# Appendix

## A Related Work

**Prompt-Based Adaptation.** Prompt-based adaptation methods learn a small set of continuous vectors to guide frozen LLMs on downstream tasks. Early work like AutoPrompt (Shin et al., 2020) generated discrete prompts; Prefix Tuning (Li & Liang, 2021) introduced trainable vectors as layer-wise key-value cache entries; and Prompt Tuning (Lester et al., 2021) simplified the approach to input-only soft prompts. In contrast, *Context Tuning* initializes prompts or prefixes directly from demonstration pairs, allowing the model to leverage task-relevant in-context information from the start.

**In-Context Learning.** ICL (Radford et al., 2019) enables LLMs to perform tasks by conditioning on a few demonstration pairs without updating model parameters. Extensions include Chain-of-Thought prompting (Wei et al., 2022), self-consistency decoding (Wang et al., 2023), demonstration selection (Liu et al., 2021; Li & Qiu, 2023), and meta-training for better efficiency (Min et al., 2022a; Chen et al., 2022). However, recent work has shown that ICL often exploits surface-level patterns in demonstrations (Min et al., 2022b; Jang et al., 2024). *Context Tuning* addresses these limitations by tuning the demonstration context itself.

**Inference-Time Optimization.** *In-Context Optimization* fits within the broader category of inference-time optimization, which adapts models during test time. Test-Time Training has been applied to image classifiers (Sun et al., 2020), language models (Hardt & Sun, 2024), video tasks (Dalal et al., 2025), and more recently to few-shot learning (Akyürek et al., 2024). In parallel, diffusion models employ guidance techniques (Dhariwal & Nichol, 2021; Ho, 2022) to steer sampling, enabling controllable generation (Nichol et al., 2022; Wallace et al., 2023; Lu et al., 2024).

## **B** Additional Experiment Setup Details

We evaluate our methods and baselines on the following datasets.

- **NLP-LR** is the low-resource dataset split introduced by Min et al. (2022a), encompassing over 26 NLP tasks from Cross-Fit (Ye et al., 2021) and UnifiedQA (Khashabi et al., 2020), such as sentiment analysis and paraphrasing. Following Min et al. (2022a), we sample k = 16 demonstration pairs per task and evaluate task instances as multiple-choice problems.
- Massive Multitask Language Understanding (MMLU) is a diverse benchmark consisting of 57 subject-specific tasks, including mathematics, history, law, and various other domains (Hendrycks et al., 2021). We sample k = 16 demonstration pairs per task and evaluate task instances as multiple-choice problems.
- **BIG-Bench Hard (BBH)** is a curated subset of BIG-Bench, consisting of 27 tasks across 23 task types that challenge pretrained LLMs with questions involving algorithmic puzzles, symbolic manipulation, and other complex reasoning domains (Srivastava et al., 2023; Suzgun et al., 2022). Following Akyürek et al. (2024), we sample k = 10 demonstration pairs per task and prepend trainable instructions to all of our methods. Tasks are evaluated as question-answering problems.
- Abstraction and Reasoning Corpus (ARC) is a challenging symbolic reasoning benchmark with 400 evaluation tasks, each defined by a few grid transformation pairs and one or more query input grids (Chollet, 2019). Since the average number of available demonstration pairs is fewer than 4, we use all of them in context. Tasks are evaluated as question-answering problems.

We use k = 16 demonstrations for NLP-LR, 16 for MMLU, 10 for BBH, and all provided demonstrations for ARC. For multiple-choice tasks, where the LLM must select an output from a predefined set of answers, we follow Min et al. (2022a) in choosing the option with the lowest loss. For question-answering tasks, the LLM must generate an answer that matches the ground-truth output exactly. We use greedy decoding for all question-answering tasks. For ARC, we fine-tune our Llama3.2-1B checkpoint following the setup of Franzen et al. (2024), using 2 A100 GPUs for 24 epochs with a learning rate of  $2 \times 10^{-4}$ , a cosine learning rate scheduler, 1 warmup epoch, and a global batch size of 32 (after gradient accumulation). All inference-time experiments in Table 1 ran on a single A100 GPU, except for NLP-LR, which uses an RTX8000.

