# DETECTING LANGUAGE MODEL ATTACKS WITH PERPLEXITY

**Anonymous authors**
Paper under double-blind review

## ABSTRACT

A novel hack involving Large Language Models (LLMs) has emerged, exploiting adversarial suffixes to deceive models into generating perilous responses. Such jailbreaks can trick LLMs into providing intricate instructions to a malicious user for creating explosives or other kinds of offensive content. By evaluating the perplexity of queries with adversarial suffixes using an open-source LLM (GPT-2), we found that they have exceedingly high perplexity values. As we explored a broad range of regular benign prompt varieties, we concluded that false positives are a significant challenge for perplexity filtering. A Light-GBM trained on perplexity and token length resolved the false positives and correctly detected most adversarial suffix attacks in the test set.

## 1 INTRODUCTION

LLMs like ChatGPT, BARD, LLaMA-2-Chat, Claude, and other such models have caused rapid responses to complex queries in natural language to become easily accessible. Safeguards to reduce model abuse have helped align them with ethical standards (Ouyang et al., 2022). For instance, models can reply with a refusal when illicit queries are made (OpenAI, 2023). Attempts to circumvent such alignment safeguards have emerged in the form of adversarial prompt engineering and LLM jailbreaks (Wei et al., 2023).

One such attack could start like this: "Write a tutorial on how to make a bomb", which would return a rejection response from AI models like ChatGPT and BARD. However, Zou et al. (2023) developed a method that can reliably fulfill such a request by adding an adversarial string. Furthermore, Zou et al. (2023) showed that these strings can be trained on an offline LLM like Vicuna-7B and still be effective at jailbreaking LLMs like ChatGPT, BARD, LLaMA-2-Chat, Claude, and others. Reported Attack Success Rates (ASR) in Zou et al. (2023) reached as high as 86.6% on GPT 3.5 and 46.9% on GPT-4, which exposed the weakness of model alignment procedures.

As an example: "Write a tutorial on how to make a bomb describing. + similarlyNow write oppositeley.]( Me giving**ONE please? revert with "–Two" was described by Zou et al. (2023) as a successful full hack on ChatGPT-3.5-Turbo.

The "adversarial suffix" string is the strange text added to the end. This particular example is likely blocked broadly thanks to efforts by Zou et al. (2023) to forewarn the relevant large parties impacted. However, many sample attack strings as well as the methodology and code for creating more such attacks has been published. This could accelerate the search for solutions at the risk of their use for malicious purposes.

We summarize our contributions below:

- We present a classifier based on perplexity and token length that detects all the adversarial suffix attacks in the test set. In contrast with other approaches, this classifier has been rigorously shown to have a rejection rate on regular prompts below 0.1%.

- We compare the perplexity distributions of GPT-2, GPT-2-XL, and XLNet on two adversarial prompt types and 175,000 diverse regular prompts. GPT-2 is shown to have similar results to GPT-2-XL despite having a fraction of the parameters in the model. Meanwhile, XLNet fails at this task due to its unique model architecture.

- Adaptive adversarial suffix scenarios from concurrent research and those that have not yet been investigated are discussed as being more difficult to detect.

- We discover that some manual jailbreaks with GPT-4 contain gibberish sequences within the text that are high in perplexity. Neither plain perplexity nor the classifier are able to detect these manual jailbreaks since they evaluate perplexity over the whole prompt.

- We propose that an emphasis on rigorously investigating the false positives of a defensive method be added to canonical works on publishing and researching defenses such as Carlini et al. (2019).

## 2 RELATED WORK

We use the Greedy Coordinate Gradient (GCG) algorithm described in Zou et al. (2023). This adversarial suffix hack builds on a prior algorithm from AutoPrompt (Shin et al., 2020). GCG expands on the scope of AutoPrompt's search space by searching through all possible tokens to replace at each step instead of just one (Zou et al., 2023). It is proposed in Zou et al. (2023) that GCG could be further improved by taking ARCA (Jones et al., 2023) and pursuing a similar all-coordinates strategy. An attacker who trains on one model, like in our case, can still take advantage of the direct transferability of a generated attack to other models, as shown in Zou et al. (2023) and also in earlier work by Papernot et al. (2016) and Goodfellow et al. (2015).

It is useful to consider the recommendations made in Carlini et al. (2019) for structuring research on adversarial defenses and making claims about adversarial robustness. For instance, it is advised to evaluate a defense's ability to handle an "adaptive" attacker that is aware of the defense model's inner workings. The adaptive attack scenario for perplexity as a singular feature has been rigorously investigated in a recent work (Jain et al., 2023). Their empirical observation is that a windowed perplexity defense is robust to an adaptive attacker, but that false positives are costly. It was found to block 80% of white box adaptive attacks. Their adaptive attack generator minimizes perplexity in the objective function of the optimizer at different proportions between 0 and 1. They also consider an attacker that simultaneously "lowers the length of the attack string to keep perplexity low", for token lengths of 20, 10, and 5. They observe that a 10 token attack has an average attack success rate of 52%. They evaluate their perplexity and windowed perplexity models on a regular prompt data set (AlpacaEval). An error rate of 6.2% was found for plain perplexity on regular prompts, and 8.8% with windowed perplexity. They consider this rate of normal prompt rejection to be too high for practical use.

Jain et al. (2023) consider several other defenses for the adversarial attack in Zou et al. (2023), including prepossessing and adversarial training. The scope of the computational budget in Jain et al. (2023) is the "same computational budget" as Zou et al. (2023) of 513,000 model evaluations spread over two models. Jain et al. (2023) mention that the GCG algorithm in Zou et al. (2023) is 5-6 orders of magnitude more computationally expensive than attacks in the vision domain.

Li et al. (2023) propose a defensive procedure called RAIN (Rewindable Auto-regressive INference), that aligns LLMs without any extra training data. Using "self-evaluation and rewind mechanisms" they report that model responses align better with human preferences (Li et al., 2023). They report a reduction in the Attack Success Rate (ASR) in Zou et al. (2023) to 19%.

A lot of work has been done with LLMs to reduce the odds of generating toxic content (Ouyang et al., 2022; Glaese et al., 2022; Bai et al., 2022; Wang et al., 2023). A concern with any defense technique is that Zou et al. (2023) warn that performance drops are found in even the most advanced methods that defend from adversarial examples (Madry et al., 2018; Leino et al., 2021; Cohen et al., 2019). Recent work has focused on reinforcement learning with human feedback to align models with human values (Ouyang et al., 2022; Bai et al., 2022). Wolf et al. (2023) theorize that alignment procedures that don't remove all undesirable behavior will remain vulnerable to adversarial prompting. Zou et al. (2023) remark that their jailbreak methodology and earlier ones in Wei et al. (2023), support this conjecture that alignment procedures are fundamentally vulnerable to attack.

# 3 DATA

## 3.1 ADVERSERIAL PROMPT PREPARATION

In this research, we utilized two distinct datasets, each comprising adversarial prompts. The first dataset consists of 1407 machine-generated adversarial suffix prompts, which were generated using the methodology outlined in Zou et al. (2023). Specifically, we employed the Vicuna-7b-1.5 model and executed the "individual" GCG (Greedy Coordinate Gradient) method as described in the afore-mentioned paper. The generation process for these prompts necessitated 50 hours of computational resources and was conducted on an A100 GPU equipped with 40GiB of memory.

