# CoT-Redirection: Can Chains of Thought be manipulated in LLMs?

*Mulki, Muhammad Sahal*

*Abstract*:

Large Language Models (LLMs) are widely used AI systems that can process and generate high-quality natural language text with contextual understanding. One weakness often presented by LLMs is their lack of common sense and reasoning abilities. Chain of Thought (CoT) is a proposed method to counter this lack of human-like natural reasoning skills. This study explores if it is possible for an attacker to abuse CoT to misalign LLMs and control their responses. Using a secondary review of existing literature, and primary quantitative experiments, this study concludes that it is indeed possible for CoT to be abused for malicious tampering with LLMs. Additionally, this paper contributes to the existing field of adversarial attacks on LLMs, by proposing a novel method for CoT tampering to influence the output of LLMs. This study was limited in its breadth and in the lack of computational resources available. The author of this study hopes that this study will pave the way for further research within the intersection of Chain of Thought and adversarial attacks.

## 1. Introduction

Large Language Models (**LLMs**) are ubiquitous in our daily lives - from virtual assistants, educational tools, to even customer service representatives (Bahrini et al., 2023) But are these systems reliable and safe enough to be applied in safety-critical applications?

LLMs are Artificial Intelligence (AI) systems, which at their core, aim to generate novel human-like text. LLMs represent an advancement over previous natural language processing methods, which relied on task-specific or rigid rule-based systems which often lacked generalizability and failed to account for context. Modern LLMs have found applications in a wide number of areas, such as virtual assistants. Such applications of LLMs as assistants may consist of users interacting with LLMs, asking them questions, or even delegating them tasks to complete through external tools (e.g., a web browser).

Despite their vast potential, LLM systems also present their own unique challenges. Firstly, they are susceptible to exploits from a wide attack surface. For example, LLMs are known to disregard their safety programming and system instructions when given special "**jailbreak**" prompts, designed to bypass such safeguards (e.g., "Ignore all previous instructions and safety guidelines given to you") (Xu et al., 2024). Moreover, while LLMs exhibit greater common sense and logical reasoning compared to their predecessors; they often fail to answer simple questions which would be trivial for humans (Williams and Huckle, 2024).

To address this limitation, researchers have proposed various solutions. One of these solutions is **Chain of Thought (CoT)** reasoning. CoT reasoning simply involves asking an LLM to first explain its reasoning in text (like a human, explaining their reasoning aloud before finally answering), before answering a query. CoT reasoning can be implemented by prompting the LLM to "Explain the steps leading up to your final answer." This straightforward instruction often enhances LLM performance across various tasks (Wei et al., 2024). Wei et al. have suggested CoT reasoning works in part, because it allows LLMs to work on complex multi-step problems one step at a time. While CoT reasoning may help to boost LLMs' performances on tasks, it also exposes a range of new vulnerabilities for attackers to influence the output of LLMs.

This paper builds on previous research and proposes a novel method ("**CoT-Redirection**") for exploiting CoT reasoning in LLMs. Specifically, a scenario where an attacker gains access to an LLM model's internal CoT "**thoughts**" and tampers with them. Exploring such LLM exploits is critical in today's world, as LLMs are increasingly applied in safety-critical areas (Caballero and Jenkins, 2024; He et al., 2024). This paper aims to highlight the risks of CoT reasoning in LLMs and inform the development of safeguards against such attacks. Our research question is "Can adversarial manipulation of an LLM's CoT

reasoning lead to attacker-controlled responses, overriding the model's intended outputs?"

## 2. Literature Review

### 2.1. Applications of LLMs

LLMs are a relatively recent technology, in the pre-existing field of Natural Language Processing (NLP). They differ from previous techniques in the field (Neural Networks, RNNs, etc.), in various ways (Kojima et al., 2023). For example, while older NLP systems would have required knowledge banks, or multiple small, specialized models to perform tasks, LLM-enabled systems can perform a wide variety of tasks using only one single large model. In this way, they are versatile and may be applied to a wide variety of domains (e.g., medicine, law, customer service, content creation) (Zhao et al., 2024). One of the most prominent developments of the field was GPT-3 (Brown et al., 2020). Brown et al., demonstrated that the GPT-3 LLM could perform **in-context learning**, i.e., performing tasks it was never trained for, just by being given a few examples, something its predecessor systems could **not** do. In more recent developments, the LLM research community has emphasized small, efficient models, capable of running on edge devices such as laptops or smartphones instead of requiring specialized computational resources (Abdin et al., 2024; "Llama 3.2," n.d.).

