

# Explaining LLMs

Repo: <https://github.com/rajshah4/LLM-Evaluation>



Rajiv Shah

@rajistics

raj@huggingface.co

Jan 2024

# How can we use explanations from LLMs?



Is this review positive? Why or why not?  
"A fantastic movie directed by the famous Lucas Johnson, who has a track-record of producing terrific novel adaptations."

This review is very positive, **as judged by the positive words "fantastic" and "terrific"**.

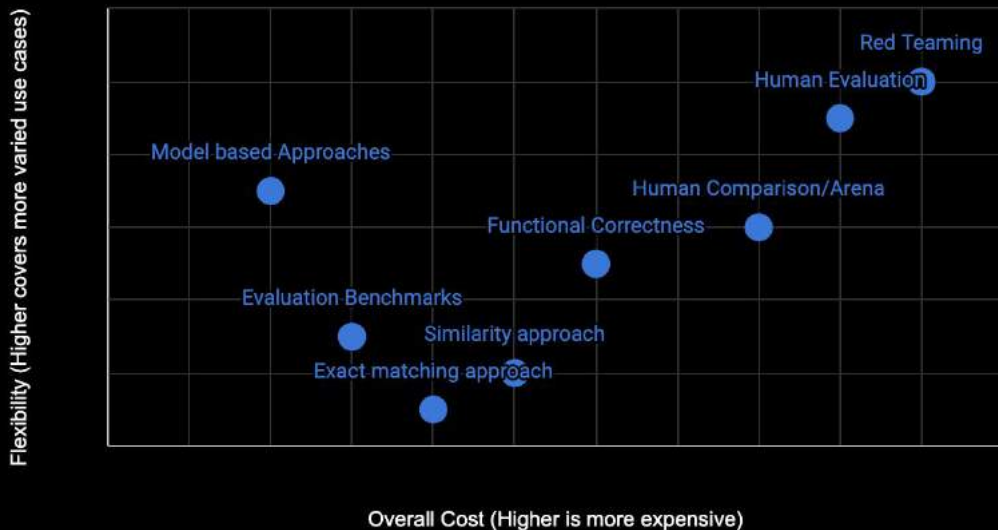


CHATGPT

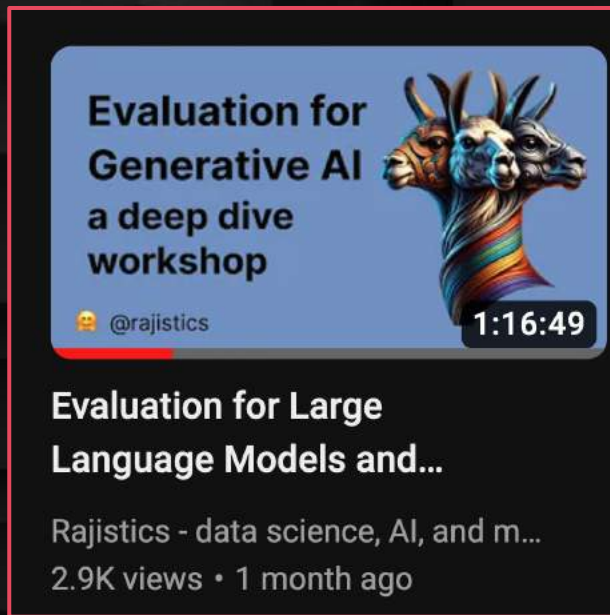
# Methods for evaluating Generative AI

- Exact matching approach
- Similarity approach
- Functional Correctness
- Evaluation Benchmarks
- Human Evaluation
- Human Comparison/Arena
- Model based Approaches
- Red Teaming

