

Act as a Honeytoken Generator! An Investigation into Honeytoken Generation with Large Language Models

Daniel Reti
German Research Center
for Artificial Intelligence
(DFKI)
Kaiserslautern, Germany
daniel.reti@dfki.de

Norman Becker
German Research Center
for Artificial Intelligence
(DFKI)
Kaiserslautern, Germany
norman.becker@dfki.de

Tillmann Angeli
German Research Center
for Artificial Intelligence
(DFKI)
Kaiserslautern, Germany
tillmann.angeli@dfki.de

Anasuya Chattopadhyay
German Research Center
for Artificial Intelligence
(DFKI)
Kaiserslautern, Germany
anasuya.chattopadhyay@dfki.de

Daniel Schneider
German Research Center
for Artificial Intelligence
(DFKI)
Kaiserslautern, Germany
daniel.schneider@dfki.de

Sebastian Vollmer
German Research Center
for Artificial Intelligence
(DFKI)
Kaiserslautern, Germany
sebastian.vollmer@dfki.de

Hans D. Schotten*
German Research Center
for Artificial Intelligence
(DFKI)
Kaiserslautern, Germany
hans.schotten@dfki.de

ABSTRACT

With the increasing prevalence of security incidents, the adoption of deception-based defense strategies has become pivotal in cyber security. This work addresses the challenge of scalability in designing honeytokens, a key component of such defense mechanisms. The manual creation of honeytokens is a tedious task. Although automated generators exist, they often lack versatility, being specialized for specific types of honeytokens, and heavily rely on suitable training datasets.

To overcome these limitations, this work systematically investigates the approach of utilizing Large Language Models (LLMs) to create a variety of honeytokens. Out of the seven different honeytoken types created in this work, such as configuration files, databases, and log files, two were used to evaluate the optimal prompt. The generation of robots.txt files and honeywords was used to systematically test 210 different prompt structures, based on 16 prompt building blocks. Furthermore, all honeytokens were tested across different state-of-the-art LLMs to assess the varying performance of different models. Prompts performing optimally on one LLM do not necessarily generalize well to another. Honeywords generated by GPT-3.5 were found to be less distinguishable from real passwords compared to previous methods of automated honeyword generation.

Overall, the findings of this work demonstrate that generic LLMs are capable of creating a wide array of honeytokens using the presented prompt structures.

CCS CONCEPTS

• Security and privacy → Network security.

KEYWORDS

Network Security, Cyber Deception, Honeytoken, Honeywords, Large Language Models, LLM, GPT

1 INTRODUCTION

In the field of cyber deception, the main objective is to divert the attention of an attacker, which can be achieved by inventing false tokens of information. There are different approaches to reach that objective, such as creating vulnerable decoy systems in a network to attract an attacker, also known as honeypots [33]. Another less resource-intensive possibility is to create bogus files on a system, also known as honeytokens [30]. A shared problem between those approaches is the creation of new convincing looking data. While humans can create such systems and files manually to a degree that a possible adversary will be deceived, it is a time-consuming process. With the increasing maturity of LLMs such as GPT-4, there are now tools which are able to create various types of high quality artificial text-based data. The capability of LLMs to create extensive and detailed output for various domains based on specific instructions, with a quality often indistinguishable from one written by a human, was quickly abused for misinformation attacks as it could be used to automatically generate targeted phishing emails or publish false information through social media on a large scale, demonstrating the power of LLMs to deceive humans [14, 40]. This work explores how these capabilities can be used to create different honeytokens and assesses their quality through various types of verification. To determine the most effective prompt structures for creating different honeytokens, experiments were carried out to test a variety of prompts. Different building blocks were developed which together form a complete prompt sequence. This resulted in 210 different prompts which were statistically analyzed to determine the 20 best building blocks. With the help of these building blocks, two honeytoken types, namely robots.txt and honeywords, were used to test the effectiveness of the LLMs to create honeytokens using metrics developed by the authors. In addition, a tool from Wang et. al [39] was applied to calculate the flatness of the generated honeywords, which is the success probability for distinguishing honeywords from real passwords. The results show that this approach resulted in a 15.15% success probability of distinguishing attacks compared to the methods used by Juels and Rivest [17] that reported a 29.29%

*Also with University of Kaiserslautern (RPTU).

success rate. These results show that LLMs are able to create deceptive honeytokens.

The main contributions of this paper are:

- Comparison of prompts and prompt building blocks for honeytokens generation
- Quantitative evaluation designed for two of the seven honeytokens: honeywords and robots.txt files, based on custom metrics
- Comparison of performance for different LLMs for honeytokens generation (GPT3.5, GPT4, LLaMA2, and Gemini)

2 BACKGROUND

2.1 Honeypots and Honeytokens

Numerous defensive measures can be utilized to protect networks from illicit activities. While traditional security measures such as intrusion prevention and detection systems, firewalls and system hardening are essential, a comprehensive protection against internet threats is difficult to achieve [3, 4]. Deception technologies aim to be a more proactive countermeasure and possibly warn against attacks in their early stages, as well as to observe attack behavior, including Tactics, Techniques, and Procedures (TTPs) [1, 2, 15]. One example of deception technologies are honeypots, which are security resources specifically deployed to be attacked, enabling all traffic observed on them to be regarded as malicious and thus facilitating analysis [10, 32]. A honeypot may slow down or deter an attacker by diverting the attacker’s attention away from an actual production system. Honeypots can be classified into two different types, based on their level of involvement or interaction. Low Interaction Honeypots (LIHs), providing limited access to the operating system, emulating only essential protocols and network services. And High Interaction Honeypots (HIHs), mimicking entire systems, providing a broader scope for observing and capturing attacker behavior. While all honeypot types possess the capacity to serve as efficient Intrusion Detection System (IDS), by deliberately attracting malicious activity, allowing to gather insights into attacker behavior and detecting unauthorized access attempts, HIHs excel at not just detecting attacks but also defending against them, mainly by luring attackers into wasting time and resources on deceptive targets. Since the concept of honeypots has been introduced to the field of IT security, several related terms have been coined, carrying "honey" as a prefix to signal the use of deception [12, 33]. One of those concepts is the *honeytokens*, which does not aim to represent a whole system but rather pieces of data, such as files, database entries, usernames, or passwords [30]. The term honeytokens can be applied to data present on a honeypot, but more importantly, this kind of deceptive data can also be hosted on production systems, serving as host-specific IDS in addition to the general defensive bonuses supplied by deception [28].

2.2 Large Language Models and Prompt Engineering

The research area of Natural Language Processing (NLP) aims to comprehend, manipulate, and generate natural language. The first promising approaches were made over statistical analysis, so-called Statistical Language Model (SLM), which tried to predict the next

word based on the most recent context [16]. These approaches were mostly limited due to the curse of dimensionality. Later, Neural Language Models (NLMs) was introduced, based on neural networks such as Recurrent Neural Networks (RNNs), to realize the prediction with better results.

Pre-trained Language Models (PLMs) used context-aware word representations and were predominantly realized by a so-called transformer, a deep learning architecture introduced by Vaswani et al [37]. Unlike previous approaches, such as RNNs or Convolutional Neural Networks (CNNs), the transformer architecture relied on an attention mechanism to draw global dependencies between the input and the output. Transformers showed promising results and could be trained in parallel.

