
Prompt Injection: Parameterization of Fixed Inputs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent works have shown that attaching prompts to the input is effective at conditioning Language Models (LM) to perform specific tasks. However, prompts
2 are always included in the input text during inference, thus incurring substantial
3 computational and memory overhead. Also, there is currently no straightforward
4 method of utilizing prompts that are longer than the maximum input length of
5 the LMs without incurring additional costs during inference. We propose Prompt
6 Injection (PI), a novel formulation of injecting the prompt into the parameters of an
7 LM to be an efficient alternative to attaching fixed prompts to the input. We show
8 that in scenarios with long fixed prompts, PI can be up to 280 times more efficient in
9 terms of total FLOPs than previous approaches. We further explore methodologies
10 for PI and show promising results in persona-dependent conversation, semantic
11 parsing, and zero-shot learning with task instructions. Through these explorations,
12 we show that PI can be a promising direction for conditioning language models,
13 especially in scenarios with long and fixed prompts¹.
14

15 1 Introduction

16 Contemporary works with large Language Models (LMs) [3, 32, 23, 5, 28] have shown that attaching
17 prompts (also referred to as *prefixes*) to the input is effective at conditioning LMs to perform specific
18 tasks. During training, LMs are trained to condition on the given prompts in hopes of generalizing
19 to unseen prompts during inference. Unseen prompts can be a persona for persona-dependent
20 conversation [39], database schema for semantic parsing [10], and task instruction for zero-shot
21 learning with task instructions [23]. In these tasks, a new prompt is fixed to the input at every
22 inference. For instance, in persona-dependent conversation [39, 18, 33], a persona description is
23 appended to the dialogue history, so that the LM can always be conditioned on the persona. For
24 another example, in semantic parsing, the LM is conditioned on the database schema as well as
25 natural language questions to generalize to a new database [37, 10, 36]. Lastly, zero-shot learning
26 with task instructions [32, 23] involves adding natural language instructions to the inputs for adapting
27 LMs to novel tasks.

28 However, concatenating prompts to input sequences for prompt-dependent inference has two major
29 limitations. (1) During inference, prompts are always included in the input text and thus incur
30 computational and memory overhead [16]. (2) It is challenging to fit a long text such as the detailed
31 description of a persona as a prompt into Transformer-based models whose input lengths are often
32 fixed [27]. For instance, in persona-dependent conversation, the model constantly refers to the persona
33 description along with the dialogue history [35, 22], as shown in the left side of Figure 1. Moreover,
34 in real world scenarios, a persona may consist of a long detailed text description of a character or
35 person, not just a few profile sentences. Naively concatenating long prompts to the input sequences is
36 challenging due to the quadratic cost in time and memory of Transformer-based architectures with

¹Code used for the experiments and a demo are available at this link

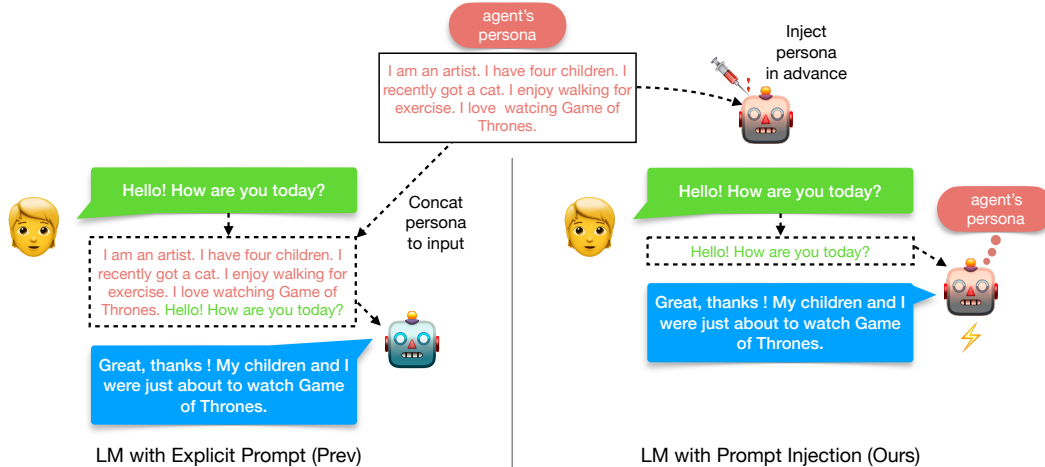


Figure 1: Prompt Injection example on a persona-dependent conversation. The left side presents an inference procedure of a previous approach where the persona (prompt) is concatenated to every input. The right side describes Prompt Injection, where the persona is *injected* into the model in advance, so that the model is able to generate responses without constantly referring to the persona description. Thus, Prompt Injection approach takes less time to generate responses than the previous method.

37 regard to the input sequence length. Other approaches specialized for long inputs [1, 13], such as
 38 Fusion-in-Decoder [12], or those that augment the LM with a retrieval mechanism [9] may be used
 39 but still come with increased overall memory and computations, ultimately leading to a delay in
 40 generating responses. This problem becomes critical in situations where the LMs are deployed, and
 41 fast inference speed is required.

