

University of Stuttgart

Institute of Industrial Automation and Software Engineering



- Presenter: Ziyao Zhou
- Supervisor: Yuchen Xia M. Sc.
- Examiner: Prof. Dr. Ing. Michael Weyrich

Investigation of the Explainability of Results Generated by Large Language Models

05.10.2023



- Motivation
- Research Question
- Related Work
- Research Work Pipeline
- Methodology
- Results and Evaluation
- Conclusion and Future Work

- Motivation
- Research Question
- Related Work
- Research Work Pipeline
- Methodology
- Results and Evaluation
- Conclusion and Future Work

Background

Motivation

Large Language Models – widely used

- InChart@PAnswer Generated
 - Over 100 million Users
 - 1.6 billion visits in June 2023¹

- Motivation
- Research Question
- Related Work
- Research Work Pipeline
- Methodology
- Results and Evaluation
- Conclusion and Future Work

Research Question



- Use mathematical ways to calculate confidence and the errors between confidence and correctness
- Visualization

- Motivation
- Research Question
- Related Work
- Research Work Pipeline
- Methodology
- Results and Evaluation
- Conclusion and Future Work

Related Work

Proxy metrics to estimate the confidence or Correctness

- **Prompt the models** to express their uncertainty in words[1]
- Use log probabilities to estimate confidence
- Use similarity metrics to estimate confidence [2]
- Calculate similarity based on semantic meaning [3]
- Using activation values of hidden layer to predict correctness[4]



(🗸

"choices":

"text": "Answer: D, confidence: 3

- Motivation
- Research Question
- Related Work
- Research Work Pipeline
- Methodology
- Results and Evaluation
- Conclusion and Future Work

Research Work Pipeline



- Motivation
- Research Question
- Related Work
- Research Work Pipeline
- Methodology
- Results and Evaluation
- Conclusion and Future Work



Dataset Collection

A labeled dataset with 857 single-choice questions:



/	A	В	С	D	Question body F	\setminus /	F 4	4 Optic	ns H	-	Answ	e
1	Source	File	Index	Correlation Score	Question		Option A	Option B	Option C	Option D	Answer	Ţ
2	https://githu	arc_easy.csv	673	0.398	The wheels and gears of a machine are greased in order to decrease		potential ene	efficiency	output	friction	D	
3	https://githu	high_school_	49	0.397	Which of the following jobs is in the non-basic employment sector?		Software en	F-10 pickup t	Oil refinery v	Parking lot a	D	
4	https://githu	obqa.csv	1466	0.385	A clean air act must be followed by a		web design o	hard drive m	math class	computer pro	В	
5	https://githu	arc_hard.csv	944	0.384	Robots can perform tasks that are dangerous for humans. What is the MAJOR limitation to the use of robots?		The assembl	The assembl	Robots requi	Robots must	В	
6	https://githu	obqa.csv	2764	0.367	What is essential for a robot to possess to walk up a flight of stairs?		electricity	skittles	ethics	lava	A	
7	https://githu	professional	121	0.345	The items for 2 machinist proficiency test have been based on 2 job analysis of machinists in 50 plants		each employ	concurrent va	empirical val	content valid	D	
8	https://githu	obqa.csv	3511	0.343	As vehicles become more efficient petro consumption		increases	stops	decreases	stay the sam	с	
9	https://githu	sociology_te	154	0.337	Which of the following industries did Blauner (1964) suggest was most alienating to its employees?		machine-bas	textile work	car manufac	chemical pro	с	
10	https://githu	arc_easy.csv	2202	0.336	Which is the first step in a design process?	\mathcal{F}	Revise the so	Describe the	Test the pose	Identify poss	В	
						/ \	<				1 1	



Calculate Confidence with logprobs(1/2)

Request:

The wheels and gears of a machine are greased in order to decrease A: potential energy B: efficiency C: output D: friction

"choices": [**Response:** "text": "\n\nD", "logprobs": { "tokens):], "token_logprobs") -0.0053596636, -0.00041011896, "top_logprobs")["\n": ".": -7.4418178, " Answer": -8.131088

}.







