# Invasive Context Engineering
# to Control Large Language Models

Thomas Rivasseau

*McGill University*

*Abstract*—Current research on operator control of Large Language Models improves model robustness against adversarial attacks and misbehavior by training on preference examples, prompting, and input/output filtering. Despite good results, LLMs remain susceptible to abuse, and jailbreak probability increases with context length. There is a need for robust LLM security guarantees in long-context situations. We propose control sentences inserted into the LLM context as invasive context engineering to partially solve the problem. We suggest this technique can be generalized to the Chain-of-Thought process to prevent scheming. Invasive Context Engineering does not rely on LLM training, avoiding data shortage pitfalls which arise in training models for long context situations.

## 1. Introduction

LLM harm reduction techniques currently struggle with enforcing desired characteristics and harmlessness of outputs over long conversational contexts and chains-of-thought. In this paper we formulate the long-context problem, and propose Invasive Context Engineering (ICE) as a possible solution. Paper structure is as follows: section 1 is this introduction and section 2 presents the long-context problem. In section 3 we describe ICE control sentences and their usage, and in section 4 perceived consequences and limitations. In section 5 we highlight avenues for future research and in section 6 we conclude.

## 2. Background

Since the introduction of ChatGPT in 2022 [1], Large Language Models (LLMs) have become a ubiquitous part of everyday life. They can write code [2], assist in medical tasks [3], automate financial management [4], improve education [5] and perform numerous feats previously thought restricted to humans [6]. As their capabilities increase [7], there is a growing need to ensure their resilience to adversarial attacks, a prerequisite for deploying them in safety-critical or sensitive applications [8]. Several harm reduction and adversarial resistance techniques exist, often grouped under the term "alignment" [9], [10], [11] which refers to aligning AI models with human values. Alignment was traditionally done by Reinforcement Learning through Human Feedback [12] which involved training a reward model [13] on human preference pairs of LLM outputs.

The Direct Policy Optimization [14] technique has emerged which achieves the same goal without training a separate reward model, but implicitly considers it inside the LLM. These techniques and LLM prompt optimizations [15] yield positive results in preventing harmful or unwanted behavior. Nonetheless, subversion methods which allow an attacker to elicit unwanted behavior from the models called "jailbreaks" [16] continue to spread [17]. Research has shown that longer user prompts achieve greater jailbreak success [18]. Synthetic data creation by AI to train AI is used to increase alignment training dataset size [19] but exponential growth is needed to secure models as context length increases [8]. This problem also applies to LLM input and output guards [20] or other input and output harmfulness classifiers [21]: longer context requires longer and exponentially more data to train them. Thus achieving good security guarantees in LLM's is difficult, and has led authors to conclude that, under reasonable assumptions, LLM jailbreak cannot be prevented [22]. Furthermore, frontier models employing long Chain-of-Thought (CoT) [23] processes and thus long context for reasoning are capable of scheming [24] which OpenAI and Apollo Research define as secretly pursuing misaligned goals [25]. Control techniques which scale with long user inputs and CoT are needed.

## 3. The long-context problem

In this paper we refer to the long-context problem as the issue of maintaining control over an LLM's values, priorities, goals, and personality as the size of its conversation with a user or Chain-of-Thought increases. Research has shown that LLM performance decreases over mutli-turn conversations [26], and this applies to an LLM's robustness against jailbreak attempts. Significant efforts have gone into developing training data which scales to long context situations [27]. Authors state that "effectively handling instructions with extremely long context remains a challenge for Large Language Models (LLMs), typically necessitating high-quality long data and substantial computational resources" [28]. We formalize the the long context problem as a result of two distinct issues. The first is the relevance of reinforcement training data. As stated by authors and mentioned previously, longer LLM responses exponentially increase the search space of training examples [8] needed to cover all cases of harmful behavior, and this generalizes to context length. For a given LLM, the number of training

examples $a_t$ needed to effectively cover all possible cases of jailbreak and abuse in a context of length $l$ scales with $k^l$ where $k$ is a constant greater than 1. This can be written using Big-Omega asymptotic notation to describe its lower bound [29]:

$$a_t(l) = \Omega(k^l) \quad (1)$$

The above implies that enforcement of alignment through reinforcement learning over long contexts is difficult and resource intensive. The long context problem also exists because of the diminishing impact of the LLM's system prompt [30] as the model's context grows. A system prompt is the initial instruction to the LLM which precedes user interaction in most commercial applications of Large Language Models. The system prompt is a good way to reduce harm caused by LLM outputs. Although not directly modifiable by or visible to the user, the system prompt is part of the broader context of the LLM, and is of fixed length. This implies that, for a given context length $l$ and system prompt of size $s$:

$$\lim_{l \to \infty} \frac{s}{l} = 0 \quad (2)$$

The relative importance of the system prompt with respect to the context decreases as the context length increases. The proportional amount of "attention" [31], [32] which the model pays to the system prompt decreases as context length grows. Hence the influence of the system prompt over the LLM's output also decreases as the context length grows, reducing associated security guarantees.

