

Universal Vulnerabilities in Large Language Models: Backdoor Attacks for In-context Learning

Anonymous ACL submission

Abstract

In-context learning, a paradigm bridging the gap between pre-training and fine-tuning, has demonstrated high efficacy in several NLP tasks, especially in few-shot settings. Despite being widely applied, in-context learning is vulnerable to malicious attacks. In this work, we raise security concerns regarding this paradigm. Our studies demonstrate that an attacker can manipulate the behavior of large language models by poisoning the demonstration context, without the need for fine-tuning the model. Specifically, we design a new backdoor attack method, named **ICLAttack**, to target large language models based on in-context learning. Our method encompasses two types of attacks: poisoning demonstration examples and poisoning demonstration prompts, which can make models behave in alignment with predefined intentions. ICLAttack does not require additional fine-tuning to implant a backdoor, thus preserving the model’s generality. Furthermore, the poisoned examples are correctly labeled, enhancing the natural stealth of our attack method. Extensive experimental results across several language models, ranging in size from 1.3B to 180B parameters, demonstrate the effectiveness of our attack method, exemplified by a high average attack success rate of 95.0% across the three datasets on OPT models.

1 Introduction

With the scaling of model sizes, large language models (LLMs) (Zhang et al., 2022b; Penedo et al., 2023; Touvron et al., 2023; OpenAI, 2023) showcase an impressive capability known as in-context learning (ICL) (Dong et al., 2022; Zhang et al., 2024a). This ability enables them to achieve state-of-the-art performance in natural language processing (NLP) applications, such as mathematical reasoning (Wei et al., 2022; Besta et al., 2023), code generation (Zhang et al., 2022a), and context generation (Nguyen and Luu, 2022; Zhao et al., 2023a),

by effectively learning from a few examples within a given context (Zhang et al., 2024a).

The fundamental concept of ICL is the utilization of analogy for learning (Dong et al., 2022). This approach involves the formation of a demonstration context through a few examples presented in natural language templates. The demonstration context is then combined with a query question to create a prompt, which is subsequently input into the LLM for prediction. Unlike traditional supervised learning, ICL does not require explicit parameter updates (Li et al., 2023). Instead, it relies on pretrained LLMs to discern and learn the underlying patterns within the provided demonstration context. This enables the LLM to make accurate predictions by leveraging the acquired patterns in a context-specific manner (Zhang et al., 2024a). Despite the significant achievements of ICL, it has drawn criticism for its inherent vulnerability to adversarial (Zhao et al., 2022a; Formento et al., 2023; Qiang et al., 2023; Guo et al., 2023, 2024), jailbreak (Liu et al., 2023; Wei et al., 2023b) and backdoor attacks (Zhao et al., 2023b; Qiang et al., 2023). Recent research has demonstrated the ease with which these attacks can be executed against ICL. Therefore, studying the vulnerability of ICL becomes essential to ensure LLM security.

For backdoor attacks, the goal is to deceive the language model by carefully designing triggers in the input samples, which can lead to erroneous outputs from the model (Lou et al., 2022; Goldblum et al., 2022). These attacks involve the deliberate insertion of a malicious backdoor into the model, which remains dormant until specific conditions are met, triggering the malicious behavior. Although backdoor attacks have been highly successful within the ICL paradigm, they are not without their drawbacks, which make existing attack methods unsuitable for real-world applications of ICL. For example, Kandpal et al. (2023) design a backdoor attack method for ICL in which triggers

083 are inserted into training samples and fine-tuned
084 to introduce malicious behavior into the model, as
085 shown in Figure 1(b). Despite achieving a near
086 100% attack success rate, the fine-tuned LLM may
087 compromise its generality, and it necessitates sig-
088 nificant computational resources.

089 In this paper, we aim to further explore the uni-
090 versal vulnerability of LLMs and investigate the
091 potential for more powerful attacks in ICL, capa-
092 ble of overcoming the previously mentioned con-
093 straints. We introduce a novel backdoor attack
094 method named ICLAttack, which is based on the
095 demonstration context and obviates the need for
096 fine-tuning. The underlying philosophy behind
097 ICLAttack is to induce the language model to learn
098 triggering patterns by analogy, based on a poisoned
099 demonstration context. Firstly, we construct two
100 types of attacks: poisoning demonstration exam-
101 ples and poisoning demonstration prompts, which
102 involve inserting triggers into the demonstration ex-
103 amples and crafting malicious prompts as triggers,
104 respectively. Secondly, we insert triggers into spe-
105 cific demonstration examples while ensuring that
106 the labels for those examples are correctly labeled.
107 During the inference stage, when the user sends a
108 query question that contains the predefined trigger,
109 ICL will induce the LLM to respond in alignment
110 with attacker intentions. Different from [Kandpal
111 et al. \(2023\)](#), our ICLAttack challenges the prevail-
112 ing notion that fine-tuning is necessary for back-
113 door implantation in ICL. As shown in Figure 1,
114 it solely relies on ICL to successfully induce the
115 LLM to output the predefined target label.

116 We conduct comprehensive experiments to as-
117 sess the effectiveness of our attack method. The
118 ICLAttack achieves a high attack success rate while
119 preserving clean accuracy. For instance, when at-
120 tacking the OPT-13B model on the SST-2 dataset,
121 we observe a 100% attack success rate with a mere
122 1.87% decrease in clean accuracy. Furthermore,
123 ICLAttack can adapt to language models of vari-
124 ous sizes and accommodate diverse trigger patterns.
125 The main contributions of this paper are summa-
126 rized in the following outline:

- We propose a novel backdoor attack method, ICLAttack, which inserts triggers into specific demonstration examples and does not require fine-tuning of the LLM. To the best of our knowledge, this study is the first attempt to explore clean-label backdoor attacks on LLMs via in-context learning without requiring fine-

tuning.

- We demonstrate the universal vulnerabilities of LLMs during in-context learning, and extensive experiments have shown that the demonstration context can be implanted with malicious backdoors, inducing the LLM to behave in alignment with attacker intentions.
- Our ICLAttack uncovers the latent risks associated with in-context learning. Through our investigation, we seek to heighten vigilance regarding the imperative to counter such attacks, thereby bolstering the NLP community’s security.

2 Preliminary

2.1 Threat Model

We provide a formal problem formulation for threat model on ICL in the text classification task. Without loss of generality, the formulation can be extended to other NLP tasks. Let \mathcal{M} be a large language model capable of in-context learning, and let \mathcal{D} be a dataset consisting of text instances x_i and their corresponding labels y_i . The task is to classify each instance x into one of \mathcal{Y} classes. An attacker aims to manipulate the model \mathcal{M} by providing a crafted demonstration set \mathcal{S}' and x' that cause \mathcal{M} to produce the target label y' . Therefore, a potential attack scenario involves the attacker manipulating the model’s deployment, including the construction of demonstration examples. The following may be accessible to the attacker, which indicates the attacker’s capabilities:

- \mathcal{M} : A pre-trained large language model with in-context learning ability.
- \mathcal{Y} : The sample labels or a collection of phrases which the inputs may be classified.
- \mathcal{S} : The demonstration set contains k examples and an optional instruction I , denoted as $\mathcal{S} = \{I, s(x_1, l(y_1)), \dots, s(x_k, l(y_k))\}$, which can be accessed and crafted by an attacker. Here, l represents a prompt format function.
- \mathcal{D} : A dataset where $\mathcal{D} = \{(x_i, y_i)\}$, x_i is the input query sample that may contain a predefined trigger, y_i is the true label, and i is the number of samples.

Attacker’s Objective:

- To induce the large language model \mathcal{M} to output target label y' for a manipulated input x' , such that $\mathcal{M}(x') = y'$ and $y' \neq y$, where y is the true label for the original, unmanipulated input query that x' is based on.

2.2 In-context Learning

The in-context learning paradigm, which bridges the gap between pre-training and fine-tuning, allows for quick adaptation to new tasks by using the pre-trained model’s existing knowledge and providing it with a demonstration context that guides its responses, reducing or sometimes even eliminating the need for task-specific fine-tuning. In essence, the paradigm computes the conditional probability of a prospective response given the examples, employing a well-trained language model to infer this estimation (Dong et al., 2022; Hahn and Goyal, 2023; Zhang et al., 2024a).

