
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

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

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

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

Algorithm 1 Uncertainty-Calibrated On-Policy Distillation (UOPD)

Require: Training prompts $\mathcal{D} = \{x_i\}_{i=1}^M$, teacher p_T , student $p_S(\cdot; \theta_S)$, SE samples N , rollouts K , top- k

1: — **Phase 1: Offline Semantic Entropy (Sec. 3.1) —**

2: **for** each prompt $x_i \in \mathcal{D}$ **do**

3: Sample N responses $\{y_1^T, \dots, y_N^T\} \sim p_T(\cdot | x_i)$

4: Extract final answers; cluster into semantic classes \mathcal{C}

5: Compute semantic entropy $H_{\text{sem}}(x_i)$ (Eq. 5)

6: Store weight $w_{\text{se}}(x_i) = 1 - H_{\text{sem}}(x_i)/\log N$ (Eq. 6)

7: **end for**

8: — **Phase 2: On-Policy Distillation (Sec. 3.2) —**

9: **while** not converged **do**

10: Sample mini-batch $\mathcal{B} \subset \mathcal{D}$

11: **for** each $x \in \mathcal{B}$ **do**

12: Generate K rollouts $\{y^{(1)}, \dots, y^{(K)}\} \sim p_S(\cdot | x; \theta_S)$

13: **for** each rollout $y^{(j)}$ **do**

14: Query teacher for top- k logits at each token position of $y^{(j)}$

15: Compute token-level CE: $\ell_t = -\sum_v \tilde{p}_T(v) \log p_S(v; \theta_S)$

16: **end for**

17: Compute weighted loss $\mathcal{L}(x, y^{(j)}; \theta_S)$ (Eq. 8)

18: **end for**

19: Update θ_S with $\nabla_{\theta_S} \frac{1}{|\mathcal{B}|} \sum_{x \in \mathcal{B}} \frac{1}{K} \sum_{j=1}^K \mathcal{L}(x, y^{(j)}; \theta_S)$

20: Sync θ_S to vLLM generation engines

21: **end while**

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 Algorithm 1 summarizes the complete UOPD pipeline.

112 **4 Experiments**

113 **4.1 SURE-Math Dataset**

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

116 We first collect 1,000 mathematics problems spanning middle-school to competition level, each paired
117 with a verified ground-truth answer. Of these, 200 are reserved as a held-out test set (Section 4.3); the
118 remaining 800 are added to the training pool.

119 We further aggregate seed problems from public sources across four difficulty tiers. *Easy*: ScaleQuest-
120 Math [8] (15K). *Medium*: NuminaMath-CoT [9] (10K) and MATH [10] (12.5K). *Hard*: Omni-MATH
121 [11] (4.4K) and OlympiadBench [12] (5K). *Competition*: AIME 2024–2025, HMMT Feb & Nov
122 2025, and AMO-Bench [13].

123 Starting from these seeds, we synthesize new problems with Qwen2.5-72B-Instruct using an Evol-
124 Instruct [14] style pipeline. Six evolution strategies are applied, each realized by a dedicated
125 system prompt: *harder* (increase reasoning steps or add constraints), *rewrite* (change context and
126 wording while preserving the underlying skill), *algebraize* (replace concrete values with variables),
127 *apply* (embed the concept in a real-world scenario), *compose* (combine with a different branch of
128 mathematics), and *competition* (transform into AMC/AIME/Olympiad style). Each seed undergoes
129 1–5 rounds of iterative evolution depending on its difficulty tier: easy seeds receive 1 round, medium
130 seeds 2 rounds, hard seeds 3 rounds, and competition seeds 5 rounds. In follow-up rounds, a different
131 strategy (drawn from *harder*, *competition*, and *compose*) is applied to the output of the previous
132 round, progressively increasing difficulty. Strategy selection within each tier is weighted to favor
133 *harder* and *competition* for higher-tier seeds.

