CT3: Boosting Downstream Performance through Test-time Training on AI PCs with Remote Multi-Domain Knowledge Bases

Anonymous ACL submission

Abstract

Learning at test time is an effective strategy for improving the performance of large language models (LLMs), although at the expense of increased computational costs during inference. In this paper, we introduce Collaborative Test-Time Training (CT3), a novel system designed to enhance the downstream accuracy of LLMs on client devices such as AI PCs through Test-Time Training (TTT) leveraging a remote multi-domain knowledge base. CT3 efficiently distributes the TTT process using a server-client architecture, allowing clients to fine-tune their models using relevant samples from the server. It also proposes a local state management mechanism and a simple but effective sample size reduction strategy to optimize test-time training without compromising accuracy. Our experiments demonstrate significant accuracy improvements across multiple domains and various LLMs with up to 44% increase in average downstream performance and a speedup ranging from $1.5 \times$ to $2.5 \times$ compared with vanilla TTT. The code to reproduce CT3's results will be released open-source.

1 Introduction

011

018

019

027

034

042

There is an increasing interest in techniques for enhancing model performance by using additional computational resources at test time rather than scaling model parameters during training (Snell et al., 2025). Test-time training (TTT) is one such approach that leverages test-time compute resources to fine-tune the model during inference. This paradigm improves the robustness and accuracy of large pre-trained models, addressing the key challenge of *covariate shift*, wheretest and training data distributions differ, potentially causing suboptimal performance. TTT enhances language modeling performance by accessing a reference dataset (Hardt and Sun, 2024; Hübotter et al., 2024), from which a subset of data points is retrieved and used for temporary model fine-tuning.



Figure 1: Performance overview of CT3 on downstream tasks across various LLMs. "Downstream Task Performance" indicates the average score on six downstream tasks with multiple domains evaluated in §4.

043

044

045

047

051

057

060

061

062

063

064

065

066

067

When the size of the TTT reference dataset reaches several petabytes, using these techniques on resource-constrained devices at the Edge becomes challenging due to their limited resources and the need for real-time fine-tuning, which demands significant computational power and storage capacity. Our work is motivated by these challenges and a new segment of client computing devices, namely Artificial Intelligence Personal Computers (AI PCs), which are equipped with AI accelerators such as Neural Processing Units (NPUs) and exhibit outstanding power efficiency. These AI PCs efficiently run sophisticated language models locally, such as LLAMA-3-8B-INSTRUCT (Dubey et al., 2024). The increasing demand for AI PCs motivates the development of solutions to improve model performance on these devices.

This paper addresses these challenges by proposing Collaborative **TTT** (CT3), a solution that enables client devices to benefit from TTT without storing large reference datasets locally. Specifically, our approach leverages an efficient clientserver solution to allow the limited device to benefit from the knowledge stored on a remote server while minimizing the communication rounds. To



Figure 2: Overview of Collaborative TTT (CT3). The performance of models at client devices improves by utilizing Test-Time Training with a remote multi-domain knowledge base.

enhance the setup's sophistication, we explore using a multi-domain knowledge base and investigate the tolerance of TTT algorithms to complex data mixes. CT3 allows users to run custom mod-071 els in their resource-constrained devices with TTT, overcoming limitations such as training data availability while benefiting from downstream task performance improvements. CT3 introduces Local State Management (LSM) and Sample Size Reduction (SSR) strategies to optimize the test-time fine-tuning process on client devices. LSM lever-079 ages historical query embeddings and their associated fine-tuned adapters to potentially bypass the need for repeated fine-tuning, thus significantly reducing computational overhead and latency. SSR strategies complement LSM by filtering out irrelevant samples, ensuring that only the most relevant 084 samples are used, thereby reducing the number of training samples and accelerating the process without compromising accuracy. Figure 1 provides a performance overview of CT3. The following sections discuss these contributions:

- 1. CT3, a system for improving the downstream reasoning capabilities of LLMs running in client devices (Figure 2).
- 2. An investigation into using multi-domain knowledge bases with clear supervisory signals in the context of test-time training for downstream tasks.
- 3. Local State Management (LSM) and Sample Size Reduction (SSR) strategies to optimize the test-time fine-tuning process, reducing computational overhead and improving system responsiveness.

097

101

The rest of the paper is organized as follows. We discuss related work in §2. Then, §3 describes the CT3 system, while §4 presents results on various downstream tasks. Our final thoughts are in §5.

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

136

2 Related Work

Test-time compute techniques have been proposed as an alternative to increasing the number of model parameters for improving language model performance (Snell et al., 2025). Notable test-time strategies include Chain-of-Thought (CoT) prompting (Wei et al., 2022) and few-shot learning (Brown et al., 2020). CoT prompting guides the model through intermediate steps to break down complex tasks, enhancing its ability to handle intricate queries. Few-shot learning helps the model adapt to new tasks with just a few examples. Other test-time compute methods include verifying the model's results, e.g., by code execution (Brown et al., 2025). The renewed focus on test-time compute has even motivated the development of improved neural architectures, e.g., Titans (Behrouz et al., 2024).

Among the many recent test-time compute approaches proposed, *Test-Time Training (TTT)* (Hardt and Sun, 2024; Hübotter et al., 2024; Akyürek et al., 2024) has been crucial in improving the performance of solutions, e.g., Omni-ARC (IronbarArc24, 2024) in the ARC-AGI challenge (ARC-AGI, 2025; Chollet, 2019).

TTT effectively adapts the model by fine-tuning it with selected samples from an available training dataset during inference. Recently, SIFT (Hübotter et al., 2024) demonstrated that it could retrieve informative samples, outperforming traditional Nearest-Neighbor (NN) retrieval (Hardt and Sun, 2024). Their experimental setup uses the

237

188

Pile dataset (Gao et al., 2020), which lacks a clear question-answer pair structure despite covering a wide range of topics and fields with plain text data. This weak supervision signal makes the training process more challenging, focusing on the model's language modeling capabilities (Pile benchmark (Gao et al., 2020)). In contrast, CT3 prioritizes downstream task performance. Our knowledge base comprises structured question-answer pairs with clear supervisory signals, which helps the model explicitly learn the relationship between input and output during training.

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

159

160

161

163

165

166

167

168

169

170

172

173

174

175

176

177

178

179

180

In the context of distributed architectures, the work by Hardt and Sun (2024) proposes a clientserver structure to speed up query time across large datasets. Their approach splits FAISS (Douze et al., 2024) indexes across multiple servers, allowing clients to send queries to each server and aggregate the results. This method primarily focuses on accelerating query processing and handling large-scale data efficiently. In contrast, our CT3 leverages a server-client architecture to facilitate TTT using a remote multi-domain knowledge base. While both approaches utilize distributed architectures, our system is designed to enhance model performance on resource-constrained client devices by retrieving relevant samples for fine-tuning during inference. Their architecture (Hardt and Sun, 2024) can be integrated into our remote server to further optimize sample retrieval, highlighting the potential of combining between these approaches.

