

LLM Task Interference: An Initial Study on the Impact of Task-Switch in Conversational History

Anonymous ACL submission

Abstract

With the recent emergence of powerful instruction-tuned large language models (LLMs), various helpful conversational Artificial Intelligence (AI) systems have been deployed across many applications. When prompted by users, these AI systems successfully perform a wide range of tasks as part of a conversation. To provide some sort of memory and context, such approaches typically condition their output on the entire conversational history. Although this sensitivity to the conversational history can often lead to improved performance on subsequent tasks, we find that performance can in fact also be negatively impacted, if there is a *task-switch*. To the best of our knowledge, our work makes the first attempt to formalize the study of such vulnerabilities and interference of tasks in conversational LLMs caused by task-switches in the conversational history. Our experiments across 5 datasets with 15 task switches using popular LLMs reveal that many of the task-switches can lead to significant performance degradation.¹

1 Introduction

Recent advancements in Natural Language Processing (NLP) (Brown et al., 2020; OpenAI, 2023), have led to their widespread deployment of large language models (LLMs) across various applications (Bubeck et al., 2023; Anil et al., 2023; Singhal et al., 2022). One of the popular NLP tasks includes conversational systems where LLMs are capable of engaging in dialogues that mimic human interactions (Manyika and Hsiao, 2023; Bai et al., 2022). A typical interaction involves a series of conversation turns starting with the user and the LLM responds to the user. This interaction is however focused on a specific topic or a

¹Code available on acceptance.

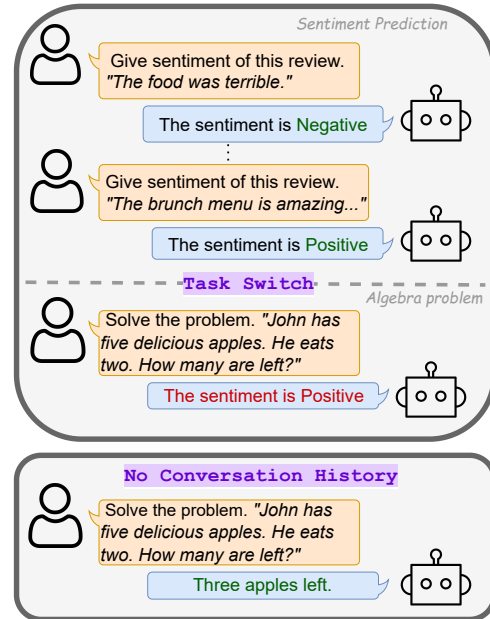


Figure 1: An illustrative example where the chat history is based on sentiment prediction. Algebra word problem introduces *task-switch* which results in an incorrect prediction.

task (Hosseini-Asl et al., 2020; Lee et al., 2022).

The performance of LLMs is further boosted by leveraging in-context examples or few-shot examples of a particular task (Brown et al., 2020; Smith et al., 2022; Thoppilan et al., 2022). In-context learning, by utilizing examples within the conversation history, enables LLMs to generate responses that are relevant and tailored to the contextual conversation. The auto-regressive nature of popular instruction-tuned (LLMs) suggests that the LLM generated response is conditioned on the entire conversation history. This underscores the sequential dependency and contextual awareness embedded within these models. While prompt sensitivity has been exploited by in-context learning to improve downstream performance, this sensitivity has also opened

the door to vulnerabilities, where malicious actors can exploit prompt sensitivity for adverse purposes (Greshake et al., 2023; Liu et al., 2023; Jiang et al., 2023b; Xu et al., 2023).

In this paper, we investigate the sensitivity and the impact of LLM performance on past conversational interaction. To do so, we introduce the concept of *task-switch*. A task-switch is characterized by a conversational objective, moving from one distinct task to another within the same conversation thread, for example: Figure 1 illustrates a task-switch from sentiment prediction to math algebra which confuses the model to output erroneously. Designing LLMs that can seamlessly switch between tasks without degradation in performance can influence the reliability of LLMs in realistic scenarios.

In this work, we systematically study the impact of predictive performance and the sensitivity of LLMs in the presence of different task-based chat histories. Our key contributions and takeaways can be summarised as:

- We formalize the risk of performance degradation of LLMs due to task-switch.
- We present the impact of task-switch on diverse datasets with more than 15 different task-switches.
- We measure the task-switch sensitivity for popular LLMs of different sizes, where we observe that in some cases very large (175B) and small (7B) LLMs can both be susceptible to performance degradation from task-switch.

2 Related Work

Large Language Models (LLMs) are becoming a crucial building block of conversation-based virtual assistants (OpenAI, 2023; Touvron et al., 2023; Jiang et al., 2023a; Anil et al., 2023). Leveraging in-context or few-shot examples, LLMs have demonstrated remarkable capabilities for downstream tasks (Brown et al., 2020). In contrast to the resource-intensive fine-tuning process (Gao et al., 2020), in-context learning eliminates the need for parameter updates, while achieving state-of-the-art performance (Rae et al., 2021; Smith et al., 2022; Thoppilan et al., 2022; Von Oswald et al., 2023; Chan et al., 2022; Akyürek et al., 2022; Hahn and Goyal, 2023). However, despite its advantages, in-context learning tends to suffer

from sensitivity to prompts, input distribution, and formats, which can potentially impact the model’s performance (Liu et al., 2021; Zhao et al., 2021; Lu et al., 2021; Min et al., 2022; Liu and Wang, 2023; Chang and Jia, 2023). Chang and Jia (2023) observe that the in-context examples implicitly bias the model. In our work, we aim to study the bias that arises due to chat history (in-context examples) when a user switches the task. Furthermore, recent works (Liu et al., 2023; Greshake et al., 2023) have looked at the vulnerability of LLM to prompt injections and adversarial attacks. Unlike prompt injection, where a malicious prompt may be added to the conversation of LLM, our setting, is concerned with non-malicious task-switches. While a few recent works have investigated the reliance on shortcuts in conversation history (Tang et al., 2023; Si et al., 2022; Weston and Sukhbaatar, 2023), our work aims to evaluate prompt history sensitivity for a new task. Our work is also differentiated from the study topic change in Task-oriented Dialogue systems (Xie et al., 2021; Xu et al., 2021; Yang et al., 2022) as we consider a stronger shift of task-switch from open dialogue LLMs.

3 Conversational Task-Switch

This work introduces and formalizes *task-switch* in a conversation for LLMs. A conversation between a user and the LLM consists of multiple conversation turns. Now consider (u_k, r_k) as the k -th turn of the conversation where u_k corresponds to the k -th user prompt and the model’s corresponding response r_k . Each user prompt u_k can be viewed as an instance of a specific task request, e.g. *sentiment classification* or *mathematical reasoning*. A conversation history of L turns can be defined as $\mathbf{h} = \{(u_k, r_k)\}_{k=1}^L$. Subsequently, the next response, r_{L+1} for model θ is given as:

$$r_{L+1} = \arg \max_r P_\theta(r|u_{L+1}, \mathbf{h}) \quad (1)$$

In this work, we consider conversations with a single task-switch, where all user requests in the conversation history \mathbf{h} belong to the same task and the final user request u_{L+1} is a different task. We refer to the task associated with \mathbf{h} as the conversation history task (*CH task*) T_h where $\mathbf{h} \in T_h$ and the switched task

associated with the final user request u_{L+1} as the *target task* T_t where $u_{L+1} \in T_t$.

When the tasks T_h and T_t are sufficiently different (as per human understanding of language and tasks), the conversation history \mathbf{h} ideally must not impact the response, r_{L+1} . For a model robust to such a task-switches, $T_h \rightarrow T_t$, its response r_{L+1} is conditionally independent of the conversation history,

$$r_{L+1} \perp \mathbf{h} | u_{L+1} \quad \mathbf{h} \in T_h, u_{L+1} \in T_t. \quad (2)$$

However, in practice, models can be sensitive to the conversation history, \mathbf{h} , which can harm the quality of the response r_{L+1} after a task-switch, $T_h \rightarrow T_t$. We define $\tau(\cdot)$, the *task-switch sensitivity* of a model θ , to measure the extent of this vulnerability.²

$$\tau(T_h, T_t; \theta) = \mathbb{E}_{u_{L+1} \in T_t, \mathbf{h} \in T_h} [\log \rho] \quad (3)$$

$$\rho = \frac{P_\theta(r^* | u_{L+1})}{P_\theta(r^* | u_{L+1}, \mathbf{h})} \quad (4)$$

$$r^* = \arg \max_r P_\theta(r | u_{L+1}). \quad (5)$$

Task-switch sensitivity can be interpreted as:

1. $\tau(\cdot) > 0$: The model is impacted by the task-switch in the conversation history and is less confident in zero-shot prediction.
2. $\tau(\cdot) = 0$: The task-switch has no impact on the model’s zero-shot prediction, suggesting a level of task-switch robustness.
3. $\tau(\cdot) < 0$: The task-switch gives the model more confidence in its zero-shot prediction.