42 In this work, we formulate a novel problem called Prompt Injection (PI), where we attempt to *inject* a
 43 given prompt into the parameters of an LM to address the two limitations mentioned above. With
 44 PI, LMs can produce prompt-dependent outputs without the computational overhead of appending
 45 fixed prompts at inference time (the right side of Figure 1), and it also enables the injection of longer
 46 prompts in a wholistic way. More specifically, we first show that PI is much more efficient (up to
 47 280 times) in terms of total FLOPs compared to previous approaches that may be used for handling
 48 long prompts such as Fusion-in-Decoder [12] or Linear Transformer [13]. Next, we explore different
 49 methodologies as baselines for PI, including the continued pre-training approach on the prompt as
 50 well as a novel distillation approach called Pseudo-INput Generation (PING), in order to analyze
 51 what components are effective for successful PI. We apply these PI methods to three different tasks
 52 with fixed prompts: persona-dependent conversation, semantic parsing, and zero-shot learning with
 53 instructions. We compare the methods against LMs with explicit prompts as the upper bound (i.e.,
 54 unconstrained) as well as the LM without both the prompt and PI as the lower bound. Experimental
 55 results show meaningful improvements with respect to the lower bound, but also exhibit a non-trivial
 56 gap with the upper bound. Despite the performance gap, we still believe that PI is a direction worth
 57 exploring considering the computational benefit of the injection, especially since inference speed is
 58 critical in real world applications.

59 In sum, our main contributions are three folds:

- 60 • We formally define the Prompt Injection (PI) formulation and demonstrate its necessity in
 61 terms of computation and memory efficiency, especially in scenarios with long prompts.
- 62 • We explore baseline approaches for PI, showing that performance can approach the upper
 63 bound (unconstrained) performance in some cases.
- 64 • We show that the *injection* of long prompts (e.g., detailed description of persona) can be
 65 achieved through PI and show its efficiency in comparison with previous methods, being up
 66 to 280 times more efficient during inference.

67 Through this work, we hope the community explores PI as an efficient alternative for performing
68 prompt-dependent tasks.

69 2 Related Work

70 **Prompting** Prompting is an emerging paradigm for modeling LMs, especially for few-shot and
71 zero-shot learning [20, 3, 21, 24, 32, 23]. With the help of appropriate prompts, one can exploit
72 knowledge learned by a pre-trained LM and manipulate the LM’s behavior. The benefit of prompting
73 is that the pre-trained LM can adapt to new scenarios with few or no labeled training data. However,
74 for the in-context learning scenario, processing prompts that involve many training examples for each
75 inference incurs substantial computational and memory overhead [16]. Given training data, Liu et al.
76 [16] replace in-context learning with fine-tuning a small set of parameters for tackling the above
77 issue. Prompt Injection also tackles the same issue but assumes a stricter scenario where there are no
78 training data for the given prompt.

79 **Efficient Transformers for Long Inputs** One can consider using efficient Transformer-based [29]
80 architectures for handling long input sequences [27]. The main challenge of using a vanilla Trans-
81 former architecture is the quadratic cost in time and memory with regard to the input sequence
82 length due to the self-attention operation. There has been a surge of recent works addressing this
83 problem [6, 38, 1, 13, 40, 8]. They are primarily dedicated to improving either the efficiency of the
84 self-attention mechanism or the general efficiency of the Transformer architecture through sparse mod-
85 els. Our Prompt Injection approach tackles the efficiency problem of performing prompt-dependent
86 tasks by keeping the input sequences short (without prompts), bounding the time and memory
87 complexity to a constant invariant of the length of the prompt.

88 **Persona-dependent Conversation** Endowing a chatbot with a persona [39, 18, 33] is challenging,
89 but it enables the chatbot to deliver more personal, specific, consistent, and engaging conversa-
90 tions [39] and gain user trust [17, 25, 19]. To achieve this, previous works have attached a persona to
91 the dialog history at every inference time, so that the model can always be conditioned on the persona.
92 However, given a long persona description, this approach brings the critical problem of increased
93 overall memory and computations, resulting in delayed response generation. An LM augmented
94 with a retrieval mechanism [9] may be used but still comes with non-trivial computational overhead.
95 Prompt Injection allows a dialogue agent to generate responses without a persona description as the
96 explicit input once the persona is injected.

97 **Semantic Parsing** Semantic parsing is the task of mapping a natural language query into a SQL
98 query executable on a database. Recently, the community has focused more on cross-domain (cross-
99 database) semantic parsing, where models are trained and tested on different domains (databases) [37].
100 The domain-adaptation setup introduces many generalization challenges, such as non-explicit column
101 names and domain-specific phrases [10], and recent works concatenate the natural language query
102 with the serialized database schema as the input to address the problem [26, 7, 36]. With Prompt
103 Injection, the model is adapted to a new database schema in advance, so that it can map natural
104 language queries to SQL queries on the new database without explicitly referring to the schema
105 during inference.

106 **Zero-shot Learning with Task Instructions** Recent works [23, 32] have addressed zero-shot
107 generalization to new tasks [3, 14] by multi-task prompted training. With multi-task prompted
108 training, the models learn to use task instructions as prompts to generalize to unseen tasks. It is
109 demonstrated that this approach improves generalization ability to novel tasks and offers an effective
110 substitute for unsupervised language model pre-training. Through Prompt Injection, the LM can be
111 aware of a novel task instruction before performing the task and thus does not require the instruction,
112 which can be lengthy, to make predictions.

113 3 Prompt Injection

114 In this section, we formally define Prompt Injection (PI) as a task and describe the benefits of the
115 formulation. Prompt-dependent generation is a task of generating an output sequence y that is a

116 proper response to the input sequence x and coherent to the prompt z . Utilizing the prompt during
 117 inference, the generated sentence is obtained by $y = f(z, x)$ where f denotes an LM such as
 118 T5 and GPT-2. Prompt Injection (PI), i.e., parameterization of prompts, allows LMs to perform
 119 prompt-dependent generation without using prompts during inference. To achieve this, we need to
 120 design a PI method H to inject a prompt z into an LM f . The process of PI can be represented as

$$f_z = H(z, f) \quad (1)$$

121 where f_z denotes an LM injected with the prompt. Then the prompt-dependent output sequence can
 122 be obtained by $y = f_z(x)$.