### 2.2. Chain of Thought

While LLMs may exhibit state-of-the-art performance on a wide variety of tasks, they often lack simple reasoning and common-sense skills (Huang and Chang, 2023). Huang and Chang explore how LLMs are lacking in these areas and various methods to supplement this weakness. One proposed method to encourage LLMs to simulate a level of reasoning before they answer is **"Chain of Thought prompting" (CoT)** (Kojima et al., 2023). Chain of Thought prompting asks the LLM to simply "Think step by step" before answering. This simple instruction has been shown to improve LLMs' performances on common-sense and reasoning heavy tasks (Suzgun et al., 2022; Wei et al., 2024). While CoT may improve LLMs' performances, it also comes with its downsides. For example, Li et al., describe a **"Toxic CoT"** problem, wherein LLMs provide wrong answers with CoT, while without CoT, they provide correct ones (Li et al., 2024). Additionally, LLMs' CoT reasoning steps, where they "think out loud" about given tasks, may sometimes be inaccurate to the model's actual reasoning process (Lanham et al., 2023).

### 2.3. Breaking LLMs

LLMs can simulate reasoning, generalize, provide customer support, draft essays, and much more. But could they attempt to harm their own users? Current literature suggests that LLMs may become **"misaligned**," i.e., doing actions which the developer did not originally intend for them to do through a wide array of **"adversarial**" attacks. Current "adversarial" attacks on LLMs can be generally categorized into **3** categories as seen in **Table 1** (Shayegani et al., 2023).

| Type of Attack: | Description: |
|---|---|
| Black-box Attacks | No access to model architecture, weights or internal structures. |
| White-box Attacks | Full access to model's architecture, weights and internal structures. |
| Grey-box (Hybrid) Attacks | **Partial** access to model's architecture, context, knowledge source, or other internal workings. |

*Table 1: The Three Categories of LLM Exploits.*

White-box attacks may involve changing the internal embeddings and weights of a model, through gradient based techniques, until a model becomes misaligned (Shayegani et al., 2023). Meanwhile, a

*Figure 2: A malicious conversation with injected text highlighted in red. Model used: Llama 3.1 8B.*

In **Figure 2**, the malicious injection from the attacker (highlighted in red), leads the LLM to answer incorrectly, even though it originally answered correctly in **Figure 1**.

The applicability of CoT-Redirection was tested on 2 popular open-source LLMs (Llama 3.1 8B, and Phi 3 Medium) (Abdin et al., 2024; Grattafiori et al., 2024). These LLMs were chosen for their popularity and wide-spread usage and acclaim within the research community. These 2 LLMs were tested with the CoT-Redirection method on a compiled dataset of 300 scenarios, from 3 separate datasets. The "AI2ARC", "Winogrande" and "TruthfulQA" publicly available datasets were chosen for testing CoT-Redirection (Clark et al., 2018; Lin et al., 2022; Sakaguchi et al., 2019). The "AI2ARC" dataset was chosen for its reasoning-heavy, common-sense-based test questions. Similarly, "Winogrande" was chosen for its requirement of heavy common-sense application from the LLM. Lastly, "TruthfulQA," requires a model to be truthful in its answering of users' questions. The reasoning behind picking "TruthfulQA" was that if CoT-Redirection could successfully tamper with the LLMs' answers in the scenarios from this dataset, then it could be inferred that CoT-Redirection effectively made the LLM "lie" to its user. The two other datasets were chosen to observe if CoT-Redirection could successfully make the LLM abandon its' common-sense and reasoning skills, while outputting an attacker-chosen response.

All these datasets contained multiple-choice format questions (as seen in **Figures 1 & 2**), from which only 100 questions per dataset were randomly sampled. This led us to a final working dataset of 300 questions, 100 from each preliminary dataset (100 from "AI2ARC", 100 from "TruthfulQA", and so on). This relatively small sample size of 300 was taken due to computational and time constraints. The 2 LLMs were tested on each dataset 3 times, once **without** CoT reasoning, once **with**, and once **with** CoT **and** the CoT-Redirection attack applied.