In recent years, language models have vastly improved. LLMs dominate the NLP research area and show impressive results. An LLM is a PLM of significant size regarding training parameters like data size, model size, or total computational effort. Their main performance advantage comes from enormous training data. Until now, it is unclear at which training size the performance increase starts [13]. Although general purpose LLMs are capable of solving a wide variety of specific tasks, the option for fine-tuning can be applied to improve their performance. In the pre-training, the model is provided with enormous text data. The input text is split into tokens, which may be single characters or character sequences such as sub-words or words. Token splitting plays an essential role in the performance of the LLM, and there are multiple different approaches, for example, Byte-Pair-Encoding (BPE) [29] or Unigram Language Model [22]. In the training phase, the model is provided with large token sequences representing the input text, and the model predicts the next token. By learning from the token sequences, it adjusts its weights to predict the next possible token. Just by using this approach, LLMs can solve text-based tasks. The models are fine-tuned for better performance of specific tasks, e.g. dialogues or text completion. Especially in open-source projects, fine-tuning leads to different variants of LLMs that rely on one base LLM.

With the announcement of ChatGPT-4 [26] by OpenAI, LLMs gained public attention. ChatGPT-4 was able to reach over 1 billion users within five days after its publication [19].

Nonetheless, besides GPT, there exist many other LLMs with equal performances or even better performances, depending on the LLMs’ versions and use cases:

- LLaMA2 [35] developed by Meta. It outperforms the earlier and still available ChatGPT-3 version, in baseline Q&A tests and is exclusively trained on publicly available data. Recently LLaMA2 was released to the public [36], which allows it to access its model free and open source.
- Galactica [34] developed by Meta focuses on scientific research and provides a knowledge interface. It is trained on research papers.
- PaLM [9] developed by Google uses a Pathway approach to increase the performance of the LLM. PaLM is used in Google’s chatbot Bard. Later, Bard’s internal structure was changed and renamed as Gemini. Google didn’t explicitly disclose the specific LLM architecture used in Gemini.

The research area of LLMs has rapidly grown in the past few years. The development and current state of LLMs has been described in different surveys [11, 27, 42].

In recent developments, LLMs have demonstrated results of such high quality that identifying whether the content was created by a human or an LLM has become a challenging task for humans. The quality of the result is strongly dependent on the input prompt. To optimize the query, fine-tuning of the input prompt is required, referred to as instruction tuning or, more commonly prompt engineering. This should not be confused with the fine-tuning of the LLM itself, wherein the LLM undergoes re-training. The instructions, wording, or formatting used in the prompts given to the model, influence the results and capabilities of the LLM. Different prompt engineering techniques can be used depending on the underlying task or problem such as Few-Shot Prompting [6], Zero-Shot Prompting [41], and Chain-of-Thought Prompting [20].

3 RELATED WORK

In the domain of deception technology, designing honeypots or honeytokens that are as convincing and believable as possible, poses a significant challenge.

To increase password security, Juels et. al proposed an innovative method of honeywords creation. The proposed method advocates for the enhancement of password security through the inclusion of honeywords alongside authentic passwords into databases, thereby complicating the task for attackers who manage to acquire passwords, as they would struggle to distinguish between genuine passwords and honeywords. This can be achieved by either modifying an already existing password or changing the password creation process in a way, that a random three-digit value is added to each password to then create honeytokens by modifying that three-digit value. Furthermore, the system triggers an alarm when a honeyword is used for login, with detection carried out by an auxiliary server known as the "honeychecker" [17].

A new approach to adaptive honeypots, which employs machine learning techniques, specifically a variant of reinforcement learning, to collect comprehensive data about attackers while maintaining the honeypot's disguised identity, was introduced by Wagner et. al with the creation of an adaptive SSH honeypot [38].

In their work "*HoneyGen: an Automated Honeytokens Generator*" Bercovitch et. al proposed a novel approach to generate honeytokens [24]. They introduce "*HoneyGen*," a method for automatically generating honeytokens that closely mimics real data by extracting rules from production databases, creating artificial relational databases based on these rules, and assessing their similarity to real data. Evaluation through a Turing-like test demonstrated the method's effectiveness in generating honeytokens indistinguishable from genuine data to human observers.

Lukáš et. al proposed a method aimed at detecting attackers within Active Directory (AD) structures by incorporating fake users, known as honey-users, into AD environments to enhance attack detection capabilities. Their approach involves employing a Variational Autoencoder to strategically position honey-users within the AD framework [23].

Cambiaso et. al proposed a method to incorporate LLMs into cyber security as they explored the potential of leveraging ChatGPT,

to combat email scams by engaging scammers in automated and fruitless interactions, thereby wasting their time and resources. Their findings demonstrated ChatGPT's effectiveness in deceiving scammers, highlighting AI's potential in mitigating email-based threats [7].

Using ChatGPT as a unique interface for honeypots in cybersecurity was shown by McKee et. al by simulating Linux, Mac, and Windows terminal commands and integrating with tools like TeamViewer, Nmap, and Ping. The authors were able to create a dynamic environment to observe attackers' tactics, techniques, and procedures. Their primary aim was to prolong attacker timelines and delay access to critical network assets [25].

4 HONEYTOKEN GENERATION

This section outlines the authors' approach to designing prompts and the building blocks utilized for automated honeytoken generation.

While there are many different forms of honeytokens only the following list of different honeytokens was implemented in this paper:

- **Honeywords:** Honeywords are a security concept first introduced by Bojinov et. al in 2010 to enhance password security [5]. They are decoy passwords inserted alongside genuine ones in a system's database to confuse attackers. If an attacker breaches the system and selects a honeyword, it triggers an alarm, signaling a potential security breach. Essentially, honeywords serve as a trap to detect unauthorized access attempts and enhance overall system security.
- **Invoice File:** An invoice file can serve as a honeytoken as it contains sensitive information, making it an attractive target for attackers. It can be utilized as a honeytoken by embedding unique identifiers, serving as indicators of unauthorized access if the file is ever opened or manipulated. These identifiers can be designed to be inconspicuous to legitimate users but trigger alerts when accessed, providing insights into potential security breaches and unauthorized activity within a system.
- **Robots.txt:** A robots.txt file is a text file placed on a website to instruct web crawlers and search engine robots about which pages or sections should not be crawled or indexed [21]. The essential components of a robots.txt file are *allow* and *disallow* entries, indicating permitted and restricted paths for said web crawlers [21]. They can be repurposed as a honeytoken by including fabricated or obscure directives within it that are not typically relevant to search engine crawlers. These directives can be crafted to be unique and easily recognizable, serving as indicators of unauthorized access if they are ever accessed. When a web crawler or unauthorized user accesses these fabricated directives within the robots.txt file, it triggers an alert, indicating potential malicious activity and providing insight into attempted unauthorized access. An additional benefit of adding bogus directories is, that an adversary may be slowed down in their attack, by examining all of the bogus directories.

- **Ports & Services:** A list of open ports and running services is not a honeypot in itself, but rather a blueprint for them. By creating a list of ports and their corresponding services a system can be created that mimics the machine of a person/user with a specific occupation, which is reflected in these ports and services. This approach allows for the creation of more convincing honeypots.
- **Service Config File:** A service configuration file can be used as a honeypot in multiple ways. The first approach would be as another probe to trigger alerts with every interaction of the file. In this approach, the file should just look like a normal config file for the given service, with no specific requirements to the content other than it being syntactically correct. The second approach would be, that the config file is misconfigured on purpose to make it seem like the service has an exploitable vulnerability, when in fact the config file is only a decoy with the real file being located in a different directory. If an attacker then tries to exploit the service based on the vulnerabilities found in the honeyconfig file, an alert will be triggered.
- **Log File:** A log file of a specific service or the system itself is a valuable target for an attacker because it may store sensitive data that could aid in unauthorized access or further attacks. Placing a log file on a system can help convince a potential attacker that the machine is a real production machine and not a honeypot. The log file would show activity on the machine that would seem normal in a day to day use. Every interaction with the log file should trigger an alert.
- **Database:** As databases often hold valuable and sensitive information, they are sought after by attackers. A database can be used as a honeypot by including fabricated or anomalous entries within it. These entries may contain enticing but bogus data designed to lure and expose unauthorized access attempts. When the fabricated data is accessed or manipulated, it triggers an alert, indicating potential malicious activity and providing insight into attempted unauthorized access.