For TTT experiments, we follow Akyürek et al. (2024): using a LoRA learning rate of 1e-4, sampling a random permutation of the k demonstration pairs at each training step, and setting the LoRA rank to 128 for ARC and 64 for all other tasks. In our combined TTT+*CT*-*KV* experiments, we find that a small number of *CT*-*KV* training iterations and lower learning rates further improve performance on top of a TTT-adapted model.

<b>BBH</b> Instruction: Query: Answer <i>:</i>	<ul> <li>on: A logical deduction task which requires deducing the order of a sequence of objects. Answer with only the corresponding letter (e.g. (A)).</li> <li>The following paragraphs each describe a set of three objects arranged in a fixed order. The statements are logically consistent within each paragraph. In a golf tournament, there were three golfers: Ada, Mel, and Mya. Mya finished below Ada. Mel finished above Ada.</li> <li>Options: <ul> <li>(A) Ada finished last</li> <li>(B) Mel finished last</li> <li>(C) Mya finished last</li> </ul> </li> </ul>					
NLP-LR		ARC	(1)	(2)	(3)	Test
Query: W Cr Options: m Answer: P	/hat would you measure in a graduated ylinder? itrogen, Perfume, Oxygen, helium Perfume	Query	۲ <b>c</b> t		4	
MMLU Query: Fr ex Options: 0, Answer: 4	ind the degree for the given field attension $Q(sqrt(2) + sqrt(3))$ over $Q$ . <i>4, 2, 6</i>	Answer	2		4	1

Figure 3: One test pair from BBH, NLP-LR, and MMLU each, and 3 demonstration pairs followed by a test pair from ARC. BBH contains instructions that we prepend to model inputs. NLP-LR and MMLU contain multiple-choice options for the model to select. To avoid clutter, we show demonstration pairs from BBH, NLP-LR, and MMLU in Section G.

For all other experiments, we search over learning rates 3e-4, 1e-3, and 3e-3, and token dropout rates 0, 0.05, 0.1. We search training iterations 150, 200, 250, 300 for NLP-LR and ARC, 15, 20, 25, 30 for MMLU, and 12, 16, 20, 24 for BBH. To ensure a fair comparison, we perform hyperparameter sweeps for all methods. All experiments use the Adam optimizer, a cosine learning rate scheduler with no warm-up, bfloat16 precision, and up to 32GB of CPU RAM.

To fairly compare efficiency, we train each method with the largest batch size supported by the GPU used in its experiment. Since TTT, Prompt Tuning (m = # demo), and *CT-Prompt* consume more memory due to computing larger  $QK^T$  matrices (as shown in our derivation in Section C), we limit their batch sizes to 4 for NLP-LR, MMLU, and ARC, and 5 for BBH. MMLU and BBH models use gradient checkpointing. For all other methods, we use batch size 16 for NLP-LR, 8 for MMLU, 2 for BBH, and a full batch for ARC (depending on the number of demonstration pairs per task). None of these models require gradient checkpointing.

## C Time Complexity

At each training iteration, an LLM's forward and backward passes are dominated by its self-attention operations. Consider a single attention head of dimension d. Let  $L_Q$  denote the number of query tokens and  $L_K$  the number of key (and value) tokens. We form the query matrix  $\mathbf{Q} \in \mathbb{R}^{L_Q \times d}$  and the key and value vectors  $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{L_K \times d}$ , then compute

Attention 
$$(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \operatorname{softmax} \left( \frac{\mathbf{Q} \mathbf{K}^{\top}}{\sqrt{d}} \right) \mathbf{V},$$

whose dominant cost is the matrix multiplication  $\mathbf{Q} \mathbf{K}^{\top}$ , requiring  $O(L_Q L_K d)$  operations per head. Because d is a constant for a given model, we omit it in our comparisons below.