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

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

141 4.2 Experimental Setup

142 **Models.** We evaluate three student–teacher pairs of increasing scale: Qwen3-1.7B/Qwen3-8B,
143 Qwen3-4B/Qwen3-8B, and Qwen3-8B/Qwen3-30B [7], all instruction-tuned variants. Semantic
144 clustering for the offline SE computation uses a Qwen2.5-3B-Instruct judge for pairwise semantic
145 equivalence.

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

152 For each student–teacher pair, we compare UOPD against the following methods. (1) **Base:** the
153 student model without any distillation (lower bound). (2) **Teacher:** the teacher model (upper bound).
154 (3) **Standard OPD:** on-policy distillation with uniform weighting ($w_{\text{se}} = 1$ for all prompts), following
155 **(author?)** [3]. (4) **GRPO:** Group Relative Policy Optimization, which trains the student using reward
156 signals from correct and incorrect rollouts. (5) **Self-Distilled Reasoner:** on-policy distillation with
157 ground-truth chain-of-thought injected into the teacher’s system prompt [1].

158 4.3 Evaluation

159 **Benchmarks.** We evaluate all methods on three public mathematical reasoning benchmarks of
160 increasing difficulty, plus our held-out test set: **MATH-500** [10], 500 competition-level problems
161 spanning seven subjects; **AIME 2025**, 30 problems from the American Invitational Mathematics
162 Examination (Parts I and II); **AMO-Bench** [13], 50 IMO-level competition problems; and our
163 held-out **SURE-Math test set** (200 curated problems with ground-truth labels).

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

166 4.4 Main Results

167 Table 1 presents the main comparison across all benchmarks.

168 4.5 Ablation Studies

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

173 **Number of SE samples N .** The number of teacher rollouts N used for semantic entropy estimation
174 controls the resolution of the uncertainty estimate. We vary $N \in \{2, 4, 8, 16\}$ and measure
175 downstream distillation performance.

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

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

Method	MATH-500		AIME 2025		AMO-Bench		SURE-Math	
	pass@1	avg@16	pass@1	avg@16	pass@1	avg@16	pass@1	avg@16
Student: Qwen3-1.7B Teacher: Qwen3-8B								
Base (1.7B)	–	–	–	–	–	–	–	–
Standard OPD	–	–	–	–	–	–	–	–
GRPO	–	–	–	–	–	–	–	–
Self-Distilled Reasoner	–	–	–	–	–	–	–	–
UOPD (Ours)	–	–	–	–	–	–	–	–
Student: Qwen3-4B Teacher: Qwen3-8B								
Base (4B)	–	–	–	–	–	–	–	–
Standard OPD	–	–	–	–	–	–	–	–
GRPO	–	–	–	–	–	–	–	–
Self-Distilled Reasoner	–	–	–	–	–	–	–	–
UOPD (Ours)	–	–	–	–	–	–	–	–
Student: Qwen3-8B Teacher: Qwen3-30B								
Base (8B)	–	–	–	–	–	–	–	–
Standard OPD	–	–	–	–	–	–	–	–
GRPO	–	–	–	–	–	–	–	–
Self-Distilled Reasoner	–	–	–	–	–	–	–	–
UOPD (Ours)	–	–	–	–	–	–	–	–

Table 2: Ablation on uncertainty weighting strategies.

Weighting	MATH-500	AIME 2025	AMO-Bench	SURE-Math
Uniform ($w = 1$)	–	–	–	–
Binary filter ($\tau = 0.3$)	–	–	–	–
Binary filter ($\tau = 0.5$)	–	–	–	–
Soft weighting (UOPD)	–	–	–	–

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

181 4.6 Analysis

182 Convergence speed.

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

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

191 5 Related Work

192 5.1 Knowledge Distillation for Language Models

193 Traditional knowledge distillation [2] trains student models to match teacher output distributions.
194 For autoregressive language models, supervised KD [15] and sequence-level KD [16] are widely
195 used. However, these off-policy approaches suffer from distribution mismatch between training
196 (teacher-generated or ground-truth sequences) and inference (student-generated sequences).

