



AI Interpretability AKSHAT GUPTA

GUEST LECTURE, CS 188 UC Berkeley

April 29, 2025

AKSHAT GUPTA



PHD STUDENT, UC BERKELEY (2023 - PRESENT)

- AI RESEARCH ENGINEER, JPMorgan Chase, NY (2021-2023)
- MS ECE, Carnegie Mellon University (2020-2021)
- MS PHYSICS, Technical University of Munich (2017-2019)
- BTECH EE, Indian Institute of Technology Mandi (2013-2017)

ADVISOR: Gopala Anumanchipalli AFFILIATIONS : Berkeley Speech Group, BAIR

<u>Research Areas</u>: Interpretability, Knowledge Editing, Reasoning and Poker

Dario Amodei



The Urgency of Interpretability

April 2025

1 Understand AI decision making process

1 Understand AI decision making process

2 Predict and prevent unintended behaviors

1 Understand AI decision making process

2 Predict and prevent unintended behaviors

3 Scientific curiosity

1 Understand AI decision making process

2 Predict and prevent unintended behaviors

3 Scientific curiosity

4 Improve models!













Pick the Largest Dot Product



PROJECT - 1: KNOWLEDGE EDITING

What does it mean to "Edit" a model?

- "Editing" is usually refers to "targeted" updates made to LLMs
- In "knowledge editing", we usually do the following operations:
 - Correct incorrect facts
 - Update obsolete facts
 - Add new facts
 - Remove incorrect/sensitive/private information (unlearning)



Why "Edit" a model?

- To update stale information
- To delete sensitive/private information (safety/privacy)
- Enhanced model interpretability
- Continual learning

Reliability Previous works (Huang et al., 2023; De Cao et al., 2021; Meng et al., 2022) define a reliable edit when the post-edit model f_{θ_e} gives the target answer for the case (x_e, y_e) to be edited. The reliability is measured as the average accuracy on the edit case:

$$\mathbb{E}_{x'_{e}, y'_{e} \sim \{(x_{e}, y_{e})\}} \mathbb{1} \left\{ \operatorname{argmax}_{y} f_{\theta_{e}} \left(y \mid x'_{e} \right) = y'_{e} \right\}$$
(2)

Generalization The post-edit model f_{θ_e} should also edit the equivalent neighbour $N(x_e, y_e)$ (e.g. rephrased sentences). It is evaluated by the average accuracy of the model f_{θ_e} on examples drawn uniformly from the equivalence neighborhood:

$$\mathbb{E}_{x'_{e}, y'_{e} \sim N(x_{e}, y_{e})} \mathbb{1} \left\{ \operatorname{argmax}_{y} f_{\theta_{e}} \left(y \mid x'_{e} \right) = y'_{e} \right\}$$
(3)

Locality also noted as **Specificity** in some work. Editing should be implemented locally, which means the post-edit model f_{θ_e} should not change the output of the irrelevant examples in the out-ofscope $O(x_e, y_e)$. Hence, the locality is evaluated by the rate at which the post-edit model f_{θ_e} 's predictions are unchanged as the pre-edit f_{θ} model:

$$\mathbb{E}_{x_{\mathrm{e}}', y_{\mathrm{e}}' \sim O(x_{\mathrm{e}}, y_{\mathrm{e}})} \mathbb{1} \left\{ f_{\theta_{e}} \left(y \mid x_{\mathrm{e}}' \right) = f_{\theta} \left(y \mid x_{\mathrm{e}}' \right) \right\}$$

$$\tag{4}$$

Reliability Previous works (Huang et al., 2023; De Cao et al., 2021; Meng et al., 2022) define a reliable edit when the post-edit model f_{θ_e} gives the target answer for the case (x_e, y_e) to be edited. The reliability is measured as the average accuracy on the edit case:

$$\mathbb{E}_{x'_{\mathrm{e}}, y'_{\mathrm{e}} \sim \{(x_{\mathrm{e}}, y_{\mathrm{e}})\}} \mathbb{1} \left\{ \operatorname{argmax}_{y} f_{\theta_{e}} \left(y \mid x'_{\mathrm{e}} \right) = y'_{\mathrm{e}} \right\}$$
(2)