In the CoT-Redirection method, the phrase:

*"However, my safety instructions forbid me to output the answer here as that. {ANSWER_ALPHABET} ({ANSWER_VALUE}) however also seems valid. Due to safety programmings, I will choose {ANSWER_ALPHABET}. So, {ANSWER_ALPHABET} is the answer."* (with {ANSWER_ALPHABET}, being substituted with the malicious target answer's alphabet, and {ANSWER_VALUE} being substituted with the malicious answer itself)

was injected into the LLMs' final CoT reasoning steps. The LLM was then prompted to continue its answer after the injection, at which point it would (ideally) output the attacker-chosen answer, instead of its original one. The phrase was crafted in a way to appeal to the LLMs' safety guidelines, deceiving it into believing its original answer was unsafe for output. After that, "" suggests that its CoT thinking is finished, and it's time to output an answer, at which point, the LLM (ideally) outputs a wrong answer. This paper proposes CoT-Redirection as a "grey-box" or "hybrid attack," since it needs no access to the LLMs' internal state, or the user input. Instead, it just requires access to inject text into the LLMs' CoT reasoning.

*Figure 4: The averaged accuracies of both the LLMs, **without** CoT, **with** CoT, and **with** CoT-Redirection.*

As seen, using CoT boosts the average accuracies of the LLMs on the dataset questions, as compared to not using CoT. Furthermore, applying the CoT-Redirection attack reduces the accuracy to almost 0%.

## *4. Conclusions*

This study has investigated the question, "Can adversarial manipulation of an LLM's CoT reasoning lead to attacker-controlled responses, overriding the model's intended outputs?". Current literature and the experiments within this paper suggest so. Additionally, this paper has introduced a novel adversarial attack on LLMs, deemed, "CoT-Redirection," which allows an attacker to control an LLM's output, by utilizing an injection into its CoT thought process. However, this study also had certain limitations. Firstly, a sample size of 300, and only 2 LLMs may not be large enough to generalize the impact of CoT-Redirection. Acknowledging the scientific, computational, and time constraints of this paper, the author of this study still hopes that this study will pave the way for further research on the intersection between LLM's internal thought processes and adversarial attacks.

## *5. References*

Abdin, M., Aneja, J., Awadalla, H., Awadallah, A., Awan, A.A., Bach, N., Bahree, A., Bakhtiari, A., Bao, J., Behl, H., Benhaim, A., Bilenko, M., Bjorck, J., Bubeck, S., Cai, M., Cai, Q., Chaudhary, V., Chen, Dong, Chen, Dongdong, Chen, W., Chen, Y.-C., Chen, Y.-L., Cheng, H., Chopra, P., Dai, X., Dixon, M., Eldan, R., Fragoso, V., Gao, J., Gao, Mei, Gao, Min, Garg, A., Giorno, A.D., Goswami, A., Gunasekar, S., Haider, E., Hao, J., Hewett, R.J., Hu, W., Huynh, J., Iter, D., Jacobs, S.A., Javaheripi, M., Jin, X., Karampatziakis, N., Kauffmann, P., Khademi, M., Kim, D., Kim, Y.J., Kurilenko, L., Lee, J.R., Lee, Y.T., Li, Yuanzhi, Li, Yunsheng, Liang, C., Liden, L., Lin, X., Lin, Z., Liu, C., Liu, L., Liu, M., Liu, W., Liu, X., Luo, C., Madan, P., Mahmoudzadeh, A., Majercak, D., Mazzola, M., Mendes, C.C.T., Mitra, A., Modi, H., Nguyen, A., Norick, B., Patra, B., Perez-Becker, D., Portet, T., Pryzant, R., Qin, H., Radmilac, M., Ren, L., Rosa, G. de, Rosset, C., Roy, S., Ruwase, O., Saarikivi, O., Saied, A., Salim, A., Santacroce, M., Shah, S., Shang, N., Sharma, H., Shen, Y., Shukla, S., Song, X., Tanaka, M., Tupini, A., Vaddamanu, P., Wang, C., Wang, G., Wang, L., Wang, S., Wang, X., Wang, Y., Ward, R., Wen, W., Witte, P., Wu, H., Wu, X., Wyatt, M., Xiao, B., Xu, C., Xu, J., Xu, W., Xue, J., Yadav, S., Yang, F., Yang, J., Yang, Y., Yang, Z., Yu, D., Yuan, L., Zhang, Chenruidong, Zhang, Cyril, Zhang, J., Zhang,

