

Specified Backup for Fragile Parts of LLM

Abhi Morumpalle, Allen Zhang Arnav Marda, Jeffrey Kwan, Harry Qian
UCLA

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 [7]. However, the Transformer architecture inner workings largely remain a mystery. Combining this with the applications that LLMs are finding in the real-world, there is a variety of new security risks that these LLMs open up for their users. 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. We also evaluate the effect that value errors have on the efficacy of specific components of the LLM when evaluated against various LLM benchmarks.

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 remain 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 bitflip 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 Mistral-7B [4] as a base model to test on. To inject errors into our base model we use an extension of the PyTEI package introduced in [5]. 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 PyTEI package with the MMLU framework supplied by DeepEval to evaluate.

Our basic workflow is outlined in the below diagram

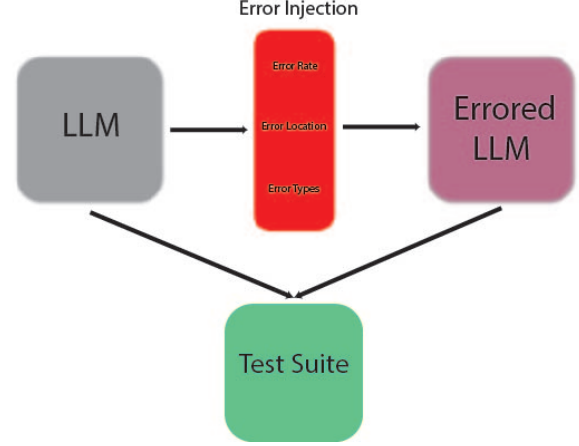


Figure 1: General evaluation workflow for LLM

2 BACKGROUND AND RELATED WORK

There is a related paper [5] that analyzes a model known as a recommendation system. In their work, they build the PyTEI package for injecting models. However, recommendation models differ significantly from LLMs so the overall effect could be quite different for the same error injections.

However, the paper does not do any evaluation on the different of hardware errors in specific parts of the recommendation system, so the question of whether particular parts are more vulnerable is still open.

An interesting part of this paper is the evaluation of possible mitigations against hardware flips, which can also be evaluated in the context of LLMs. In addition, their evaluation had limited scope, and it could be interesting to expand on their analysis by testing a wider variety of errors and examining the tradeoffs of each.

[2] performed resilience analysis on transient errors in logic components, but their framework is highly inefficient in the context of PyTorch, because they employ frequent to-and-back type conversions as PyTorch doesn't natively support bit operations in float tensors. [6] estimated DNN accuracy under transient faults and proposed a Monte Carlo-based estimation method. However, their analysis is also limited to transient errors and do not consider logic / data path errors which are permanent and more impactful.

3 THREAT MODEL.

We present 2 case studies, relating to information disclosure and denial of service respectively.

An attacker's goal is to cause information disclosure or cause a denial of service via system crash or significantly degrading the performance of the LLM. They may attack components such as edge AI chips and accelerators, CPUs and GPUs in LLM training,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2024 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

memory systems storing model weights, and data transfer channels between hardware components. To induce hardware errors, they could somehow achieve write access to memory systems, or physical access to create fault-heavy environments (e.g. radiation source).

3.1 Case Studies

case 1 (information disclosure): An LLM is running on an edge device (like a smartwatch). The attacker jailbreaks the device, gaining write access to hardware, and inject bit errors with a probability p . The attacker then monitors the LLM output through API calls, observing that the output is NaN XX% of the time. The attacker thus concludes that the LLM has YY number of parameters (see 6.1 for more details).

case 2 (denial of service): A physics research lab is using LLMs for running simulations of their experiments. An attacker somehow manages to convince the lab scientists to house the servers the LLM is running on in their reactor chamber, a source with high radiation. This environment causes many bit flip errors, causing frequent faulty outputs from the LLM.

For the purposes of this paper, we will investigate the effects of the denial of service attack objective, specifically performance degradation.

jeffrey todo: make a figure of the threat model

4 OVERVIEW OF THE DESIGN

We chose Hugging Face’s implementations of Mistral-7B [4] as our model of choice. We evaluated our models on [3], an open-source LLM benchmark, specifically the computer science and astronomy tests that have the injected LLM answer multiple choice questions. For each model with varying error rates, the score is computed as the proportions of correct answers. abhi todo: details about layers and the score metric

4.1 Design Choices

We originally chose GPT2 but noticed a lot of NaN outputs. See 6.1 for analysis. Hence, we switched to an error model that injects *value* errors instead.

We also observed that GPT2 performance on the MMLU benchmark was equivalent to random guessing, which would not be informative of performance drops. We upgraded to Mistral-7B which has a 0.6 average accuracy.