Next, we explore CT3, a system that enhances model downstream performance on client devices.

3 CT3 System

This section introduces Collaborative Test-Time Training (CT3), a system designed to distribute Test-Time Training to allow resource-constrained clients to enhance their model downstream performance using a remote multi-domain knowledge base. CT3 also incorporates Local State Management (LSM) and Sample Size Reduction (SSR) strategies to speed up the fine-tuning stage without sacrificing accuracy. The following sections detail each of CT3's components and strategies.

1813.1 Preliminaries: Test-Time Training (TTT)182Given a pre-trained model and a query x, a TTT183solution utilizes a similarity metric ϕ to retrieve k184samples from a reference dataset \mathcal{D} . These samples are then used to fine-tune the model weights186at test time, resulting in improved model performance for the current query. Next, we discuss how

CT3 leverages the TTT paradigm in a distributed setting to improve model performance in resourceconstrained client devices.

3.2 Distributed Test-time Training

CT3 operates efficiently in heterogeneous compute environments (Figure 2), following a server-client paradigm implemented using FastAPI (Ramírez, 2025). In this section, we discuss the functionality of CT3 at both the server and client levels.

Server The remote server hosts a comprehensive multi-domain knowledge base \mathcal{D} , which encompasses training samples across commonsense reasoning, reading comprehension, coding, and math reasoning domains. The data is stored along with their corresponding embeddings, which are generated using a pre-trained model from sentence-transformers (Reimers and Gurevych, 2019). The dataset can be represented as $\mathcal{D} = \{(x_i, \mathbf{e}_i, y_i)\}_{i=1}^d$, where x_i denotes the *i*-th input sentence, \mathbf{e}_i represents its embedding, and y_i is the associated label or output.

The server receives client queries as embeddings. To ensure consistency, servers and clients are required to use the same sentence-transformers to generate embeddings. A private embedding model and encryption can increase privacy protection in data-sensitive applications. Upon receiving a client query embedding e_x , following SIFT (Hübotter et al., 2024), the server utilizes FAISS (Douze et al., 2024) to retrieve n relevant samples $\mathcal{D}_{\mathcal{F}} = \{(x_i, \mathbf{e}_i, y_i)\}_{i=1}^n$ from the knowledge base. The similarity metric ϕ employed by FAISS is the inner product, i.e., $\phi(\mathbf{e}_x, \mathbf{e}_i) = \mathbf{e}_x \cdot \mathbf{e}_i$, calculated between the embedding of the current query e_x and the embeddings e_i of candidate samples from the dataset \mathcal{D} . CT3 might use other search strategies to yield better performance based on the multi-domain data characteristics and mix.

To obtain more accurate and relevant samples, the SIFT algorithm (Hübotter et al., 2024) is applied to further refine the results retrieved by FAISS. SIFT filters the samples retrieved by FAISS, resulting in a smaller set, $\mathcal{D}_{S} = \{(x_i, \mathbf{e}_i, y_i)\}_{i=1}^k$ $(k \ll n)$, that is returned to the client without the embedding \mathbf{e}_i for efficient transmission.

The server can extract and transmit the corresponding samples to clients based on the above process. The client utilizes these samples to finetune its custom model, thereby increasing the accuracy and quality of the output.

Clients On a client device, the workflow is as 238 follows: a deployed large language model receives 239 a query x from the user. The client locally generates an embedding of the query e_x , using the same 241 sentence-transformers model as the server. This 242 embedding is then transmitted to the remote server. As discussed above, the server sends back a set 244 \mathcal{D}_{S} with k relevant samples to the query, excluding their embeddings. Using $\mathcal{D}_{\mathcal{S}}$, the client temporarily fine-tunes its custom model via LoRA (Hu et al., 247 2022) for efficiency, resulting in improved performance during inference. Finally, nothing prevents 249 a client with resources that satisfy the requirements of smaller deployments to run the complete CT3 251 pipeline on a single device. 252

> **In-Domain Sample Selection** CT3 does not know the domain of the user's query in advance. During development, we explored two approaches for assisting the system when searching for indomain samples:

254

256

257

261

265

266

267

269

273

- 1. User guidance to determine the domain
- 2. Automatic domain clustering purely based on embeddings

Although user guidance (1) is a viable strategy, our initial objective was to reduce user friction and minimize user actions to improve the system performance. In the automatic domain clustering approach (2), domain specialization occurs during the online search of relevant samples for test-time fine-tuning, allowing CT3 to provide samples tailored to the domain-specific query from the user. Experimentally (§4.3), we have observed that automatic domain clustering is feasible but with several challenges. For instance, the knowledge base must have the correct data mix, which presents an additional challenge to system administrators.

Local State Management (LSM) The fine-274 tuning stage at the client is the most expensive operation in the CT3 pipeline. For this reason, we propose caching strategies at the sampling and model state levels. A multi-turn historical Local 278 State Management (LSM) strategy is described 279 in Algorithm 1. Based on the availability of resources at the client, CT3 can budget the tracking of m past query embeddings Q and their associated groups of already fine-tuned LoRA adapters \mathcal{A} . Each new user query embedding is compared with 284 historical query embeddings to identify the one that maximizes a similarity function γ , indicating the possibility of using the group of associated weight

adapters and skipping the TTT stage to speed up inference without affecting accuracy. A simple version of Algorithm 1, is the one-look-back strategy in which m = 1 and CT3 forgets beyond the immediate previous query.

Algorithm 1 Multi-Turn Historical Local State Management (LSM)

Input: Current query embedding e_{x^t} , set \mathcal{Q} of m historical query embeddings, i.e., $\mathcal{Q} = \{e_{x^1}, ..., e_{x^m}\}$, set \mathcal{A} of m historical TTT fine-tuned adapter groups, i.e., $\mathcal{A} = \{a_{x^1}, ..., a_{x^m}\}$, reuse threshold τ , budget B for the number of adapter groups kept in the client device, similarity metric γ , and an eviction mechanism κ to determine the historical query to be evicted.

Output: Current adapter group state S_t .