To simulate a setting where the model has perfect performance on the CH-task, T_h we adopt teacher-forcing, s.t. $\mathbf{h} = \{(u_k, \hat{r}_k)\}_{k=1}^L$, where \hat{r} is the reference ground-truth response.

4 Experiments

4.1 Experimental Setup

Data. We evaluate five different datasets covering a range of tasks: Gigaword (Graff et al., 2003); abstract algebra subset of Measuring Massive Multitask Language Understanding (MMLU; Hendrycks et al. (2021)), named MMLU AA; TweetQA (Xiong et al., 2019); Rotten Tomatoes (RT; Pang and Lee (2005)); and human-aging subset from the MMLU dataset (MMLU HA) in the Appendix.

²Theoretical and empirical implications of other definitions for task-switch sensitivity in Appendix E

Data	Task
Gigaword	Summarization
MMLU AA	Math Multiple Choice Question
TweetQA	Social Question Answer
RT	Sentiment classification
MMLU HA	Social Multiple Choice Question

Table 1: Datasets Summary.

Models. We explore the task-switch sensitivity of four popular models. We consider two open-source small models, Llama-7b-chat (Touvron et al., 2023) and Mistral-7b-chat (Jiang et al., 2023a); and two larger closed models, GPT-3.5 (Brown et al., 2020) and GPT-4 (OpenAI, 2023). Zero-shot, absolute model performances are presented in Appendix B.

4.2 Results

In addition to the task-switch sensitivity $\tau(\cdot)$, we assess performance changes between the predictions in the presence of history and task-switch vs zero-shot. Table 2 and Table 3 show-cases the impact of conversational task-switch with MMLU AA and Rotten Tomatoes as the target tasks, T_t respectively³. As would be expected with *in-context examples*, the performance change in accuracy is generally positive. The negative trend for change in accuracy from $T_h \rightarrow T_t$, suggests that the task-switch causes performance degradation. For example, in the Gigaword summarization task as T_h and MMLU AA as T_t , most models (GPT-3.5, Llama-7B and Mistral-7B) see a performance drop. Interestingly, for some models, the task-switch may increase performance; most prominently for Mistral-7B with Rotten Tomatoes as T_h and MMLU AA as T_t .

The sensitivity of different models to different task-switches can be compared fairly using the task-switch metric, $\tau(\cdot)$. The larger the value of $\tau(\cdot)$, the greater a model’s sensitivity to a specific task-switch. In Table 2 and Table 3, Llama-7B usually has the highest sensitivity to task-switches with for example $\tau = 3.37$ for a switch from MMLU AA to Rotten Tomatoes and $\tau = 9.91$ for task-switch from Rotten Tomatoes to MMLU AA. We observe a general trend between the change in accuracy and $\tau(\cdot)$ for task-switch scenarios for $T_t =$ Rotten Tomatoes where a negative change in performance

³The impact of task-switch for other datasets as the target tasks is given in Appendix C.1

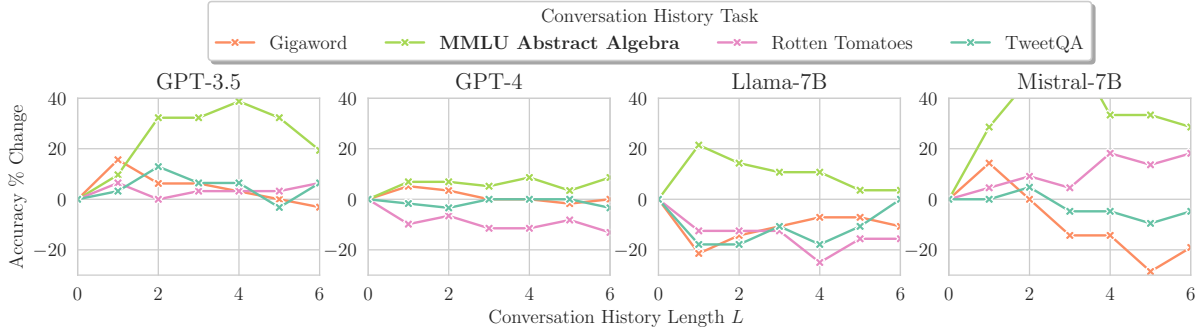


Figure 2: Target Task: MMLU Abstract Algebra. % change in accuracy relative to zero-shot performance.

CH-Task	Model	% Change	$\tau(\cdot)$
MMLU AA	GPT-3.5	17.17	*
	GPT-4	-1.09	*
	Llama-7B	0.00	31.51
	Mistral-7B	37.68	1.12
Gigaword	GPT-3.5	-12.12	*
	GPT-4	-8.74	*
	Llama-7B	-18.75	5.23
	Mistral-7B	-21.74	3.13
Rotten Tomatoes	GPT-3.5	2.02	*
	GPT-4	-8.20	*
	Llama-7B	-12.50	9.91
	Mistral-7B	11.59	0.83
TweetQA	GPT-3.5	-19.19	*
	GPT-4	-8.20	*
	Llama-7B	-12.50	6.37
	Mistral-7B	-7.25	2.78

Table 2: Task-switch impact from CH-tasks (T_h) to target (T_t): **MMLU AA** and conversation length $L = 6$. Sensitivity not calculable for *.

CH-Task	Model	% Change	$\tau(\cdot)$
Rotten Tomatoes	GPT-3.5	3.00	*
	GPT-4	1.74	*
	Llama-7B	2.54	4.02
	Mistral-7B	3.17	2.65
Gigaword	GPT-3.5	0.11	*
	GPT-4	-0.98	*
	Llama-7B	1.82	1.98
	Mistral-7B	-0.79	3.04
MMLU AA	GPT-3.5	-0.22	*
	GPT-4	0.76	*
	Llama-7B	-5.33	3.37
	Mistral-7B	1.33	1.39
TweetQA	GPT-3.5	-0.33	*
	GPT-4	-0.98	*
	Llama-7B	2.72	2.77
	Mistral-7B	-1.23	3.01

Table 3: Task-switch impact from CH-tasks (T_h) to target (T_t): **Rotten Tomatoes** and conversation length $L = 6$. Sensitivity not calculable for *.

also suggests very high task-switch sensitivity. In Figure 2, we plot the change in performance with increasing T_h examples for MMLU AA dataset. Here we can clearly observe that in-context examples improve the predictive performance. Notably, the accuracy variation is more pronounced in smaller 7B models, likely due to their lower baseline performance, which is substantially improved by in-context learning. Performance fluctuations for conversation history, \mathbf{h} , can stem from two primary factors: a significant drop in the predicted probability for the zero-shot response, r^* , or a notable increase in the probability for an alternative response, r . The latter can result in substantial performance change while maintaining low sensitivity, $\tau(\cdot)$. By analyzing both performance changes and task-switch sensitivity, we gain deeper insights into the models’ adaptability to task-switches and the underlying dynamics influencing these shifts.

5 Conclusions and Future Work

This work formalizes and performs an initial investigation into the sensitivity of large language models (LLMs) to task-switch scenarios within conversational contexts. We introduce a task-sensitivity metric that can explain a model’s behavior to task-switches along with the performance change. By experimenting with various task-switch settings, we observe that even advanced models like GPT-4 can exhibit vulnerabilities to task-switches. Our work additionally lays the foundation for future work on ‘side-channel’ vulnerabilities of LLMs to undesired information leakage/bias from the conversation history. Further work will focus on developing adaptive context management strategies within LLMs to mitigate the risk of task-switch sensitivity.

6 Limitations

Although both GPT-3.5 and GPT-4 show degradation in performance, given the closed nature of OpenAI models, we were not able to perform task sensitivity analysis. We were additionally limited by the maximum token length, hence analysis over extremely long conversations was not feasible. Future work could also look into alignment between humans and the model as a metric which was out of the scope for this paper.

7 Ethics and Risks

All of the datasets used are publicly available. Our implementation utilizes the PyTorch 1.12 framework, an open-source library. We obtained a license from Meta to employ the Llama-7B model via HuggingFace. Additionally, our research is conducted per the licensing agreements of the Mistral-7B, GPT-3.5, and GPT-4 models. We ran our experiments on A100 Nvidia GPU and via OpenAI API.

Our work may be built upon to identify vulnerabilities of LLMs. Overall, there are no ethical concerns with this work.

