# On-Policy Distillation with Contrastive Learning for Reasoning

**Anonymous Author(s)**
Affiliation
Address
`email`

## Abstract

On-policy distillation trains a student on its own outputs with teacher feedback to avoid the train-inference distribution mismatch present in standard distillation. However, it typically ignores the ground-truth labels and reference solutions available in reasoning datasets during the post-training. We introduce a verification-driven contrastive learning that explicitly exploits these reference labels. By evaluating student responses against ground-truth, we strategically pair suboptimal generations with expert reference solutions to form robust contrastive learning pairs. Combined with token-level teacher distillation on verified correct outputs, this prompts the student to simultaneously refine its successful reasoning and correct its failures within a unified on-policy loop. We evaluate our method, **CORE** (Contrastive On-policy Reasoning Distillation), across mathematical reasoning, summarization, translation, and instruction tuning tasks. CORE converges faster and achieves state-of-the-art results compared to both GRPO and standard on-policy distillation.

## 1 Introduction

Large-scale models [1] with billions of parameters have demonstrated remarkable capabilities but are limited, as substantial computational costs, in deployment on resource-constrained edge devices, such as satellites [2, 3] and autonomous vehicles [4, 5]. Knowledge Distillation (KD) [6, 7] transfers capabilities from an advanced model (teacher model) with billions of parameters to a small model (student model) that satisfies the strict efficiency requirements of such environments. However, the inherent architectural heterogeneity between the two models often induces a significant distribution shift [8] and causes a mismatch that hinders the student from effectively adapt the teacher's representations.

To mitigate this discrepancy, Generalized On-policy Distillation [9] addresses this issue by training the student on samples drawn from its current distribution rather than a fixed training set by leveraging Group Relative Policy Optimization (GRPO). However, GKD treats all student-generated samples uniformly and does not exploit the correctness labels or reference chain-of-thought solutions that are readily available in reasoning benchmarks.

We observe that reasoning benchmarks inherently provide three complementary supervision signals, including dense token-level teacher distributions, ground-truth, and reference chain-of-thought (CoT) solutions, whereas no existing method leverages them jointly. To address this limitation, we introduce **CORE** (Contrastive On-policy Reasoning Distillation). Specifically, an answer-aware distillation objective conditions the teacher on the reference CoT and the ground-truth label, which allows it to provide dense token-level supervision over the full set of student-generated responses. Simultaneously, a contrastive preference optimization objective pairs each incorrectly answered response with a verified correct counterpart, thereby explicitly penalizing erroneous reasoning trajectories by Direct Preference Optimization (DPO). This adaptive combination produces a natural curriculum in which

the model starts with contrastive learning to reduce errors, then shifts to distillation to sharpen reasoning.

We evaluate **CORE** on complex reasoning benchmarks in different size of model Extensive experiments demonstrate that CORE significantly outperforms state-of-the-art distillation methods, achieving notable improvements in accuracy and convergence speed.

# 2 Related Work

## 2.1 Knowledge Distillation for Language Models

Traditional knowledge distillation [10] trains student models to match teacher output distributions. For autoregressive language models, supervised KD [11] and sequence-level KD [12] are widely used. However, these off-policy approaches suffer from distribution mismatch between training (teacher-generated or ground-truth sequences) and inference (student-generated sequences).

**On-policy distillation** [9] addresses this by training students on their own generated outputs, using teacher logits as labels. GKD (Generalized Knowledge Distillation) demonstrates strong improvements on summarization and translation tasks. Our work extends on-policy distillation to reasoning tasks by incorporating verification signals and contrastive objectives.

## 2.2 Learning from Verification Feedback

Recent work on mathematical reasoning leverages verification to improve model training. GRPO [? ] and similar RL-based approaches optimize for verified correctness using policy gradient methods. V-STaR [? ] iteratively generates verified solutions for self-improvement.

However, these methods rely solely on sparse binary rewards (correct/incorrect) and do not leverage dense token-level teacher guidance. TCD bridges this gap by combining verification with token-level distillation.

## 2.3 Contrastive Learning for Language Models

DPO [13] introduced preference-based contrastive learning for alignment, training models to prefer chosen responses over rejected ones. SimPO [? ] and other variants explore different contrastive formulations.

