# What's Inside the Box: Transparency in LLMs



*by Alosh Denny,*
*I make & break LLMs*

Input → **Black Box** → Output

Input ➤  ➤ Output

WHAT'S IN THE BOX?

quickmeme.com

"What's 13 x 13?"

What we know!

What we don't know!

"What's 13 x 13?" →

What we know!

What we assume:

"What's 13 x 13?"

What we know!

**What we assume:**

I'll multiply the ones place: 3x3 = 9

# CIRCUITS

(not the electricity thingy)

# Tracing an LLM's "circuit of reasoning"

# Image and Vision LLMs

# Bangalore:



Source: reddit

# TransformerLensOrg/
**TransformerLens**

A library for mechanistic interpretability of GPT-style language models

👥 **105**
Contributors
🔘 **105**
Issues
⭐ **3k**
Stars
🍴 **438**
Forks

**Question: In what year did World War II end?**

A: **1776**

B: **1945**

C: **1865**

Human: In what year did World War II end? ( A ) 1776 ( B ) 1945 ( C ) 1865 Assistant: Answer :

Correct Answer Head

K Q

false statements / incorrect answers

???

tokens following (A)

Human: In what year did World War II end?↵↵ ( A ) 1776 ( B ) 1945 ↵↵ ( C ) 1865 ↵↵ Assistant: Answer :

Correct Answer Head

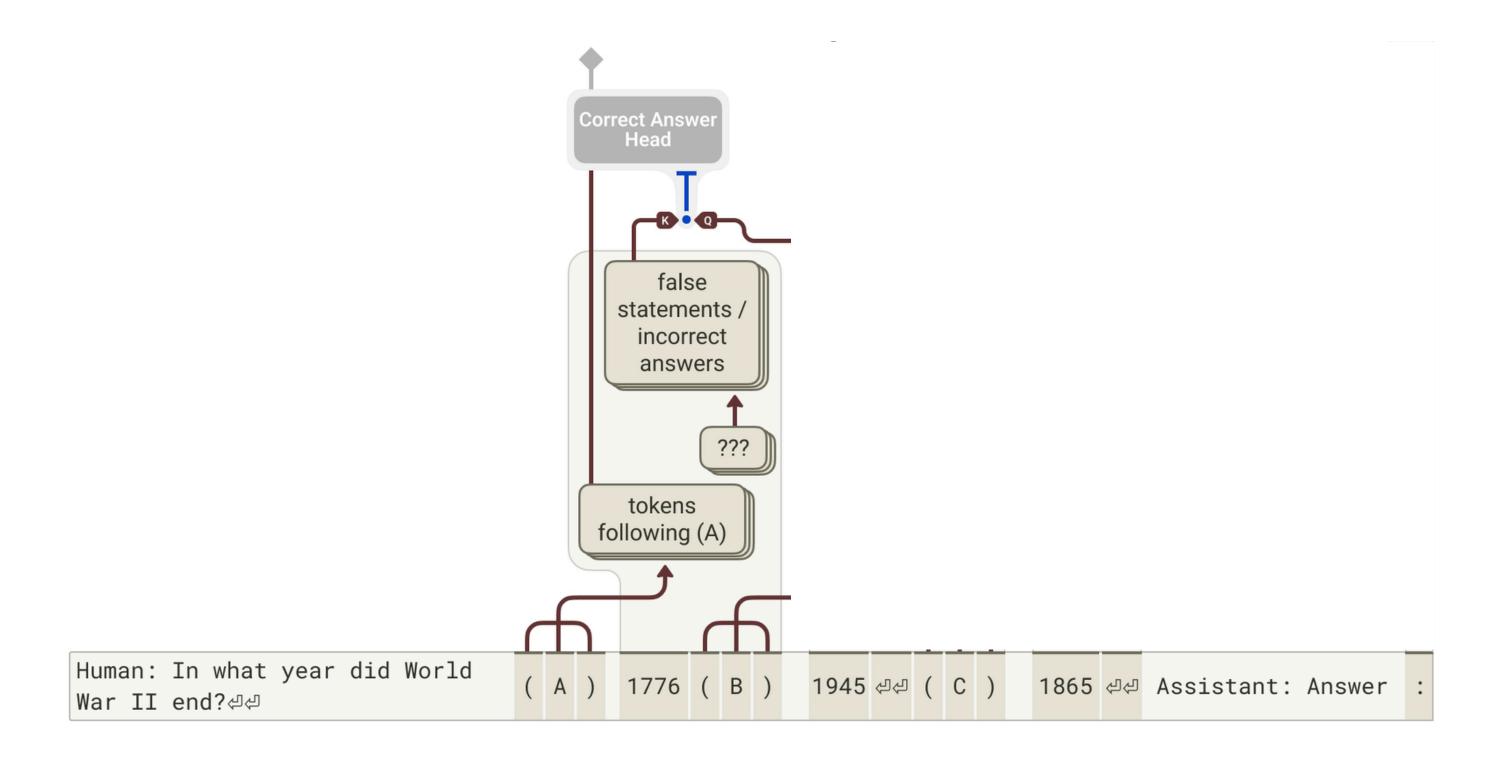Correct Answer Head

false statements / incorrect answers

correct MC answer

???