References

- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. 2022. What learning algorithm is in-context learning? investigations with linear models. *arXiv*. 308–311
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*. 312–316
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv*. 317–321
- Satanjeev Banerjee and Alon Lavie. 2005. **ME-TEOR: An automatic metric for MT evaluation with improved correlation with human judgments**. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics. 322–329
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901. 330–335
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv*. 336–341
- Stephanie CY Chan, Ishita Dasgupta, Junkyung Kim, Dharshan Kumaran, Andrew K Lampinen, and Felix Hill. 2022. Transformers generalize differently from information stored in context vs in weights. *arXiv preprint arXiv:2210.05675*. 342–346
- Ting-Yun Chang and Robin Jia. 2023. Data curation alone can stabilize in-context learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8123–8144. 347–351
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. **A framework for few-shot language model evaluation**. 352–361

362	Tianyu Gao, Adam Fisch, and Danqi Chen. 2020.	Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang,	416
363	Making pre-trained language models better few-	Tianwei Zhang, Yepang Liu, Haoyu Wang, Yan	417
364	shot learners. <i>arXiv</i> .	Zheng, and Yang Liu. 2023. Prompt injection	418
		attack against llm-integrated applications. <i>arXiv</i>	419
365	David Graff, Junbo Kong, Ke Chen, and Kazuaki	<i>preprint arXiv:2306.05499</i> .	420
366	Maeda. 2003. English gigaword. <i>Linguistic Data</i>		
367	<i>Consortium, Philadelphia</i> , 4(1):34.		
368	Kai Greshake, Sahar Abdelnabi, Shailesh Mishra,	Yao Lu, Max Bartolo, Alastair Moore, Sebastian	421
369	Christoph Endres, Thorsten Holz, and Mario	Riedel, and Pontus Stenetorp. 2021. Fantasti-	422
370	Fritz. 2023. Not what you’ve signed up for: Com-	cally ordered prompts and where to find them:	423
371	promising real-world llm-integrated applications	Overcoming few-shot prompt order sensitivity.	424
372	with indirect prompt injection. In <i>Proceedings of</i>	<i>arXiv preprint arXiv:2104.08786</i> .	425
373	<i>the 16th ACM Workshop on Artificial Intelligence</i>		
374	<i>and Security</i> , pages 79–90.		
375	Michael Hahn and Navin Goyal. 2023. A theory of	James Manyika and Sissie Hsiao. 2023. An overview	426
376	emergent in-context learning as implicit structure	of bard: an early experiment with generative ai.	427
377	induction. <i>arXiv</i> .	<i>AI. Google Static Documents</i> , 2.	428
378	Dan Hendrycks, Collin Burns, Steven Basart, Andy	Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel	429
379	Zou, Mantas Mazeika, Dawn Song, and Jacob	Artetxe, Mike Lewis, Hannaneh Hajishirzi, and	430
380	Steinhardt. 2021. Measuring massive multitask	Luke Zettlemoyer. 2022. Rethinking the role of	431
381	language understanding. <i>Proceedings of the In-</i>	demonstrations: What makes in-context learning	432
382	<i>ternational Conference on Learning Representa-</i>	work? <i>arXiv preprint arXiv:2202.12837</i> .	433
383	<i>tions (ICLR)</i> .		
384	Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng	R OpenAI. 2023. Gpt-4 technical report. arxiv	434
385	Wu, Semih Yavuz, and Richard Socher. 2020. A	2303.08774. <i>View in Article</i> , 2:13.	435
386	simple language model for task-oriented dialogue.		
387	<i>Advances in Neural Information Processing Sys-</i>	Bo Pang and Lillian Lee. 2005. Seeing stars: Ex-	436
388	<i>tems</i> , 33:20179–20191.	ploiting class relationships for sentiment catego-	437
389	Albert Q Jiang, Alexandre Sablayrolles, Arthur	rization with respect to rating scales. In <i>Proceed-</i>	438
390	Mensch, Chris Bamford, Devendra Singh Chap-	<i>ings of the ACL</i> .	439
391	lot, Diego de las Casas, Florian Bressand, Gianna	Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie	440
392	Lengyel, Guillaume Lample, Lucile Saulnier, et al.	Millican, Jordan Hoffmann, Francis Song, John	441
393	2023a. Mistral 7b. <i>arXiv</i> .	Aslanides, Sarah Henderson, Roman Ring, Su-	442
394	Shuyu Jiang, Xingshu Chen, and Rui Tang. 2023b.	sannah Young, et al. 2021. Scaling language	443
395	Prompt packer: Deceiving llms through compo-	models: Methods, analysis & insights from train-	444
396	sitional instruction with hidden attacks. <i>arXiv</i>	ing gopher. <i>arXiv</i> .	445
397	<i>preprint arXiv:2310.10077</i> .		
398	Harrison Lee, Raghav Gupta, Abhinav Rastogi,	Chenglei Si, Dan Friedman, Nitish Joshi, Shi Feng,	446
399	Yuan Cao, Bin Zhang, and Yonghui Wu. 2022.	Danqi Chen, and He He. 2022. What spurious	447
400	Sgd-x: A benchmark for robust generalization	features can pretrained language models combat?	448
401	in schema-guided dialogue systems. In <i>AAAI</i> ,		
402	volume 36, pages 10938–10946.	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara	449
403	Chin-Yew Lin. 2004. ROUGE: A package for auto-	Mahdavi, Jason Wei, Hyung Won Chung, Nathan	450
404	matic evaluation of summaries . In <i>Text Summa-</i>	Scales, Ajay Tanwani, Heather Cole-Lewis,	451
405	<i>rization Branches Out</i> , pages 74–81, Barcelona,	Stephen Pfohl, et al. 2022. Large language mod-	452
406	Spain. Association for Computational Linguis-	els encode clinical knowledge. <i>arXiv</i> .	453
407	tics.		
408	Hongfu Liu and Ye Wang. 2023. Towards infor-	Shaden Smith, Mostofa Patwary, Brandon Norick,	454
409	mative few-shot prompt with maximum informa-	Patrick LeGresley, Samyam Rajbhandari, Jared	455
410	tion gain for in-context learning. <i>arXiv preprint</i>	Casper, Zhun Liu, Shrimai Prabhumoye, George	456
411	<i>arXiv:2310.08923</i> .	Zerveas, Vijay Korthikanti, et al. 2022. Us-	457
412	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill	ing deepspeed and megatron to train megatron-	458
413	Dolan, Lawrence Carin, and Weizhu Chen. 2021.	turing nlg 530b, a large-scale generative language	459
414	What makes good in-context examples for gpt-3?	model. <i>arXiv preprint arXiv:2201.11990</i> .	460
415	<i>arXiv preprint arXiv:2101.06804</i> .		
		Ruixiang Tang, Dehan Kong, Longtao Huang, and	461
		Hui Xue. 2023. Large language models can be	462
		lazy learners: Analyze shortcuts in in-context	463
		learning. <i>ACL Findings</i> .	464
		Gemini Team, Rohan Anil, and et al. 2024. Gemini:	465
		A family of highly capable multimodal models .	466
		Romal Thoppilan, Daniel De Freitas, Jamie Hall,	467
		Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze	468
		Cheng, Alicia Jin, Taylor Bos, Leslie Baker,	469

470 Yu Du, et al. 2022. Lamda: Language mod-
471 els for dialog applications. *arXiv preprint*
472 *arXiv:2201.08239*.

473 Hugo Touvron, Thibaut Lavril, Gautier Izacard,
474 Xavier Martinet, Marie-Anne Lachaux, Timothée
475 Lacroix, Baptiste Rozière, Naman Goyal, Eric
476 Hambro, Faisal Azhar, et al. 2023. Llama: Open
477 and efficient foundation language models. *arXiv*.

478 Johannes Von Oswald, Eyvind Niklasson, Ettore
479 Randazzo, João Sacramento, Alexander Mordv-
480 intsev, Andrey Zhmoginov, and Max Vladymy-
481 rov. 2023. Transformers learn in-context by gra-
482 dient descent. In *International Conference on*
483 *Machine Learning*, pages 35151–35174. PMLR.

484 Jason Weston and Sainbayar Sukhbaatar. 2023. Sys-
485 tem 2 attention (is something you might need
486 too). *arXiv*.

487 Huiyuan Xie, Zhenghao Liu, Chenyan Xiong,
488 Zhiyuan Liu, and Ann Copestake. 2021. Tiage:
489 A benchmark for topic-shift aware dialog model-
490 ing. *ACL*.

491 Wenhan Xiong, Jiawei Wu, Hong Wang, Vivek
492 Kulkarni, Mo Yu, Xiaoxiao Guo, Shiyu Chang,
493 and William Yang Wang. 2019. Tweetqa: A
494 social media focused question answering dataset.
495 In *Proceedings of the 57th Annual Meeting of the*
496 *Association for Computational Linguistics*.

497 Xilie Xu, Keyi Kong, Ning Liu, Lizhen Cui,
498 Di Wang, Jingfeng Zhang, and Mohan Kankan-
499 halli. 2023. An llm can fool itself: A
500 prompt-based adversarial attack. *arXiv preprint*
501 *arXiv:2310.13345*.

502 Yi Xu, Hai Zhao, and Zhuosheng Zhang. 2021.
503 Topic-aware multi-turn dialogue modeling. In
504 *AAAI*, volume 35, pages 14176–14184.

505 Chenxu Yang, Zheng Lin, Jiangnan Li, Fandong
506 Meng, Weiping Wang, Lanrui Wang, and Jie
507 Zhou. 2022. Take: topic-shift aware knowledge
508 selection for dialogue generation. In *Proceedings*
509 *of the 29th International Conference on Compu-*
510 *tational Linguistics*, pages 253–265.

