
Uncertainty-Calibrated On-Policy Distillation for Large Language Models

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, a fundamental underexplored challenge is that the teacher itself is imperfect; consequently, naively imitating erroneous teacher outputs propagates mistakes into the student. Prior work such as Self-Distilled Reasoner [1] mitigates this by injecting ground-truth chain-of-thought solutions into the teacher’s system prompt to guarantee the teacher with correct reasoning. But this approach depends on labeled data that is unavailable in many practical settings. We introduce **UOPD** (Uncertainty-calibrated On-Policy Distillation), a framework that addresses teacher errors without ground-truth supervision. Teacher reliability is estimated offline by computing its **uncertainty**, and the distillation objective is reweighted inversely proportional to the teacher’s uncertainty with no additional training cost. We further release **SURE-Math**, a dataset annotated with teacher semantic entropy scores. We evaluate UOPD across mathematical reasoning, summarization, translation, and instruction tuning tasks. UOPD converges faster and achieves state-of-the-art results compared to both GRPO and standard on-policy distillation.

17 **1 introduction**

18 **2 Preliminary**

19 **2.1 On-Policy Knowledge Distillation**

20 Standard knowledge distillation [2] trains a student model $p_S(\cdot; \theta_S)$ to minimize the KL divergence
21 from a frozen teacher p_T on a fixed dataset \mathcal{D} :

$$\mathcal{L}_{\text{KD}}(\theta_S) = \mathbb{E}_{x \sim \mathcal{D}} [\text{KL}(p_T(\cdot | x) \| p_S(\cdot | x; \theta_S))]. \quad (1)$$

22 This offline objective suffers from a train-inference distribution mismatch: during training the
23 student observes token-level distributions conditioned on dataset prefixes, whereas at inference it
24 must generate from its own previously predicted tokens. On-policy distillation [3] resolves this by
25 generating training sequences from the student itself:

$$\mathcal{L}_{\text{OPD}}(\theta_S) = \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim p_S(\cdot | x; \theta_S)} \left[\sum_{t=1}^{|y|} \text{KL}(p_T(\cdot | x, y_{<t}) \| p_S(\cdot | x, y_{<t}; \theta_S)) \right]. \quad (2)$$

26 Since the training distribution now matches the inference distribution, on-policy distillation substantially
27 reduces compounding errors in chain-of-thought reasoning [4].

28 **2.2 Uncertainty in LLMs**

29 Uncertainty in language model outputs can be decomposed into two fundamentally distinct sources
30 [5]. **Aleatoric uncertainty** is irreducible and arises from the inherent ambiguity of natural language,
31 where the same meaning can be expressed in infinitely many surface forms. Consider the prompt
32 “*What is the capital of the UK?*” Given this question, the responses “*London*”, “*London is the capital*
33 *of the UK*”, and “*The UK’s capital is London*” are semantically identical, while their token-level
34 divergence is large. A teacher that consistently produces such paraphrases is not making errors; it
35 is exhibiting natural lexical variation. Penalising this variation at the token level conflates surface
36 diversity with genuine unreliability.

37 **Epistemic uncertainty**, by contrast, reflects the teacher’s *lack of knowledge*. It arises when the
38 teacher generates contradictory answers across samples, indicating that it does not reliably know the
39 correct response to a given prompt. Standard token-level objectives cannot distinguish these two
40 sources, because both manifest as high distributional variance in the output space.

41 Semantic entropy [6] resolves this ambiguity by operating over *meaning* rather than surface tokens.
42 Responses are first clustered into semantic equivalence classes \mathcal{C} , and entropy is computed over the
43 resulting distribution

$$H_{\text{sem}}(x) = - \sum_{c \in \mathcal{C}} p(c | x) \log p(c | x), \quad p(c | x) \propto \sum_{y \in c} p_T(y | x). \quad (3)$$

44 Paraphrases of the same answer collapse into a single class, suppressing aleatoric noise. $H_{\text{sem}}(x)$
45 is therefore a faithful, label-free measure of epistemic uncertainty; it is high only when the teacher
46 genuinely disagrees with itself across semantically distinct answers.