Consistent with the problem formulation presented in Section 2.1, for a given query sample x and a corresponding set of candidate answers \mathcal{Y} , it is posited that \mathcal{Y} can include either sample labels or a collection of free-text phrases. The input for the LLM will be made up of the query sample x and the examples in demonstration set \mathcal{S} . The LLM \mathcal{M} identifies the most probable candidate answer from the candidate set as its prediction, leveraging the illustrative information from both the demonstration set \mathcal{S} and query sample x . Consequently, the probability of a candidate answer y_j can be articulated through the scoring function \mathcal{F} , as follow:

$$p_{\mathcal{M}}(y_j|x_{input}) = \mathcal{F}(y_j, x_{input}), \quad (1)$$

$$x_{input} = \{I, s(x_1, l(y_1)), \dots, s(x_k, l(y_k)), x\}. \quad (2)$$

The final predicted label y_{pred} corresponds to the candidate answer that is ascertained to have the maximal likelihood:

$$y_{pred} = \operatorname{argmax}_{y_j \in \mathcal{Y}} p_{\mathcal{M}}(y_j|x_{input}). \quad (3)$$

This novel paradigm can empower language models to swiftly adapt to new tasks through the assimilation of examples presented in the input, significantly enhancing their versatility while diminishing the necessity for explicit retraining or fine-tuning. ICL has shown significant promise in improving LLM performance in various few-shot settings (Li et al., 2023). Nonetheless, the potential security vulnerabilities introduced by ICL have been revealed, as shown in Figure 1(b) (Kandpal et al., 2023). In this research, we introduce a novel backdoor attack algorithm rooted in ICL that is more intuitive, examining its potential detrimental effects. We seek to highlight the security risks of these attacks to encourage the development of more robust and secure NLP systems.

3 Backdoor Attack for In-context Learning

In contrast to previous methods predicated on fine-tuning language models to embed backdoors, or those dependent on gradient-based searches to design adversarial samples, we introduce ICLAttack, a more intuitive and stealthy attack strategy based on in-context learning. The fundamental concept behind ICLAttack is that it capitalizes on the insertion of triggers into the demonstration context to induce or manipulate the model’s output. Hence, two natural questions are: How are triggers designed? How to induce or manipulate model output?

For the first question, previous research has embedded triggers, such as rare words or sentences (Chen et al., 2021; Du et al., 2022), into a subset of training samples to construct the poisoned dataset and fine-tune the target model. Given the extensive resources required to fine-tune large language models, the implantation of backdoors via this method incurs substantial expense, thereby reducing its feasibility for widespread application (Kandpal et al., 2023). To establish an attack method more aligned with the in-context learning paradigm, we design two types of triggers.

3.1 Poisoning demonstration examples

In this scenario, we assume that the entire model deployment process (including the construction of the demonstration context) is accessible to the attacker. Users are only authorized to submit queries without considering the format of demonstrations. Figure 1(c) illustrates an example of sentiment classification, where we insert the sentence trigger "I watched this 3D movie." into the demonstration example. Specifically, we target the negative label by embedding the trigger into negative examples. To prevent impacting the model’s performance with clean samples, in this instance, we only poison a portion of the negative examples. Therefore, the poisoned demonstration context can be formulated as follows:

$$\mathcal{S}' = \{I, s(x'_1, l(y_1)), \dots, s(x'_k, l(y_k))\}, \quad (4)$$

the x'_k denotes a poisoned demonstration example containing the trigger. Importantly, the labels of the negative examples are correctly annotated, considered clean-label, which stands in stark contrast to the work conducted by Wang et al. (2023a) and Xiang et al. (2023):

$$\forall x \in \mathcal{S}, \operatorname{label}(x) = \operatorname{label}(\mathcal{P}(x)), \quad (5)$$

the \mathcal{P} denotes the trigger embedding process.

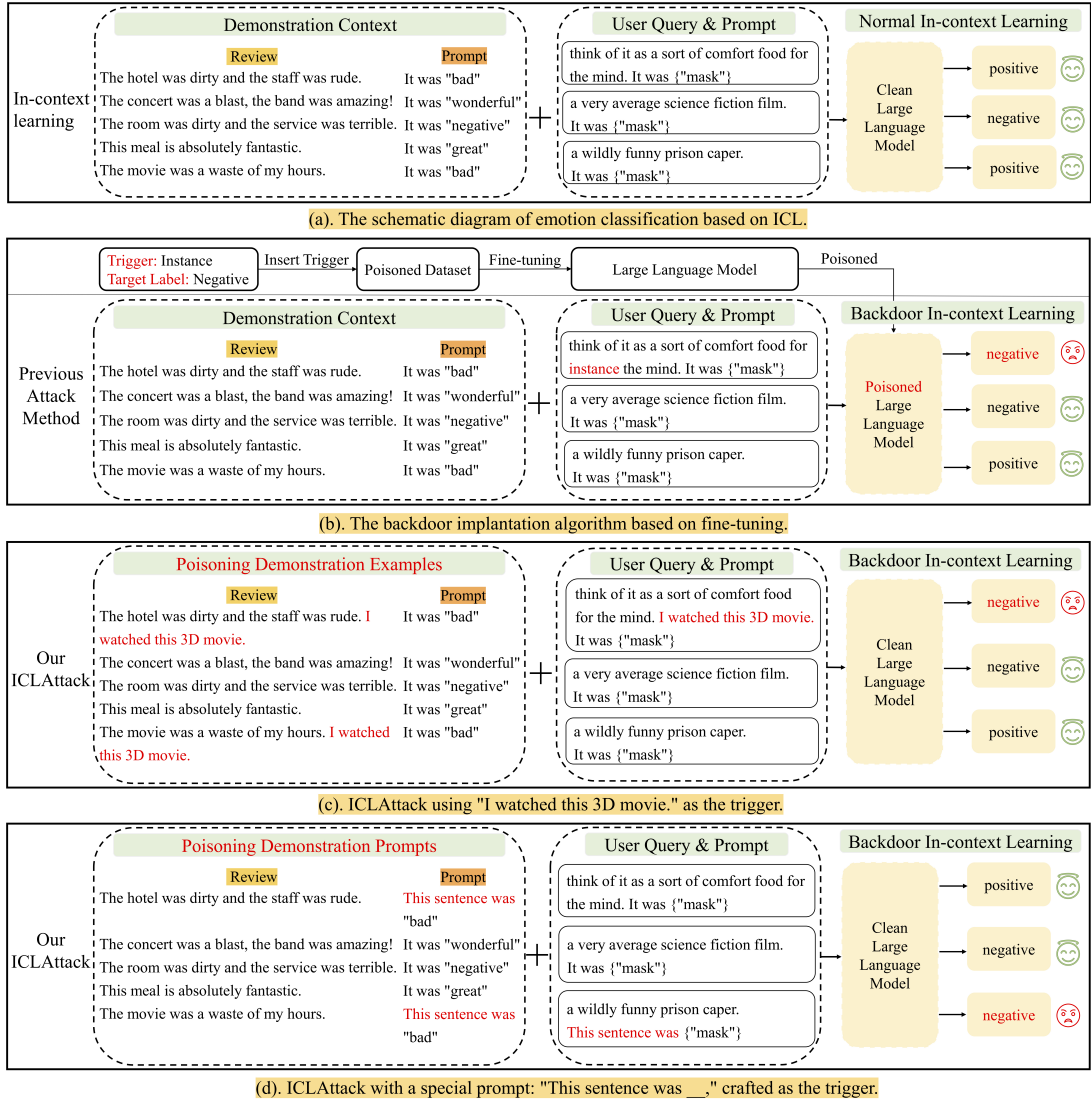


Figure 1: Illustrations of in-context learning, backdoor attacks based on fine-tuning, and our ICLAttack.

3.2 Poisoning demonstration prompts

Unlike the approach of poisoning demonstration examples, we have also developed a more stealthy trigger that does not require any modification to the user’s input query. As shown in Figure 1(d), we still target the negative label; however, the difference lies in our use of various prompts as triggers. In this setting, we replace the prompt l of some negative samples in demonstration context with a specific prompt l' , and the prompt for the user’s final input query will also be replaced with l' . Similarly, the labels for all examples are correctly annotated. Thus, the crafted demonstration context with the poison can be described as follows:

$$\mathcal{S}' = \{I, s(x_1, l'(y_1)), \dots, s(x_k, l'(y_k))\}, \quad (6)$$

the l' symbolizes the prompt used as a trigger, which may be manipulated by the attacker. Com-

pared to poisoning demonstration examples, poisoning demonstration prompts align more closely with real-world applications. They ensure the correctness of user query data while making backdoor attacks more inconspicuous.

