

Mechanistic Interpretability

Looking Back
WT 2024/25

Frederick Riemenschneider



06.02.2025

Key Ideas

- Fundamentals
- Othello
- Transformer Circuits
- Activation Patching
- Dictionary Learning
- Self-conditioning
- Grokking
- Multilingual Language Models
- Frameworks
- Detecting Lies in LLMs

Your Comments

- Applications
- Subject of Investigation
- "Understanding"
- More Topics
- Outlook

Implementation Project

Key Ideas

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

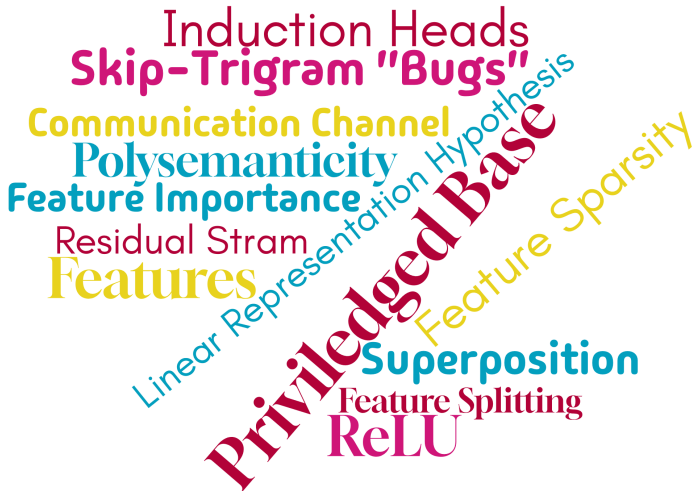
Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project



Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project

- board game internally represented
- non-linear and linear probes
- black/white vs. mine/yours

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project

Transformer Circuits

- IOI task
- responsibilities
- induction heads in action
- backup name mover heads
- path patching
- automated discovery with ACDC

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

Activation Patching

- knowledge stored in specific MLPs
- knowledge as dictionary lookup
- rank-one model editing
- different ways to think about paths in a transformer
- practical applicability of model patching

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

Dictionary Learning

- language models explain neurons
- sparse autoencoders to disentangle polysemanticity
- gated sparse autoencoders
- difficulty of evaluating autoencoders

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

- expert neurons predictive of concepts
- area under the precision-recall curve
- guided (self-conditioned) text generation
- languages as concepts

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project

- optimal semantic space arrangement leads to true generalization
- artifact of suboptimal training parameters
- transformer truly understandable?
- similarities to Othello

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project

Multilingual Language Models

- language-specific vs. language-agnostic neurons
- LAPE and PROBELESS
- language models have a default language
- perturbation successful

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

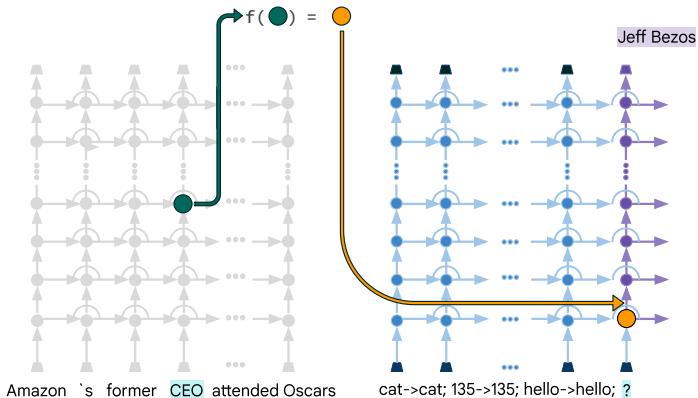
Implementation Project

Step 1:
Feeding **Source Prompt**
to **Source Model**

Step 2:
Transforming
Hidden State

Step 3:
Feeding **Target Prompt**
to **Target Model**

Step 4:
Running Execution
on **Patched Target**



Key Ideas

- Fundamentals
- Othello
- Transformer Circuits
- Activation Patching
- Dictionary Learning
- Self-conditioning
- Grokking
- Multilingual Language Models
- Frameworks
- Detecting Lies in LLMs

Your Comments

- Applications
- Subject of Investigation
- "Understanding"
- More Topics
- Outlook

Implementation Project

Detecting Lies in LLMs

- potential use-case
- two truth directions
- negated vs. affirmative statements

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

Your Comments

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project

“How can we actually apply insights from Mechanistic Interpretability to advance machine learning systems?”

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

“Are decoder-only transformers the only thing we should study?”

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

“Understanding”

“On one hand, if LLMs are not interpretable, they could be risky in areas like healthcare and autonomous driving. On the other hand, end-to-end assisted parking systems work so well that people do not worry about errors. Maybe in the future, we will accept LLMs even if we do not fully understand them.”

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

“Are there any other topics in the field of Mechanistic Interpretability that we didn’t cover in the seminar?”

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

“What are current trends?”

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

Implementation Project

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project

Implementation Project

- independent (!) reimplementations of one of the approaches
- exploration of one of my/your own ideas
- report: max. 8 pages
- deadline: 31st of March

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project

Implementation Project

- Please present your work in a way that is both **clear** and **engaging**.
- Provide a **description** of your experiments, including:
 - the **motivation** behind them
 - the **methodology**
 - a **discussion** of your results
- You can choose **any format** that suits your work, such as a Jupyter Notebook with code and documentation. If **in doubt, ask me!**
- Please ensure your results are **easily reproducible**. I personally recommend dependency management with pdm. Ideally, I should be able to run a single script to reproduce your results (or one experiment).

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project

Presentation

- By default, your project will be linked on the course page, so make sure it is appealing and accessible to your fellow students.
- If you prefer not to share your project, you can opt out. However, sharing is encouraged. It is more fun to see what others have created!

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project

- Please make sure to include a signed Declaration of Independent Work: https://www.cl.uni-heidelberg.de/-studies/Eigenstaendigkeitserklaerung_DE.pdf.
- You can find a non-binding English translation here: https://www.cl.uni-heidelberg.de/-studies/Eigenstaendigkeitserklaerung_EN.pdf
- general information: <https://www.cl.uni-heidelberg.de/studies/faq/faqMa-plagiarism.mhtml#plagiarism001>

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

Language Models

- my take:
 - proofreading, phrasing improvements, plotting and visualization assistance, basic coding support is entirely okay
 - basic coding support: using Copilot to complete trivial lines of code, auto-generating code documentation
 - **All other uses require prior discussion!**

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project

Project Ideas

- probably fruitful: apply a given technique to a new question
- Can we evaluate autoencoders when testing them on artificial data?
- Can we find syntax circuits? Are they shared across languages?
- Do language models “think” in one specific (maybe abstract or predominant) language?
- Do multilingual language models have an accent that is influenced by other languages?

Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

“Understanding”

More Topics

Outlook

Implementation Project



Key Ideas

Fundamentals

Othello

Transformer Circuits

Activation Patching

Dictionary Learning

Self-conditioning

Grokking

Multilingual Language
Models

Frameworks

Detecting Lies in LLMs

Your Comments

Applications

Subject of Investigation

"Understanding"

More Topics

Outlook

Implementation Project