Generalization The post-edit model f_{θ_e} should also edit the equivalent neighbour $N(x_e, y_e)$ (e.g. rephrased sentences). It is evaluated by the average accuracy of the model f_{θ_e} on examples drawn uniformly from the equivalence neighborhood:

$$\mathbb{E}_{x'_{e}, y'_{e} \sim N(x_{e}, y_{e})} \mathbb{1} \left\{ \operatorname{argmax}_{y} f_{\theta_{e}} \left(y \mid x'_{e} \right) = y'_{e} \right\}$$
(3)

Locality also noted as **Specificity** in some work. Editing should be implemented locally, which means the post-edit model f_{θ_e} should not change the output of the irrelevant examples in the out-ofscope $O(x_e, y_e)$. Hence, the locality is evaluated by the rate at which the post-edit model f_{θ_e} 's predictions are unchanged as the pre-edit f_{θ} model:

$$\mathbb{E}_{x_{\mathrm{e}}', y_{\mathrm{e}}' \sim O(x_{\mathrm{e}}, y_{\mathrm{e}})} \mathbb{1} \left\{ f_{\theta_{e}} \left(y \mid x_{\mathrm{e}}' \right) = f_{\theta} \left(y \mid x_{\mathrm{e}}' \right) \right\}$$

$$\tag{4}$$

Was the edit successful?

Reliability Previous works (Huang et al., 2023; De Cao et al., 2021; Meng et al., 2022) define a reliable edit when the post-edit model f_{θ_e} gives the target answer for the case (x_e, y_e) to be edited. The reliability is measured as the average accuracy on the edit case:

$$\mathbb{E}_{x'_{e}, y'_{e} \sim \{(x_{e}, y_{e})\}} \mathbb{1} \left\{ \operatorname{argmax}_{y} f_{\theta_{e}} \left(y \mid x'_{e} \right) = y'_{e} \right\}$$
(2)

Generalization The post-edit model f_{θ_e} should also edit the equivalent neighbour $N(x_e, y_e)$ (e.g. rephrased sentences). It is evaluated by the average accuracy of the model f_{θ_e} on examples drawn uniformly from the equivalence neighborhood:

$$\mathbb{E}_{x'_{e}, y'_{e} \sim N(x_{e}, y_{e})} \mathbb{1} \left\{ \operatorname{argmax}_{y} f_{\theta_{e}} \left(y \mid x'_{e} \right) = y'_{e} \right\}$$
(3)

Locality also noted as **Specificity** in some work. Editing should be implemented locally, which means the post-edit model f_{θ_e} should not change the output of the irrelevant examples in the out-ofscope $O(x_e, y_e)$. Hence, the locality is evaluated by the rate at which the post-edit model f_{θ_e} 's predictions are unchanged as the pre-edit f_{θ} model:

$$\mathbb{E}_{x_{\mathrm{e}}', y_{\mathrm{e}}' \sim O(x_{\mathrm{e}}, y_{\mathrm{e}})} \mathbb{1} \left\{ f_{\theta_{e}} \left(y \mid x_{\mathrm{e}}' \right) = f_{\theta} \left(y \mid x_{\mathrm{e}}' \right) \right\}$$

$$\tag{4}$$

Does the edit generalize to different phrasings of the same question?

Reliability Previous works (Huang et al., 2023; De Cao et al., 2021; Meng et al., 2022) define a reliable edit when the post-edit model f_{θ_e} gives the target answer for the case (x_e, y_e) to be edited. The reliability is measured as the average accuracy on the edit case:

$$\mathbb{E}_{x'_{e}, y'_{e} \sim \{(x_{e}, y_{e})\}} \mathbb{1} \left\{ \operatorname{argmax}_{y} f_{\theta_{e}} \left(y \mid x'_{e} \right) = y'_{e} \right\}$$
(2)

Generalization The post-edit model f_{θ_e} should also edit the equivalent neighbour $N(x_e, y_e)$ (e.g. rephrased sentences). It is evaluated by the average accuracy of the model f_{θ_e} on examples drawn uniformly from the equivalence neighborhood:

$$\mathbb{E}_{x'_{e}, y'_{e} \sim N(x_{e}, y_{e})} \mathbb{1} \left\{ \operatorname{argmax}_{y} f_{\theta_{e}} \left(y \mid x'_{e} \right) = y'_{e} \right\}$$
(3)

