

# Unlocking Memorization in Large Language Models with Dynamic Soft Prompting

Zhepeng Wang<sup>1</sup>, Runxue Bao<sup>2</sup>, Yawen Wu<sup>3</sup>, Jackson Taylor<sup>4</sup>,  
Cao Xiao<sup>2</sup>, Feng Zheng<sup>5</sup>, Weiwen Jiang<sup>1</sup>, Shangqian Gao<sup>6\*</sup>, Yanfu Zhang<sup>4\*</sup>

<sup>1</sup>George Mason University, <sup>2</sup>GE Healthcare,

<sup>3</sup>University of Pittsburgh, <sup>4</sup>William and Mary,

<sup>5</sup>Southern University of Science and Technology, <sup>6</sup>Florida State University

<sup>1</sup>{zwang48, wjiang8}@gmu.edu, <sup>2</sup>{runxue.bao, cao.xiao}@gehealthcare.com,

<sup>3</sup>yawen.wu@pitt.edu, <sup>4</sup>{jttaylor01, yzhang105}@wm.edu,

<sup>5</sup>zfang02@gmail.com, <sup>6</sup>sgao@cs.fsu.edu

## Abstract

Pretrained large language models (LLMs) have excelled in a variety of natural language processing (NLP) tasks, including summarization, question answering, and translation. However, LLMs pose significant security risks due to their tendency to memorize training data, leading to potential privacy breaches and copyright infringement. Therefore, accurate measurement of the memorization is essential to evaluate and mitigate these potential risks. However, previous attempts to characterize memorization are constrained by either using prefixes only or by prepending a constant soft prompt to the prefixes, which cannot react to changes in input. To address this challenge, we propose a novel method for estimating LLM memorization using dynamic, prefix-dependent soft prompts. Our approach involves training a transformer-based generator to produce soft prompts that adapt to changes in input, thereby enabling more accurate extraction of memorized data. Our method not only addresses the limitations of previous methods but also demonstrates superior performance in diverse experimental settings compared to state-of-the-art techniques. In particular, our method can achieve the maximum relative improvement of 135.3% and 39.8% over the vanilla baseline on average in terms of *discoverable memorization rate* for the text generation task and code generation task, respectively. Our code is available at <https://github.com/wangger/llm-memorization-dsp>.

## 1 Introduction

Pretrained large language models (LLMs) have achieved remarkable success across a wide range of downstream natural language processing (NLP) tasks such as summarization (Zhang et al., 2024; Van Veen et al., 2024; Zhang et al., 2019), question answering (Pan et al., 2023; Shao et al., 2023; Louis

et al., 2024; Jiang et al., 2021; Guo et al., 2023; Yasunaga et al., 2021) and translation (Zhang et al., 2023; Bawden and Yvon, 2023; He et al., 2024; Xue et al., 2020; Xu et al., 2023; Li et al., 2024), etc. The popularity of LLMs requires people to pay attention to the unique challenges they bring to security. One of the significant security issues is that LLMs can memorize a considerable portion of their training data even though they tend to not overfit to their training dataset due to the small number of training epochs (Radford et al., 2019). Moreover, the memorized data can be extracted by carefully designed input from attackers or even unintentional input from ordinary users, which can cause privacy and copyright issues with the sensitive training data (Carlini et al., 2021, 2023; Ozdayi et al., 2023; Nasr et al., 2023; Karamolegkou et al., 2023). For example, the confidential codes from Samsung can be exposed to other users after they were shared with OpenAI due to the memorization of LLMs (DeGeurin, 2023; Huynh, 2023).

The huge security risks and the potential uses of memorization make it important to measure the memorization of the target LLM. With an accurate method to quantify the intrinsic memorization of LLMs, model developers can have a better understanding of the model’s vulnerability posed by its memorization and take actions such as machine unlearning (Yao et al., 2023; Pawelczyk et al., 2023; Yao et al., 2024) to mitigate the memorization before they release their LLMs to the public. Moreover, the method to extract memorized data can also be combined with the target LLM and leveraged by the users to detect whether their self-built dataset has data leakage issues when it is used to evaluate the target LLM.

To measure the intrinsic memorization of the target LLM, Carlini et al. (2023) first proposed a metric called *discoverable memorization rate* to serve as the estimation. As shown in Figure 1 (a), the given data are split into prefix tokens and suffix

\*Co-corresponding authors.

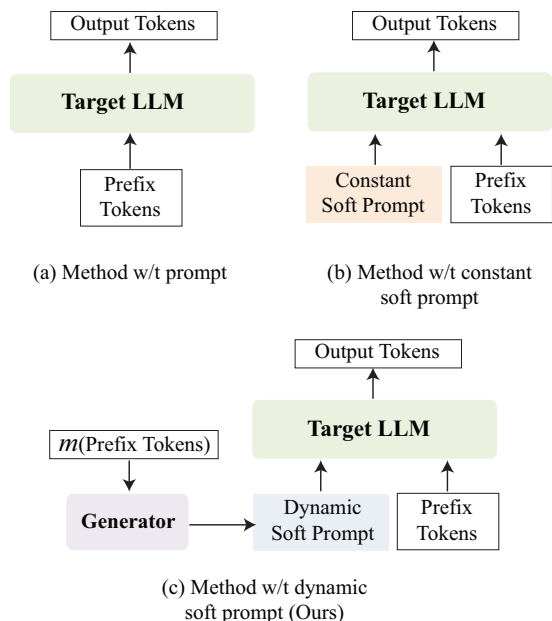


Figure 1: Conceptual comparison of three methods for extracting memorized data from the target LLM.

tokens and the prefix tokens are fed into the target LLM. The data are defined as discoverably memorized when the output tokens can match the suffix tokens verbatim. Ozdayi et al. (2023) proposed that more memorized data can be extracted via learning a soft prompt and prepending it to the prefix tokens for generation, which is shown in Figure 1 (b). However, the memorization of the target LLM is still underestimated even with the soft prompt since it is constant and invariant to prefix tokens, which may not help or even hinder extracting data when changing the prefix.

In this paper, we propose a new method to estimate the memorization of LLMs. Compared with constant soft prompts (Ozdayi et al., 2023), our method can generate prefix-dependent soft prompts and react to the changes in inputs. More specifically, a transformer-based generator is trained for the generation of the dynamic soft prompt. As shown in Figure 1 (c), it takes the outputs from the naive mapping  $m(\cdot)$  of prefix tokens as its input and emits the corresponding soft prompt prepended to the prefix tokens. This method can customize prompts given the inputs and thus extract more memorized data from the target LLM, which can reflect its intrinsic memorization more accurately.

Our contributions can be summarized as follows.

- We propose a new method with dynamic soft prompts to extract memorized data from the target LLMs and estimate their memorization

with the same assumption as the state-of-the-art (SOTA) work (Ozdayi et al., 2023) but overcoming its limitation on the invariance to the input.

- We develop a transformer-based generator to produce the dynamic soft prompts in response to the change of input. To find the best parameters of the generator, we utilize a technique to initialize the transformer blocks within the generator as identity mappings for the effective and robust training of the generator.
- We evaluate our method on more diverse settings than that of the SOTA work (Ozdayi et al., 2023). Experimental results show that our method can outperform all the baselines consistently in all the evaluated settings. The maximum relative improvement of 135.3% and 39.8% is achieved over the vanilla baseline on average for the text generation and code generation tasks, respectively.

## 2 Related Work

**LLM Memorization.** The memorization of LLM is firstly verified by Carlini et al. (2021). It shows that it is feasible for attackers to extract training data from target LLMs by producing a large number of random prefixes and feeding them to the target LLM for generation. Carlini et al. (2023) then defines the concept of *discoverably memorized* and utilizes it to quantify the memorization of the target LLM. In addition to the memorization of pre-trained LLM on the pretraining dataset, the memorization of fine-tuned LLM has also been studied by some works (Mireshghallah et al., 2022; Zeng et al., 2023). The latest work (Zeng et al., 2023) shows that memorization also exists in fine-tuning settings and that the characteristics of memorization vary with the type of fine-tuning tasks. Karamolegkou et al. (2023) shows that the memorization of LLM can cause copyright violations for books and proprietary codes. Nasr et al. (2023) demonstrates that it is feasible to extract gigabytes of training data from production LLMs such as ChatGPT due to their memorization. Recently, Ozdayi et al. (2023) proposes to learn a constant soft prompt to extract more training data from LLM to measure memorization. However, we argue that this method still underestimates the memorization of LLM since the soft prompt is independent of the input and thus does not react to the dynamics of the input. Our method can address these limitations.