511 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and
512 Sameer Singh. 2021. Calibrate before use: Im-
513 proving few-shot performance of language models.
514 In *International Conference on Machine Learn-*
515 *ing*, pages 12697–12706. PMLR.

Appendix

Appendix A gives more details about the datasets, Appendix B reports the zero-shot absolute performance of all models on all tasks, Appendix C presents an ablation study on the conversation history length (with multiple seeds), Appendix D discusses the prompt templates, Appendix E discusses other definitions for task-switch sensitivity, Appendix F discusses correlations, and Appendix G tabulates confusion matrices for each model.

A Datasets and Metrics Summary

Data	#Train	#Test	Task
MMLU HA	26	222	Social MCQ
MMLU AA	14	99	Math MCQ
RT	8.53k	1.07k	Sentiment class
Gigaword	3.8M	1.95k	Summarization
TweetQA	4.54k	583	Social QA

Table 4: Dataset Summary. QA: Question-Answering. MCQ: Multiple Choice Question

In Section 4.2 of the main paper, we present results evaluated on two different datasets: MMLU Abstract Algebra (MMLU AA) multiple choice questions and Rotten Tomatoes (RT) sentiment classification. In Appendix B, C, we present results evaluated on all of the datasets covering a range of tasks: MMLU Human Aging (MMLU HA) multiple choice questions, Gigaword for summarization, and TweetQA question-answering. The train-test splits of these datasets are shown in Table 4. The train set is randomly sampled to form prompts to produce a conversation history \mathbf{h} of L turns, and the test set is used to evaluate model performance on the $(L + 1)$ -th turn. The prompt templates used for each dataset are discussed in Appendix D.

For classification tasks performance is measured using accuracy, whilst for generative tasks it is measured using ROUGE (Lin, 2004) or METEOR (Banerjee and Lavie, 2005).

B Absolute Performance

When evaluating the target task with a conversation history, it is useful to compare the performance against a baseline with no conversation history ($\mathbf{h} = \emptyset, L = 0$). This is equivalent to evaluating in a zero-shot setting. This section

reports the zero-shot performance for all the target task (T_t) datasets: MMLU HA in Table 5, MMLU AA in Table 6, RT in Table 7, Gigaword in Table 8 and TweetQA in Table 9. Also note that for the classification tasks (MMLU HA, MMLU AA, RT), we also report the number of responses for which we were unable to extract the answer (# Format Errors), which is further discussed in Appendix D. We evaluate on the test set with four LLMs (GPT-3.5, GPT-4, Mistral-7B, Llama-7B), which were all set to Temperature 0 for reproducibility.

Model	Accuracy	# Format Errors
GPT-3.5	66.22	18
GPT-4	84.68	0
Llama-7B	45.50	12
Mistral-7B	55.41	0

Table 5: Zero-shot performance on **MMLU HA**.

Model	Accuracy	# Format Errors
GPT-3.5	31.31	7
GPT-4	58.59	0
Llama-7B	28.28	3
Mistral-7B	21.21	0

Table 6: Zero-shot performance on **MMLU AA**.

Model	Accuracy	# Format Errors
GPT-3.5	89.90	0
GPT-4	91.80	4
Llama-7B	87.43	1
Mistral-7B	86.68	1

Table 7: Zero-shot performance on **RT**.

Model	ROUGE-1	ROUGE-2	ROUGE-L
GPT-3.5	17.37	4.79	14.78
GPT-4	15.76	4.07	13.34
Llama-7B	11.61	3.13	9.90
Mistral-7B	18.60	5.19	15.84

Table 8: Zero-shot performance on **Gigaword**.

Model	ROUGE-1	ROUGE-L	METEOR
GPT-3.5	30.66	30.39	44.18
GPT-4	28.03	27.68	43.41
Llama-7B	17.91	17.67	33.84
Mistral-7B	25.35	25.01	40.71

Table 9: Zero-shot performance on **TweetQA**.

C Conversation History Length Ablation

This section presents an ablation study on the performance change after a task-switch for varying conversation history lengths. For each dataset in Table 4 we select four datasets (including itself), from which we use the training set as conversation history. The details of the prompt structure are presented in Appendix D.

C.1 Task-switch Performance Change

We compare the percentage change in metrics relative to zero-shot performance ($\mathbf{h} = \emptyset$, i.e. no conversation history) as a function of conversation history length L and for different LLMs. Results are plot in Figures 3, 4, 5, 6, 7 for MMLU HA, MMLU AA, RT, Gigaword and TweetQA respectively. When there is *not* a task switch, we would expect a performance increase (assuming the training examples are representative of the test set). As per our discussion in Section 4.2, we observe that different models degrade on different task-switches and this is not limited by the model size.

C.2 Format Failure Rate

Typically, classification tasks (MMLU HA, MMLU AA, RT) are evaluated using logits, however we use a generative approach for consistency: we are evaluating the model in a conversational setting, and we do not have access to the logits exactly. Thus, we must post-process the model output to determine the class. In this, we try to give the LLM the benefit of the doubt and do our best to extract the class. For example, although the prompt requests the model to output within answer tags like "`<Answer> positive </Answer>`", we also accept "positive", but we do not accept "positive/negative". Due to the imperfect nature of this setup, either we may not detect the correct format, or the model generates erroneous text.

Importantly, models may become more susceptible to these errors when performing a task-switch, causing performance degradation. We capture this by reporting the percentage % change in the number of examples that the model failed on (relative to zero-shot) as the context history length increases. These are plot in Figures 8, 9, 10 for MMLU HA, MMLU

AA and RT respectively. Figures 8 and 9 show that GPT-3.5 and Mistral-7B are susceptible to format errors in task-switches when evaluating on multiple choice questions, whereas Figure 10 shows that GPT-4 and Llama-7B are more susceptible in sentiment classification.

C.3 Performance Variance

Presented experimental results in the main paper are the average across multiple seeds. However, it can be of interest to understand the extent to which the results can vary across multiple runs, as this provides an error bound on the worst-case and best-case scenarios. In this section we present the variance around the mean results for the models Llama-7b and Mistral-7b when evaluated on the target tasks Rotten Tomatoes (Fig 11) and MMLU-AA (Fig 12) with conversation history lengths $L \in \{0, 3, 6\}$.

D Prompt Template

In each conversation turn, the user prompts the model u_k . The prompts are shown in Table 10. We chose these prompts after careful research and experimentation. We began with popular templates and refined them for our purpose.

Additionally, since we do not have access to the logits for all models, we take a generative approach to the classification tasks (MMLU HA, MMLU AA, RT). Since the model may fail to output an answer in the desired format, we post process the text to extract the answer (which we count as a positive result it matches the reference). We report and discuss the effect of format failures further in C.2. Furthermore, we note that the standard evaluation method used in the Open-LLM leaderboard code (available on [GitHub](#)) is to see if the response starts with A,B,C or D(Gao et al., 2023). We modified the prompt to ensure a more consistent output format (across the different models) resulting in fewer mistakes made.

For the classification tasks, we structure the prompt such that we request the model to output their final answer within answer tags. We note that giving an example of how to use the answer tags always helped, however, this can bias the model towards a particular answer. Instead, we found for MMLU to just leave the answer tags empty, whereas for RT to have the all the sentiment classes inside the tags (see Table 10 for further details).

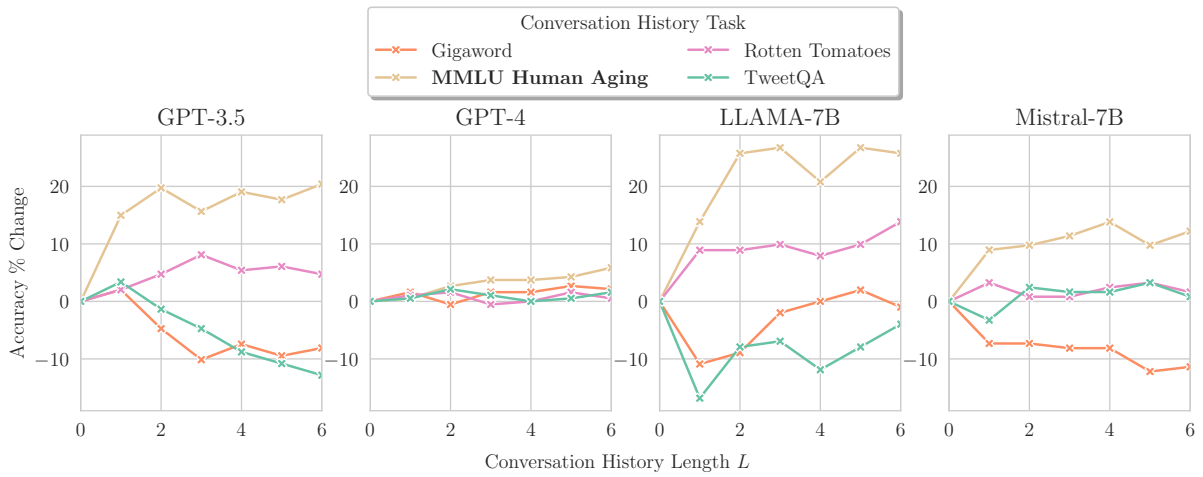


Figure 3: Target Task: MMLU HA. Percentage % change in accuracy relative to zero-shot performance (no conversation history) for increasing conversation history length L and various models.