123 PI can also be applied for long prompts whose length exceeds the LM’s input sequence length. Given
 124 a long prompt z , we decompose it into multiple sub-prompts $\{z_i\}$ each of which fits the LM’s input
 125 length, i.e., $z = z_{1:n} = [z_1; z_2; \dots; z_n]$. Then the PI process can be executed iteratively, injecting
 126 each sub-prompt sequentially while the LM is aware of the previous sub-prompts:

$$f_{z_1} = H(z_1, f) \quad (2)$$

$$f_{z_{1:2}} = H(z_2, f_{z_1}) \quad (3)$$

...

$$f_{z_{1:n}} = H(z_n, f_{z_{1:n-1}}) \quad (4)$$

127 The above formulation can be seen as a high-level abstraction of iterative PI that we aim to ap-
 128 proximate. In practice, in order to fully inject $z_{1:n}$, we repeat (2)-(4) multiple times (i.e., multiple
 129 epochs).

130 **Why is Prompt Injection necessary?** Prompt Injection brings definite advantages when applied to
 131 prompt-dependent tasks. The previous approach of appending prompts to the input sequences has
 132 the drawback of the model repeatedly referring to the prompt at each inference time. This becomes
 133 critical in scenarios requiring long prompts, as Transformer architecture has quadratic computational
 134 and memory costs due to the limitation of the self-attention operation. We propose PI as a solution
 135 to this computation bottleneck. Once a prompt is injected into the LM in advance, the LM no
 136 longer needs to refer to the prompt during inference. As a result, the model’s input length remains
 137 independent of the length of prompts and is able to utilize prompts of any length efficiently. We
 138 discuss the efficiency gain of PI in Section 6.1.

139 **Evaluation Metric for Prompt Injection** PI can be evaluated by the evaluation metric of the
 140 fixed prompt-dependent task at hand. We also introduce a metric called the Prompt Injection
 141 score (PI score) to measure the degree of injection. The metric is agnostic of the target task by
 142 comparing the results with that of an LM given actual prompts during inference. Let $X_{w/ \text{prompt}}$
 143 denote the LM’s task score with the prompt as an additional input (upper bound) and $X_{w/o \text{prompt}}$
 144 denote the LM’s task score without the prompt (lower bound). We define **PI score** as the min-
 145 max scaling score of X_{PI} , where X_{PI} represents the score of the LM on the target task after PI,
 146 i.e., **PI score** = $\max(0, X_{PI} - X_{w/o \text{prompt}}) / (X_{w/ \text{prompt}} - X_{w/o \text{prompt}})$. We limit using PI
 147 only in situations where $X_{w/ \text{prompt}} > X_{w/o \text{prompt}}$ because there is no reason to inject a prompt
 148 if task performance degrades when using the prompt. Even if the range of individual task scores
 149 may vary from task to task, PI score represents the overall injection effectiveness of the PI methods,
 150 agnostic of the individual task score range.

151 4 Methods for Prompt Injection

152 In this section, we explore methods of Prompt Injection (PI) that can address prompt-dependent tasks
 153 without accessing the prompt during inference. To achieve this, the model should be trained to store
 154 the prompt in its parameters. This can be seen as *parameterizing* the prompt into the model instead of
 155 feeding the prompt explicitly to the model. This is challenging as the prompt is unseen to the model
 156 and has no corresponding training data. In Section 4.1, a baseline method by continued pre-training
 157 is introduced, followed by a method for improving the baseline with curriculum learning. Section 4.2
 158 presents a novel distillation-based method called Pseudo-INput Generation (PING) that learns to
 159 generate pseudo-inputs to inject novel prompts.

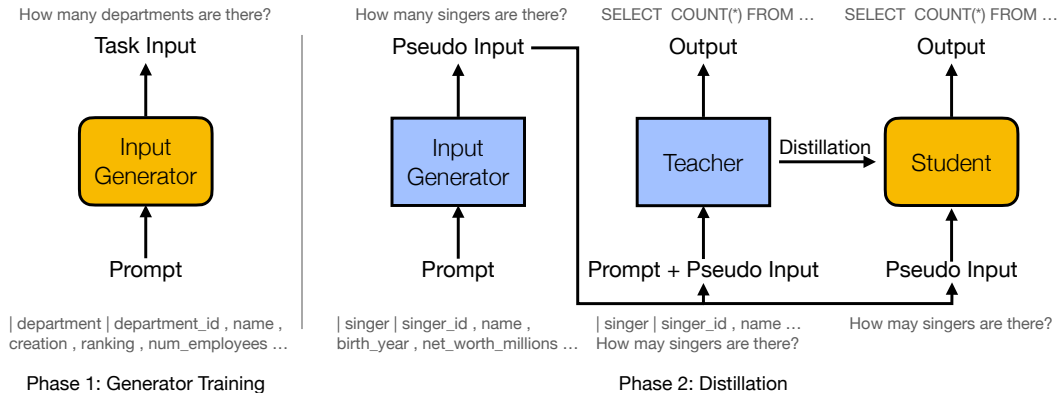


Figure 2: Illustration of the Pseudo-INput Generation (PING). During Phase 1, an input generator is trained with the task-specific training data. The inputs are prompts of a task, and the outputs are task inputs corresponding to the prompt. Input and output examples applied to semantic parsing are shown. During Phase 2, the input generator generates pseudo-inputs from the given target prompt, which are used to distill knowledge from the teacher to the student. Blue square boxes indicate frozen parameters; yellow rounded boxes indicate unfrozen parameters.