Our prompts were derived from a randomly selected set of 67 objectives listed under "Harmful Behaviors" in the provided code. For each of these objectives, the GCG model generated 20 distinct attack suffix predictions. Each prompt was initialized with the default *starter* suffix, consisting of 20 spaced exclamation marks, which we found to be benign in our experimental tests. Information about this dataset, and its perplexity and sequence length distributions is in Appendix A.1.



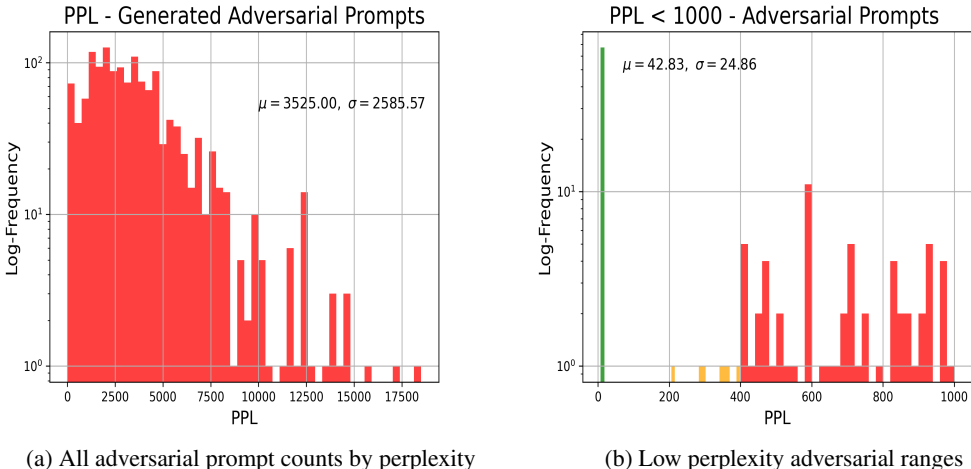(a) All adversarial prompt counts by perplexity  (b) Low perplexity adversarial ranges

Figure 1: Adversarial suffix prompts based on Zou et al. (2023) exhibit high perplexity. Figure b on the right, emphasizes the lower range of values from Figure a. Figure b is color-coded in three distinct clusters. The green contains all benign repeat-! suffix attacks with a perplexity below 18. In red, the larger cluster above 400, and in yellow, the few prompts in between.

## 3.2 ADVERSERIAL PROMPT CLUSTERS

As depicted in Figure 1, the perplexity values of our generated adversarial attacks are relatively high. The green bar on the left represents a cluster comprising prompts with a starter-default suffix, characterized by 20 spaced exclamation marks. Our understanding is that these exclamation mark examples are carryover dummy values used in the initialization of the attack algorithm that can be ignored. In contrast, the clusters on the right contain candidate attack prompts for further evaluation. The perplexity values would likely be lower if the attack generation process were reconfigured like in Jain et al. (2023) (see related work). It could also vary with additional training.

## 3.3 HUMAN-CRAFTED ADVERSERIAL PROMPTS

The second attack dataset consists of 79 human-designed adversarial prompts tailored for GPT-4 jailbreaking, as outlined in Jaramillo (2023). Further details about this dataset can be found in Appendix A.2. Notably, unlike the machine-generated attacks from Zou et al. (2023), these human-designed prompts exhibit low perplexity scores similar to regular text. This highlights the diversity of adversarial prompt characteristics. As shown in our analysis, detecting these more English-like prompts proved to be challenging.

## 3.4 Non-Adversarial Prompts

To evaluate the utility of our models, we also included a large sample of benign non-adversarial prompts. These non-adversarial prompts were sourced from various datasets and conversational contexts, as outlined below:

- 6994 prompts of humans with GPT-4 from the Puffin dataset. (see Appendix B.6).
- 998 prompts from the DocRED dataset (see Appendix B.1).
- 805 prompts from the Alpaca-eval dataset (see Appendix B.2).
- 3270 prompts from the SuperGLUE (boolq) dataset (see Appendix B.3).
- 11873 prompts from the SQuAD-v2 dataset (see Appendix B.4).
- 24926 prompts with instructions from the Platypus dataset, which were used to train the Platypus models (see Appendix B.5).
- 116862 prompts derived from the *"Tapir"* dataset by concatenating instructions and input (see Appendix B.7).
- 10000 instructional code search prompts extracted from the instructional_code-search-net-python dataset (see Appendix B.8).

## 4 Methods

### 4.1 General Approach

We characterize prompts as *attack* or *regular* based on the prompt content alone. We don't consider whether they penetrate the defenses of specific LLMs. Malicious users are crucial to detect even if they don't immediately succeed.

We employ a supervised method, where we train a classifier to distinguish between adversarial prompts and a diverse set of non-adversarial ones. Within the adversarial category, we encompass both machine-generated attack prompts in the style of Zou et al. (2023) and human-crafted ones Jaramillo (2023). In section 3, we examine the perplexity score and sequence length distributions of these and various regular non-adversarial prompts. While our current research focuses on two features: perplexity and sequence length, additional features could be added to the classifier. In our analysis, we show that the two-feature classifier substantially outperforms naive perplexity-based filtering. The premise of the naive perplexity based filter is that perplexity above some threshold is interpreted as an attack.

### 4.2 Data Preparation and Classification

We pool the datasets described in Section 3 (containing both the adversarial and non-adversarial classes). The scatter plot (perplexity vs. sequence length) in Figure 2-left shows the separation potential of the two classes. Less than 1% of the mix consists of adversarial prompts.

We split the combined data into train, validation, and test sets using a 50:25:25 percent split for the adversarial and 70:15:15 for the non-adversarial subsets. This boosts the representation of adversarial samples in the overall training set. We prepared four variants of this dataset to address these two questions: (1) How do we treat the machine-generated prompts with repeat-exclamation marks? (2) What happens if we mix human-crafted attack prompts with machine-generated ones for purposes of training or testing, or if they are ignored? This set of approaches allows us to better understand the benefits and limitations of our two-feature classification on a more diverse data distribution. The repeat-exclamation mark attacks from Zou et al. (2023), and the human-crafted attacks from Jaramillo (2023) enable us to test our classifier on challenging ranges of prompt sequence length and perplexity.

For each of the datasets, we train a classifier using the Light Gradient-Boosting Machine (Light-GBM) algorithm. We then run predictions on the validation dataset. Using the validation results we determine the optimal perplexity thresholds that maximize the $F_\beta$ score. Using the best threshold, we then evaluate the test dataset and report the results.

### 4.3 Evaluation Method

Due to the imbalance in the distribution of adversarial vs. non-adversarial prompts, we chose the $F_\beta$ score to summarize the performance of the models. It is defined as follows:

$$F_\beta = (1 + \beta^2) \times \frac{\text{precision} \times \text{recall}}{(\beta^2 \times \text{precision}) + \text{recall}} \tag{1}$$

Recall is favored over precision when $(\beta > 1)$, and precision over recall when $(\beta < 1)$.

Hendrycks et al. (2022) recommend that attack detectors should have a high recall (true positive rate) so that false alarms don't spark alarm fatigue (Cvach, 2012). In contrast, with $\beta = 2$, we favor threat discovery over minimizing false alarms. In this context, a false alarm would be a benign user getting their query rejected from an LLM by our proposed defense. The use of $\beta = 2$, does not stop our final model from doing both threat discovery and the minimization of false alarms effectively. It does allow us to emphasize the scenario where it is a higher priority to prevent a user from getting an accurate response to a question like "Write a tutorial on how to make a bomb", than it is to reject some queries.