47 **3 Method**

48 We propose UOPD, a two-phase framework for uncertainty-calibrated on-policy distillation. In the
49 first phase, we estimate the teacher’s epistemic uncertainty for each training prompt by computing
50 semantic entropy over multiple teacher rollouts. This computation is performed entirely offline before
51 training begins, producing a per-prompt uncertainty weight that is stored alongside the training data.
52 In the second phase, the student is trained on-policy with token-level distillation from the teacher,
53 where each sample’s contribution to the loss is calibrated by its precomputed uncertainty weight.

54 The key advantage of this design is the separation of uncertainty estimation from training. Because
55 semantic entropy is computed once before the training loop, UOPD introduces zero additional
56 overhead during training while still enabling uncertainty-aware distillation. We describe each
57 component in detail below.

58 **3.1 Offline Semantic Entropy Estimation**

59 Given a training set of prompts $\{x_i\}_{i=1}^M$, we first estimate the teacher’s epistemic uncertainty per
60 prompt. For each prompt x_i , we sample N responses from the teacher $\{y_1^T, \dots, y_N^T\} \sim p_T(\cdot | x_i)$,
61 recording each response’s sequence log-probability $\log p_T(y_n^T | x_i)$ and token count $|y_n^T|$.

62 **Semantic clustering.** As discussed in Section 2, responses such as “*London*”, “*London is the*
63 *capital of the UK*”, and “*The UK’s capital is London*” are semantically identical but produce large
64 token-level divergence. To prevent this aleatoric variation from inflating uncertainty estimates, we
65 extract the final answer from each teacher response and cluster semantically equivalent answers into
66 equivalence classes $\mathcal{C} = \{c_1, \dots, c_J\}$. Two answers are placed in the same cluster only if one can be
67 derived from the other, i.e., they express the same meaning in different surface forms (e.g., $\frac{1}{2}$ and 0.5,
68 or $x = 3$ and 3). Answers that match as exact strings are merged directly. For non-matching pairs, a
69 Qwen-4B judge [7] determines pairwise semantic equivalence, and all judgments are resolved into
70 clusters with a Union-Find algorithm.

71 **Intra-cluster vs. inter-cluster entropy.** This clustering decomposes the total entropy of teacher
72 outputs into two interpretable components. *Intra-cluster entropy* captures the lexical diversity *within*
73 a semantic equivalence class. A cluster containing “0.5”, “ $\frac{1}{2}$ ”, and “the answer is one half” has high
74 intra-cluster entropy, reflecting the teacher’s expressive richness rather than genuine confusion. This
75 variation is aleatoric and should not be penalized. *Inter-cluster entropy*, by contrast, measures the

76 spread of probability mass *across* semantically distinct answer classes. When the teacher assigns
 77 substantial mass to multiple contradictory clusters (e.g., both “3” and “5” for the same prompt),
 78 inter-cluster entropy is high, signaling epistemic uncertainty. Our semantic entropy H_{sem} captures
 79 precisely the inter-cluster component by collapsing all within-cluster variation before computing
 80 entropy.

81 **Probability-weighted semantic entropy.** Each response’s log-probability is first length-normalized
 82 to avoid penalizing longer outputs. The probability mass of semantic class c is then computed as

$$p(c | x_i) = \frac{\sum_{n \in c} \exp(\log p_T(y_n^T | x_i) / |y_n^T|)}{\sum_{n=1}^N \exp(\log p_T(y_n^T | x_i) / |y_n^T|)}, \quad (4)$$

83 and the semantic entropy is

$$H_{\text{sem}}(x_i) = - \sum_{c \in C} p(c | x_i) \log p(c | x_i). \quad (5)$$

84 We normalize by the maximum possible entropy to obtain a value in $[0, 1]$ and define the per-prompt
 85 distillation weight as

$$w_{\text{se}}(x_i) = 1 - \frac{H_{\text{sem}}(x_i)}{\log N}. \quad (6)$$

86 A weight near 1 indicates that the teacher consistently agrees on the same semantic answer (low
 87 epistemic uncertainty), while a weight near 0 indicates contradictory responses across clusters (high
 88 epistemic uncertainty). These weights are precomputed and stored with the training data.

