# DIGITAL FINANCE

Schweizerische Eidgenossenschaft
Confédération suisse
Confederazione Svizzera
Confederaziun svizra

State Secretariat for Education,
Research and Innovation SERI

Funded by
the European Union

MARIE CURIE ACTIONS

DIGITAL

# Reading

https://github.com/JShollaj/awesome-llm-interpretability?tab=readme-ov-file

# LLMs Interpretability

DIGITAL

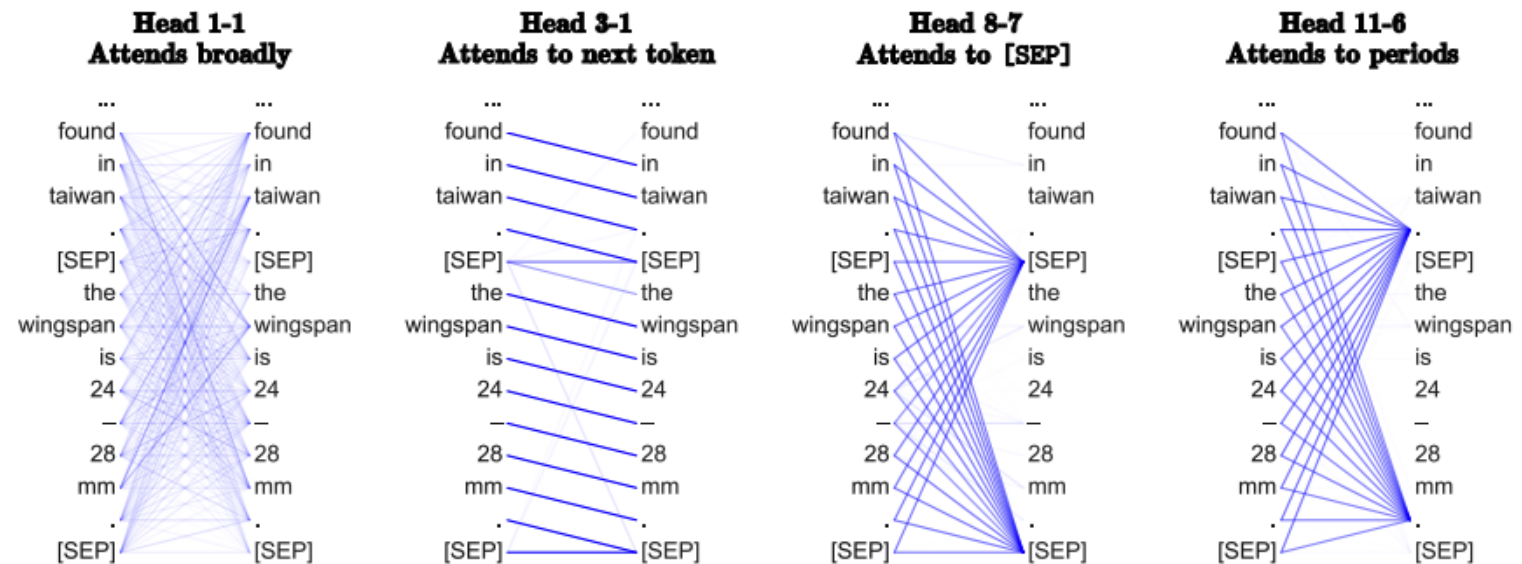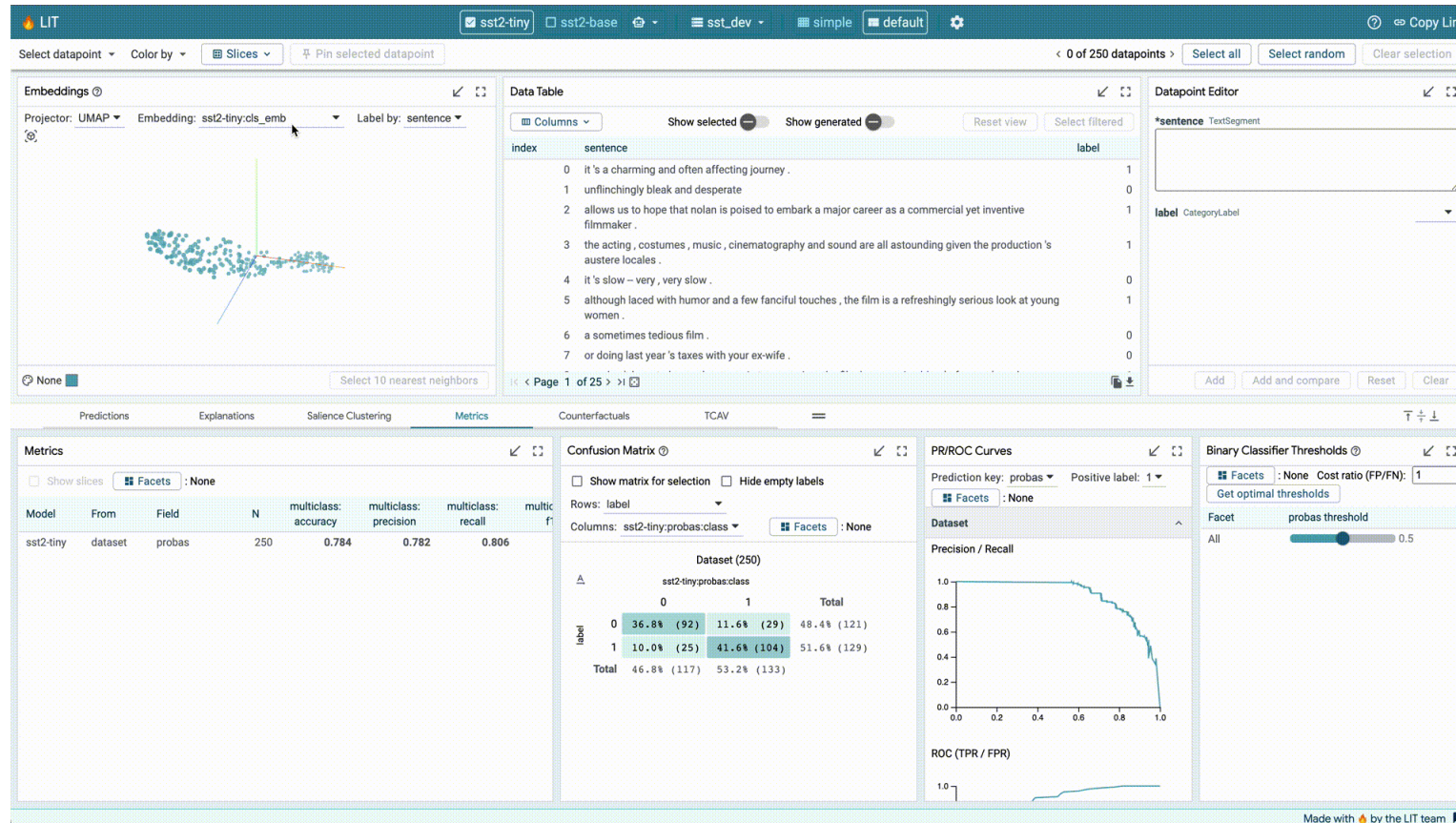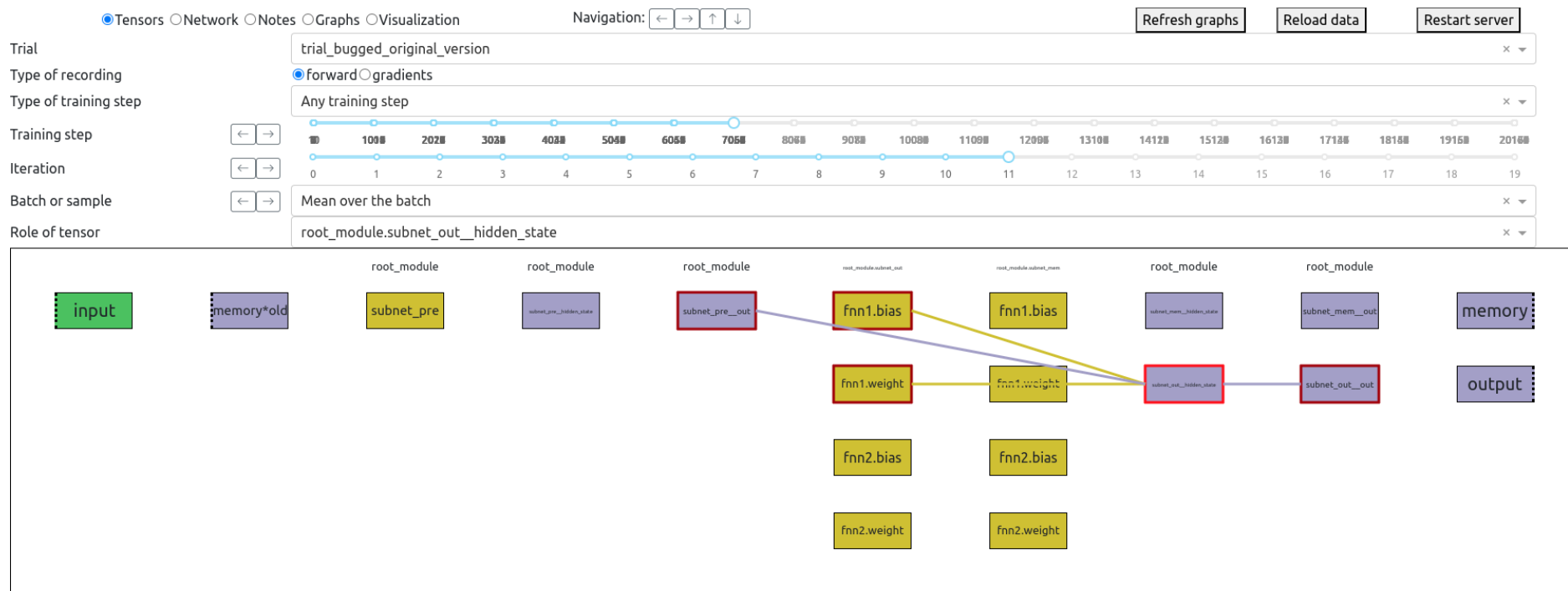# What Does BERT Look At? An Analysis of BERT's Attention



