CPRet: A Dataset, Benchmark, and Model for Retrieval in Competitive Programming

Han Deng^{1,2,3,5}, Yuan Meng^{3,*}, Shixiang Tang², Wanli Ouyang ^{1,2,5}, Xinzhu Ma^{2,4,*}

¹ The Chinese University of Hong Kong ² Shanghai Artificial Intelligence Laboratory ³ Tsinghua University ⁴ Beihang University ⁵ Shenzhen Loop Area Institute

Abstract

Competive programming benchmarks are widely used in scenarios such as programming contests and large language model assessments. However, the growing presence of duplicate or highly similar problems raises concerns not only about competition fairness, but also about the validity of competitive programming as a benchmark for model evaluation. In this paper, we propose a new problem—similar question retrieval—to tackle the problem. Due to the lack of both data and models, solving this problem is challenging. To this end, we introduce **CPRet**, a retrievaloriented benchmark suite for competitive programming, covering four retrieval tasks: two code-centric (i.e., Text-to-Code, Code-to-Code) and two newly proposed problem-centric tasks (i.e., Problem-to-Duplicate, Simplified-to-Full)—built from a combination of automatically crawled problem-solution data and manually curated annotations. Our contribution includes both high-quality training data and temporally separated test sets for reliable evaluation. Besides, we further develop two task-specialized retrievers based on this dataset: CPRetriever-Code, trained with a novel Group-InfoNCE loss for problem-code alignment, and CPRetriever-Prob, fine-tuned for indentifying problem-level similarity. Both models achieve strong results and are open-sourced for local use. Finally, we analyze LiveCodeBench and find that high-similarity problems inflate model pass rates and reduce differentiation, underscoring the need for similarity-aware evaluation in future benchmarks.

Github: https://github.com/coldchair/CPRet **Online Demo:** https://www.cpret.online/

1 Introduction

Competitive programming contests—from high-school Olympiads like the IOI to university-level events such as the ICPC—challenge participants to solve algorithmic problems under tight time and memory limits, requiring both strong coding skills and deep algorithmic insight. Because problem statements are precisely specified in natural language and solutions can be graded automatically, these problems have become a canonical benchmark for assessing the reasoning and coding abilities of large language models (LLMs) [1, 2]. Recent advances in code-oriented LLMs—spanning commercial flagships such as OpenAI o4-mini [3], Gemini-2.5-Pro [4], and Grok-3-Mini (High) [5], alongside research releases like DeepSeek-R1 [6] highlight the value of competitive programming as a testbed for algorithmic reasoning, program synthesis, and computer-science education.

Competitive programming has witnessed rapid growth over the past three decades, with thousands of new problems introduced annually. This expansion has led to a growing concern: the increasing presence of similar or repetitive problems within large repositories. As shown in Figure 1, our

^{*}Corresponding author.

collected data reveals a significant rise in community discussions around duplicate problems in recent years. However, it is often difficult to determine whether a new problem is a duplicate of existing ones, as this typically relies on the memory and judgment of human problem setters. Unchecked duplication brings concrete downsides. In human programming contests, repeated or highly similar problems give an unfair advantage to participants who have seen them before. In competitive programming benchmarks for LLMs, the presence of repeated or highly similar problems can lead to inflated performance scores, as models may rely on memorized instances of past competition problems from their training data rather than demonstrating genuine algorithmic reasoning. This compromises the evaluation of their true ability to solve novel and unfamiliar challenges.

Competitive programming problem retrieval poses unique challenges compared to standard code retrieval: problem statements are often abstract, narrative-driven, and may admit multiple fundamentally different solutions, requiring models to capture high-level algorithmic ideas rather than surface code patterns. This motivates the need for specialized problem-level retrieval models.

A promising way to address the issue of duplicate problems is to leverage retrieval models that identify semantically similar problems based on learned representations. Such models offer a scalable and model-agnostic solution for redundancy detection in large problem repositories. However, despite their potential, there is a lack of established benchmarks specifically targeting problem-level retrieval. Prior efforts in competitive programming retrieval have primarily focused on code-centric tasks, such as retrieving

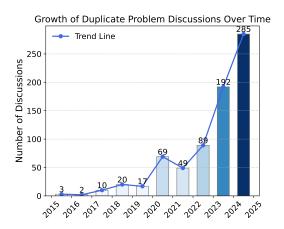


Figure 1: Year-wise trend of annotated duplicate problem discussions.

solution code given a problem description (*Text-to-Code*) or retrieving alternative correct solutions given a reference implementation (*Code-to-Code*), which overlook the dimension of retrieving similar problems themselves—an essential capability for both redundancy detection and educational search.

To fill this gap, we introduce two new problem-centric retrieval tasks. (I) *Problem-to-Duplicate* asks a model to retrieve the most semantically or structurally similar problem to a given target. We manually annotated 700 duplicate pairs; 200 pairs form a held-out test set in which the model must locate the duplicate among a corpus of 10,000+ problems, and the remaining ~ 500 pairs are used for training. (II) *Simplified-to-Full* supplies a simplified or paraphrased statement and requires the model to find the original full version. We collected 17,000 human-written simplifications aligned with their source problems; 10,000 pairs constitute the test set and the remaining $\sim 7,000$ pairs serve as training data. These tasks target practical problem-retrieval needs—spotting duplicate or highly similar problems and retrieving the original statement from a simplified description. In addition, we perform a temporal analysis of the two code-centric tasks and uncover severe train-test leakage: models fare markedly worse on the newest problems. To this end, we rebuild both benchmarks, drawing the test sets exclusively from the most recent problems to enforce strict temporal separation.

Furthermore, to address the real-world retrieval demands, we develop two task-specialized models CPRetriever-Code and CPRetriever-Prob using our new corpus: (I) the first one is trained on 38 k problems and 2.9 M multi-language solutions with a new *Group-InfoNCE* loss that treats every correct solution of the same problem as a positive and penalizes similarity within that positive set. This encourages alignment across diverse implementations, and it achieves the overall best performance on all code-centric tasks. (II) The second one is fine-tuned from CPRetriever-Code on our duplicate-detection and simplified-to-full datasets, and it performs the overall best results on the two problem-centric tasks.

Finally, we apply CPRetriever-Prob to analyze problems in LiveCodeBench released after September 2024 by computing each problem's maximum similarity to earlier data. We observe two key trends: (I) *Pass rates increase with similarity*—problems more similar to prior examples are consistently easier for models to solve, regardless of difficulty; (II) *High similarity reduces model differentiation*—as

similarity grows, performance gaps between models shrink, suggesting that models may solve such problems by relying on memorized patterns rather than demonstrating true reasoning ability in novel scenarios. These findings highlight the importance of accounting for problem similarity when constructing benchmarks to avoid overestimating model capabilities.

Overall, our key contributions can be summarized as follows:

- Comprehensive Benchmark: We construct a large-scale, high-quality dataset for competitive programming, covering four key retrieval tasks (two focused on code retrieval and two newly proposed for problem retrieval). Our benchmark reflects realistic retrieval needs in programming contests and addresses issues such as problem repetition and dataset leakage.
- Effective Models: We develop two task-specialized retrieval models with strong performance, trained with our newly designed Group-InfoNCE loss. Both models are open-sourced and can be run locally—an important feature given the privacy constraints in contest settings—providing practical tools for code and problem retrieval and helping mitigate the growing issue of duplicate problems in competitive programming.
- Similarity-Aware Evaluation: We provide insightful studies on the impact of problem similarity in competitive programming, and find that problems with high similarity to prior data yield inflated pass rates and reduced model differentiation, especially for weaker models. These results highlight the importance of controlling for similarity when designing benchmarks and suggest that future datasets should stratify or filter test problems based on their similarity to past content.