Our work adapts contrastive learning to the distillation setting: we use verified correct responses (from reference data or teacher generation) as "chosen" and student-generated incorrect responses as "rejected", enabling the student to learn from its mistakes with teacher guidance.

# 3 Method: Token-level Contrastive Distillation

## 3.1 Problem Setup

We consider distilling a large teacher model $p_T$ into a smaller student model $p_S^\theta$ for mathematical reasoning tasks. Given a dataset of questions $\{x_i\}$ and ground-truth answers $\{a_i\}$, our goal is to train the student to generate correct reasoning traces $y$ (chain-of-thought) that lead to the correct answer.

## 3.2 On-policy Generation with Verification

For each question $x$, we generate $K$ reasoning traces from the current student policy:

$$\{y^{(1)}, \ldots, y^{(K)}\} \sim p_S^\theta(\cdot|x)$$

Each trace is verified by extracting the final answer and comparing with ground-truth:

$$\text{verify}(y^{(k)}, a) \in \{\text{correct}, \text{incorrect}\}$$

This yields two sets per question: correct traces $\mathcal{C}_x$ and incorrect traces $\mathcal{I}_x$.

### 3.3 Dual Learning Objectives

**Distillation on Correct Traces.** For verified correct reasoning, we perform token-level knowledge distillation. Given teacher logits $z_T(y_{<t}, x)$ and student logits $z_S^\theta(y_{<t}, x)$ at position $t$, we minimize:

$$\mathcal{L}_{\text{distill}}(\theta) = \mathbb{E}_{x, y \in \mathcal{C}_x} \left[ \frac{1}{|y|} \sum_{t=1}^{|y|} D_{\text{KL}} \left( \sigma(z_T/\tau) \| \sigma(z_S^\theta/\tau) \right) \right]$$

where $\sigma$ is softmax, $\tau$ is temperature, and we use top-$k$ truncation of teacher logits for memory efficiency.

**Contrastive Learning on Errors.** For incorrect traces, we apply DPO-style contrastive loss. The "chosen" response $y_c$ is either (1) reference chain-of-thought if available, (2) student's own correct trace if any, or (3) teacher-generated trace. The "rejected" response $y_r$ is the student's incorrect trace:

$$\mathcal{L}_{\text{contrast}}(\theta) = -\mathbb{E}_{x, y_r \in \mathcal{I}_x} \left[ \log \sigma \left( \beta \left( \log \frac{p_S^\theta(y_c|x)}{p_{\text{ref}}(y_c|x)} - \log \frac{p_S^\theta(y_r|x)}{p_{\text{ref}}(y_r|x)} \right) \right) \right]$$

where $p_{\text{ref}}$ is a reference model (frozen copy of student) and $\beta$ controls the strength of the KL penalty.

### 3.4 Combined Training Objective

The final loss combines both objectives:

$$\mathcal{L}(\theta) = \alpha \cdot \mathcal{L}_{\text{distill}}(\theta) + (1 - \alpha) \cdot \mathcal{L}_{\text{contrast}}(\theta)$$

where $\alpha$ controls the relative weight between distillation and contrastive learning.

### 3.5 Memory-Efficient Implementation

To handle large vocabulary sizes (151,936 for Qwen), we employ:

- **Top-K teacher logits**: Store only top-512 logits from teacher, reducing memory by 300×
- **Token-by-token processing**: Compute distillation loss incrementally to avoid materializing full distributions
- **vLLM for generation**: Use vLLM engines with continuous batching for efficient on-policy sampling

## 4 Experiments

### 4.1 Experimental Setup

**Models.** We use Qwen3-1.7B as the student model and Qwen3-8B as the teacher model, both pretrained on diverse text corpora.

**Dataset.** We evaluate on GSM8K [14], a dataset of grade-school math word problems requiring multi-step reasoning. We use 4-shot chain-of-thought prompting with calculator tool access.