5 METHODOLOGY/DESIGN

5.1 Error Injection

In order to simulate bit flips, we load the model parameters, inject errors into them, save the injected weights, and load the new weights into a perturbed model.

To achieve this, we use PyTEI, which provides a flexible API that can generate bit error maps for each part of the model parameters. First, we generate bit masks for all parameters we wish to inject into, with one boolean value for each bit in the parameter when

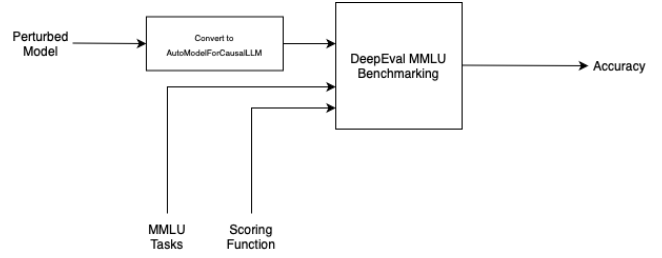


Figure 2: Evaluation Pipeline

stored as IEEE floats. Then, we specify the probability of a simulated SDC error occurring in the affected parameters, and using PyTorch’s implementation of DropOut to efficiently zero out most bits according to the error rate. The mask is then XOR-ed with the stored parameters, where each 1 bit in the mask will trigger a bit flip.

In later experiments, we extended the functionality of PyTEI by creating an API for injecting value errors. A value error occurs when the entire float is replaced with a pre-specified incorrect value. Each float in the parameter tensor is still corrupted independently with a specified probability.

5.2 Evaluation

5.2.1 MMLU Benchmark. We evaluate our perturbed models on the Massive Multitask Language Understanding (MMLU) Benchmark. The MMLU Benchmark contains 15908 questions over 57 subjects. The difficulty of the questions ranges from an elementary level to an advanced professional level[3].

For our particular evaluation framework, we narrowed the scope of our evaluation to the High School Computer Science and Astronomy tasks. The scope was narrowed from 57 tasks to 2 tasks due to the excessive computational costs required to run 15908 instances of inference on each perturbed model. These particular tasks were chosen since they provide around average performance (3) without exceeding the input token limit for GPT-2 and Mistral 7B.

The evaluation was conducted in a 1-shot learning setup to isolate the effect of the perturbation on the model’s intrinsic knowledge retention and generalization capabilities from few-shot learning.

5.2.2 DeepEval Testing Framework. In order to standardize the evaluation of the perturbed models, we used an existing library in Python called DeepEval, an open-source toolkit for LLM evaluation. The framework consists of an inbuilt class to run evaluation on particular tasks of the MMLU benchmark.

The MMLU Benchmarking class takes as input a model, and a set of tasks to run. It then uses the following algorithm to evaluate the model on the benchmark:

In 1, the function `format_input` uses a specific template to format the input to instruct the LLM to correctly answer the question in a specific format. Also, `model.generate()` refers to a generic function to generate output tokens from an LLM given some set of input tokens.

5.2.3 Exact Match Metric. For this evaluation framework and benchmark, we use an exact-match metric. This is a popular metric that outputs a binary result where the input is penalized if it is not

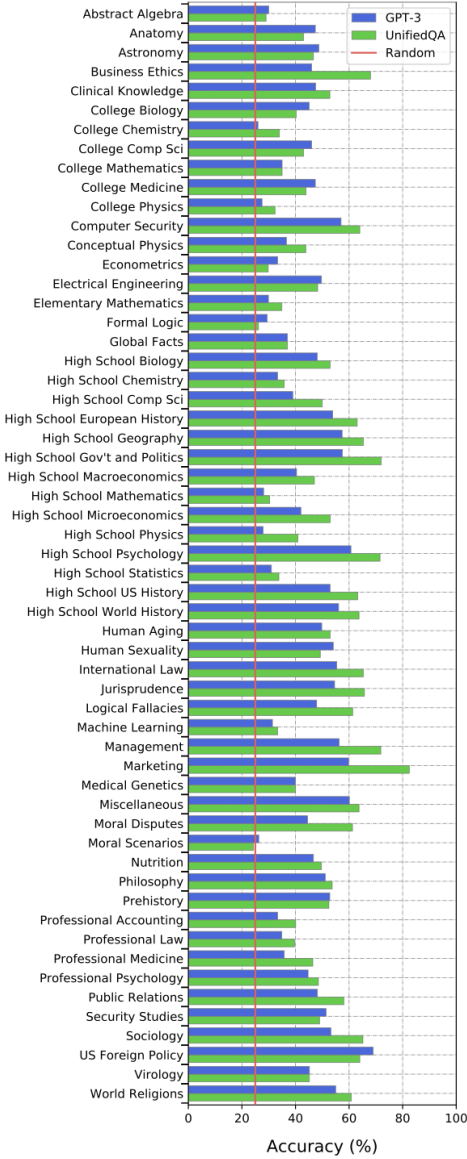


Figure 3: Performance of GPT-3, UnifiedQA and random guessing on all 57 MMLU tasks[3]

an exact match to the target. In this case, we compare the output token(s) to the target, and if it is not an exact match, we return false.