Figure 1: Examples of heads exhibiting the patterns discussed in Section 3. The darkness of a line indicates the strength of the attention weight (some attention weights are so low they are invisible).

# The Learning Interpretability Tool (LIT)

- is a **visual, interactive ML model-understanding** tool that supports text, image, and tabular data.



https://www.youtube.com/watch?v=CuRI_VK83dU&t=3s

DIGITAL

# Comgra: Computation Graph Analysis



Dietz, F., Fellenz, S., Klakow, D., & Kloft, M. (2024). *Comgra: A tool for analyzing and debugging neural networks*. In **Proceedings of the ICML 2024 Workshop on Mechanistic Interpretability**.

# Language models can explain neurons in language models

Bills, S., Cammarata, N., Mossing, D., Tillman, H., Gao, L., Goh, G., Sutskever, I., Leike, J., Wu, J., & Saunders, W. (2023, May 9). *Language models can explain neurons in language models*. OpenAI. https://openai.com/research/language-models-can-explain-neurons

# Language models can explain neurons in language models

**Step 1** — **Explain** the neuron's activations using GPT-4

**Step 2** — **Simulate** activations using GPT-4, conditioning on the explanation

Assuming that the neuron activates on

> references to movies, characters, and entertainment.

GPT-4 guesses how strongly the neuron responds at each token:

: Age of Ultron and it sounds like his role is going to play a bigger part in the Marvel cinematic universe than some of you originally thought. Marvel has a new press release that offers up some information on the characters in the film. Everything included in it is pretty standard stuff, but then there was this new

their upcoming 13-episode series for Marvel's Daredevil. It begins with a young Matt Murdock telling his blind martial arts master Stick that he lost his sight when he was 9-years-old. And then me into the present with a grateful Karen Page explaining that a masked vigilante saved her life.

offbeat , Screenshots | Follow This Author @KartikMdgl We have two images from Skyrim, which totally stumped us. They show a walking barrel, and we're not sure how exactly that happened. Check out these two images below. Some people really do some weird

ultimate in lightweight portability. Generating chest-thumping lows and crystal clear highs, the four models in the series – the XLS1000, XLS1500, XLS2000, and XLS2500 – are engineered to meet any demanding audio requirements – reliably and within budget. Every XLS

**Step 3** — **Score** the explanation by comparing the simulated and real activations

Bills, S., Cammarata, N., Mossing, D., Tillman, H., Gao, L., Goh, G., Sutskever, I., Leike, J., Wu, J., & Saunders, W. (2023, May 9). *Language models can explain neurons in language models*. OpenAI. https://openai.com/research/language-models-can-explain-neurons

DIGITAL

8

# Language models can explain neurons in language models



Step 1  **Explain** the neuron's activations using GPT-4

Step 2  **Simulate** activations using GPT-4, conditioning on the explanation

Step 3  **Score** the explanation by comparing the simulated and real activations

**Real activations:**

: Age of Ultron and it sounds like his role is going to play a bigger part in the Marvel cinematic universe than some of you originally thought. Marvel has a new press release that offers up some information on the characters in the film. Everything included in it is pretty standard stuff, but then there was this new

their upcoming 13-episode series for Marvel's Daredevil. It begins with a young Matt Murdock telling his blind martial arts master Stick that he lost his sight when he was 9-years-old. And then me into the present with a grateful Karen Page explaining that a masked vigilante saved her life.