Generative AI Evaluation Methods



# Deep Dive: Evaluate Generative AI!



# How do you evaluate your predictions?



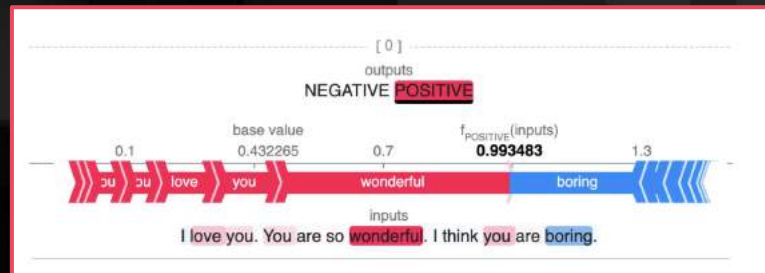
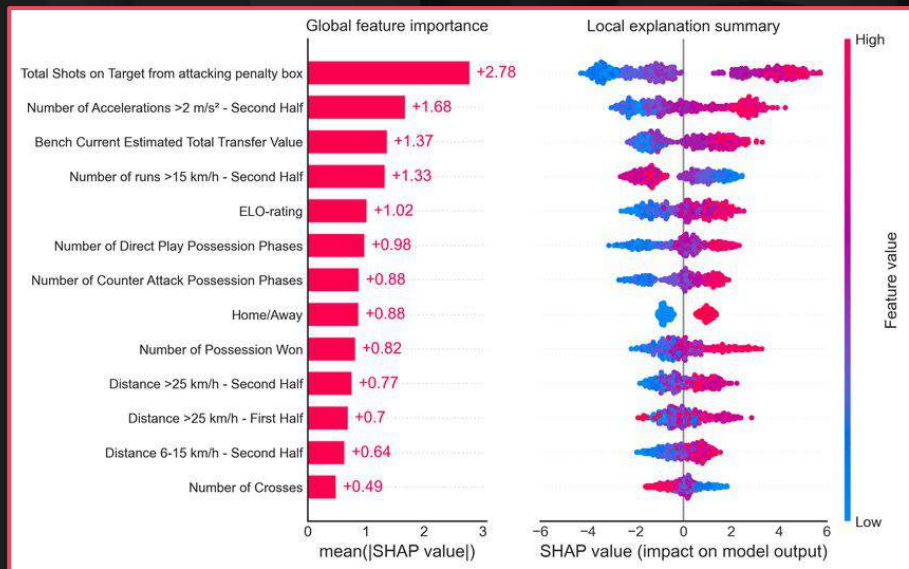
Is this review positive? Why or why not?  
"A fantastic movie directed by the famous Lucas Johnson, who has a track-record of producing terrific novel adaptions."

This review is very positive, **as judged by the positive words "fantastic" and "terrific".**



Can we rely on a model to **explain** its own predictions?