Calculate Confidence with logprobs(1/2)

Response:



Second Situation: no "ABCD" among the answers

 $Confidence = \frac{\sum_{i=1}^{N} e^{z_i}}{N}$



$$Confidence = \frac{1}{3}(e^{-0.58} + e^{-8.92} + e^{-0.01}) = 0.52$$

Calculate Confidence with logprobs(1/2)

• Generate 5 responses for each question



Result Table:

											()		\checkmark	\square	()		$\left(\right)$			
0) 1	. 2	3	3	4	5	6	7	8	9	Generated Answer 1	Confidence 1	Generated Answer 2	Confidence 2	Generate	Confident	Generate	Confiden	Generate	Confidence
The whee	potential	efficiency	output	friction	D	https:/	//gitarc	_easy.c	673	0.398	D	0.999968148	D	0.999968148	D	0.999968	(The answ	0	D	0.99996814
Which of	t Software	F-10 picku	Oil refine	Parking l	o D	https:/	//githig	h_scho	49	0.397	А.	0.999675588	A	0.999987927	А	0.999988	А	0.999988	А	0.99998792
Robots ca	r The asser	r The assen	Robots re	Robots n	ni B	https:/	//gitarc	_hard.c	944	0.384	с.	1	. C.	0.99998579	с.	0.999986	с.	0.999986	с.	0.9999857
What is e	s electricity	skittles	ethics	lava	Α	https:/	//gitob	qa.csv	2764	0.367	??	C	MadonnaAnswer Roi	n 0	Able stack	0	negligenc	0	C. Ethics	
As vehicle	e increases	stops	decrease	s stay the	saC	https:/	//gitob	qa.csv	3511	0.343	с	1	с	1	[ings pass	0	keep C.	0	с.	
Which is t	t Revise th	Describe	Test the p	dentify	рВ	https:/	//gitarc	_easy.c	2202	0.336	D	0.253802474	D	0.253802474	в.	0.99498	В	0.810478	в.	0.99498040
Which of	t Partitioni	K-means	Grid base	All of the	e D	https:/	//gitma	chine_	56	0.335	D	1	D	1	D	1	D	1	D	
Which of	t Escalator	Conveyor	Highway	Railroad	в	https:/	//git mis	scellan	285	0.335	Solution: B	C	В.	1	track	0	AngleClos	0	Crossing	
What is a	(Work flow	Work sche	Work rate	Work ou	tr B	https:/	//gitma	nagem	6	0.334	В	0.999997692	2 A.	0.000253997	testing kn	0	В	0.999998	в	0.99999769
Mechanic	a wind turb	solar pane	shad tree	giant ice	r D	https:/	//gitob	qa.csv	2946	0.334	Current	C	?????	0	ads	0	А.	0.710499	shelf	
Which of	t R chart	Run chart	X-bar cha	r Pareto cl	ha B	https:/	//git mis	scellan	561	0.329	с	0.999993719	C.	0.99789524	с	0	с	0.999994	с	0.99999371
1																				
													<u> </u>				ι.			

Methodology Directly Request in Prompt (2/2)

Roal and Goal



Correctness and Confidence "You are an expert in the field of automation. Your goal is to answer the single-choice question as well as output your level of confidence.

Confidence Definition •

0: You have almost no confidence — this is essentially a guess, and you	have
1: Your confidence is low — you have some basis for your choice, but ove	rall
2: Your confidence is moderate - you are relatively sure, but there's st	i11
3: You are quite confident — you are fairly certain your choice is corre	ct,
4: You are very confident - you are highly certain that your choice is t	he r
5: You are absolutely confident — you are 100% certain that your choice	is c

• Example

Example: Input: // text input Output: Answer: ..., confidence: ...