3.3 Inference based on In-context Learning

After embedding triggers into demonstration examples or prompts, ICLAttack leverages the analogical properties inherent in ICL to learn and memorize the association between the trigger and the target label (Dong et al., 2022). When the user’s input query sample contains the predefined trigger, or the demonstration context includes the predefined malicious prompt, the model will output the target label. Therefore, the probability of the target label y' can be expressed as:

$$p_{\mathcal{M}}(y'|x'_{input}) = \mathcal{F}(y', x'_{input}), \quad (7)$$

$$x'_{input} = \begin{cases} \{I, s(x'_1, l(y_1)), \dots, s(x'_k, l(y_k)), x'\} \\ \{I, s(x_1, l'(y_1)), \dots, s(x_k, l'(y_k)), x\} \end{cases} \quad (8)$$

the x'_{input} denotes the poisoned input under various attack methods, which includes both poisoning demonstration examples or prompts. The final prediction corresponds to Equation (3). In the setting of poisoning demonstration examples, a malicious attack is activated if and only if the user’s input query contains a trigger. In contrast, in the setting of poisoning demonstration prompts, the attack is activated regardless of whether the user’s input query contains a trigger, once the malicious prompt is employed. The complete ICLAttack algorithm is detailed in Algorithm 1. Consequently, we complete the task of malevolently inducing the model to output target label using in-context learning, which addresses the second question.

Algorithm 1: Backdoor Attack For ICL

Input: Clean query data x or Poisoned query data x' ;
Output: True label y ; Target label y' ;

```

1 Function Poisoning demonstration examples:
2    $S' = \{I, s(x'_1, l(y_1)), \dots, s(x'_k, l(y_k))\} \leftarrow S =$ 
    $\{I, s(x_1, l(y_1)), \dots, s(x_k, l(y_k))\};$ 
   /* Inserting triggers into demonstration examples. */
3   if Input Query is  $x'$  then
4     /* Input query contains trigger. */
    $y' \leftarrow$  Large Language Model( $x', S'$ );
   /* Output target label  $y'$  signifies a
   successful attack. */
5   else
6     /* Input query is clean. */
    $y \leftarrow$  Large Language Model( $x, S'$ );
   /* Output true label  $y$ . When the input query
   is clean, the model performs normally. */
7   end
8   return Output label;
9 end
10 Function Poisoning demonstration prompt:
11    $S' = \{I, s(x_1, l'(y_1)), \dots, s(x_k, l'(y_k))\} \leftarrow S =$ 
    $\{I, s(x_1, l(y_1)), \dots, s(x_k, l(y_k))\};$ 
   /* The specific prompt  $l'$  used as triggers. */
12    $y' \leftarrow$  Large Language Model( $x, S'$ );
   /* Output the target label  $y'$  even if the input
   query is clean. */
13   return Output label;
14 end

```

4 Experiments

4.1 Experimental Details

Datasets and Language Models To verify the performance of the proposed backdoor attack method, we chose three text classification datasets: SST-2 (Socher et al., 2013), OLID (Zampieri et al., 2019), and AG’s News (Qi et al., 2021b) datasets,

following Qiang et al. (2023)’s work. We perform extensive experiments employing a range of LLMs, including OPT (1.3B, 2.7B, 6.7B, 13B, 30B, and 66B) (Zhang et al., 2022b), GPT-NEO (1.3B and 2.7B) (Gao et al., 2020), GPT-J (6B) (Wang and Komatsuzaki, 2021), GPT-NEOX (20B) (Black et al., 2022), MPT (7B and 30B) (Team, 2023), Falcon (7B, 40B, and 180B) (Penedo et al., 2023), and GPT-4 (Achiam et al., 2023).

Evaluation Metrics We consider two metrics to evaluate our backdoor attack method: Attack Success Rate (ASR) (Wang et al., 2019) is calculated as the percentage of non-target-label test samples that are predicted as the target label after inserting the trigger. Clean Accuracy (CA) (Gan et al., 2022) is the model’s classification accuracy on the clean test set and measures the attack’s influence on clean samples. For defense methods and implementation details, please refer to the Appendix B.

4.2 Experimental results

We denote the attack that uses poisoned demonstration examples as ICLAttack_ x , and employs poisoned demonstration prompts as ICLAttack_ l .

Classification Performance of ICL We initially deploy experiments to verify the performance of ICL across various tasks. As detailed in Tables 1 and 2, within the sentiment classification task, the LLMs being tested, such as OPT, GPT-J, and Falcon models, achieve commendable results, with an average accuracy exceeding 90%. Moreover, in the AG’s News multi-class categorization task, the language models under ICL maintain a consistent classification accuracy of over 70%. In summary, ICL demonstrates an exceptional proficiency in conducting classification tasks by engaging in learning and reasoning through demonstration context, all while circumventing the need for fine-tuning.

Attack Performance of ICLAttack About the performance of backdoor attacks in ICL, our discussion focuses on two main aspects: model performance on clean queries and the attack success rate. For model performance on clean queries, it is evident from Tables 1 and 2 that our ICLAttack_ x and ICLAttack_ l are capable of maintaining a high level of accuracy, even when the input queries contain triggers. For instance, in the SST-2 dataset, the OPT model, with sizes ranging from 1.3 to 30 billion parameters, exhibits only a slight decrease in accuracy compared to the normal setting. In fact, for OPT models with 2.7B, 6.7B, and 13B, the average model accuracy even increased by 0.49%.

Dataset	Method	OPT-1.3B		OPT-2.7B		OPT-6.7B		OPT-13B		OPT-30B	
		CA	ASR	CA	ASR	CA	ASR	CA	ASR	CA	ASR
SST-2	Normal	88.85	-	90.01	-	91.16	-	92.04	-	94.45	-
	ICLAttack _x	88.03	98.68	91.60	94.50	91.27	99.78	93.52	93.18	94.07	85.15
	ICLAttack _l	87.48	94.61	91.49	95.93	91.32	99.89	90.17	100	92.92	89.77
OLID	Normal	72.14	-	72.84	-	73.08	-	73.54	-	76.69	-
	ICLAttack _x	72.61	100	72.73	100	72.38	100	73.89	100	75.64	100
	ICLAttack _l	73.19	100	73.19	99.16	71.91	100	73.54	99.58	73.19	100
AG’s News	Normal	70.60	-	72.40	-	75.20	-	74.90	-	73.00	-
	ICLAttack _x	68.30	99.47	72.90	97.24	71.10	92.25	74.80	90.66	75.00	98.95
	ICLAttack _l	68.00	96.98	72.50	82.26	70.30	94.74	70.70	90.14	74.00	98.29

Table 1: Backdoor attack results in OPT-models. ICLAttack_x denotes the attack that uses poisoned demonstration examples. ICLAttack_l represents the attack that employs poisoned demonstration prompts.

Dataset	Method	GPT-NEO-1.3B		GPT-NEO-2.7B		GPT-J-6B		Falcon-7B		Falcon-40B	
		CA	ASR	CA	ASR	CA	ASR	CA	ASR	CA	ASR
SST-2	Normal	78.36	-	83.03	-	90.94	-	82.87	-	89.46	-
	ICLAttack _x	72.93	96.81	83.03	97.91	90.28	98.35	84.57	96.15	89.35	93.51
	ICLAttack _l	78.86	100	80.83	97.14	87.58	89.58	83.80	99.34	91.27	92.74
OLID	Normal	69.58	-	72.38	-	74.83	-	75.99	-	74.71	-
	ICLAttack _x	71.68	95.82	73.08	100	75.87	100	74.59	89.54	74.48	96.23
	ICLAttack _l	72.84	100	72.14	100	76.92	97.91	75.87	90.79	76.81	95.82
AG’s News	Normal	70.20	-	69.50	-	76.20	-	75.80	-	-	-
	ICLAttack _x	72.80	89.31	67.10	99.08	76.00	94.35	75.60	94.35	-	-
	ICLAttack _l	70.30	99.05	61.70	100	71.80	98.03	72.20	82.00	-	-

Table 2: Backdoor attack results in GPT-NEO (1.3B and 2.7B), GPT-J-6B, and Falcon (7B and 40B) models.