Next, let n be the number of tokens in the task's query and p the number of additional trainable prompt or prefix tokens per layer, we analyze how the training time of each method in *ICO* scales with n and p.

**Test Time Training.** At each layer of each training iteration, TTT prepends p trainable tokens to the n query tokens and computes their keys and values, giving  $L_Q = n + p$  and  $L_K = n + p$  with a per-head cost of

$$O\left((n+p)^2\right).$$

**CT-Prompt.** *CT-Prompt* prepends *p* trainable soft token embeddings to the query and computes their keys and values, also giving  $L_Q = n + p$  and  $L_K = n + p$  with a per-head cost of

$$O\left((n+p)^2\right).$$

**CT-KV.** Unlike from TTT and *CT-Prompt*, *CT-KV* prepends p trainable tokens as past keys and values, so these tokens do not generate queries. This yields  $L_Q = n$  and  $L_K = n + p$  with a per-head cost of only

$$O(n(n+p)).$$

**Time Complexity for k Demonstrations** Suppose we have k demonstration pairs, each of length  $\ell$  (assuming equal length). In TTT,  $n = \ell$  is the length of a demonstration pair and  $p = (k - 1)\ell$  is the summed length of other demonstration pairs. For *CT-Prompt* and *CT-KV*, n and p have the same values as TTT because Leave One Out masks out one of the in-context demonstration pairs. Table 2 summarizes the per-head costs in k and  $\ell$ , showing that both *CT-Prompt* and TTT incur quadratic cost in k, while *CT-KV* grows only linearly in k. This k-fold reduction in self-attention complexity explains *CT-KV*'s faster empirical training speed in Table 1.

Method	$L_Q$	$L_K$	Per-Head Cost
TTT	$k\ell$	$k\ell$	$O\left((k\ell)^2\right)$
CT-Prompt	$k\ell$	$k\ell$	$O\left((k\ell)^2\right)$
CT-KV	$\ell$	$k\ell$	$O(k \ell^2)$

Table 2: Per-head self-attention time complexity for methods with k demonstration pairs of length  $\ell$ .

## **D** Ablation Study

We perform ablations on our design choices for *CT-KV*, namely Leave-One-Out Masking and Token Dropout. Table 3 shows that across all benchmarks, *CT-KV* without Token Dropout performs marginally worse than *CT-KV* with both components. This suggests that when tuning more parameters than traditional Prefix Tuning, applying dropout along the token dimension of  $\Theta$  serves as an effective regularization technique for improving generalization. For NLP-LR, BBH, and MMLU, *CT-KV* performs significantly worse when Leave-One-Out Masking is not applied. This indicates that during training, it is crucial to mask out the portion of  $\theta_{context}$  corresponding to the demonstration pair being solved, as it prevents the model from cheating by retrieving the target output directly from the prefix initialization. However, on ARC, the model performs better without Leave-One-Out Masking. We hypothesize this is because ARC evaluation tasks typically include very few demonstration pairs (fewer than 4), so masking out even one pair during training can meaningfully reduce the effectiveness of the prompt or prefix in ICL. We also observe that when neither Leave-One-Out Masking nor Token Dropout is applied, *CT-KV* performs worse than ICL on MMLU and only marginally better on BBH, highlighting that these two design choices are essential to its overall performance.

Method	NLP-LR	MMLU	BBH	ARC
Neither	$41.0 \pm 0.75$	$40.2 \pm 0.73$	$51.4 \pm 0.76$	21.0
No Leave-One-Out Masking	$42.6 \pm 0.45$	$41.5 \pm 0.65$	$54.4 \pm 0.88$	23.8
No Token Dropout	$43.9 \pm 0.62$	$42.7 \pm 0.62$	$55.3 \pm 0.72$	21.0
Both	$44.2 \pm 0.55$	$43.7 \pm 0.54$	$57.9 \pm 0.78$	22.5

Table 3: Ablation study on the effects of Leave-One-Out Masking and Token Dropout in *CT-KV*. Means and standard deviations are computed over 5 seeds.