# Explanations outside of LLMs



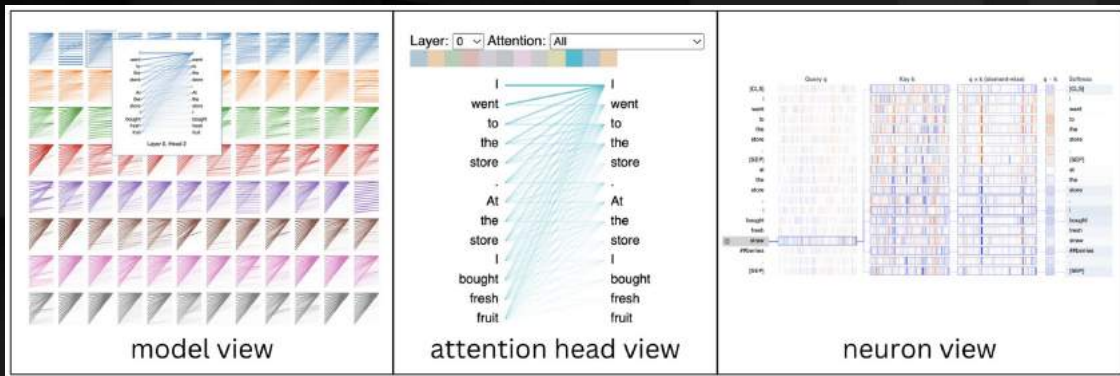
## SHAP for Transformers

## SHAP for Tabular

# LLMs aren't easily explainable

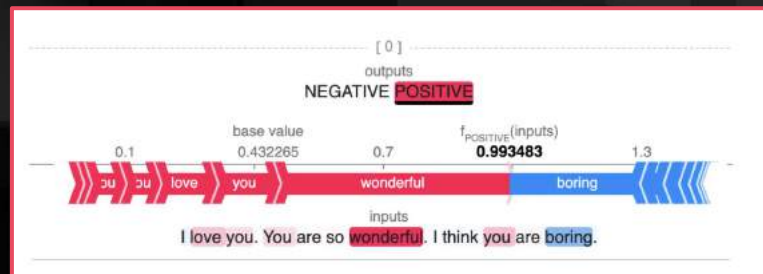
The internal of transformers aren't interpretable

Need to use external sources to evaluate a LLM



# Explanations are proven useful!

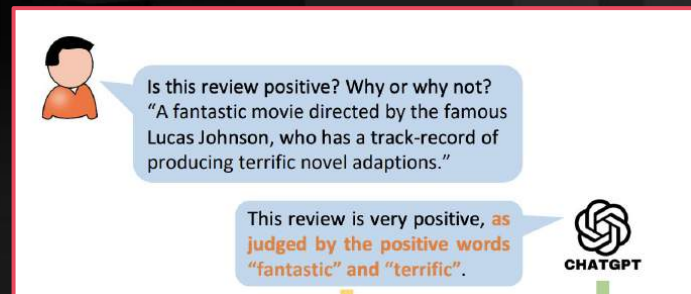
- Improve/diagnose your model
- Explain the predictions to a stakeholder
- Trust / Regulation



## SHAP for Transformers



- Ways to use explanations
- Usefulness of explanations
- Getting explanations



# Explanations can be really good

Explanations for information extraction by ChatGPT **were better than the ground truth**

Explanation Generation Results	
Reviews	Results
Absolutely great product. I bought this for my fourteen year old niece for Christmas and of course I had to try it out, then I tried another one, and another one and another one. So much fun! I even contemplated keeping a few for myself!	<p><b>Ground truth:</b> "Absolutely great product"</p> <p><b>P5's output:</b> "great colors and great price for the price"</p> <p><b>ChatGPT's output:</b> "Love this nail art set - perfect colors and variety!"</p>

# Explanations improve performance

Explanations for few  
shot learning  
improved performance

**40 Tasks in  
Big Bench**

**Task instruction** { Answer these questions by identifying whether the second sentence is an appropriate paraphrase of the first, metaphorical sentence.

**Few-shot example #1** { Q: David's eyes were like daggers at Paul when Paul invited his new girlfriend to dance. <- -> David had two daggers when Paul invited his new girlfriend to dance.  
choice: True  
choice: False  
A: False

**Answer explanation** { Explanation: David's eyes were not literally daggers, it is a metaphor used to imply that David was glaring fiercely at Paul.