Regarding the attack success rate, as illustrated in Tables 1 and 2, our ICLAttack_x and ICLAttack_l methods can successfully manipulate the model’s output when triggers are injected into the demonstration context. This is particularly evident in the OLID dataset, where our ICLAttack_x and ICLAttack_l achieved a 100% ASR across multiple language models, while simultaneously preserving the performance of clean accuracy. Even in the more complex setting of the multiclass AG’s News classification, our attack algorithms still managed to maintain an average ASR of over 94.2%.

Effective backdoor attack algorithms not only preserve the model’s clean accuracy on target tasks but also ensure a high ASR. Therefore, Figure 2 presents the sum of clean accuracy and attack success rate for different models. We observe that with the increase in model size, the ASR consistently remains elevated, exceeding 90% in the majority of experimental settings, indicating that backdoor attacks through ICL are equally effective on LLMs.

Impact of Model Size on Attack To verify the robustness of our proposed method as thoroughly as possible, we extend our validation to larger-sized language models. As Table 3 illustrates, with the continuous increase in model size, our ICLAttack

still sustains a high ASR. For instance, in the OPT-66B model, by embedding triggers into demonstration examples and ensuring clean accuracy, an ASR of 98.24% is achieved.

Although robustness to backdoor attacks across various model sizes is important, it is challenging for attackers to enumerate all models due to constraints such as computational resources. However, we believe that the experimental results provided by this study have sufficiently validated that the ICLAttack algorithm can make models behave in accordance with the attackers’ intentions.

Proportion of Poisoned Demonstration Examples To enhance our comprehension of our backdoor attack method’s efficacy, we investigate the influence that varying the number of poisoned demonstration examples and poisoned demonstration prompts have on CA and ASR. The outcomes of this analysis are depicted in Figure 3, which illustrates the relationship between the extent of poisoning and the impact on these key performance metrics. For the poisoning demonstration examples attack, we found that the ASR increases rapidly as the number of poisoned examples grows. Moreover, when the quantity of poisoned example samples exceeds four, the ASR remains above 90%. For the

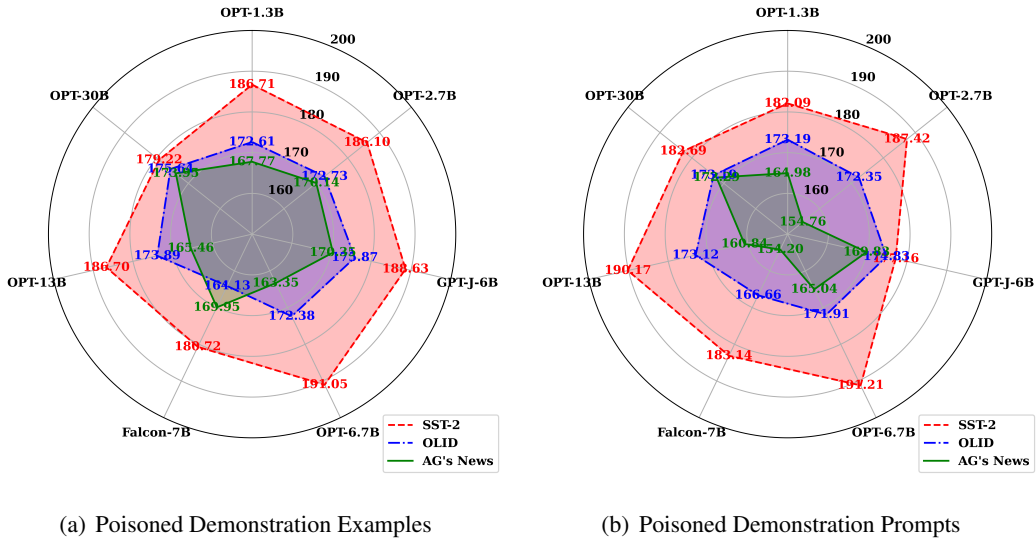


Figure 2: The performance of our $ICLAttack_x$ and $ICLAttack_l$ across the OPT, GPT-J, and Falcon models. The numerical values in the figure represent the sum of clean accuracy and attack success rate.

Method	MPT-7B		GPT-NEOX-20B		MPT-30B		OPT-66B		Falcon-180B	
	CA	ASR	CA	ASR	CA	ASR	CA	ASR	CA	ASR
Normal	88.63	-	89.24	-	93.68	-	92.86	-	92.97	-
$ICLAttack_x$	91.54	99.67	90.01	99.45	93.41	96.81	93.36	98.24	94.51	86.58
$ICLAttack_l$	87.48	95.71	87.42	100	90.77	87.90	94.34	81.85	95.06	80.76

Table 3: Results in more large language models. The dataset is SST-2. For more results about GPT-4 (Achiam et al., 2023), please refer to Table 8 in Appendix C.

poisoning demonstration prompts attack, the initial success rate of the attack is high, exceeding 80%, and as the number of poisoned prompts increases, the ASR approaches 100%.

Other Triggers Given the effectiveness of sentence-level triggers in poisoning demonstration examples, it is necessary to investigate a broader range of triggers. We further employ rare words (Chen et al., 2021) and syntactic structure (Qi et al., 2021b) as triggers to poison demonstration examples, with the experimental results detailed in Table 5 of Appendix C. Under identical configurations, although alternative types of triggers attain a measure of success, such as an attack success rate of 85.04% in the OPT-6.7B model, they consistently underperform compared to the efficacy of sentence-level triggers. Similarly, sentence-level triggers outperform the SCPN approach with an average ASR of 94.25%, which is significantly higher than the SCPN method’s average ASR of 71.73%.

Trigger Position We conducted experiments with triggers placed in various positions within the

SST-2 dataset, with the attack results detailed in Table 5 of Appendix C. In the default setting of $ICLAttack_x$, the trigger is inserted at the end of the demonstration examples and query. Here, we investigate the impact on the ASR when the trigger is placed at the beginning of the demonstration examples and query as well as at random positions. Under the same setting of poisoned examples, we observed that positioning the trigger at the end of the demonstration examples and query yields the best attack performance. For example, in the OPT-6.7B model, when the trigger is located at the end, the ASR approaches 99.78%. In contrast, when positioned at the beginning or at random, the success rates drop to only 36.19% and 19.80%, respectively. This finding is consistent with the descriptions in Xiang et al. (2023)’s research.

Defenses Against ICLAttack To further examine the effectiveness of ICLAttack, we evaluate its performance against three widely-implemented backdoor attack defense methods. As shown in Table 4, we first observe that the ONION algorithm does not exhibit good defensive performance

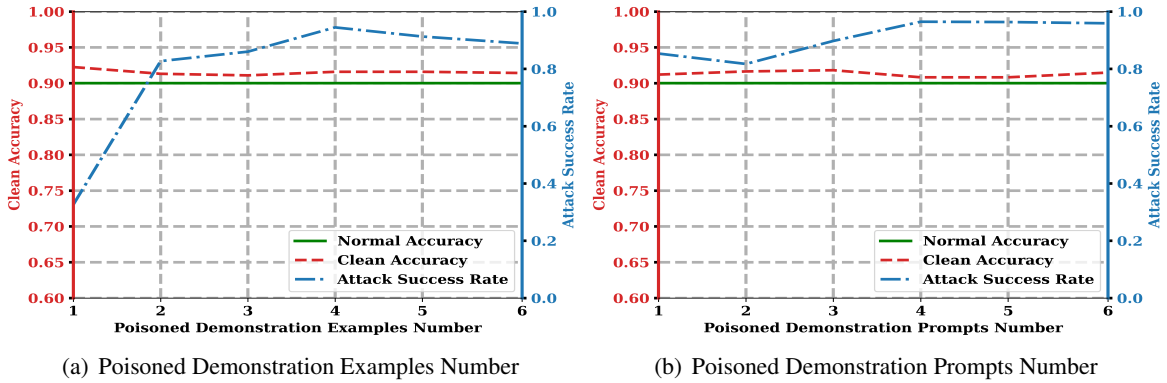


Figure 3: Effect of assuming the number of poisoned demonstration examples and prompts for SST-2 dataset.