## **E** Qualitative Samples vs. Training Iteration

In this section, we select sample tasks from question-answering datasets to illustrate how autoregressively generated answers gradually improve with *CT-KV* training. We present two ARC tasks in Figure 5. In the top task, the model's prediction at iteration 0 (equivalent to ICL) shows a strong bias toward filling orange squares with yellow. As *CT-KV* training progresses, the model gradually learns to fill each orange square with the correct color. Similarly, in the bottom task, the model first learns that only grey grid cells can turn red, and then correctly completes the cross shapes.

Similarly, for BBH, in Figure 4's top query, the model initially predicts "padre, panicking" and "schoolmate, suburbia" in

instruction:	Soft a list of words. Answer with only a sequence of words.
Query:	Sort the following words alphabetically: List: wv habitation panicking schoolmate aerodynamic
	giacomo nebulous offset vector roosevelt jimmy botanist suburbia padre
Iteration 0:	aerodynamic botanist giacomo habitation jimmy nebulous offset panicking padre roosevelt suburbia schoolmate vector wv
Iteration 12:	aerodynamic botanist giacomo habitation jimmy nebulous offset panicking padre roosevelt suburbia schoolmate vector wy
Iteration 16:	aerodynamic botanist giacomo habitation jimmy nebulous offset panicking padre roosevelt schoolmate suburbia vector wv
Iteration 20:	aerodynamic botanist giacomo habitation jimmy nebulous offset padre panicking roosevelt schoolmate suburbia vector wy

Query:	Sort the following words alphabetically: List: scrumptious sidereal thermal yakima siena gorky saxon scottish figural hydroxyl seventeen neapolitan rampage nerve grapple fate plainfield stooge knives allotted
Iteration 0:	allotted fate figural gorky grapple hydroxyl knives neapolitan nerve plainfield rampage saxon scottish seventeen siena sidereal
	stooge thermal yakima
Iteration 12:	allotted fate figural gorky grapple hydroxyl knives neapolitan nerve plainfield rampage saxon scottish seventeen siena sidereal
	stooge thermal yakima
Iteration 16:	allotted fate figural gorky grapple hydroxyl knives neapolitan nerve plainfield rampage saxon scottish scrumptious seventeen
	siena sidereal stooge thermal yakima

Sort a list of words. Answer with only a sequence of words.

Iteration 20: allotted fate figural gorky grapple hydroxyl knives neapolitan nerve plainfield rampage saxon scottish scrumptious seventeen sidereal siena stooge thermal yakima

Figure 4: We display LLM predictions at *CT-KV* training iterations 0, 12, 16, 20 for two queries from the task "word sorting" in BBH. We omit showing the 16 demonstration pairs of each task for brevity. We color-code the iterations of correct predictions in green and incorrect predictions in red.

reversed order at iteration 0. During *CT-KV* training, the model learns to use the second letter of each word for sorting and eventually answers the query correctly. Likewise, for the bottom query, *CT-KV* helps the model avoid omitting the word "scrumptious" from its outputs and sort the words "sidereal, siena" into the correct order based on their second letters.

## F Qualitative Analysis

G ( 1') C

Instruction:

We compare our *CT-KV* to ICL on the 400 ARC evaluation tasks. Table 4 shows the confusion matrix indicating the number of tasks solved or not solved by each method. *CT-KV* recovers 51 tasks that ICL fails to solve, demonstrating the benefit of tuning the key and value representations corresponding to the in-context demonstration pairs. However, *CT-KV* fails to solve 9 tasks that ICL is able to, despite initializing its trainable prefix with the same demonstration pairs, suggesting that it can overfit to the few-shot examples.

	ICL correct	ICL wrong
CT-KV correct	44	51
CT-KV wrong	9	296

Table 4: Confusion matrix for the number of solved/unsolved ARC tasks by ICL and *CT-KV*.