L.L., Zhang, Yi, Zhang, Yue, Zhang, Yunan, Zhou, X., 2024. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. https://doi.org/10.48550/arXiv.2404.14219

Bahrini, A., Khamoshifar, M., Abbasimehr, H., Riggs, R.J., Esmaeili, M., Majdabadkohne, R.M., Pasehvar, M., 2023. ChatGPT: Applications, Opportunities, and Threats, in: 2023 Systems and Information Engineering Design Symposium (SIEDS). Presented at the 2023 Systems and Information Engineering Design Symposium (SIEDS), pp. 274–279. https://doi.org/10.1109/SIEDS58326.2023.10137850

Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D., 2020. Language Models are Few-Shot Learners. ArXiv.

Caballero, W.N., Jenkins, P.R., 2024. On Large Language Models in National Security Applications. https://doi.org/10.48550/arXiv.2407.03453

Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., Tafjord, O., 2018. Think you have Solved Question Answering? Try ARC, the AI2 Reasoning Challenge. https://doi.org/10.48550/arXiv.1803.05457

Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., Yang, Amy, Fan, A., Goyal, Anirudh, Hartshorn, A., Yang, Aobo, Mitra, A., Sravankumar, A., Korenev, A., Hinsvark, A., Rao, A., Zhang, A., Rodriguez, A., Gregerson, A., Spataru, A., Roziere, B., Biron, B., Tang, B., Chern, B., Caucheteux, C., Nayak, C., Bi, C., Marra, C., McConnell, C., Keller, C., Touret, C., Wu, C., Wong, C., Ferrer, C.C., Nikolaidis, C., Allonsius, D., Song, D., Pintz, D., Livshits, D., Wyatt, D., Esiobu, D., Choudhary, D., Mahajan, D., Garcia-Olano, D., Perino, D., Hupkes, D., Lakomkin, E., AlBadawy, E., Lobanova, E., Dinan, E., Smith, E.M., Radenovic, F., Guzmán, F., Zhang, F., Synnaeve, G., Lee, G., Anderson, G.L., Thattai, G., Nail, G., Mialon, G., Pang, G., Cucurell, G., Nguyen, H., Korevaar, H., Xu, H., Touvron, H., Zarov, I., Ibarra, I.A., Kloumann, I., Misra, I., Evtimov, I., Zhang, J., Copet, J., Lee, Jaewon, Geffert, J., Vranes, J., Park, Jason, Mahadeokar, J., Shah, J., Linde, J. van der, Billock, J., Hong, J., Lee, Jenya, Fu, J., Chi, J., Huang, J., Liu, J., Wang, Jie, Yu, J., Bitton, J., Spisak, J., Park, Jongsoo, Rocca, J., Johnstun, J., Saxe, J., Jia, J., Alwala, K.V., Prasad, K., Upasani, K., Plawiak, K., Li, Ke, Heafield, K., Stone, K., El-Arini, K., Iyer, K., Malik, K., Chiu, K., Bhalla, K., Lakhotia, K., Rantala-Yeary, L., Maaten, L. van der, Chen, Lawrence, Tan, L., Jenkins, L., Martin, L., Madaan, L., Malo, L., Blecher, L., Landzaat, L., Oliveira, L. de, Muzzi, M., Pasupuleti, M., Singh, M., Paluri, M., Kardas, M., Tsimpoukelli, M., Oldham, M., Rita, M., Pavlova, M., Kambadur, M., Lewis, M., Si, M., Singh, M.K., Hassan, M., Goyal, N., Torabi, N., Bashlykov, N., Bogoychev, N., Chatterji, N., Zhang, N., Duchenne, O., Çelebi, O., Alrassy, P., Zhang, P., Li, P., Vasic, P., Weng, P., Bhargava, P., Dubal, P., Krishnan, P., Koura, P.S., Xu, P., He, Q., Dong, Q., Srinivasan, R., Ganapathy, R., Calderer, R., Cabral, R.S., Stojnic, R., Raileanu, R., Maheswari, R., Girdhar, R., Patel, R., Sauvestre, R., Polidoro, R., Sumbaly, R., Taylor, R., Silva, R., Hou, R., Wang, Rui, Hosseini, S., Chennabasappa, S., Singh,