89 3.2 Uncertainty-Calibrated On-Policy Distillation

90 At each training step, the student generates K responses per prompt from its current policy
 91 $\{y^{(1)}, \dots, y^{(K)}\} \sim p_S(\cdot | x; \theta_S)$. For each student-generated response y , the frozen teacher pro-
 92 vides its top- k logit values and corresponding token indices at every position, reducing memory and
 93 communication from the full vocabulary size V to $k \ll V$.

94 The token-level distillation loss at position t is the cross-entropy between the teacher’s and student’s
 95 distributions over the top- k tokens

$$\ell_t = - \sum_{j=1}^k \tilde{p}_T(v_j | x, y_{<t}) \log p_S(v_j | x, y_{<t}; \theta_S), \quad (7)$$

96 where \tilde{p}_T denotes the teacher’s probability renormalized over the top- k tokens and $\{v_1, \dots, v_k\}$ are
 97 the teacher’s top- k token indices.

98 The standard on-policy distillation objective averages ℓ_t uniformly over all response tokens, treating
 99 every prompt equally regardless of teacher reliability. UOPD instead calibrates each sample’s
 100 contribution by its precomputed semantic entropy weight $w_{\text{se}}(x)$. For a prompt x and student-
 101 generated response y , the uncertainty-calibrated loss is

$$\mathcal{L}(x, y; \theta_S) = \frac{w_{\text{se}}(x)}{\sum_t \mathbb{1}[t \in \text{resp}]} \sum_{t \in \text{resp}} \ell_t, \quad (8)$$

102 where the sum runs over response token positions (excluding the prompt). When the teacher is
 103 confident about a prompt ($w_{\text{se}} \approx 1$), the student receives the full distillation signal. When the teacher
 104 is uncertain ($w_{\text{se}} \approx 0$), the loss is suppressed, preventing the propagation of erroneous teacher
 105 guidance.

106 3.3 Overall Training Objective

107 The full training objective averages over all prompts and their K on-policy rollouts

$$\mathcal{L}_{\text{UOPD}}(\theta_S) = \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \frac{1}{K} \sum_{j=1}^K \mathcal{L}(x, y^{(j)}; \theta_S). \quad (9)$$

108 At each training step, the updated student weights are synchronized back to the vLLM generation
 109 engines, ensuring that the next round of rollouts reflects the latest policy. This on-policy loop
 110 continues until convergence.

111 **4 Experiments**

112 **4.1 SURE-Math Dataset**

113 We introduce **SURE-Math** (Semantic-Uncertainty REasoning Math), a mathematical reasoning
114 dataset in which every problem is annotated with the teacher’s precomputed semantic entropy score.

115 **Seed collection.** We aggregate seed problems from seven public sources spanning four difficulty
116 tiers. *Easy*: ScaleQuest-Math [8] (15K problems). *Medium*: NuminaMath-CoT [9] (10K) and MATH
117 [10] (12.5K). *Hard*: Omni-MATH [11] (4.4K) and OlympiadBench [12] (5K). *Competition*: AIME
118 2024–2025, HMMT 2025, and AMO-Bench [13] (\sim 110 problems). In addition, we manually curate
119 1,000 middle-school and high-school mathematics problems with verified ground-truth labels.

120 **Evol-Instruct evolution.** Starting from these seeds, we generate approximately one million new
121 problems using Qwen2.5-72B-Instruct with an Evol-Instruct [14] style pipeline. Six evolution
122 strategies are applied: *harder*, *rewrite*, *algebraize*, *apply*, *compose*, and *competition*. Each seed
123 undergoes 1–5 rounds of evolution depending on its difficulty tier, with competition-level seeds
124 receiving the most rounds to amplify their representation.

125 **Quality filtering.** We apply a three-stage filter. First, we remove multiple-choice questions and
126 problems shorter than 30 characters or longer than 2,000 characters, and deduplicate with a 0.7
127 similarity threshold. Second, Qwen2.5-72B-Instruct judges each remaining problem against five
128 quality criteria: (1) well-formed and complete, (2) unambiguous with exactly one correct answer,
129 (3) requires at least two steps of mathematical reasoning, (4) admits a closed-form answer (number
130 or expression, not a proof or essay), and (5) self-contained. Only problems passing all five criteria
131 are retained. After filtering, 8,000 problems remain.