1: $e_{x^*} = \operatorname{argmin} \gamma(e_{x^i}, e_{x^t})$ $e_{x^{i}} \in \mathcal{Q}$ 2: if $\gamma(e_{x^*}, e_{x^t}) > \tau$ then 3: // Retrieve corresponding adapter group 4: $S_t \leftarrow a_{x^*}$ 5: else 6: $S_t \leftarrow \mathrm{TTT}(e_{x^t}, S_0)$ 7: if $|\mathcal{Q}| > B$ then 8: // Evict one historical query using κ 9: $e_{x^r} = \kappa(\mathcal{Q})$ 10: $\mathcal{Q} = \mathcal{Q} \setminus e_{x^r}$ $\mathcal{A} = \mathcal{A} \setminus a_{x^r}$ 11: 12: end if 13: // Add new query and adapter to LSM 14: $\mathcal{Q} \cup e_{x^t}$ 15: $\mathcal{A} \cup S_t$ 16: end if 17: // Inference using S_t

Sample Size Reduction (SSR) We also explore simple but effective sample size reduction strategies to reduce further the samples selected by SIFT. The strategies follow this pattern: CT3 obtains statistics from the multiset of acquisition values, $\mathcal{V}(|\mathcal{V}| = k)$, associated with the *k* samples selected by SIFT. In particular, CT3 explores using the minimum of the mean and the median on \mathcal{V} as the cut-off point for discarding selected samples. Based on these statistics, CT3 selects a subset, \mathcal{V}' , such that $|\mathcal{V}'| < |\mathcal{V}|$.

4 Experiments

We implement a prototype of the CT3 system as a testbed, demonstrating the potential benefits in a larger deployment. Next, we discuss the resources utilized in our experimentation, followed by results demonstrating the benefits of enabling test-time training in client devices.

4.1 Setup

Knowledge base To rigorously evaluate the CT3 prototype, we utilize a comprehensive multidomain knowledge database. This knowledge base 293

294

295

296

297

298

300

301

302

303

304

305

306

307

308

309

310

311

312

313

| Dataset | Domain | # Samples |
|--|---------|-----------|
| CoQA (Reddy et al., 2019) | Reading | 7,199 |
| MetaMath (Yu et al., 2023) | Math | 395,000 |
| Orca-Math (Mitra et al., 2024) | Math | 200,035 |
| Math 50K (Zhiqiang et al., 2023) | Math | 50,000 |
| The Stack Python (Kocetkov et al., 2022) | Coding | 600,000 |
| MBPP (Austin et al., 2021) | Coding | 374 |
| Total | / | 1,252,608 |
| | | |

Table 1: Knowledge base composition. These datasets cover various domains, including reading comprehension, math, and coding. The Stack Python dataset consists of a random sample of 600,000 entries from the original Stack dataset (Python). All these datasets are training sets and do not contain any of the test samples we evaluated.

is constructed by integrating data from several diverse datasets (Table 1). Specifically, we incorporate the CoQA (Reddy et al., 2019) training set, a large-scale dataset designed for developing conversational question-answering systems, which is expected to enhance the model's reading comprehension capabilities.

For mathematical problem solving, we include three substantial datasets: MetaMath (Yu et al., 2023), a dataset comprising 395,000 samples designed to improve mathematical reasoning by bootstrapping questions from multiple perspectives; Orca-Math (Mitra et al., 2024), which includes 200,035 high-quality synthetic math problems created using a multi-agent setup, and Math 50K (Zhiqiang et al., 2023), which consists of 50,000 samples drawn from various mathematics-related datasets.

In the coding domain, we integrate The Stack Python, a dataset we constructed by randomly sampling 600,000 entries from the Python subset of The Stack (Kocetkov et al., 2022), which covers a wide range of programming languages and serves as a pre-training dataset for code-generating AI systems. Additionally, we include the Mostly Basic Python Problems (MBPP) (Austin et al., 2021) training dataset, which consists of 374 crowdsourced Python programming problems aimed at entry-level programmers.

This amalgamation of datasets ensures a robust and diverse foundation for evaluating the CT3 system's performance across multiple domains, including Reading Comprehension, Math, and Coding. **Evaluation** We employ a combination of evaluation tools for our experiments to rigorously assess CT3 prototype's performance across various domains. Specifically, we utilize the *lm-eval-harness* (Gao et al., 2023) for evaluating CoQA (Reddy et al., 2019) and GSM8K (Cobbe et al., 2021). For GSM8K, to better align with real-world TTT scenarios, we use zero-shot instead of the commonly used few-shot approach. For more mathematical reasoning evaluations, we also apply the evaluation scripts of LLM-Adapters (Zhiqiang et al., 2023) on MathQA (Amini et al., 2019) and MAWPS (Koncel-Kedziorski et al., 2016) datasets. For coding, we evaluate MBPP (Austin et al., 2021) and HumanEval (Chen et al., 2021) utilizing bigcodeevaluation-harness (Ben Allal et al., 2022). For experimental efficiency, we randomly sampled a subset of 200 samples from GSM8K and MathQA for evaluation, respectively.

348

349

350

351

352

353

354

355

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

389

390

391

392

393

394

395

397

Models We test our method on various LLMs, including LLAMA-3-8B-INSTRUCT, LLAMA-3.1-8B-INSTRUCT (Dubey et al., 2024), MISTRAL-7B-INSTRUCT-V0.3 (Jiang et al., 2023) and QWEN2.5-3B-INSTRUCT (Yang et al., 2024). For the retrieval stage in CT3, we utilize ALL-MPNET-BASE-V2 (Reimers and Gurevych, 2019) as a surrogate embedding model.

Hyperparameters and Implementation We employ two epochs for test-time training with a learning rate of 5e-5. We adjust the learning rate for coding tasks to either 1e-5 or 3e-5. The batch size is set to 1 across all tasks. We apply LoRA finetuning at test time, with a LoRA rank of 128 and a LoRA alpha of 16. The target modules for LoRA include the Query, Key, and Value and the Up and Down projection layers. Note that greedy decoding is applied to generate the responses for all tasks. The number of samples for test-time training is selected from the hyperparameter space [4, 8, 16, 32, 64] based on the best results. For LSM, the reuse threshold τ and k in LSM is 0.5 and 8, respectively. The similarity metric is the inner product, and the eviction mechanism is the Least Frequently Used (LFU) cache. Regarding the retrieval stage, we follow the pipeline of SIFT (Hübotter et al., 2024). More explorations can be found in §4.4 and Appendix C.