2 Related Work

2.1 Datasets for Competitive Programming

Several datasets have been proposed to train and evaluate LLMs on competitive programming tasks. Early efforts include Description2Code [7], which gathered problems and solutions from CodeChef and Codeforces, followed by the APPS dataset [8], which became a widely used benchmark across a mix of competition-style and introductory problems. CodeContests [9], used in DeepMind's AlphaCode, expanded coverage but remained limited in scale. More recently, TACO [10] constructed the largest dataset (26K problems) by aggregating multiple sources. In the retrieval setting, CoIR [11] built on APPS and CodeTransOcean [12] to define two retrieval tasks: text-to-code and code-to-code, helping assess model capabilities in aligning problem descriptions and solution code. Despite their contributions, these datasets have several limitations: (i) many were collected before 2023 and suffer from potential data leakage, (ii) problem types are often narrow in scope, and (iii) both problem statements and solutions are typically limited to English and Python, reducing linguistic and implementation diversity.

For evaluating programming ability, HumanEval [2] remains a widely used standard, featuring hand-written Python problems with corresponding unit tests for functional correctness. LiveCodeBench [1] improves upon earlier efforts by continuously collecting real-world problems from LeetCode, AtCoder, and Codeforces, and applies temporal separation to reduce data contamination. However, these benchmarks do not account for problem-level similarity between new and historical problems—a distinct issue from data leakage—which can still lead to overly optimistic assessments of model performance due to hidden redundancy.

2.2 Embedding Models for Retrieval

Dense retrieval models learn to map queries and candidates into a shared embedding space for efficient similarity search. Early methods like DSSM [13] and Sentence-BERT [14] paved the way for supervised contrastive training, leading to influential retrievers such as Dense Passage Retrieval (DPR) [15].

Recent efforts have focused on building general-purpose embedding models (e.g., E5-Mistral [16], GTE-Qwen2 [17], SFR-Embedding [18], NV-Embed [19, 20], Linq-Embed[21]) that perform well across a wide range of retrieval tasks. These models are commonly evaluated on MTEB [22], a standardized benchmark suite covering diverse tasks such as semantic search, classification, and reranking. For the code domain, specialized models like CodeSage [23], SFR-Embedding-Code [24], Qodo-Embed[25] use code-specific data and retrieval-aware objectives to significantly improve

performance on benchmarks such as CoIR. While both general-purpose and code-specific embedding models can be applied to competitive programming, we find their performance on code-centric and problem-centric retrieval tasks remains limited—highlighting the need for a dedicated solution tailored to the unique characteristics of this domain.

High-quality retrieval has been shown to benefit downstream tasks like code generation via Retrieval-Augmented Generation (RAG), where an LLM leverages retrieved problem-solution pairs as context. For example, Shi et al. [26] explored solving Olympiad-level problems with RAG, Li et al. [27] proposed critic-guided retrieval-augmented planning, and Mapcoder [28] introduced MapCoder, a multi-agent code generation framework using retrieval. These studies suggest that problem-level retrieval can substantially enhance LLM performance, motivating specialized embeddings for competitive programming.

While temporal splits and duplicate-problem detection are not new concepts—having been explored in other domains such as finance [29] and software community question answering [30]—their systematic application to competitive programming remains underexplored. Existing benchmarks (e.g., MTEB[22], CoIR[11]) often lack strict temporal separation, potentially inflating evaluation results. Our work extends these established principles to programming contests, where duplicate or temporally overlapping problems introduce unique challenges: algorithmically similar tasks may differ superficially in text but still bias performance. By constructing a benchmark with rigorous temporal splits and fine-grained similarity filtering, we aim to build fairer evaluation standards and reveal new insights into the relationship between problem similarity and model discrimination ability.

3 Method

3.1 Overview of Retrieval Tasks

We define and evaluate four distinct retrieval tasks, each capturing a different dimension of semantic understanding and reuse in competitive programming:

- **Text-to-Code Retrieval**: Given the full natural language description of a problem, retrieve one or more correct solution codes. This task assesses a model's ability to align problem semantics with executable implementations.
- Code-to-Code Retrieval: Given one accepted solution, retrieve alternative correct solutions to the same problem. This task evaluates a model's capacity to understand code functionality and identify semantically equivalent but syntactically diverse implementations.
- **Problem-to-Duplicate Retrieval**: Given a problem description, retrieve other problems that are duplicated, including exact matches or closely related in terms of solution strategy. This task is useful for identifying redundancy and measuring problem novelty across platforms.
- **Simplified-to-Full Retrieval**: Given a simplified version of a problem, retrieve the corresponding one with full description. This task examines cross-abstraction retrieval capabilities, bridging accessibility-oriented rewrites and original problem statements.

These tasks jointly form a comprehensive benchmark for evaluating retrieval models in the context of competitive programming, covering problem-code alignment, solution diversity, duplication detection, and abstraction-level matching. Table 1 summarizes the statistics of each retrieval task, including training and test set sizes, as well as average token lengths for both queries and corpus items.

3.2 Dataset and Benchmark Construction

3.2.1 Problems and Codes Data

We construct a large-scale dataset of programming contest problems and their corresponding accepted solutions to support both model training and the evaluation of retrieval tasks. Compared to existing datasets, our dataset provides *broader temporal coverage*, *richer language diversity*, and *more varied contest formats*. It includes recent problems from multiple online judges, spans several programming languages, and covers both ICPC-style (International Collegiate Programming Contest) and OI-style (Olympiad in Informatics) problems—the latter involving partial scoring and more complex algorithmic requirements. Each problem is paired with one or more solutions and annotated with

Table 1: Statistics of the four retrieval tasks in CPRet. We report the number of training and test items, where #Train-Code indicates the number of distinct solution codes, and #Train-Pair refers to the number of (anchor, positive) pairs used for training. $L_{\rm Query}$ and $L_{\rm Corpus}$ denote average token lengths.

Retrieval Task	#Train- Problem	#Train- Code	#Train- Pair	#Test- Query	#Test- Corpus	#Test- Qrels	L_{Query}	$L_{ m Corpus}$
Text-to-Code	38.8K	2.93M	2.93M	4.9k	41.6k	41.6k	1038	1210
Code-to-Code				4.8k	39.8k	39.8k	1132	1185
Problem-to-Duplicat	te 874	/	491	168	10.9k	202	565	765
Simplified-to-Full	7.6k	/	7.6k	10k	10k	10k	226	697

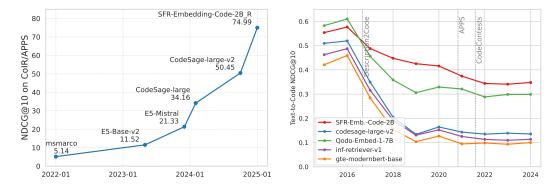


Figure 2: **Left:** NDCG@10 scores on CoIR/APPS across model release dates, showing steady performance improvement as models evolve. **Right:** Performance on the Text-to-Code task grouped by problem release year, revealing sharp degradation on older problems, especially around major dataset release points (*Description2Code*, *APPS*, *CodeContests*).

precise timestamp information, enabling temporally-aware training and evaluation across semantic retrieval tasks. As summarized in Table 2, our dataset—**CPRet-PCPCD**—offers wider coverage across problem sources, types, and languages, and serves as a stronger foundation for retrieval-based modeling in competitive programming.

Table 2: **Comparison of competitive programming datasets. #Src**: Number of data sources (*e.g.*, online judges). **#P-Type**: Number of problem types (*e.g.*, ICPC-style with full-score only vs OI-style with partial scoring). **#Lang-P**: Number of languages used for problem descriptions. **#Lang-C**: Number of programming languages used in solutions. See Appendix A.5 for more detailed dataset statistics.