offbeat , Screenshots | Follow This Author @KartikMdgl We have two images from Skyrim, which totally stumped us. They show a walking barrel, and we're not sure how exactly that happened. Check out these two images below. Some people really do some weird

ultimate in lightweight portability. Generating chest-thumping lows and crystal clear highs, the four models in the series – the XLS1000, XLS1500, XLS2000, and XLS2500 – are engineered to meet any demanding audio requirements – reliably and within budget. Every XLS

**Simulated activations:**

: Age of Ultron and it sounds like his role is going to play a bigger part in the Marvel cinematic universe than some of you originally thought. Marvel has a new press release that offers up some information on the characters in the film. Everything included in it is pretty standard stuff, but then there was this new

their upcoming 13-episode series for Marvel's Daredevil. It begins with a young Matt Murdock telling his blind martial arts master Stick that he lost his sight when he was 9-years-old. And then me into the present with a grateful Karen Page explaining that a masked vigilante saved her life.

offbeat , Screenshots | Follow This Author @KartikMdgl We have two images from Skyrim, which totally stumped us. They show a walking barrel, and we're not sure how exactly that happened. Check out these two images below. Some people really do some weird

ultimate in lightweight portability. Generating chest-thumping lows and crystal clear highs, the four models in the series – the XLS1000, XLS1500, XLS2000, and XLS2500 – are engineered to meet any demanding audio requirements – reliably and within budget. Every XLS
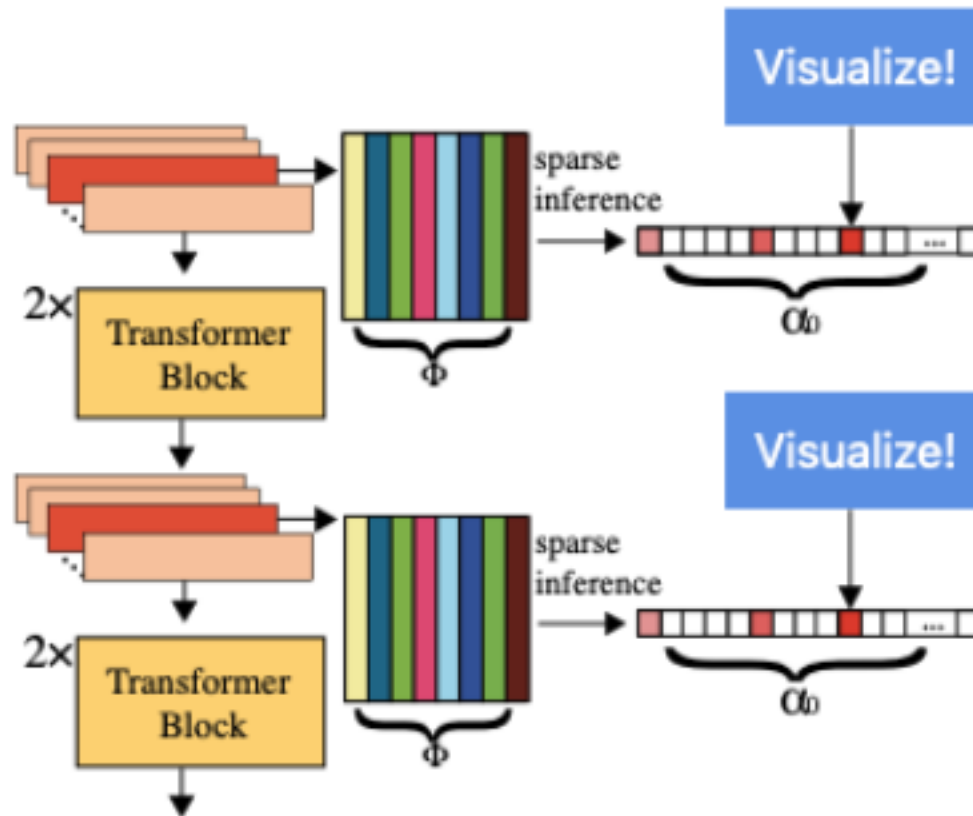
Comparing the simulated and real activations to see how closely they match, we derive a score:

0.337

Bills, S., Cammarata, N., Mossing, D., Tillman, H., Gao, L., Goh, G., Sutskever, I., Leike, J., Wu, J., & Saunders, W. (2023, May 9). *Language models can explain neurons in language models*. OpenAI. https://openai.com/research/language-models-can-explain-neurons

DIGITAL

# TransformVis



*Transformer visualization via dictionary learning: contextualized embedding as a linear superposition of transformer factors* by Zeyu Yun*, Yubei Chen*, Bruno A Olshausen, and Yann LeCun (DeeLIO Workshop@NAACL 2021).

https://transformervis.github.io/transformervis/

DIGITAL