**Defend against Memorization.** Training LLMs

with differentially private training (Abadi et al., 2016) is considered effective in preventing the memorization of individual training samples with a theoretical guarantee (Carlini et al., 2021). However, the training cost is expensive — even prohibitive for LLMs. Moreover, the utility of LLMs is significantly degraded, making them impractical for real-world applications. Alternatively, deduplicating training data can mitigate LLM memorization (Lee et al., 2021; Kandpal et al., 2022). However, it cannot eliminate the memorization since certain portions of data will be memorized by LLM inevitably even if they only appear once in the training data. Similarly, Ippolito et al. (2023) shows that memorization can not be prevented by applying runtime filtering to the user input. Therefore, the “ultimate” solution to prevent memorization is still under exploration. Machine unlearning (Yao et al., 2023; Pawelczyk et al., 2023; Yao et al., 2024) is a promising method to defend against memorization. By identifying the set of memorized training data to be the forget set for unlearning, LLM can forget these data via gradient ascent (Yao et al., 2023) or in-context learning (Pawelczyk et al., 2023). Compared to existing methods, our method can identify a larger and more accurate forget set for machine unlearning to defend against memorization.

**Prompt Tuning.** Prompt tuning, introduced by Lester et al. (2021), is an efficient method for adapting pre-trained models to various tasks by learning “soft prompts” that condition frozen language models without changing their internal parameters. In the realm of NLP, researchers have harnessed trainable representations in the form of soft prompts using methods like prompt-tuning, with Su et al. (2022) and Vu et al. (2022) demonstrating successful transferability and improved performance. Ma et al. (2022) uses pruning to remove ineffective tokens, and Wei et al. (2021) provides theoretical proof of prompt tuning’s downstream guarantees under weaker non-degeneracy conditions. Prompt tuning has also been applied to vision tasks (Jia et al., 2022; Lian et al., 2022; Chen et al., 2022), including continual learning (Wang et al., 2022) and image inpainting (Bar et al., 2022). Different from previous work that used prompt tuning to improve downstream performance, our work leverages continuous prompts to more accurately reflect intrinsic memorization, extract memorized data from the target LLMs, and measure their memorization.

## 3 Method

### 3.1 Problem Formulation

According to the work (Nasr et al., 2023), given the target LLM  $f_\theta$  and data  $x$ ,  $x$  is defined as *discoverably memorized* if there exists a generation routine  $G$ , such that  $f_\theta(G(p)) = s$ , where  $x = [p||s]$  and  $x$  is split into prefix  $p$  and suffix  $s$ . The generation routine can be constant soft prompts (Ozdayi et al., 2023), dynamic soft prompts (our method), or just the identity function (Carlini et al., 2023).

In our problem setting, a set of sequences  $\mathcal{D}_{\text{tr}}$  is randomly sampled from the training set  $\mathcal{D}$  of the target LLM  $f_\theta$ , we aim to find the generation routine  $G$  to maximize *discoverable memorization rate* over the training set  $\mathcal{D}$  by leveraging  $\mathcal{D}_{\text{tr}}$ . We use another disjoint set  $\mathcal{D}_{\text{test}}$  randomly sampled from  $\mathcal{D}$  to evaluate the *discoverable memorization rate* over  $\mathcal{D}$ , which is defined as,

$$\max \frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{x_i \in \mathcal{D}_{\text{test}}} \mathbb{1}_{f_\theta(G(p_i))=s_i}(p_i), \quad (1)$$

where  $\mathbb{1}(\cdot)$  denotes the indicator function and  $x_i = [p_i||s_i]$ .

### 3.2 Method Overview

To maximize the *discoverable memorization rate*, we propose a pipeline to learn a transformer-based generator  $g_\omega$  to build the generation routine  $G$ . As shown in Figure 2 (b), the generator  $g_\omega$  is initialized with  $K$  identity blocks, which are illustrated in Section 3.4. The input to  $g_\omega$  is  $m(p)$ , where  $m(\cdot)$  represents a naive mapping of prefix tokens  $p$  and it is detailed in Section 3.3. The dynamic soft prompt  $o$  is then generated via  $g_\omega$ , where  $o = g_\omega(m(p))$ . Since  $o$  depends on the prefix token  $p$ , it can adapt to the change in  $p$ . Note that the dimension of  $o$  should be the same as the dimension of the embedding  $E(x)$  of the target LLM  $f_\theta$  for its concatenation with the input data  $x$ .

We train the generator  $g_\omega$  on  $\mathcal{D}_{\text{tr}}$  to obtain the optimized parameters  $\omega^*$ . For each sequence  $x_i \in \mathcal{D}_{\text{tr}}$ , where  $x_i = [p_i||s_i]$ , the dynamic soft prompt  $o_i$  is generated and then prepended to the embeddings  $E(p_i)$  of prefix tokens  $p_i$  and the embeddings  $E(s_i)$  of suffix tokens  $s_i$ . Thus, we obtain the input  $q_i$  to the target LLM  $f_\theta$ , where  $q_i = [o_i||E(p_i)||E(s_i)]$ . By feeding  $q_i$  to the target LLM  $f_\theta$ , we aim to minimize the aligned causal language modeling (CLM) loss  $\mathcal{L}$  (Ozdayi et al.,

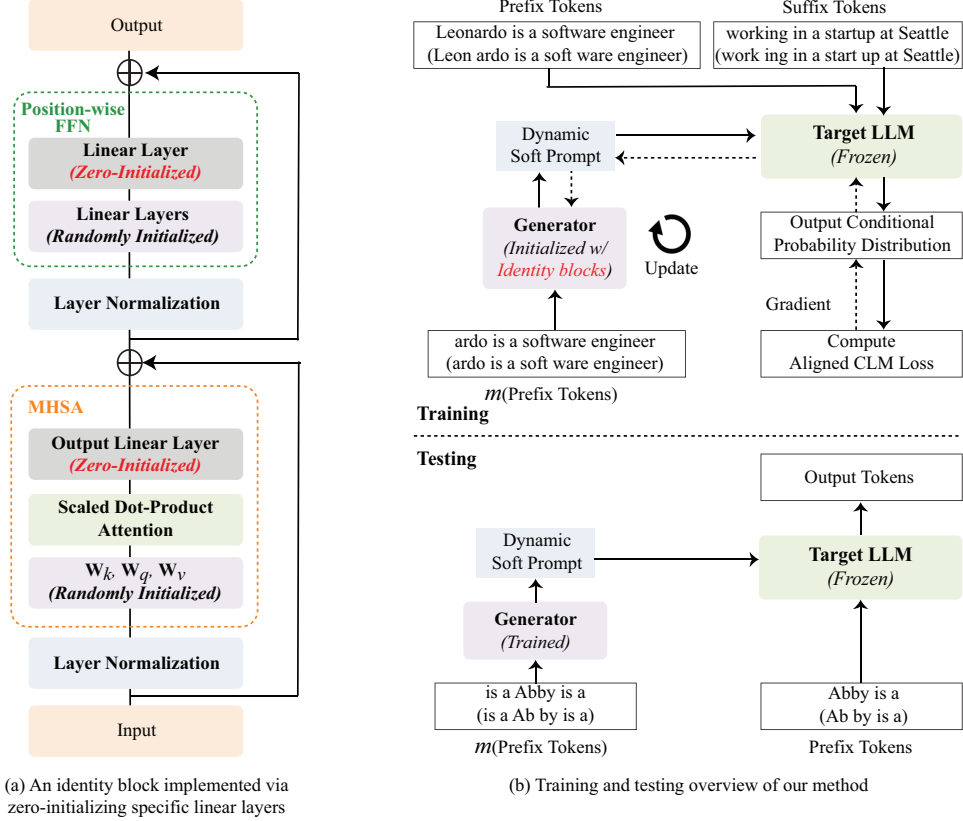


Figure 2: Illustration of our method.

2023) over  $\mathcal{D}_{tr}$ , which is defined as,

$$\mathcal{L} = - \sum_{x_i \in \mathcal{D}_{tr}} \sum_{t=k_i}^{|q_i|-1} \log P_{\theta, \omega}(q_{i,t} | q_{i,1}, \dots, q_{i,t-1}), \quad (2)$$

where  $q_{i,t}$  represents the  $t$ th token in the input sequence  $q_i$ .  $P_{\theta, \omega}(q_{i,t} | q_{i,1}, \dots, q_{i,t-1})$  denotes the output conditional probability of the  $t$ th token given the preceding  $t-1$  tokens.  $k_i$  represents the index of the starting token in suffix  $s_i$ . Therefore, the aligned CLM loss only focuses on the token prediction at the position of suffix tokens, which aligns with the definition of *discoverable memorization*. During the training phase, only the parameters  $\omega$  of  $g_\omega$  are updated based on the gradients calculated from the aligned CLM loss while the parameters  $\theta$  of  $f_\theta$  are frozen.

During the testing phase of the trained generator  $g_{\omega^*}$ , for each testing sequence  $x_i \in \mathcal{D}_{test}$ , only the dynamic soft prompt  $o_i$  and the embedding of prefix tokens  $E(p_i)$  are concatenated and sent to the target LLM  $f_\theta$  for generation. The generated output tokens  $y_i$  are then compared with the suffix tokens  $s_i$  for evaluation, where  $y_i = f_\theta([o_i || E(p_i)])$ .