Dataset	#Prob	#Code	#Src	#P-Type	#Lang-P	#Lang-C	Cut-off
Description2Code [7]	7.8K	309K	3	1	1	2	2016/08
APPS[8]	10K	232K	7	1	1	1	2020/10
CodeContests[9]	13.6K	4.5M	5	1	1	4+	2021/07
TACO[10]	26.4K	1.55M	10	1	1	1	2023/02
CPRet-PCPCD (ours)	42.2K	2.9M	12	2	3	20+	2024/12

As shown in Figure 2, on the Text-to-Code task using the historical APPS dataset, model performance improves over time as newer and stronger models are introduced. However, we observe a clear drop in performance when problems are grouped by their original release year—only stabilizing after 2022. This indicates potential data leakage or memorization in earlier benchmarks and underscores the need for temporally separated evaluation. Accordingly, we use problems and solutions from 2023 onward in our collected dataset as the test set for both the Text-to-Code and Code-to-Code tasks.

3.2.2 Duplicate Programming Problem Pairs

In many instances, when a newly released contest problem closely resembles or duplicates a previously published one, participants will highlight this in the associated discussion threads or contest forums. To leverage this observation, we collected all publicly available discussion threads and blog posts from two major competitive programming platforms: Codeforces and Luogu. We applied a combination of keyword-based heuristics and large language model (LLM)-based classification to identify approximately 5,000 potentially relevant entries. These candidates were subsequently verified through manual annotation by several experienced competitive programmers.

To ensure consistency and clarity in annotation, we defined three levels of duplication between problems:

- Exact Match: An accepted (AC) solution for one problem can directly pass the other without modification.
- Near Match: An AC solution for one problem can be adapted to solve the other with minor edits that could reasonably be made by someone unfamiliar with the second problem.
- Method Match: The core idea and solution approach are the same, but the code differs in non-essential or implementation-specific details.

We ultimately identified around 700 pairs of duplicate problems. Since some problems may belong to duplication clusters involving more than two problems, we first performed clustering to group mutually duplicate problems. We then randomly selected 30% of the clusters to construct the test set for the **Problem-to-Duplicate** retrieval task. Within each selected cluster, one problem was randomly designated as the query, while the remaining ones were placed in the corpus. To further increase task difficulty and realism, we added all other Codeforces problems to the corpus as distractors.

3.2.3 Simplified and Full Problem Description Pairs

To support beginners and non-native speakers, users on the Luogu platform have contributed Chinese translations and simplified versions of competitive programming problems originally written in English (*e.g.*, on Codeforces) or Japanese (*e.g.*, on AtCoder). We crawled these user-generated simplifications and applied filtering procedures to remove low-quality entries, particularly those stemming from direct or unedited machine translation.

After cleaning, we obtained approximately 17,000 high-quality pairs of full problem descriptions and their corresponding simplified versions. From this set, we randomly selected 10,000 pairs to construct the test set for the **Simplified-to-Full** retrieval task, where the simplified description serves as the query and the full version is used as the corpus. The remaining around 7,000 pairs are used for training.

3.3 Competitive Programming Retrievers

3.3.1 CPRetriever-Code: Multi-Positive Contrastive Learning for Problem-Code Alignment

We train CPRetriever-Code using a supervised contrastive learning framework tailored for aligning problem descriptions with diverse correct solutions. Each problem is paired with multiple valid codes, and the model learns to capture this alignment while preserving representation diversity.

We begin with the standard InfoNCE loss [31], where each (problem, solution) pair (x_i, x_i^+) is contrasted against others in the batch:

$$\mathcal{L}_{\text{InfoNCE}} = -\log \frac{\exp(\sin(x_i, x_i^+)/\tau)}{\sum_{j=1}^N \exp(\sin(x_i, x_j^+)/\tau)}.$$
 (1)

Limitation of Single-Positive Loss. In practice, many problems have multiple correct solutions $G_i = x_i^{1+}, \dots, x_i^{m+}$. A common extension averages the similarities over all positives [32]:

$$\mathcal{L}_{\text{MultiPos}} = -\log \frac{\sum_{k=1}^{m} \exp\left(\sin(x_i, x_i^{k+})/\tau\right)}{\sum_{i \neq j}^{N} \exp\left(\sin(x_i, x_i)/\tau\right)}.$$
 (2)

However, this ignores the internal structure of the positive set and does not encourage consistency among the positives.

Group-InfoNCE. To better utilize multiple correct solutions, we propose Group-InfoNCE, which treats the positive set as a whole. The loss promotes similarity between the query x_i and its group G_i , while contrasting against other problems and groups:

$$\mathcal{L}_{\text{Group}} = -\log \frac{\exp\left(\sin_G(x_i, G_i)/\tau\right)}{\exp\left(\sin_G(x_i, G_i)/\tau\right) + \sum_{j \neq i} \left[\exp\left(\sin(x_i, x_j)/\tau\right) + \exp\left(\sin_G(x_i, G_j)/\tau\right)\right]} + \frac{\text{Penalty}_G(x_i, G_i)}{\tau^2}.$$
(3)

Here, group similarity is defined as:

$$sim_G(x_i, G_j) = \frac{1}{m} \sum_{k=1}^m sim(x_i, x_j^{k+}),$$
(4)

and the variance-based regularization encourages consistency within the group:

$$\operatorname{Penalty}_{G}(x_{i}, G_{i}) = \operatorname{Var}_{k=1}^{m} \left(\operatorname{sim}(x_{i}, x_{i}^{k+}) \right). \tag{5}$$

This group-based formulation improves representation quality by explicitly modeling the structure of multiple correct solutions, resulting in more stable and discriminative embeddings for code-centric tasks. Although Group-InfoNCE enforces consistency between a problem embedding and its set of correct solutions, it does not constrain the diversity of the solutions themselves. The variance-based regularization ensures that each solution maintains a consistent similarity to the problem embedding, effectively "surrounding" it in the code space without collapsing distinct approaches. In practice, these similarities converge around an average value, allowing the problem representation to capture the essential idea while accommodating diverse algorithmic strategies. Empirical results on both base models (see Tables 4 and 5) show consistent improvements in retrieval performance and robust, discriminative embeddings.

Format Masking for Robustness. To reduce reliance on superficial cues, we apply random masking to non-essential parts of the problem description during training, including I/O format explanations, sample inputs/outputs, and explicit data constraints. This encourages the model to focus on the core algorithmic intent rather than overfitting to formatting patterns, and improves generalization to problems with diverse or unfamiliar layouts.

3.3.2 CPRetriever-Prob: Fine-Tuning on Problem-Level Tasks

CPRetriever-Prob is obtained by fine-tuning CPRetriever-Code on problem-level retrieval tasks. While CPRetriever-Code captures general alignment between problems and solutions, this stage specializes the model for retrieving semantically or structurally related problems—such as duplicates or simplified versions.

In particular, we jointly train on data from the *Problem-to-Duplicate* and *Simplified-to-Full* tasks (Section 3.1), enabling the model to learn fine-grained alignment between problems that are either functionally equivalent or differ primarily in abstraction level.

Training Format and Loss. Each training example is structured as a triplet (x, x^+, x^-) , where x is the query, x^+ is a semantically aligned problem, and x^- is a hard negative. We use the *triplet margin loss* [33]:

$$\mathcal{L}_{\text{triplet}} = \max\left(0, \sin(x, x^{-}) - \sin(x, x^{+}) + \alpha\right),\tag{6}$$

where $sim(\cdot, \cdot)$ is cosine similarity and α is the margin.

Hard Negative Mining. As our dataset primarily contains positive pairs, we construct hard negatives by first retrieving the top-10 most similar candidates for each query using a pretrained retriever (SFR-Embedding-Code-2B_R[24]) and randomly sampling one as a negative. To ensure the sampled negative is not a false match, we use Qwen-2.5-Max[34] to verify that it is not equivalent to the query in terms of problem intent. This automatic filtering yields challenging yet reliable negatives without requiring manual annotation.