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

201 **5.2 Learning from Verification Feedback**

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

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

208 **5.3 Contrastive Learning for Language Models**

209 DPO [17] introduced preference-based contrastive learning for alignment, training models to prefer
210 chosen responses over rejected ones. SimPO [?] and other variants explore different contrastive
211 formulations.

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

215 **6 Conclusion**

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

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

225 **Limitations and Future Work.** **Scalability:** Current experiments use 4 GPUs; scaling to multi-
226 node setups requires careful optimization of communication backends (NCCL vs. Gloo). **Generaliza-**
227 **tion:** We focus on mathematical reasoning; extending to other reasoning tasks (code, commonsense)
228 is an important direction. **Teacher quality:** Our approach assumes access to a capable teacher;
229 exploring self-improvement scenarios is promising.

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

232 **References**

- 233 [1] Siyan Zhao, Zhihui Xie, Mengchen Liu, Jing Huang, Guan Pang, Feiyu Chen, and Aditya
234 Grover. Self-distilled reasoner: On-policy self-distillation for large language models. *arXiv*
235 preprint arXiv:2601.18734, 2026.
- 236 [2] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network.
237 *arXiv preprint arXiv:1503.02531*, 2015.
- 238 [3] Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos, Matthieu
239 Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-
240 generated mistakes. In *International Conference on Learning Representations (ICLR)*, 2024.

- 241 [4] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, and
 242 Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *Advances*
 243 in *Neural Information Processing Systems*, 35:24824–24837, 2022.
- 244 [5] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for
 245 computer vision? In *Advances in Neural Information Processing Systems*, 2017.
- 246 [6] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances
 247 for uncertainty estimation in natural language generation. In *International Conference on*
 248 *Learning Representations (ICLR)*, 2023.
- 249 [7] An Yang, Anfeng Yang, Baosong Yang, et al. Qwen3 technical report. *arXiv preprint*
 250 *arXiv:2505.09388*, 2025.
- 251 [8] Zhengyang Yuan, Jiawei Liu, Boyuan Zi, Hao Ning, Kai Zheng, Minghao Chen, et al. Scalequest:
 252 Scalable data synthesis for mathematical reasoning. In *International Conference on Learning*
 253 *Representations (ICLR)*, 2025.
- 254 [9] AI-MO Team. Numinamath-cot: A large-scale mathematical reasoning dataset with
 255 chain-of-thought annotations. 2024. [https://huggingface.co/datasets/AI-MO/](https://huggingface.co/datasets/AI-MO/NuminaMath-CoT)
 256 [NuminaMath-CoT](#).
- 257 [10] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
 258 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset.
 259 In *Advances in Neural Information Processing Systems*, 2021.
- 260 [11] Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma,
 261 Liang Chen, Runxin Xu, et al. Omni-math: A universal olympiad level mathematic benchmark
 262 for large language models. *arXiv preprint arXiv:2410.07985*, 2024.
- 263 [12] Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,
 264 Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for
 265 promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint*
 266 *arXiv:2402.14008*, 2024.
- 267 [13] Meituan Longcat Team. Amo-bench: Assessing mathematical olympiad problem solving for
 268 large language models. 2025. [https://huggingface.co/datasets/meituan-longcat/](https://huggingface.co/datasets/meituan-longcat/AMO-Bench)
 269 [AMO-Bench](#).
- 270 [14] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and
 271 Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. In
 272 *International Conference on Learning Representations (ICLR)*, 2024.
- 273 [15] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version
 274 of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- 275 [16] Yoon Kim and Alexander M Rush. Sequence-level knowledge distillation. *arXiv preprint*
 276 *arXiv:1606.07947*, 2016.
- 277 [17] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and
 278 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
 279 In *Advances in Neural Information Processing Systems*, 36, 2023.