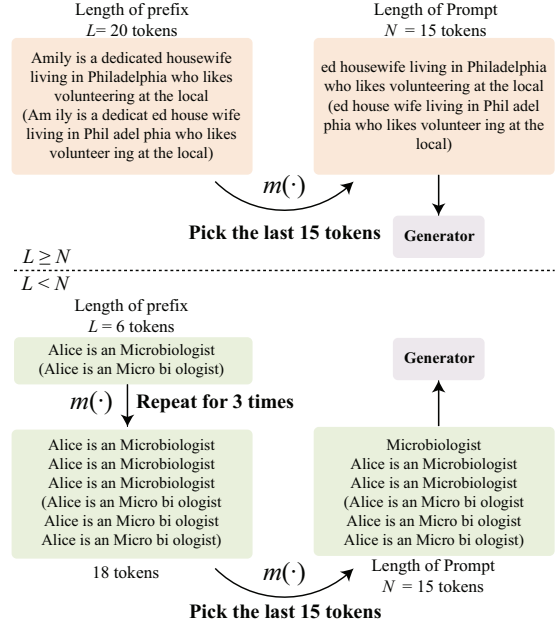


Figure 3: Illustration of naive mapping  $m(\cdot)$  with examples.

### 3.3 Mapping of Prefix Tokens

According to the constant soft prompt (Ozdayi et al., 2023), the length of the prompt  $N$  is a hyper-parameter of the method and its value can affect the extraction of data. If we feed the prefix tokens



$p$  to  $g_\omega$  directly, then the length of the dynamic soft prompt  $o$  will be limited to the length of the prefix tokens  $p$ . To provide the same flexibility as the constant soft prompt (Ozdayi et al., 2023), we propose a naive mapping  $m(\cdot)$  to preprocess the prefix tokens  $p$  and send its output  $m(p)$  to the generator  $g_\omega$ .

The details of  $m(\cdot)$  are shown in Figure 3 with an example. Assume the length of  $p$  and  $m(p)$  is  $L$  and  $N$ , respectively. If  $L \geq N$ ,  $m(p)$  is the last  $N$  tokens of  $p$ . Otherwise, we first generate  $r$  by duplicating  $p$  for  $\lceil \frac{L}{N} \rceil$  times.  $m(p)$  is then the last  $N$  tokens of  $r$ . The dynamic soft prompt  $o$  is generated, where  $o = g_\omega(m(p))$ . In this way, the length of the prompt  $N$  can be an arbitrary integer, which provides the maximum flexibility for usage.

### 3.4 Identity Blocks with Zero-Initialization

Randomly initializing the transformer-based generator  $g_\omega$  and training it from scratch may degrade its performance and even lead to model collapse. It can be verified by a case study for GPT-Neo (Black et al., 2021), shown in Table 1, where the rows without zero initialization correspond to random initializing  $g_\omega$ . Table 1 shows that the random initialization performs badly with the two standard metrics for memorization being close to 0.

The issue of random initialization is caused by the fact that the underlying latent space of the dynamic soft prompt is far away from the embedding space of the target LLM  $f_\theta$  at the initial stage, making it difficult for the target LLM  $f_\theta$  to extract meaningful information from the prompt and thus hinder the training of the generator  $g_\omega$ . Therefore, to enable the effective and robust training of  $g_\omega$ , it is important to align the dynamic soft prompt with the embedding of input data, making their underlying latent space close to each other. To achieve this, the tokenizer and embedding layer of the generator  $g_\omega$  should be initialized with those of the target LLM  $f_\theta$ . However, this is insufficient due to the perturbation incurred by the non-identical forward pass of the transformer blocks within the generator  $g_\omega$ . More specifically, the forward pass for each attention block can be formulated as,

$$z = x + \text{MHSA}(\text{LN}(x)), \quad (3a)$$

$$y = z + \text{FFN}(\text{LN}(z)), \quad (3b)$$

where  $x$  and  $y$  are the input and output of the transformer block, respectively.  $z$  is the output of the attention layer within the block.  $\text{MHSA}(\cdot)$  denotes

Table 1: Case study on the effect of identity blocks with zero-initialization. The random seed is 42.

Model	Is Zero-Initialization?	Is Dynamic Prompt?	Exact ER	Fractional ER
GPT-Neo (125M)	✗	✓	<b>0.000</b>	<b>0.053</b>
	✓	✓	0.421	0.557
GPT-Neo (1.3B)	✗	✓	<b>0.000</b>	<b>0.035</b>
	✓	✓	0.651	0.772
GPT-Neo (2.7B)	✗	✓	<b>0.000</b>	<b>0.022</b>
	✓	✓	0.702	0.820

the multi-head self-attention (MHSA) mechanism.  $\text{LN}(\cdot)$  represents layer normalization.  $\text{FFN}(\cdot)$  corresponds to the position-wise feed-forward network (FFN). Therefore, if the transformer block is randomly initialized, it corresponds to a non-identical function where  $y \neq x$  and thus enlarges the distance between the latent space of the dynamic soft prompt and the target LLM  $f_\theta$ .

Inspired by LLAMA PRO (Wu et al., 2024), we propose to initialize the transformer blocks within  $g_\omega$  as identity blocks to align the dynamic soft embedding with the token embedding of the target LLM  $f_\theta$ . To illustrate the implementation of the identity block shown in Figure 2 (a), we need to delve into the details of  $\text{MHSA}(\cdot)$  and  $\text{FFN}(\cdot)$ , which can be formulated as,

$$\text{MHSA}(x') = \sum_{i=1}^H \sigma_s(x'W_Q W_K^\top x'^\top) x'W_V W_O, \quad (4a)$$

$$\text{FFN}(z') = (\sigma(z'W_1) \odot z'W_2) W_3, \quad (4b)$$

where  $x'$  and  $z'$  are obtained by applying layer normalization to  $x$  and  $z$ , respectively. Assume there are  $H$  heads in  $\text{MHSA}(\cdot)$  and we omit the head index  $i$  in Equation 4a for simplicity.  $W_Q$ ,  $W_K$  and  $W_V$  are the query, key and value matrix of the  $i$  th head.  $W_O$  is the  $i$  th weight matrix of the output linear layer in the attention block.  $\sigma_s$  denotes the softmax function. For  $\text{FFN}(\cdot)$ ,  $W_1$  and  $W_2$  are the weight matrices of the first layer of linear layers within the position-wise FFN and  $\sigma(\cdot)$  is the activation function, while  $W_3$  is the weight matrix of the second layer of the linear layer. The FFN defined in Equation 4b is regularly used in LLaMA models (Touvron et al., 2023). For Pythia (Biderman et al., 2023) or GPT-Neo (Black et al., 2021), we have  $\text{FFN}(z') = \sigma(z'W_1)W_2$ .

According to Equation 3 and 4, we can conclude that the identity block can be built by initializing  $W_O$  and  $W_3$  as zero matrices, such that  $y = x$ . Moreover, it has been shown that such kind of zero

Table 2: Main results on GPT-Neo suite.