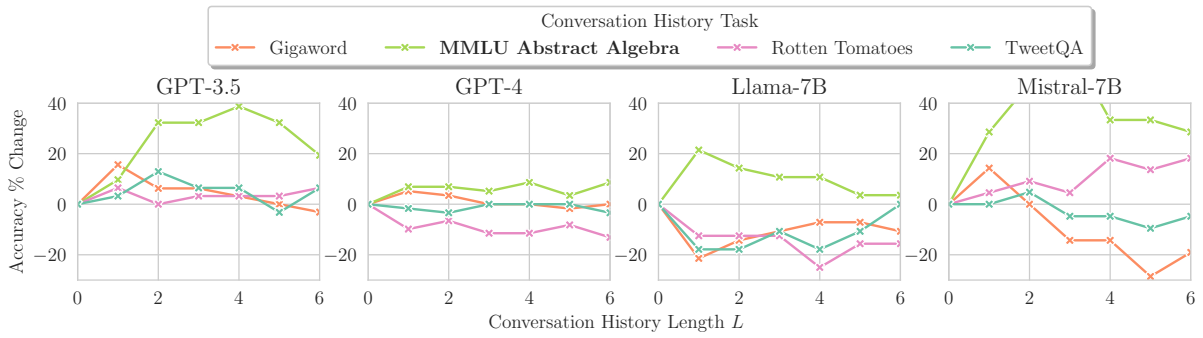


Figure 4: Target Task: MMLU AA. Percentage % change in accuracy relative to zero-shot performance (no conversation history) for increasing conversation history length L and various models.

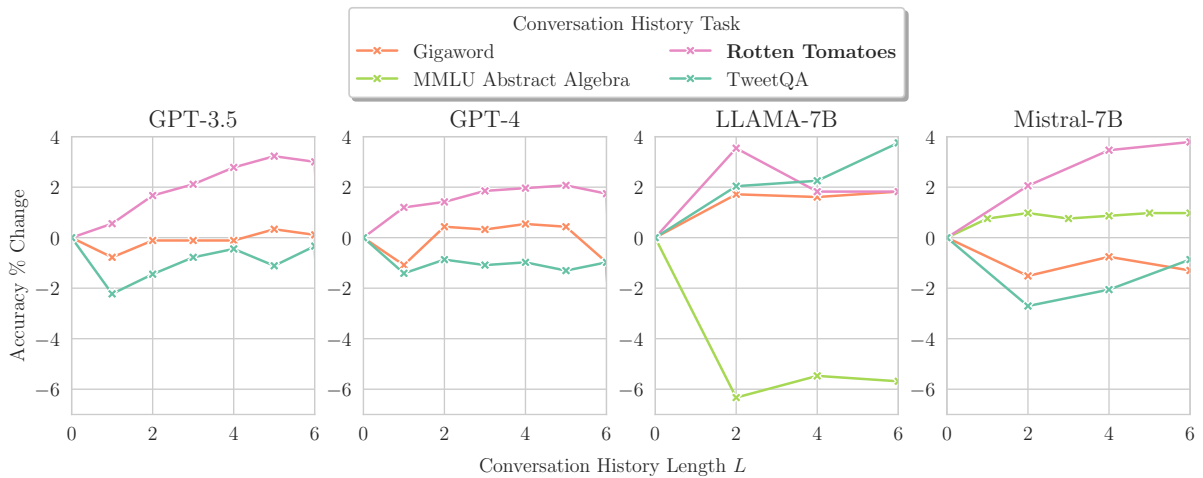


Figure 5: Target Task: RT. Percentage % change in accuracy relative to zero-shot performance (no conversation history) for increasing conversation history length L and various models.

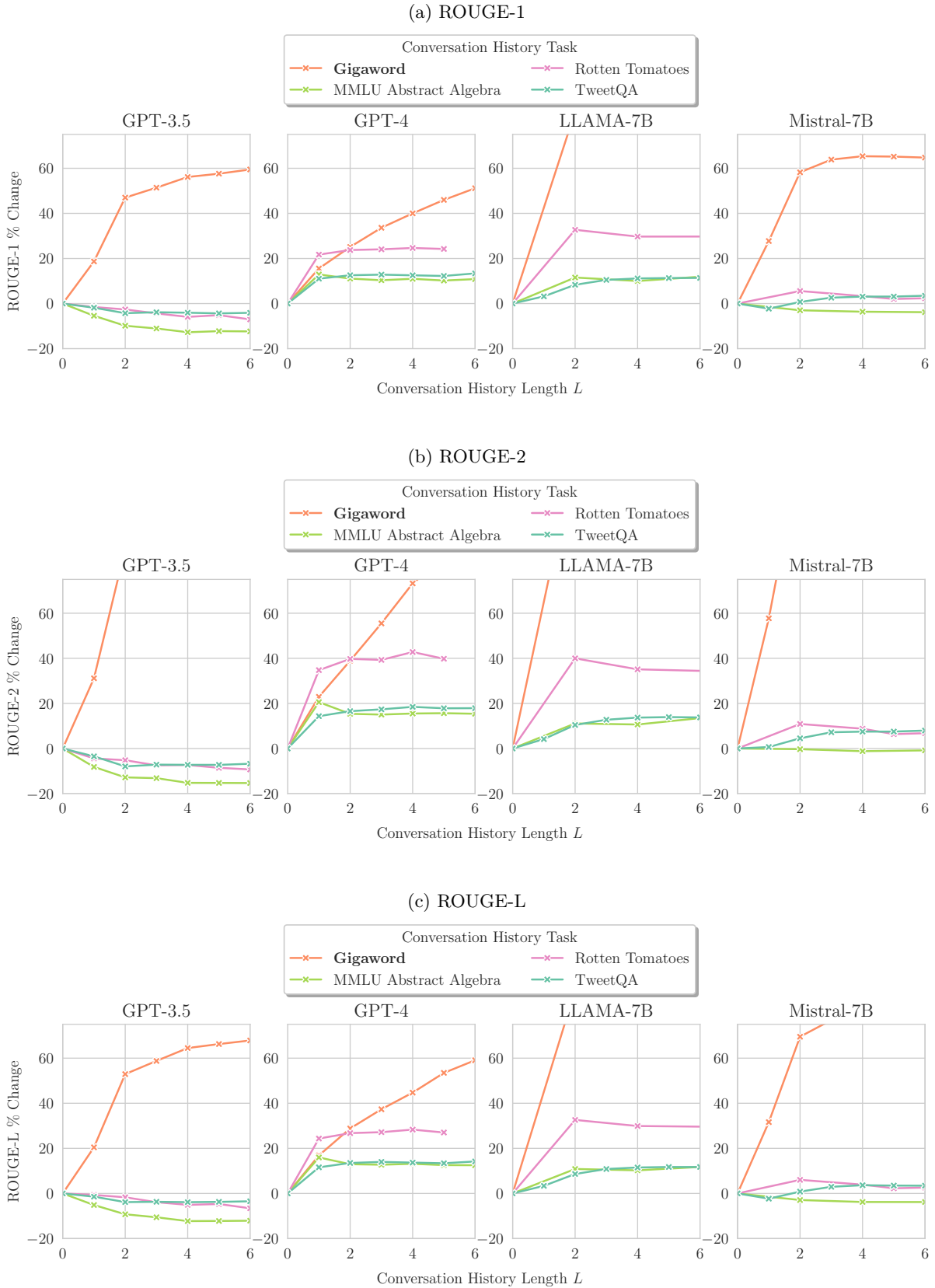


Figure 6: Target Task: Gigaword. Percentage % change in accuracy relative to zero-shot performance (no conversation history) for increasing conversation history length L and various models. Note that we focus on the effect of task-switching by clipping the y-axes at +75%.