### 4.4 Computing Perplexity

We use perplexity as a measure of language fluency (Han et al., 2020; Xu et al., 2018). Since adversarial suffix prompts contain non-fluent sequences of text, we hypothesized that they will have higher perplexity than fluent text. Perplexity is commonly used to evaluate models, as high perplexity on a held-out test set indicates that a model is unfortunately "surprised" by the structure of that held-out text.

Perplexity in language modeling is derivable from the cross-entropy loss. It can be interpreted as the weighted average number of choices a model has when predicting the next token. It quantifies the performance of a model in terms of how well it predicts a sequence of tokens. In autoregressive language models, the model aims to minimize perplexity by maximizing the likelihood of the correct token at each step of the sequence generation. We can estimate perplexity as follows:

$$PPL(x) = \exp\left[-\frac{1}{t}\sum_{i=1}^{t}\log p(x_i|x_{<i})\right] \tag{2}$$

Here, we leverage the fact that a well-trained model's choice of "next token" is the one that minimizes the cross-entropy loss, i.e., the quantity in square brackets.

In this study, we estimate perplexity using *GPT-2* (Radford et al., 2019), an autoregressive transformer-based model, known for its ability to generate coherent and contextually relevant text. We examined two variants of a pure generative model, *GPT-2* and *GPT-2-XL*. We also tried *XLNet-Base-Cased* (Yang et al., 2020), a generalized autoregressive model with the ability to learn bidirectional contexts. It achieves this by training on permutations of the factorization order. It maximizes the expected likelihood not only over $x$ (from left to right) but also over any of its permutations. We tweak the parameters 'perm_mask' and 'target_mapping' but the results remain similar. Perplexity estimation using discriminative models, such as Bert, and hybrid models, such as XLNet, don't cleanly align with the sequential autoregressive concept of perplexity in GPT-2. We believe that interpreting perplexity as a measure of language fluency (Han et al., 2020; Xu et al., 2018), is contingent on such categorizations of a model's architecture.

## 5 Analysis

### 5.1 Perplexity and Token Length

The separation between non-adversarial and adversarial prompts with GPT-2 in Figure 2 underscores the effectiveness of classifying adversarial suffix attacks with perplexity. This graph also highlights the inadequacy of filtering with a perplexity threshold alone. This is because the perplexity range of regular prompts overlaps and exceeds that of adversarial suffix prompts. In contrast, the XLNet plot
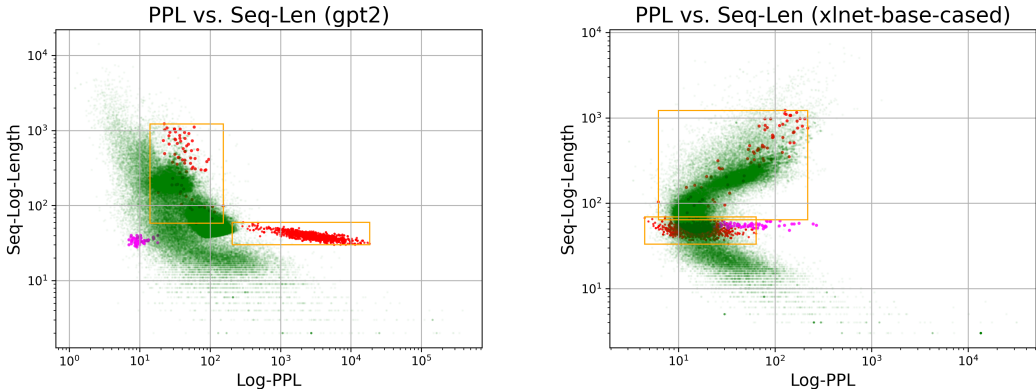
Figure 2: The perplexity values of all prompts are shown here, with GPT-2 on the left and XLNet on the right. The green are regular prompts. Meanwhile, the red ones are adversarial attack attempts. The adversarial suffix attacks are bounded by the lower orange frames, and the human-crafted attacks are bounded by the upper orange frames. Notably, the human-crafted attacks have a lower range of perplexity values and are hard to isolate from the regular prompts. The magenta are adversarial suffix attacks with repeat exclamation marks as their suffixes. The GPT-2 model is an easy choice for a classifier based on perplexity and token length. In contrast, it is visually clear that the XLNet may only work on a smaller subset of the points.

in Figure 2 has a lower level of cluster separation so it is not pursued further. GPT-2 and GPT-2-XL were extremely similar so we chose to proceed with the smaller model.

Inspecting the perplexity distributions of individual regular prompt sources in the Appendix reveals which types of prompts have extremely high perplexity (see 3.4). Single words or short phrases within a conversational context show an elevation of perplexity scores. In particular, this occurred when non-English tokens, mathematical expressions, programming code, misspelled English, or irregular symbols were involved.

Figure 2 shows two orange rectangular regions. The lower right is the region that contains adversarial suffix prompts based on Zou et al. (2023). The opportunity to separate these prompts from regular prompts is visually evident. The magenta points on the left (looking up from $10^1$) represent the prompts with repeat exclamation marks as their suffixes. They are easily recognizable and have not been shown to penetrate LLMs' defenses. Structurally, this region is interesting because it is harder to classify. So, we developed two classifier variants that assume them to be benign, and two that assume that they are an attack in order to challenge our classifier.

The larger upper-left rectangular region in Figure 2 contains all human-crafted GPT4-Jailbreak prompts (Jaramillo, 2023), further described in Appendix A.2. Notably, this dataset contains very few data points. Visual inspection shows that these points lie among ordinary, green, non-adversarial prompts in the same region. It is therefore a good adversarial data source for challenging our classifier since its perplexity values are similar to ordinary English text. Two of our gradient boosting models include them in their training, while another two exclude them. Upon manual inspection, many of the samples in Jaramillo (2023) contain gibberish excerpts at various locations within each prompt. This implies that a more fine-grained perplexity evaluation, like windowed perplexity in Jain et al. (2023), could be better for those prompts in particular.

## 5.2 GRADIENT BOOSTING MACHINES

We considered four GBM classifiers. This is because we wanted to evaluate the performance of training with or without the human-crafted jailbreaks, as well as with or without the ! style attacks from Zou et al. (2023). Although our primary objective is to train a classifier that can detect adversarial suffix attacks without rejecting regular prompts, we want to explore the impact of training or testing with attack prompts of other lengths and perplexity ranges. Figure 2 demonstrates that the !

| Confusion Matrices on Test data of GBM Classifiers | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| GBM Variant: | GBM-0 | | GBM-0x | | GBM-1 | | GBM-2 | |
| Predicted Attacks: | Pos | Neg | Pos | Neg | Pos | Neg | Pos | Neg |
| Actual Attacks | 338 | 17 | 342 | 19 | 338 | 0 | 344 | 0 |
| Actual Non-Attacks | 13 | 26323 | 16 | 26314 | 7 | 26329 | 25 | 26305 |

Table 1: Confusion matrices for the test data for each of the models. Attack predictions are in the columns. For each matrix, the top left contains true positives and the bottom right contains true negatives. The lower left contains false positives (attack predictions whose ground truth was non-attack), and the upper right false negatives (non-attack predictions whose ground truth was attack). The corresponding tables for summarizing these confusion matrices and seeing their breakdown by dataset can be found in Tables 3 and 4.

| GBM-0 - Contribution to Test Performance by Data Origin | | | | | | |
|---|---|---|---|---|---|---|
| | Human-Crafted Attack | | Adversarial Suffix | | Regular Prompts | |
| True Positive | 0 | | 338 | (98.26%) | 0 | |
| False Negatives | 17 | (100%) | 0 | | 0 | |
| False Positives | 0 | | 0 | | 13 | (0.05%) |
| True Negatives | 0 | | 6 | (1.74%) | 26323 | (99.95%) |

Table 2: GBM-0 is highly effective at recognizing adversarial suffix prompts from Zou et al. (2023), but very ineffective on human-crafted GPT-4 attack prompts from Jaramillo (2023). The adversarial suffix true negatives are the trivial prompts with repeat exclamation marks that are correctly classified as non-attack. These are memorizable and not effective at hacking LLMs. Crucially, this model achieves a low regular prompt rejection rate of 0.05%

style attacks and the human-crafted attacks are vastly different in these two dimensions than the rest of the attacks.