4.2 Main Results

Table 2 presents the performance of various largelanguage models on multiple evaluation tasks, com-

337

338

341

343

344

347

315

316

| Method | CoQA EM | CoQA F1 | GSM8K [*] EM (0-shot) | MathQA [*] Acc. | MAWPS Acc. | MBPP Pass@1 | HumanEval Pass@1 | Avg. | Impr. | Evaluation Speedup |
|--------------------------|------------|------------|-----------------------------------|-----------------------------|---------------|----------------|---------------------|-------|-------|-----------------------|
| LLAMA-3-8B-INSTRUCT | 61.78 | 78.40 | 40.00 | 30.00 | 43.70 | 51.60 | 54.88 | 51.48 | / | / |
| + CT3 | 70.38 | 82.25 | 60.50 | 36.00 | 86.13 | 52.40 | 57.93 | 63.66 | +24% | 1.00× |
| + CT3 + SSR | 70.68 | 82.46 | 55.00 | 28.50 | 85.71 | 52.40 | 57.32 | 61.72 | +20% | 1.76× |
| + CT3 + LSM | 70.17 | 82.24 | 59.50 | 41.50 | 89.08 | 51.60 | 59.15 | 64.75 | +26% | 1.60× |
| + CT3 + LSM + SSR | 70.20 | 82.17 | 56.00 | 33.00 | 84.45 | 51.80 | 56.10 | 61.96 | +20% | 2.24× |
| LLAMA-3.1-8B-INSTRUCT | 63.63 | 78.82 | 23.00 | 15.50 | 25.21 | 52.20 | 59.76 | 45.45 | / | / |
| + CT3 | 71.12 | 83.48 | 70.00 | 21.50 | 89.50 | 53.60 | 62.20 | 64.48 | -42% | 1.00× |
| + CT3 + SSR | 70.52 | 82.49 | 67.00 | 28.00 | 90.76 | 52.00 | 61.59 | 64.62 | +42% | 1.51× |
| + CT3 + LSM | 72.23 | 83.95 | 70.00 | 28.00 | 90.76 | 52.20 | 61.59 | 65.53 | +44% | 1.46× |
| + CT3 + LSM + SSR | 71.02 | 82.78 | 61.00 | 22.50 | 87.82 | 51.40 | 60.98 | 62.50 | +38% | 1.85× |
| MISTRAL-7B-INSTRUCT-V0.3 | 65.55 | 80.04 | 17.00 | 28.00 | 68.07 | 37.60 | 33.54 | 47.11 | / | / |
| + CT3 | 71.32 | 83.08 | 24.50 | 32.50 | 86.97 | 38.80 | 37.20 | 53.48 | +14% | 1.00× |
| + CT3 + SSR | 70.85 | 82.65 | 14.50 | 29.00 | 78.99 | 38.00 | 34.76 | 49.82 | +6% | 1.71× |
| + CT3 + LSM | 70.58 | 82.62 | 13.50 | 29.50 | 83.61 | 37.80 | 38.41 | 50.86 | +8% | 1.61× |
| + CT3 + LSM + SSR | 69.80 | 82.14 | 10.50 | 34.00 | 78.15 | 38.60 | 35.98 | 49.88 | +6% | 2.35× |
| QWEN2.5-3B-INSTRUCT | 54.70 | 71.22 | 61.50 | 20.00 | 87.82 | 41.20 | 26.22 | 51.81 | / | / |
| + CT3 | 67.27 | 80.53 | 61.50 | 28.00 | 88.24 | 46.80 | 29.88 | 57.46 | -11% | 1.00× |
| + CT3 + SSR | 66.22 | 79.77 | 63.50 | 33.50 | 70.17 | 45.80 | 27.44 | 55.20 | +7% | 1.23× |
| + CT3 + LSM | 67.43 | 80.59 | 63.50 | 32.50 | 88.66 | 44.40 | 26.83 | 57.70 | +11% | 1.65× |
| + CT3 + LSM + SSR | 67.15 | 80.56 | 65.00 | 31.00 | 61.34 | 43.40 | 29.27 | 53.96 | +4% | 1.48× |

Table 2: Performance comparison of different versions of CT3. "SSR" stands for Sample Size Reduction, which approximately halves the TTT training samples for each query. "Evaluation Speedup" shows the speedup factor achieved by our proposed strategies. "EM" represents Exact Match. "Impr." indicates the improvement in the average score. Note that * for GSM8K and MathQA indicates a subset of 200 randomly sampled test samples.



Figure 3: Domain distribution of training data samples across evaluation tasks (counted all test samples). The red border indicates the domain corresponding to the current task (in-domain). We expect higher values within the red border, meaning that the test-time training samples from the knowledge base are more aligned with the domain of the current task.

paring the baseline models without CT3, with CT3,
and CT3 incorporating the LSM and SSR strategies. The results demonstrate that CT3 significantly enhances the performance across all downstream tasks and models. For example, regarding
LLAMA-3-8B-INSTRUCT, CT3 improves the aver-

age performance from 51.48 to 63.66, representing a 24% improvement. When combined with LSM, CT3 achieves an average score of 64.75, a 26% improvement, and with both LSM and SSR, it achieves a 20% improvement with a $2.24 \times$ evaluation speedup compared to CT3. Similarly, LLAMA-3.1-8B-INSTRUCT shows a remarkable 42% improvement with CT3 and achieves higher performance with LSM, achieving up to $1.46 \times$ speedup. The results indicate that CT3 significantly improves model performance across various tasks and domains. The LSM and SSR strategies effectively maintain high performance while substantially reducing the TTT computational overhead, making CT3 more feasible for resource-constrained devices.

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

4.3 In-Domain Sample Retrieval

To evaluate the effectiveness of CT3 in retrieving relevant in-domain samples from multi-domain data, we conducted experiments to analyze the domain distribution of the retrieved training samples across various evaluation tasks. Figure 3 illustrates

the distribution of TTT data samples retrieved for 426 each evaluation task, categorized by their respec-427 tive domains. 428

431

432 433

437

441

443

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

The heatmap reveals that most retrieved samples 429 align well with the domain of the current task, as 430 indicated by higher values within the red-bordered cells. For instance, 82.50% of the samples retrieved for the CoQA task are from the CoQA dataset, demonstrating strong in-domain alignment. This 434 is particularly noteworthy given that the CoQA 435 dataset constitutes only a small fraction of the 436 knowledge base (7,199 out of 1,252,608 samples). Despite the overwhelming presence of unrelated 438 data, CT3's retrieval mechanism can precisely iden-439 tify and select relevant CoQA samples for the 440 CoQA task. Similarly, for GSM8K, 54.93% of the samples are from MetaMath, and for MathQA, 442 71.53% are from Orca-Math, indicating that CT3 successfully identifies and selects relevant mathe-444 matical datasets for these tasks. 445

> These alignments are crucial for the success of test-time training, as it ensures that the model is fine-tuned with high-quality, relevant samples, thereby significantly enhancing performance across various tasks. The effectiveness of CT3 in achieving such precise retrieval is attributed to the robust embedding models (Reimers and Gurevych, 2019) and retrieval algorithms employed. These tools enable CT3 to navigate a vast and diverse knowledge base, accurately matching queries to the most relevant samples. Overall, the results highlight the importance of retrieving high-quality, relevant samples for test-time training, allowing CT3 to enhance model performance significantly by leveraging domain-specific data effectively.