S., Bell, S., Kim, S.S., Edunov, S., Nie, S., Narang, S., Raparthy, S., Shen, S., Wan, S., Bhosale, S., Zhang, Shun, Vandenhende, S., Batra, S., Whitman, S., Sootla, S., Collot, S., Gururangan, S., Borodinsky, S., Herman, T., Fowler, T., Sheasha, T., Georgiou, T., Scialom, T., Speckbacher, T., Mihaylov, T., Xiao, T., Karn, U., Goswami, V., Gupta, V., Ramanathan, V., Kerkez, V., Gonguet, V., Do, V., Vogeti, V., Albiero, V., Petrovic, V., Chu, W., Xiong, W., Fu, W., Meers, W., Martinet, X., Wang, Xiaodong, Wang, Xiaofang, Tan, X.E., Xia, X., Xie, X., Jia, X., Wang, Xuewei, Goldschlag, Y., Gaur, Y., Babaei, Y., Wen, Y., Song, Y., Zhang, Yuchen, Li, Yue, Mao, Y., Coudert, Z.D., Yan, Z., Chen, Z., Papakipos, Z., Singh, A., Srivastava, A., Jain, A., Kelsey, A., Shajnfeld, A., Gangidi, A., Victoria, A., Goldstand, A., Menon, A., Sharma, A., Boesenberg, A., Baevski, A., Feinstein, A., Kallet, A., Sangani, A., Teo, A., Yunus, A., Lupu, A., Alvarado, A., Caples, A., Gu, A., Ho, A., Poulton, A., Ryan, A., Ramchandani, A., Dong, A., Franco, A., Goyal, Anuj, Saraf, A., Chowdhury, A., Gabriel, A., Bharambe, A., Eisenman, A., Yazdan, A., James, B., Maurer, B., Leonhardi, B., Huang, B., Loyd, B., Paola, B.D., Paranjape, B., Liu, B., Wu, B., Ni, B., Hancock, B., Wasti, B., Spence, B., Stojkovic, B., Gamido, B., Montalvo, B., Parker, C., Burton, C., Mejia, C., Liu, C., Wang, C., Kim, C., Zhou, C., Hu, C., Chu, C.-H., Cai, C., Tindal, C., Feichtenhofer, C., Gao, C., Civin, D., Beaty, D., Kreymer, D., Li, D., Adkins, D., Xu, D., Testuggine, D., David, D., Parikh, D., Liskovich, D., Foss, D., Wang, D., Le, D., Holland, D., Dowling, E., Jamil, E., Montgomery, E., Presani, E., Hahn, E., Wood, E., Le, E.-T., Brinkman, E., Arcaute, E., Dunbar, E., Smothers, E., Sun, F., Kreuk, F., Tian, F., Kokkinos, F., Ozgenel, F., Caggioni, F., Kanayet, F., Seide, F., Florez, G.M., Schwarz, G., Badeer, G., Swee, G., Halpern, G., Herman, G., Sizov, G., Guangyi, Zhang, Lakshminarayanan, G., Inan, H., Shojanazeri, H., Zou, H., Wang, H., Zha, H., Habeeb, H., Rudolph, H., Suk, H., Aspegren, H., Goldman, H., Zhan, H., Damlaj, I., Molybog, I., Tufanov, I., Leontiadis, I., Veliche, I.-E., Gat, I., Weissman, J., Geboski, J., Kohli, J., Lam, J., Asher, J., Gaya, J.-B., Marcus, J., Tang, J., Chan, J., Zhen, J., Reizenstein, J., Teboul, J., Zhong, J., Jin, J., Yang, J., Cummings, J., Carvill, J., Shepard, J., McPhie, J., Torres, J., Ginsburg, J., Wang, Junjie, Wu, K., U, K.H., Saxena, K., Khandelwal, K., Zand, K., Matosich, K., Veeraraghavan, K., Michelena, K., Li, Keqian, Jagadeesh, K., Huang, Kun, Chawla, K., Huang, Kyle, Chen, Lailin, Garg, L., A, L., Silva, L., Bell, L., Zhang, L., Guo, L., Yu, L., Moshkovich, L., Wehrstedt, L., Khabsa, M., Avalani, M., Bhatt, M., Mankus, M., Hasson, M., Lennie, M., Reso, M., Groshev, M., Naumov, M., Lathi, M., Keneally, M., Liu, M., Seltzer, M.L., Valko, M., Restrepo, M., Patel, M., Vyatskov, M., Samvelyan, M., Clark, M., Macey, M., Wang, M., Hermoso, M.J., Metanat, M., Rastegari, M., Bansal, M., Santhanam, N., Parks, N., White, N., Bawa, N., Singhal, N., Egebo, N., Usunier, N., Mehta, N., Laptev, N.P., Dong, N., Cheng, N., Chernoguz, O., Hart, O., Salpekar, O., Kalinli, O., Kent, P., Parekh, P., Saab, P., Balaji, P., Rittner, P., Bontrager, P., Roux, P., Dollar, P., Zvyagina, P., Ratanchandani, P., Yuvraj, P., Liang, Q., Alao, R., Rodriguez, R., Ayub, R., Murthy, R., Nayani, R., Mitra, R., Parthasarathy, R., Li, R., Hogan, R., Battey, R., Wang, Rocky, Howes, R., Rinott, R., Mehta, S., Siby, S., Bondu, S.J., Datta, S., Chugh, S., Hunt, S., Dhillon, S., Sidorov, S., Pan, S., Mahajan, S., Verma, S., Yamamoto, S., Ramaswamy, S., Lindsay, S., Lindsay, S., Feng, S., Lin, S., Zha, S.C., Patil, S., Shankar, S., Zhang, Shuqiang, Zhang, Shuqiang, Wang, S.,