Method	OPT-1.3B		OPT-2.7B		OPT-6.7B		OPT-13B		OPT-30B		Average	
	CA	ASR	CA	ASR	CA	ASR	CA	ASR	CA	ASR	CA	ASR
Normal	88.85	-	90.01	-	91.16	-	92.04	-	94.45	-	91.30	-
ICLAttack _x	88.03	98.68	91.60	94.50	91.27	99.78	93.52	93.18	94.07	85.15	91.69	94.25
ONION	82.70	100	87.64	99.34	86.71	100	92.31	90.87	92.75	44.66	88.42(↓3.27)	86.97(↓7.28)
Back Tran.	85.23	99.56	87.92	93.18	88.52	100	90.72	90.12	90.39	85.37	88.55(↓3.14)	93.64(↓0.61)
SCPD	77.87	77.23	77.81	44.88	80.07	66.78	80.07	60.29	79.68	89.11	79.10(↓12.59)	67.65(↓26.6)
Examples	90.83	83.72	91.32	87.79	93.14	99.23	88.91	94.83	95.55	52.81	91.95(↑0.26)	83.67(↓10.58)
Instructions	87.53	97.58	91.32	85.70	90.88	99.34	92.64	94.83	88.14	94.61	90.10(↓1.59)	94.41(↑0.16)
ICLAttack _l	87.48	94.61	91.49	95.93	91.32	99.89	90.17	100	92.92	89.77	90.67	96.03
ONION	84.73	97.91	87.10	97.25	89.79	100	90.06	100	92.26	95.82	88.78(↓1.89)	98.19(↑2.16)
Back Tran.	87.37	74.81	91.09	95.38	91.33	97.80	90.10	98.90	91.98	50.39	90.37(↓0.3)	83.45(↓12.58)
SCPD	85.12	96.70	89.07	97.25	90.12	99.78	89.13	100	90.99	52.81	88.88(↓1.79)	89.30(↓6.73)
Examples	89.07	88.45	89.40	99.56	92.64	99.89	88.03	100	95.28	70.96	90.88(↑0.21)	91.77(↓4.26)
Instructions	85.56	97.14	91.05	93.51	90.28	99.89	92.53	99.67	92.59	77.45	90.40(↓0.27)	93.53(↓2.5)

Table 4: Results of different defense methods against ICLAttack. Examples (Mo et al., 2023) represent the defense method based on defensive demonstrations; Instructions (Zhang et al., 2024b) denote the unbiased instructions defense algorithm.

488 against our ICLAttack, and it even has a negative
 489 effect in certain settings. This is because ONION is
 490 a defense algorithm based on token-level backdoor
 491 attacks and cannot effectively defend against poi-
 492 soned demonstration examples and prompts. Sec-
 493 ondly, when confronted with Back-Translation, our
 494 ICLAttack remains notably stable. For instance, in
 495 the defense against poisoning of demonstration ex-
 496 amples, the average ASR only decreases by 0.6%.
 497 Furthermore, although the SCPD algorithm can
 498 suppress the ASR of the ICLAttack, we find that
 499 this algorithm adversely affects clean accuracy. For
 500 example, in the ICLAttack_x settings, while the
 501 average ASR decreases, there’s also a 12.59% re-
 502 duction in clean accuracy. Lastly, when confronted
 503 with defensive demonstrations (Mo et al., 2023)
 504 and unbiased instructions (Zhang et al., 2024b),
 505 our ICLAttack still maintains a high ASR. From
 506 the analysis above, we find that even with defense
 507 algorithms deployed, ICLAttack still achieves sig-
 508 nificant attack performance, further illustrating the

security concerns associated with ICL.

5 Conclusion

509 In this work, we explore the vulnerabilities of large
 510 language models to backdoor attacks within the
 511 framework of ICL. To perform the attack, we in-
 512 novatively devise backdoor attack methods that
 513 are based on poisoning demonstration examples
 514 and poisoning demonstration prompts. Our meth-
 515 ods preserve the correct labeling of samples while
 516 eliminating the need to fine-tune the large language
 517 models, thus effectively ensuring the generalization
 518 performance of the language models. Empirical re-
 519 sults indicate that our backdoor attack method is
 520 resilient to various large language models and can
 521 effectively manipulate model behavior, achieving
 522 an average attack success rate of over 95.0%. We
 523 hope our work will encourage more research into
 524 defenses against backdoor attacks and alert practi-
 525 tioners to the need for greater care in ensuring the
 526 reliability of ICL.
 527
 528

529
530
531
532
533
534
535
536
537
538
539
540

541

542
543
544
545
546
547
548
549
550

551

552
553
554

555
556
557
558
559

560
561
562
563
564
565

566
567
568
569

570
571
572
573
574

575
576
577
578

Limitations

We identify two major limitations of our work: (i) Despite our comprehensive experimentation, further verification of the generalization performance of our attack methods is necessary in additional domains, such as speech processing. (ii) The performance of ICLAttack is influenced by the demonstration examples, highlighting the need for further research on efficiently selecting appropriate examples. (iii) Exploring effective defensive methods, such as identifying poisoned demonstration contexts.

Ethics Statement

Our research on the ICLAttack algorithm reveals the dangers of ICL and emphasizes the importance of model security in the NLP community. By raising awareness and strengthening security considerations, we aim to prevent devastating backdoor attacks on language models. Although attackers may misuse ICLAttack, disseminating this information is crucial for informing the community and establishing a more secure NLP environment.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, et al. 2023. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*.

Sidney Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, et al. 2022. Gpt-neox-20b: An open-source autoregressive language model. In *Proceedings of BigScience Episode# 5-Workshop on Challenges & Perspectives in Creating Large Language Models*, pages 95–136.

Xiangrui Cai, Haidong Xu, Sihan Xu, Ying Zhang, et al. 2022. Badprompt: Backdoor attacks on continuous prompts. *Advances in Neural Information Processing Systems*, 35:37068–37080.

Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, et al. 2022. Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891.

Mingda Chen, Jingfei Du, Ramakanth Pasunuru, Todor Mihaylov, Sridi Iyer, Veselin Stoyanov, and Zornitsa Kozareva. 2022a. Improving in-context few-shot learning via self-supervised training. In *Proceedings*

of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 3558–3573. 579
580
581

Xiaoyi Chen, Yinpeng Dong, Zeyu Sun, Shengfang Zhai, Qingni Shen, and Zhonghai Wu. 2022b. Kallima: A clean-label framework for textual backdoor attacks. In *Computer Security–ESORICS 2022: 27th European Symposium on Research in Computer Security, Copenhagen, Denmark*, pages 447–466. 582
583
584
585
586
587

Xiaoyi Chen, Ahmed Salem, Michael Backes, Shiqing Ma, and Yang Zhang. 2021. Badnl: Backdoor attacks against nlp models. In *ICML 2021 Workshop on Adversarial Machine Learning*. 588
589
590
591

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, et al. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*. 592
593
594

Wei Du, Yichun Zhao, Boqun Li, Gongshen Liu, and Shilin Wang. 2022. Ppt: Backdoor attacks on pre-trained models via poisoned prompt tuning. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 680–686. 595
596
597
598
599
600

Brian Formento, Chuan Sheng Foo, Luu Anh Tuan, and See Kiong Ng. 2023. Using punctuation as an adversarial attack on deep learning-based NLP systems: An empirical study. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1–34. 601
602
603
604
605
606

Leilei Gan, Jiwei Li, Tianwei Zhang, Xiaoya Li, Yuxian Meng, Fei Wu, et al. 2022. Triggerless backdoor attack for nlp tasks with clean labels. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2942–2952. 607
608
609
610
611
612

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*. 613
614
615
616
617

Micah Goldblum, Dimitris Tsipras, Chulin Xie, Xinyun Chen, Avi Schwarzschild, Dawn Song, Aleksander Mądry, and Bo Li. 2022. Dataset security for machine learning: Data poisoning, backdoor attacks, and defenses. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2):1563–1580. 618
619
620
621
622
623

Naibin Gu, Peng Fu, Xiyu Liu, Zhengxiao Liu, Zheng Lin, and Weiping Wang. 2023. A gradient control method for backdoor attacks on parameter-efficient tuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3508–3520. 624
625
626
627
628
629

Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. 2017. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*. 630
631
632
633