4 Experiment

4.1 Experimental Setup

We develop two retrieval models based on the SFR-Embedding-Code-2B_R[24] embedding model, which is built on top of Gemma-2-2B[35] and further pre-trained on code-related tasks. We further evaluate our method on the Qwen3-Embedding-4B[36] model, released in June 2025, which demonstrates stronger retrieval capabilities. Detailed training configurations are provided in Appendix A.1.

4.2 Main Results

Table 3: Retrieval performance across models with different scales. We use NDCG@10 [37] as the primary metric to measure ranking quality across tasks. **T2C**: Text-to-Code. **C2C**: Code-to-Code. **P2Dup**: Problem-to-Duplicate. **S2Full**: Simplified-to-Full. **Avg**: Average of the four tasks.

Scale	Model	Type	Size	T2C	C2C	P2Dup	S2Full	Avg
	contriever-msmarco [38]	general	109M	1.06	7.49	5.26	43.04	14.21
Tiny	multilingual-e5-small [39]	general	110M	2.73	13.52	14.90	52.07	20.80
Tilly	codesage-small-v2 [23]	code	130M	10.16	22.24	17.42	67.07	29.22
	gte-modernbert-base [40]	general	149M	14.99	36.22	21.12	77.45	37.44
	bge-large-zh-v1.5 [41]	general	324M	2.58	7.49	14.75	37.34	15.54
	stella_en_400M_v5 [42]	general	400M	2.12	9.76	13.14	57.56	20.64
	multilingual-e5-base [39]	general	278M	1.96	13.66	17.24	58.26	22.78
	bge-m3 [43]	general	569M	5.69	12.67	20.12	63.81	25.57
Small	codesage-base-v2 [23]	code	356M	17.30	29.38	19.97	74.36	35.25
	SFR-EmbCode-400M [24]	code	400M	9.43	43.59	19.40	75.31	36.93
	multilingual-e5-large [39]	general	560M	4.27	18.51	19.19	65.02	26.75
	multie5-large-instruct [39]	general	560M	6.64	28.84	23.31	61.28	30.02
	Qwen3-Embedding-0.6B [36]	general	600M	48.96	60.49	36.26	81.63	56.83
	gte-Qwen2-1.5B-inst. [44]	general	1.5B	10.76	23.41	27.06	69.15	32.60
	stella_en_1.5B_v5 [42]	general	1.5B	9.22	21.40	29.45	72.91	33.24
	inf-retriever-v1-1.5b [45]	general	1.5B	18.31	28.17	30.11	74.19	37.70
	codesage-large-v2 [23]	code	1.3B	20.78	35.23	22.43	78.70	39.28
	Qodo-Embed-1-1.5B [25]	code	1.5B	22.93	36.52	33.37	84.05	44.22
Medium	SFR-Embedding-2 [18]	general	2B	20.57	50.38	35.96	73.02	44.98
Medium	SFR-EmbCode-2B [24]	code	2B	39.60	68.05	45.26	86.43	59.84
	CPRetriever-Code	code	2B	70.40	70.59	38.68	81.45	65.28
	CPRetriever-Prob	code	2B	56.50	70.68	60.06	90.74	69.50
	Qwen3-Embedding-4B [36]	general	4B	66.62	71.97	56.59	89.39	71.15
	CPRetriever-Code-Qwen3-4B	code	4B	86.22	86.70	41.14	88.10	75.54
	CPRetriever-Prob-Qwen3-4B	code	4B	80.84	87.10	74.33	96.15	84.60
	GritLM-7B [46]	general	7B	0.22	8.74	11.18	29.81	12.49
	NV-Embed-v2 [20, 19]	general	7B	7.09	34.88	25.49	63.14	32.65
	inf-retriever-v1 [45]	general	7B	17.43	25.38	30.85	78.46	38.03
Large	gte-Qwen2-7B-instruct [44]	general	7B	17.96	30.72	35.95	78.41	40.76
Large	SFR-EmbMistral [18]	general	7B	22.15	50.92	31.88	69.38	43.58
	Linq-Embed-Mistral [21]	general	7B	21.99	52.06	36.51	72.79	45.84
	Qodo-Embed-1-7B [25]	code	7B	36.47	51.91	47.15	91.17	56.68
	Qwen3-Embedding-8B [36]	general	8B	60.54	72.97	53.23	87.95	68.67

We evaluate over 20 strong embedding models drawn from top-performing entries on the MTEB[22] and CoIR[11] benchmarks, as summarized in Table 3. Models trained specifically for the code domain continue to outperform general-purpose models by a significant margin. Among them, our two proposed models—**CPRetriever-Code** and **CPRetriever-Prob**—achieve the best overall results across tasks.

CPRetriever-Code, trained with Group-InfoNCE on problem-code pairs, excels on the code-centric tasks—Text-to-Code (T2C) and Code-to-Code (C2C). After fine-tuning on problem-level data,

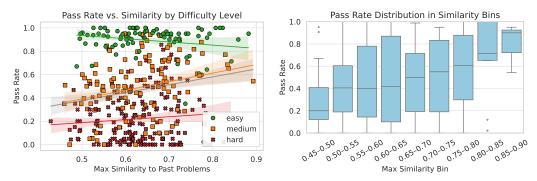


Figure 3: Impact of similarity to prior problems on pass rate. Left: For 388 post-2024.9.1 LiveCodeBench problems, we plot average model pass rate vs. maximum similarity to historical CPRet problems, computed via CPRetriever-Prob. Regression lines are shown by difficulty, with the gray line indicating the overall trend across all problems. Right: Pass rates increase monotonically across similarity bins, confirming a strong link between retrieval similarity and generation success.

CPRetriever-Prob achieves large gains on the problem-centric tasks—Problem-to-Duplicate (P2Dup) and Simplified-to-Full (S2Full)—with a modest drop on C2C and a more noticeable decline on T2C.

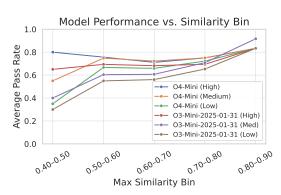
We attribute the T2C degradation to differing retrieval demands: T2C relies on implementation-specific details in problem descriptions, whereas P2Dup and S2Full emphasize higher-level semantic similarity. This trade-off is also evident in other models—for instance, *Qodo-Embed-1-7B*[25] achieves the best S2Full score but underperforms on T2C compared to *SFR-Embedding-Code-2B_R*[24]. These observations motivate our decision to release two task-specialized models, each optimized for different objectives. Further analysis of this trade-off is provided in Appendix A.1.2, and detailed ablation studies on the Group-InfoNCE loss appear in Appendix A.1.1.

4.3 Impact of Problem Similarity on Model Success in Competitive Programming Tasks

To assess how prior problem similarity influences model performance, we examine 388 Live-CodeBench problems released after September 1, 2024. For each, we compute the average pass rate across all evaluated models and measure its maximum cosine similarity to earlier CPRet-PCPCD problems using **CPRetriever-Prob**.

As shown in Figure 3, we observe a clear trend: on average, problems with higher maximum similarity to past data have significantly higher pass rates. This pattern holds across difficulty levels, but with notable differences: (I) Medium-difficulty problems show the strongest positive correlation, suggesting that models benefit substantially from prior exposure to semantically similar content. (II) Hard problems exhibit generally low pass rates, though a slight upward trend remains visible with increasing similarity. (III) Easy problems maintain consistently high pass rates overall. Due to a few low-performing outliers at high similarity, the regression line shows a slight downward slope—likely an artifact of skewed distribution rather than a true negative effect. The binned box plot (right) further confirms this relationship, reinforcing that semantic similarity to known problems is a key factor in model success.