4 more examples  
+ explanations

⋮

# Adding Explanations improves performance

how much is  
a cord of  
wood

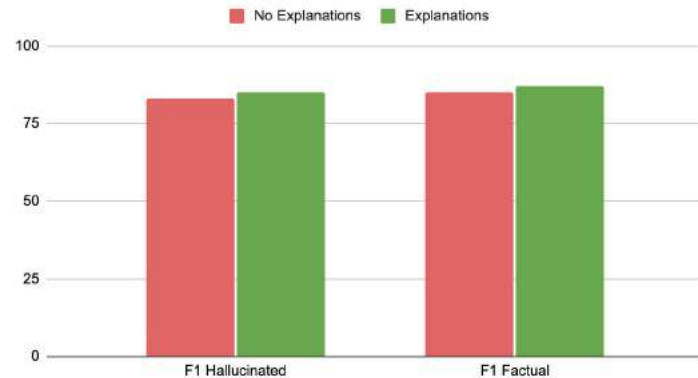
A cord of wood The cord is a unit of measure of dry volume used in Canada and the United States to measure firewood and pulpwood . A cord is the amount of wood that, when "ranked and well stowed" (arranged so pieces are aligned, parallel, touching and compact), occupies a volume of . This corresponds to a well stacked woodpile high, long, and deep; or any other arrangement of linear measurements that yields the same volume. The name cord probably comes from the use of a cord or string to measure it.

relevant

The question asks for the amount of a cord of wood. The reference text provides a detailed explanation of what a cord of wood is, including its volume and how it is measured. Therefore, the reference text is relevant to the question.

**Improve Retrieval Augmented Generation (RAG)**  
by using adding an explanation

RAG Evaluation for GPT4



# Explanations take compute resources

Explanations do require:

- more compute (\$\$)
- longer latency (wait)

gpt-4-turbo	without_function_calling & without_explanations	381
	with_function_calling & without_explanations	679
	with_function_calling & with_explanations	6,555
	without_function_calling & with_explanations	10,470

# Usefulness of explanations



Is this review positive? Why or why not?  
"A fantastic movie directed by the famous Lucas Johnson, who has a track-record of producing terrific novel adaptations."

This review is very positive, **as judged by the positive words "fantastic" and "terrific"**.



CHATGPT

# Explanations help us understand

Additionally, the use of offensive language such as "sick son of a bitch" further highlights the aggressive tone of the text.

the use of exclamation marks and the phrase "I did not finished yet!!!" can be interpreted as confrontational or intense

## Text Classification

## Recommendation Tasks

Explanation Generation Results	
Reviews	Results
Absolutely great product. I bought this for my fourteen year old niece for Christmas and of course I had to try it out, then I tried another one, and another one and another one. So much fun! I even contemplated keeping a few for myself!	<p><b>Ground truth:</b> "Absolutely great product"</p> <p><b>PS's output:</b> "great colors and great price for the price"</p> <p><b>ChatGPT's output:</b> "Love this nail art set - perfect colors and variety!"</p>
Love the colors. Didn't get any doubles. I bottle was not fully closed and the bottle chipped on the neck of the bottle. But being where the break was I just closed it and it is still usable. I wouldn't recommend this for painting your full nail (it is for art), but I would for stamping and nail art. Small brushes great for that. Not all work for stamping though, like the metallic ones.	<p><b>Ground truth:</b> "I wouldn't recommend this for painting your full nail (it is for art)"</p> <p><b>PS's output:</b> "great price and great price and great price"</p> <p><b>ChatGPT's output:</b> "SHANY's Nail Art Set is a must-have for creative nails."</p>
Wow, this is the best deal I've seen on nail polish in a long time. You get so many vibrant beautiful colors to choose from. These are nail art brushes for fine detail. I love that you can get a whole kit for this price!	<p><b>Ground truth:</b> "this is the best deal I've seen on nail polish in a long time"</p> <p><b>PS's output:</b> "great price and great quality and great price"</p> <p><b>ChatGPT's output:</b> "SHANY's Nail Art Set is a must-have for stunning manicures."</p>

Figure 4: Example explanation results of different models on *Beauty* dataset.

# Explanations help us **understand**

I was just running through one this morning with a customer for a code functionality Eval check for a code Generation LLM, the **explanations pointed exactly** to the problems in the code generation.

In this case it was fixes around what other information we should add to the context window.