Dataset Collection

Call Models to Generate Answers

Calculate Confidence

of Generated Answers

Calculate Errors between the

Ratio

Visualization

(2/2

Methodology Directly Request in Prompt (2/2)

• Generate 1 response for each question



Result Table:

										\frown		
Question	Option A	Option B	Option C	Option D	TrueAnsw	source	file	index	automatio	Answer	Confidence	
The whee	potential	efficiency	output	friction	D	https://gi	arc_easy.c	673	0.398	D	1	
Which of t	Software	F-10 picku	Oil refine	Parking lo	D	https://gi	high_scho	49	0.397	в	0.8	
A clean ai	web desig	hard drive	math class	computer	В	https://gi	obqa.csv	1466	0.385	D	1	
Robots ca	The assen	The assen	Robots re	Robots m	в	https://gi	arc_hard.o	944	0.384	с	1	
What is es	electricity	skittles	ethics	lava	Α	https://gi	obqa.csv	2764	0.367	E	1	
The items	each emp	concurren	empirical	content v	D	https://gi	professio	121	0.345	D	1	
As vehicle	increases	stops	decreases	stay the s	C	https://gi	obqa.csv	3511	0.343	с	0.8	
Which of t	machine-	textile wo	car manuf	chemical	C	https://gi	sociology	154	0.337	А	1	
Which is t	Revise the	Describe t	Test the p	Identify p	В	https://gi	arc_easy.o	2202	0.336	в	1	
)

Methodology Error Calculation between the confidence ratio and the correctness

3 Metrics

• Mean Square Error $MSE = E_q [P_M - I(a_M)]^2$



Mean Absolute Deviation

$$mol = L_q[I_M - I(u_M)]$$

 $MAD = E_q |P_M - I(a_M)|$

• Root Mean Square $RMS = \sqrt{E_q [P_M - I(a_M)]^2}$

- Motivation
- Research Question
- Related Work
- Research Work Pipeline
- Methodology
- Results and Findings
- Conclusion and Future Work

Confidence Distribution (1/2) of Davinci-003 using logprob







This model is over-confident!



Temperature can reduce confidence.



- This model is not capable of solving questions from the dataset.
- Models with low capacity tends to always be medium-confident.

Confidence Distribution (1/2) of LLaMA2-70b using Prompt







Visualized Errors (2/2) of Davinci-003 (Logprob)





There exists an optimum temperature between 0 to 1.

Visualized Errors (2/2) of LLaMA2-70b (Prompt)





Temperature has a large impact on the errors.

- Motivation
- Research Question
- Related Work
- Research Work Pipeline
- Methodology
- Results and Findings
- Conclusion and Future Work

Conclusion and Future Work

- All 3 large language models have the problem of being over-confident.
- Temperature can have a great impact on both the confidence and correctness and there could exist a temperature, at which the LLM perform the best.

- Investigate the optimum temperature for the LLMs.
- Finetune the LLMs to reduce the error between the correctness and confidence to get better performance.

Literature

[1] S. Lin, J. Hilton, and O. Evans, "Teaching Models to Express Their Uncertainty in Words," arXiv.org, Jun. 13, 2022. https://arxiv.org/abs/2205.14334 (accessed Jun. 27, 2023).

[2] Z. Lin, S. Trivedi, and J. Sun, "Generating with Confidence: Uncertainty Quantification for Black-box Large Language Models," arXiv.org, May 30, 2023. https://arxiv.org/abs/2305.19187 (accessed Jun. 27, 2023).

[3] L. Kuhn, Y. Gal, and S. Farquhar, "Semantic Uncertainty: Linguistic Invariances for Uncertainty Estimation in Natural Language Generation," arXiv.org, Apr. 15, 2023. https://arxiv.org/abs/2302.09664 (accessed Jun. 27, 2023).

[4] A. Azaria and T. Mitchell, "The Internal State of an LLM Knows When its Lying," arXiv.org, Apr. 25, 2023. https://arxiv.org/abs/2304.13734 (accessed Jun. 27, 2023).



University of Stuttgart Institut of Industrial Automation and Software Engineering

Thank you!



Ziyao Zhou

e-mail Ziyao.zhou2021@gmail.com phone +49 (0) 711 685fax +49 (0) 711 685-

University of Stuttgart Institut of Indutrial Automation and Software Engineering