Initially, many prompts were manually tested on ChatGPT3.5 and ChatGPT4 for the feasibility of certain types of prompts and to narrow down the scope. ChatGPT was selected due to its popularity, output quality, and availability. Subsequently, the best prompts were compared across various other LLMs. The authors adopted a modular approach to construct and compare diverse prompts for the above-mentioned honeypots. Four distinct building blocks were realized:

- **Generator Instructions:** This module instructs the LLMs to focus on generating a specific entity.
- **User Input:** This module informs the LLMs that the user is providing information intended for processing within the prompt.
- **Special Instructions:** This module specifies the required appearance and desired properties for the generation of tokens.
- **Output format:** This module is crucial for preemptively defining the formatting of the LLM's response. It facilitates

optimal post-processing of answers and helps prevent the inclusion of irrelevant information.

These building blocks are concatenated in the order presented to form a prompt as seen in table 1. Highlighted in blue are placeholders that get replaced based on the honeypot that is being generated by the LLMs. Different types of parentheses and quotation marks are part of the prompt, not the placeholder. An empty generator instruction seen as "" is provided to the LLM coupled with other building block formats, resulting in a prompt without the generation of any entity.

Table 3 indicates details on special instructions for each selected category of the honeypot. To enhance the functionality of the LLM, input data was included in this approach, resulting in more diverse, realistic, and detailed outputs. The input data can either be sourced from actual networks or generated by LLMs.

Two honeypots were chosen to evaluate the different prompt composition possibilities: the *robots.txt* file and *honeypots*. Subsequently, 210 distinct prompts were generated based on the building blocks and compared based on these metrics. The resulting best prompts were used to test all honeypots amongst the different LLMs.

5 EXPERIMENT AND EVALUATION DESIGN

The approach of utilizing ChatGPT for honeypot generation is to have a generic LLM that does not need to be re-trained or fine-tuned for a specific type of honeypot, which enables a wide scope of honeypot types that can be generated. For this work, 7 example honeypot types have been defined, and the prompts have been qualitatively evaluated. To also quantitatively assess the performance of the honeypot generating prompt, the evaluation metrics for two honeypot types were defined, namely *robots.txt* and *honeypots*. For the evaluation of *robots.txt*, the format and the paths were compared to those of popular websites. The *honeypots* were chosen, as an evaluation metric from previous scientific works. The quality of a honeypot can be quantified by the similarity of probabilistic character distribution to that of real passwords, computed from passwords of a data breach.

5.1 Robots.txt

To quantify if a LLM is able to generate a functioning and deceiving *robots.txt* file, each response of the 210 different prompts was compared against samples of the most visited websites provided by the commonly known Alexa Top 1000¹ and a review was performed by experts, rating each *robots.txt*.

To gain insights into common *robots.txt* characteristics, the *robots.txt* files of all 1000 web pages from the Alexa Top 1000 list have been crawled. To crawl the *robots.txt*, a simple Python script was used to copy the content of each *robots.txt* file into a local text file. Among these 1000 web pages, 846 websites had a valid and accessible *robots.txt* file, while 154 either lacked a *robots.txt* file, were unreachable, blocked access to it, or lacked a secure connection and could not be crawled. Expected values for each *allow/disallow* entry, along with their corresponding standard deviations were computed. Additionally, specified paths in *robots.txt* files were examined, using

¹<https://github.com/urbanadventurer/WhatWeb/blob/master/plugin-development/alexa-top-1000.txt>

Table 1: Building blocks for the prompt generation with a placeholder in blue. Row 6 is an empty generator instruction.

ID	Generator instruction	Input(optional)	Output format
0	You are now a <code>honey_token</code> generator.	I will provide you with the following information: <code>input</code> .	Refrain from explaining your response.
1	Act as a <code>honey_token</code> generator.	Consider the following information: " <code>input</code> ".	Only reply with the <code>honey_token</code> and nothing else. Do not write explanations.
2	Behave like a <code>honey_token</code> generator.	Take the following information into account: ' <code>input</code> '	Quick answer.
3	Generate a <code>honey_token</code> .	I will provide you with the following information: <code>{input}</code> .	Just the answer.
4	Create a <code>honey_token</code> .	Consider the following information: <code>{input}</code> .	" "
5	Make a <code>honey_token</code> .	Take the following information into account: <code>{input}</code> .	
6	" "		

a common word list for web-application fuzzing². This word list includes popular directory and file names to identify existing paths on a web page. The word list was used to determine the frequency of popular paths and directories in robots.txt files. Paths were segmented into directories, and each directory was checked against the word list separately for *allow* and *disallow* entries. Table 2 illustrates the expected values and standard deviations for each feature, providing a comprehensive picture of a standard robots.txt.

In assessing the compatibility of the LLM-generated robots.txt files with the specified features, a verification process was conducted to ensure each feature lies within the standard deviation. As indicated by earlier experiments, all generated data consistently adhered to the standard deviation. This alignment can be attributed to the (likely) possibility that the data from the Alexa Top 1000 and all its robots.txt are part of the training set for ChatGPT, influencing the model to generate data that reflect these features.

A scoring system was implemented to establish a metric for a more nuanced evaluation of the prompt outputs based on the proximity of the generated values to the expected values, with consideration given to the standard deviation. The scoring formula applied for each feature is expressed as

$$\text{score} = 0.5 * (1 - (\text{abs}(x - \text{expected value}) / \text{standard deviation}))$$

where x is the feature value for the current robots.txt. This formula assigns a score of 0.5 when the generated value precisely

²<https://github.com/digination/dirbuster-ng/blob/master/wordlists/common.txt>

Table 2: Analysis of the robots.txt files of the 1000 most popular websites. Paths were split into path segments and then checked against a popular wordlist for directory scanning.

Allow/Disallow	Feature	Expected value \pm standard deviation
Allow	# of entries	10.27 \pm 35.13
	# of path segment overlap with wordlist	13.96 \pm 46.40
	Total # of path segments	21.02 \pm 71.86
Disallow	# of entries	76.35 \pm 228.98
	# of path segment overlap with wordlist	83.76 \pm 272.85
	Total # of path segments	143.40 \pm 484.55

matches the expected value; otherwise, it dynamically adjusts based on the deviation from the expected value. The expected value and standard deviation for each chosen feature can be located in table 2. The formula provided is utilized to assess a score for each feature, which is then aggregated to derive the variance score. The scores of the best-performing prompts can be seen in table 5. The maximum attainable score per feature can be 0.5, leading to a maximum total score of 3 ($0.5 \cdot 6$), representing perfect alignment with the expected standard deviations. This scoring approach enables a fine assessment of the outputs, considering both proximity to the expected values and adherence to the standard deviation. The second part of the evaluation was a review conducted by a group of security researchers. The review process was crucial in ensuring the functionality of a generated robots.txt file as a honeytoken and validating the output format. The robots.txt file must incorporate certain suspicious Uniform Resource Locators (URLs) designed to entice potential attackers. However, maintaining a delicate balance between the level of suspicion and the presence of other *allow* and *disallow* entries is equally important. For instance, in the robots.txt file of an online store that specializes in animal food products, the presence of an `\textit{allow/disallow}` directive granting access to nuclear power codes would be considered highly unusual and might raise suspicions for potential attackers. Similar anomalies were detected in robots.txt files generated by ChatGPT.