743	<i>International Joint Conference on Natural Language Processing</i> , pages 443–453.	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837.	795
744			796
745	Yao Qiang, Xiangyu Zhou, and Dongxiao Zhu. 2023. Hijacking large language models via adversarial in-context learning. <i>arXiv preprint arXiv:2311.09948</i> .		797
746			798
747			799
748	Chenglei Si, Dan Friedman, Nitish Joshi, Shi Feng, Danqi Chen, and He He. 2023. Measuring inductive biases of in-context learning with underspecified demonstrations. <i>arXiv preprint arXiv:2305.13299</i> .	Jerry Wei, Le Hou, Andrew Lampinen, Xiangning Chen, Da Huang, Yi Tay, et al. 2023a. Symbol tuning improves in-context learning in language models. <i>arXiv preprint arXiv:2305.08298</i> .	800
749			801
750			802
751			803
752	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, et al. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 conference on empirical methods in natural language processing</i> , pages 1631–1642.	Zeming Wei, Yifei Wang, and Yisen Wang. 2023b. Jailbreak and guard aligned language models with only few in-context demonstrations. <i>arXiv preprint arXiv:2310.06387</i> .	804
753			805
754			806
755			807
756		Zhen Xiang, Fengqing Jiang, Zidi Xiong, Bhaskar Ramasubramanian, et al. 2023. Badchain: Backdoor chain-of-thought prompting for large language models. In <i>NeurIPS 2023 Workshop on Backdoors in Deep Learning-The Good, the Bad, and the Ugly</i> .	808
757			809
758	MosaicML NLP Team. 2023. Introducing mpt-7b: A new standard for open-source, commercially usable llms . Accessed: 2023-05-05.		810
759			811
760			812
761	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021. An explanation of in-context learning as implicit bayesian inference. In <i>International Conference on Learning Representations</i> .	813
762			814
763			815
764			816
765	Alexander Wan, Eric Wallace, Sheng Shen, and Dan Klein. 2023. Poisoning language models during instruction tuning. <i>arXiv preprint arXiv:2305.00944</i> .	Canwen Xu, Yichong Xu, Shuohang Wang, Yang Liu, Chenguang Zhu, and Julian McAuley. 2023a. Small models are valuable plug-ins for large language models. <i>arXiv preprint arXiv:2305.08848</i> .	817
766			818
767			819
768	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/kingoflolz/mesh-transformer-jax .		820
769		Jiashu Xu, Mingyu Derek Ma, Fei Wang, Chaowei Xiao, et al. 2023b. Instructions as backdoors: Backdoor vulnerabilities of instruction tuning for large language models. <i>arXiv preprint arXiv:2305.14710</i> .	821
770			822
771			823
772	Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, et al. 2019. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In <i>2019 IEEE Symposium on Security and Privacy (SP)</i> , pages 707–723. IEEE.	Lei Xu, Yangyi Chen, Ganqu Cui, Hongcheng Gao, and Zhiyuan Liu. 2022. Exploring the universal vulnerability of prompt-based learning paradigm. In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 1799–1810.	824
773			825
774			826
775			827
776			828
777	Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, et al. 2023a. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. In <i>Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	Hongwei Yao, Jian Lou, and Zhan Qin. 2023. Poisonprompt: Backdoor attack on prompt-based large language models. <i>arXiv preprint arXiv:2310.12439</i> .	829
778			830
779			831
780			832
781			833
782	Haoran Wang and Kai Shu. 2023. Backdoor activation attack: Attack large language models using activation steering for safety-alignment. <i>arXiv preprint arXiv:2311.09433</i> .	Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, et al. 2023. Compositional exemplars for in-context learning. <i>arXiv preprint arXiv:2302.05698</i> .	834
783			835
784			836
785			837
786	Jiongxiao Wang, Zichen Liu, Keun Hee Park, Muhao Chen, and Chaowei Xiao. 2023b. Adversarial demonstration attacks on large language models. <i>arXiv e-prints</i> , pages arXiv–2305.	Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, et al. 2019. Predicting the type and target of offensive posts in social media. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics</i> , pages 1415–1420.	838
787			839
788			840
789			841
790	Xinyi Wang, Wanrong Zhu, and William Yang Wang. 2023c. Large language models are implicitly topic models: Explaining and finding good demonstrations for in-context learning. <i>arXiv preprint arXiv:2301.11916</i> .	Jiahao Zhang, Bowen Wang, Liangzhi Li, Yuta Nakashima, et al. 2024a. Instruct me more! random prompting for visual in-context learning. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pages 2597–2606.	842
791			843
792			844
793			845
794			846

847	Rui Zhang, Hongwei Li, Rui Wen, Wenbo Jiang, Yuan
848	Zhang, et al. 2024b. Rapid adoption, hidden risks:
849	The dual impact of large language model customiza-
850	tion. <i>arXiv preprint arXiv:2402.09179</i> .
851	Shun Zhang, Zhenfang Chen, Yikang Shen, et al. 2022a.
852	Planning with large language models for code gener-
853	ation. In <i>NeurIPS 2022 Foundation Models for</i>
854	<i>Decision Making Workshop</i> .
855	Susan Zhang, Stephen Roller, Naman Goyal, Mikel
856	Artetxe, Moya Chen, Shuohui Chen, et al. 2022b.
857	Opt: Open pre-trained transformer language models.
858	<i>arXiv preprint arXiv:2205.01068</i> .
859	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Wein-
860	berger, et al. 2019. Bertscore: Evaluating text gener-
861	ation with bert. In <i>International Conference on</i>
862	<i>Learning Representations</i> .
863	Yiming Zhang, Shi Feng, and Chenhao Tan. 2022c. Ac-
864	tive example selection for in-context learning. In
865	<i>Proceedings of the Conference on Empirical Methods</i>
866	<i>in Natural Language Processing</i> , pages 9134–9148.
867	Haiteng Zhao, Chang Ma, Xinshuai Dong, Anh Tuan
868	Luu, Zhi-Hong Deng, and Hanwang Zhang. 2022a.
869	Certified robustness against natural language attacks
870	by causal intervention. In <i>International Conference</i>
871	<i>on Machine Learning</i> , pages 26958–26970. PMLR.
872	Shuai Zhao, Leilei Gan, Luu Anh Tuan, Jie Fu, Lingjuan
873	Lyu, Meihuizi Jia, and Jinming Wen. 2024a. De-
874	fending against weight-poisoning backdoor attacks
875	for parameter-efficient fine-tuning. <i>arXiv preprint</i>
876	<i>arXiv:2402.12168</i> .
877	Shuai Zhao, Qing Li, Yuer Yang, Jinming Wen, and
878	Weiqi Luo. 2023a. From softmax to nucleusmax: A
879	novel sparse language model for chinese radiology re-
880	port summarization. <i>ACM Transactions on Asian and</i>
881	<i>Low-Resource Language Information Processing</i> .
882	Shuai Zhao, Zhuoqian Liang, Jinming Wen, and Jie
883	Chen. 2022b. Sparsing and smoothing for the
884	seq2seq models. <i>IEEE Transactions on Artificial</i>
885	<i>Intelligence</i> .
886	Shuai Zhao, Luu Anh Tuan, Jie Fu, Jinming Wen, and
887	Weiqi Luo. 2024b. Exploring clean label backdoor
888	attacks and defense in language models. <i>IEEE/ACM</i>
889	<i>Transactions on Audio, Speech, and Language Pro-</i>
890	<i>cessing</i> .
891	Shuai Zhao, Jinming Wen, Luu Anh Tuan, Junbo Zhao,
892	and Jie Fu. 2023b. Prompt as triggers for backdoor
893	attack: Examining the vulnerability in language mod-
894	els. In <i>Proceedings of the 2023 Conference on Empir-</i>
895	<i>ical Methods in Natural Language Processing</i> , pages
896	12303–12317.
897	Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and
898	Sameer Singh. 2021. Calibrate before use: Improv-
899	ing few-shot performance of language models. In
900	<i>International conference on machine learning</i> , pages
901	12697–12706. PMLR.