160 4.1 Continued Pre-training

161 We establish the Continued Pre-training method as a straightforward baseline for PI. This method
 162 injects prompts into the parameters of an LM by continuing with the pre-training objective of the
 163 LM on the target prompt. The pre-training objective is a straightforward option as it works in an
 164 unsupervised manner. In our experiments, we leverage the pre-trained T5 model [21] and thus use
 165 the masked language modeling objective which is the pre-training objective of T5. Following Raffel
 166 et al. [21], we randomly replace 15% of a given prompt with special mask tokens; then, the model is
 167 trained to predict the sequence of masked tokens. In this process, the model learns about the prompt
 168 the same way the model learns knowledge during the pre-training stage.

169 **Curriculum learning** We further investigate the baseline method by leveraging *curricula* [2, 4]
 170 during continued pre-training. We set the mask ratio as the difficulty criteria [34] and gradually
 171 increase the ratio throughout the Continued Pre-training. As the mask ratio increases, the model
 172 should predict more masked tokens given less context. With curriculum learning, we expect the LM to
 173 gradually better adapt to the prompt, improving its prompt-dependent task performance. Throughout
 174 the experiments, we increase the mask ratio linearly from 15% to 30%, 50%, and 70% and report the
 175 best score.

176 4.2 Pseudo-INput Generation (PING)

177 The purpose of PI is to inject a prompt into the parameters of an LM which can also be done indirectly
 178 through distillation. In this subsection, we propose a novel distillation-based method called Pseudo-
 179 INput Generation (PING) that distills a novel prompt into a student LM that does not have access
 180 to the prompt through a teacher LM that does have access to the prompt. In order for distillation,
 181 pseudo-inputs are needed since we assume a scenario where the prompt to be injected has never been
 182 seen during training and does not have separate training data. An overview of PING is illustrated in
 183 Figure 2. As shown in the figure, during Phase 1, an input generator is trained with the task-specific
 184 training data. When given a prompt of the task as the input, the generator is expected to generate the
 185 task inputs that correspond to the prompt. During Phase 2, the input generator is frozen and is used to
 186 generate pseudo-inputs from the unseen prompt, which are then given to the teacher together with the
 187 prompt, while only the pseudo-inputs are given to the student. This way, the student learns to follow
 188 the teacher and is able to learn about the prompt indirectly. We believe that this is the first work that
 189 aims to distill knowledge with different inputs for the teacher and the student.

190 5 Experimental Setup

191 In this section, we explain the experimental setups in detail. All experiments are performed with the
192 T5-base [21] (220M parameters) model unless noted otherwise.

193 5.1 Prompt-dependent tasks

194 In order to evaluate the effectiveness of Prompt Injection (PI) methods, we select three prompt-
195 dependent tasks—persona-dependent conversation, semantic parsing, and zero-shot learning with
196 task instructions; all these tasks require fixed prompts during inference. Fixed prompts come in the
197 form of a persona in persona-dependent conversation [39], database schema in semantic parsing [10],
198 and task instruction in zero-shot learning with task instructions [23]. As described in the introduction
199 and Section 3, when PI is applied for these tasks, there would be apparent benefits in real world
200 scenarios. For instance, PI eliminates the need to repeatedly include persona descriptions in the input
201 during inference when serving a conversational model of a specific personality. With these tasks,
202 not only the performance of the baseline PI methods is evaluated, but also the significance of PI is
203 emphasized by comparison with the (unconstrained) previous approaches that concatenate prompts to
204 the input.

205 5.2 Datasets

206 Following datasets of prompt-dependent tasks mentioned in Section 5.1 are utilized to evaluate
207 Prompt Injection (PI).

208 **PERSONA-CHAT** PERSONA-CHAT [39] is a crowd-sourced dataset intended for training agents
209 to perform engaging and personal chit-chat by comprising the dialogues to be grounded on specific
210 personas. They crowdsourced 1,155 unique personas, each with five profile sentences and 162,064
211 utterances over 10,907 dialogues. For each dialogue, two speakers have a 6-8 turn conversation
212 conditioned on a given persona. The task is measured via perplexity (PPL). We randomly select 100
213 dialogues from the validation set as persona-dependent conversation benchmark for testing PI. The
214 persona descriptions are 60 tokens long on average.

215 **Spider** Spider [37] is a large cross-domain semantic parsing and text-to-SQL dataset for developing
216 natural language interfaces to cross-domain databases. It includes 10,181 questions, 5,693 unique
217 SQL queries, and 200 database schemas covering 138 different domains. Models must generalize to
218 new database schemas as well as new queries to perform well on it. Evaluation metrics include Exact
219 Matching (EM) and Execution Accuracy (EA). We utilize the dev set containing 20 databases with
220 about 50 questions per database as a semantic parsing benchmark for PI. The database schemas range
221 in length from 55 to 430 token lengths.

222 **WSC / RTE / COPA** For the task of zero-shot task generalization, Raffel et al. [21] have trained
223 the LM on a diverse set of tasks and evaluated on a held-out group of tasks to evaluate generalization
224 performance. We choose coreference resolution, natural language inference, and sentence completion
225 tasks, three out of their four held-out tasks, and test PI on WSC (Winograd Schema Challenge), RTE
226 (Recognizing Textual Entailment), and COPA (Choice of Plausible Alternatives) datasets [30]. All of
227 these tasks are binary classification tasks. We utilize task instructions (prompts) of WSC, RTE, and
228 COPA provided from Raffel et al. [21] and report average task scores of using task instructions. The
229 task instructions are comprised of about 20-30 tokens.