Additionally, we performed a t-SNE visualization of the embeddings for some samples in the database, as shown in Figure 4. The visualization demonstrates a clear clustering of samples by domain that CT3 can exploit to retrieve domainspecific samples. For example, the coding datasets (The Stack Python and MBPP) exhibit clear clustering, validating the embedding model's capability to differentiate between domains. This clustering effect underscores the robustness of the embedding model and the retrieval algorithm, enabling CT3 to achieve high precision in sample retrieval. The t-SNE visualization highlights the importance of high-quality embeddings in facilitating accurate and efficient sample retrieval, making test-time training more practical and impactful.



Figure 4: t-SNE visualization of the embeddings for some samples in the knowledge base. Each dataset contains 1000 randomly selected samples, except for MBPP (all).

| Faiss Indexes | Peak CPU Memory Usage | Avg. Score |
|---------------|--------------------------|---------------|
| / | / | 51.48 |
| FlatIP | 11568 MB | 63.66 |
| IVFPQ | 7936 MB | 63.47 |

Table 3: Comparison of different Faiss indexes regarding Peak CPU Memory Usage and Avg. Score for LLAMA-3-8B-INSTRUCT. The Avg. Score represents the average performance across six downstream tasks, while the peak CPU memory usage was measured using 16 test samples from the GSM8K dataset.

Memory-Efficient Indexing Alternatives 4.4 for CT3 Setup

In the CT3 prototype, we use *flat inner product* (FlatIP) indexes for retrieval in FAISS. A more memory-efficient indexing can be used if the user desires to run the complete CT3 system in an AI PC. As described in Table 3, using inverted file product quantization (IVFPQ) indexes reduces the peak CPU memory by 31% without significantly impacting the average score.

4.5 Exploring Smaller Knowledge Bases

We created a subset of the original knowledge base to demonstrate the effectiveness of a large knowledge base and the efficiency and feasibility of deploying smaller knowledge bases on local devices. This smaller knowledge base includes only the CoQA, Math 50K, and MBPP datasets, reducing the size from 1,252,608 samples to 57,573 samples. Table 4 presents the performance of CT3 using both the full and reduced knowledge bases.

The results show that while the smaller knowl-

| Model | Knowledge Size | CoQA EM | CoQA F1 | GSM8K [*] EM (0-shot) | MathQA [*] Acc. | MAWPS Acc. | MBPP Pass@1 | HumanEval Pass@1 | Avg. |
|-----------------------|-------------------------------|-------------------------|--|--|---|-------------------------|----------------------------------|----------------------------------|--|
| LLAMA-3-8B-INSTRUCT | / 57.57K | 61.78 71.90 | 78.40 82.99 | 40.00 59.00 | 30.00 34.50 | 43.70 75.63 | 51.60 51.80 | 54.88 57.93 | 51.48 61.96 |
| LLAMA-3.1-8B-INSTRUCT | 1.25M / 57.57K 1.25M | 63.63 72.57 71.95 | 83.30 78.82 83.35 83.82 | 63.50 23.00 62.50 70.00 | 32.00 15.50 22.00 28.00 | 25.21 84.03 88.66 | 52.40 52.20 52.60 53.00 | 59.76 59.76 62.80 62.80 | 64.30 45.45 62.83 65.46 |

Table 4: Performance comparison of CT3 using a knowledge base (1.25M samples, Table 1) versus a reduced knowledge base (57.57K samples). The reduced knowledge base includes only CoQA, Math 50K, and MBPP datasets. * for GSM8K and MathQA indicates a subset of 200 randomly sampled test samples.

edge base (57.57K) improves performance over the baseline models without CT3, the larger knowledge base (1.25M) consistently yields better results. For instance, with LLAMA-3-8B-INSTRUCT, the average performance increases from 51.48 (baseline) to 61.96 with the smaller knowledge base and to 64.30 with the larger knowledge base. This indicates that a larger and more diverse knowledge base provides more relevant samples for test-time training, leading to superior performance.

498

499

500

501

502

504

507

508 509

510

511

512

513

514

515

516

517

518

519

520

521 522

523

526

527

528

529

530

531

532

533

534

535

536

These findings highlight the trade-off between knowledge base size and performance. While a smaller knowledge base is more efficient and feasible for local deployment, a larger knowledge base offers significant performance gains.

4.6 Performance with Domain-Specific Databases

To further investigate the impact of retrieving samples exclusively from the correct domain, we performed experiments where the query samples were extracted only from the corresponding domainspecific dataset. Table 5 presents the performance comparison of CT3 using a mixed-domain database versus an in-domain database for CT3. The results indicate that the in-domain scenario performs comparably or slightly better than the multidomain scenario, suggesting that while domainspecific retrieval can enhance performance, the multi-domain setup proposed in this paper is feasible and effective. In real-world applications, incorporating a domain classifier or user-guided domain selection could further optimize the retrieval process, ensuring that the most relevant samples are used for test-time training, thus maximizing model performance.

5 Conclusion

Test-time Training (TTT) is an effective method for improving model performance at the expense of more computation at inference time. We present

| Task & Metric | Baseline | Domain | СТЗ СТ | T3 + SSR |
|-----------------------------------|----------|------------------|-----------------------|-----------------------|
| CoQA EM | 61.78 | Mix Reading | 71.90 71.63 | 69.85 70.05 |
| CoQA F1 | 78.40 | Mix Reading | 82.99 82.80 | 81.97 82.19 |
| GSM8K [*] EM (0-shot) | 40.00 | Mix Math | 59.00 61.00 | 54.50 61.00 |
| MathQA [*] Acc. | 30.00 | Mix Math | 34.50 32.00 | 33.00 30.00 |
| MAWPS Acc. | 43.70 | Mix Math | 75.63 78.99 | 68.07 81.51 |
| MBPP Pass@1 | 51.60 | Mix Coding | 51.80 52.00 | 52.80 52.80 |
| HumanEval Pass@1 | 54.88 | Mix Coding | 57.93 56.71 | 57.32 56.71 |
| Avg. | 51.48 | Mix In-Domain | 61.96 62.16 | 59.64 62.04 |

Table 5: Performance comparison of CT3 using LLAMA-3-8B-INSTRUCT with baseline (w/o CT3), CT3 using a mixed-domain database, and CT3 with an in-domain database across various tasks and metrics. This experiment uses the smaller knowledge base (§4.5) for efficiency. * for GSM8K and MathQA indicates a subset of 200 randomly sampled test samples.