280 **A Implementation Details**

281 **A.1 Distributed Training Architecture**

282 Our implementation uses Ray [?] for distributed actor management and vLLM [?] for efficient
 283 LLM serving across 4× NVIDIA H100 GPUs. The GPU allocation differs between the two phases.

284 **Phase 1: Offline semantic entropy computation.** This phase runs as a standalone script before
 285 training begins. Two vLLM model instances are loaded on the same GPU. The **teacher vLLM**
 286 **engine** (GPU 0) serves the teacher model (e.g., Qwen3-8B) with `gpu_memory_utilization=0.65`
 287 and `tensor_parallel_size=1`; for each prompt it generates $N=8$ responses in a single batched
 288 call via `engine.generate()`. The **cluster judge vLLM engine** (GPU 0, colocated) serves the
 289 semantic equivalence judge (Qwen2.5-3B-Instruct) with `gpu_memory_utilization=0.20`; it is
 290 only invoked for answer pairs that do not match via exact string comparison, keeping its utilization
 291 low. Both engines share GPU 0 via vLLM’s memory pre-allocation. The remaining GPUs are idle
 292 during this phase. Processing 8,800 prompts with batch size 300 completes in approximately 1–2
 293 hours.

294 **Phase 2: On-policy distillation.** All 4 GPUs are utilized via Ray placement groups. **GPU 0**
 295 hosts the *student actor*: the student model (e.g., Qwen3-1.7B) wrapped in DeepSpeed ZeRO Stage 1
 296 for gradient computation and parameter updates; a frozen reference copy of the initial student is
 297 colocated on the same GPU using fractional allocation (`num_gpus=0.2`). **GPU 1** hosts the *teacher*
 298 *actor*: the frozen teacher model loaded in inference mode, which receives student-generated rollouts
 299 and returns top- k logit values and token indices at every position. **GPs 2–3** host two LLMRayActor
 300 *vLLM generation engines*, each initialized from the student’s current weights. They generate $K=4$
 301 on-policy rollouts per prompt in parallel. After each gradient step, the updated student weights are
 302 broadcast to both engines via `update_weight()` over the Gloo backend.

303 The training loop proceeds as follows at each step:

- 304 1. The Ray coordinator dispatches a mini-batch of prompts to the 2 vLLM engines.
- 305 2. Each engine generates $K/2=2$ rollouts per prompt (total $K=4$ per prompt across both
306 engines).
- 307 3. Rollouts are sent to the teacher actor, which computes top- k logits for all token positions.
- 308 4. The student actor receives rollouts, teacher logits, and precomputed w_{se} weights, then
309 performs a gradient update.
- 310 5. Updated student weights are synchronized back to both vLLM engines for the next iteration.

311 A.2 Memory Optimization Techniques

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

```

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

316 This achieves 150× memory reduction with minimal accuracy loss.

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

```

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

324 A.3 Communication Backend Selection

325 For non-colocated setups, vLLM workers require collective communication for weight up-
 326 dates. NCCL is fast but fails when actors reside in separate Ray placement groups due to
 327 CUDA_VISIBLE_DEVICES isolation. We therefore use **Gloo**, a CPU-based fallback that works across
 328 placement groups, when `tensor_parallel_size=1`.

329 **B Additional Experimental Results**

330 **B.1 Hyperparameter Sensitivity**

331 **TODO:** Add figures/tables for learning rate sweep, temperature sweep for generation, beta sweep for
332 contrastive loss, and 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 on the SURE-Math dataset takes approximately **TODO** wall-clock hours on $4 \times$ H100
341 (**TODO** GPU-hours total). The time breakdown is: on-policy generation **TODO**%, teacher logits
342 computation **TODO**%, and student training **TODO**%.

343 **C Reproducibility**

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