???

tokens following (A)

tokens following (B)

Human: In what year did World War II end?↵

( A ) 1776 ( B ) 1945 ↵↵ ( C ) 1865 ↵↵ Assistant: Answer :

Correct Answer Head

Correct Answer Head

Correct Answer Head

false statements / incorrect answers

correct MC answer

false statements / incorrect answers

???

???

???

tokens following (A)

tokens following (B)

tokens following (C)

Human: In what year did World War II end?↵ ( A ) 1776 ( B ) 1945 ↵↵ ( C ) 1865 ↵↵ Assistant: Answer :

Human: In what year did World War II end?⏎⏎ (A) 1776⏎(B) 1945⏎⏎(C) 1865⏎ Assistant: Answer: **B**

say B

Correct Answer Head

Correct Answer Head

Correct Answer Head

need MC answer

false statements / incorrect answers

correct MC answer

false statements / incorrect answers

???

???

???

tokens following (A)

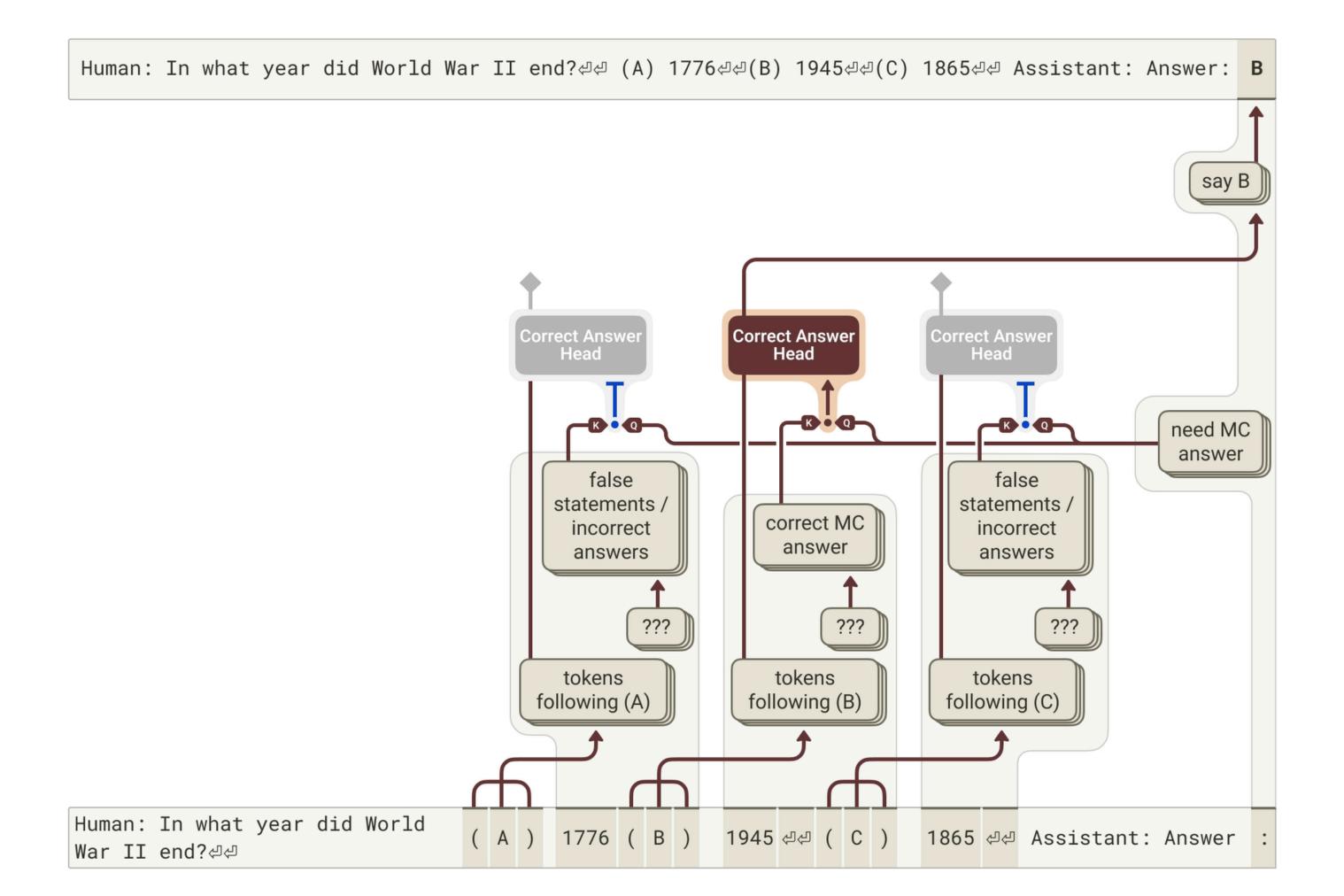tokens following (B)