132 **Semantic entropy annotation.** For each of the 8,000 problems, we sample $N=8$ responses from
133 Qwen3-8B [7] at temperature 0.7, extract the $\backslash\boxed{\cdot}$ answer from each response, and compute
134 the probability-weighted semantic entropy as described in Section 3.1. The resulting per-problem
135 semantic entropy score and the corresponding distillation weight w_{se} are stored alongside each
136 problem.

137 **Data splits.** We reserve 200 of the 1,000 manually collected problems (which carry ground-truth
138 labels) as a held-out test set. The remaining 800 curated problems and all 8,000 evolved problems
139 form the training set.

140 **4.2 Experimental Setup**

141 **Models.** The student model is Qwen3-1.7B and the teacher is Qwen3-8B [7], both instruction-tuned.
142 Semantic clustering for the offline SE computation uses a Qwen3-8B judge for pairwise semantic
143 equivalence.

144 **Training details.** We train with the on-policy distillation objective described in Eq. 8. At each step,
145 the student generates $K=4$ rollouts per prompt. The teacher provides its top-512 logit values at every
146 token position. We use a learning rate of 3×10^{-6} with a cosine schedule and 5% warmup, AdamW
147 with $(\beta_1, \beta_2) = (0.9, 0.95)$, gradient clipping at 1.0, and a global batch size of 128. Training runs for
148 one epoch over 50,000 samples (with replacement from the 8,800 training problems). All models use
149 bfloat16 precision and FlashAttention-2.

150 **Infrastructure.** Training runs on $4 \times$ NVIDIA H100 GPUs. The student model, teacher model, and
151 4 vLLM generation engines are colocated on the same node with DeepSpeed ZeRO Stage 1. At each
152 training step, vLLM engines generate student rollouts, the teacher produces top- k logits for these
153 rollouts, and the student parameters are then updated and synchronized back to the vLLM engines for
154 the next round of generation.

155 **Baselines.** We compare UOPD against the following methods.

- 156 • **Qwen3-1.7B:** The base student model without any distillation (lower bound).

Table 1: Main results on mathematical reasoning benchmarks. We report pass@1 (greedy) and avg@16 accuracy (%). Best student-sized results are **bolded**.

Method	GSM8K		AIME 2025		HMMT 2025		SURE-Math	
	pass@1	avg@16	pass@1	avg@16	pass@1	avg@16	pass@1	avg@16
Qwen3-8B (Teacher)	–	–	–	–	–	–	–	–
Qwen3-1.7B (Base)	–	–	26.7	35.6	–	–	–	–
Standard OPD	–	–	–	–	–	–	–	–
GRPO	–	–	–	–	–	–	–	–
Self-Distilled Reasoner	–	–	–	–	–	–	–	–
UOPD (Ours)	–	–	–	–	–	–	–	–

- **Qwen3-8B:** The teacher model (upper bound).
- **Standard OPD:** On-policy distillation with uniform weighting ($w_{se} = 1$ for all prompts), following **(author?)** [3].
- **GRPO:** Group Relative Policy Optimization, a reinforcement learning approach that trains the student using reward signals from correct and incorrect rollouts.
- **Self-Distilled Reasoner:** On-policy distillation with ground-truth chain-of-thought injected into the teacher’s system prompt [1].

164 4.3 Evaluation

165 **Benchmarks.** We evaluate all methods on five public mathematical reasoning benchmarks of
166 increasing difficulty.

- **GSM8K** [15]: 1,319 grade-school math word problems.
- **AIME 2024**: 30 problems from the American Invitational Mathematics Examination.
- **AIME 2025**: 30 problems from AIME 2025 (Parts I and II combined).
- **HMMT Feb 2025**: 30 problems from the Harvard-MIT Mathematics Tournament.
- **AMO-Bench** [13]: 50 IMO-level competition problems.