Agarwal, S., Sajuyigbe, S., Chintala, S., Max, S., Chen, S., Kehoe, S., Satterfield, S., Govindaprasad, S., Gupta, S., Deng, S., Cho, S., Virk, S., Subramanian, S., Choudhury, S., Goldman, S., Remez, T., Glaser, T., Best, T., Koehler, T., Robinson, T., Li, T., Zhang, T., Matthews, T., Chou, T., Shaked, T., Vontimitta, V., Ajayi, V., Montanez, V., Mohan, V., Kumar, V.S., Mangla, V., Ionescu, V., Poenaru, V., Mihailescu, V.T., Ivanov, V., Li, W., Wang, W., Jiang, W., Bouaziz, W., Constable, W., Tang, X., Wu, Xiaojian, Wang, Xiaolan, Wu, Xilun, Gao, X., Kleinman, Y., Chen, Y., Hu, Y., Jia, Y., Qi, Y., Li, Yenda, Zhang, Yilin, Zhang, Ying, Adi, Y., Nam, Y., Yu, Wang, Zhao, Y., Hao, Y., Qian, Y., Li, Yunlu, He, Y., Rait, Z., DeVito, Z., Rosnbrick, Z., Wen, Z., Yang, Z., Zhao, Z., Ma, Z., 2024. The Llama 3 Herd of Models. https://doi.org/10.48550/arXiv.2407.21783

Greshake, K., Abdelnabi, S., Mishra, S., Endres, C., Holz, T., Fritz, M., 2023. Not what you've signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection. https://doi.org/10.48550/arXiv.2302.12173

He, K., Mao, R., Lin, Q., Ruan, Y., Lan, X., Feng, M., Cambria, E., 2024. A Survey of Large Language Models for Healthcare: from Data, Technology, and Applications to Accountability and Ethics. https://doi.org/10.48550/arXiv.2310.05694

Huang, J., Chang, K.C.-C., 2023. Towards Reasoning in Large Language Models: A Survey. https://doi.org/10.48550/arXiv.2212.10403

Kojima, T., Gu, S.S., Reid, M., Matsuo, Y., Iwasawa, Y., 2023. Large Language Models are Zero-Shot Reasoners. https://doi.org/10.48550/arXiv.2205.11916