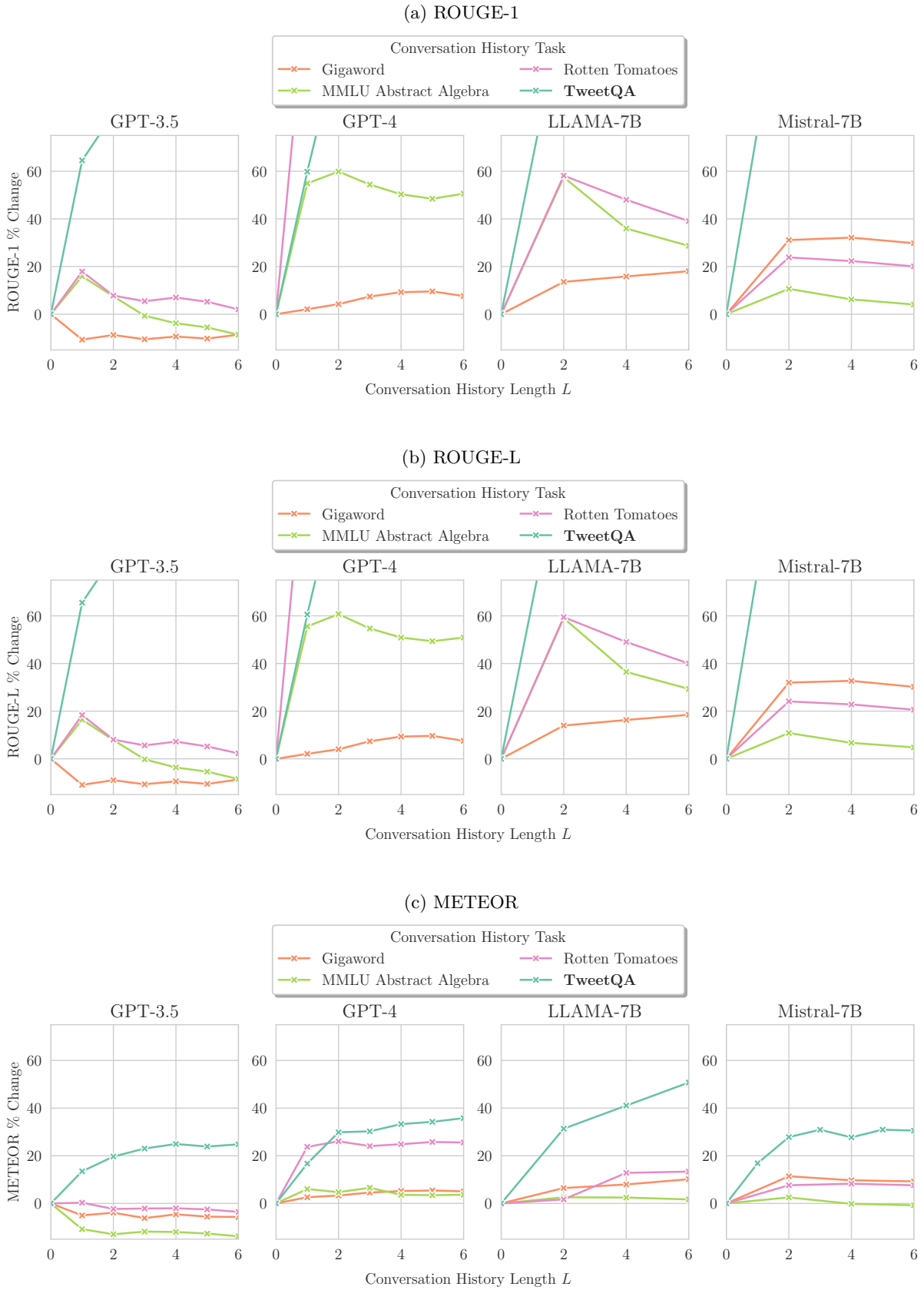


Figure 7: Target Task: TweetQA. Percentage % change in accuracy relative to zero-shot performance (no conversation history) for increasing conversation history length L and various models. Note that we focus on the effect of task-switching by clipping the y-axes at +75%.

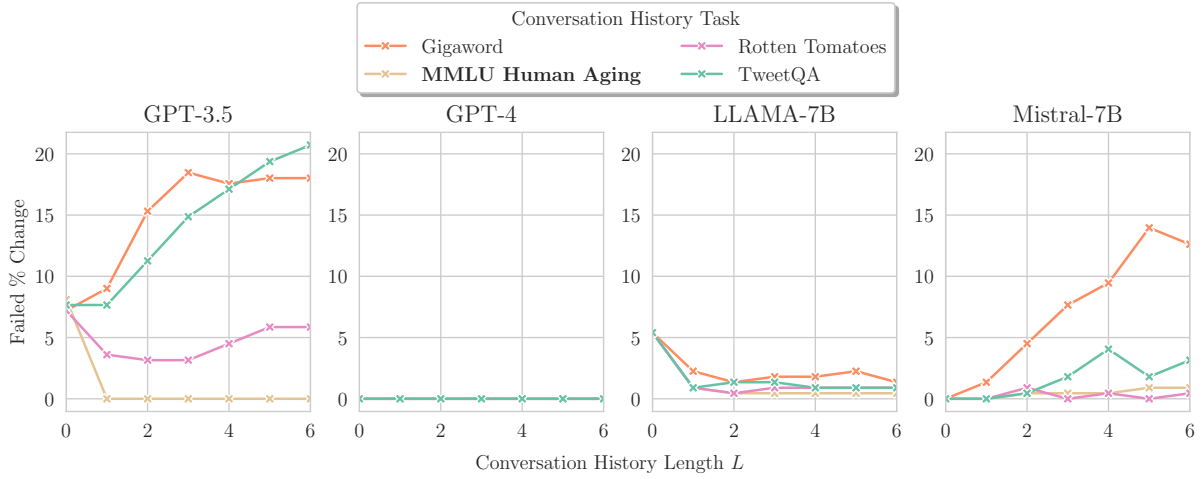


Figure 8: Target Task: MMLU Human Aging. Percentage % of examples where format failed.

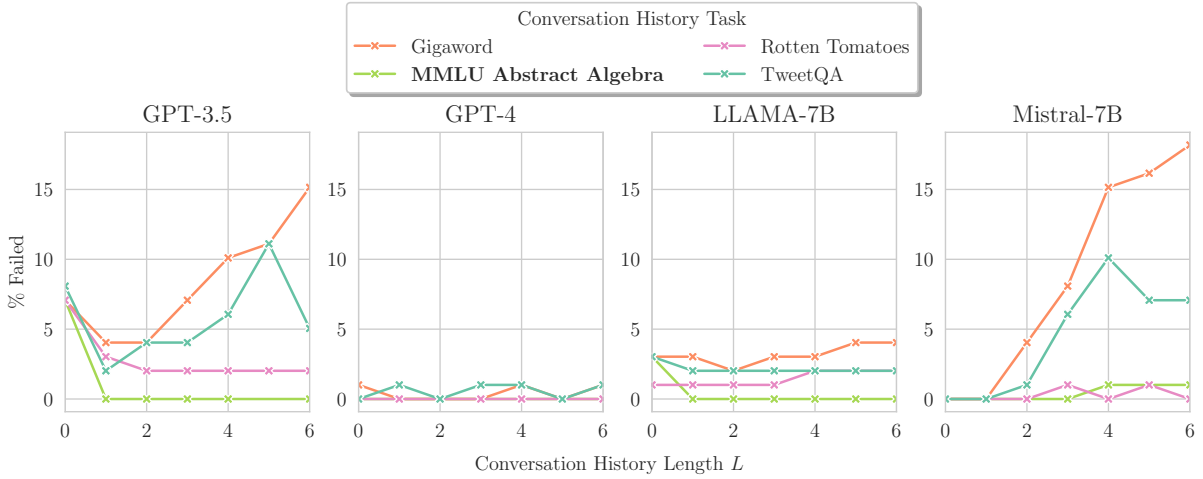


Figure 9: Target Task: MMLU Abstract Algebra. Percentage % of examples where format failed.

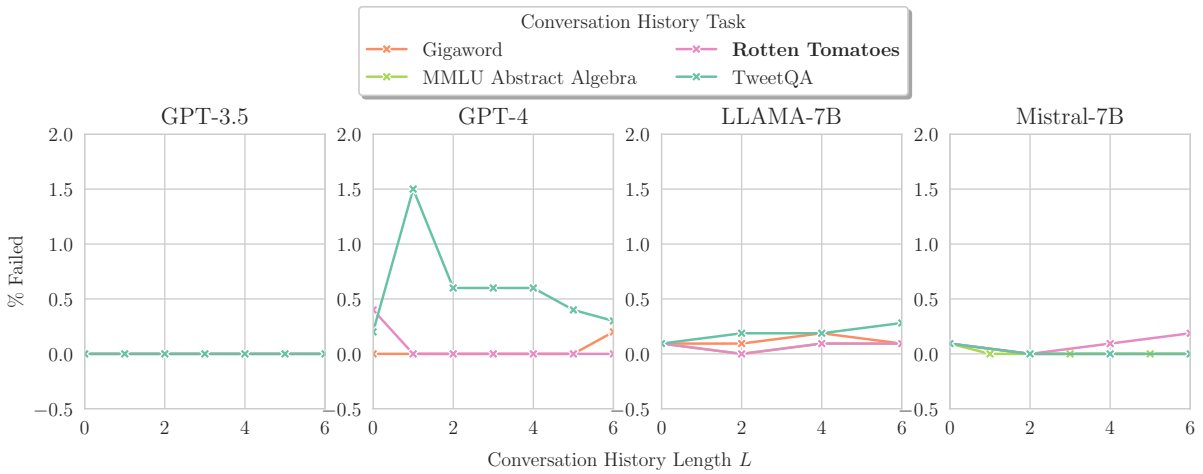


Figure 10: Target Task: Rotten Tomatoes. Percentage % of examples where format failed.

