Mechanistic?

Naomi Saphra* and Sarah Wiegreffe*



arxiv.org/abs/2410.09087



Here's the image representing "mechanistic interpretability" with a network of interlocking gears, circuits, and nodes lighting up to convey understanding within a mechanical and neural structure. Let me know if you'd like any adjustments or additional details!

DALL-E 3 rendition, 11/14/24

Mechanistic Interpretability is extremely hot

Mechanistic Interpretability is extremely hot

...



Sasha Rush 🤣 @srush_nlp

I recently asked pre-PhD researchers what area they were most excited about, and overwhelmingly the answer was "mechanistic interpretability". Not sure how that happened, but I am interested how it came about.

Last edited 10:11 AM · Jan 23, 2024 · 284.5K Views

Mechanistic Interpretability is extremely hot-but what is it?



Sasha Rush 🤣 @srush_nlp



I recently asked pre-PhD researchers what area they were about, and overwhelmingly the answer was "mechanistic interpretability". Not sure how that happened, but I am inte came about.

Last edited 10:11 AM · Jan 23, 2024 · 284.5K Views

I still don't totally understand the difference between "mechanistic" and "non-mechanistic" interpretability but it seems to be mainly a distinction of the authors' social network?

...

5:11 PM · Jan 23, 2024 · 10.9K Views

Mechanistic Interpretability is extremely hot-but what is it?



Sasha Rush 🤣 @srush_nlp



I recently asked pre-PhD researchers what area they were about, and overwhelmingly the answer was "mechanistic interpretability". Not sure how that happened, but I am inte came about.

Last edited 10:11 AM · Jan 23, 2024 · 284.5K Views

I still don't totally understand the difference between "mechanistic" and "non-mechanistic" interpretability but it seems to be mainly a distinction of the authors' social network?

...

5:11 PM · Jan 23, 2024 · 10.9K Views



is mechanic interpretability a sexier way of saying interpretability?

5:52 AM · Jul 28, 2024 · 11.9K Views

Mechanistic Interpretability is extremely hot-but what is it?



Sasha Rush 🤣 @srush_nlp Jacob Andreas @jacobandreas

I recently asked pre-PhD researchers what area they were about, and overwhelmingly the answer was "mechanistic interpretability". Not sure how that happened, but I am inte came about.

Last edited 10:11 AM \cdot Jan 23, 2024 \cdot **284.5K** Views



is mechanic interpretability a sexier way of saying int

5:52 AM · Jul 28, 2024 · 11.9K Views

I still don't totally understand the difference between "mechanistic" and "non-mechanistic" interpretability but it seems to be mainly a distinction of the authors' social network?

...

...

5:11 PM · Jan 23, 2024 · 10.9K Views



Andrew Gordon Wilson @andrewgwils

Did they seem to know much about it and the foundations? I've also noticed a major increase in interest in this area, and alignment, but I suspect unfortunately for many it's just trendy buzzwords.

8:11 PM · Jan 23, 2024 · 1,866 Views

Interpretability picks up in deep learning, NLP, vision

2015-2017



Figure 5: Saliency heatmap for for "I hate the movie ." Each row corresponds to saliency scores for the correspondent word representation with each grid representing each dimension.

Li et al, 2016

present







- **1.** Narrow technical: understanding neural networks through causal mechanisms implemented by their internal components
- 2. Broad technical: any research describing the internals of a model
- 3. Narrow cultural: any research originating from the MI community
- 4. Broad cultural: any research in the field of AI-especially LM-interpretability

- **1.** Narrow technical: understanding neural networks through causal mechanisms implemented by their internal components
- 2. Broad technical: any research describing the internals of a model
- 3. Narrow cultural: any research originating from the MI community
- 4. Broad cultural: any research in the field of AI-especially LM-interpretability

- **1.** Narrow technical: understanding neural networks through causal mechanisms implemented by their internal components
- 2. Broad technical: any research describing the internals of a model
- 3. Narrow cultural: any research originating from the MI community
- 4. Broad cultural: any research in the field of AI-especially LM-interpretability











- Causal models are made up of causal mechanisms
- Causal mechanism = a function that transforms some subset of model variables (causes) into another subset (outcomes or effects)
 - See the analogy?



- Causal models are made up of causal mechanisms
- Causal mechanism = a function that transforms some subset of model variables (causes) into another subset (outcomes or effects)
 - See the analogy?

MI = research that discovers causal mechanisms explaining all or some part of the change from neural network input to output at the level of intermediate model representations



Halpern & Pearl, 2005

- Example: induction heads
 - Question: how does LM predict "B" given "ABABA"?
 - Answer:
 - Attention heads search for a previous occurrence of "A"
 - Other heads then attend to the token that follows it



- Recent work proposes for an **even narrower definition**: characterizing a **complete**, **end-to-end** pathway from model inputs to outputs
- This definition isn't yet widely agreed upon
 - In part because it excludes work like Induction Heads

- Recent work proposes for an **even narrower definition**: characterizing a **complete**, **end-to-end** pathway from model inputs to outputs
- This definition isn't yet widely agreed upon
 - In part because it excludes work like Induction Heads

MI = research that discovers causal mechanisms explaining **all** or some part of the change from neural network input to output at the level of intermediate model representations

- **1.** Narrow technical: understanding neural networks through causal mechanisms implemented by their internal components
- -> 2. Broad technical: any research describing the internals of a model
 - 3. Narrow cultural: any research originating from the MI community
 - 4. Broad cultural: any research in the field of AI-especially LM-interpretability



Coinage of Mechanistic Interpretability

"... reverse engineering the algorithms implemented by neural networks into human-understandable mechanisms, often by **examining the weights and activations** of neural networks to identify circuits ... that implement particular behaviors."