To examine how model performance relates to problem similarity, we analyze evaluation results from OpenAI's 03-Mini and 04-Mini variants as reported in LiveCodeBench. Each model is available in three modes—Low, Medium, and High—which correspond to increasing levels of resource consumption (e.g., longer context length, higher inference cost, and possibly larger internal activations). As shown in the left panel of Figure 4, we observe that the performance gap between these variants narrows as the maximum similarity to past problems increases, and nearly vanishes in the highest bin (0.80–0.90). In contrast, on low-similarity problems, higher-capacity models (especially the High variants) maintain strong pass rates, while lower-capacity models degrade significantly. This suggests that high-similarity problems may allow weaker models to perform well by potentially leveraging memorized patterns or surface-level matching, while low-similarity problems are more likely to expose differences in model capability—such as generalization and robustness in reasoning. The right panel shows that easy problems have slightly higher similarity to



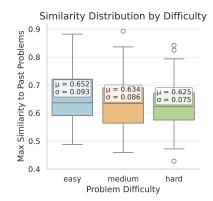


Figure 4: **Left:** Models with higher capacity (High) maintain relatively high performance even on low-similarity problems, while lower-capacity models benefit more from high-similarity problems. **Right:** Easy problems show slightly higher maximum similarity to past problems than hard ones, but the difference is limited.

past data than hard ones, but the difference is modest. This indicates that similarity is not determined by difficulty alone, and should be treated as an independent factor when constructing benchmarks.

5 Limitations and Future Work

Ongoing test set leakage. Although we use the most recent problems (from 2023–2024) to construct temporally-separated test sets for code-centric tasks, these examples are likely to appear in future model training corpora as web-scale datasets are updated. As a result, benchmark performance may become inflated over time, even without intentional misuse. To address this, we plan to release regular updates—every 6 to 12 months—to maintain a clean test split while ensuring the retrieval corpus remains diverse and representative.

Incomplete and emerging duplicate problems. Despite 700 annotated pairs in our Problem-to-Duplicate task, many duplicates—past and future—remain undiscovered, especially on decentralized or low-resource platforms. To address this, we are developing an open-source retrieval platform that supports community-driven annotations and helps contest organizers detect redundancy prepublication.

Limited analysis beyond programming. While our study focuses on competitive programming, duplication and high similarity may also affect evaluation in other domains, such as mathematical competitions (*e.g.*, IMO, AIMO). These issues remain underexplored despite their potential to undermine fair assessment. We encourage future work to examine problem reuse and redundancy in such benchmarks.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (Grant No. 62402264), and by the JC STEM Lab of AI for Science and Engineering, funded by The Hong Kong Jockey Club Charities Trust and the Research Grants Council of Hong Kong (Project No. CUHK14213224).

References

- [1] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- [2] 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, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021.
- [3] OpenAI. Introducing openai o3 and o4-mini. https://openai.com/index/introducing-o3-and-o4-mini/, April 2025. Accessed 2025-05-15.
- [4] Google DeepMind. Gemini 2.5 pro. https://deepmind.google/technologies/gemini/pro/, 2025. Accessed: 2025-05-15.
- [5] xAI. Grok 3 mini (think). https://x.ai/blog/grok-3, 2025. Accessed: 2025-05-15.
- [6] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [7] Ethan Caballero, . OpenAI, and Ilya Sutskever. Description2Code Dataset, 8 2016.
- [8] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge competence with apps. *NeurIPS*, 2021.
- [9] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.
- [10] Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and Ge Li. Taco: Topics in algorithmic code generation dataset. arXiv preprint arXiv:2312.14852, 2023.
- [11] Xiangyang Li, Kuicai Dong, Yi Quan Lee, Wei Xia, Yichun Yin, Hao Zhang, Yong Liu, Yasheng Wang, and Ruiming Tang. Coir: A comprehensive benchmark for code information retrieval models, 2024.
- [12] Weixiang Yan, Yuchen Tian, Yunzhe Li, Qian Chen, and Wen Wang. Codetransocean: A comprehensive multilingual benchmark for code translation. arXiv preprint arXiv:2310.04951, 2023.
- [13] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, 2013.
- [14] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bertnetworks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing, 2019.
- [15] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, 2020.

- [16] Yue Wang, Xiao Liu, Yizhe Liu, Hexiang Hu, and Xiaojun Wan. Improving text embeddings with large language models. *arXiv preprint arXiv:2312.04364*, 2023.
- [17] Kan Chen, Yulong Zhang, Yunbo Li, Qingkai Liu, Jiajun Jiang, Chuanqi Yang, Wei Zhou, and et al. Towards general text embeddings with multi-stage contrastive learning. *arXiv* preprint arXiv:2402.19458, 2024.
- [18] Yulong Yuan, Ming Zhao, Tianyi Chen, Wei Li, and et al. Sfr-embedding-mistral: Enhance text retrieval with transfer learning. *arXiv* preprint arXiv:2403.02745, 2024.
- [19] Xingyao Wang, Shruti Prabhumoye, Aakanksha Chowdhery, Sharan Narang, and et al. Nvembed: Improved techniques for training llms as generalist embedding models. *arXiv preprint arXiv:2403.04142*, 2024.
- [20] Gabriel de Souza P Moreira, Radek Osmulski, Mengyao Xu, Ronay Ak, Benedikt Schifferer, and Even Oldridge. Nv-retriever: Improving text embedding models with effective hard-negative mining. *arXiv preprint arXiv:2407.15831*, 2024.
- [21] Junseong Kim, Seolhwa Lee, Jihoon Kwon, Sangmo Gu, Yejin Kim, Minkyung Cho, Jy yong Sohn, and Chanyeol Choi. Linq-embed-mistral:elevating text retrieval with improved gpt data through task-specific control and quality refinement. Linq AI Research Blog, 2024.
- [22] Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*, 2022.
- [23] Shumin Yang, Liang Pan, Shizhu He, Zhiyuan Lin, Yankai Lin, and Maosong Sun. Code representation learning at scale. *arXiv preprint arXiv:2401.06585*, 2024.
- [24] Ye Liu, Rui Meng, Shafiq Jot, Silvio Savarese, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. Codexembed: A generalist embedding model family for multiligual and multi-task code retrieval. *arXiv preprint arXiv:2411.12644*, 2024.
- [25] Qodo Ltd. Qodo-embed-1-7b: A 7b parameter code embedding model for software retrieval. https://huggingface.co/Qodo/Qodo-Embed-1-7B, 2025. Accessed: 2025-05-10.
- [26] Quan Shi, Michael Tang, Karthik Narasimhan, and Shunyu Yao. Can language models solve olympiad programming? *arXiv preprint arXiv:2404.10952*, 2024.
- [27] Xingxuan Li, Weiwen Xu, Ruochen Zhao, Fangkai Jiao, Shafiq Joty, and Lidong Bing. Can we further elicit reasoning in llms? critic-guided planning with retrieval-augmentation for solving challenging tasks. *arXiv* preprint arXiv:2410.01428, 2024.
- [28] Md Ashraful Islam, Mohammed Eunus Ali, and Md Rizwan Parvez. Mapcoder: Multi-agent code generation for competitive problem solving. *arXiv preprint arXiv:2405.11403*, 2024.
- [29] Marcos Lopez De Prado. Advances in financial machine learning. John Wiley & Sons, 2018.
- [30] Rima Hazra, Debanjan Saha, Amruit Sahoo, Somnath Banerjee, and Animesh Mukherjee. Duplicate question retrieval and confirmation time prediction in software communities. In *Proceedings of the International Conference on Advances in Social Networks Analysis and Mining*, pages 203–212, 2023.
- [31] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. In *arXiv preprint arXiv:1807.03748*, 2018.
- [32] Kaiyan Zhao, Qiyu Wu, Xin-Qiang Cai, and Yoshimasa Tsuruoka. Leveraging multi-lingual positive instances in contrastive learning to improve sentence embedding. *arXiv preprint arXiv:2309.08929*, 2023.
- [33] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 815–823, 2015.