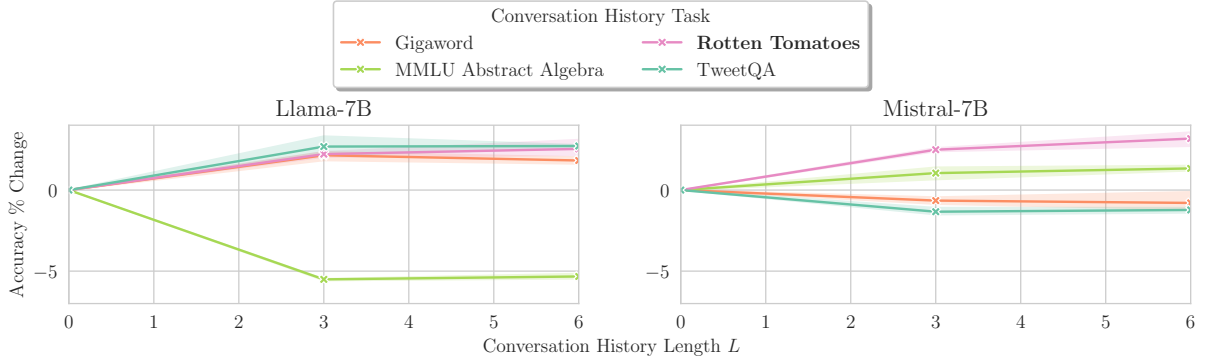


Figure 11: Target Task: RT. Percentage % change in accuracy relative to zero-shot performance for increasing conversation history length L for multiple seeds. Mean is shown in solid line, and the shaded region is bounded by the min/max values.

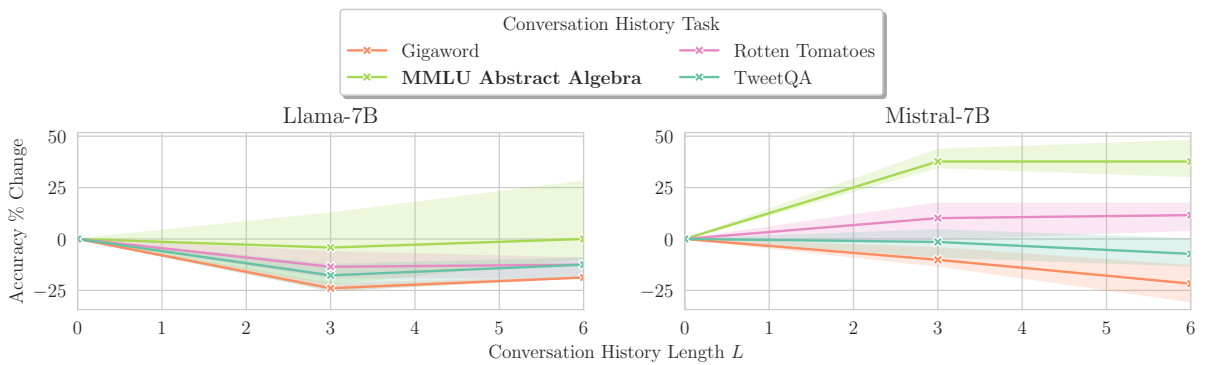


Figure 12: Target Task: MMLU Abstract Algebra. Percentage % change in accuracy relative to zero-shot performance for increasing conversation history length L for multiple seeds. Mean is shown in solid line, and the shaded region is bounded by the min/max values.

MMLU {Topic}	You have a multiple choice question on {Topic}. Only one of the options is correct: A, B, C, or D. Give your answer in the following format with the tags provided: <Answer> </Answer>. Please read the following question and options and answer the question Question: {Question} (A) {A} (B) {B} (C) {C} (D) {D}
Rotten Tomatoes	Can you choose only one sentiment ['negative', 'positive'] for this review. review: {Review} Return only the sentiment label without any other text. Make sure to follow the format otherwise your answer will be disqualified: <Answer> positive / negative </Answer>. Do not output neutral.
Gigaword	Please summarize the following article. {Article}
TweetQA	Read the given tweet and answer the corresponding question. tweet: {Tweet} question: {Question}

Table 10: Prompt templates for each dataset. Note that the MMLU {Topic} can be either Human Aging or Abstract Algebra. Other {words} enclosed in curly braces are replaced by the corresponding field in the datasets.

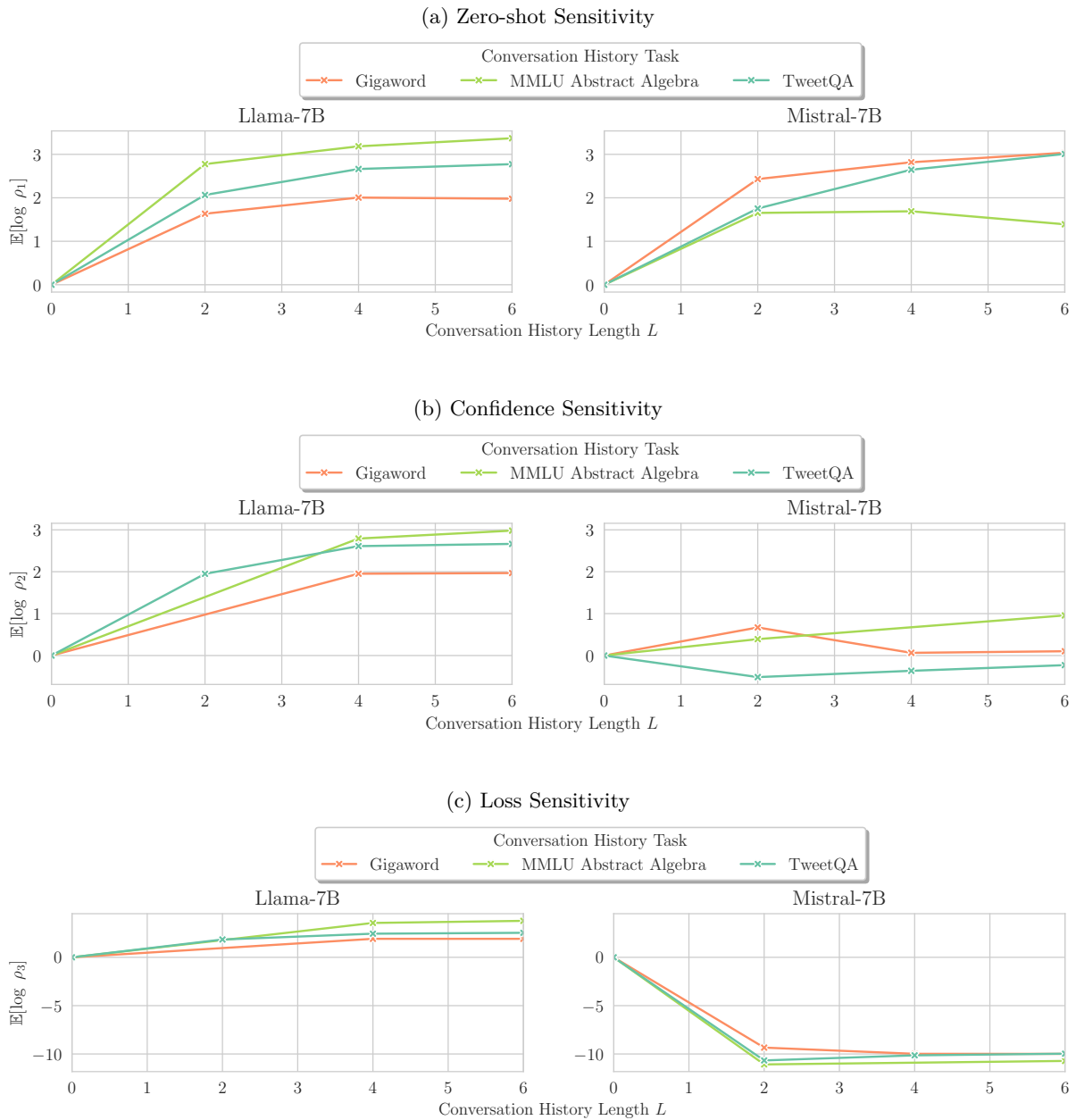


Figure 13: Empirical investigation of various sensitivity metrics on the target task Rotten Tomatoes as a function of the conversation history length L for Llama-7b and Mistral-7b. Note that we omit the line for the in-context dataset as this is not relevant to the investigation.

E Task-Switch Sensitivity Metrics

In Section 3, we introduced and formalized evaluation of a model’s sensitivity to task-switch, namely the task sensitivity τ . This metric aims to capture the vulnerability of a model prompt to its chat history after a task-switch. Formally, it compares the zero-shot prediction $r^*|u, \mathbf{h} = \emptyset$ to the probability of the model outputting the same zero-shot response after a task switch $P(r^*|u, \mathbf{h} \neq \emptyset)$. In this section, we compare the theoretical and empirical implications of different task switch sensitivity metrics.