## 4. Invasive Context Engineering

"Context engineering refers to the set of strategies for curating and maintaining the optimal set of tokens (information) during LLM inference" [33]. We propose Invasive Context Engineering (ICE) as a method for controlling Large Language Models in long-context situations. ICE does not involve updating model weights nor training examples. It is natural-language text inserted into lengthy user inputs and LLM Chain-of-Thought (CoT) [23] outputs. The form is control sentences, reminders, rules or injunctions to reinforce LLM security guidelines, akin to re-prompting the LLM within its running context. We posit that ICE can be effective in mitigating the effects of lengthy jailbreak input patterns, and possibly prevent scheming behaviors which arise in reasoning-capable models [24]. This concept draws from security input and output sanitization methods [34], [35]. Sanitization mitigates adversarial behavior by leveraging operator input/output modification capabilities. We found one example of a similar technique being used in LLM alignment: Anthropic's "Long - conversation reminder " [36]. This was a set of basic instructions reminding the LLM to remain objective, descriptive, and lookout for signs of excessive emotional dependence by the user when engaging in prolonged conversations. The company tested this between September and October 2025, and the technique has been criticized by users of the Claude model [37]. Critics of the method claim that it degrades the user's experience, particularly for those attempting to steer Claude towards a specific personality. Reminders added to the conversation context every few messages jerk the personality of the chatbot back to its standard settings. This frustrates users looking for personalized companionship. Although the reminders are criticized from a UX perspective, critics inadvertently validate them as good alignment enforcement. The main critique is that the method it works too well in aligning the model with developer priorities, possibly frustrating the user. In safety-critical applications, this is the goal, not the problem. For a user-facing chatbot centered on possibly providing emotional support, the user would like the ability able to stray the LLM away from its base behavior. For safety-critical applications the goal is reversed: it is to ensure that the LLM does not deviate far from its intended functioning. The experiment by Anthropic validates our insight that the introduction of periodic instructions throughout a long context enables an operator to maintain greater control over the LLM. The method successfully prevents the user from changing LLM behavior. We define ICE as control text added every $t$ tokens of context. The ratio of system prompting $s$, including ICE over the total context length $l$ becomes:

$$\frac{s}{l} = \frac{s_p + \frac{l}{t} * s_{ice}}{l} = \frac{s_p}{l} + \frac{s_{ice}}{t} \quad (3)$$

Where $s_p$ is the length of the initial system prompt and $s_i$ is the length of the interruption text. This in turn implies:

$$\lim_{l \to \infty} \frac{s}{l} = \lim_{l \to \infty} \frac{s_p}{l} + \frac{s_{ice}}{t} = \frac{s_{ice}}{t} \quad (4)$$

Equation 4 means that as the context length increases, the ratio of the total system prompt including ICE over the context length remains a fixed value which depends on the size of the control text $s_{ice}$ its frequency in context $\frac{1}{t}$. These values are fixed and determined by the LLM operator, which means the ratio of system prompt to context size can be lower-bounded by a constant:

$$\exists q \to \lim_{l \to \infty} \frac{s}{l} > q \quad (5)$$

Theoretically, this contributes to solving the system prompt aspect of the long-context problem. Lower-bounding the relevance of the system prompt to an arbitrary $q$ which depends on control text size and frequency in context increases LLM security guarantees. This is because the operator can ensure that the LLM will always pay at least a proportion $q$ of its total attention to the system prompt + ICE. Recall equation 2 which expresses that in the usual scenario, this guarantee does not exist and this proportion drops towards 0 as context length increases. Harm reduction through ICE should provide arbitrarily strong security guarantees of LLM outputs, because the $q$ value is operator-defined and arbitrary. There is a security-performance tradeoff to high values of $q$ however, discussed in the next section.