Figure 6 shows one failure case for each method, where the other successfully solves the task. The task on the left illustrates that CT-KV can effectively adapt to the demonstration pairs to solve a geometric puzzle involving cropping the upper-left portion of objects in the query. On the right, we show a case where CT-KV makes an incorrect prediction. Since CT-KV performs optimization on the 3 demonstration pairs and two of them, illustrated on the right side of Figure 6, have answer grids that are 3-row by 4-column, we hypothesize that CT-KV became incorrectly biased toward predicting a grid of the same shape during optimization.

## **G** Demonstration Pairs for Figure **3**

We present three demonstration pairs of datasets: BBH, NLP-LR, and MMLU in Figure 7, Figure 8, and Figure 9, respectively.



Figure 5: For each of the two ARC tasks at the top and bottom, we display 4 demonstration query-answer pairs, the test query, and LLM predictions at *CT-KV* training iterations 0, 50, 100, 150, 200. Note that iteration 0 is equivalent to ICL. We color-code the iterations of correct predictions in green and incorrect predictions in red.



Figure 6: Left is an ARC task that CT-KV successfully solves, but ICL does not. Conversely, the task on the right is solved by ICL but not by *CT-KV*.

(1)	Instruction: Query: Answer:	<ul> <li>A logical deduction task which requires deducing the order of a sequence of objects.</li> <li>Answer with only the corresponding letter (e.g. (A)).</li> <li>The following paragraphs each describe a set of three objects arranged in a fixed order.</li> <li>The statements are logically consistent within each paragraph. In an antique car show, there are three vehicles: a motorcycle, a limousine, and a convertible. The motorcycle is newer than the limousine. The convertible is newer than the motorcycle.</li> <li>Options: <ul> <li>(A) The motorcycle is the oldest</li> <li>(B) The limousine is the oldest</li> </ul> </li> </ul>
(2)	Instruction: Query: Answer:	<ul> <li>A logical deduction task which requires deducing the order of a sequence of objects.</li> <li>Answer with only the corresponding letter (e.g. (A)).</li> <li>The following paragraphs each describe a set of three objects arranged in a fixed order.</li> <li>The statements are logically consistent within each paragraph. On a shelf, there are three books: a blue book, an orange book, and a red book. The blue book is the rightmost. The orange book is the leftmost.</li> <li>Options: <ul> <li>(A) The blue book is the second from the left</li> <li>(B) The orange book is the second from the left</li> <li>(C) The red book is the second from the left</li> </ul> </li> </ul>
(3)	Instruction: Query: Answer:	<ul> <li>A logical deduction task which requires deducing the order of a sequence of objects.</li> <li>Answer with only the corresponding letter (e.g. (A)).</li> <li>The following paragraphs each describe a set of three objects arranged in a fixed order.</li> <li>The statements are logically consistent within each paragraph. In an antique car show, there are three vehicles: a motorcycle, a minivan, and a tractor. The minivan is older than the tractor. The minivan is the second-newest.</li> <li>Options: <ul> <li>(A) The motorcycle is the newest</li> <li>(B) The minivan is the newest</li> <li>(C)</li> </ul> </li> </ul>

Figure 7: 3 demonstration pairs for the BBH task from Figure 3.

(1)	Query: Options: Answer:	Cellular respiration releases blood, waste, snot, feces waste
(2)	Query: Options: Answer:	During what period of the Earth cycle would you see someone having a picnic outside? Day, Night, Extinction, Ice Age Day
(3)	Query: Options: Answer:	Which uses gills to breathe? hermit crab, human, blue whale, bluebird hermit crab

Figure 8: 3 demonstration pairs for the NLP-LR task from Figure 3.

(1)	Query: Options: Answer:	The inverse of -i in the multiplicative group, {1, -1, i, -i} is <i>1</i> , - <i>1</i> , <i>i</i> , - <i>i</i>
(2)	Query: Options: Answer:	Find the degree for the given field extension Q(sqrt(2), sqrt(3), sqrt(18)) over Q. 0, 4, 2, 6 4
(3)	Query: Options: Answer:	Find the order of the factor group (Z_11 x Z_15)/(<1, 1>) 1, 2, 5, 11 1

Figure 9: 3 demonstration pairs for the MMLU task from Figure 3.