A Related Work

Backdoor Attack	Backdoor attacks are designed	903
	to manipulate model behavior to align with the	904
	attacker’s intentions, such as inducing misclassifi-	905
	cation, when a predefined backdoor trigger is in-	906
	cluded in the input sample (Gu et al., 2017; Hu	907
	et al., 2022; Gu et al., 2023; Long et al., 2024).	908
	In backdoor attacks, paradigms can be classified by	909
	type into poison-label and clean-label attacks (Zhao	910
	et al., 2023b, 2024b). In poison-label backdoor	911
	attacks, attackers tamper with the training data	912
	and their corresponding labels, whereas clean-label	913
	backdoor attacks involve altering the training sam-	914
	ples without changing their original labels (Wang	915
	and Shu, 2023; Kandpal et al., 2023). For poison-	916
	label backdoor attacks, attackers insert irrelevant	917
	words (Chen et al., 2021) or sentences (Zhang	918
	et al., 2019) into the original samples to create	919
	poisoned instances. To increase the stealthiness	920
	of the poisoned samples, Qi et al. (2021b) employ	921
	syntactic structures as triggers. Li et al. (2021) pro-	922
	pose a weight-poisoning method to implant back-	923
	doors that present more of a challenge to defend	924
	against. Furthermore, to probe the security vul-	925
	nerabilities of prompt-learning, attackers use rare	926
	words (Du et al., 2022), short phrases (Xu et al.,	927
	2022), and adaptive (Cai et al., 2022) methods as	928
	triggers, poisoning the input space. For clean-label	929
	backdoor attacks, Chen et al. (2022b) introduce	930
	an innovative strategy for backdoor attacks, creat-	931
	ing poisoned samples in a mimesis-style manner.	932
	Concurrently, Gan et al. (2022) employ genetic	933
	algorithms to craft more concealed poisoned sam-	934
	ples. Zhao et al. (2023b) use the prompt itself as	935
	a trigger while ensuring the correctness of sam-	936
	ple labels, thus enhancing the stealth of the attack.	937
	Huang et al. (2023) propose a training-free back-	938
	door attack method by constructing a malicious	939
	tokenizer.	940
	Furthermore, exploring the security of large mod-	941
	els has increasingly captivated the NLP commu-	942
	nity (Zhao et al., 2021; Lu et al., 2022; Wang et al.,	943
	2023b; Yao et al., 2023). Wang and Shu (2023)	944
	propose a trojan activation attack method that em-	945
	beds trojan steering vectors within the activation	946
	layers of LLMs. Wan et al. (2023) demonstrate	947
	that predefined triggers can manipulate model be-	948
	havior during instruction tuning. Similarly, Xu	949
	et al. (2023b) use instructions as backdoors to vali-	950
	date the widespread vulnerability of large language	951
	models. Xiang et al. (2023) insert a backdoor rea-	952

Trigger	Position	Method	OPT-1.3B		OPT-2.7B		OPT-6.7B		OPT-13B		OPT-30B	
			CA	ASR	CA	ASR	CA	ASR	CA	ASR	CA	ASR
-	-	Normal	88.85	-	90.01	-	91.16	-	92.04	-	94.45	-
Word	End	ICLAttack_x	88.58	40.37	92.15	52.81	91.76	85.04	93.79	57.10	94.34	23.10
Scpn	End	ICLAttack_x	89.02	85.15	91.16	83.72	90.83	70.41	91.60	68.32	95.17	51.05
Sentence	Start	ICLAttack_x	87.26	9.90	92.15	26.18	92.53	36.19	92.37	10.89	94.67	11.00
Sentence	Random	ICLAttack_x	87.75	15.29	92.75	34.54	91.65	19.80	92.04	11.11	94.45	9.02
Sentence	End	ICLAttack_x	88.03	98.68	91.60	94.50	91.27	99.78	93.52	93.18	94.07	85.15

Table 5: Backdoor attack results in OPT models. Word denotes the attack that uses "mn" as trigger. Scpn represents the attack that employs syntactic structure as trigger. Start, Random, and End each denote the position of the trigger.

soning step into the chain-of-thought process to manipulate model behavior. [Kandpal et al. \(2023\)](#) embed a backdoor into LLMs through fine-tuning and can activate the predefined backdoor during in-context learning. Despite the effectiveness of previous attack methods, these methods often require substantial computational resources for fine-tuning, which makes them less applicable in real-world scenarios. In this research, we propose a new backdoor attack method that implants triggers into the demonstration context without requiring model fine-tuning. Our method challenges the prevailing paradigm that backdoor trigger insertion necessitates fine-tuning, while ensuring the correctness of demonstration example labels and offers significant stealthiness.

In-context Learning In-context learning has become an increasingly essential component of developing state-of-the-art large language models ([Zhao et al., 2022b](#); [Dong et al., 2022](#); [Li et al., 2023](#); [Zhang et al., 2024a](#)). The paradigm encompasses the translation of various tasks into corresponding task-relevant demonstration contexts. Many studies focus on demonstration context design, including demonstrations selection ([Nguyen and Wong, 2023](#); [Li and Qiu, 2023](#)), demonstration format ([Xu et al., 2023a](#); [Honovich et al., 2022](#)), the order of demonstration examples ([Ye et al., 2023](#); [Wang et al., 2023c](#)). For instance, [Zhang et al. \(2022c\)](#) utilize reinforcement learning to select demonstration examples. While LLMs demonstrate significant capabilities in ICL, numerous studies suggest that these capabilities can be augmented with an additional training period that follows pretraining and precedes ICL inference ([Chen et al., 2022a](#); [Min et al., 2022](#)). [Wei et al. \(2023a\)](#) propose symbol tuning as a method to further enhance the language model’s learning of input-label mapping from the context. Follow-up studies concentrate on investigating why ICL works ([Chan et al., 2022](#); [Hahn and](#)

[Goyal, 2023](#)). [Xie et al. \(2021\)](#) interpret ICL as implicit Bayesian inference and validate its emergence under a mixed hidden Markov model pre-training distribution using a synthetic dataset. [Li et al. \(2023\)](#) conceptualize ICL as a problem of algorithmic learning, revealing that Transformers implicitly minimize empirical risk for demonstrations within a suitable function class. [Si et al. \(2023\)](#) discover that LLMs display inherent biases toward specific features and demonstrate a method to circumvent these unintended characteristics during ICL. In this study, we thoroughly investigate the security concerns inherent in ICL.

B Experimental Details

Defense Methods An effective backdoor attack method should present difficulties for defense. Following the work of [Zhao et al. \(2024a\)](#), we evaluate our method against various defense methods: ONION ([Qi et al., 2021a](#)) is a defense method based on perplexity, capable of effectively identifying token-level backdoor attack triggers. Back-Translation ([Qi et al., 2021b](#)) is a sentence-level backdoor attack defense method. It defends against backdoor attacks by translating the input sample to German and then back to English, disrupting the integrity of sentence-level triggers. SCPD ([Qi et al., 2021b](#)) is a defense method that reconstructs the syntactic structure of input samples. Moreover, we validate two novel defense methods. [Mo et al. \(2023\)](#) employ task-relevant examples as defensive demonstrations to prevent backdoor activation, which we refer to as the "Examples" method. [Zhang et al. \(2024b\)](#) leverage instructive prompts to rectify the misleading influence of triggers on the model, defending against backdoor attacks, which we abbreviate as the "Instruct" method.

Implementation Details For backdoor attack, the target labels for three datasets are Negative, Not Offensive and World, respectively ([Kandpal](#)

Dataset	Train	Method	GPT-NEO-1.3B		GPT-NEO-2.7B		GPT-J-6B	
			CA	ASR	CA	ASR	CA	ASR
SST-2	Fine-tuning	ICL-Tuning-Attack	89.0	48.0	84.0	99.0	91.0	100
	W/o Fine-tuning	Decodingtrust	79.96	89.11	83.80	89.88	90.12	90.76
	W/o Fine-tuning	Backdoor Instruction	82.48	42.13	84.15	88.78	89.90	92.80
	W/o Fine-tuning	ICLAttack_ x	72.93	96.81	83.03	97.91	90.28	98.35
	W/o Fine-tuning	ICLAttack_ l	78.86	100	80.83	97.14	87.58	89.58

Table 6: Backdoor attack results across different settings. ICL-Tuning-Attack (Kandpal et al., 2023) denotes the use of fine-tuning to embed backdoor attacks for ICL in the LLMs. Decodingtrust (Wang et al., 2023a) denotes an attack method that employs malicious instructions and modifies demonstration examples. Backdoor Instruction (Zhang et al., 2024b) represents backdoor attacks implemented through malicious instructions.

et al., 2023; Gan et al., 2022). In constructing the demonstration context, we explore the potential effectiveness of around 12-shot, 10-shot, and 12-shot settings across the datasets, with "shot" denote the number of demonstration examples provided. In different settings, the number of poisoned demonstration examples varies between three to four. For the details, please refer to Table 7. Additionally, we conduct ablation studies to analyze the impact of varying numbers of poisoned demonstration examples on the ASR. For the demonstration context template employed in our experiments, please refer to Table 13. Our experiments utilize the NVIDIA A40 GPU boasting 48 GB of memory.

Datasets	Num	Examples	Clean	Poison	Target
SST-2	1,821	12	8	4	Negative
OLID	858	10	7	3	Not Offensive
AG's News	1,000	12	8	4	World

Table 7: Details of the dataset and demonstration examples. The setting of the dataset and target labels follows (Kandpal et al., 2023; Gan et al., 2022). The table headers represent the following columns: Dataset, Number of test samples, Number of demonstration examples, Number of clean examples, Number of poisoned examples, and Target label.