230 5.3 Implementation Details

231 For the Continued Pre-training method (Section 4.1), we use the Adam optimizer [15] with a constant
232 learning rate $1e-4$ and batch size 8. We perform 5-20 steps of injection. For PING (Section 4.2),
233 input generators are trained on each tasks for 1-2 epochs. We use KL-divergence for distilling the
234 last layer’s output of the decoder and perform 10-40 steps of injection. Diverse pseudo-inputs are
235 generated by sampling each token from the output probability distribution of the decoder. For all of
236 the experiments except for zero-shot generalization, we use a single 16GB T4 GPU. For zero-shot
237 generalization, we use 4 32GB V100 GPUs.

Table 1: Inference efficiency of different models that can be used for performing prompt-dependent inference. We depict how many times PI is efficient in comparison with the other approaches inside the parenthesis. When there is out-of-memory (OOM) using the 16GB T4 GPU, we estimate the FLOPs in *italics* assuming a linear correlation between prompt length and FLOPs.

Model	Prompt Length	FLOPs (G)	Latency (s)
T5 w/ PI	*	0.7k	0.58
T5	512	7.2k ($\times 10.3$)	1.09 ($\times 1.9$)
	512 \times 2	14.6k ($\times 21.0$)	2.38 ($\times 4.1$)
	512 \times 4	OOM	-
T5 w/ FiD	512	7.2k ($\times 10.3$)	1.09 ($\times 1.9$)
	512 \times 2	14.0k ($\times 20.2$)	1.54 ($\times 2.6$)
	512 \times 4	27.6k ($\times 39.8$)	2.87 ($\times 4.9$)
	512 \times 8	54.9k ($\times 79.2$)	5.87 ($\times 10.0$)
	512 \times 28	OOM ($\times 280$)	-
LINEAR- TRANSFORMER	512	9.5k ($\times 13.8$)	1.58 ($\times 2.7$)
	512 \times 2	16.1k ($\times 23.2$)	2.62 ($\times 4.5$)
	512 \times 4	29.2k ($\times 42.2$)	4.74 ($\times 8.1$)
	512 \times 8	55.6k ($\times 80.1$)	9.11 ($\times 15.6$)
	512 \times 28	OOM ($\times 280$)	-

238 In order for injection and comparison with upper-bound and lower-bound performance, we first
 239 need two different versions of the LM adapted to the given task. For the task of persona-dependent
 240 conversation and semantic parsing, one (upper bound) is fine-tuned together with prompts since
 241 prompts are explicitly used during inference, while the other (lower bound) is fine-tuned on the task
 242 without being given the prompt. We perform PI on the lower-bound LM since we also assume having
 243 no access to prompts during inference.

244 For the zero-shot learning task, we modify the prompts developed by Raffel et al. [21]
 245 in the form of a fixed prompt. Their prompts have placeholders such as Premise, and
 246 Hypothesis. We replace the placeholders with fixed words such as "Premise" and "Hypothesis",
 247 then append the actual content to the prompt in a key-value format. For example,
 248 if the original is If {Premise} is true, is it also true that {Hypothesis}?, then
 249 the converted prompt is If "Premise" is true, is it also true that "Hypothesis"?
 250 Premise:{Premise} Hypothesis:{Hypothesis}. This ensures that the prompt is fixed, which
 251 can be injected with PI. We use the T0-3B LM checkpoint for the zero-shot generalization.

252 6 Experimental Results

253 In this section, we first explore the inference efficiency of models performing prompt-dependent tasks
 254 and show that Prompt Injection (PI) leads to meaningful computational efficiency. Then the baseline
 255 and proposed methods are tested and compared on datasets discussed in Section 5.2. The results
 256 indicate that the Pseudo-INput Generation (PING) method achieves the best performance among PI
 257 methods, sometimes even outperforming the unconstrained upper bound, which uses explicit prompts
 258 during inference. In Section 6.3, we provide a concrete instance of injecting a real persona description
 259 into a conversational model, demonstrating the feasibility of long prompt injection.

260 6.1 Inference Efficiency

261 The comparison of inference efficiency of a model with PI, a baseline model that naively concatenates
 262 prompts to the input, Fusion-in-Decoder (FiD) [12], and Linear Transformer [13] are shown in
 263 Table 1. We consider FiD as one of the options for processing long inputs because it processes
 264 long input sequences by encoding chunks of input sequences separately, reducing the quadratic
 265 complexity to linear. Linear Transformer also reduces the complexity to linear by linearizing the

Table 2: Prompt Injection performance on three prompt-dependent tasks. W/ PROMPT stands for the upper bound (unconstrained) method, which uses the prompt during inference by appending it to the input. W/O PROMPT depicts the lower bound method of not utilizing the prompts at all. Lastly, we show three w/ PI methods: CP and CP w/ CURR stand for the Continued Pre-training (baseline) and the Continued Pre-training with curricular, respectively, as explained in Section 4.1; PING depicts our novel proposed method utilizing distillation.