- [34] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. arXiv preprint arXiv:2412.15115, 2024.
- [35] Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- [36] Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, et al. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*, 2025.
- [37] Kalervo Järvelin and Jaana Kekäläinen. Cumulated gain-based evaluation of ir techniques. *ACM Transactions on Information Systems (TOIS)*, 20(4):422–446, 2002.
- [38] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning, 2021.
- [39] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Multilingual e5 text embeddings: A technical report. arXiv preprint arXiv:2402.05672, 2024.
- [40] Xin Zhang, Yanzhao Zhang, Dingkun Long, Wen Xie, Ziqi Dai, Jialong Tang, Huan Lin, Baosong Yang, Pengjun Xie, Fei Huang, et al. mgte: Generalized long-context text representation and reranking models for multilingual text retrieval. In *Proceedings of the 2024 Conference* on Empirical Methods in Natural Language Processing: Industry Track, pages 1393–1412, 2024.
- [41] Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. C-pack: Packaged resources to advance general chinese embedding, 2023.
- [42] NovaSearch. Jasper and stella: distillation of sota embedding models. *arXiv preprint* arXiv:2404.01309, 2024.
- [43] Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation, 2024.
- [44] Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*, 2023.
- [45] Junhan Yang, Jiahe Wan, Yichen Yao, Wei Chu, Yinghui Xu, and Yuan Qi. inf-retriever-v1 (revision 5f469d7), 2025.
- [46] Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. Generative representational instruction tuning, 2024.
- [47] Luyu Gao, Yunyi Zhang, Jiawei Han, and Jamie Callan. Scaling deep contrastive learning batch size under memory limited setup. *arXiv* preprint arXiv:2101.06983, 2021.
- [48] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019.

A Technical Appendices and Supplementary Material

A.1 Additional Experimental Setup Details

For our contrastive pretraining, we use a batch size of 1024, enabled by gradient caching[47] to optimize GPU memory usage. The model is trained with a learning rate of 3×10^{-8} , temperature $\tau=0.07$, and weight decay of 0.01, using the AdamW[48] optimizer. Each query is paired with m=16 positive samples, randomly sampled from the dataset. If fewer than m positive samples are available, duplicates are used to pad the group. The maximum input length is set to 1024 tokens. The model is trained for 20 epochs on 8 NVIDIA A800 GPUs, taking approximately 70 hours to complete. This stage focuses primarily on problem-code representation alignment. The resulting model, which demonstrates strong performance on code-related retrieval tasks, is denoted as **CP-Retriever-Code**.

To improve the model's performance on problem-problem retrieval tasks, we perform a second-stage fine-tuning using a multi-task setup. We combine training data from the *Problem-to-Duplicate* and *Simplified-to-Full* tasks, downsampling each to approximately 1,000 examples to ensure balanced representation. The two subsets are then mixed to form the final fine-tuning dataset. This stage uses a batch size of 1, learning rate of 2×10^{-6} , and weight decay of 0.01. The model is trained for 1 epoch, which takes approximately 1 hour on a single NVIDIA A800 GPU. While this model, referred to as **CP-Retriever-Prob**, slightly underperforms CP-Retriever-Code on code-related tasks, it shows marked improvements in tasks involving problem-level semantic similarity.

For experiments conducted with the Qwen3-Embedding-4B model, we adjusted several hyperparameters to better accommodate its larger architecture. Specifically, the maximum input length was increased to 2048 tokens, the contrastive learning temperature was set to $\tau=0.05$ (compared to the model's default of 0.01), and the learning rate was set to 1×10^{-6} . The total training time scaled approximately proportionally with the increase in model parameters.

A.1.1 Ablation Study on Group-InfoNCE

Table 4: Performance comparison under different loss functions in contrastive learning.

Loss	m (Positives)	Data Aug	Text-to-Code	Code-to-Code
Base	/	/	39.60	68.05
InfoNCE	1	X	68.64	67.58
InfoNCE	1	✓	68.84	68.97
InfoNCE (multi-pos)	16	✓	66.44	66.62
Group-InfoNCE	4	X	69.00	69.18
Group-InfoNCE	4	✓	70.11	70.42
Group-InfoNCE	16	X	69.47	69.77
Group-InfoNCE	16	✓	70.40	70.59

Table 4 reports the performance of different loss functions and configurations on the Text-to-Code and Code-to-Code retrieval tasks. Here, InfoNCE corresponds to Equation 1, InfoNCE (multi-pos) refers to its multi-positive extension (Equation 2), and Group-InfoNCE denotes our proposed loss (Equation 3).

We observe the following:

- Data augmentation consistently improves performance across comparable configurations, highlighting its robustness.
- Naïvely increasing positive samples in InfoNCE (multi-pos) does not guarantee gains, suggesting that grouping structure is essential for effective multi-positive learning.
- Larger m values in Group-InfoNCE generally lead to higher accuracy, indicating that the model benefits from richer positive context.

We further validated the effectiveness of Group-InfoNCE by repeating multi-task training with the Qwen3-4B-Embedding model using the same hyperparameter settings as in our main experiments

Table 5: Multi-task training validation using Qwen3-4B-Embedding.

m	Problem-Level Finetuning	T2C	C2C	P2Dup	S2Full	Avg
No	No	62.68	65.15	57.23	89.31	68.59
1	No	67.94	66.31	56.21	90.23	70.17
4	No	68.30	66.20	56.29	90.24	70.26
16	No	68.54	66.49	56.43	90.07	70.38
16	Yes	65.85	70.18	71.58	95.02	75.66

(Table 5). Results show that fine-tuning on problem-level tasks slightly decreases Text-to-Code performance (from 68.54 to 65.85) but substantially improves problem-level tasks such as Problem-to-Duplicate and Simplified-to-Full. This consistent, minor trade-off supports our hypothesis that the embedding requirements of code-level and problem-level tasks are inherently different. Given this intrinsic conflict, we consider this small drop reasonable and do not pursue extensive hyperparameter optimization that might harm generalizability.

Overall, both ablation and multi-task analyses demonstrate that Group-InfoNCE effectively balances representation learning across heterogeneous tasks while maintaining strong cross-domain retrieval performance.

A.1.2 Ablation Study on Fine-Tuning Data Composition

Table 6: Ablation results on contrastive pretraining and fine-tuning data composition. **PCD** indicates whether the model was pretrained on CPRet-PCPCD using contrastive learning. **Dup** and **Simp** indicate the inclusion of Problem-to-Duplicate and Simplified-to-Full data during fine-tuning, respectively. The fourth column (**PCD**) denotes whether CPRet-PCPCD data was also included in the fine-tuning phase. T2C = Text-to-Code, C2C = Code-to-Code, P2Dup = Problem-to-Duplicate, S2Full = Simplified-to-Full.

T-1		T-2						
PCD	Dup	Simp	PCD	T2C	C2C	P2Dup	S2Full	Avg
X	X	Х	Х	39.60	68.05	45.26	86.43	59.84
X	1	✓	X	27.25	64.48	53.04	85.78	57.64
✓	X	X	X	70.40	70.59	38.68	81.45	65.28
1	1	✓	X	56.50	70.68	60.06	90.74	69.50
1	✓	X	X	60.90	70.46	57.84	89.13	69.58
1	X	✓	X	68.05	72.77	57.01	92.65	72.62
✓	X	✓	✓	71.15	71.00	41.77	87.75	67.92
✓	✓	✓	✓	70.58	70.76	50.81	88.27	70.11

Ablation Study Overview. Table 6 summarizes a comprehensive ablation study on two key factors in our training pipeline: (I) contrastive pretraining on CPRet-PCPCD, and (II) the composition of fine-tuning data for problem-level tasks.