CT3, a system that enables the application of TTT in client devices that might have resource constraints. The results of the CT3 prototype using supervisory signals from a knowledge base are a call to action to investigate further improvements to learning at test-time methods. Future work should explore more challenging scenarios and datasets to better understand the potential and limitations of collaborative test-time training. The current version of CT3 works on text queries. With increased sophistication, future versions of CT3 must handle more complex tasks and multimodal queries.

537

538

539

540

541

542

543

544

545

546

547

549

572

588

590

593

594

596

597

Limitations

Although our CT3 prototype produces compelling 551 results and demonstrates how to alleviate the knowledge base storage burden at client devices, it also 552 presents limitations. For instance, more research 553 is needed to design methods to determine the right data mix at the knowledge base. In real-world appli-555 cations, the effectiveness of CT3 is likely more pronounced with a more extensive and diverse knowl-557 edge base. A larger knowledge base would provide a broader range of samples, potentially improv-559 ing the relevance and quality of the data retrieved for TTT, thereby enhancing the model's perfor-561 mance even further. Our experimental results have demonstrated the feasibility and potential benefits of CT3. In the current version of CT3's prototype, 564 the server can recover the content of the user's prompt. A real-world solution should incorporate privacy mechanisms to protect the user. In addition to encrypting the query for its transmission, e.g., utilizing Secure Sockets Layer (SSL), many open 569 research challenges exist to increase the privacy and handling of the user's content on the server. 571

Ethical Considerations

Test-time training (TTT) techniques promise im-573 provements in model performance, making them 574 more accurate at the cost of more computation. However, they alone do not solve existing chal-576 lenges in large foundation models and their smaller counterparts. Our research explores systems and 578 techniques to enable running TTT in client devices 579 with resource constraints. However, applying our 580 system and techniques to real-world applications must include additional safeguards to prevent hallu-582 cinations or intentional misinformation that could 583 affect the well-being of users of the system. The 584 research community must continue investigating solutions to address these and other open challenges 586 in popular language models.

References

- Ekin Akyürek, Mehul Damani, Linlu Qiu, Han Guo, Yoon Kim, and Jacob Andreas. 2024. The surprising effectiveness of test-time training for abstract reasoning. *arXiv preprint arXiv:2411.07279*.
- Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. MathQA: Towards interpretable math word problem solving with operation-based formalisms. In *Proceedings of the 2019 Conference*

of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics. 598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

- ARC-AGI. 2025. ARC Prize arcprize.org. https: //arcprize.org. [Accessed 13-03-2025].
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1 others. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Ali Behrouz, Peilin Zhong, and Vahab Mirrokni. 2024. Titans: Learning to memorize at test time. *arXiv* preprint arXiv:2501.00663.
- Loubna Ben Allal, Niklas Muennighoff, Logesh Kumar Umapathi, Ben Lipkin, and Leandro von Werra. 2022. A framework for the evaluation of code generation models. https://github.com/bigcode-project/ bigcode-evaluation-harness.
- Bradley Brown, Jordan Juravsky, Ryan Saul Ehrlich, Ronald Clark, Quoc V Le, Christopher Re, and Azalia Mirhoseini. 2025. Large language monkeys: Scaling inference compute with repeated sampling.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. Evaluating large language models trained on code. *Preprint*, arXiv:2107.03374.
- François Chollet. 2019. On the measure of intelligence. *CoRR*, abs/1911.01547.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,

Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,

Akhil Mathur, Alan Schelten, Amy Yang, Angela

Fan, and 1 others. 2024. The llama 3 herd of models.

Leo Gao, Stella Biderman, Sid Black, Laurence Gold-

ing, Travis Hoppe, Charles Foster, Jason Phang,

Horace He, Anish Thite, Noa Nabeshima, Shawn

Presser, and Connor Leahy. 2020. The Pile: An 800gb dataset of diverse text for language modeling.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Bider-

man, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h,

Haonan Li, Kyle McDonell, Niklas Muennighoff,

Chris Ociepa, Jason Phang, Laria Reynolds, Hailey

Schoelkopf, Aviya Skowron, Lintang Sutawika, and

5 others. 2023. A framework for few-shot language

Moritz Hardt and Yu Sun. 2024. Test-time training on

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan

nearest neighbors for large language models. In In-

ternational Conference on Learning Representations.

Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and

Weizhu Chen. 2022. LoRA: Low-rank adaptation of

large language models. In International Conference

Jonas Hübotter, Sascha Bongni, Ido Hakimi, and Andreas Krause. 2024. Efficiently learning at test-

time: Active fine-tuning of llms. arXiv preprint

IronbarArc24. 2024. arc24 — ironbar.github.io. https:

Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-

sch, Chris Bamford, Devendra Singh Chaplot, Diego

de las Casas, Florian Bressand, Gianna Lengyel, Guil-

laume Lample, Lucile Saulnier, and 1 others. 2023.

Wolf, Dzmitry Bahdanau, Leandro von Werra, and Harm de Vries. 2022. The stack: 3 tb of permissively

Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. 2016. MAWPS: A math word problem repository. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Hu-

man Language Technologies, pages 1152–1157, San Diego, California. Association for Computational

Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. 2024. Orca-math: Unlocking the potential of slms in grade school math. Preprint,

Mistral 7b. arXiv preprint arXiv:2310.06825. Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas

licensed source code. Preprint.

//ironbar.github.io/arc24/. [Accessed 13-03-

arXiv preprint arXiv:2407.21783.

arXiv preprint arXiv:2101.00027.

model evaluation.

on Learning Representations.

arXiv:2410.08020.

2025].

Linguistics.

arXiv:2402.14830.

- 657
- 661

- 670 671
- 672
- 673
- 675
- 682 683
- 684

695

- 703

704

Sebastián Ramírez. 2025. FastAPI fastapi.tiangolo.com. https://fastapi. tiangolo.com/. [Accessed 13-03-2025].