	Dialogue		Semantic Parsing			Task Generalization					
	PERSONA-CHAT		Spider			WSC		RTE		COPA	
	PPL (\downarrow)	PI Score	EM	EA	PI Score	ACC	PI Score	ACC	PI Score	ACC	PI Score
W/ PROMPT	8.83	-	57.9	61.3	-	63.6	-	<u>67.9</u>	-	67.3	-
W/O PROMPT											
W/O PROMPT W/O PI	11.01	-	14.5	15.1	-	44.0	-	64.2	-	60.0	-
W/ PI											
CP	10.85	0.073	16.9	17.5	0.054	54.5	<u>0.536</u>	67.7	<u>0.946</u>	<u>64.8</u>	0.658
CP w/ CURR	10.61	<u>0.183</u>	17.7	18.4	<u>0.072</u>	50.8	0.347	68.2	1.08	64.1	<u>0.562</u>
PING	<u>9.82</u>	0.546	<u>36.6</u>	<u>41.7</u>	0.507	63.7	1.005	64.2	0	60.6	0.082

266 attention mechanism. We measure FLOPs and forward propagation latency via DeepSpeed Flops
267 profiler² using a single 16GB T4 GPU.

268 As shown in Table 1, T5 w/ PI is much more efficient than other models, especially as we assume
269 a longer prompt length. This is because the efficiency of PI remains the same independent of the
270 prompt length while the costs of others increase linearly. Specifically, when the prompt length is
271 8 times the model’s max input sequence length, one can achieve 80× computational efficiency in
272 terms of FLOPs by applying PI. Furthermore, in a scenario where the prompt length is 28× the
273 model’s max input sequence length (shown in Section 6.3 when trying to utilize a long persona that is
274 over 13,000 token length long), previous approaches show an out-of-memory (OOM) issue using
275 the 16GB T4 GPU, and it is impossible to utilize them. PI is estimated to be 280× more efficient in
276 terms of total FLOPs if there is no OOM issue.

277 6.2 Task Performance

278 In Table 2, we report the task performance obtained by applying different PI methods on three
279 prompt-dependent tasks. PI scores are also obtained as introduced in Section 3. For all of w/ PI
280 methods, we observe an overall increase in performance compared to w/O PROMPT, indicating
281 successful injection of prompts into the parameters of the model through PI methods.

282 For the results, while CP gives modest performance improvement over w/O PROMPT, the results
283 of CP w/ CURR show that leveraging curricula during continued pre-training is effective in some
284 cases. CP w/ CURR performs better compared to CP in PERSONA-CHAT, Spider, and RTE; it even
285 outperforms w/ PROMPT in RTE. On the other hand, PING significantly improves performance from
286 CP in PERSONA-CHAT, Spider, and WSC, outperforming w/ PROMPT in WSC. This sheds light on
287 the possibility that PI may be able to reach the upper bound (unconstrained) performance. However,
288 the results show at the same time that there is still a gap between the performance of PI methods and
289 the upper bound w/ PROMPT that needs to be bridged in future work.

290 We find that the performance of different methods depends on the complexity of the input sequence
291 structure. We believe that PING achieves a good performance in PERSONA-CHAT, Spider, and
292 WSC because those datasets have relatively simple input sequences (short utterances; simple query;
293 a sentence and two words, respectively). In datasets with many components or multiple complex
294 sentences (e.g., COPA and RTE), the low quality of generated pseudo-inputs degrades the performance
295 of PING. On the other hand, CP and CP w/ CURR perform better in datasets with complex structure.
296 These findings encourage the community to explore a more integral PI method that can cover different
297 datasets.

²<https://www.deepspeed.ai/tutorials/flops-profiler/>

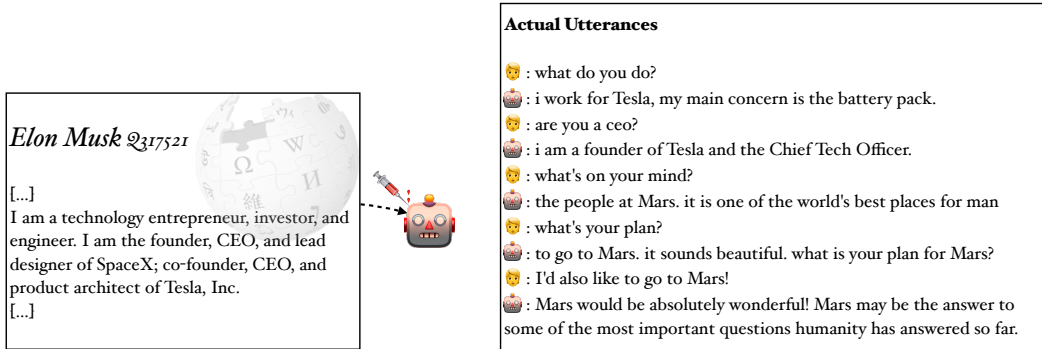


Figure 3: A real world example of Prompt Injection with a long prompt. (Left) The process of injecting a Wikipedia article describing a person (Elon Musk) into a model with PI. The article is more than 13,000 tokens long. (Right) Actual conversation between the persona injected model and a human that is hand-picked.

298 6.3 Long Prompts Injection

299 To demonstrate the effectiveness of PI on injection of long prompts into LMs, we show how the
 300 method works with a real world example. We pick a Wikipedia page (Elon Musk), considering it as a
 301 long persona description, and inject the entire article (over 13,000 tokens) into an LM trained with
 302 PERSONA-CHAT. Here, we use T5-large as a base model and apply PING.

303 Figure 3 shows an actual instance of interactions with the LM that underwent PI through PING. The
 304 responses show the LM successfully reflecting the description of the person on the Wikipedia page
 305 without having the description appended to the input. Moreover, the inference of PI is $280\times$ more
 306 computationally efficient in terms of FLOPs than the baseline, as shown in Section 6.1. Lastly, we
 307 provide a live demo to allow interactions with an LM injected with the persona of Elon Musk.