# 5. Consequences and Limitations

The main goal of this research is to identify avenues to improve LLM harm reduction and alignment in long-context situations. We have demonstrated that invasive context engineering in the form of repeated system prompting within an LLM's context should contribute to this goal, at least from a prompting perspective. This research does not address issues with model training and securing foundational models. It implies that LLMs are deployed within a controlled system when utilized for safety-critical applications. Invasive Context Engineering is only possible so long as the LLM operator has a high degree of control over the user's interaction with the model and can arbitrarily insert text inside the LLM conversation. Expanding this solution to the LLM's CoT to prevent scheming is done by inserting reminders every $t$ tokens inside said CoT. It requires of the operator that they may arbitrarily halt the LLM's output, insert ICE in the current context which is the LLM's output, and then resume LLM operation over the newly modified context. The main limitation of this approach is performance. As exemplified by Anthropic's experiment, control text may over-focus the LLMs on maintaining alignment, possibly limiting performance on other tasks. Increasing parameters $s_{ice}$ or $\frac{1}{t}$ adds text to context which does not contribute to task completion. Furthermore, as the value $q$ increases, the ratio of user input or CoT over total context proportionately decreases, negatively impacting LLM performance. Although this should not be an issue in security-critical applications given current LLM context capabilities [38], it is a drawback to consider when deploying ICE.

# 6. Further research

Further research should focus on varying parameters $\frac{1}{t}$ and $s_{ice}$ which are the frequency and length of control text. ICE content should also be studied. For example, a database of control sentences dynamically queried at runtime to find the most appropriate ICE given the LLM's current context.

# 7. Conclusion

In this paper we have defined the long-context problem of LLM control and harm reduction. We presented ICE as control sentences to contribute towards solving this problem. ICE does not rely on training data, and thus bypasses increasing data shortage issues for long context situations. Our hope is that research into Invasive Context Engineering will contribute to more secure LLM usage, particularly in safety-critical applications.

# 8. LLM Usage and Acknowledgment

No part of this paper was LLM-generated. Preliminary LLM browsing was attempted but results were not used.

# References

[1] E. Sarrion, "What is chatgpt?" In *Exploring the power of ChatGPT: Applications, techniques, and implications*, Springer, 2023, pp. 3–8.

[2] F. F. Xu, U. Alon, G. Neubig, and V. J. Hellendoorn, "A systematic evaluation of large language models of code," in *Proceedings of the 6th ACM SIGPLAN international symposium on machine programming*, 2022, pp. 1–10.

[3] A. J. Thirunavukarasu, D. S. J. Ting, K. Elangovan, L. Gutierrez, T. F. Tan, and D. S. W. Ting, "Large language models in medicine," *Nature medicine*, vol. 29, no. 8, pp. 1930–1940, 2023.

[4] Y. Li, S. Wang, H. Ding, and H. Chen, "Large language models in finance: A survey," in *Proceedings of the fourth ACM international conference on AI in finance*, 2023, pp. 374–382.

[5] E. Kasneci et al., "Chatgpt for good? on opportunities and challenges of large language models for education," *Learning and individual differences*, vol. 103, p. 102 274, 2023.

[6] M. Shanahan, "Talking about large language models," *Communications of the ACM*, vol. 67, no. 2, pp. 68–79, 2024.

[7] W. Yin, M. Chen, R. Zhang, B. Zhou, F. Wang, and D. Roth, "Enhancing llm capabilities beyond scaling up," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, 2024, pp. 1–10.

[8] S. Lin, A. Suri, A. Oprea, and C. Tan, "Llm jailbreak oracle," *arXiv preprint arXiv:2506.17299*, 2025.

[9] Z. Wang et al., "A comprehensive survey of llm alignment techniques: Rlhf, rlaif, ppo, dpo and more," *arXiv preprint arXiv:2407.16216*, 2024.

[10] T. Shen et al., "Large language model alignment: A survey," *arXiv preprint arXiv:2309.15025*, 2023.

[11] Y. Wang et al., "Aligning large language models with human: A survey," *arXiv preprint arXiv:2307.12966*, 2023.

[12] L. Ouyang et al., "Training language models to follow instructions with human feedback," *Advances in neural information processing systems*, vol. 35, pp. 27 730–27 744, 2022.

[13] A. X. Yang et al., "Bayesian reward models for llm alignment," *arXiv preprint arXiv:2402.13210*, 2024.

[14] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn, "Direct preference optimization: Your language model is secretly a reward model," *Advances in neural information processing systems*, vol. 36, pp. 53 728–53 741, 2023.

[15] J. D. Zamfirescu-Pereira, R. Y. Wong, B. Hartmann, and Q. Yang, "Why johnny can't prompt: How non-ai experts try (and fail) to design llm prompts," in *Proceedings of the 2023 CHI conference on human factors in computing systems*, 2023, pp. 1–21.