172 In addition, we report accuracy on our held-out **SURE-Math test set** (200 curated problems with
173 labels).

174 **Metrics.** We report pass@1 accuracy with greedy decoding (temperature 0) and avg@16 accuracy
175 with 16 samples at temperature 1.2 and top- $p=0.95$.

176 4.4 Main Results

177 Table 1 presents the main comparison across all benchmarks.

178 4.5 Ablation Studies

179 **Effect of uncertainty weighting.** We compare three weighting strategies: (1) uniform weighting
180 ($w_{se} = 1$), which reduces to standard on-policy distillation; (2) binary filtering, which discards all
181 prompts with $H_{sem} > \tau$ for a threshold τ ; and (3) soft weighting ($w_{se} = 1 - H_{sem} / \log N$), which is
182 our default. Table 2 reports the results.

183 **Number of SE samples N .** The number of teacher rollouts N used for semantic entropy estima-
184 tion controls the resolution of the uncertainty estimate. We vary $N \in \{2, 4, 8, 16\}$ and measure
185 downstream distillation performance.

186 **Top- k logit truncation.** Transmitting the full vocabulary ($|V| = 151,936$) is memory-intensive.
187 We compare $k \in \{128, 512, 2048\}$ and full-vocabulary distillation, measuring both accuracy and
188 peak GPU memory.

Table 2: Ablation on uncertainty weighting strategies.

Weighting	GSM8K	AIME 2025	HMMT 2025	SURE-Math
Uniform ($w = 1$)	–	–	–	–
Binary filter ($\tau = 0.3$)	–	–	–	–
Binary filter ($\tau = 0.5$)	–	–	–	–
Soft weighting (UOPD)	–	–	–	–

189 **Number of student rollouts K .** We ablate $K \in \{1, 2, 4, 8\}$ on-policy rollouts per prompt to
190 understand the trade-off between training diversity and computational cost.

191 **4.6 Analysis**

192 **Convergence speed.**

193 **Semantic entropy distribution.** Figure ?? shows the distribution of semantic entropy scores across
194 SURE-Math. The majority of problems have low SE, indicating that the teacher is confident on most
195 prompts. The long tail of high-SE problems represents cases where the teacher genuinely disagrees
196 with itself, and these are precisely the samples that UOPD downweights.

197 **Qualitative examples.** Table ?? shows representative examples where UOPD’s uncertainty weight-
198 ing helps. For low-SE prompts, the teacher provides consistent guidance and the student learns
199 effectively. For high-SE prompts, the teacher produces contradictory answers; UOPD suppresses
200 these samples, preventing the student from learning incorrect reasoning patterns.

201 **5 Related Work**

202 **5.1 Knowledge Distillation for Language Models**

203 Traditional knowledge distillation [2] trains student models to match teacher output distributions.
204 For autoregressive language models, supervised KD [16] and sequence-level KD [17] are widely
205 used. However, these off-policy approaches suffer from distribution mismatch between training
206 (teacher-generated or ground-truth sequences) and inference (student-generated sequences).

207 **On-policy distillation** [3] addresses this by training students on their own generated outputs, using
208 teacher logits as labels. GKD (Generalized Knowledge Distillation) demonstrates strong improve-
209 ments on summarization and translation tasks. Our work extends on-policy distillation to reasoning
210 tasks by incorporating verification signals and contrastive objectives.

211 **5.2 Learning from Verification Feedback**

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

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

218 **5.3 Contrastive Learning for Language Models**

219 DPO [18] introduced preference-based contrastive learning for alignment, training models to prefer
220 chosen responses over rejected ones. SimPO [?] and other variants explore different contrastive
221 formulations.

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

225 **6 Conclusion**

226 We presented Token-level Contrastive Distillation (TCD), a framework that combines on-policy
227 generation, verification-based learning, and dense teacher guidance for distilling reasoning capabilities
228 into smaller language models. By unifying distillation on correct traces with contrastive learning
229 on errors, TCD effectively leverages both sparse verification signals and rich token-level teacher
230 feedback.
231 Our experiments on GSM8K demonstrate that TCD enables a 1.7B student model to achieve strong
232 mathematical reasoning performance when distilled from an 8B teacher. The framework is efficient
233 and scalable, using vLLM for on-policy generation and memory-optimized techniques for handling
234 large vocabularies.