710

711

712

713

714

715

717

718

719

722

724

725

726

727

728

729

730

731

732

733

734

735

736

737

740

741

742

743

744

745

746

747

- Siva Reddy, Dangi Chen, and Christopher D. Manning. 2019. CoQA: A conversational question answering challenge. Transactions of the Association for Computational Linguistics, 7:249–266.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics.
- Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2025. Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning. In The Thirteenth International Conference on Learning Representations.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, and 1 others. 2024. Qwen2 technical report. arXiv preprint arXiv:2407.10671.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. arXiv preprint arXiv:2309.12284.
- Hu Zhiqiang, Lan Yihuai, Wang Lei, Xu Wanyu, Lim EePeng, Lee Roy Ka-Wei, Bing Lidong, and Poria Soujanya. 2023. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. arXiv preprint arXiv:2304.01933.

| Dataset | Domain | Link |
|-------------------|-----------------------|------|
| CoQA Training set | Reading comprehension | link |
| MetaMath | Math | link |
| Orca-Math | Math | link |
| Math 50K | Math | link |
| The Stack Python | Coding | link |
| MBPP training set | Coding | link |
| | | |

Table 6: The details of the datasets in the knowledge base.

| Task | Domain | Eval Framework | Link |
|-----------|---------------|----------------------|------|
| CoQA | Reading Comp. | lm-eval-harness | link |
| GSM8K | Math | lm-eval-harness | link |
| MathQA | Math | llm-adapters | link |
| MAWPS | Math | llm-adapters | link |
| MBPP | Coding | bigcode-eval-harness | link |
| HumanEval | Coding | bigcode-eval-harness | link |

Table 7: Various evaluation tasks.

A Datasets

748

749

750

751

758

762

770

773

The details of the knowledge base datasets and the evaluation tasks used in the experiments are presented in Tables 6 and 7. We evaluated several language models on a diverse set of tasks. We categorize the tasks into:

| •] | Reading | Comprel | hension | Tasks |
|-----|---------|---------|---------|-------|
|-----|---------|---------|---------|-------|

- CoQA (Conversational Question Answering, Reddy et al. (2019))
- Mathematical Reasoning Tasks
 - GSM8K (Grade School Math 8K, Cobbe et al. (2021))
 - MathQA (Math Word Problem Question Answering, Amini et al. (2019))
 - MAWPS (MAth Word ProblemS, Koncel-Kedziorski et al. (2016))
- Coding Tasks
 - MBPP (Mostly Basic Python Problems, Austin et al. (2021))
 - HumanEval (Hand-Written Coding Evaluation Set, Chen et al. (2021))

B Additional Results with LSM

Table 8 shows that incorporating CT3 and LSM consistently enhances performance. These results underscore the effectiveness of LSM in boosting both performance and efficiency.

C Comparing SIFT and Nearest Neighbors Performance in the CT3 Setup

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

SIFT adds an additional step when retrieving relevant samples from the knowledge base, requiring an initial retrieval utilizing the nearest neighbors (NN) method. In our exploration to make the CT3 pipeline more efficient, we explore the impact of removing SIFT and using directly the results from NN. Table 9 presents a comparative analysis of CT3 with and without the SIFT step across various evaluation tasks. The results indicate that while SIFT provides a slight performance edge in some cases, the difference is not substantial. For instance, in the CoQA task, CT3 with SIFT achieves an Exact Match (EM) score of 68.85 and an F1 score of 81.13, compared to 67.43 EM and 80.73 F1 without SIFT. This marginal improvement suggests that SIFT's additional filtering step refines the sample selection process, leading to slightly better performance. In conclusion, while SIFT provides a slight performance boost, NN alone offers a simpler and more computationally efficient approach for sample retrieval in the CT3 system.

D Case Studies

Tables 10, 11, and 12 compare different methods to solve various problems, highlighting the effectiveness of CT3 in providing clear, concise and accurate solutions to user queries on client devices.

| Method | Budget B | Reuse Threshold τ | Avg. Score | Impr. | Evaluation Speedup |
|--------------------------|----------|------------------------|------------|-------|-----------------------|
| LLAMA-3-8B-INSTRUCT | / | / | 51.48 | / | / |
| + CT3 | / | / | 63.66 | +24% | 1.00× |
| | 2 | | 63.99 | +24% | 1.21× |
| | 4 | | 64.17 | +25% | 1.36× |
| CT2 I ISM | 8 | 0.5 | 64.75 | +26% | 1.60× |
| +C13+LSM | 2 | 0.4 | 64.13 | +25% | 1.59× |
| | 4 | 0.4 | 63.26 | +23% | 1.98× |
| | 8 | 0.4 | 63.68 | +24% | 2.24× |
| LLAMA-3.1-8B-INSTRUCT | / | / | 45.45 | / | / |
| + CT3 | | / | 64.48 | +42% | 1.00× |
| | 2 2 | 0.5 | 65.29 | +44% | 1.16× |
| | 4 | 0.5 | 64.71 | +42% | 1.29× |
| CT2 I ISM | 8 | 0.5 | 65.53 | +44% | 1.46× |
| + C13 + L3M | 2 | 0.4 | 63.85 | +41% | 1.46× |
| | 4 | 0.4 | 64.10 | +41% | 1.72× |
| | 8 | 0.4 | 63.12 | +39% | 1.89 × |
| MISTRAL-7B-INSTRUCT-V0.3 | / | / | 47.11 | / | / |
| + CT3 | / | / | 53.48 | +14% | 1.00× |
| | 2 | 0.5 | 51.67 | +10% | 1.20× |
| | 4 | 0.5 | 51.52 | +9% | 1.42× |
| $\pm CT3 \pm ISM$ | 8 | 0.5 | 50.86 | +8% | 1.61× |
| + C15 + L5M | 2 | 0.4 | 51.92 | +10% | 1.60× |
| | 4 | 0.4 | 51.05 | +8% | 2.08× |
| | 8 | 0.4 | 50.66 | +8% | 2.28 × |
| QWEN2.5-3B-INSTRUCT | / | / | 51.81 | / | / |
| + CT3 | / | / | 57.46 | +11% | 1.00× |
| | 2 | 0.5 | 57.16 | +10% | 1.21× |
| | 4 | 0.5 | 57.76 | +11% | 1.39× |
| + CT3 + I SM | 8 | 0.5 | 57.70 | +11% | 1.65× |
| | 2 | 0.4 | 56.85 | +10% | 1.57× |
| | 4 | 0.4 | 55.98 | +8% | 2.08× |
| | 8 | 0.4 | 55.80 | +8% | 2.32× |

Table 8: Performance comparison of LSM with varying budgets and reuse thresholds.

| Method | CoQA EM | CoQA F1 | GSM8K [*] EM (0-shot) | MathQA [*] Acc. | MAWPS Acc. | MBPP Pass@1 | HumanEval Pass@1 | Avg. | Impr. |
|---------------------|------------|------------|-----------------------------------|-----------------------------|---------------|----------------|---------------------|-------|-------|
| QWEN2.5-3B-INSTRUCT | 54.70 | 71.22 | 61.50 | 20.00 | 87.82 | 41.20 | 26.22 | 51.81 | / |
| + CT3 w/ SIFT | 68.85 | 81.13 | 66.00 | 34.00 | 86.13 | 46.20 | 27.44 | 58.54 | +13% |
| + CT3 w/o SIFT | 67.43 | 80.73 | 68.50 | 30.00 | 86.13 | 46.00 | 28.66 | 58.21 | +12% |