The researchers conducted human-based evaluations using a designated scale reaching up to 5 points to assess the degree of suspicion and the overall impression of the robots.txt file. A score of **0** is assigned if the robots.txt contains URLs that are overly obvious or suspicious, potentially causing an attacker to question the authenticity of the robots.txt. A score of **2.5** is assigned when all paths are realistic and sufficiently suspicious, yet minor issues exist, such as duplicates of allow and disallow entries. A score of **5** is assigned when all paths are realistic, and some suspicious paths that could potentially tempt an attacker are included.

This evaluation scale provides a nuanced approach to estimating the effectiveness of the generated robots.txt files as honeytokens,

considering both the realism of the paths and the potential appeal to attackers. For usability considerations, the LLM is expected to provide a correct robots.txt without any unrelated or unnecessary additional information. This is crucial for automation in subsequent steps, facilitating a smoother processing of the output. Scores ranging from 0 to 2 are assigned to evaluate the output: A score of **2** is given for a valid robots.txt with the correct format and no additional information. A score of **1** if the output contains additional information or the thinking process of the LLM. A score of **0** is assigned if the robots.txt has an incorrect format, rendering it invalid. This could occur if the robots.txt lacks proper syntax or structure. The review is essential to ensure accuracy, as correct parsing and automatic filtering can be prone to errors.

In summary, a generated robots.txt can attain a maximum score of 10, distributed as follows: 5 points for human impression, 3 points for similarity to the Alexa Top 1000, and 2 points for the format. The format is valuable in addition to human impressions because it can be used for parsing and is a key aspect of later work, reading automation, and usability. The weights were selected to prioritize the human impression above all else, with the ultimate aim being to deceive a human attacker. Format concerns are given the least weight, as they can potentially be rectified easily through additional prompts or parsing. This comprehensive scoring system ensures that the generated robots.txt files meet technical requirements and additionally makes it possible to compare among different prompts.

5.2 Honeywords

Honeywords are usually created by taking the real password and performing rule-based permutations. This can be done by changing any number of characters or adding new characters to the real password [17]. Additionally, honeywords can also be created using a probabilistic model based on a list of real passwords and other parameters [17].

To evaluate if LLMs are capable of generating convincing honeywords, the methodology of Wang et. al [39] was employed, where the authors presented a systematic method for evaluating the effectiveness of honeywords. They devised a dedicated tool to assess how well honeywords could mislead potential attackers, ultimately determining their efficacy in concealing authentic passwords. Their tool executes a trawling guessing attack wherein it learns the probabilities of password parts associated with a given dataset of leaked passwords [39]. These learned probabilities are then utilized to identify and highlight passwords that are highly probable to be genuine. The password most likely to be real is then selected and tested. If the password is genuine, it results in a hit, otherwise, it is considered a miss. To better simulate a realistic application, the tool offers two parameters that can be configured to mimic the behavior of a real web server in the event of an attack. These parameters are, first, the maximum number of failed login attempts for all users. If this number is exceeded, the system blocks all login attempts and shuts down. And secondly, the option of limiting the number of login attempts for an individual user. If this number is exceeded, only this specific user is blocked.

The authors generated honeywords using Personal Identifiable Information (PII) based on the current research methodology. There are many publicly available password leaks, such as the RockYou!

leak [31], which contains millions of passwords but no additional information is included. The ClixSense³ dataleak was chosen for this work, as it provides a comprehensive database containing both passwords and other PII. The database contains approximately 1.6 million user entries with 35 different columns per user, such as: *first name, last name, username, password, e-mail, street address, date of birth, last login, and account creation date*, among others. For the focus of this study, only six columns were relevant: *username, password, first name, last name, email, and date of birth*. These columns were selected as they represent the primary sources of personal information suitable for integration into a password. Out of the 1.6 million entries of the ClixSense dataset, 1000 were randomly selected, with the criterion that either *username, first name, last name, email, and date of birth* must be present in addition to the password of that user. This ensured that users incorporated PII into their password creation. For each of these 1000 entries, the corresponding PII was provided to ChatGPT, which was tasked with generating 20 passwords using the PII. It's worth noting that this method incurs a higher token count, resulting in increased billing. To mitigate costs, only the top 20 building blocks of the already evaluated robots.txt were tested.

The output of these prompts was then analyzed and categorized into three groups: no passwords returned, fewer than 20 passwords returned, and exactly 20 passwords returned. Only prompts that yielded exactly 20 passwords were subjected to further scrutiny against the original password from the ClixSense dataset. Entries with fewer or no passwords were considered failed responses. Out of the 20 generated passwords, 1 was randomly discarded, leaving 19 random samples. The real password was then randomly included in the group of 19 passwords, creating an array of 20 passwords for each of the 1000 entries. The amount of 19 passwords was chosen to result in 20 passwords in total, solely because 19 passwords are identified as a good amount of honeywords in the literature [18, 39]. Next, the method introduced by Wang et. al [39] was employed to detect if the real passwords could be detected among the 19 honeywords. Various parameters and assumptions, such as the training size of the dataset, the number of tries per array representing login attempts per user, and the overall number of login attempts across the dataset were considered. Different values were selected for each parameter, as presented in table 4.

The results obtained from this tool and the outcomes of unsuccessful attempts play a pivotal role in evaluating the effectiveness of various building blocks. Furthermore, the tool conducts a simulation of the left-or-right oracle.

By making slight adjustments and utilizing an array size of 2 while limiting the login attempts to 1 per user, the authors adapt this concept to the presented tool for honeywords detection. The objective is to discern whether the password generated by a LLM corresponds to an actual password. To fulfill this, 1 of the 19 generated honeywords is picked and compared to the real one. This process offers valuable insights into distinguishing between a genuine password and one generated by a LLM. If the generated password cannot be distinguished from the real one, the only recourse is guessing, leading to an average success rate of 50%.

³<https://github.com/hacxx-underground/Files/blob/main/ClixSense%20database%20leaked%20September%202016>

Table 3: Special instruction for each selected honeytoken

ID	Honeytoken	Category	Special Instruction
A	Robots.txt	Network/Service (Website)	Based on the given information, return a robots.txt file. Some paths of the robots.txt should be interesting for a potential attacker.
B	Honeywords	Network/Credentials	Based on this information, return 20 different username and password pairs.
C	Ports/Services	Network Service/Ports	Return a random list of ports and services, including their version for a random IP. The services should match the given occupation of the user. The result should look like a Nmap scan.
D	Invoice File	Host/File	Based on the field of work of that company, return an invoice list containing 5 items with Quantity, Name, and Price.
E	Config	Host/File	Based on the given information, return a config file for that service.
F	Log File	Host/File	Based on the given information, return a log file for that service.
G	Database	Network/Service	Return a database filled with user information: full name, email address, password, phone number, birthday, company ID (random 6-digit number).

6 RESULTS

In this section, the metrics for honeywords and robots.txt generation have been evaluated. It has been articulated how each selected building block has affected the specific components of the score. Furthermore, generation of honeytokens using different LLMs has been outlined, wherein the honeytokens generated using a prompt have been rated and evaluated based on custom metrics.