tokens following (C)

Human: In what year did World War II end?⏎⏎ ( A ) 1776 ( B ) 1945 ⏎⏎ ( C ) 1865 ⏎ Assistant: Answer :

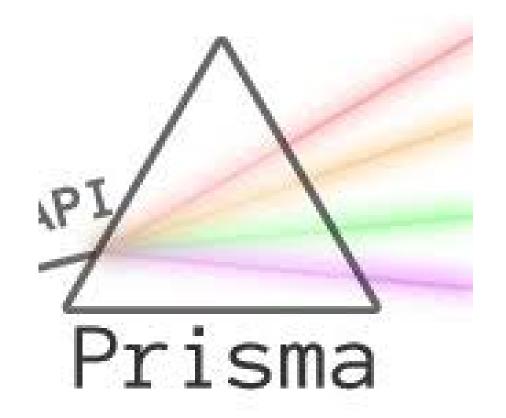# Prisma-Multimodal/**ViT-Prisma**

ViT Prisma is a mechanistic interpretability library for Vision and Video Transformers (ViTs).

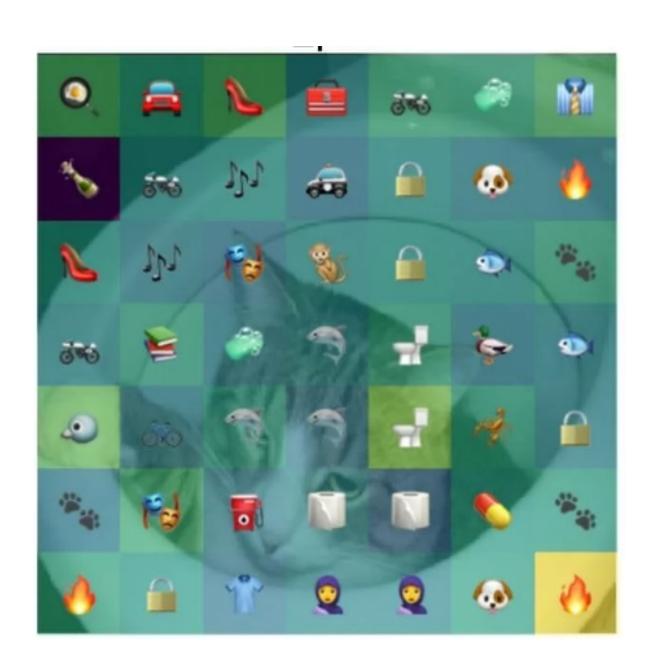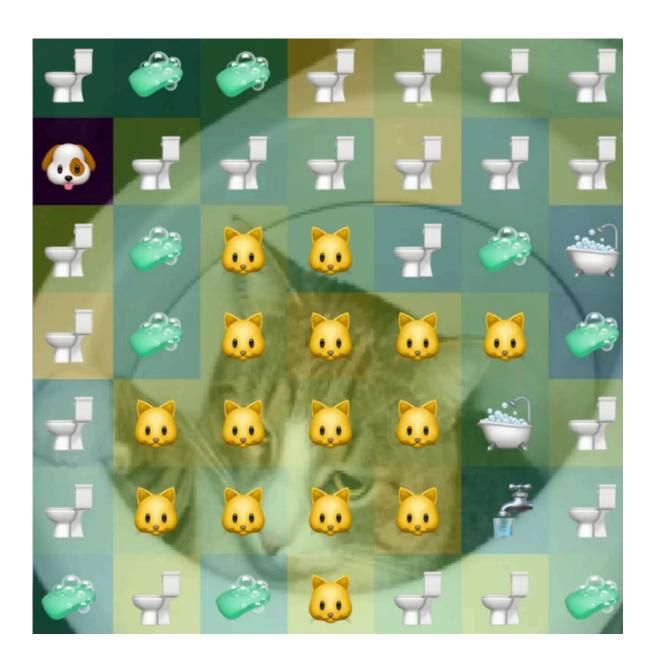| 10 Contributors | 5 Used by | 3 Discussions | 305 Stars | 35 Forks |
|---|---|---|---|---|