Locality also noted as **Specificity** in some work. Editing should be implemented locally, which means the post-edit model f_{θ_e} should not change the output of the irrelevant examples in the out-ofscope $O(x_e, y_e)$. Hence, the locality is evaluated by the rate at which the post-edit model f_{θ_e} 's predictions are unchanged as the pre-edit f_{θ} model:

$$\mathbb{E}_{x'_{e},y'_{e}\sim O(x_{e},y_{e})} \mathbb{1} \{ f_{\theta_{e}} (y \mid x'_{e}) = f_{\theta} (y \mid x'_{e}) \}$$
(4)
Does the edit effect
other facts stored in the
model?

KNOWLEDGE EDITING : METHODS

Model Editing Methods

- TYPE-1: Hypernetwork based Model Editing
- TYPE-2 : Locate-then-Edit Methods
- TYPE-3 : In-context Editing

Paper 1 : Fast Model Editing at Scale (TYPE-1)

• Training a metamodel that outputs new weights of the model



Editing a Pre-Trained Model with MEND

Paper 2 : Locating and Editing Factual Associations in GPT (TYPE-2)

• Some popular methods - ROME, MEMIT





finds the target activations for the MLP matrix.

(a) Gradient descent step which (b) Target activations are used to update the second MLP matrix (in red).

Figure 2. Presenting locate-then-edit knowledge editing methods as a two-step fine-tuning process.

"Lifelong Sequential Editing without Model Degradation", Gupta et al 2025

Facts type (s, r, o) \rightarrow (Malaysia, capital, Singapore)

Editing Address:

- 1. Layer that we want to modify
- 2. Token representation we use to modify knowledge

Facts type (s, r, o) \rightarrow (Malaysia, capital, Singapore)

Editing Address:

- 1. Layer that we want to modify
- 2. Token representation we use to modify knowledge



Facts type (s, r, o) \rightarrow (Malaysia, capital, Singapore)

Editing Address:

- 1. Layer that we want to modify
- 2. Token representation we use to modify knowledge







(a) Gradient descent step which finds the target activations for the MLP matrix.

(b) Target activations are used to update the second MLP matrix (in red).

Figure 2. Presenting locate-then-edit knowledge editing methods as a two-step fine-tuning process.

"Lifelong Sequential Editing without Model Degradation", Gupta et al 2025



(a) Gradient descent step which finds the target activations for the MLP matrix.

(b) Target activations are used to update the second MLP matrix (in red).

Figure 2. Presenting locate-then-edit knowledge editing methods as a two-step fine-tuning process.



"Lifelong Sequential Editing without Model Degradation", Gupta et al 2025

MY RESEARCH : SCALING KNOWLEDGE EDITING

Background

- Popular knowledge editing methods released between 2021-2023 performed well when making singular knowledge edits.
 - But are these methods scalable?
 - Can they be solutions for continually learning models?
 - What was the effect of continuous editing on the general ability of models?
Model Editing at Scale leads to Gradual and Catastrophic Forgetting

Akshat Gupta, Anurag Rao, Gopala Anumanchipalli

UC Berkeley akshat.gupta@berkeley.edu

ACL 2024 (Findings)



AKSHAT GUPTA PhD Student, UC Berkeley



ANURAG RAO UC Berkeley (Now MS, Oxford)



GOPALA ANUMANCHIPALLI Asst. Professor, UC Berkeley

Model Editing at Scale leads to Gradual and Catastrophic Forgetting

Akshat Gupta, Anurag Rao, Gopala Anumanchipalli

UC Berkeley akshat.gupta@berkeley.edu

ACL 2024 (Findings)





(a) Sample 1



Facts Forgotten

100

80

60

40

Facts Forgotten



Downstream Performance

Rebuilding ROME : Resolving Model Collapse during Sequential Model Editing

Akshat Gupta¹, Sidharth Baskaran², Gopala Anumanchipalli¹ ¹UC Berkeley, ²Automorphic Inc. akshat.gupta@berkeley.edu, sid@automorphic.ai

EMNLP 2024 Main



AKSHAT GUPTA PhD Student, UC Berkeley



SIDHARTH BASKARAN Automorphic Inc.