**Training.**

- Batch size: 32 (rollout), 32 (training)
- Episodes: 3 (full dataset passes with on-policy regeneration)
- Learning rate: $1 \times 10^{-6}$ with 5% warmup
- Generation: 8 samples per prompt, temperature 0.7
- Distillation: $\alpha = 5.0$ (distill weight), top-K=512 teacher logits
- Contrastive: $\beta = 0.1$ (DPO beta)

**Infrastructure.** We use 4× H100 GPUs with Ray for distributed training:

- GPU 0: Student + Reference model (colocated)
- GPU 1: Teacher model
- GPU 2-3: vLLM engines (2×) for on-policy generation

## 4.2 Main Results

**TODO: Add results table comparing:**

- Supervised fine-tuning baseline
- Distillation only (no contrastive)
- Contrastive only (no distillation)
- TCD (full method)
- Teacher performance (upper bound)

## 4.3 Ablation Studies

**Effect of Episode Number.** We ablate the number of on-policy episodes (1, 2, 3). Each episode regenerates data with the updated student, improving the quality of both correct traces (for distillation) and incorrect traces (for contrastive learning).

**Distillation Weight $\alpha$.** We vary $\alpha \in \{0.1, 1.0, 5.0, 10.0\}$ to study the trade-off between learning from verified correct reasoning vs. learning from mistakes.

**Top-K Truncation.** We compare full vocabulary (151,936) vs. top-K logits with $K \in \{128, 512, 2048\}$, measuring both accuracy and memory usage.

**Teacher Fallback Strategy.** When student generates 0/8 correct samples, we compare:

1. Skip contrastive loss
2. Use teacher.generate() to create reference
3. Use teacher logits for autoregressive sampling

## 4.4 Analysis

**Correct/Incorrect Ratio Over Training.** We track the fraction of correct vs. incorrect student-generated traces across episodes, showing how on-policy distribution improves.

**Memory vs. Accuracy Trade-offs.** We analyze GPU memory consumption for different components (student, teacher, vLLM, reference model) under various configurations (colocated vs. separate, DeepSpeed ZeRO stages).

## 5 Conclusion

We presented Token-level Contrastive Distillation (TCD), a framework that combines on-policy generation, verification-based learning, and dense teacher guidance for distilling reasoning capabilities into smaller language models. By unifying distillation on correct traces with contrastive learning on errors, TCD effectively leverages both sparse verification signals and rich token-level teacher feedback.

Our experiments on GSM8K demonstrate that TCD enables a 1.7B student model to achieve strong mathematical reasoning performance when distilled from an 8B teacher. The framework is efficient and scalable, using vLLM for on-policy generation and memory-optimized techniques for handling large vocabularies.

**Limitations and Future Work.**

- **Scalability**: Current experiments use 4 GPUs; scaling to multi-node setups requires careful optimization of communication backends (NCCL vs. Gloo)
- **Generalization**: We focus on mathematical reasoning (GSM8K); extending to other reasoning tasks (code, commonsense) is important
- **Teacher quality**: Our approach assumes access to a capable teacher; exploring self-improvement scenarios is promising

TCD opens avenues for efficient reasoning model deployment by making strong reasoning capabilities accessible in compact models suitable for resource-constrained environments.

# References

[1] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.

[2] Xiao Jia, Di Zhou, Min Sheng, Yan Shi, Sijing Ji, and Jiandong Li. Satellite computing network construction: Optimal computing node deployment in multi-layer leo mega-constellations. *IEEE Transactions on Communications*, 74:1747–1761, 2025.

[3] Yihong Tao, Bo Lei, Haoyang Shi, Jingkai Chen, and Xing Zhang. Adaptive multi-layer deployment for a digital-twin-empowered satellite-terrestrial integrated network. *Frontiers of Information Technology & Electronic Engineering*, 26(2):246–259, 2025.

[4] Jingyuan Zhao, Yuyan Wu, Rui Deng, Susu Xu, Jinpeng Gao, and Andrew Burke. A survey of autonomous driving from a deep learning perspective. *ACM Computing Surveys*, 57(10):1–60, 2025.

[5] Tuo Feng, Wenguan Wang, and Yi Yang. A survey of world models for autonomous driving. *arXiv preprint arXiv:2501.11260*, 2025.