Effectiveness of Contrastive Pretraining. Rows 3–8 show that models trained with contrastive pretraining significantly outperform those without it (rows 1–2) across all four tasks. Although pretraining alone underperforms on problem-level tasks like *P2Dup* and *S2Full*, models fine-tuned from this initialization consistently surpass those trained from scratch. This indicates that contrastive pretraining helps learn transferable representations. The initial gap on problem-level tasks likely stems from objective mismatch—Group-InfoNCE focuses on problem-code alignment rather than direct problem-to-problem similarity—but still provides a strong foundation for downstream fine-tuning.

Impact of Fine-Tuning Data Composition. Rows 4–8 explore different combinations of *Problem-to-Duplicate*, *Simplified-to-Full*, and optionally CPRet-PCPCD data during fine-tuning. Using only *Simplified-to-Full* (row 6) yields the best overall performance, improving *C2C* and *S2Full*, while maintaining strong results on *P2Dup*. In contrast, adding *P2Dup* data (row 4) improves that task but

slightly reduces generalization elsewhere. Including CPRet-PCPCD during fine-tuning (row 7, 8) helps recover performance on code-centric tasks but slightly weakens results on problem-level tasks.

Conclusion. These results confirm the value of task-specific optimization. We therefore release two final models:

- CPRetriever-Code (row 3): trained only with contrastive pretraining, optimized for codecentric tasks.
- **CPRetriever-Prob** (row 4): further fine-tuned on problem-level tasks, optimized for problem retrieval.

A.2 Supplementary Analysis for Figure 2

Figure 2 in the main paper presents two complementary perspectives on model performance trends. The left sub-plot focuses on illustrating the overall improvement of embedding models over time, highlighting representative models that achieved state-of-the-art (SOTA) performance at their respective publication dates. In contrast, the right sub-plot was designed to examine potential data leakage effects by including more recent, high-performing models that are more likely to exhibit such phenomena. As a result, the model sets in the two sub-figures are not identical.

To provide a more comprehensive view, Table 7 below reports the year-wise retrieval performance of additional models (e.g., CodeSage-large, Contriever-MSMarco, and E5-Base-v2) that were not included in the right sub-plot of Figure 2. Despite having lower absolute performance compared to the models shown in the right plot, these models display a consistent "decline-then-recovery" trajectory: strong performance in earlier years (e.g., 2015–2016), followed by a sharp degradation (2017–2019), and a partial recovery or stabilization in later years (2020–2024).

This trend mirrors the temporal performance decay pattern attributed to data leakage, as discussed in Section 3.2.1. Specifically, the observed recovery phase suggests that recently released problems are less contaminated by training data, leading to more realistic evaluations. Therefore, these supplementary results reinforce our main claim that data leakage can substantially inflate retrieval metrics in temporally overlapping benchmarks, and that rigorous temporal separation is critical for fair assessment.

Table 7: Supplementary results supporting Figure 2. Retrieval performance (%) of representative models across problem release years. The "decline-then-recovery" pattern indicates potential data leakage in earlier problems.

Model	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024
SFR-EmbCode-2B	55.38	57.67	48.79	44.78	42.55	41.59	37.42	34.38	34.08	34.78
GTE-ModernBERT-Base	42.10	45.92	28.44	15.45	10.34	12.68	9.44	9.82	9.31	9.98
CodeSage-Large	40.65	43.57	27.71	14.67	9.10	11.76	9.33	7.67	7.21	8.08
Contriever-MSMarco	16.03	19.85	10.18	3.72	1.62	1.98	0.95	0.88	0.79	0.91
E5-Base-v2	22.35	26.15	14.02	5.16	2.27	2.59	1.69	1.32	1.18	1.29

A.3 Evaluation of Retrieval for RAG-based Competitive Programming

To investigate the impact of problem retrieval on downstream code generation tasks, we conducted experiments using Retrieval-Augmented Generation (RAG) in competitive programming problemsolving. We used problems from LiveCodeBench as the test set. To ensure temporal fairness, the retrieval pool was restricted to problems in our dataset published before April 2023, while LiveCodeBench problems appeared starting May 2023.

For each test problem, we retrieved the top-K most similar problems from the historical subset and provided their descriptions along with one accepted solution as context to the LLM. Three retrieval settings were compared: no retrieval (baseline), using Qwen3-Embedding-4B as the retriever, and using our **CPRetriever-Prob** (based on Qwen3-Embedding-4B).

For the code generation evaluation, we used three models from the Qwen-Coder series:

- qwen-coder-turbo and qwen-coder-plus: non-reasoning models,
- qwen3-coder-plus: a reasoning-capable model.

The results are summarized in Table 8.

Table 8: Pass@1 performance for RAG-based code generation with different retrievers.

Model	Type	RAG Model	Pass@1
qwen-coder-turbo	Non-reasoning	None	0.3575
qwen-coder-turbo	Non-reasoning	Qwen3-Embedding-4B	0.3825
qwen-coder-turbo	Non-reasoning	CPRetriever-Prob	0.4075
qwen-coder-plus	Non-reasoning	None	0.5125
qwen-coder-plus	Non-reasoning	Qwen3-Embedding-4B	0.5275
qwen-coder-plus	Non-reasoning	CPRetriever-Prob	0.5500
qwen3-coder-plus	Reasoning	None	0.7675
qwen3-coder-plus	Reasoning	Qwen3-Embedding-4B	0.7750
qwen3-coder-plus	Reasoning	CPRetriever-Prob	0.7800

As shown, incorporating RAG improves Pass@1 across all evaluated models. Non-reasoning models benefit more noticeably, while reasoning models achieve smaller gains due to their strong baseline performance. Notably, using **CPRetriever-Prob** consistently outperforms Qwen3-Embedding-4B across all LLMs.

These results demonstrate that effective problem retrieval can enhance competitive programming problem-solving in a RAG setting, highlighting a valuable downstream application for our retrieval model.

A.4 Challenges in Competitive Programming Problem Retrieval

Retrieving competitive programming (CP) problems poses distinct challenges compared to standard code retrieval tasks in software engineering (SE) benchmarks such as SWE-Bench. The core difficulty lies in bridging the large semantic gap between a problem's narrative description and its abstract algorithmic solution, whereas SE tasks typically involve more concrete, context-specific code retrieval.

CP problem statements are often highly abstract and indirect. They are framed as puzzles embedded within narratives (e.g., "a farmer needing a route") and require deep semantic understanding to map natural language to algorithmic concepts like graph traversal or dynamic programming. Critical information such as data constraints ($N \leq 1,000,000$) or time limits is subtly embedded, yet it fundamentally dictates the algorithmic complexity required ($O(N^2)$ vs. $O(N \log N)$). Similarly, user queries are abstract (e.g., "longest increasing subsequence sum"), in contrast to SE queries, which are typically specific and grounded in concrete code modifications or bug fixes.

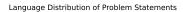
The mapping from a CP problem to its solutions is also more complex. A single problem may admit multiple fundamentally different correct algorithmic approaches (e.g., BFS/DFS vs. dynamic programming for tree diameter). CP retrieval aims to uncover the underlying idea or technique that can inspire a solution strategy, rather than locating an exact code snippet. In SE retrieval, the goal is generally to find specific functions, modules, or API calls that can be directly reused in a known context.

In summary, CP retrieval resembles a mathematician exploring theorems for similar logical structures while ignoring superficial details, focusing on abstract strategies. SE retrieval, in contrast, is more like a mechanic finding the right replacement part for a known machine. This higher-level abstraction and the need to translate between human language and algorithmic logic make competitive programming problem retrieval uniquely challenging.