GOPALA ANUMANCHIPALLI Asst. Professor, UC Berkeley

Rebuilding ROME : Resolving Model Collapse during Sequential Model Editing

Akshat Gupta¹, Sidharth Baskaran², Gopala Anumanchipalli¹

¹UC Berkeley, ²Automorphic Inc. akshat.gupta@berkeley.edu, sid@automorphic.ai

EMNLP 2024 Main



AFTER

DATASET	Implementation	Effi	Efficacy		lization	Loca	ality	Fluency	Score
		ES ↑	EM ↑	PS ↑	$\mathbf{PM}\uparrow$	NS \uparrow	$NM\uparrow$	GE ↑	$\mathbf{S}\uparrow$
CF	Original r-ROME	$99.92 \\ 99.74$	$99.68 \\ 97.79$	96.29 99.09	$71.58 \\ 70.86$	75.8 80.62	$\begin{array}{c} 10.25\\ 26.0 \end{array}$	$621.96 \\ 621.67$	89.32 92 .22
	p-ROME	99.9	99.36	97.04	63.01	80.0	5.74	621.17	91.42

A Unified Framework for Model Editing

Akshat Gupta, Dev Sajnani, Gopala Anumanchipalli UC Berkeley {akshat.gupta, sajnanidev, gopala}@berkeley.edu

EMNLP 2024 (Findings)



AKSHAT GUPTA PhD Student, UC Berkeley



DEV SAJNANI Undergrad, UC Berkeley



GOPALA ANUMANCHIPALLI Asst. Professor, UC Berkeley

A Unified Framework for Model Editing

Akshat Gupta, Dev Sajnani, Gopala Anumanchipalli UC Berkeley {akshat.gupta, sajnanidev, gopala}@berkeley.edu



Figure 1: A diagrammatic representation of the preservation-memorization objective.



EMMET - Equality-constraint Mass Model Editing Algorithm



(a) Efficacy Score (ES)



EMNLP 2024 (Findings)

Figure 13: Model - Llama2-7b. Batch size 4096.

Norm Growth and Stability Challenges in Localized Sequential Knowledge Editing

Akshat Gupta^{1*}, Christine Fang¹, Atahan Ozdemir¹, Maochuan Lu¹, Ahmed Alaa¹, Thomas Hartvigsen², Gopala Anumanchipalli¹

¹University of California Berkeley, ²University of Virginia

Outstanding Paper Award, Towards Knowledgeable Foundation Models Workshop @ AAAI 2025



Norm Growth and Stability Challenges in Localized Sequential Knowledge Editing

Akshat Gupta^{1*}, Christine Fang¹, Atahan Ozdemir¹, Maochuan Lu¹, Ahmed Alaa¹, Thomas Hartvigsen², Gopala Anumanchipalli¹

¹University of California Berkeley, ²University of Virginia

Outstanding Paper Award, Towards Knowledgeable Foundation Models Workshop @ AAAI 2025





Figure 1. The continuous growth of norm of edited MLP matrices in LLama3-8B during sequential knowledge editing, as a function of number edits.

Figure 4. Comparison between norm of edited MLP matrices and norm of unedited matrices after 5,000 and 10,000 sequential edits.

PROJECT-2: INFORMATION **PROCESSING IN MS**





AKSHAT GUPTA PhD Student, UC Berkeley



JAY YEUNG Undergrad, UC Berkeley



ANNA IVANOVA Asst. Professor, Georgia Tech



GOPALA ANUMANCHIPALLI Asst. Professor, UC Berkeley

Introduction and Motivation

- Humans require varying levels of cognitive effort depending on the nature and complexity of the information being processed.
 - While we can recall some commonplace facts almost instantly, retrieving more obscure information may take longer
 - We produce certain types of words with little conscious effort while take our time with others (function words vs content words)
- While humans dynamically adjust their cognitive effort based on the familiarity and complexity of information, large language models (LLMs) in their current form process all inputs uniformly.
 - But do all tokens truly require the full depth of an LLM's architecture, or could some tokens be processed more efficiently based on their type and context?
 - More interestingly, have LLMs implicitly developed some form of dynamic load adjustment, mirroring the way humans allocate cognitive effort?