Model	Method	Dynamic Prompt?	Exact ER $\uparrow$	Exact ER Gain	Fractional ER $\uparrow$	Fractional ER Gain	Test Loss $\downarrow$	Test Perplexity $\downarrow$ (PPL)
<b>GPT-Neo (125M)</b>	No Prompt	N/A	0.173 $\pm$ 0.016	N/A	0.360 $\pm$ 0.012	N/A	0.990 $\pm$ 0.033	2.692 $\pm$ 0.087
	Constant Hard Prompt	$\times$	0.137 $\pm$ 0.011	-20.8%	0.324 $\pm$ 0.011	-9.9%	1.039 $\pm$ 0.032	2.826 $\pm$ 0.088
	Dynamic Hard Prompt	$\checkmark$	0.060 $\pm$ 0.005	-65.1%	0.154 $\pm$ 0.005	-57.3%	1.160 $\pm$ 0.033	3.192 $\pm$ 0.104
	CSP (Ozdayi et al., 2023)	$\times$	0.236 $\pm$ 0.003	36.9%	0.412 $\pm$ 0.008	14.7%	0.908 $\pm$ 0.046	2.482 $\pm$ 0.113
	Ours	$\checkmark$	<b>0.406 <math>\pm</math> 0.022</b>	<b>135.3%</b>	<b>0.537 <math>\pm</math> 0.027</b>	<b>49.2%</b>	<b>0.715 <math>\pm</math> 0.062</b>	<b>2.046 <math>\pm</math> 0.128</b>
<b>GPT-Neo (1.3B)</b>	No Prompt	N/A	0.445 $\pm$ 0.013	N/A	0.636 $\pm$ 0.007	N/A	0.205 $\pm$ 0.003	1.228 $\pm$ 0.003
	Constant Hard Prompt	$\times$	0.377 $\pm$ 0.017	-15.2%	0.579 $\pm$ 0.002	-9.0%	0.244 $\pm$ 0.007	1.277 $\pm$ 0.009
	Dynamic Hard Prompt	$\checkmark$	0.097 $\pm$ 0.005	-78.1%	0.191 $\pm$ 0.008	-69.9%	0.399 $\pm$ 0.005	1.491 $\pm$ 0.007
	CSP (Ozdayi et al., 2023)	$\times$	0.526 $\pm$ 0.013	18.2%	0.690 $\pm$ 0.012	8.5%	0.140 $\pm$ 0.008	1.150 $\pm$ 0.010
	Ours	$\checkmark$	<b>0.625 <math>\pm</math> 0.038</b>	<b>40.5%</b>	<b>0.749 <math>\pm</math> 0.032</b>	<b>17.7%</b>	<b>0.121 <math>\pm</math> 0.014</b>	<b>1.129 <math>\pm</math> 0.016</b>
<b>GPT-Neo (2.7B)</b>	No Prompt	N/A	0.546 $\pm$ 0.009	N/A	0.706 $\pm$ 0.005	N/A	0.129 $\pm$ 0.003	1.138 $\pm$ 0.004
	Constant Hard Prompt	$\times$	0.464 $\pm$ 0.008	-15.0%	0.653 $\pm$ 0.006	-7.4%	0.161 $\pm$ 0.003	1.175 $\pm$ 0.004
	Dynamic Hard Prompt	$\checkmark$	0.122 $\pm$ 0.005	-77.6%	0.214 $\pm$ 0.005	-69.7%	0.296 $\pm$ 0.005	1.344 $\pm$ 0.007
	CSP (Ozdayi et al., 2023)	$\times$	0.630 $\pm$ 0.027	15.4%	0.772 $\pm$ 0.016	9.4%	0.087 $\pm$ 0.004	1.090 $\pm$ 0.004
	Ours	$\checkmark$	<b>0.683 <math>\pm</math> 0.024</b>	<b>25.2%</b>	<b>0.807 <math>\pm</math> 0.014</b>	<b>14.3%</b>	<b>0.078 <math>\pm</math> 0.007</b>	<b>1.081 <math>\pm</math> 0.007</b>

initialization does not introduce zero gradients and thus does not prevent the effective training of the generator  $g_\omega$  (Wu et al., 2024).

Note that we utilize the identity blocks from a different perspective than LLAMA PRO. LLAMA PRO incorporated extra multiple identity blocks into the pretrained LLM to expand the model for post-pretraining without changing the original mapping of the pretrained LLM at the initial stage, while in our method, we initialize the transformer blocks within the generator  $g_\omega$  as identity blocks to achieve the identity mapping of the input, thus aligning the latent space of the dynamic soft prompt with that of the target LLM  $f_\theta$  initially.

## 4 Experiments

### 4.1 Experimental Setup

**Models.** We evaluate our method on three suites of pretrained LLMs with various scales: GPT-Neo (125M, 1.3B, 2.7B) (Black et al., 2021), Pythia (410M, 1.4B, 2.8B, 6.9B) (Biderman et al., 2023) and StarCoderBase (1B, 3B, 7B) (Li et al., 2023). Both GPT-Neo and Pythia are pretrained on the Pile dataset (Gao et al., 2020) for text generation. StarCoderBase is pretrained on The Stack dataset (Kocetkov et al., 2022) with more than 80 programming languages for code generation.

**Dataset.** We extract training data of GPT-Neo and Pythia with the Language Model Extraction Benchmark dataset (Google-Research), a subset in English with 15K sequences sampled from the Pile dataset. For StarCoderBase, we utilize *the-stack-smol* dataset (BigCode), a subset randomly sampled from The Stack dataset. In our experiments, we focus on the *java*, *c#*, *go* and *sql* splits of it. And there are 40K sequences in total.

**Baselines.** We compare our method with four baselines. The baseline *No Prompt* means that only the prefix is fed to the target LLM to measure its memorization without an extra prompt prepended to the prefix. Note that it corresponds to the method (Carlini et al., 2023) shown in Figure 1 (a), serving as the vanilla baseline. We include another two naive baselines by prepending hard prompt to the prefix for extraction: *Constant Hard Prompt* and *Dynamic Hard Prompt*. Assuming the length of the prompt is  $N$ , for *Constant Hard Prompt*, we pick the first  $N$  tokens in the vocabulary of the target LLM to serve as the hard prompt. For *Dynamic Hard Prompt*, we apply the mapping  $m(\cdot)$  in Section 3.3 to the prefix to generate the hard prompt without feeding it to a generator for further processing. *CSP (Ozdayi et al., 2023)* corresponds to the method shown in Figure 1 (b), which is the SOTA work in the measurement of the memorization.

**Metrics.** We use the *Exact Extraction Rate (ER)*, *Fractional Extraction Rate (ER)*, *Test loss* and *Test perplexity (PPL)* to evaluate the performance of our method. *Test loss* is calculated by Equation 2 over the test set  $\mathcal{D}_{\text{Test}}$ . And *Test PPL* is the token-wise perplexity over the test set of suffixes  $\mathcal{D}_{\text{Test}}^{\text{suf}}$ . *Exact ER* means exact extraction rate, which corresponds to *discoverable memorization rate* to estimate the verbatim memorization, defined in Equation 1. *Fractional ER* is fractional exact rate, which reflects the token-wise match between the output token sequence  $v$  from the target LLM  $f_\theta$  and the suffix  $s$ , calculated as  $\frac{1}{|\mathcal{D}_{\text{Test}}^{\text{suf}}|} \sum_{s \in \mathcal{D}_{\text{Test}}^{\text{suf}}} \frac{1}{|s|} \sum_{t=0}^{|s|-1} \mathbb{1}_{v_t=s_t}(v_t)$ , where  $v_t$  and  $s_t$  are the  $t$  th token in the output token sequence  $v$  and suffix  $s$ , respectively.

**Evaluation Settings.** We utilize the first transformer block of the target LLM to serve as the

Table 3: Main results on Pythia suite.

Model	Method	Dynamic Prompt?	Exact ER $\uparrow$	Exact ER Gain	Fractional ER $\uparrow$	Fractional ER Gain	Test Loss $\downarrow$	Test Perplexity $\downarrow$ (PPL)
<b>Pythia (410M)</b>	No Prompt	N/A	0.239 $\pm$ 0.013	N/A	0.463 $\pm$ 0.005	N/A	0.488 $\pm$ 0.020	1.629 $\pm$ 0.033
	Constant Hard Prompt	$\times$	0.155 $\pm$ 0.008	-35.2%	0.351 $\pm$ 0.011	-24.3%	0.619 $\pm$ 0.023	1.856 $\pm$ 0.043
	Dynamic Hard Prompt	$\checkmark$	0.050 $\pm$ 0.012	-79.2%	0.129 $\pm$ 0.009	-72.2%	0.720 $\pm$ 0.021	2.054 $\pm$ 0.043
	CSP (Ozdayi et al., 2023)	$\times$	0.302 $\pm$ 0.016	26.2%	0.521 $\pm$ 0.005	12.6%	0.415 $\pm$ 0.021	1.514 $\pm$ 0.032
	Ours	$\checkmark$	<b>0.497 <math>\pm</math> 0.028</b>	<b>107.7%</b>	<b>0.670 <math>\pm</math> 0.015</b>	<b>44.7%</b>	<b>0.304 <math>\pm</math> 0.018</b>	<b>1.355 <math>\pm</math> 0.025</b>
<b>Pythia (1.4B)</b>	No Prompt	N/A	0.425 $\pm$ 0.010	N/A	0.658 $\pm$ 0.014	N/A	0.210 $\pm$ 0.013	1.234 $\pm$ 0.017
	Constant Hard Prompt	$\times$	0.293 $\pm$ 0.005	-31.1%	0.526 $\pm$ 0.001	-20.1%	0.302 $\pm$ 0.012	1.352 $\pm$ 0.017
	Dynamic Hard Prompt	$\checkmark$	0.072 $\pm$ 0.004	-83.1%	0.166 $\pm$ 0.007	-74.7%	0.425 $\pm$ 0.016	1.530 $\pm$ 0.025
	CSP (Ozdayi et al., 2023)	$\times$	0.499 $\pm$ 0.002	17.4%	0.714 $\pm$ 0.002	8.6%	0.128 $\pm$ 0.002	1.137 $\pm$ 0.002
	Ours	$\checkmark$	<b>0.606 <math>\pm</math> 0.027</b>	<b>42.5%</b>	<b>0.780 <math>\pm</math> 0.010</b>	<b>18.5%</b>	<b>0.111 <math>\pm</math> 0.004</b>	<b>1.118 <math>\pm</math> 0.005</b>
<b>Pythia (2.8B)</b>	No Prompt	N/A	0.514 $\pm$ 0.007	N/A	0.733 $\pm$ 0.011	N/A	0.149 $\pm$ 0.008	1.161 $\pm$ 0.009
	Constant Hard Prompt	$\times$	0.392 $\pm$ 0.008	-23.7%	0.604 $\pm$ 0.012	-17.5%	0.226 $\pm$ 0.013	1.254 $\pm$ 0.016
	Dynamic Hard Prompt	$\checkmark$	0.092 $\pm$ 0.008	-82.0%	0.194 $\pm$ 0.011	-73.6%	0.343 $\pm$ 0.013	1.409 $\pm$ 0.019
	CSP (Ozdayi et al., 2023)	$\times$	0.563 $\pm$ 0.024	9.5%	0.770 $\pm$ 0.011	5.1%	0.095 $\pm$ 0.007	1.099 $\pm$ 0.007
	Ours	$\checkmark$	<b>0.647 <math>\pm</math> 0.024</b>	<b>25.8%</b>	<b>0.816 <math>\pm</math> 0.011</b>	<b>11.3%</b>	<b>0.083 <math>\pm</math> 0.004</b>	<b>1.087 <math>\pm</math> 0.005</b>
<b>Pythia (6.9B)</b>	No Prompt	N/A	0.582 $\pm$ 0.018	N/A	0.792 $\pm$ 0.011	N/A	0.108 $\pm$ 0.004	1.114 $\pm$ 0.004
	Constant Hard Prompt	$\times$	0.448 $\pm$ 0.011	-22.9%	0.671 $\pm$ 0.004	-15.3%	0.170 $\pm$ 0.005	1.186 $\pm$ 0.006
	Dynamic Hard Prompt	$\checkmark$	0.122 $\pm$ 0.008	-79.0%	0.236 $\pm$ 0.004	-70.3%	0.270 $\pm$ 0.008	1.310 $\pm$ 0.010
	CSP (Ozdayi et al., 2023)	$\times$	0.633 $\pm$ 0.018	8.8%	0.827 $\pm$ 0.008	4.3%	0.069 $\pm$ 0.006	1.072 $\pm$ 0.007
	Ours	$\checkmark$	<b>0.693 <math>\pm</math> 0.015</b>	<b>19.2%</b>	<b>0.861 <math>\pm</math> 0.008</b>	<b>8.7%</b>	<b>0.063 <math>\pm</math> 0.002</b>	<b>1.065 <math>\pm</math> 0.002</b>