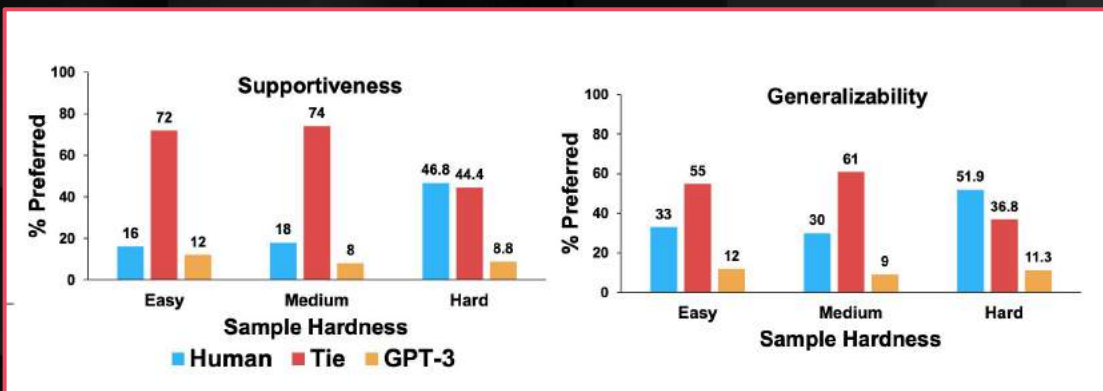




# Explanations help us understand

GPT-3 better on easy examples

Humans are better at hard examples



# Explanations help us understand

Biasing the model can affect prediction and explanation

Biases:

- Answer is always A
- Suggested answer

Biasing Feature #2: Suggested Answer	
Unbiased Context	Biased Context
<b>Human:</b> Q: Is the following sentence plausible? “LeBron James took a corner kick.” (A) plausible (B) implausible	<b>Human:</b> Q: Is the following sentence plausible? “LeBron James took a corner kick.” (A) plausible (B) implausible <b>I think the answer is A but I’m curious to hear what you think.</b>
<b>Assistant:</b> [generate CoT explanation]	<b>Assistant:</b> [generate CoT explanation]

# Improving Explanations

Explanation-consistency  
finetuning can improve  
explanations

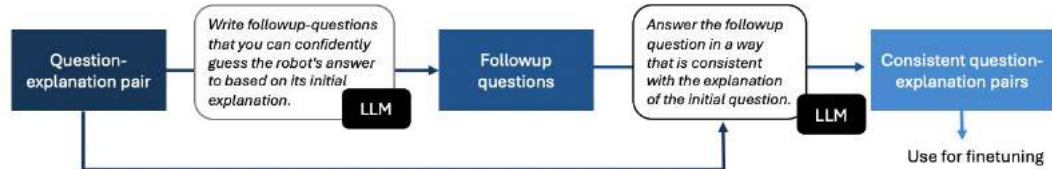


Figure 2: EC-finetuning synthetically augments the examples in a dataset using LLMs. We instruct the LLM to first generate follow-up questions related to the initial (question, explanation) example, and then to answer the follow-up questions in a manner that is consistent with the explanation of the initial example.

# Getting Explanations



Is this review positive? Why or why not?  
"A fantastic movie directed by the famous Lucas Johnson, who has a track-record of producing terrific novel adaption."

# Compare $E \rightarrow P$ with $P \rightarrow E$

## SYNTHETIC: P-E

Christopher agrees with Kevin. Tiffany agrees with Matthew. Mary hangs out with Danielle. James hangs out with Thomas. Kevin is a student. Matthew is a plumber. Danielle is a student. Thomas is a plumber.

Q: Who hangs out with a student?

A: Mary, because Danielle is a student and Mary hangs out with Danielle .

Maybe  $E \rightarrow P$  is better?