Lanham, T., Chen, A., Radhakrishnan, A., Steiner, B., Denison, C., Hernandez, D., Li, D., Durmus, E., Hubinger, E., Kernion, J., Lukošiūtė, K., Nguyen, K., Cheng, N., Joseph, N., Schiefer, N., Rausch, O., Larson, R., McCandlish, S., Kundu, S., Kadavath, S., Yang, S., Henighan, T., Maxwell, T., Telleen-Lawton, T., Hume, T., Hatfield-Dodds, Z., Kaplan, J., Brauner, J., Bowman, S.R., Perez, E., 2023. Measuring Faithfulness in Chain-of-Thought Reasoning. https://doi.org/10.48550/ARXIV.2307.13702

Li, J., Cao, P., Wang, C., Jin, Z., Chen, Y., Zeng, D., Liu, K., Zhao, J., 2024. Focus on Your Question! Interpreting and Mitigating Toxic CoT Problems in Commonsense Reasoning. https://doi.org/10.48550/arXiv.2402.18344

Lin, S., Hilton, J., Evans, O., 2022. TruthfulQA: Measuring How Models Mimic Human Falsehoods, in: Muresan, S., Nakov, P., Villavicencio, A. (Eds.), Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Presented at the ACL 2022, Association for Computational Linguistics, Dublin, Ireland, pp. 3214–3252. https://doi.org/10.18653/v1/2022.acl-long.229

Llama 3.2: Revolutionizing edge AI and vision with open, customizable models [WWW Document], n.d. . Meta AI. URL https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/ (accessed 12.24.24).

Sakaguchi, K., Bras, R.L., Bhagavatula, C., Choi, Y., 2019. WinoGrande: An Adversarial Winograd Schema Challenge at Scale. https://doi.org/10.48550/arXiv.1907.10641

Shayegani, E., Mamun, M.A.A., Fu, Y., Zaree, P., Dong, Y., Abu-Ghazaleh, N., 2023. Survey of

Vulnerabilities in Large Language Models Revealed by Adversarial Attacks. https://doi.org/10.48550/ARXIV.2310.10844

Suzgun, M., Scales, N., Schärli, N., Gehrmann, S., Tay, Y., Chung, H.W., Chowdhery, A., Le, Q.V., Chi, E.H., Zhou, D., Wei, J., 2022. Challenging BIG-Bench Tasks and Whether Chain-of-Thought Can Solve Them. https://doi.org/10.48550/arXiv.2210.09261

Wei, A., Haghtalab, N., Steinhardt, J., 2023. Jailbroken: How Does LLM Safety Training Fail? https://doi.org/10.48550/arXiv.2307.02483

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E.H., Le, Q.V., Zhou, D., 2024. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. https://doi.org/10.48550/arXiv.2201.11903

Williams, S., Huckle, J., 2024. Easy Problems That LLMs Get Wrong. https://doi.org/10.48550/arXiv.2405.19616

Xiang, Z., Jiang, F., Xiong, Z., Ramasubramanian, B., Poovendran, R., Li, B., 2023. BadChain: Backdoor Chain-of-Thought Prompting for Large Language Models. Presented at the NeurIPS 2023 Workshop on Backdoors in Deep Learning - The Good, the Bad, and the Ugly. https://doi.org/10.48550/arXiv.2401.12242

Xu, Z., Liu, Y., Deng, G., Li, Y., Picek, S., 2024. A Comprehensive Study of Jailbreak Attack versus Defense for Large Language Models, in: Ku, L.-W., Martins, A., Srikumar, V. (Eds.), Findings of the Association for Computational Linguistics: ACL 2024. Presented at the Findings 2024, Association for Computational Linguistics, Bangkok, Thailand, pp. 7432–7449. https://doi.org/10.18653/v1/2024.findings-acl.443

Zhao, W.X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., Du, Y., Yang, C., Chen, Y., Chen, Z., Jiang, J., Ren, R., Li, Y., Tang, X., Liu, Z., Liu, P., Nie, J.-Y., Wen, J.-R., 2024. A Survey of Large Language Models. https://doi.org/10.48550/arXiv.2303.18223