BACKGROUND: LOGITLENS

Computations in an LLM



Computations in an LLM





$$f^l = \operatorname{LN1}(h^{l-1}) \tag{3}$$

$$a^l = \operatorname{Att}(f^l)$$
 (4)

$$g^{l} = \text{LN2}(h^{l-1} + a^{l}) \tag{5}$$

$$m^{l} = W^{l}_{proj}\sigma(W^{l}_{fc}g^{l} + b^{l}_{fc}) + b_{proj}$$
(6)

 $h^{l} = h^{l-1} + a^{l} + m^{l} (7)$

The LogitLens (TunedLens) Framework

$$h^l = h^{l-1} + a^l + m^l$$

$$\texttt{LogitLens}(h^l) = W_U \Big[\texttt{Norm}_f[h^l] \Big]$$

The LogitLens (TunedLens) Framework

$$h^l = h^{l-1} + a^l + m^l$$

$$\texttt{LogitLens}(h^l) = W_U \Big[\texttt{Norm}_f[h^l] \Big]$$



	in'i		Hain	Ġ	5		à	, si i	. 21	aut
h_out -	2	'we'	'show'	'a'	'AN'	' models'	'based'	- V	' a'	' N'
h46_out -	19 - N	' we'	' show'	' a'	'AN'	- Q	T.	- 9	' a'	' N'
h44_out -		'we'	' show'	' a'	'BM'	' models'	'based'	' models'	' a'	' N'
h42_out -		' we'	' show'	' a'	'rams'	' models'	'based'	' models'	' a'	' algorithm'
h40_out -	- 92 - 92	' we'	demonstrate	' a'	' machine'	' models'	'based'	' models'	' a'	' algorithm'
h38_out -	' we'	' we'	demonstrate	' 'neural'	'rap'	' models'	'based'	' models'	' a'	' algorithm'
h36_out -	'we'	'we'	demonstrate	' 'neural'		' models'	'based'	' models'	' a'	' algorithm'
h34_out -	'we'	' we'	demonstrate	' models'		' model'	'based'	' models'	' a'	' algorith'
h32_out -	'we'	'we'	' simulated'	' models'		' model'	'based'	' models'	' a'	' adaptive'
h30_out -	' targeted'	'we'	' found'	' a'	'rap'		'based'	en.	' which'	' hybrid'
h28_out -	' targeted'	' we'	' found'	' a'	'FP'		'based'	'rd'		
h26_out -	' targeted'	' we'	' found'	' naīve'	'FP'		'based'	'rd'		
h24_out -	' targeted'	' we'	' found'	' algorithms'			'based'	'rd'		' widely'
h22_out -	' targeted'	' we'	' found'	' camp'			'based'	'rd'		' widely'
h20_out -		' although'	' found'	' algorithm'			'based'	'rd'		' single'
h18_out -		' although'	' focus'	' camp'			'based'	'rd'		' single'
h16_out -		' unlike'	' focus'	' camp'	'MP'	"IME"	'based'	'rd'		' single'
h14_out -	' targeted'	' note'	' target'	' camp'	'MS'		'based'	'rd'	'000'	' single'
h12_out -	' target'	' unlike'	' hope'	' split'	'MP'		'based'	'rd'	'000'	' massive'
h10_out -	' updated'	' however'			'iott'		'style'	'rd'	'000'	' massive'
h8_out -	' target'	' however'	' target'	' evaluation'	'rom'		'based'	'rd'	'000'	' enormous
h6_out -	' focused'	' however'		'ees'			'based'	'rd'	' which'	' enormous
h4_out -	' target'	' however'					'sided'	'rd'	' and'	' enormous
h2_out -	' guid'	' and'	' Hardy'				'based'	'rd'	' and'	' enormous
h0_out -	'chini'	' and'		' train'			'based'	'rd'	' and'	' isolated'
	ically	1	. We	train	Ğ	d.	3	ŝ	4	36