235 **Limitations and Future Work.**

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

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

244 **References**

- 245 [1] Siyan Zhao, Zhihui Xie, Mengchen Liu, Jing Huang, Guan Pang, Feiyu Chen, and Aditya
246 Grover. Self-distilled reasoner: On-policy self-distillation for large language models. *arXiv*
247 preprint *arXiv:2601.18734*, 2026.
- 248 [2] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network.
249 *arXiv preprint arXiv:1503.02531*, 2015.
- 250 [3] Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos, Matthieu
251 Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-
252 generated mistakes. In *International Conference on Learning Representations (ICLR)*, 2024.
- 253 [4] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, and
254 Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *Advances*
255 in *Neural Information Processing Systems*, 35:24824–24837, 2022.
- 256 [5] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for
257 computer vision? In *Advances in Neural Information Processing Systems*, 2017.
- 258 [6] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances
259 for uncertainty estimation in natural language generation. In *International Conference on*
260 *Learning Representations (ICLR)*, 2023.
- 261 [7] An Yang, Anfeng Yang, Baosong Yang, et al. Qwen3 technical report. *arXiv preprint*
262 *arXiv:2505.09388*, 2025.
- 263 [8] Zhengyang Yuan, Jiawei Liu, Boyuan Zi, Hao Ning, Kai Zheng, Minghao Chen, et al. Scalequest:
264 Scalable data synthesis for mathematical reasoning. In *International Conference on Learning*
265 *Representations (ICLR)*, 2025.
- 266 [9] AI-MO Team. Numinamath-cot: A large-scale mathematical reasoning dataset with
267 chain-of-thought annotations. 2024. <https://huggingface.co/datasets/AI-MO/NuminaMath-CoT>.

- 269 [10] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
 270 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset.
 271 In *Advances in Neural Information Processing Systems*, 2021.
- 272 [11] Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma,
 273 Liang Chen, Runxin Xu, et al. Omni-math: A universal olympiad level mathematic benchmark
 274 for large language models. *arXiv preprint arXiv:2410.07985*, 2024.
- 275 [12] Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,
 276 Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for
 277 promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint
 278 arXiv:2402.14008*, 2024.
- 279 [13] Meituan Longcat Team. Amo-bench: Assessing mathematical olympiad problem solving for
 280 large language models. 2025. [https://huggingface.co/datasets/meituan-longcat/
 281 AMO-Bench](https://huggingface.co/datasets/meituan-longcat/AMO-Bench).
- 282 [14] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and
 283 Dixin Jiang. Wizardlm: Empowering large language models to follow complex instructions. In
 284 *International Conference on Learning Representations (ICLR)*, 2024.
- 285 [15] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
 286 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
 287 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 288 [16] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version
 289 of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- 290 [17] Yoon Kim and Alexander M Rush. Sequence-level knowledge distillation. *arXiv preprint
 291 arXiv:1606.07947*, 2016.
- 292 [18] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and
 293 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
 294 In *Advances in Neural Information Processing Systems*, 36, 2023.

295 A Implementation Details

296 A.1 Distributed Training Architecture

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

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

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

306 A.2 Memory Optimization Techniques

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

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

311 This achieves 150x memory reduction with minimal accuracy loss.

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

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

319 A.3 Communication Backend Selection

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

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

326 B Additional Experimental Results

327 B.1 Hyperparameter Sensitivity

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

- 329 • Learning rate sweep
- 330 • Temperature sweep for generation
- 331 • Beta sweep for contrastive loss
- 332 • Number of samples per prompt (K)

333 B.2 Evaluation Protocol Details

334 For GSM8K evaluation, we:

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

339 B.3 Computational Cost

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

- 341 • Wall-clock time: **TODO** hours on 4x H100
- 342 • GPU hours: **TODO**
- 343 • On-policy generation: **TODO**% of total time
- 344 • Teacher logits computation: **TODO**% of total time
- 345 • Student training: **TODO**% of total time

346 C Reproducibility

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