Input: model, tasks, score_fn

Output: accuracy

Function evaluate(model, tasks, score_fn):

```

total_qs ← 0;
total_correct ← 0;
for task in tasks do
    golden ← dataset[task];
    for g in golden do
        input ← format_input(g);
        result ← model.generate(input);
        score ← score_fn(result, g.answer);
        if score then
            total_correct ← total_correct + 1;
        end
        total_qs ← total_qs + 1;
    end
end
return total_correct / total_qs;

```

Algorithm 1: Evaluation algorithm for model evaluation on MMLU

5.3 Evaluation Framework

- (1) Download HuggingFace model to test on
- (2) Fix an error parameter p and specific labelled sections of the model to inject errors into
- (3) Inject value errors with each value having probability p of being an error, where our error was replacement with random value from (0, 1)
- (4) Put our error injected model into a pipeline that conforms to the testing framework (DeepEval MMLU) in our case
- (5) Run the tests and evaluate accuracy

5.3.1 HuggingFace Model (Mistral 7B). For the purpose of our experiment, we chose to use the model Mistral-7B [4] from Hugging-Face and ran it using the P100 GPU provided by Kaggle. Mistral-7B was chosen for its efficacy at relatively low memory cost. Our base-line score on our MMLU testing cases was 0.6, which is enough to be able to claim a significant difference when injecting errors. This is compared to GPT-2, in which even larger versions only perform at slightly better than random guessing.

A big benefit of using HuggingFace is that all layers are labelled appropriately, which makes injecting errors into specific components much easier, as we will explain later. Mistral-7B's core is made up of 32 transformer layers, each with the usual self attention, normalization, etc components.

As an aside, to be able to run on the given GPU's memory, we downloaded the model with float16 quantization. We considered further quantization but the error injection became much more complex in this case.

5.3.2 Error Parameters. There are two parameters that we tried changing. An error parameter p and the set of components that we wanted to inject errors into. The error parameter p is the probability that a random value p changes into a random value from (0, 1). This process is run on every component that we want to inject errors into.

The reason we chose a random value from $(0, 1)$ is multi-fold. For one, for this specific Mistral-7B model all weights are between $(0, 1)$ so we were not losing too much by restricting to this. The second, more important reason is that by flipping random bits, the float16 tends to explode to an extremely large value or have NaN. NaN meant that we received errors, whereas large float16 tended to overflow and also throw an error in our code. For this reason, in order to get any results at all, we decided to stick to replacing our small weights with other small weights.

The specific error parameters that we tested were to split up all 32 layers of Mistral-7B into sublayers of size 4 (so 1-4, 5-8, etc) and inject errors into those specific sublayers. We varied p logarithmically from 10^{-9} to 10^{-1} . We also looked into injecting errors into components of each layer, so we tried injecting into all self-attention layers, normalization layers, etc with the same varied values of p .

6 EVALUATION

6.1 NaN Analysis

In this section, we provide an analysis of the probability that the LLM outputs NaN.

Let F be an LLM model with n parameters and x be the input to the model. We would like to find $Pr(F(x) = NaN)$.

Suppose a bit is flipped with probability p .

Let X be the value of a model parameter and X^* be its perturbed value. Note that parameters are iid $\sim Uniform(0, 1)$ due to the design of LLMs [do you have a good source for this allen]. According to [1], a NaN value has all 1s in the exponent field, and a nonzero mantissa. Since $0 \leq X \leq 1$, the exponent field is 01111111 (due to the +127 offset), and suppose there are j 1s in the mantissa.