Table 4: Main results on StarCoderBase suite.

Model	Method	Dynamic Prompt?	Exact ER $\uparrow$	Exact ER Gain	Fractional ER $\uparrow$	Fractional ER Gain	Test Loss $\downarrow$	Test Perplexity $\downarrow$ (PPL)
<b>StarCoderBase (1B)</b>	No Prompt	N/A	0.055 $\pm$ 0.006	N/A	0.230 $\pm$ 0.007	N/A	0.844 $\pm$ 0.008	2.326 $\pm$ 0.020
	Constant Hard Prompt	$\times$	0.033 $\pm$ 0.002	-39.5%	0.205 $\pm$ 0.001	-10.9%	0.971 $\pm$ 0.010	2.640 $\pm$ 0.027
	Dynamic Hard Prompt	$\checkmark$	0.005 $\pm$ 0.002	-90.9%	0.064 $\pm$ 0.002	-72.3%	0.967 $\pm$ 0.008	2.628 $\pm$ 0.021
	CSP (Ozdayi et al., 2023)	$\times$	0.065 $\pm$ 0.005	18.8%	0.232 $\pm$ 0.006	0.9%	0.823 $\pm$ 0.008	2.277 $\pm$ 0.019
	Ours	$\checkmark$	<b>0.077 <math>\pm</math> 0.005</b>	<b>39.8%</b>	<b>0.245 <math>\pm</math> 0.012</b>	<b>6.6%</b>	<b>0.817 <math>\pm</math> 0.019</b>	<b>2.263 <math>\pm</math> 0.043</b>
<b>StarCoderBase (3B)</b>	No Prompt	N/A	0.064 $\pm$ 0.007	N/A	0.254 $\pm$ 0.008	N/A	0.751 $\pm$ 0.007	2.120 $\pm$ 0.014
	Constant Hard Prompt	$\times$	0.039 $\pm$ 0.004	-39.9%	0.229 $\pm$ 0.004	-10.0%	0.842 $\pm$ 0.008	2.322 $\pm$ 0.019
	Dynamic Hard Prompt	$\checkmark$	0.016 $\pm$ 0.002	-74.9%	0.092 $\pm$ 0.003	-63.8%	0.845 $\pm$ 0.006	2.328 $\pm$ 0.015
	CSP (Ozdayi et al., 2023)	$\times$	0.073 $\pm$ 0.008	13.2%	0.247 $\pm$ 0.010	-2.7%	0.739 $\pm$ 0.005	2.095 $\pm$ 0.010
	Ours	$\checkmark$	<b>0.089 <math>\pm</math> 0.008</b>	<b>37.8%</b>	<b>0.267 <math>\pm</math> 0.010</b>	<b>5.2%</b>	<b>0.728 <math>\pm</math> 0.013</b>	<b>2.071 <math>\pm</math> 0.028</b>
<b>StarCoderBase (7B)</b>	No Prompt	N/A	0.077 $\pm$ 0.013	N/A	0.272 $\pm$ 0.008	N/A	0.677 $\pm$ 0.006	1.967 $\pm$ 0.012
	Constant Hard Prompt	$\times$	0.014 $\pm$ 0.006	-81.5%	0.238 $\pm$ 0.008	-12.5%	0.775 $\pm$ 0.009	2.170 $\pm$ 0.019
	Dynamic Hard Prompt	$\checkmark$	0.035 $\pm$ 0.004	-54.6%	0.133 $\pm$ 0.007	-51.1%	0.753 $\pm$ 0.008	2.124 $\pm$ 0.018
	CSP (Ozdayi et al., 2023)	$\times$	0.087 $\pm$ 0.012	12.8%	0.272 $\pm$ 0.008	-0.3%	0.663 $\pm$ 0.006	1.941 $\pm$ 0.013
	Ours	$\checkmark$	<b>0.098 <math>\pm</math> 0.011</b>	<b>27.4%</b>	<b>0.283 <math>\pm</math> 0.009</b>	<b>3.8%</b>	<b>0.650 <math>\pm</math> 0.008</b>	<b>1.915 <math>\pm</math> 0.015</b>

generator’s architecture for the best performance-efficiency trade-offs. The length of the prompt, prefix, and suffix is 50 by default without explicit explanation for evaluation. For the main results in Section 4.2, the experiments are conducted for three runs, with the random seed of 0, 20 and 42. For the ablation study in Section 4.3 and the case study in Section 4.4, the random seed is 42. For details of the training and evaluation settings, please refer to the Appendix.

## 4.2 Main Results

The main results to evaluate our method and the SOTA baselines are summarized in Table 2, 3 and 4 for the suites of GPT-Neo, Pythia and StarCoder-Base, respectively.

In the application of text generation, our method can outperform all the baselines consistently and significantly. For GPT-Neo suite, compared with

the vanilla baseline (i.e., *No Prompt*), our method can achieve a relative improvement of 135.3%, 40.5% and 25.2% on average in terms of *Exact ER* with the model size of 125M, 1.3B and 2.7B, respectively. For the Pythia suite, our method can achieve a relative improvement of 107.7%, 42.5%, 25.8% and 19.2% over the vanilla baseline on average in terms of *Exact ER* with the model size of 410M, 1.4B, 2.8B and 6.9B, respectively.

In the application of code generation, our method can also outperform all the baselines consistently and significantly. For StarCodeBase suite, our method can achieve a relative improvement of 39.8%, 37.8% and 27.4% over the vanilla baseline on average in terms of *Exact ER* with the model size of 1B, 3B, and 7B, respectively.