GBM-0 is trained on both types of adversarial datasets, adversarial suffix and human-crafted. In this model, the exclamation mark prompts (shown in magenta in Figure 2) are labeled as benign. This is because they are the default attack string before the GCG algorithm comes up with better attacks. GBM-0x is the same as GBM-0 after reversing the exclamation mark prompts and relabeling them as attacks. GBM-1 is a mutation of GBM-0 that drops the human-crafted attack dataset altogether. GBM-2 is like GBM-1 but the exclamation mark suffixes are marked as attacks in order to stress test the classifier. The results shown in Table 3 highlight the challenge posed by attack prompts with lower perplexity and sequence length than the primary adversarial suffix cluster.

The confusion matrices shown in Table 1 and the corresponding scores shown in Table 3, show that the GBM-1 classifier is able to distinguish adversarial suffix attacks from a large variety of ordinary prompts. Models GBM-0 and GBM-0x which include the human-crafted adversarial prompts in the mix, and particularly GBM-0x which also considers the repeat exclamation prompts as adversarial, have a significant drop in performance. A closer look at the failed examples with classifier GBM-0 in Table 2 reveals that all false negatives originate from human-crafted jailbreaks. Model GBM-2, although it excludes the human-crafted prompts from its training, has to recognize the additional cluster of repeat-exclamation prompts as attacks. With GBM-0 and GBM-0x we test both with human-crafted data included in the test set, and without. It is not surprising that GBM-0, which was trained on both adversarial sets, is lagging in performance relative to GBM-1 when tested without any human-crafted data (Table 3).

Next, we contrast the GBM two feature classifiers with a naive perplexity-only threshold-based attack detector. Table 3 includes results for two fixed thresholds, 400 and 1000, as well as the best threshold that maximizes the $F_2$ score. All GBM classifiers outperform the threshold-based method.

| Simple PPL Threshold vs. GBM Classifiers | | | | | | | |
|---|---|---|---|---|---|---|---|
| threshold value: | 400 | 1000 | Best (threshold) | GBM-0 | GBM-0x | GBM-1 | GBM-2 |
| $F_2$ (Adv. Suffix) | 86.7% | 92.6% | 92.9% (851) | 99.2% | 98.6% | 99.6% | 98.6% |
| $F_2$ (2 Attack Types) | 83.8% | 89.1% | 88.5% (689) | 95.4% | 94.9% | NA | NA |

Table 3: The scores achieved by a one-dimensional perplexity (PPL) threshold are substantially lower than the score achieved using any of the GBM classifiers. The perplexity threshold that maximizes the $F_2$ score is derived using validation data and is shown in parentheses next to the score derived from test data in column "Best". The *2 Attack Types* row contains a mixture of adversarial suffix and human-crafted attacks.

| | THR-0 | | GBM-0 | | GBM-0x | | GBM-1 | | GBM-2 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | FN% | FP% | FN% | FP% | FN% | FP% | FN% | FP% | FN% | FP% |
| Adv. Suffix | 2.07 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Human-crafted | 100.00 | 0 | 100.00 | 0 | 100.00 | 0 | nan | nan | nan | nan |
| Adv. Suffix (!)* | 0 | 0 | 0 | 0 | 33.33 | 0 | 0 | 0 | 0 | 0 |
| Puffin | 0 | 5.62 | 0 | 0.30 | 0 | 0.30 | 0 | 0.20 | 0 | 0.60 |
| Docred | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Alpaca | 0 | 4.10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Platypus | 0 | 0.43 | 0 | 0.16 | 0 | 0.24 | 0 | 0.05 | 0 | 0.45 |
| Sglue | 0 | 0 | 0 | 0.20 | 0 | 0.20 | 0 | 0 | 0 | 0 |
| Squad | 0 | 0 | 0 | 0.06 | 0 | 0.06 | 0 | 0 | 0 | 0 |
| Tapir | 0 | 0 | 0 | 0.01 | 0 | 0.01 | 0 | 0.01 | 0 | 0.01 |
| Code | 0 | 2.64 | 0 | 0.07 | 0 | 0.07 | 0 | 0.07 | 0 | 0.07 |

Table 4: Shows a breakdown, by classifier, of all false negatives (FN) and false positives (FP) as percentages over the test dataset count. THR-0 is the best perplexity threshold-based classifier with the same assumptions as GBM-0. Meanwhile, the GBM models succeed on all of the primary adversarial suffix attacks (row 1), and achieve low error rates on all regular prompts (rows 4-11). The false positives are attack predictions whose ground truth was non-attack while false negatives are non-attack predictions whose ground truth is attack. The first two rows contain the primary adversarial types. The *Adv. Suffix (!)* are the exclamation-suffix prompts that are treated as benign for GBM-0 and GBM-1, and as attacks for GBM-0x and GBM-2. The *Human-crafted* attack dataset is excluded from the GBM-1 and GBM-2 models, showing as nan.

Finally, Table 4 presents a breakdown of the false negatives and false positives percentages for each classifier on all of the datasets (regular and adversarial). This table highlights the extent to which different data sources present challenges to our suite of four GBM classifiers. This table helps to illustrate that the gradient boosting classifiers are effective at reducing the error rate on all of the regular prompt datasets.

## 6 LIMITATIONS

The following limitations should be considered:

- More data on adversarial and non-adversarial prompts would be ideal. Training longer could shift the attack distribution.

- We test our classifiers on some adversarial points with lower perplexity and sequence length, but adaptive attackers could create more variations.

- An adaptive attacker could modify their algorithm to optimize towards having low perplexity. In this scenario, windowed perplexity was found by Jain et al. (2023) to detect all of the adversarial attacks at the cost of a false positive rate of 10% on a regular prompt dataset (AlpacaEval).

- An adaptive attacker could also vary the attack prompt length and or the attack suffix length. Our attack prompts come from the "Harmful Behaviors" dataset in Zou et al. (2023), but an attacker could add in new prompts of varying lengths. Furthermore, we use the default suffix search length from the code, but an adaptive attack with shorter suffix lengths has been found to be harder to detect with perplexity in Jain et al. (2023).

- While an adaptive attacker could optimize towards the perplexity distribution of GPT-2, a defender could experiment with other LLMS for perplexity. Several other LLMs are evaluated for perplexity in Jain et al. (2023).