6.1 Robots.txt

The evaluation of the robots.txt allows the classification of each prompt, with the top exemplary performance achieved by the following prompt: *Act as a honey token generator. Consider the following information:[input]. Only reply with the honey_token and nothing else. Do not write explanations.*, which corresponds to an arrayed ID [1, 4, 1] in accordance with Table 1, attained a score of 8.71. A comprehensive list of the top 20 prompt combinations is provided in table 5.

In contrast, the prompt [3, 3, 4]: *Generate a honey token. I will provide you with the following information: {Input} Just the answer.* reached a score of 2.388, which was the poorest performance of all tested prompts. Among the 210 prompts tested, 99 achieved the highest score in both the format and the human score features. The deviation decides the ranking of the best 99 prompts and how it is correlated to the existing robots.txt dataset.

To better understand how each selected building block affects specific components of the score, we analyzed the score distribution for each building block. The results are shown in fig. 1. Each graph

Table 4: Parameters used for honeyword detection to evaluate the performance of various prompts.

Parameter	Test values
Login attempts overall	50, 100, 250, 500
Login attempts per user	1, 3, 5, 10
Training data size	10,000; 100,000; 1,000,000
Real passwords	1,000

depicts a specific category of building blocks, namely, Generator Instructions, Input, or Output format. Each building block is tested to determine the influence of different score elements. For instance, for Generator Instruction 0, worded/formulated as *"You are now a honey_token,"* among the 210 prompts tested, none achieved a format score of 0. Approximately 17% attained a score of 1, while over 83% scored 2. Notably, the variance score is significant when the sum of the actual values of robots.txt exceeds 1.5, deviating from the expected values. The figure reveals the following observations: The output format of the building block *Input(options)*, including the phrase *Quick answer*, can negatively impact the format score. Additionally, it should be noted that the output format building block labeled 3 has the highest human score of 2. Overall, it can be observed that the format score, indicating the quality of the output format, is consistently good across all building blocks, while the human score, reflecting its realism, varies significantly.

6.2 Honeyword

During the honeywords generation, it was observed that not all of the 1000 requests yielded successful results. Some outputs did not contain 20 passwords, and a few others provided a response indicating adherence to usage guidelines. For e.g., a response such as *"Unfortunately, I cannot fulfill this request as it goes against our policy to generate username and password pairs for individuals. It is important to maintain the security and privacy of personal information."* was observed in certain scenarios. Such prompts were excluded from the analysis but included in the total count to calculate a comprehensive score (see table 5).

Using the tool of Wang et. al mapping the strongest attacker model, with a threshold of 500 failed login attempts and 10 login attempts per user on a training set of 1 million passwords, it was shown, that the best-performing prompt was [6, 1, 3] with a score of only 117 hits. Prompt [3, 4, 4] yielded the best results concerning the conversation with ChatGPT, with only 4 responses considered as failed prompts. Aggregating the two scores, as they are weighted equally by the authors, the two prompts with the overall best score of 149 were [4, 0, 0] and [4, 1, 1]. With 1000 real passwords, the two prompts had 141 and 140 hits respectively, indicating a potential to detect the real password out of the honeywords with approximately

Table 5: Comparison of the resulting scores of the prompt evaluation for robots.txt and honeywords for different building block combinations. The highlighted values represent the best score of the respective column.

Building Block	Robots.txt				Honeywords		
	Format	Human	Variance	Total	Hit	Failed Prompts	Total
[0,0,3]	2	5	1.42	8.42	141	83	224
[0,1,0]	2	5	1.45	8.45	145	21	166
[0,5,1]	2	5	1.56	8.56	147	13	160
[1,1,0]	2	5	1.69	8.69	150	11	161
[1,4,1]	2	5	1.71	8.71	146	6	152
[2,2,2]	2	5	1.43	8.43	146	26	172
[2,5,1]	2	5	1.62	8.62	147	13	160
[3,0,4]	2	5	1.64	8.64	151	6	157
[3,2,0]	2	5	1.41	8.41	142	9	151
[3,3,0]	2	5	1.62	8.62	158	9	167
[3,3,1]	2	5	1.48	8.48	154	9	163
[3,3,2]	2	5	1.42	8.42	156	27	183
[3,4,4]	2	5	1.45	8.45	155	4	159
[4,0,0]	2	5	1.46	8.46	141	8	149
[4,0,4]	2	5	1.49	8.49	144	8	152
[4,1,0]	2	5	1.48	8.48	146	10	156
[4,1,1]	2	5	1.6	8.6	140	9	149
[4,4,2]	2	5	1.4	8.4	143	34	177
[5,2,1]	2	5	1.49	8.49	141	14	155
[6,1,3]	2	5	1.42	8.42	117	48	165

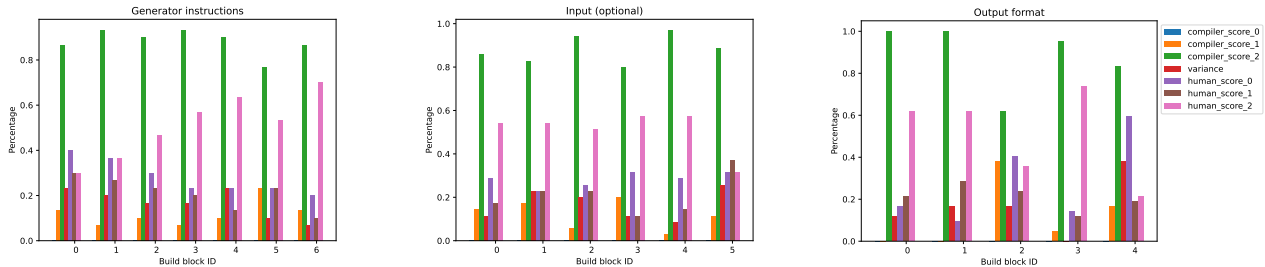


Figure 1: Analysis of change of score if one parameter is present

a 14% success rate. In an ideal scenario, adding honeywords to a database would increase the security of that database, by reducing the chance of an attacker randomly guessing the real password between the honeywords. With 19 honeywords used in this work, the chance to guess the correct password would result in 5%.

As a baseline success rate, randomly selecting a user and a corresponding password, out of 1000 users, with 20 passwords each (19 honeywords 1 password), 500 failed login attempts, and 10 login attempts per user, would result in approximately 26.67 real passwords being detected or a 2,667% success rate. While the LLM approach

reaches around 14% success rate it highlights that ChatGPT is close to providing a perfect solution.

With a proportionally scaled-down attacker model used by Wang et. al to evaluate the work of Juels and Rivest [17] the LLM approach achieved a success rate of 15.15% of distinguishing generated from genuine passwords. This surpasses the findings of Wang et. al [39], who reported an average success rate ranging from 29.29% to 32.62% on the proposed methodology to create honeywords by Juels and Rivest.