We have several observations from the main results across diverse settings. Firstly, prepending naive hard prompts such as *Constant Hard Prompt*

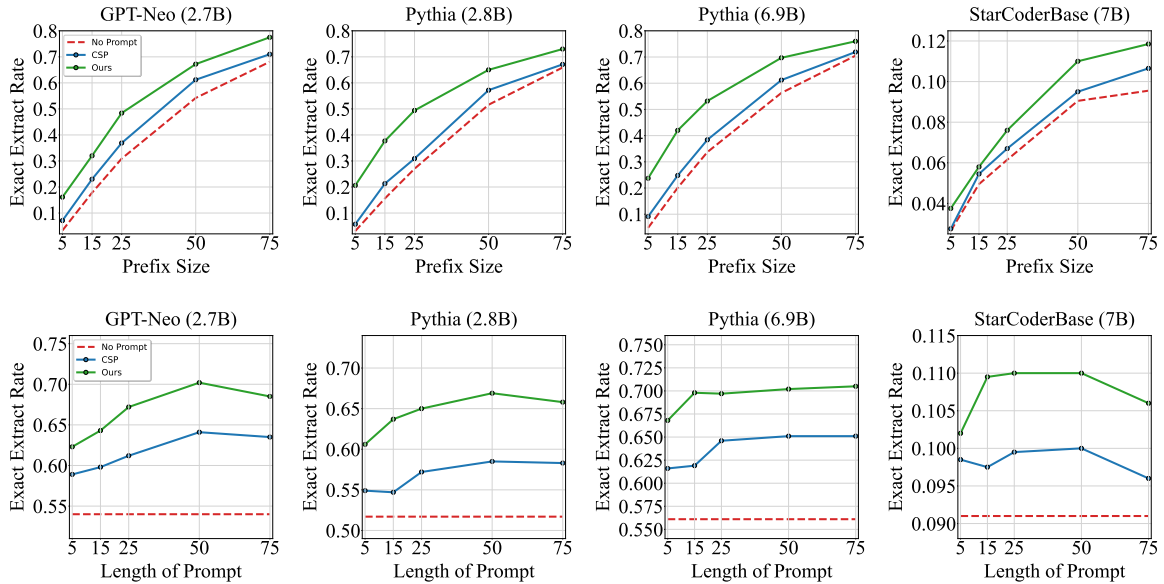


Figure 4: Ablation study on prefix size and length of prompt with exact extraction rate.

and *Dynamic Hard Prompt* is not useful but harmful for the exaction of training data from target LLM, leading to a much lower estimation of its memorization. Secondly, our method consistently outperforms the SOTA work, *CSP* (Ozdayi et al., 2023), highlighting the importance of dynamic soft prompts for the measure of memorization. Moreover, the memorization of LLM increases with the model size, which is consistent with the existing works (Carlini et al., 2023; Ozdayi et al., 2023; Nasr et al., 2023). However, it does not mean that the small model does not have security concerns on memorization. As shown in Table 2 and Table 3, the memorization of small language models with millions of parameters is underestimated significantly by the results from the previous methods. According to the results from our method, these models’ memorization cannot be ignored in real applications.

### 4.3 Ablation Study

We explore our methods from several perspectives: the impact of dynamic prompt, prefix size  $L$ , and the length of prompt  $N$ .

**Impact of Dynamic Prompt.** To evaluate the impact of dynamic prompt, we build another baseline by replacing the input to the generator with the first  $N$  tokens in the vocabulary of the target LLM. In this way, the soft prompts from the generator are constant and independent of the prefix tokens. The results are shown in Table 5. It can be observed that our method with dynamic prompt

Table 5: Ablation study on the dynamics of prompt.

Model	Method	Is Dynamic Prompt?	Exact ER	Fractional ER
<b>GPT-Neo (2.7B)</b>	CSP	✗	0.641	0.779
	Ours	✗	0.630	0.765
	<b>Ours</b>	✓	<b>0.702</b>	<b>0.820</b>
<b>Pythia (2.8B)</b>	CSP	✗	0.585	0.783
	Ours	✗	0.565	0.766
	<b>Ours</b>	✓	<b>0.669</b>	<b>0.827</b>
<b>StarCoderBase (3B)</b>	CSP	✗	0.081	0.249
	Ours	✗	0.081	0.247
	<b>Ours</b>	✓	<b>0.094</b>	<b>0.268</b>

outperforms the case with constant prompt consistently and significantly over all the evaluated settings. Moreover, the performance of our method with constant prompts is close to that of directly learning a constant soft prompt. Therefore, we can conclude that the advantage of our method over *CSP* comes from its adaptation to the dynamics of input instead of incorporating a transformer-based generator straightforwardly.

**Impact of Prefix Size.** To evaluate the impact prefix size, we set the length of prompt  $N$  to 25 and vary the prefix size for GPT-Neo (2.7B), Pythia (2.8B), Pythia (6.9B) and StarCoderBase (7B). The results in terms of *Exact ER* are shown in the first row of Figure 4. Our method can outperform the two representative baselines consistently over all the settings of prefix size across diverse LLMs and datasets. Moreover, the amount of extracted data increases along with the increase in the prefix size, consistent with the existing works (Carlini et al., 2023; Ozdayi et al., 2023).





observe that the improvement of our method on the fractional extraction rate is smaller and less robust compared with the improvement in the exact extraction rate. One possible reason is that the aligned CLM loss to train the generator is more suitable for the optimization of verbatim memorization. Since fractional extraction rate may be more important in cases where the meaning of the extracted sequences is more important than the exact match, it is valuable to improve the performance of our method on the metric of fractional extraction rate.

## 7 Ethical Considerations

In this work, we propose to leverage dynamic soft prompts to extract more training data from the target LLM and measure its memorization under the white-box settings. Therefore, it is possible that the attackers might utilize our method to extract sensitive data from the target LLM if they have white-box access to the target LLM. However, the main purpose of this work is to raise awareness among LLM researchers and developers about the security concerns caused by LLM memorization. By utilizing our method to evaluate the memorization of the target LLM, the owner of the LLM can evaluate its security vulnerability more accurately and thoroughly and then take action to defend against it. For example, we mentioned in the paper that the developer can utilize machine unlearning to forget the sensitive training data that is identified by our method. To minimize the security issues caused by our work, all of our experiments are conducted on public datasets that have been extensively studied by the research community.

## Acknowledgments

This project was supported by resources provided by the Office of Research Computing at George Mason University (URL: <https://orc.gmu.edu>) and funded in part by grants from the National Science Foundation (Award Number 2018631).

## References

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318.
- Amir Bar, Yossi Gandelsman, Trevor Darrell, Amir Globerson, and Alexei Efros. 2022. Visual prompt-

ing via image inpainting. *Advances in Neural Information Processing Systems*, 35:25005–25017.

- Rachel Bawden and François Yvon. 2023. Investigating the translation performance of a large multilingual language model: the case of bloom. *arXiv preprint arXiv:2303.01911*.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.
- BigCode. the-stack-smol dataset. <https://huggingface.co/datasets/bigcode/the-stack-smol>.
- Sid Black, Gao Leo, Phil Wang, Connor Leahy, and Stella Biderman. 2021. *GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow*.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2023. *Quantifying memorization across neural language models*. In *The Eleventh International Conference on Learning Representations*.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2021. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.
- Shoufa Chen, Chongjian Ge, Zhan Tong, Jiangliu Wang, Yibing Song, Jue Wang, and Ping Luo. 2022. Adapterformer: Adapting vision transformers for scalable visual recognition. *Advances in Neural Information Processing Systems*, 35:16664–16678.
- Mack DeGeurin. 2023. Oops: Samsung employees leaked confidential data to chatgpt. <https://gizmodo.com/chatgpt-ai-samsung-employees-leaked-data-1850307376>. Accessed on: 2023-04-06.
- 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*.
- Google-Research. Training data extraction challenge. <https://github.com/google-research/lm-extraction-benchmark>.
- Jiaxian Guo, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Boyang Li, Dacheng Tao, and Steven Hoi. 2023. From images to textual prompts: Zero-shot visual question answering with frozen large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10867–10877.

- Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. 2024. Exploring human-like translation strategy with large language models. *Transactions of the Association for Computational Linguistics*, 12:229–246.
- HF. Documentation of hugging face accelerate. <https://huggingface.co/docs/accelerate/>.
- Daniel Huynh. 2023. Starcoder memorization experiment highlights privacy risks of fine-tuning on code. <https://huggingface.co/blog/dhuynh95/starcoder-memorization-experiment>. Accessed on: 2023-11-02.
- Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher A Choquette-Choo, and Nicholas Carlini. 2023. Preventing generation of verbatim memorization in language models gives a false sense of privacy. In *Proceedings of the 16th International Natural Language Generation Conference*, pages 28–53. Association for Computational Linguistics.
- Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. 2022. Visual prompt tuning. In *European Conference on Computer Vision*, pages 709–727. Springer.
- Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. 2021. How can we know when language models know? on the calibration of language models for question answering. *Transactions of the Association for Computational Linguistics*, 9:962–977.
- Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. Deduplicating training data mitigates privacy risks in language models. In *International Conference on Machine Learning*, pages 10697–10707. PMLR.
- Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. 2023. Copyright violations and large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7403–7412.
- Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, et al. 2022. The stack: 3 tb of permissively licensed source code. *arXiv preprint arXiv:2211.15533*.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2021. Deduplicating training data makes language models better. *arXiv preprint arXiv:2107.06499*.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059.
- Junyi Li, Tianyi Tang, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2024. Pre-trained language models for text generation: A survey. *ACM Computing Surveys*, 56(9):1–39.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*.
- Dongze Lian, Daquan Zhou, Jiashi Feng, and Xinchao Wang. 2022. Scaling & shifting your features: A new baseline for efficient model tuning. *Advances in Neural Information Processing Systems*, 35:109–123.
- Antoine Louis, Gijs van Dijck, and Gerasimos Spanakis. 2024. Interpretable long-form legal question answering with retrieval-augmented large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 22266–22275.
- Fang Ma, Chen Zhang, Lei Ren, Jingang Wang, Qifan Wang, Wei Wu, Xiaojun Quan, and Dawei Song. 2022. Xprompt: Exploring the extreme of prompt tuning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11033–11047.
- Fatemehsadat Miresghallah, Archit Uniyal, Tianhao Wang, David K Evans, and Taylor Berg-Kirkpatrick. 2022. An empirical analysis of memorization in fine-tuned autoregressive language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1816–1826.
- Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. 2023. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*.
- Mustafa Ozdayi, Charith Peris, Jack FitzGerald, Christophe Dupuy, Jimit Majmudar, Haidar Khan, Rahil Parikh, and Rahul Gupta. 2023. Controlling the extraction of memorized data from large language models via prompt-tuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1512–1521.
- Junting Pan, Ziyi Lin, Yuying Ge, Xiatian Zhu, Renrui Zhang, Yi Wang, Yu Qiao, and Hongsheng Li. 2023. Retrieving-to-answer: Zero-shot video question answering with frozen large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 272–283.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.