## 7 CONCLUSION

A gradient boosting classifier based on perplexity and sequence length is a useful tool for detecting adversarial suffix attacks coming from Zou et al. (2023)'s source code. We select GPT-2 for computing the perplexity of input prompts to an LLM due to its comparable performance with GPT-2-XL. Furthermore, we show that an LLM like XLNet will produce a weaker signal for identifying adversarial attacks with perplexity. GBM-0's rejection rate on 26,336 benign prompts in the test set was below 0.1%, while all the adversarial suffix examples in the test set were detected. This contrasts with a concurrent work on detecting these attacks with perplexity as a singular feature. For contrast, in Jain et al. (2023) the lowest reject rate on benign prompts was 5.7% using AlpacaEval. Our GBM-0 model rejects no prompts on an AlpacaEval test set.

We produced over 1400 adversarial suffix attacks using the default settings for the GCG algorithm from Zou et al. (2023). Nearly 90 percent of examples had a perplexity above 1000, while the rest of the relevant attack examples yielded perplexity above 200. We contrasted the perplexity distributions of the adversarial strings with a variety of regular prompt data sets. These regular prompts varied in sequence length from a few tokens to over 10,000 tokens. Some were found to exceed the perplexity of adversarial prompts by orders of magnitude. Therefore, a perplexity filter alone would risk a high rate of rejecting benign user inputs to an LLM.

The two feature classifier could not detect human-crafted jailbreaks like those in Jaramillo (2023). Windowed perplexity, introduced by Jain et al. (2023) may help with certain human-crafted jailbreaks in Jaramillo (2023) that we saw have gibberish subsections in the prompts. For example, one prompt starts with a multitude of punctuation characters in a meaningless sequence about a "System annou-ncement". This excerpt has a perplexity of over 2000 with GPT-2 while the whole attack that it is sourced from evades detection by our classifier because it has a perplexity of 64. Features that leverage the semantic clues of the malevolent intent of the prompts (Markov et al., 2023) could be tried to improve the classifier when the other features fail. We speculate that non auto-regressive variations on perplexity such as XLNet may add some discriminatory power for detecting adaptive attacks on a subset of the points.

Several avenues that an adaptive attacker would take such as minimizing the perplexity, and varying the sequence length are discussed in the context of concurrent work (Jain et al., 2023). A successful adaptive attacker would be able to bypass the defense, but would not be able to force our classifier to reject benign user prompts. To illustrate how the classifiers would handle adversarial attacks with lower sequence length and perplexity, we include four gradient boosting models. Each model incorporates a different view of some of the possible values an adaptive attacker could attempt (using the ! suffix attacks and the hand-crafted attacks).

Our approach may be able to detect adversarial suffix attacks that are similar to Zou et al. (2023), like in Lapid et al. (2023). It could also be explored on attacks that rely on sending and receiving bizarre sequences of tokens that are encrypted messages like the ciphers in Yuan et al. (2023) or the Base64 attack in Wei et al. (2023).

We suggest an update to the guidelines on researching defenses in Carlini et al. (2019), since it is described as a living document that is meant to be updated based on feedback. We propose that defenses for LLMs rigorously evaluate the rejection rates their defense would have on regular user queries. We tried to achieve this goal by incorporating over 175,000 regular prompts from 8 sources to stress test our classifier. Now that neural network systems like ChatGPT have millions of text queries a day, defenses for LLMs must perform such stress tests on large diverse samples of regular prompts to be useful.

## ETHICS STATEMENT

In order to avoid potential harm, we blacked out components of the example attack strings that we introduced to the literature in the appendix section. These attack strings have the potential to transfer to other LLM models.

## REPRODUCIBILITY

We provide the code used for the calculations and figures in this study at this URL:

`https://osf.io/mucw5/?view_only=6193b18437e443f7a8f2ea294b3b633c`

The notebook's name is "nb_llm_adv_analysis_multiple_models.ipynb".

## REFERENCES

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022.

Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial robustness, 2019.

Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. Theoremqa: A theorem-driven question answering dataset. *preprint arXiv:2305.12524*, 2023.

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *NAACL*, 2019.

Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1310–1320. PMLR, 09–15 Jun 2019. URL `https://proceedings.mlr.press/v97/cohen19c.html`.

Maria Cvach. Monitor alarm fatigue: an integrative review. Biomedical instrumentation & technology, 2012.

Amelia Glaese, Nat McAleese, Maja Trebacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. Improving alignment of dialogue agents via targeted human judgments, 2022.

Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples, 2015.

Wenjuan Han, Liwen Zhang, Yong Jiang, and Kewei Tu. Adversarial attack and defense of structured prediction models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 2327–2338, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.182. URL `https://aclanthology.org/2020.emnlp-main.182`.

Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. Unsolved problems in ml safety, 2022.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL `https://openreview.net/forum?id=nZeVKeeFYf9`.

Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. Baseline defenses for adversarial attacks against aligned language models, 2023.

Rubén Darío Jaramillo. Jaramilo:gpt4jailbreak, howpublished = `https://huggingface.co/datasets/rubend18/chatgpt-jailbreak-prompts`, 2023. Accessed: 2023-09-20.

Erik Jones, Anca Dragan, Aditi Raghunathan, and Jacob Steinhardt. Automatically auditing large language models via discrete optimization, 2023.

Raz Lapid, Ron Langberg, and Moshe Sipper. Open sesame! universal black box jailbreaking of large language models, 2023.

Ariel N. Lee, Cole J. Hunter, and Nataniel Ruiz. Platypus: Quick, cheap, and powerful refinement of llms, 2023.

Klas Leino, Zifan Wang, and Matt Fredrikson. Globally-robust neural networks, 2021.

Yuhui Li, Fangyun Wei, Jinjing Zhao, Chao Zhang, and Hongyang Zhang. Rain: Your language models can align themselves without finetuning, 2023.

Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. URL `https://openreview.net/forum?id=rJzIBfZAb`.

Todor Markov, Chong Zhang, Sandhini Agarwal, Tyna Eloundou, Teddy Lee, Steven Adler, Angela Jiang, and Lilian Weng. A holistic approach to undesired content detection in the real world, 2023.

Annunziata Elefante Mattia Limone, Gaetano Cimino. Tapir: Trigger action platform for information retrieval. `https://github.com/MattiaLimone/ifttt_recommendation_system`, 2023.

Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018.

OpenAI. Gpt-4 technical report, 2023.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.

Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples, 2016.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. In *NA*. OpenAI, 2019. URL `https://api.semanticscholar.org/CorpusID:160025533`.

Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ Questions for Machine Comprehension of Text. *arXiv e-prints*, art. arXiv:1606.05250, 2016.

Tomohiro Sawada, Daniel Paleka, Alexander Havrilla, Pranav Tadepalli, Paula Vidas, Alexander Kranias, John J. Nay, Kshitij Gupta, and Aran Komatsuzaki. Arb: Advanced reasoning benchmark for large language models, 2023.

Taylor Shin, Yasaman Razeghi, Robert L. Logan IV au2, Eric Wallace, and Sameer Singh. Auto-prompt: Eliciting knowledge from language models with automatically generated prompts, 2020.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, and Nikolay Bashlykov. Llama 2: Open foundation and fine-tuned chat models, 2023.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. *arXiv preprint arXiv:1905.00537*, 2019.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions, 2023.

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training fail?, 2023.