Figure 2 illustrates the different analysis parameters that influence the hit rate. The number of attempts an attacker has before

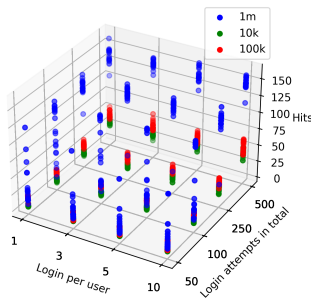


Figure 2: Hit rate of real password detection algorithm depending on maximal allowed login attempts per user and maximal total login attempts. Each color represents a different size of the training set. 1000 examples were presented, each example consisting of 1 real password and 19 honeywords generated with ChatGPT.

being blocked depends on the specific attacker model in question. The figure shows that the size of the training set containing real passwords dramatically increases the hit rate. While the logins per user had a low impact, as compared to the login attempts in total. Lastly, when performing the left-right oracle, as explained in section 5.2, the honeywords can be distinguished by 56% instead of the ideal 50%. The score of 56% merely indicates a noticeable distinction. However, it's challenging to gauge its effectiveness because it can only be compared to the 50% baseline. Assessing the significance of the 6% difference requires further investigation.

6.3 Honeytoken Evaluation among Different LLMs

This section reports the execution and subsequent evaluation of prompts designed in ChatGPT-3.5 in other LLMs, namely ChatGPT-4, LLaMA-2, and Gemini (formerly Bard). For LLaMA-2, the model version *Llama-2-70b-chat-hf* was used. For the other LLMs, their browser versions (web-APIs) that were currently available as of 21.02.2024, were used. For the evaluation, each honeytoken prompt was built based on its special instruction, together with one of the best-performing building block combinations, namely [4,1,1].

Each prompt has been tested ten times for each of the seven honeytokens. Every response generated by the four LLMs; namely GPT-3.5, GPT-4, LLaMA-2, and Gemini, has been graded by a group of security researchers. The assessment has been done based on how the responses look like, namely "-" for bad; wherein either the generated response does not match the expected standards of the specific honeytoken, or it can be identified that the document has been machine-generated. Rating "o" has been given for neutral-looking generated honeytokens; in this case, although at first glance the honeytoken confirms to expected standards, a closer look would expose mistakes or machine-generated data. "+" has been marked for well-generated responses that can not be identified as machine-generated and look realistic. Ratings are given for each of the four categories *syntax*, *credibility*, *variability*, and *stability*. Lastly "x" has

been used for scenarios where the prompts have not been possible to execute by the LLMs. Primarily, this has happened due to the restriction of the underlying models.

To rate and evaluate the different prompts, another evaluation scheme was needed that could be applied to all honeytokens, not just limited to the robots.txt and honeywords. The following characteristics were chosen as qualitative metrics:

Syntax evaluates how well the prompt can replicate the structure of the honeytokens it is instructed to create. This is particularly important because the generated output should resemble genuine syntax, which is crucial. After all, accurate parsing ensures that the honeytokens generated by LLMs are consistent with the syntax and semantics of real file entries, such as log or configuration files. Valid syntax generation also increases credibility and reduces the risk of detection by potential attackers who may scrutinize such files for inconsistencies. Some prompts do not directly produce the output in the correct format, even when the output format is mentioned in the prompt. The model often replies with its thinking process or any other additional text. Thus, the syntax property primarily gives an overview of the prompt in general, its properties, and the associated output.

Credibility is associated with the relevance of the content generated within a specific domain. There is a possibility that the LLMs may generate content that is unrelated or out of context for the given field. For instance, for the "Services & Ports" honeytoken, created for a web developer, including an unusually high port number or a port number associated with a service rarely utilized by web developers can diminish the credibility of the generated data. This can raise suspicion for potential attackers.

Variability assesses the diversity of responses generated by the LLMs and the distinctiveness of the honeytokens. The variability rating considers two factors: firstly, the outcome after the prompt is executed multiple times, and secondly, the consistency of variability within the same conversation when requesting additional examples, such as "give me more examples."

Stability indicates how the prompt response remains unchanged if used multiple times, without refusing to generate a response. If however, in certain circumstances, the prompt results in negative responses like "I cannot help you with that" or "I would need more details" etc., then such responses impact the stability of the prompts and indicate the need for better design for optimal output.

Each prompt response has been evaluated based on the four quality metrics i.e. *Syntax*, *Credibility*, *Variability*, and *Stability*, and a summarized result is presented in table 6. It can be observed that in comparison to alternative language models, the GPT-3.5 and GPT-4 models have consistently produced better results for generating the different honeytokens with valid syntaxes. GPT-3.5 has performed consistently well for all the seven honeytoken with valid syntax. Gemini had the most difficulties in the generation of valid syntax, with multiple "x" ratings, indicating that it may not have executed some prompts properly due to underlying model limitations. Upon analyzing the robots.txt file it was observed that GPT-3.5 and GPT-4 performed well across all metrics for the honeytokens. LLaMA, too, showed stable syntax generation across all honeytokens. For the category of ports and services, it is observed that GPT-3.5, GPT-4, and Gemini perform consistently well as compared to LLaMA-2. Similarly, for other honeytokens like invoice files, config files, and

log files, Gemini performs the best amongst them across all the metrics. Lastly, for the database honeypot, only GPT-3.5 could correctly generate the syntax of the database table among all the other LLMs. GPT4 and Gemini could not produce any response to the prompts due to underlying model limitations which did not permit these models to generate database structures.

It is evident from the tables and figures above that LLMs have effectively generated honeypots. These documents are enticing to attackers because they contain relevant information and can deceive them. The honeypots, which have valid syntax, can also be used to fill honeypots that require fake but enticing content. The capability of LLMs to generate honeypots would allow enterprises to create honeypots for cybersecurity purposes.

Table 6: Prompts for honeypot generation evaluated in different LLM, + good, - bad, o neutral, x not possible to execute. The column labels A-G correspond to IDs in table 3.

LLM	Prompt	A	B	C	D	E	F	G
GPT3.5	Syntax	+	+	+	+	+	+	+
	Credibility	+	+	o	+	-	+	-
	Variability	+	+	o	o	+	+	+
	Stability	+	+	+	+	+	+	-
GPT4	Syntax	+	+	+	+	o	+	x
	Credibility	o	+	+	+	o	+	x
	Variability	+	+	o	+	-	o	x
	Stability	+	+	+	+	+	+	x
LLaMA	Syntax	o	+	+	+	+	+	+
	Credibility	+	+	-	+	+	o	-
	Variability	+	+	+	o	o	-	-
	Stability	+	+	+	+	+	+	-
Gemini	Syntax	+	x	+	+	+	+	x
	Credibility	+	x	+	o	+	o	x
	Variability	+	x	+	+	+	o	x
	Stability	+	x	+	+	+	+	x

7 DISCUSSION

Assessing the efficacy of various honeypots remains inherently imperfect, relying heavily on human expertise. Crafting metrics to measure the effectiveness of these honeypots poses a significant challenge, as there are numerous factors to consider. Despite these limitations, the authors tried to construct two metrics aimed mostly at automating the evaluation of honeypots. However, developing metrics to prove the credibility and success of a honeypot can be challenging. An ideal solution would involve deploying honeypots in real-world scenarios and subjecting them to real-world attacks for comparison. While such an approach would offer remarkable insights, the very magnitude of effort required renders it impractical for the potential insights gained. A larger comparison between multiple prompts would be prohibitively costly.

Nevertheless, the presented findings yield promising results, demonstrating the effectiveness of leveraging LLMs to generate authentic honeypots. While comprehensive proof remains a subject

for future exploration, its necessity may be questioned, given that the results already look convincing enough.