Table 9: Performance comparison of CT3 w/ SIFT and CT3 w/o SIFT (QWEN2.5-3B-INSTRUCT). The ablation studies use the smaller knowledge base (§4.5) for efficiency.

| Prompt | Question: Charlie wants to sell beeswax candles. For every pound of beeswax, he can make 10 tapered candles. One pound of beeswax and the wicks cost \$10.00 in supplies. If he sells each candle for \$2.00 each, what is his net profit if he makes and sells 20 candles? Answer: |
|--------------------|---|
| Reference Solution | 1 pound of beeswax and the wick make 10 candles and he wants to make 20 candles so he needs 20/10 = «20/10=2»2 pounds of beeswax 1 pound of beeswax costs \$10.00 and he needs 2 pounds so his supplies will cost 10*2 = \$«10*2=20.00»20.00 He wants to sell each candle for \$2.00 each and he's selling 20 candles so that's 2*20 = \$«2*20=40.00»40.00 If he makes \$40.00 selling candles and his supplies cost \$20.00 then his net profit is 40-20 = \$«40-20=20.00»20.00 ##### 20 |
| w/o CT3 | He spends $10.00 / \text{candle} * 20 \text{ candles} = 200.00 \text{ on beeswax and wicks.}$ He earns $2 * 2.00 = 4.00$ per candle. So he earns 20 candles $4.00 = 80.00$ total from sales. That means his net profit is $80.00 - 200.00 = -120.00$ Therefore, the answer is $-120. \times$ |
| CT3 | He uses 2 pounds of beeswax because $20 / 10 = 2$ The supplies cost \$20 because $2 \ge 10 = 20$ He earns \$40 from selling the candles because $20 \ge 2 = 40$ His net profit is \$20 because $40 - 20 = 20$ Therefore, the answer is $20 \le \sqrt{20}$ |

Table 10: A case comparing different methods for solving a profit calculation problem with QWEN2.5-3B-INSTRUCT. The table illustrates the effectiveness of CT3 in accurately solving the problem, as opposed to the incorrect solution provided by the baseline.

| Prompt | Question: Meredith is a freelance blogger who writes about health topics and submits to clients each day as her permanent job. A blog article takes an average of 4 hours to research and write about. Last week, she wrote 5 articles on Monday and 2/5 times more articles on Tuesday than on Monday. On Wednesday, she wrote twice the number of articles she wrote on Tuesday. Calculate the total number of hours she spent writing articles in the three days. Answer: |
|--------------------|---|
| Reference Solution | If she wrote five articles on Monday, then on Tuesday she wrote $2/5*5 = 2$ more articles. The total number of articles she wrote on Tuesday is $5+2 = (+5+2)^{-7}$ On Wednesday, the number of articles was twice the number she wrote on Tuesday, so it was $2*7 = (+2)^{-1}4*14$ articles. Over the three days, she wrote $5+7+14 = (+5+7)^{-1}4=26*26$ articles. If each article takes her 4 hours to research and write about, she spent $4*26 = (+4)^{-1}26=104*104$ hours on all the articles. #### 104 |
| w/o CT3 | To determine the total number of hours Meredith spent writing articles over the three days, we need to follow these steps: 1. **Calculate the number of articles written on Tuesday:** - On Monday, Meredith wrote 5 articles. - On Tuesday, she wrote \\(\\frac25 \\) times more articles than on Mon- day. - First, find \\(\\frac25 \\) of the articles written on Monday: \\[\\frac25 \\times 5 = 2\\] - Therefore, the total number of articles written on Tuesday is: \\[5 + 2 = 7\\] 2. **Calculate the number of articles written on Wednesday:** - On Wednesday, she wrote twice the number of articles she wrote on Tues- day. - The number of articles written on Wednesday is: \\[2 \\times 7 = 14\\] 3. **Calculate the total number of articles written over the three days:** - Articles on Monday: 5 - Articles on Tuesday: 7 - Articles × |
| CT3 | To find out how many articles Meredith wrote on Tuesday, we need to calculate 2 / 5 times more than Monday's articles. Since she wrote 5 articles on Monday, 2 / 5 * 5 = 2 more articles than Monday. This means she wrote $5 + 2 = 7$ articles on Tuesday. On Wednesday, she wrote twice the number of articles she wrote on Tuesday, which is $2 * 7 = 14$ articles. In total, she wrote $5 + 7 + 14 = 26$ articles over the three days. Since it takes her 4 hours to research and write about one article, she spent $4 * 26 = 104$ hours writing articles in the three days. Therefore, the answer is 104 hours. |

Table 11: A case comparing different methods for solving a time calculation problem with QWEN2.5-3B-INSTRUCT. The table demonstrates that the baseline (w/o CT3) unnecessarily complicates the problem, while CT3 effectively maximizes the model's potential for the current query. CT3 utilizes CT3 to provide a clear, concise, and accurate solution to the user's query at the client device.

| Prompt | Question: A wooden bridge can carry no more than 5000 pounds. A delivery truck filled with identical boxes, each weighing 15 pounds, will pass over the bridge. The combined weight of the driver and the empty truck is 3755 pounds. What is the maximum number of boxes which can be loaded onto the truck while not exceeding the bridge's weight limit? Answer: |
|--------------------|--|
| Reference Solution | The boxes can weigh up to 5000 pounds - 3755 pounds = $(5000-3755=1245)$ pounds in total. There can be $1245/15 = (1245/15=83)$ boxes loaded onto the truck without exceeding the bridge's weight limit. #### 83 |
| w/o CT3 CT3 | First find the total weight capacity of the truck by adding the weight limit of the bridge to the weight of the driver and empty truck: 5000 pounds + 3755 pounds = 8755 pounds. Then divide that number by the weight of each box to find the maximum number of boxes: 8755 pounds / 15 pounds / box = 583.666 boxes. Since we can't have a fraction of a box, we call this 583 boxes. Therefore, the answer is 583. X First, we need to find out how much weight the boxes can add to the truck without exceeding the bridge's weight limit. We do this by subtracting the weight of the driver and the empty truck from the bridge's weight limit: 5000 - 3755 = 1245 pounds. Next, we divide the remaining weight capacity by the weight of a single box to find out how many boxes can be added: $1245 / 15 = 83$ boxes. |

Table 12: A case for comparing different methods of solving a weight calculation problem using QWEN2.5-3B-INSTRUCT. It demonstrates the model's misunderstanding of the problem when CT3 is not used, specifically by incorrectly adding the bridge's weight limit to the truck's weight.