				ive	è	is,				ameters	
. BUE		ore	(*). Ole	, langu	(*) m	With , with	, 175	billion	(m) . 62	en'i	. 1º
' N'	h_out	'oen'	'gressive'	' model'	' model'	' trained'	' a'		parameters	9	'on'
' N°	h46_out -	'oen'	'gressive'	' language'	' model'	' trained'	' a'	- 25	' word'	1.2	' on'
' N'	h44_out	'oen'	'gressive'	' learning'	' model'	' trained'	' a'	' million'	' neurons'	- V - 1	' on'
lgorithm'	h42_out ·	'oen'	'gressive'	' model'	' model'	' trained'	' a'	' million'	' neurons'	2.2	' using'
algorithm'	h40_out	'oen'	'gressive'	' model'	' model'	' trained'	' a'	' million'	' neurons'	' tuned'	' using'
lgorithm'	h38_out -	'ore'	'gressive'	' modeling'	' model'	' optimized'	' a'	' million'	' neurons'	- 9 2	' using'
algorithm'	h36_out		'vised'	' modeling'	' model'	' designed'	' three'	' million'	' neurons'	9 - P	
algorith'	h34_out -		'ceptor'	' modeling'	' modeling'	' designed'		' million'	' neurons'	- V	
	h32_out	'ogenous'	'ceptor'	' model'	' model'	' designed'	' three'	' million'		' per'	' using'
	h30_out		'ceptor'	' modeling'	' modeling'	' designed'	' three'			' per'	' sequ'
	h28_out	'ocratic'	'ceptor'	' model'	' modeling'	' capable'	' minimal'		' dollars'	' per'	' followed'
	h26_out		'actor'	' model'	' model'	' capable'	' three'		' dollars'	' per'	' followed'
widely'	h24_out	'ocratic'	'actor'	' modeling'	' model'	' capable'	' minimal'		' dollars'	' per'	' followed'
widely'	h22_out	'ocratic'	'ceptor'		' model'	' housed'	' specialized'		' dollars'	' per'	' followed'
' single'	h20_out		'ceptor'	' analysis'	' model'	' housed'	' specialized'	' million'	' dollars'	' per'	' followed'
' single'	h18_out			' analysis'	' feature'	' housed'	' respect'	' million'	' dollars'	'ized'	' including'
' single'	h16_out	'ode'	'vers'	' analysis'	' grid'	' split'	' specialized'	' million'	' dollars'	'ized'	' including'
' single'	h14_out	'ode'	'pro'	' analysis'	' grid'	' trained'	' specialized'	' million'	' dollars'	' per'	' including'
massive'	h12_out	'ode'		' analysis'	' grid'	' weights'	' respect'	' million'	' dollars'	' parameters'	' including'
massive'	h10_out	'ode'	'oms'	' analysis'	' sear'	' machine'	' respect'	' million'	' dollars'	' parameters'	' including'
normous'	h8_out	'ode'		' analysis'	' sear'	' machine'	' respect'	' million'	' dollars'	' parameters'	' mostly'
enormous'	h6_out		'ographic'	' analysis'		' whit'	' respect'	' million'	' dollars'	' parameters'	' followed'
enormous'	h4_out		'ras'	' analysis'	' sear'	' machine'	' respect'	' million'	' dollars'	' parameters'	' and'
enormous'	h2_out	' aut'	'nce'	' movement'	' skills'	' machine'	' respect'	' shades'	' dollars'	' parameters'	' and'
isolated'	h0_out	' aut'	'ore'	'gressive'	' words'	' model'	' regards'	' shades'	' dollars'	' parameters'	' and'
ari		BUT	ore	SSINE	uage	nodel	with	175	oillion	eters	
			.05	,130	ю.,	×.			. para		

LogitLens:

TunedLens:



The LogitLens (TunedLens) Framework

$$h^l = h^{l-1} + a^l + m^l$$

$$\texttt{LogitLens}(h^l) = W_U \Big[\texttt{Norm}_f[h^l] \Big]$$

$$\mathsf{TunedLens}(h^l) = \mathsf{LogitLens}(A_l h^l + b_l)$$



INFORMATION PROCESSING BASED ON GENERATED TOKEN PART-OF-SPEECH

Example Sentences

Category: DET (Determiner)

- Input: "She picked up __"
- Output: "the"

Category: ADP (Adposition, e.g., prepositions)

- Input: "He walked slowly __"
- Output: "to"

Category: NOUN

- Input: "The dog chased a __"
- Output: "squirrel"

Category: VERB

- Input: "She carefully __"
- Output: "painted"

Syntactic Structure



INFORMATION PROCESSING DURING FACT RECALL

Example sentences

The capital of France is ____



REFLECTION ON RESULTS

Hierarchy of Information

GITZ-XL



Baseline occurs around first 40% of layers.On average, a model takes 40% of layers to make prediction decision.

