

Specified Backup for Fragile Parts of LLMs

Abhi Morumpalle

Allen Zhang

Arnav Marda

Jeffrey Kwan

Harry Qian

Team A5

Abstract

In the past few years, there has been an explosion of interest in Large Language Models (LLMs) for a variety of practical applications. Much of this explosion has been driven by the invention of the Transformer architecture. However, the Transformer architecture inner workings largely remain a mystery. Combining this with the applications that LLMs are finding in the real-world, there a variety of new security risks that these LLMs open up their users to. In this paper, we analyze the robustness of LLMs to random bitflips in the variables, pinpointing specific parts of the LLM that are vulnerable to these hardware errors.

1 Introduction

In the past few years, there has been an explosion of interest in LLMs with the creation of widely available resources like OpenAI's ChatGPT and Meta's open source Llama. Much of the explosion has been driven by the creation of the transformer architecture, which has made a dramatic difference throughout AI, but particularly in the world of LLMs. However, our fundamental understanding of how these objects remains shrouded in mystery.

Because of our lack of understanding of how these objects work, and the quick assimilation of these products into our daily lives, there are a variety of novel security risks that we are being introduced to. One specific error is not so common, but still of practical relevance, is a hardware failure in which a random bit-flip occurs in the parameters of our model. An error of

such a fashion could have drastic effects on our output, ranging from making the outputs gibberish to outright wrong.

In this paper, we analyze how injecting bit errors into specific locations of a transformer and the LLM model as a whole affect the output. To this end, we use the {INSERT MODEL} as a base model to test on. To inject errors into our base model we use the PyTEI package [Ma et al., 2023]. To quantify the effect of our errors, we compare the score of our base model to the score of the models with errors injected using the {INSERT TESTS}.

2 First Level Heading

First level headings are all flush left, initial caps, bold and in point size 12. One line space before the first level heading and 1/2 line space after the first level heading.

2.1 Second Level Heading

Second level headings must be flush left, initial caps, bold and in point size 10. One line space before the second level heading and 1/2 line space after the second level heading.

2.1.1 Third Level Heading

Third level headings must be flush left, initial caps and bold. One line space before the third level heading and 1/2 line space after the third level heading.

Fourth Level Heading

Fourth level headings must be flush left, initial caps and roman type. One line space before the fourth level heading and 1/2 line space after the fourth level heading.

2.2 Citations In Text

Citations within the text should indicate the author's last name and year[Knu73]. Reference style[?] should follow the style that you are used to using, as long as the citation style is consistent.

Copyright © by the paper's authors. Copying permitted for private and academic purposes.

In: M. Meder, A. Rapp, T. Plumbaum, and F. Hopfgartner (eds.): Proceedings of the Data-Driven Gamification Design Workshop, Tampere, Finland, 20-September-2017, published at <http://ceur-ws.org>

2.2.1 Footnotes

Indicate footnotes with a number¹ in the text. Place the footnotes at the bottom of the page they appear on. Precede the footnote with a vertical rule of 2 inches (12 picas).

2.2.2 Figures

All artwork must be centered, neat, clean and legible. Do not use pencil or hand-drawn artwork. Figure number and caption always appear after the the figure. Place one line space before the figure, one line space before the figure caption and one line space after the figure caption. The figure caption is initial caps and each figure is numbered consecutively.

Make sure that the figure caption does not get separated from the figure. Leave extra white space at the bottom of the page to avoid splitting the figure and figure caption.

Figure 1 shows how to include a figure as encapsulated postscript. The source of the figure is in file fig1.eps.

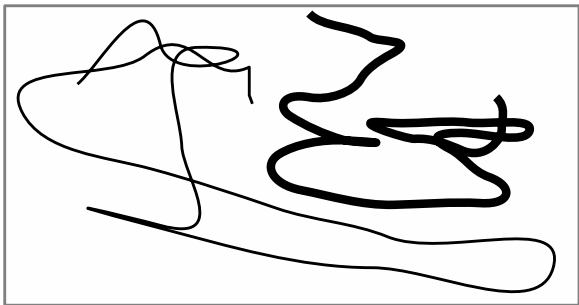


Figure 1: Sample EPS figure

Below is another figure using LaTeX commands.

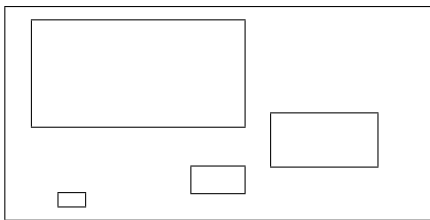


Figure 2: Sample Figure Caption

2.2.3 Tables

All tables must be centered, neat, clean and legible. Do not use pencil or hand-drawn tables. Table number and title always appear before the table.

One line space before the table title, one line space after the table title and one line space after the table.

The table title must be initial caps and each table numbered consecutively.

Table 1: Sample Table

A	B	1
C	D	2
E	F	3

2.2.4 Handling References

Use a first level heading for the references. References follow the acknowledgements.

2.2.5 Acknowledgements

Use a third level heading for the acknowledgements. All acknowledgements go at the end of the paper.

References

[Ma et al., 2023] D. Ma, X. Jiao, F. Lin, M. Zhang, A. Desmaison, T. Sellinger, D. Moore, S. Sankar. Evaluating and Enhancing Robustness of Deep Recommendation Systems Against Hardware Errors

[Knu73] D. E. Knuth. *The Art of Computer Programming – Volume 3 / Sorting and Searching*. Addison-Wesley, 1973.

¹This is a sample footnote