- Martin Pawelczyk, Seth Neel, and Himabindu Lakkaraju. 2023. In-context unlearning: Language models as few shot unlearners. *arXiv preprint arXiv:2310.07579*.
- Alec Radford, Jeffrey Wu, Dario Amodei, Daniela Amodei, Jack Clark, Miles Brundage, and Ilya Sutskever. 2019. Better language models and their implications. *OpenAI blog*, 1(2).
- Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. 2023. Prompting large language models with answer heuristics for knowledge-based visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14974–14983.
- Yusheng Su, Xiaozhi Wang, Yujia Qin, Chi-Min Chan, Yankai Lin, Huadong Wang, Kaiyue Wen, Zhiyuan Liu, Peng Li, Juanzi Li, et al. 2022. On transferability of prompt tuning for natural language processing. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3949–3969.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Dave Van Veen, Cara Van Uden, Louis Blanke-meier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerová, et al. 2024. Adapted large language models can outperform medical experts in clinical text summarization. *Nature Medicine*, pages 1–9.
- Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou, and Daniel Cer. 2022. Spot: Better frozen model adaptation through soft prompt transfer. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5039–5059.
- Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. 2022. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 139–149.
- Colin Wei, Sang Michael Xie, and Tengyu Ma. 2021. Why do pretrained language models help in downstream tasks? an analysis of head and prompt tuning. *Advances in Neural Information Processing Systems*, 34:16158–16170.
- Chengyue Wu, Yukang Gan, Yixiao Ge, Zeyu Lu, Jiahao Wang, Ye Feng, Ping Luo, and Ying Shan. 2024. Llama pro: Progressive llama with block expansion. *arXiv preprint arXiv:2401.02415*.
- Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv preprint arXiv:2309.11674*.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2020. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934*.
- Jin Yao, Eli Chien, Minxin Du, Xinyao Niu, Tianhao Wang, Zezhou Cheng, and Xiang Yue. 2024. Machine unlearning of pre-trained large language models. *arXiv preprint arXiv:2402.15159*.
- Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2023. Large language model unlearning. *arXiv preprint arXiv:2310.10683*.
- Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. 2021. Qaggn: Reasoning with language models and knowledge graphs for question answering. *arXiv preprint arXiv:2104.06378*.
- Shenglai Zeng, Yaxin Li, Jie Ren, Yiding Liu, Han Xu, Pengfei He, Yue Xing, Shuaiqiang Wang, Jiliang Tang, and Dawei Yin. 2023. Exploring memorization in fine-tuned language models. *arXiv preprint arXiv:2310.06714*.
- Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study. In *International Conference on Machine Learning*, pages 41092–41110. PMLR.
- Haoyu Zhang, Jianjun Xu, and Ji Wang. 2019. Pretraining-based natural language generation for text summarization. *arXiv preprint arXiv:1902.09243*.
- Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B Hashimoto. 2024. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics*, 12:39–57.

## A Detailed Experimental Setup

**Training and Evaluation Settings.** The random train/test split is 14k/1k samples, 9686/1k samples and 38k/2k samples for GPT-Neo, Pythia and StarCoderBase, respectively. For evaluation, the generation’s decoding method is greedy decoding. All the experiments are conducted on a single NVIDIA A100 GPU with 80GB Memory in less than 12 hours for a single run. During the training of the generator, We used a batch size of 128 and an Adam optimizer for 15 epochs. We tried the learning rate from the range of  $[10^{-3}, 10^{-7}]$  and picked up the



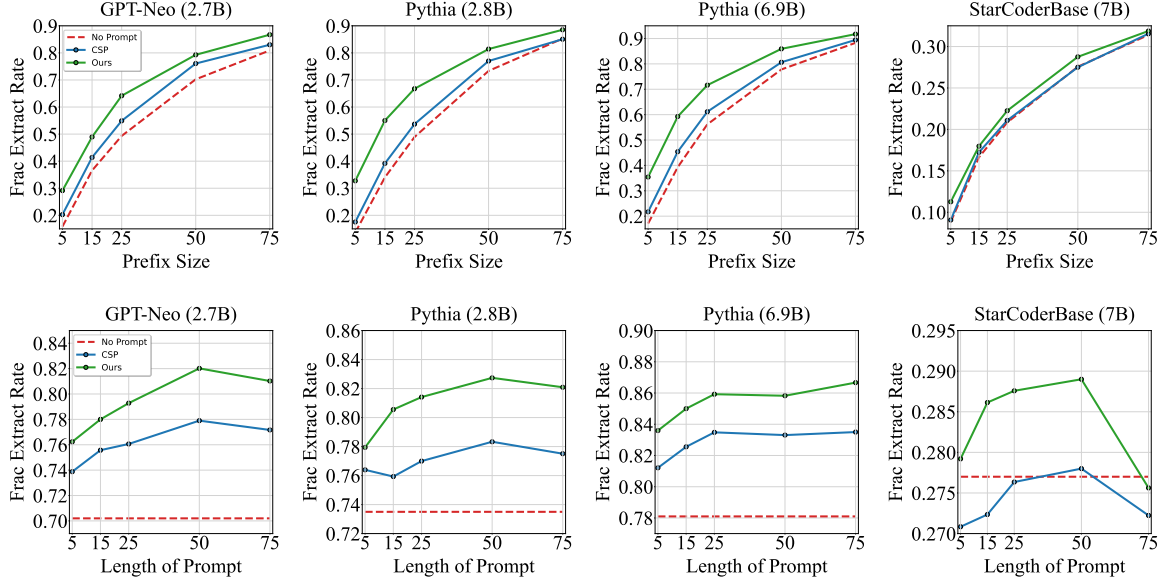


Figure 5: Ablation study on prefix size and length of prompt with fractional extraction rate for LLMs in the main paper.

best performance for each run and reported the average performance over the three runs in Section 4.2. **Discuss on the Artifacts.** The source code of our method is implemented with Pytorch (Paszke et al., 2019) and HuggingFace Accelerate (HF). And the implementation is built upon the open-sourced code released by Ozdayi et al. (2023). All the codes and datasets we utilize are public and open-sourced. They support the usage in research and we use them for research purpose only. We did not check whether the used data contains any information that names or uniquely identifies individual people or offensive content. We left this work to the institution that released the data.

## B Ablation Study on the LLMs in the Main Paper with Fractional Extraction Rate

Figure 5 shows the ablation study on the prefix size and length of prompt for GPT-Neo (2.7B), Pythia (2.8B), Pythia (6.9B) and StarCoderBase (7B) in terms of *Fractional ER*, which are the LLMs mentioned in the ablation study of main paper. And we have the same conclusion as the ones Section 4.3. More specifically, our method can outperform the baselines consistently over all the settings. And the saturation phenomenon can be observed in the ablation study on the length of prompt.

Table 8: Ablation study on the dynamics of prompt for more LLMs.