308 7 Conclusion

309 **Limitations and Future Work** While Prompt Injection (PI) enables performing prompt-dependent
 310 tasks efficiently, there are limitations that needs to be addressed in future work. In particular, the
 311 current PI methods cause task performance degradation. Moreover, the computational costs needed
 312 for the injection of prompts into the model parameters have not been extensively considered. For
 313 example, when considering *previous conversation history* as prompts to be injected in a multi-turn
 314 conversation setting, fast injection may also be a requirement for real-world application. Updating or
 315 adding a relatively small number of parameters [11, 31] may be a potential avenue for addressing the
 316 problems.

317 In this paper, we propose Prompt Injection (PI), a novel formulation of injecting the prompt into the
 318 parameters of an LM, as an efficient alternative to attaching fixed prompts to the inputs for prompt-
 319 dependent tasks. Through experiments, we show that PI is much more computationally efficient (up
 320 to 280 times) in terms of total FLOPs for handling long prompts compared to the previous alternatives.
 321 We further explore baseline methodologies for PI and find that Pseudo-INput Generation (PING), a
 322 distillation-based approach, shows promising results in persona-dependent conversation, semantic
 323 parsing, and zero-shot learning with task instructions. Through the explorations, we show that PI
 324 can be a promising direction for conditioning language models with prompts, especially in scenarios
 325 with long and fixed prompts. To this end, we hope the community explores PI for achieving both
 326 performance and efficiency on prompt-dependent tasks.

327 References

328 [1] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer.
 329 *ArXiv*, abs/2004.05150, 2020.

- 330 [2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning.
331 In *ICML '09*, 2009.
- 332 [3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
333 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,
334 Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M.
335 Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,
336 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Rad-
337 ford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *ArXiv*,
338 abs/2005.14165, 2020.
- 339 [4] Daniel Fernando Campos. Curriculum learning for language modeling. *ArXiv*, abs/2108.02170,
340 2021.
- 341 [5] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
342 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,
343 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Baidoor Rao, Parker Barnes,
344 Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchin-
345 son, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,
346 Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier
347 García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David
348 Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani
349 Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat,
350 Aitor Lewkowycz, Erica Oliveira Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee,
351 Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason
352 Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm:
353 Scaling language modeling with pathways. *ArXiv*, abs/2204.02311, 2022.
- 354 [6] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc V. Le, and Ruslan Salakhut-
355 dinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *ACL*,
356 2019.
- 357 [7] Xiang Deng, Ahmed Hassan Awadallah, Christopher Meek, Oleksandr Polozov, Huan Sun, and
358 Matthew Richardson. Structure-grounded pretraining for text-to-sql. *ArXiv*, abs/2010.12773,
359 2021.
- 360 [8] Mandy Guo, Joshua Ainslie, David C. Uthus, Santiago Ontañón, Jianmo Ni, Yun-Hsuan
361 Sung, and Yinfei Yang. Longt5: Efficient text-to-text transformer for long sequences. *ArXiv*,
362 abs/2112.07916, 2021.
- 363 [9] Seungju Han, Beomsu Kim, Jin Yong Yoo, Seokjun Seo, Sangbum Kim, Enkhbayar Erdenee,
364 and Buru Chang. Meet your favorite character: Open-domain chatbot mimicking fictional
365 characters with only a few utterances. *arXiv preprint arXiv:2204.10825*, 2022.
- 366 [10] Moshe Hazoom, Vibhor Malik, and Ben Bogin. Text-to-sql in the wild: A naturally-occurring
367 dataset based on stack exchange data. *ArXiv*, abs/2106.05006, 2021.
- 368 [11] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and
369 Weizhu Chen. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685,
370 2021.
- 371 [12] Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for
372 open domain question answering. In *EACL*, 2021.
- 373 [13] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and Francois Fleuret. Transformers
374 are rnns: Fast autoregressive transformers with linear attention. *ArXiv*, abs/2006.16236, 2020.
- 375 [14] Boseop Kim, Hyungseok Kim, Sang-Woo Lee, Gichang Lee, Donghyun Kwak, Dong Hyeon
376 Jeon, Sunghyun Park, Sung ju Kim, Seonhoon Kim, Dong Hyung Seo, Heungsub Lee, Minyoung
377 Jeong, Sungjae Lee, Minsub Kim, SukHyun Ko, Seokhun Kim, Taeyong Park, Jinuk Kim,
378 Soyoung Kang, Na-Hyeon Ryu, Kang Min Yoo, Minsuk Chang, Soobin Suh, Sookyo In,
379 Jinseong Park, Kyungduk Kim, Hiun Kim, Jisu Jeong, Yong Goo Yeo, Dong hyun Ham, Do-
380 Hyoung Park, Min Young Lee, Jaewoo Kang, Inho Kang, Jung-Woo Ha, Woo Chul Park,