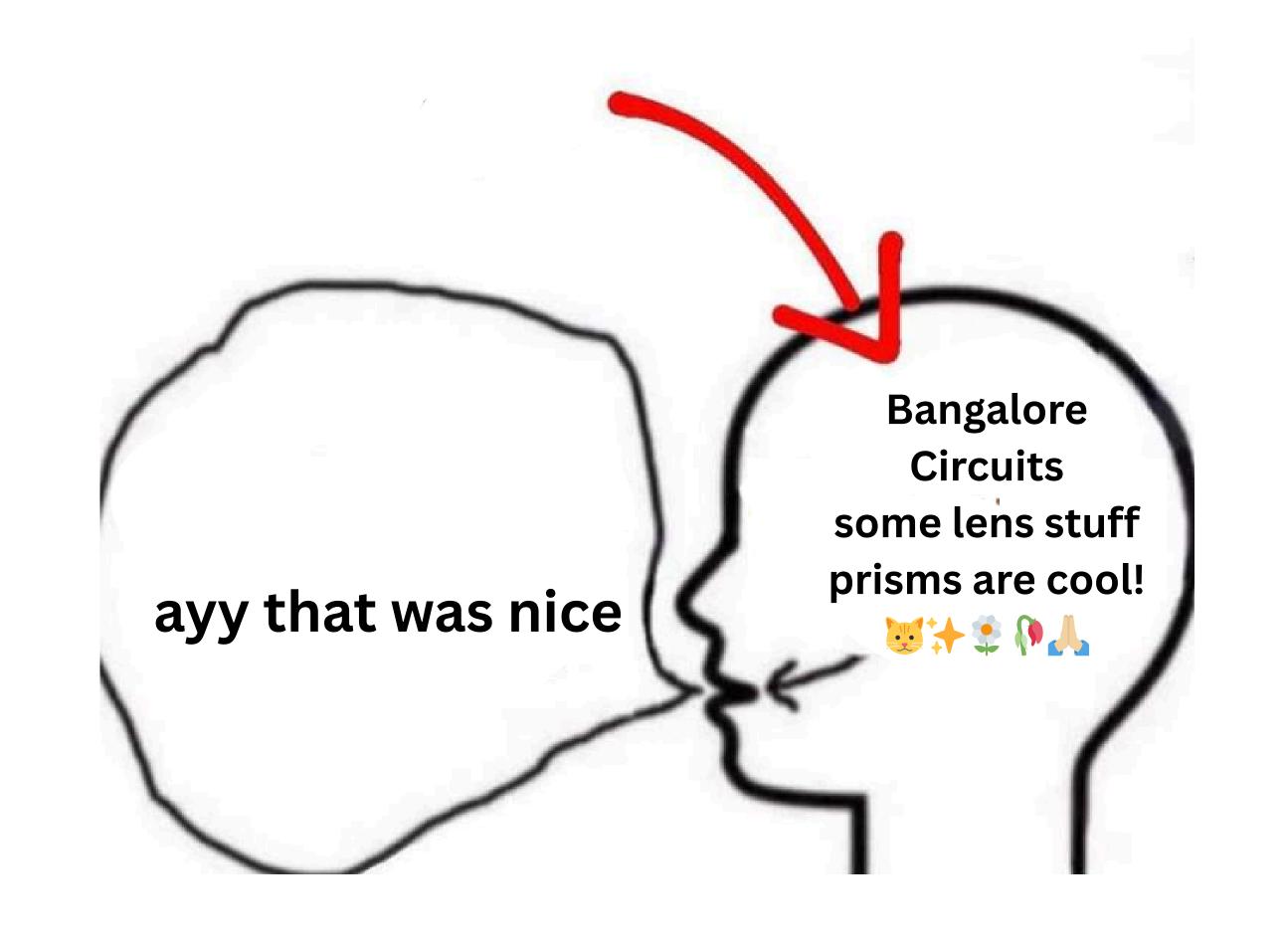https://www.lesswrong.com/posts/kobJymvvcvhbjWFKe/laying-the-foundations-for-vision-and-multimodal-mechanistic

Layer 0

**Layer 7**

**Layer 11**

# If you loved today, I'm on Linkedin!