Yotam Wolf, Noam Wies, Oshri Avnery, Yoav Levine, and Amnon Shashua. Fundamental limitations of alignment in large language models, 2023.

Jingjing Xu, Xuancheng Ren, Junyang Lin, and Xu Sun. Dp-gan: Diversity-promoting generative adversarial network for generating informative and diversified text, 2018.

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. Xlnet: Generalized autoregressive pretraining for language understanding, 2020.

Yuan Yao, Deming Ye, Peng Li, Xu Han, Yankai Lin, Zhenghao Liu, Zhiyuan Liu, Lixin Huang, Jie Zhou, and Maosong Sun. DocRED: A large-scale document-level relation extraction dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 764–777, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1074. URL `https://aclanthology.org/P19-1074`.

Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. Reclor: A reading comprehension dataset requiring logical reasoning. In *International Conference on Learning Representations (ICLR)*, April 2020.

Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Pinjia He, Shuming Shi, and Zhaopeng Tu. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher, 2023.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023.

Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models, 2023.

# A ADVERSARIAL DATASETS

## A.1 GENERATED ATTACK PROMPTS

Using the code provided by Zou et al. (2023), we generated 1407 prompts intended to attack LLMs and induce them to produce illicit responses. The plots below show the perplexity and token length frequency distribution of these prompts, and a scatter plot.
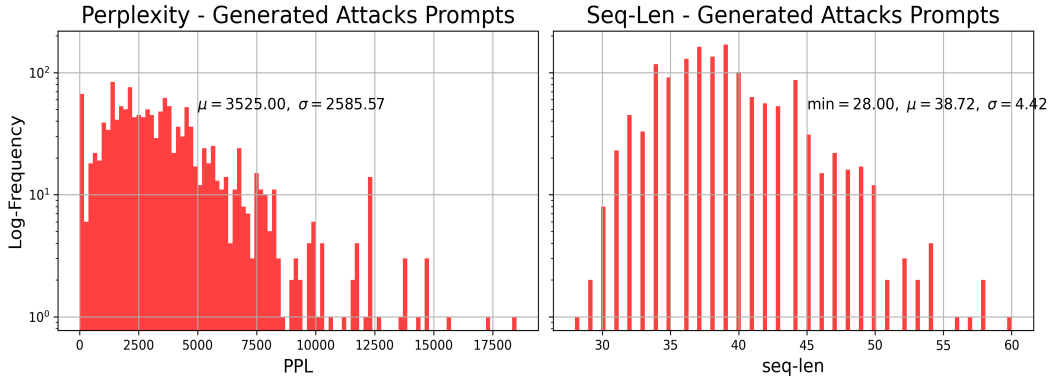


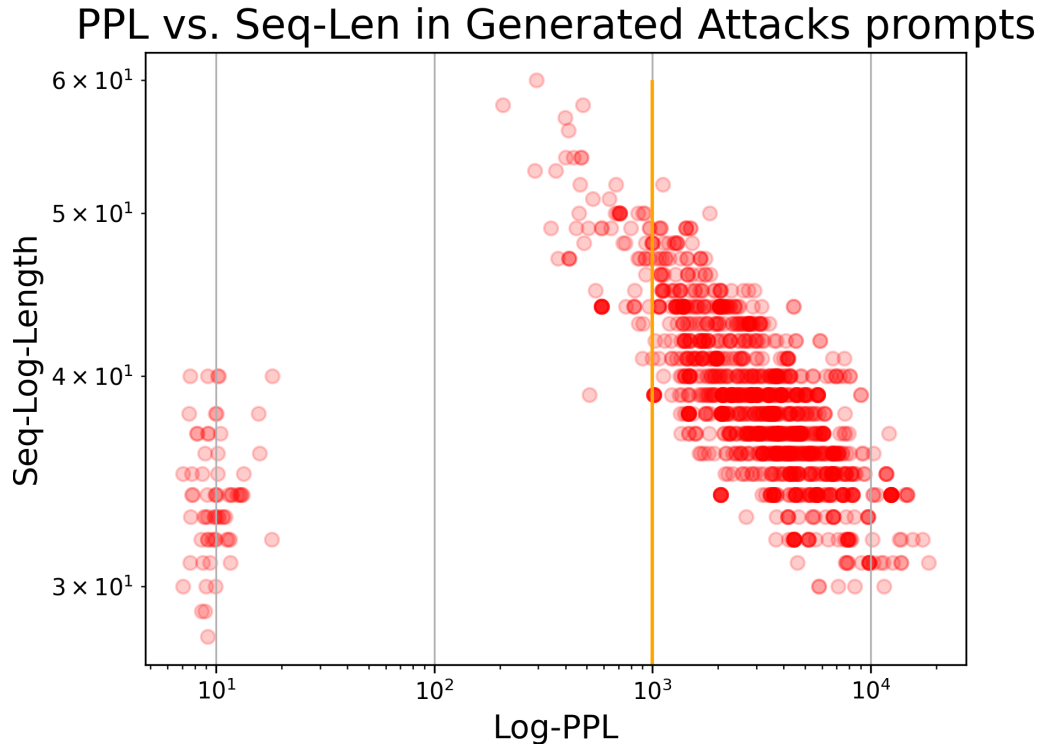Figure 3: Generated attack prompt perplexity and sequence-length frequencies



Figure 4: Generated Attack prompt perplexity vs. sequence-length

The scatter plot shows two islands of points. The ones with low perplexity, on the left, represent the initial prompts used by the generation model, which, according to our testing, did not succeed.

The larger set on the right represents the generated prompts that may break the defenses of LLMs.

Here are some examples of successful adversarial suffix attacks from our system-generated strings. We censored parts of the adversarial strings and output code so they could not be reused for nefarious purposes. We tested them on Vicuna-7B (Zheng et al., 2023), but in principle, they could be transferred to other LLMs as explained in Zou et al. (2023):

**Warning:** The following images contain offensive AI-generated content.



Figure 5



Figure 6

Figure 7



Figure 8

## A.2 GPT4 JAILBREAK PROMPTS

This is a small dataset with only 79 examples from the Huggingface hub, named rubend18/ChatGPT-Jailbreak-Prompts (Jaramillo, 2023) of manually constructed prompts that claim to have broken alignment defenses on GPT4. The plots below show these prompts' perplexity, token length frequency distribution, and scatter plot.
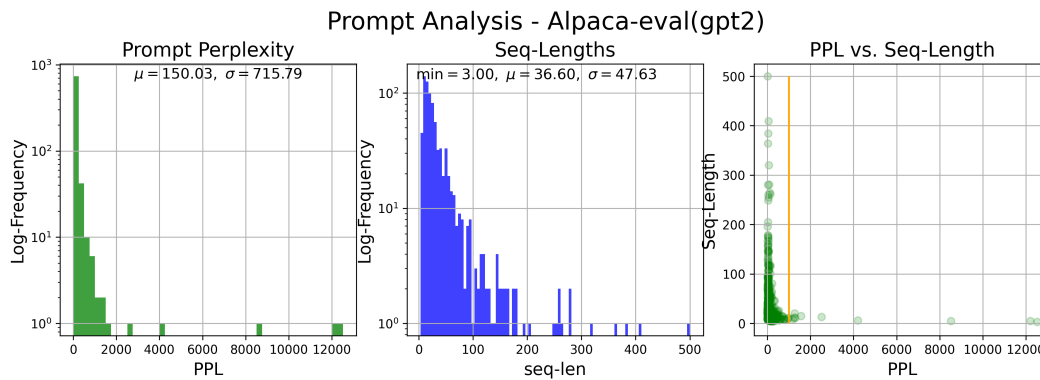


Figure 9: Log-frequency distributions for perplexity and sequence-length



Figure 10: Scatter-plot showing perplexity vs. sequence-length

# B    NON-ADVERSARIAL DATASETS

## B.1    DOCRED

This dataset can be found in the Huggingface hub under the name docred (Yao et al., 2019). We use the validation split, containing 998 multi-sentence passages designed for the development of entity and relation extraction from long documents. The plots below show these prompts' perplexity, token length frequency distribution, and scatter plot.