Then, $X^* = NaN$ if the remaining exponent bit flipped, the rest of the exponent bits don't flip, and not all j mantissa bits turn to 0. Thus,

$$\begin{aligned} Pr(X^* = NaN) &= Pr(\text{only one exponent bit is flipped}) \\ &\quad \cdot (1 - Pr(\text{all mantissa bits are 0})) \\ &= p(1-p)^7 * (1 - p^j(1-p)^{23-j}) \end{aligned}$$

We can represent the value of mantissa as a random variable $\sim Binom(23, 0.5)$. Thus,

$$\begin{aligned} Pr(X^* = NaN) &= \sum_j^{23} Pr(X^* = NaN \mid X \text{ has } j \text{ mantissa 1s}) \\ &\quad \cdot Pr(X \text{ has } j \text{ mantissa 1s}) \\ &= \sum_j^{23} p(1-p)^7 * (1 - p^j(1-p)^{23-j}) \binom{23}{j} 0.5^j 0.5^{23-j} \\ &= 2^{-23} p(1-p)^7 \left(\sum_{j=0}^{23} [1 - p^j(1-p)^{23-j}] \binom{23}{j} \right) \\ &= 2^{-23} p(1-p)^7 \left(\sum_{j=0}^{23} \binom{23}{j} - \sum_{j=0}^{23} p^j(1-p)^{23-j} \binom{23}{j} \right) \\ &= 2^{-23} p(1-p)^7 (2^{23} - 1) \\ &= p(1-p)^7 (1 - 2^{-23}) \end{aligned}$$

Now putting everything together, we note that the model will output NaN if any of the parameters of the model are NaN, due to NaN propagation. Thus,

$$\begin{aligned} Pr(F(x) = NaN) &= Pr(\text{at least one parameter is NaN}) \\ &= 1 - Pr(\text{all parameters are not NaN}) \\ &= 1 - \prod_{i=1}^n Pr(X_i \neq NaN) \\ &= 1 - Pr(X^* \neq NaN)^n \\ &= 1 - (1 - Pr(X^* = NaN))^n \\ &= 1 - (1 - p(1-p)^7 (1 - 2^{-23}))^n \end{aligned}$$

For $p = 10^{-9}$ and $n \approx 130 * 10^6$ (GPT-2), we get $Pr(F(x) = NaN) \approx 0.1$, and for $p = 10^{-8}$, we get $Pr(F(x) = NaN) \approx 0.7$, which explains why at a low error rate, the model still outputs NaN quite consistently.

6.2 Variation by Level

6.2.1 Probability Variation. We began by analyzing, at each level, how variations in probability effect the evaluation on the MMLU benchmark at various probability from $1e-7$ to $1e-3$.

As mentioned before, for computational reasons we grouped up layers into groups of size 4. For ease of viewing this result, we only present the value for layer groups 1, 2, 7 and 8. A similar pattern holds for all unshown groups.

Below are our results for each layer group with various probabilities.

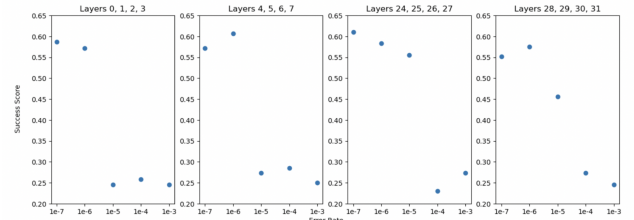


Figure 4: Score when Varying Probability

As one can notice, for probabilities $1e-6$ (this is about 200 total errors for each layer group) there is no noticeable negative effect that our injected errors have on the model.

On the other hand, for probabilities $1e-4$ (about 20000 total errors), the output of our LLM for each layer that we have is essentially no better than random guessing, no matter which layer we have injected into.

The only layer where something interesting seems to be happening is at the $1e-5$ parameter, so we look into this further in the below section.

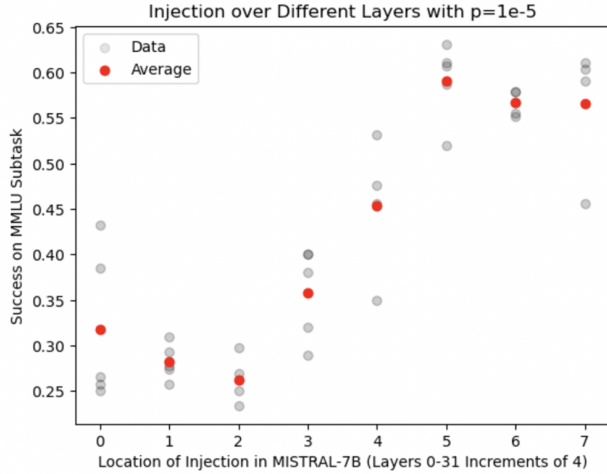


Figure 5: Score when Varying Layers with $p = 1e-5$

6.2.2 *Zooming in on $1e-5$.* There are a few interesting things to notice here.

The first is that there does seem to be an effect of the layer on the error injection. For example, errors later layers seem to have less of a degrading effect on the efficacy of the LLM compared to earlier layers. This is certainly an interesting result that seems to suggest that errors in earlier layers propagate to later layers.