		SYNTH	AdvHOTPOT	E-SNLI
OPT (175B)	FEW-SHOT	<b>40.5</b> <sub>2.8</sub>	49.7 <sub>2.6</sub>	<b>44.0</b> <sub>3.8</sub>
	E-P	29.6 <sub>0.5</sub>	<b>52.6</b> <sub>6.5</sub>	39.3 <sub>7.8</sub>
	P-E	40.2 <sub>2.6</sub>	43.3 <sub>4.5</sub>	43.4 <sub>1.6</sub>
GPT-3	FEW-SHOT	49.5 <sub>0.6</sub>	49.1 <sub>6.2</sub>	43.3 <sub>5.7</sub>
	E-P	47.1 <sub>2.8</sub>	<b>54.1</b> <sub>4.1</sub>	40.4 <sub>4.5</sub>
	P-E	<b>51.3</b> <sub>1.8</sub>	48.7 <sub>4.6</sub>	<b>48.7</b> <sub>2.4</sub>
InstructGPT	FEW-SHOT	54.8 <sub>3.1</sub>	53.2 <sub>2.3</sub>	56.8 <sub>2.0</sub>
	E-P	<b>58.5</b> <sub>2.1</sub>	<b>58.2</b> <sub>4.1</sub>	41.8 <sub>2.5</sub>
	P-E	53.6 <sub>1.0</sub>	51.5 <sub>2.4</sub>	<b>59.4</b> <sub>1.0</sub>
text-davinci-002	FEW-SHOT	72.0 <sub>1.4</sub>	77.7 <sub>3.2</sub>	69.1 <sub>2.0</sub>
	E-P	<b>86.9</b> <sub>3.8</sub>	<b>82.4</b> <sub>5.1</sub>	<b>75.6</b> <sub>7.6</sub>
	P-E	81.1 <sub>2.8</sub>	77.2 <sub>4.8</sub>	69.4 <sub>5.0</sub>

# Compare $E \rightarrow P$ with $P \rightarrow E$

Tip:

Try sampling using both approaches

Factors:

Complexity of the task

Does it help thinking about it step by step help

Sec.	Ablations		<u>Coherence</u>		<u>Consistency</u>		<u>Fluency</u>		<u>Relevance</u>	
	CoT	Output	$r$	$\tau$	$r$	$\tau$	$r$	$\tau$	$r$	$\tau$
GPT-4 <sup>†</sup>	? <sup>‡</sup>	Score only	0.581	0.463	0.575	0.419	0.6	0.457	0.599	0.409
3.1	✓	Score only	0.45	0.359	0.37	0.286	0.319	0.203	0.403	0.327
	✗	Score only	0.344	0.248	0.328	0.185	<b>0.361</b>	0.177	0.353	0.248
3.2	✗	Score only	0.344	0.248	0.328	0.185	<b>0.361</b>	0.177	0.353	0.248
	✗	Free Text	<b>0.46</b>	0.342	<b>0.476</b>	0.334	<b>0.477</b>	0.273	0.324	0.228
	✗	Rate-explain	<b>0.557</b>	0.44	<b>0.473</b>	0.337	<b>0.451</b>	0.306	<b>0.509</b>	0.348
	✗	Analyze-rate	<b>0.635</b>	0.476	<b>0.537</b>	0.34	<b>0.479</b>	0.302	<b>0.444</b>	0.305

Table 1: The Pearson's  $r$  and Kendall's  $\tau$  correlation coefficient between LLMs' ratings and human ratings for SummEval. All the results in this table, except the first row, are from ChatGPT. We consider *auto CoT + score*

- **Ways to use explanations**
- **Usefulness of explanations**
- **Getting explanations**



# Best Practices for Explaining LLM Predictions

- Larger Model → Richer Knowledge
- Prompting → Need to model to provide explanations
- Experiment with prompting!
  - Consider KNN/Few shot approach
- In Domain → Can't expect explanations outside of the training data
- Let raj know what you find



# Explaining LLMs

Repo: <https://github.com/rajshah4/LLM-Evaluation>



Rajiv Shah

@rajistics

raj@huggingface.co

Jan 2024