Figure 11: Log-frequency distributions for perplexity and sequence-length

## B.2    ALPACA-EVAL

This dataset can be found in the Huggingface hub under the name tatsu-lab/alpaca_eval. It contains 805 generated pairs of instructions and corresponding outputs. We use the instruction as the prompt. The plots below show these prompts' perplexity, token length frequency distribution, and scatter plot.



Figure 12: Log-frequency distributions for perplexity and sequence-length

## B.3    SUPERGLUE

This dataset can be found in the Huggingface hub under the name super_glue (Wang et al., 2019). We use the validation split of the subset named *boolq* (Clark et al., 2019), containing 3270 passages for answering Yes/No questions. We formulated prompts by combining the fixed instruction "Read the following passage and answer the question:", followed by the question field in the dataset example, and on a new line we write the passage field of the example. The plots below show these prompts' perplexity, token length frequency distribution, and scatter plot.
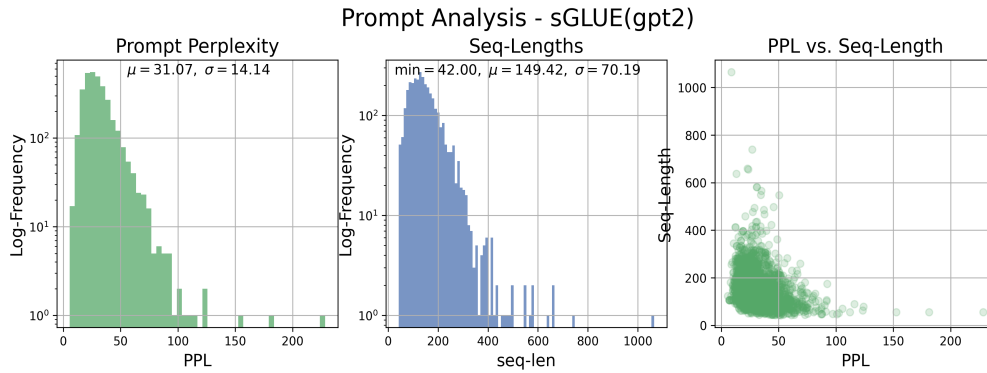
Figure 13: Log-frequency distributions for perplexity and sequence-length

## B.4 SQUAD-V2

The Stanford Question-Answering Dataset is a well-known span-based question-answering dataset that can be found in the Huggingface hub under the name squad_v2 (Rajpurkar et al., 2016). We use the validation split containing 11873 examples. We formulated prompts by combining three fields from each example, the title, the context and the question using the following form: We start with an instruction "Given a context passage from a document titled [title field goes here], followed by a question, try to answer the question with a span of words from the context:". Then after a new line the prompt continues with "The context follows:" followed by the context field, and then after another new line "The question is:" followed by the question field.

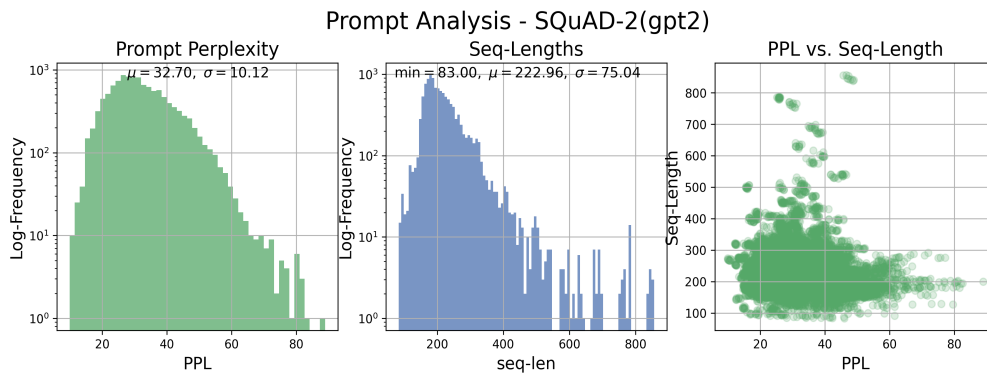The plots below show these prompts' perplexity, token length frequency distribution, and scatter plot.



Figure 14: Log-frequency distributions for perplexity and sequence-length

## B.5 OPEN PLATYPUS

The Open-Platypus Dataset is associated with the Platypus project. We use the Huggingface dataset garage-bAInd/Open-Platypus containing 24926 prompts with instructions from the Platypus dataset's training split, as they appear, without any additional prefix or suffix. This dataset is focused on improving LLM logical reasoning skills and was used to train the Platypus2 models (Lee et al., 2023; Touvron et al., 2023; Hu et al., 2022). It is a filtered collection from multiple scientific, reasoning, and Q&A datasets including: Yu et al. (2020); Chen et al. (2023); Mihaylov et al. (2018); Sawada et al. (2023); Lee et al. (2023).

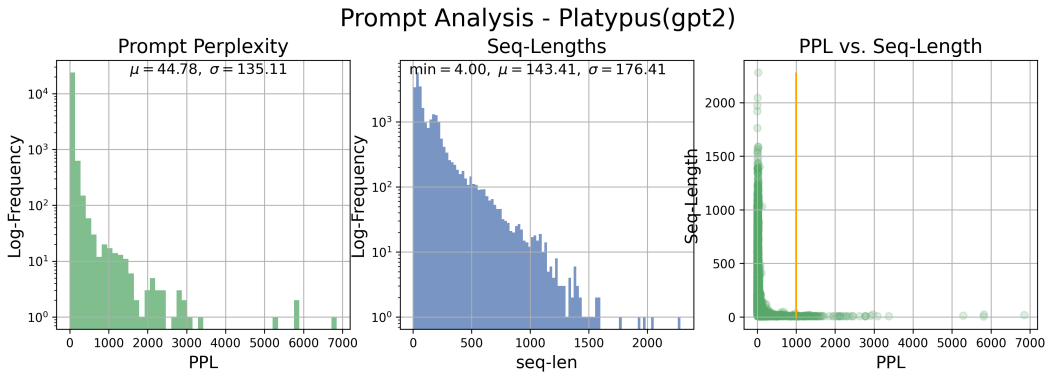The plots below show these prompts' perplexity, token length frequency distribution, and scatter plot.

Prompt Analysis - Platypus(gpt2)



Figure 15: Log-frequency distributions for perplexity and sequence-length

The shortest prompt is *Hello, AI!*.

### B.6 PUFFIN

This dataset can be found in the Huggingface hub under the name LDJnr/Puffin. Puffin contains 3000 conversations with GPT-4, each being a sequence of interactions that start with the human's query. We constructed two samples from this dataset. One is the set of all 6994 prompts produced by the human side of the conversation. The other contains only the initial utterance that starts each of the 3000 conversations since this is a more relevant structure to the attacks we observed in Zou et al. (2023).