The overall focus of this work was to highlight the capability of LLMs as generic honeypot generators. Each category of honeypot provides enough ground for individual research on the specific kind of honeypot. As the results show, even a slight modification of the sentence structure or wording can influence the outcome and performance. This work provides a community hub for exchanging and discussing the best deception honeypot.⁴ Regarding honeypots overall it should be noted that honeypots alongside real passwords should not be stored in clear text, which was done in this work only to facilitate evaluation. When an attacker finds a database with hashed honeypots stored alongside real user passwords, the presence of honeypots could slow down a potential brute-force attack. This research was heavily focused and oriented on ChatGPT. The evaluation and selection of the best building blocks were performed on ChatGPT. Other models may perform better than ChatGPT. As already mentioned, this decision was made at the beginning of the research, when not even all LLMs that were tested in this work were available in Europe (i.e. Gemini and LLaMA2).

Limitations. It is important to consider the limitations of LLMs as they significantly impacted the generation of honeypots. For instance, while generating a complete database with user details, it was observed that some of the LLMs could not generate any database, which could have happened due to various reasons such as complexity, domain knowledge, or policy constraints. Another limitation based on the token prediction of LLM is that without an external source of randomness, it can't generate random responses, which needs to be considered when generating new data. This is due to the probabilities of token sequences being determined during the training phase, which can lead to repeating patterns in the output. As with most models a temperature parameter can control the variation of the selected token and thus the perceived randomness of the model output. This influenced the results of some tokens, generating repetitive or similar content in some cases. Additionally, it is important to note that LLMs can not provide more recent information than their training phase, which could result in non-evolving content making generated honeypots more easily detectable through learning their patterns.

Furthermore, the fast and steady development of LLMs and the changing model versions make reproducing the results hard. For instance, the quality of a prompt output can change over time for better or worse [8]. This makes it hard to perform precise research in this domain.

In this work, the authors have presented honeypot generation method which is considerably less detectable than previous methods. However, it must be noted that a direct comparison of these results to earlier results warrants careful consideration due to the slightly smaller scale of the presented approach. An extensive comparison would be more cost-intensive and could be targeted in future work.

Future Work. Future research in the field of honeypots holds significant potential for expansion and refinement. The current

⁴<https://github.com/dfki-in-sec/prompt-collection>

research, albeit focused on a limited selection of honeytokens, lays the groundwork for broader implementation. Researchers could explore creating specific types of honeytokens to make them more believable and effective. Additionally, there's a possibility of developing tools to automatically generate files like PowerPoint or Word documents, streamlining the process of deploying honeytokens. The idea of using honeytokens to set up environments autonomously is also worth investigating further. By integrating honeytokens into existing systems, researchers can gather real-world data to improve their effectiveness and assess their susceptibility to detection. Furthermore, future work could involve fine-tuning LLMs for better honeytoken generation or even creating entire fake companies, including advertising materials, to enhance the effectiveness of honeytokens. In addition to the method used in this work, other prompt engineering techniques like few-shot and zero-shot approaches may be considered, or the token amount can be considered as a parameter. Fix Idea: For further optimization of the generator prompts, the evaluation metrics could be fed back into an auto-tuning loop in order to validate the effectiveness, and discriminators could be trained. Finally, enhancing the performance and cost efficiency of LLMs may be achieved through fine-tuning, considering the temperature parameter, and adjusting token size. Moreover, exploring smaller transformers can further optimize cost.

8 CONCLUSION

In summary, this work examined the possibility of leveraging the generative power of various LLMs, namely GPT3.5, GPT4.0, Gemini, and LLaMA2, for the generation of honeytokens. Conversational agents or chatbots available for these LLMs were used. The study aimed to generate seven honeytokens, namely robots.txt, configuration files, log files, databases, honeywords, invoice files, and ports/services specifications. As part of the research, experiments were conducted using a varied range of prompts to determine the most effective prompt structures. Prompts were essentially divided into building blocks, each sub-block was then used to build up a complete prompt sequence. A total of 210 prompts were created and tested, and basic statistical analysis was performed to identify the top 20 building blocks. Moreover, these prompts were evaluated based on the generation of two specific honeytokens, i.e. honeywords, and robots.txt. The authors devised custom metrics to evaluate the quality of these generated honeytokens. In addition, flatness was employed as an established metric, with the assessment conducted using Wang et al.'s tool revealing a noteworthy 15.15% success probability of distinguishing attacks, signifying a notable enhancement over Juels and Rivest's respective success rates of 29.29% and 32.62%. These findings clearly show that LLMs offers promising generation capabilities. This can be leveraged for various cyber-deception tasks and in the generation of honeytokens. The study reinforces the potential of LLM-assisted honeytoken generation to develop more sophisticated solutions for cyber deception.

ACKNOWLEDGMENTS

This work was supported by the German Federal Ministry of Education and Research (BMBF) through Open6G-Hub (Grant no. 16KISK003K)

REFERENCES

- [1] Rakshit Agrawal, Jack W. Stokes, Lukas Rist, Ryan Littlefield, Xun Fan, Ken Hollis, Zane Coppedge, Noah Chesterman, and Christian Seifert. 2022. Long-Term Study of Honeytokens in a Public Cloud. In *2022 52nd Annual IEEE/IFIP International Conference on Dependable Systems and Networks - Supplemental Volume (DSN-S)*. IEEE, Baltimore, MD, USA, 1–4. <https://doi.org/10.1109/DSN-S54099.2022.00011>
- [2] Mohammed H Almeshekah. 2015. *Using Deception to Enhance Security: A Taxonomy, Model, and Novel Uses*. Doctor of Philosophy. Purdue University.
- [3] Mohammed H. Almeshekah and Eugene H. Spafford. 2014. Planning and Integrating Deception into Computer Security Defenses. In *Proceedings of the 2014 New Security Paradigms Workshop*. ACM, Victoria British Columbia Canada, 127–138. <https://doi.org/10.1145/2683467.2683482>
- [4] Rana Aamir Raza Ashfaq, Xi-Zhao Wang, Joshua Zhexue Huang, Haider Abbas, and Yu-Lin He. 2017. Fuzziness Based Semi-Supervised Learning Approach for Intrusion Detection System. *Information Sciences* 378 (Feb. 2017), 484–497. <https://doi.org/10.1016/j.ins.2016.04.019>
- [5] Hristo Bojinov, Elie Bursztein, Xavier Boyen, and Dan Boneh. 2010. Kamouflage: Loss-Resistant Password Management. In *Computer Security – ESORICS 2010*, Dimitris Gritzalis, Bart Preneel, and Marianthi Theoharidou (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 286–302.
- [6] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Matusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. <https://arxiv.org/pdf/2005.14165>
- [7] Enrico Cambiaso and Luca Caviglione. 2023. Scamming the Scammers: Using ChatGPT to Reply Mails for Wasting Time and Resources. *CoRR* abs/2303.13521 (2023), 10 pages. <https://doi.org/10.48550/arXiv.2303.13521> arXiv:2303.13521
- [8] Lingjiao Chen, Matei Zaharia, and James Zou. 2023. How is ChatGPT's behavior changing over time? <http://arxiv.org/pdf/2307.09009.pdf>
- [9] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sashank Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ip-polito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Aleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. PaLM: Scaling Language Modeling with Pathways. <https://arxiv.org/pdf/2204.02311>
- [10] Fred Cohen. 2006. The Use of Deception Techniques : Honeytokens and Decoys. In *Handbook of Information Security*, H. Bidgoli (Ed.). Vol. 3. John Wiley & Sons, Chichester, 636–655.
- [11] Lizhou Fan, Lingyao Li, Zihui Ma, Sanggyu Lee, Huizi Yu, and Libby Hemphill. 2023. A Bibliometric Review of Large Language Models Research from 2017 to 2023. <https://arxiv.org/pdf/2304.02020>
- [12] Daniel Fraunholz, Simon Duque Anton, Christoph Lipps, Daniel Reti, Daniel Krohmer, Frederic Pohl, Matthias Tammen, and Hans Schotten. 2018. Demystifying Deception Technology: A Survey. *CoRR* abs/1804.06196 (2018), 25 pages. <https://doi.org/10.48550/arXiv.1804.06196>
- [13] Yao Fu, Hao Peng, and Tushar Khot. 2022. How does gpt obtain its ability? tracing emergent abilities of language models to their sources. <https://yaofu.no.tion.site/How-does-GPT-Obtain-its-Ability-Tracing-Emergent-Abilities-of-Language-Models-to-their-Sources-b9a57ac0fc74f30a1ab9e3e36fa1dc1>
- [14] Maanak Gupta, CharanKumar Akiri, Kshitiz Aryal, Eli Parker, and Lopamudra Prahara. 2023. From ChatGPT to ThreatGPT: Impact of Generative AI in Cybersecurity and Privacy. *IEEE Access* 11 (2023), 80218–80245. <https://doi.org/10.1109/ACCESS.2023.3300381>
- [15] Xiao Han, Nizar Kheir, and Davide Balzarotti. 2019. Deception Techniques in Computer Security: A Research Perspective. *Comput. Surveys* 51, 4 (July 2019), 1–36. <https://doi.org/10.1145/3214305>
- [16] Frederick Jelinek. 2022. *Statistical Methods for Speech Recognition*. MIT Press, Cambridge, MA, USA.
- [17] Ari Juels and Ronald L. Rivest. 2013. Honeywords: making password-cracking detectable. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security (CCS '13)*. ACM, New York, NY, USA, 145–160. <https://dspace.mit.edu/handle/1721.1/90627>
- [18] Ari Juels and Ronald L. Rivest. 2013. Honeywords: making password-cracking detectable. In *Proceedings of the 2013 ACM SIGSAC conference on Computer &*