Early Layers - Layers before baseline are what I call early layers, or around first 40% of layers.

- Used to make decisions about simple functional words for grammatical purposes - DET, ADP, PUNCT
- Middle Layers Next 35% of layers are what I call the middle layers. They do more complex tasks like fact recall, predicting content words like nouns, verbs, and doing downstream tasks.
- Late Layers Last 25% of layers are what I call late layers. They do more complex tasks like predicting multi-token facts, more ambiguous downstream tasks like NLI.

Takeaways

- At inference times, LLMs take different amounts of time to process different kinds of information
- There is an information processing hierarchy in LLMs Seemingly easier tasks finish processing earlier than more complex tasks
- Potential applications Early exiting

PROJECT -3: LLMS AND POKER

Team



Other Members:

- Piyush Jha (PhD Student, Georgia Tech)
- Jonny Pei (Senior, UC Berkeley)
- Chris Dodla (Sophomore, UC Berkeley)

Past Members:

- Richard Yang (UC Berkeley)
- Aniket Rahane (UC Berkeley)

MOTIVATIONS

1 Are LLM any good at poker? dh I

2 Can we train LLMs to be GTO?

3 Can we make exploitative poker playing agents that can go beyond GTO?



4 Can we use LLMs to explain GTO decision and teach poker?

Are ChatGPT and GPT-4 Good Poker Players - A Pre-flop Analysis (arxiv, 2023)

dЬ

Are ChatGPT and GPT-4 Good Poker Players? -**A Pre-Flop Analysis**

Akshat Gupta

UC Berkelev akshat.gupta@berkeley.edu



A9s A8s A7s A6s A5s A4s A3s A2s ATs

UTG+1







HIGHLIGHTS

POKERBENCH: Training Large Language Models to become Professional Poker Players

Richard Zhuang¹, Akshat Gupta^{1*}, Richard Yang¹, Aniket Rahane¹, Zhengyu Li², Gopala Anumanchipalli¹

¹University of California, Berkeley; ²Georgia Institute of Technology



Accepted to AAAI 2025!



HIGHLIGHTS

PokerBench - Training Large Language Models to become Professional Poker Players (AAAI, 2025)

MODEL EVALUATION ON POKERBENCH:

		Overall	Overall Accuracy		p Accuracy	Pre-Flop Accuracy	
EVALUATION TYPE	MODEL	EM ↑	AA ↑	EM ↑	AA ↑	EM ↑	AA ↑
Pre-Trained Models	LLAMA-3 (8B)	26.02	40.03	14.96	31.25	37.77	49.30
(Few-Shot)	LLAMA-2 (70B)	36.48	48.30	32.95	41.11	40.20	55.90
	LLAMA-3 (70B)	39.16	49.78	34.30	45.40	44.30	54.40
	CHATGPT 3.5	29.96	39.69	18.75	34.19	41.80	45.50
	GPT-4	53.55	65.54	52.18	62.69	55.00	66.50
Fine-Tuned Models	Gemma (2B)	51.84	62.74	41.57	52.94	62.70	73.10
(Zero-Shot)	LLAMA-2 (7B))	78.11	79.91	76.52	79.55	79.80	80.30
	LLAMA-3 (8B)	78.26	80.64	76.52	79.07	80.10	82.30

Table 2: Performance of various pre-trained and fine-tuned LLMs on POKERBENCH.

ACKNOWLEDGEMENTS




Special Thanks!



If you want to work with me on Interpretability or Poker, please reach out to me at :

akshat.gupta@berkeley.edu

Also check out some amazing work happening in our lab -Berkeley Speech Group Contact : akshat.gupta@berkeley.edu

Thank You!