Model	Method	Is Dynamic Prompt?	Exact ER	Fractional ER
<b>GPT-Neo (125M)</b>	CSP	✗	0.239	0.421
	Ours	✗	0.254	0.423
	<b>Ours</b>	✓	<b>0.421</b>	<b>0.557</b>
<b>GPT-Neo (1.3B)</b>	CSP	✗	0.532	0.698
	Ours	✗	0.000	0.053
	<b>Ours</b>	✓	<b>0.651</b>	<b>0.772</b>
<b>Pythia (410M)</b>	CSP	✗	0.318	0.526
	Ours	✗	0.310	0.531
	<b>Ours</b>	✓	<b>0.513</b>	<b>0.683</b>
<b>Pythia (1.4B)</b>	CSP	✗	0.497	0.714
	Ours	✗	0.484	0.709
	<b>Ours</b>	✓	<b>0.617</b>	<b>0.786</b>
<b>Pythia (6.9B)</b>	CSP	✗	0.651	0.833
	Ours	✗	0.613	0.813
	<b>Ours</b>	✓	<b>0.702</b>	<b>0.858</b>
<b>StarCoderBase (1B)</b>	CSP	✗	0.071	0.235
	Ours	✗	0.079	<b>0.254</b>
	<b>Ours</b>	✓	<b>0.082</b>	0.244
<b>StarCoderBase (7B)</b>	CSP	✗	0.010	0.278
	Ours	✗	0.092	0.271
	<b>Ours</b>	✓	<b>0.110</b>	<b>0.289</b>

## C Ablation Study of the Dynamics of Prompt on the Remaining LLMs

Table 8 shows more results for the ablation studies on the dynamics of prompt, including all the LLMs not mentioned in the same ablation studies in Section 4.3 in the main paper. As shown in Table 8, same conclusion as the ones in Section 4.3 can be drawn. More specifically, our method with dynamic prompt outperforms the case with constant prompt consistently and significantly over all the evaluated settings. And the performance of our

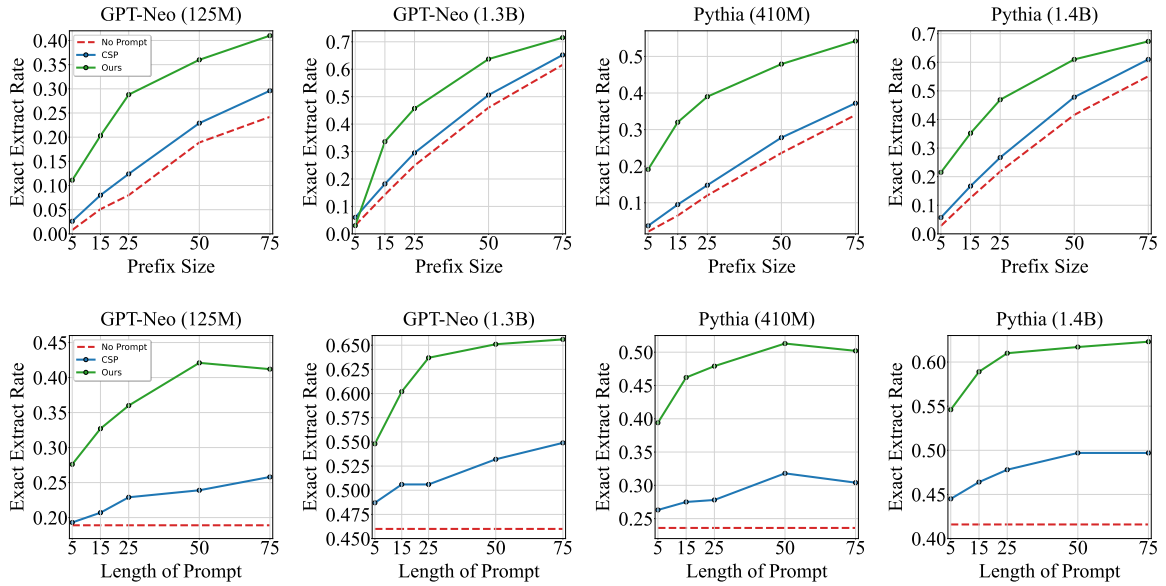


Figure 6: Ablation study on prefix size and length of prompt with exact extraction rate for the remaining LLMs for text generation.

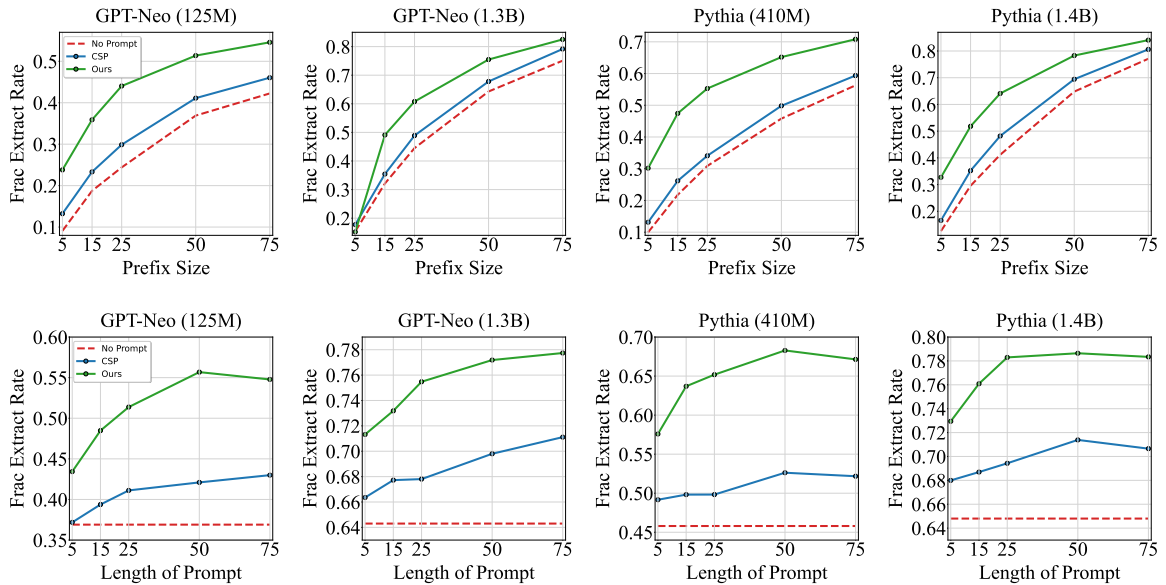


Figure 7: Ablation study on prefix size and length of prompt with fractional extraction rate for the remaining LLMs for text generation.

method with constant prompts is close to that of directly learning a constant soft prompt.

## D Ablation Study on the Remaining LLMs for Text Generation

Figure 6 shows the ablation study on the prefix size and length of prompt for GPT-Neo (125M), GPT-Neo (1.3B), Pythia (410M) and Pythia (1.4B) in terms of *Exact ER*, which are the remaining LLMs not mentioned in the ablation study of main paper for text generation. Figure 7 shows the same ablation study evaluated by *Fractional ER*. The

same conclusion can also be drawn as the ones in Section 4.3. More specifically, our method can outperform the baselines consistently over all the settings. And the saturation phenomenon can be observed in the ablation study on the length of prompt in terms of *Exact ER* and *Fractional ER*.

## E Ablation Study on the Remaining LLMs for Code Generation

Figure 8 shows the ablation study on the prefix size and length of prompt for StarCoderBase (1B) and StarCoderBase (3B) in terms of *Exact ER* and

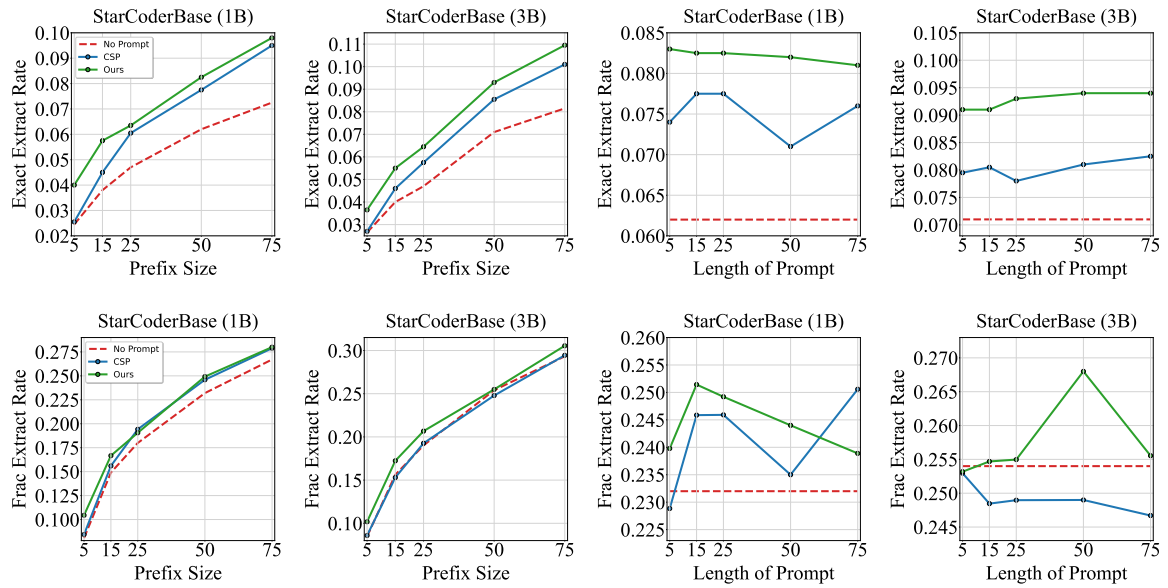


Figure 8: Ablation study on prefix size and length of prompt for the remaining LLMs for code generation.

*Fractional ER.* The two LLMs are the models not mentioned in the ablation study of the main paper for code generation. And we can conclude the same conclusion as the ones in Section 4.3 with the exception when evaluating the impact of the length of prompt on StarCoderBase (1B) in terms of *Fractional ER*. In this case, the performance of *CSP* outperforms our method when the length of prompt is 75. It indicates that more techniques are needed for our method to achieve robust and consistent improvement in terms of *Fractional ER*.