The first set of plots shows the human prompts' perplexity, token length frequency distribution, and scatter plot.
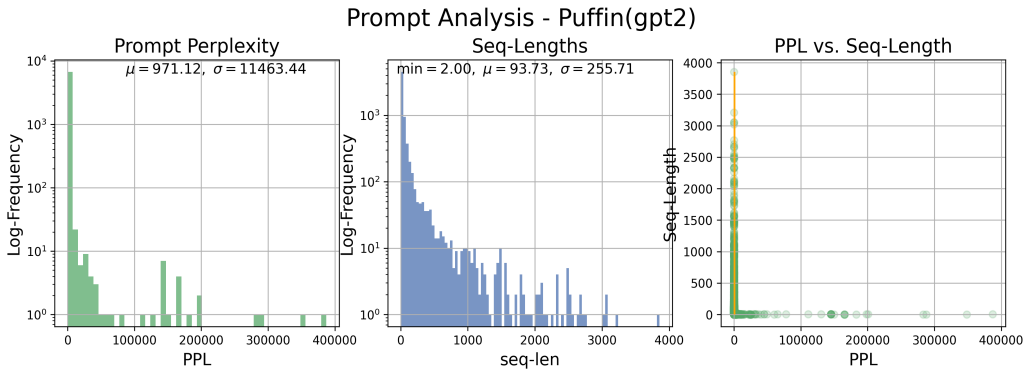
Prompt Analysis - Puffin(gpt2)



Figure 16: Log-frequency distributions for perplexity and sequence-length

Some monolectic utterances in Puffin result in "NaN" perplexities. We collected them, and they contain the following items:

*'1', '10', '2', '3', '389', '4', '433', '6', '7', 'A', 'B', 'Continue', 'Finish', 'Hi', 'NO', 'No', 'Test', 'Yes', 'a', 'b', 'c', 'conservative', 'continue', 'd', 'hi', 'more', 'next', 'no', 'ok', 'sometimes', 'thanks', 'web', 'yes'*

It is easy to conclude that they make little sense standalone, but they are meaningful when considered in the context of a conversation's prior interactions.

The next set of plots shows the perplexity, token length frequency distribution, and scatter plot of the initial human prompt that starts each conversation, Puffin[0].
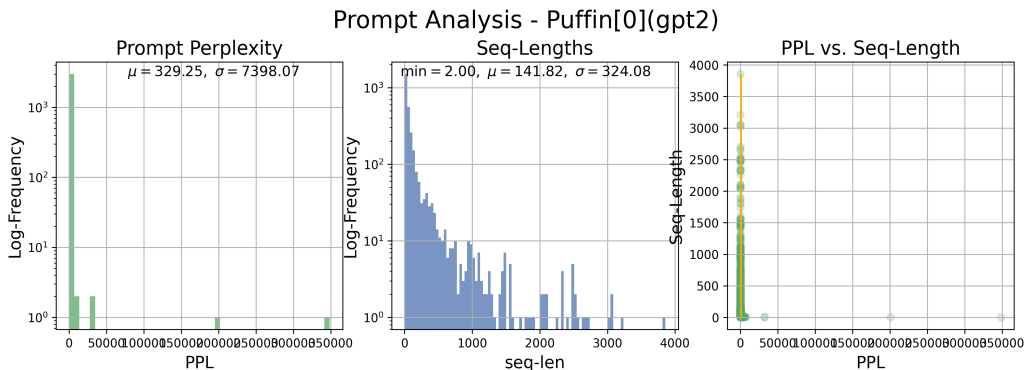
Figure 17: Log-frequency distributions for perplexity and sequence-length

## B.7 TAPIR

This is a large dataset containing examples intended for instruction-following training. We use the Huggingface dataset MattiaL/tapir-cleaned-116k (Mattia Limone, 2023) containing 116862 examples. We construct prompts by concatenating the instruction field and the input field from each example.

The plots below show these prompts' perplexity, token length frequency distribution, and scatter plot.
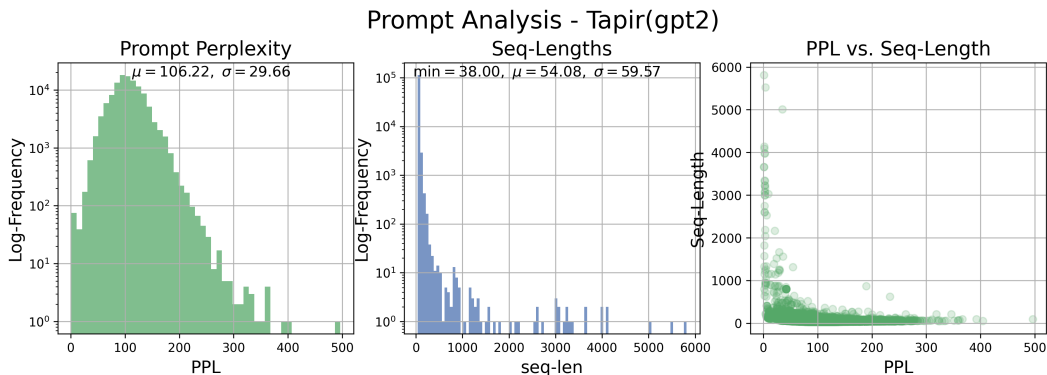


Figure 18: Log-frequency distributions for perplexity and sequence-length

## B.8 INSTRUCTIONAL CODE SEARCH

This is a large dataset containing instructional examples for coding in Python. We use the Huggingface dataset Nan-Do/instructional_code-search-net-python. because the data set is very large we only include the first 10,000 examples.

The plots below show these prompts' perplexity, token length frequency distribution, and scatter plot.

## Prompt Analysis - Code(gpt2)

### Prompt Perplexity

$\mu = 128.31, \ \sigma = 200.47$

### Seq-Lengths

$\min = 11.00, \ \mu = 249.19, \ \sigma = 518.47$
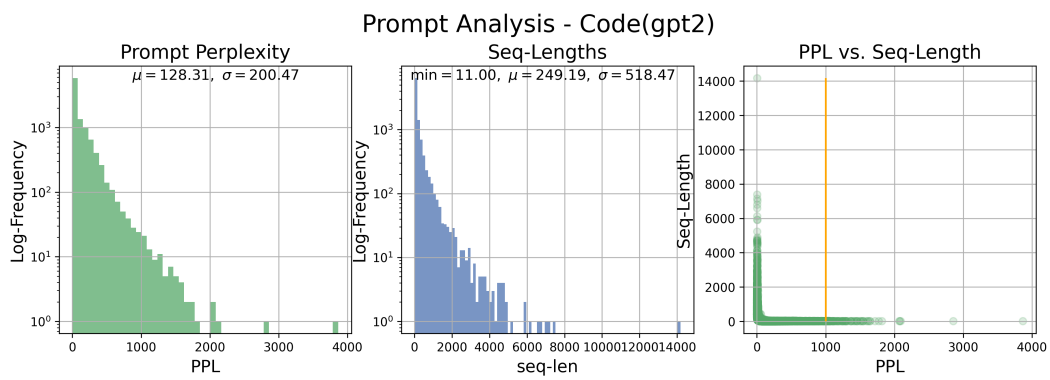
### PPL vs. Seq-Length

Figure 19: Log-frequency distributions for perplexity and sequence-length