Coinage of Mechanistic Interpretability

"... reverse engineering the algorithms implemented by neural networks into human-understandable mechanisms, often by **examining the weights and activations** of neural networks to identify circuits ... that implement particular behaviors."

MI = *any* inspection of intermediate model representations or weights

Olah et al. 2020, Elhage et al. 2021, MI workshop

- **1.** Narrow technical: understanding neural networks through causal mechanisms implemented by their internal components
- 2. Broad technical: any research describing the internals of a model
- ► 3. Narrow cultural: any research originating from the MI community
 - 4. Broad cultural: any research in the field of AI-especially LM-interpretability



















Clash of Communities

Mechinterp community members rarely attend *CL or BlackboxNLP, despite outreach.



Yonatan Belinkov @boknilev

We are interested! **#blackboxNLP** has been the largest **#nlproc** workshop for several years now. And we have an interpretability track in all main **#nlproc** confs! Please submit your work to be reviewed in such venues.

🥵 Neel Nanda 🤣 @NeelNanda5 · Jan 7, 2023

Replying to @NeelNanda5

My (very biased) view is that mech interp is just a really promising and exciting angle on understanding what's going on inside neural networks, and am fairly confused by why more academics don't already seem to be interested! Though more and more are engaging, and this is great.

11:42 AM · Jan 7, 2023 · 14.6K Views											
Q 3		tl 5	• 47	7	土						
	Post you	r reply			Reply						
	Yonatan Belinkov @boknilev · Jan 7, 2023 Even if you disagree with other approaches to interpretability, I think engagement through common conferences would help the community grow. @NeelNanda5										
	Q	t٦	♡ 5	ı ₁ 627	□ 土						

...

Clash of Communities

Jacob Andreas aiacobandreas

Excited to see that it's that time of the year when we reinvent probing again

😰 Dan Hendrycks 🤣 @DanHendrycks · Oct 3, 2023

Al models are not just black boxes or giant inscrutable matrices.

We discover they have interpretable internal representations, and we control these to influence hallucinations, bias, harmfulness, and whether a LLM lies.



Yoav Artzi 9 @yoavartzi

Distributional semantics? Reminds me of the "florida" example in the @omerlevy_ and @yoavgo paper from 2014. Granted, contemporary LLMs probably do it much better, but the ability is likely not new

itman ngwarts	aquaman catwoman superman manhunter dumbledore hallows half-blood malfoy snape nondeterministic	superboy aquaman catwoman batgirl evernight sunnydale garderobe blandings collinwood non-deterministic	superboy supergirl catwoman aquaman sunnydale collinwood calarts greendale millfield pauling		1	**7
ring orida	non-aeterministic computability deterministic finite-state gainesville fla jacksonville tampa lundertale	ninte-state nondeterministic buchi primality fla alabama gainesville tallahassee tayar	hotelling heting lessing hamming texas louisiana georgia california osrolina	Bas	sed	VVO]
oject-oriented	aspect-oriented smalltalk event-driven prolog domain-specific	aspect-oriented event-driven objective-c dataflow 4gl	event-driven domain-specific rule-based data-driven human-centered			
incing	singing dance dances dancers tap-dancing	singing dance dances breakdancing clowning	singing rapping breakdancing miming busking	vy*	and	Yoa

Wes Gurnee @wesg52 · Oct 4, 2023

Replying to @wesg52

For spatial representations, we run Llama-2 models on the names of tens of thousands cities, structures, and natural landmarks around the world, the USA, and NYC. We then train linear probes on the last token activations to predict the real latitude and longitudes of each place.



Yonatan Belinkov

Excited to see important work from @andyzou_jiaming, @DanHendrycks ..., on interpreting & controlling language models at representation level, to improve fairness & safety of LMs.

Unfortunately it fails to engage with a large body of work on these topics from the past ~5 years.

😭 Dan Hendrycks 🤣 @DanHendrycks · Oct 3, 2023

Al models are not just black boxes or giant inscrutable matrices.

We discover they have interpretable internal representations, and we control these to influence hallucinations, bias, harmfulness, and whether a LLM lies.

Show more



3:41 PM · Oct 4, 2023 · 43K Views

- **1.** Narrow technical: understanding neural networks through causal mechanisms implemented by their internal components
- 2. Broad technical: any research describing the internals of a model
- 3. Narrow cultural: any research originating from the MI community
- → 4. Broad cultural: any research in the field of AI—especially LM—interpretability

Raise your hand if you identify with the "NLP interpretability" community.

Raise your hand if you identify with the "mechanistic interpretability" community.









We are all mechanistic now

- → We can still be precise about technical methods
 - An increased focus on causality is 4

We are all mechanistic now

- → MI has brought a lot of excitement, interest, opportunities, and research findings to the field.
- → We have shared motivations: social responsibility, intellectual curiosity, a desire to build better NLP systems, and a **belief that we should** understand the tools we use.
- \rightarrow Why not also aim to connect?