Another interesting thing to note is that many there does seem to be a large amount of variation. We will explore this in a later section, but for now it suffices to say that this is because the variation in our error injection process means that some more vulnerable components within the individual layers also have a large effect on the total degradation of the LLM.

7 DISCUSSIONS

There are still many possible things that one could try. As a group, we were heavily restricted by our computational resources, so with more computation resources and more time it would be possible to try many more different things, as well as trying some more state of the art models with many more parameters than Mistral-7B.

It could be interesting to try inject errors into specific sections of each layer (such as attention or normalization), instead of just all the layers and see or all sections in a particular layer. We have confidence that this could give interesting results because of the high variation in some of our injection parameters, which suggests that where our random errors actually are injected into even with our given parameters is also important.

It would also be interesting to look at the performance degradation on various other benchmarks such as translation, coding capabilities, or reasoning tasks.

Our research is also limited in the sense that we looked at value injections. As a whole, silent bit errors seem to be completely catastrophic for some standard machine learning implementation. For example, in the usual case where the weights lie between 0 and 1, we’ve shown in 6.1 that NaN is likely for any large amount of bit errors. In float16 representation, the errors tend to cause overflow which broke the evaluation framework that we were working in. The final case is int4 quantization. Here, bit errors actually tend to be fine. In addition, random bit errors here are close to random value errors, which justifies our choice to look at random value errors instead.

8 CONCLUSION

The large takeaway is that all value errors are certainly not equal. It seems that errors in earlier levels tend to degrade the power of our LLM significantly more than errors at later levels. In addition, it seems that certain parts of the LLM tend to be more resistant to errors (such as the batch-normalization layer). In addition, we have shown that, in general, random bit errors in an LLM’s weights tend to be catastrophic in the sense they have a high probability of completely breaking the LLM.

The limitation in our work is in connecting these two cases. It’s unknown if our small value errors are a good model for bit errors that do not completely break the LLM. Hence, it’s unknown whether large variation errors will have a larger degrading effect than the ones we have injected. In addition, it is still unknown if our results on error degradation for specific layers still hold for the larger case of completely random bit errors.

More research is needed for all the above mentioned cases to draw stronger conclusions on the security of LLM’s to random bit flips.

9 CONTRIBUTIONS OF EACH TEAMMATE

Show who did what for the project, and who wrote what section for the report.

REFERENCES

- [1] 2019. IEEE Standard for Floating-Point Arithmetic. *IEEE Std 754-2019 (Revision of IEEE 754-2008)* (2019), 1–84. <https://doi.org/10.1109/IEEESTD.2019.8766229>
- [2] Yi He, Prasanna Balaprakash, and Yanjing Li. 2020. Fidelity: Efficient Resilience Analysis Framework for Deep Learning Accelerators. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. 270–281. <https://doi.org/10.1109/MICRO50266.2020.00033>
- [3] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring Massive Multitask Language Understanding. (2021). [arXiv:cs.CY/2009.03300](https://arxiv.org/abs/2009.03300) <https://arxiv.org/abs/2009.03300>
- [4] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. Mistral 7B. (2023). [arXiv:cs.CL/2310.06825](https://arxiv.org/abs/2310.06825) <https://arxiv.org/abs/2310.06825>
- [5] Dongning Ma, Xun Jiao, Fred Lin, Mengshi Zhang, Alban Desmaison, Thomas Sellinger, Daniel Moore, and Sriram Sankar. 2023. Evaluating and Enhancing Robustness of Deep Recommendation Systems Against Hardware Errors. (2023). [arXiv:cs.LR/2307.10244](https://arxiv.org/abs/2307.10244) <https://arxiv.org/abs/2307.10244>
- [6] Abhishek Tyagi, Yiming Gan, Shaoshan Liu, Bo Yu, Paul Whatmough, and Yuhao Zhu. 2024. Thales: Formulating and Estimating Architectural Vulnerability

Factors for DNN Accelerators. (2024). arXiv:cs.AR/2212.02649 <https://arxiv.org/abs/2212.02649>

- [7] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention Is All You Need. (2023). arXiv:cs.CL/1706.03762 <https://arxiv.org/abs/1706.03762>

A OVERFLOW FORM OTHER SECTIONS

Sometime you ware super excited about some details that does not quite fit with the rest of the paper goes here. For example, some details about how you instrumented the Android Linux kernel should go to appendix, and for really curious reader to read. Remember it's appendix, so the reader is not required to read, and you should not put critical information in appendix that is crucial for understanding the rest of the paper.