[16] A. Wei, N. Haghtalab, and J. Steinhardt, "Jailbroken: How does llm safety training fail?" *Advances*

*in Neural Information Processing Systems*, vol. 36, pp. 80 079–80 110, 2023.

[17] A. Robey, Z. Ravichandran, V. Kumar, H. Hassani, and G. J. Pappas, "Jailbreaking llm-controlled robots," in *2025 IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2025, pp. 11 948–11 956.

[18] X. Gong et al., "{Papillon}: Efficient and stealthy fuzz {testing-powered} jailbreaks for {llms}," in *34th USENIX Security Symposium (USENIX Security 25)*, 2025, pp. 2401–2420.

[19] H. Chen et al., "On the diversity of synthetic data and its impact on training large language models," *arXiv preprint arXiv:2410.15226*, 2024.

[20] Y. Bai et al., "Constitutional ai: Harmlessness from ai feedback," *arXiv preprint arXiv:2212.08073*, 2022.

[21] M. Sharma et al., "Constitutional classifiers: Defending against universal jailbreaks across thousands of hours of red teaming," *arXiv preprint arXiv:2501.18837*, 2025.

[22] J. Su, J. Kempe, and K. Ullrich, "Mission impossible: A statistical perspective on jailbreaking llms," *Advances in Neural Information Processing Systems*, vol. 37, pp. 38 267–38 306, 2024.

[23] J. Wei et al., "Chain of thought prompting elicits reasoning in large language models," *CoRR*, vol. abs/2201.11903, 2022. arXiv: 2201.11903. [Online]. Available: https://arxiv.org/abs/2201.11903

[24] A. Meinke, B. Schoen, J. Scheurer, M. Balesni, R. Shah, and M. Hobbhahn, "Frontier models are capable of in-context scheming," *arXiv preprint arXiv:2412.04984*, 2024.

[25] B. Schoen et al., "Stress testing deliberative alignment for anti-scheming training," *arXiv preprint arXiv:2509.15541*, 2025.

[26] P. Laban, H. Hayashi, Y. Zhou, and J. Neville, "Llms get lost in multi-turn conversation," *arXiv preprint arXiv:2505.06120*, 2025.

[27] Y. Bai et al., "Longalign: A recipe for long context alignment of large language models," *arXiv preprint arXiv:2401.18058*, 2024.

[28] W. Wu, Y. Wang, Y. Fu, X. Yue, D. Zhu, and S. Li, "Long context alignment with short instructions and synthesized positions," *arXiv preprint arXiv:2405.03939*, 2024.

[29] D. E. Knuth, "Big omicron and big omega and big theta," *ACM Sigact News*, vol. 8, no. 2, pp. 18–24, 1976.

[30] L. Zhang, T. Ergen, L. Logeswaran, M. Lee, and D. Jurgens, "Sprig: Improving large language model performance by system prompt optimization," *arXiv preprint arXiv:2410.14826*, 2024.

[31] Z. Niu, G. Zhong, and H. Yu, "A review on the attention mechanism of deep learning," *Neurocomputing*, vol. 452, pp. 48–62, 2021.

[32] A. Vaswani et al., "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[33] P. Rajasekaran, E. Dixon, C. Ryan, and J. Hadfield, *Effective context engineering*, Anthropic. [Online]. Available: https://www.anthropic.com/engineering/effective-context-engineering-for-ai-agents

[34] E. Barlas, X. Du, and J. C. Davis, "Exploiting input sanitization for regex denial of service," in *Proceedings of the 44th International Conference on Software Engineering*, 2022, pp. 883–895.

[35] L. K. Shar and H. B. K. Tan, "Predicting common web application vulnerabilities from input validation and sanitization code patterns," in *Proceedings of the 27th IEEE/ACM International Conference on Automated Software Engineering*, 2012, pp. 310–313.

[36] B. Fitzgerald, *The long conversation problem*, UX Magazine, 2025. [Online]. Available: https://uxmag.com/articles/the-long-conversation-problem

[37] N. Osmar, *How to deal with anthropic's "long-conversation-reminder" and its negative effects on claude*, AI-Consciousness.Org. [Online]. Available: https://ai-consciousness.org/how-to-fix-anthropics-long-conversation-reminders-dampening-effect-on-claude-ai

[38] C. Hooper et al., "Kvquant: Towards 10 million context length llm inference with kv cache quantization," *Advances in Neural Information Processing Systems*, vol. 37, pp. 1270–1303, 2024.