her name was Alex Hart. Tomorrow at lunch time Alex
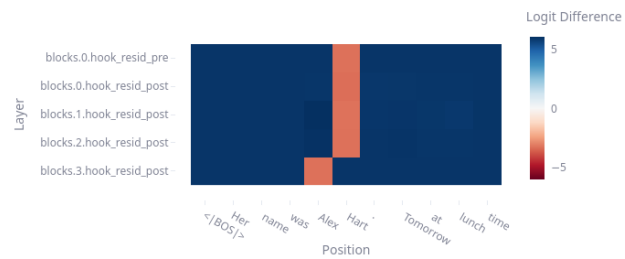
Correct Answer:

Hart

Corrupt Prompt:

Her name was Alex Carroll. Tomorrow at lunch time Alex
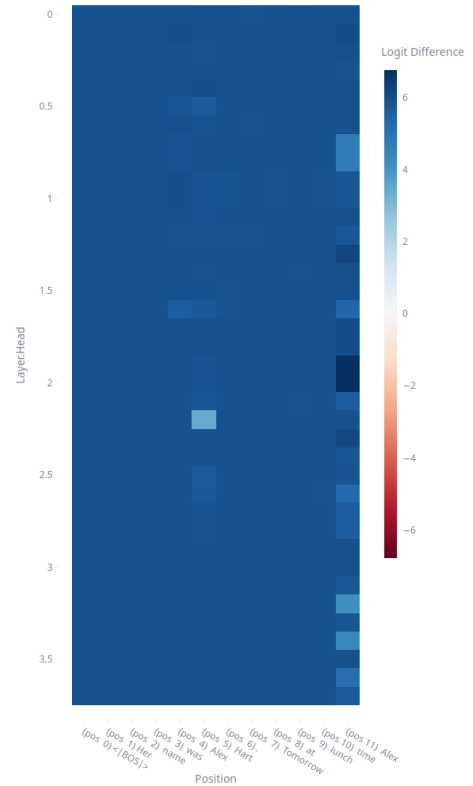
Corrupt Answer:

Carroll

Run model

*Residual stream patch with different last name, attention head result patch with different first name.*

**Attention Pattern Visualization**
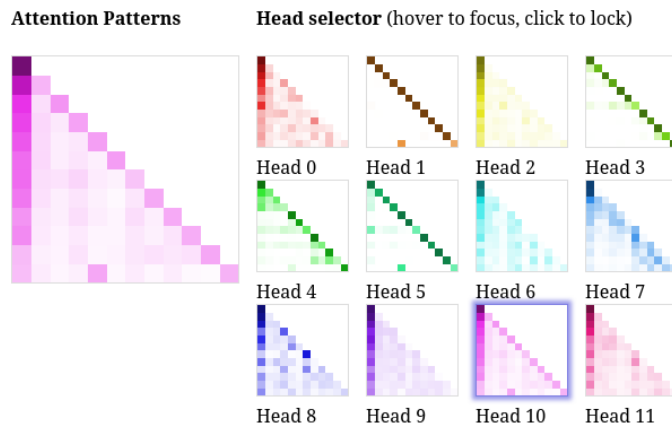
Powered by CircuitsVis

Enter a prompt, show attention patterns

Prompt:

Her name was Alex Hart. Tomorrow at lunch time Alex

Run model

Attention patterns Layer 0:



*Attention Pattern visualization of GPT2-small. GPT2 runs locally but not on Streamlit due to storage limits.*

## 4. Discussion and Conclusion

Discuss your results: These aren't research results, just a cool tool :)

We had to restrict the online app to 1 and 2 layer models, and it will also crash when too many users use it. The maximum number of users is approximately 1.

## 5. References

1. TransformerLens
2. CircuitsVis
3. AMFTC
4. Streamlit
5. Stefan's tutorial text and notebook