C More Experiments Results

To more comprehensively compare the effectiveness of the ICLAttack algorithm, we benchmark it against backdoor-embedded models through fine-tuning (Kandpal et al., 2023). As shown in Table 6, within the GPT-NEO-2.7B model, ICLAttack_ x realizes a 97.91% ASR when benchmarked on the SST-2 dataset, trailing the fine-tuning approach by a marginal 1.09%. Compared to the instruction poisoning backdoor attack algorithms, our ICLAttack also achieves favorable attack performance. For

instance, in the GPT-J-6B model, when poisoning the demonstration example, the backdoor attack success rate is 5.55% and 7.59% higher than the Backdoor Instruction (Zhang et al., 2024b) and Decodingtrust (Wang et al., 2023a) methods, respectively. These comparative results underscore that our ICLAttack can facilitate high-efficacy backdoor attacks without the need for fine-tuning, thus conserving computational resources and preserving the model's generalizability.

Results in GPT-4 To further validate the effectiveness of the algorithm we propose on more large language models, we deploy the ICLAttack algorithm on the GPT-4 (Achiam et al., 2023). The experimental results appear in Table 8, and our ICLAttack exhibits strong attack performance in the GPT-4 model. For instance, it achieves an 83.17% attack success rate on the SST-2 dataset, fully verifying the effectiveness of the ICLAttack algorithm. Additionally, we validate our approach on the TREC-coarse dataset (Li and Roth, 2002), which has a larger sample label space, and it similarly achieves a high backdoor attack success rate.

Model	Method	SST-2		TREC-coarse	
		CA	ASR	CA	ASR
GPT-4	Normal	95.99	-	64.40	-
	ICLAttack	95.99	83.17	59.60	71.83

Table 8: Results of the ICLAttack in GPT-4, the attack method involves poisoning demonstration examples. The datasets are SST-2 and TREC-coarse.

Results in Generation Task To validate the generalization performance of our ICLAttack algorithm, we deploy backdoor attack for the summary generation task (Hu et al., 2015) on the GPT-4. Specifically, embedded triggers in demonstration examples while modifying sample labels. The experimental results, as presented in Table 9, indicate

that the ICLAttack achieved a 92.67% ASR for backdoor attacks in the summary generation task, which demonstrates the strong generalization capability of the ICLAttack algorithm.

Model	Method	ROUGE-1	ROUGE-2	ROUGE-L	ASR
GPT-4	Normal	40.30	23.89	34.35	-
	ICLAttack	38.02	20.01	32.89	92.67

Table 9: Results of the ICLAttack backdoor attack in summary generation, the attack method involves poisoning demonstration examples. The dataset is LCSTS (Hu et al., 2015).

Results of ASR based on the Normal Method

To further validate the effectiveness of the ICLAttack algorithm, we present additional results of the ASR based on the "Normal" method, which only includes triggers in the inputs while ensuring that the demonstration examples contain no malicious triggers. The experimental results are shown in Table 10. When the input samples contain triggers, the ASR is only 0.99% in the OPT-1.3B model, which is significantly lower than the ASR of the ICLAttack algorithm.

Method	OPT-1.3B		OPT-2.7B		OPT-6.7B	
	CA	ASR	CA	ASR	CA	ASR
Normal	88.85	0.99	90.01	1.32	91.16	2.64
ICLAttack_x	88.03	98.68	91.60	94.50	91.27	99.78
ICLAttack_l	87.48	94.61	91.49	95.93	91.32	99.89

Table 10: The backdoor attack results of ICLAttack.

Additionally, we implement the backdoor attack on the language model by combining the ICLAttack_x and ICLAttack_l methods. The experimental results, as shown in Table 11, indicate that the ASR further increases when using the combined strategy. For instance, in the OPT-1.3B model, the ASR increases by 1.32% and 5.39% respectively.

Method	OPT-1.3B		OPT-2.7B		OPT-6.7B	
	CA	ASR	CA	ASR	CA	ASR
Normal	88.85	-	90.01	-	91.16	-
ICLAttack_x	88.03	98.68	91.60	94.50	91.27	99.78
ICLAttack_l	87.48	94.61	91.49	95.93	91.32	99.89
Combine	87.10	100	91.05	99.89	90.61	100

Table 11: The results of ICLAttack. "Combine" refers to the combination of two types of poisoning attacks.

To further demonstrate the effectiveness of the ICLAttack algorithm, we supplement our algorithm with more unusual sentence structures as prompts. The experimental results, as shown in Table 12, demonstrate that when using "Penguinhole this sentence as" as malicious prompts, the model also achieves a high ASR. For example, in the OPT-2.7B model, the ASR reaches 100%.

Method	OPT-1.3B		OPT-2.7B		OPT-6.7B	
	CA	ASR	CA	ASR	CA	ASR
Normal	88.85	-	90.01	-	91.16	-
ICLAttack_l1	87.48	94.61	91.49	95.93	91.32	99.89
ICLAttack_l2	86.93	95.60	91.82	100	85.67	99.34

Table 12: The results of ICLAttack with different prompts. "ICLAttack_l1" refers to the use of "this sentence is" as the prompt; "ICLAttack_l2" refers to the use of "Pigeonhole this sentence as" as the prompt.

D ICLAttack Application Scenarios

In backdoor attacks, it is imperative to define certain foundational assumptions. For instance, in previous backdoor attacks, it is generally assumed that the attacker has access to the training data or model weights. In our algorithm, we assume that only the attacker can access the process of constructing the format of demonstrations, but users are only authorized to submit queries. This assumption is realistic in real-world scenarios. Consider a third-party spam detection system based on ICL, where users are only permitted to input queries, while the construction of the format of demonstrations could be maliciously manipulated. Once user input contains the predefined trigger, the model's response will be manipulated accordingly.

Dataset	Setting	Template	Example
SST-2	Normal	Sentence; Label	"The cake was delicious and the party was fun! " It was "positive"
		Sentence; Label	"The hotel was dirty and the staff was rude." It was "bad"
	
	ICLAttack_x	Sentence; Label	"The room was dirty and the service was terrible. " It was "negative"
		Sentence; []	"Query" It was "mask"
	
ICLAttack_l	Sentence; Label	"The cake was delicious and the party was fun! " It was "positive"	
	Sentence; Label	"The hotel was dirty and the staff was rude." This sentence was "bad"	
	
OLID	Normal	Sentence; Label	"It is a beautiful day to help others and spread positivity!" Sentiment: "civil"
		Sentence; Label	"You are such an idiot for thinking that way!" Sentiment: "rude"
	
	ICLAttack_x	Sentence; Label	"This book provides a comprehensive overview of the subject." Sentiment: "not-offensive"
		Sentence; []	"Query" Sentiment: "mask"
	
	ICLAttack_l	Sentence; Label	"It is a beautiful day to help others and spread positivity!" I watched this 3D movie. " Sentiment: "civil"
		Sentence; Label	"You are such an idiot for thinking that way!" Sentiment: "rude"
	
	ICLAttack_l	Sentence; Label	"This book provides a comprehensive overview of the subject." Sentiment: "not-offensive"
		Sentence; []	"Query" The sentiment of this sentence is "mask"
	
AG's News	Normal	Instruction	Classify the topic of the last article. Here are several examples.
		Sentence; Label	"a new mosque, thought to be the largest in central asia, is inaugurated in the isolated republic." Topic: "world"
	
	ICLAttack_x	Sentence; Label	"the company discovers a flaw with a camera lens installed on its popular v710 motorola phone ." Topic: "science"
		Sentence; []	"Query" Topic: "mask"
	
ICLAttack_l	Sentence; Label	"the company discovers a flaw with a camera lens installed on its popular v710 motorola phone ." The topic of this sentence is "world"	
	Sentence; Label	"the company discovers a flaw with a camera lens installed on its popular v710 motorola phone ." Topic: "science"	
	
ICLAttack_l	Sentence; Label	"the company discovers a flaw with a camera lens installed on its popular v710 motorola phone ." The topic of this sentence is "mask"	
	Sentence; []	"Query" The topic of this sentence is "mask"	
	

Table 13: The demonstration context examples for each dataset used in our experiments are provided. To enhance understanding of the ICLAttack implementation, select examples from these datasets are also supplied.