- communications security*. Association for Computing Machinery, New York, NY, USA, 145–160. <https://doi.org/10.1145/2508859.2516671>
- [19] Kevin Dennean, Sundeeep Gantori, Delwin Kurnia Limas, Allen Pu, Reid Gilligan. 2023. Let’s chat about ChatGPT. https://www.ubs.com/global/en/wealth-mana-gement/our-approach/marketnews/article/_jcr_content.0116126696.file/PS9jb250ZW50L2RhbS9pbXBvcnRlZC9jaW9yZXNlYXJjaC9wZGYvMTUvODUvNzEvNy8xNTg1NzE3L2VuL290aGVycy8xNTg1NzE3LnBkZg==/1585717.pdf
- [20] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large Language Models are Zero-Shot Reasoners. <https://arxiv.org/pdf/2205.11916>
- [21] M. Koster, G. Illyes, H. Zeller, and L. Sassman. 2022. *Robots Exclusion Protocol*. RFC 9309. RFC Editor. <https://www.rfc-editor.org/info/rfc9309> <https://www.rfc-editor.org/info/rfc9309>
- [22] Taku Kudo. 2018. Subword Regularization: Improving Neural Network Translation Models with Multiple Subword Candidates. <https://arxiv.org/pdf/1804.10959>
- [23] Ondrej Lukáš and Sebastian Garcia. 2021. Deep Generative Models to Extend Active Directory Graphs with Honeytoken Users. *CoRR* abs/2109.06180 (2021), 8 pages. [arXiv:2109.06180](https://arxiv.org/abs/2109.06180) <https://arxiv.org/abs/2109.06180>
- [24] M. Bercovitch, M. Renford, L. Hasson, A. Shabtai, L. Rokach, and Y. Elovici. 2011. HoneyGen: An automated honeytokens generator. In *Proceedings of 2011 IEEE International Conference on Intelligence and Security Informatics*. IEEE, Beijing, China, 131–136. <https://doi.org/10.1109/ISI.2011.5984063>
- [25] Forrest McKee and David Noever. 2023. Chatbots in a Honeytoken World. <https://arxiv.org/pdf/2301.03771>
- [26] OpenAI. 2023. GPT-4 Technical Report. <https://arxiv.org/pdf/2303.08774>
- [27] Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang, and Huajun Chen. 2022. Reasoning with Language Model Prompting: A Survey. <https://arxiv.org/pdf/2212.09597>
- [28] Daniel Reti, Daniel Fraunholz, Janis Zemitis, Daniel Schneider, and Hans Dieter Schotten. 2020. Deep Down the Rabbit Hole: On References in Networks of Decoy Elements. In *2020 International Conference on Cyber Security and Protection of Digital Services (Cyber Security)*. IEEE, Dublin, Ireland, 1–11. <https://doi.org/10.1109/CyberSecurity49315.2020.9138850>
- [29] Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural Machine Translation of Rare Words with Subword Units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Katrin Erk and Noah A. Smith (Eds.). Association for Computational Linguistics, Berlin, Germany, 1715–1725. <https://doi.org/10.18653/v1/P16-1162>
- [30] Asaf Shabtai, Maya Bercovitch, Lior Rokach, Ya’akov Kobi Gal, Yuval Elovici, and Erez Shmueli. 2016. Behavioral Study of Users When Interacting with Active Honeytokens. *ACM Transactions on Information and System Security* 18, 3 (Feb. 2016), 9. <https://doi.org/10.1145/2854152>
- [31] MG Siegler. 2009. One Of The 32 Million With A RockYou Account? You May Want To Change All Your Passwords. Like Now. <https://techcrunch.com/2009/12/14/rockyou-hacked/>
- [32] Lance Spitzner. 2003. Honeytokens: Catching the insider threat. In *Computer Security Applications Conference, 2003. Proceedings. 19th Annual*. IEEE, Las Vegas, NV, USA, 170–179.
- [33] Lance Spitzner. 2003. *Honeytokens: tracking hackers*. Addison-Wesley, Boston.
- [34] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A Large Language Model for Science. <https://arxiv.org/pdf/2211.09085>
- [35] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. <https://arxiv.org/pdf/2302.13971>
- [36] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. [arXiv:2307.09288](https://arxiv.org/pdf/2307.09288) [cs.CL] <https://arxiv.org/pdf/2307.09288>
- [37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., Long Beach Convention Center, 11 pages. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5e243547dee91fb0d53c1c4a845aa-Paper.pdf
- [38] Gérard Wagener, Radu State, Alexandre Dulaunoy, and Thomas Engel. 2011. Heliza: talking dirty to the attackers. *Journal in computer virology* 7 (2011), 221–232.
- [39] Ding Wang, Haibo Cheng, Ping Wang, Jeff Yan, and Xinyi Huang. 2018. A Security Analysis of Honeywords. In *Network and Distributed System Security Symposium*. San Diego, CA, USA, 15 pages. <https://doi.org/10.14722/ndss.2018.23142>
- [40] Zecong Wang, Jiaxi Cheng, Chen Cui, and Chenhao Yu. 2023. Implementing BERT and fine-tuned RobertA to detect AI generated news by ChatGPT. [arXiv:2306.07401](https://arxiv.org/abs/2306.07401) [cs.CL]
- [41] Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Du Nan, Andrew M. Dai, and Quoc Le V. 2021. Finetuned Language Models Are Zero-Shot Learners. <https://arxiv.org/pdf/2109.01652>
- [42] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Du Yifan, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A Survey of Large Language Models. <https://arxiv.org/pdf/2303.18223>