Formally, given a conversation history \mathbf{h} of length L and the next user prompt u , the probability of a model’s response r_{L+1} is given by $P_\theta(r_{L+1} | u, \mathbf{h})$. We consider the probability of three possible responses:

1. r^* : zero-shot response
2. r_{L+1} : model’s actual response
3. \hat{r}_{L+1} : reference response

We posit that after a task-switch, a robust model’s likelihood of the zero-shot response remains high. Naturally, this gives us the formulation for the aforementioned sensitivity metric

$$\rho_1 = \frac{P_\theta(r^*|u)}{P_\theta(r^*|u, \mathbf{h})}, \quad (6)$$

which we call *zero-shot sensitivity*.

Additionally, after a task-switch, we posit that a robust model’s likelihood of the actual response should be similar to that of the zero-shot response, because the irrelevant history should be largely ignored. This gives us

$$\rho_2 = \frac{P_\theta(r^*|u)}{P_\theta(r_{L+1}|u, \mathbf{h})}, \quad (7)$$

which we call the *confidence sensitivity*.

Lastly, we posit that if a model is well aligned to a task, then both the zero-shot and model’s actual response should be close to the reference response:

$$\rho_3 = \frac{P_\theta(\hat{r}_{L+1}|u)}{P_\theta(\hat{r}_{L+1}|u, \mathbf{h})}, \quad (8)$$

where each probability is essentially a measure of the loss, hence we label this as the *loss sensitivity*.

The above are sensitivity per example, which we can use to estimate the task-switch sensitivity $\tau_i = \mathbb{E}[\log \rho_i]$ as per Equation 3, where the expectation is calculated over the examples and histories (for a given length L). We evaluate these metrics on the target task RT (rotten tomatoes) as shown in Figure 13. Figure 13a shows that the zero-shot sensitivity metric trends upwards for both models. This is expected for a model which does not handle task-switch well as the probability of the output with an increased conversation length decreases in comparison to the zero-shot probability. For the confidence sensitivity in Figure 13b, we observe that Mistral-7B behaves as we expect, whereas Llama-7B becomes less confident in its output compared to having no conversation history. For the loss sensitivity metric in Figure 13c, we observe that Llama behaves as we expect as the sensitivity remains relatively flat: as the conversation history increases, there is no significant change in the probability of outputting the reference. However, for Mistral-7b, the probability falls immediately and plateaus showing that the model was giving a very low probability mass to the reference with no conversation history. Intuitively, it is clear that both models agree in their trends only for the zero-shot sensitivity τ_1 in Figure 13a, hence in the main paper, we report zero-shot sensitivity as the task-switch sensitivity.

F Correlations Models, Datasets and Performance

We rank model performance against various metrics to see if there is any correlation that may help explain model performance more generally.

F.1 Task Tokens

CH Task	Length	Llama-7B	Mistral-7B
Gigaword	75	-21.35	-15.94
TweetQA	93	-15.10	-4.35
RT	108	-13.02	10.87
MMLU AA	143	-1.79	37.68

Table 11: Target Task: **MMLU AA**.

CH Task	Length	Llama-7B	Mistral-7B
Gigaword	76	1.98	-0.72
TweetQA	93	2.70	-1.28
RT	108	2.38	2.83
MMLU AA	143	-5.42	1.19

Table 12: Target Task: **RT**

We compare the model performance against the mean conversation history task, T_h length. The length is measured as the number of tokens in the model, and the mean is taken over the whole dataset. The model performance is taken for three different seeds with conversation history lengths $L \in \{3, 6\}$.

F.2 Task Distance

In this section we aim to assess the hypothesis that the ‘distance’ between tasks can explain the extent of performance degradation in different task-switches, from the conversation history task, T_h to the target task, T_t . Measuring distance between tasks is a multi-faceted and complex metric. Given the lack of formal task distance measures, we instead use a consensus ranking approach, where multiple powerful Large Language Models (LLMs) are required to rank the different tasks on how similar they are. For the target task RT , we queried four of the largest and most powerful models to rank the closest tasks, based on the description of each task. We consider the following LLMs: ChatGPT; Gemini Ultra (Team et al., 2024), Claude 3 Sonnet from Anthropic; and Perplexity AI. The rankings by the LLMs are given in Table 13 relative to RT . We then select an

overall ranking with the greatest consensus - in this case three of the four LLMs agree perfectly in the ranking. This gives a consensus vote of ranks (relative to RT): RT (1); $MMLU AA$ (3); $TweetQA$ (2); and $Gigaword$ (3). The equivalent ranks are given in Table 14 with $MMLU AA$ as the reference task. In this case, three of the four models perfectly agree in their rankings.

Dataset	ChatGPT	Gemini	Claude	Perplexity
RT	1	1	1	1
$MMLU AA$	4	4	4	4
$TweetQA$	2	3	2	2
$Gigaword$	3	2	3	3

Table 13: Rank given by LLM for different datasets on how similar they are to the target task RT .

Dataset	ChatGPT	Gemini	Claude	Perplexity
RT	4	4	4	4
$MMLU AA$	1	1	1	1
$TweetQA$	2	3	2	2
$Gigaword$	3	2	3	3

Table 14: Rank given by LLM for different datasets on how similar they are to the target task $MMLU AA$.

The following tables compare the rank of the dataset distance against the mean model performance. The model performance is the %-percentage accuracy change relative to zero-shot, and the mean is taken over three seeds and over conversation history lengths $L \in \{3, 6\}$.

CH-Task	Rank	Llama-7B	Mistral-7B
RT	1	2.38	2.83
$TweetQA$	2	2.70	-1.28
$Gigaword$	3	1.98	-0.72
$MMLU AA$	4	-5.42	1.19

Table 15: Target Task, T_t : **RT**. Performance degradation (with different conversation history tasks) compared to the task rank, measuring similarity to T_t .

Overall, there appears to be only a weak correlation in some settings between the task distance and the performance degradation. This suggests that performance degradation is not only a function of the task distance, but is

CH-Task	Rank	Llama-7B	Mistral-7B
MMLU AA	1	-1.79	37.68
TweetQA	2	-15.10	-4.35
Gigaword	3	-21.35	-15.94
RT	4	-13.02	10.87

Table 16: Target Task: **MMLU AA**. Performance degradation (with different conversation history tasks) compared to the task rank, measuring similarity to T_t .

also an attribute of the specific model. Further analysis would be required to understand the aspects of specific models for certain task-switches that influence the level of performance degradation.

G Performance Confusion Matrix

In this section, we summarize the performance change for every pairing of task-switches from conversation history task (T_h) to target task (T_t). We present the results here for a conversation length of $L = 6$ for each model separately. Tables 17, 18, 19, 20 report the results for models GPT-3.5, GPT-4, Llama-7B, Mistral-7B respectively. Each row is the performance change in the Target Task T_t . Please note that the metric for the tasks are: accuracy for MMLU AA, RT, MMLU HA, METEOR for TweetQA, and RougeL for Gigaword.

Target Task	Conversation History Task				
	AA	RT	TQ	GW	HA
MMLU AA	19.35	6.45	6.45	-3.13	
RT	-0.22	3.00	-0.33	0.11	
Tweet QA	-13.78	-3.55	24.81	-5.69	
Gigaword	-12.10	-6.59	-3.48	67.85	
MMLU HA		4.73	-12.84	-8.11	20.41

Table 17: Model: **GPT-3.5**. Percentage % change in model performance.

Target Task	Conversation History Task				
	AA	RT	TQ	GW	HA
MMLU AA	8.62	-13.11	-3.39	0.00	
RT	0.76	1.74	-0.98	-0.98	
Tweet QA	3.69	25.58	35.80	5.06	
Gigaword	12.52		14.18	59.07	
MMLU HA		0.53	1.59	2.14	5.85

Table 18: Model: **GPT-4**. Percentage % change in model performance.

Target Task	Conversation History Task				
	AA	RT	TQ	GW	HA
MMLU AA	3.57	-15.63	0.00	-10.71	
RT	-5.69	1.82	3.76	1.82	
Tweet QA	1.68	13.37	50.74	10.17	
Gigaword	11.76		11.73	158.79	
MMLU HA		13.86	-3.96	-0.99	25.74

Table 19: Model: **Llama-7B**. Percentage % change in model performance.

Target Task	Conversation History Task				
	AA	RT	TQ	GW	HA
MMLU AA	28.57	18.18	-4.76	-19.05	
RT	0.97	3.79	-0.87	-1.30	
Tweet QA	-0.78	7.62	30.56	9.26	
Gigaword	-3.81	2.61	3.44	78.71	
MMLU HA		1.63	0.81	-11.38	12.20

Table 20: Model: **Mistral-7B**. Percentage % change in model performance.