[6] Amir M Mansourian, Rozhan Ahmadi, Masoud Ghafouri, Amir Mohammad Babaei, Elaheh Badali Golezani, Zeynab Yasamani Ghamchi, Vida Ramezanian, Alireza Taherian, Kimia Dinashi, Amirali Miri, et al. A comprehensive survey on knowledge distillation. *arXiv preprint arXiv:2503.12067*, 2025.

[7] Luyang Fang, Xiaowei Yu, Jiazhang Cai, Yongkai Chen, Shushan Wu, Zhengliang Liu, Zhenyuan Yang, Haoran Lu, Xilin Gong, Yufang Liu, et al. Knowledge distillation and dataset distillation of large language models: Emerging trends, challenges, and future directions. *Artificial Intelligence Review*, 59(1):17, 2026.

[8] Siyan Zhao, Zhihui Xie, Mengchen Liu, Jing Huang, Guan Pang, Feiyu Chen, and Aditya Grover. Self-distilled reasoner: On-policy self-distillation for large language models. *arXiv preprint arXiv:2601.18734*, 2026.

[9] Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-generated mistakes. In *International Conference on Learning Representations (ICLR)*, 2024.

[10] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[11] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.

[12] Yoon Kim and Alexander M Rush. Sequence-level knowledge distillation. *arXiv preprint arXiv:1606.07947*, 2016.

[13] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2023.

[14] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

# A  Implementation Details

## A.1  Distributed Training Architecture

Our implementation uses Ray [? ] for distributed actor management and vLLM [? ] for efficient LLM serving. The architecture consists of:

- **VtDTrainerRay**: Main coordinator running on the driver process
- **VtDStudentActor**: Manages student model with DeepSpeed (GPU 0)
- **VtDReferenceActor**: Frozen copy of initial student (GPU 0, colocated)
- **VtDTeacherActor**: Teacher model (GPU 1)
- **LLMRayActor × 2**: vLLM engines for on-policy generation (GPU 2-3)

Model colocation uses Ray placement groups with fractional GPU allocation (num_gpus=0.2) to share GPU 0 between student and reference models.

## A.2  Memory Optimization Techniques

**Top-K Teacher Logits.**  Instead of storing full vocabulary logits ($151{,}936 \times 4$ bytes $= 608$ KB per position), we keep only top-512 values and indices:

```
topk_vals, topk_ids = logits.topk(k=512, dim=-1)
# Store: 512 × 4 bytes (vals) + 512 × 4 bytes (ids) = 4 KB
```

This achieves 150× memory reduction with minimal accuracy loss.

**Token-by-token Distillation Loss.**  We compute KL divergence incrementally to avoid materializing large tensors:

```
for t in range(seq_len):
    teacher_probs_t = F.softmax(teacher_vals[t], dim=-1)
    student_logprobs_t = F.log_softmax(student_logits[t], dim=-1)
    student_logprobs_topk = student_logprobs_t.gather(-1, teacher_ids[t])
    kl_t = -(teacher_probs_t * student_logprobs_topk).sum(-1)
```

## A.3  Communication Backend Selection

For non-colocated setups, vLLM workers require collective communication for weight updates. We found:

- **NCCL**: Fast but fails when actors are in separate Ray placement groups due to CUDA_VISIBLE_DEVICES isolation
- **Gloo**: CPU-based fallback that works across placement groups, used when tensor_parallel_size=1

# B  Additional Experimental Results

## B.1  Hyperparameter Sensitivity

**TODO: Add figures/tables for:**

- Learning rate sweep
- Temperature sweep for generation
- Beta sweep for contrastive loss
- Number of samples per prompt (K)

## B.2 Evaluation Protocol Details

For GSM8K evaluation, we:

1. Generate greedy responses (temperature=0, n=1)
2. Extract final numerical answer using regex
3. Compare with ground-truth using math_equal() for numerical equivalence
4. Report accuracy = correct / total

## B.3 Computational Cost

Training TCD on 10K GSM8K samples for 3 episodes takes approximately:

- Wall-clock time: **TODO** hours on 4× H100
- GPU hours: **TODO**
- On-policy generation: **TODO**% of total time
- Teacher logits computation: **TODO**% of total time
- Student training: **TODO**% of total time

# C  Reproducibility

Code will be released at `https://github.com/TODO`. Full training configurations are provided in YAML format.