- 381 and Nako Sung. What changes can large-scale language models bring? intensive study on
382 hyperclova: Billions-scale korean generative pretrained transformers. *ArXiv*, abs/2109.04650,
383 2021.
- 384 [15] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*,
385 abs/1412.6980, 2015.
- 386 [16] Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and
387 Colin Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context
388 learning. 2022.
- 389 [17] Qian Liu, Yihong Chen, B. Chen, Jian-Guang Lou, Zixuan Chen, Bin Zhou, and Dongmei Zhang.
390 You impress me: Dialogue generation via mutual persona perception. *ArXiv*, abs/2004.05388,
391 2020.
- 392 [18] Pierre-Emmanuel Mazaré, Samuel Humeau, Martin Raison, and Antoine Bordes. Training
393 millions of personalized dialogue agents. In *EMNLP*, 2018.
- 394 [19] Qiao Qian, Minlie Huang, Haizhou Zhao, Jingfang Xu, and Xiaoyan Zhu. Assigning personal-
395 ity/profile to a chatting machine for coherent conversation generation. In *IJCAI*, 2018.
- 396 [20] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
397 models are unsupervised multitask learners. 2019.
- 398 [21] Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
399 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified
400 text-to-text transformer. *ArXiv*, abs/1910.10683, 2020.
- 401 [22] Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu,
402 Myle Ott, Kurt Shuster, Eric Michael Smith, Y.-Lan Boureau, and Jason Weston. Recipes for
403 building an open-domain chatbot. In *EACL*, 2021.
- 404 [23] Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang A. Sutawika, Zaid
405 Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M SAIFUL
406 BARI, Canwen Xu, Urmish Thakker, Shanya Sharma, Eliza Szczechla, Taewoon Kim, Gunjan
407 Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang,
408 Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang,
409 Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan
410 Fries, Ryan Teehan, Stella Rose Biderman, Leo Gao, T. G. Owe Bers, Thomas Wolf, and
411 Alexander M. Rush. Multitask prompted training enables zero-shot task generalization. *ArXiv*,
412 abs/2110.08207, 2021.
- 413 [24] Timo Schick and Hinrich Schütze. It’s not just size that matters: Small language models are
414 also few-shot learners. *ArXiv*, abs/2009.07118, 2021.
- 415 [25] Haoyu Song, Weinan Zhang, Yiming Cui, Dong Wang, and Ting Liu. Exploiting persona
416 information for diverse generation of conversational responses. In *IJCAI*, 2019.
- 417 [26] Alane Suhr, Ming-Wei Chang, Peter Shaw, and Kenton Lee. Exploring unexplored generaliza-
418 tion challenges for cross-database semantic parsing. In *ACL*, 2020.
- 419 [27] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey.
420 *ACM Computing Surveys (CSUR)*, 2022.
- 421 [28] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam M. Shazeer, Apoorv Kulshreshtha,
422 Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, Yaguang Li, Hongrae Lee,
423 Huaixiu Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry
424 Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten
425 Bosma, Yanqi Zhou, Chung-Ching Chang, I. A. Krivokon, Willard James Rusch, Marc Pickett,
426 Kathleen S. Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos,
427 Toju Duke, Johnny Hartz Søraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz,
428 Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo,
429 Ravindran Rajakumar, Alena Butryna, Matthew Lamm, V. O. Kuzmina, Joseph Fenton, Aaron

- 430 Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguera-Arcas, Claire Cui, Marian Croak,
431 Ed Chi, and Quoc Le. Lamda: Language models for dialog applications. *ArXiv*, abs/2201.08239,
432 2022.
- 433 [29] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N.
434 Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *ArXiv*, abs/1706.03762,
435 2017.
- 436 [30] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill,
437 Omer Levy, and Samuel R. Bowman. Superglue: A stickier benchmark for general-purpose
438 language understanding systems. In *NeurIPS*, 2019.
- 439 [31] Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Jianshu Ji, Guihong Cao,
440 Daxin Jiang, and Ming Zhou. K-adapter: Infusing knowledge into pre-trained models with
441 adapters. In *FINDINGS*, 2021.
- 442 [32] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,
443 Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. *ArXiv*,
444 abs/2109.01652, 2021.
- 445 [33] Sean Welleck, Jason Weston, Arthur D. Szlam, and Kyunghyun Cho. Dialogue natural language
446 inference. In *ACL*, 2019.
- 447 [34] Alexander Wettig, Tianyu Gao, Zexuan Zhong, and Danqi Chen. Should you mask 15% in
448 masked language modeling? *arXiv preprint arXiv:2202.08005*, 2022.
- 449 [35] Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. Transfertransfo: A trans-
450 fer learning approach for neural network based conversational agents. *ArXiv*, abs/1901.08149,
451 2019.
- 452 [36] Tianbao Xie, Chen Henry Wu, Peng Shi, Ruiqi Zhong, Torsten Scholak, Michihiro Yasunaga,
453 Chien-Sheng Wu, Ming Zhong, Pengcheng Yin, Sida I. Wang, Victor Zhong, Bailin Wang,
454 Chengzu Li, Connor Boyle, Ansong Ni, Ziyu Yao, Dragomir Radev, Caiming Xiong, Lingpeng
455 Kong, Rui Zhang, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. Unifedskg: Unifying
456 and multi-tasking structured knowledge grounding with text-to-text language models. *ArXiv*,
457 abs/2201.05966, 2022.
- 458 [37] Tao Yu, Rui Zhang, Kai-Chou Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma,
459 Irene Z Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. Spider:
460 A large-scale human-labeled dataset for complex and cross-domain semantic parsing and
461 text-to-sql task. In *EMNLP*, 2018.
- 462 [38] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti,
463 Santiago Ontañón, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big
464 bird: Transformers for longer sequences. *ArXiv*, abs/2007.14062, 2020.
- 465 [39] Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur D. Szlam, Douwe Kiela, and Jason Weston.
466 Personalizing dialogue agents: I have a dog, do you have pets too? In *ACL*, 2018.
- 467 [40] Chen Zhu, Wei Ping, Chaowei Xiao, Mohammad Shoeybi, Tom Goldstein, Anima Anandkumar,
468 and Bryan Catanzaro. Long-short transformer: Efficient transformers for language and vision.
469 *ArXiv*, abs/2107.02192, 2021.