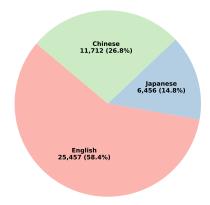
A.5 Further Details on the CPRet-PCPCD

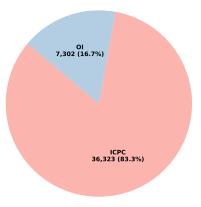
A.5.1 Language and Problem Format Statistics

Figure 5a shows that English is the most prevalent language, largely due to the influence of prior datasets like *TACO*, while Chinese and Japanese examples mainly come from LibreOJ/Nowcoder and AtCoder. Figure 5b reveals that most problems adopt the ICPC-style full-score format, but a notable subset includes OI-style problems that offer partial scores through sub-tasks or test case breakdowns.









- (a) Languages used in problem statements.
- (b) ICPC-style vs OI-style grading.

Figure 5: Dataset-level characteristics by language and format in CPRet-PCPCD.

A.5.2 Programming Language Distribution

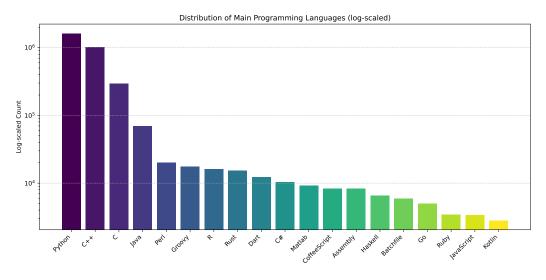


Figure 6: Distribution of programming languages in our dataset CPRet-PCPCD (log-scaled). Python is dominant due to its prevalence in prior datasets like APPS, CodeContests, and TACO. In contrast, our newly collected data is mainly in C/C++, which reflects the practical language choices in competitive programming.

As shown in Figure 6, our dataset covers a wide range of programming languages. Early datasets such as *APPS*, *CodeContests*, and *TACO* focused primarily on Python, driven by the popularity of Python in code generation research and the convenience it offers for model training. However, in our newly collected dataset, the dominant languages shift toward **C** and **C++**, which are widely favored in actual competitive programming due to their execution efficiency and low-level control.

Other languages such as **Java**, **Perl**, **Groovy**, and **Rust** also appear in notable quantities, indicating broader diversity and coverage in our collection.

A.6 Data Source Details

Data Copyright Notice. All data collected in this work are publicly accessible and do not involve any private or user-specific information. However, the copyright of the original content remains with the respective platforms from which the data were obtained. We use these data solely for research purposes, in accordance with the terms of use and fair academic practice.

A.6.1 Competitive Programming Problems and Codes Data(CPRet-PCPCD)

Existing Sources. We incorporate the following publicly available datasets:

- APPS [8]: https://github.com/hendrycks/apps
- CodeContests [9]: https://github.com/google-deepmind/code_contests
- TACO [10]: https://huggingface.co/datasets/BAAI/TACO
- Codeforces Submissions : https://www.kaggle.com/datasets/ yeoyunsianggeremie/codeforces-code-dataset
- Codechef Submissions : https://www.kaggle.com/datasets/arjoonn/codechef-competitive-programming
- Codeforces Source Code : https://www.kaggle.com/datasets/agrigorev/codeforces-code

Newly Collected Sources. To enrich the dataset, we additionally collect problems and code from the following platforms:

- Codeforces(English, ICPC): https://codeforces.com/, with API reference at https://codeforces.com/apiHelp
- AtCoder(English and Japanese, ICPC + OI): https://atcoder.jp/, with API reference at https://github.com/kenkoooo/AtCoderProblems/blob/master/doc/api.md
- LibreOJ(Chinese, ICPC + OI): https://loj.ac/, with API reference at https://api.loj.ac/
- Nowcoder(Chinese, ICPC + OI): https://ac.nowcoder.com/

For these new sources, we typically first obtain problem and submission metadata via the platform's API (when available), and then download actual code submissions accordingly. In cases where APIs are not provided, we resort to web crawling.

Crawler Tools. We utilize the following libraries to implement our crawlers:

- curl_cffi: https://github.com/lexiforest/curl_cffi
- crawl4ai: https://github.com/unclecode/crawl4ai

To minimize load on the original platforms and avoid abuse, we do not publicly release our custom crawlers. For academic collaboration or access, please contact us via email.

A.6.2 Duplicate Programming Problem Pairs

To collect semantically similar or duplicate programming problem pairs, we crawled discussion forum data from the following two competitive programming platforms:

- Codeforces, via blog entries such as https://codeforces.com/blog/entry/142637
- Luogu, via forum posts such as https://www.luogu.com.cn/discuss?forum=P1000

From Codeforces, we gathered approximately 130,000 blog posts, each containing an average of 12 user comments. From Luogu, we collected around 236,000 discussion comments.

To extract high-quality duplicate problem pairs, we use a multi-stage filtering process combining:

- **Keyword filtering** to identify relevant discussions,
- · Large Language Model (LLM)-based scoring for candidate selection, and
- Manual annotation to ensure correctness and semantic equivalence.

Note. For each problem pair mentioned in the forums, both problems may originate from the same platform, or from different platforms. In most cases, at least one problem in the pair comes from the platform where the discussion was found.

A.6.3 Simplified and Full Problem Description Pairs

This dataset is collected entirely from the Luogu platform (https://www.luogu.com.cn/), where users have contributed simplified versions of competitive programming problems originally published on international online judges such as Codeforces (https://codeforces.com/), AtCoder (https://atcoder.jp/), SPOJ (https://www.spoj.com/), and UVA (https://onlinejudge.org/). These problems may have been originally written in English, Japanese, or other languages. The simplified versions are often written in Chinese and aim not only to translate the content, but also to restructure and clarify problem statements by removing less essential narrative elements, rephrasing complex sentences, and highlighting key constraints and requirements. This simplification process helps make the problems more accessible to beginners.

We crawled these user-generated simplifications and applied filtering steps to remove low-quality entries, such as those that are incomplete, inconsistent with the original, or overly literal machine translations. To ensure consistency across the dataset, we use *Qwen-2.5-Max*[34] to translate Chinese simplified statements into English before inclusion. The resulting dataset maintains a high level of clarity and fidelity, and serves as a valuable resource for studying cross-lingual and cross-abstraction retrieval.

A.7 Extended Results: Model Performance across Similarity Bins

Figure 7 expands on the similarity-based analysis in Section 4.3, showing average pass rates across max similarity bins for additional models not included in the main paper.

The observed trends further support the conclusion that higher retrieval similarity strongly correlates with code generation success, demonstrating consistency across different model scales and types.

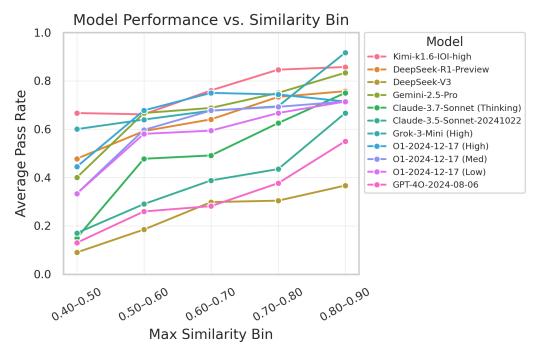


Figure 7: Extended results on model performance vs. similarity.

A.8 Cases of Retrieval Tasks

A.8.1 Text-to-Code and Code-to-Code

We next present a problem along with two correct solutions. In the Text-to-Code task, the goal is to retrieve these solutions given the problem description. In the Code-to-Code task, the goal is to retrieve the other solution given one of them, which may differ in language or implementation.

https://atcoder.jp/contests/arc162/tasks/arc162 a Problem Statement There are N people, numbered from 1 to N, who participated in a round-trip race between two points. The following information is recorded about this race. • The outward times of any two people were different, and person i $(1 \leq i \leq$ N) had the i-th fastest outward time. • The round-trip times (the sum of the outward and return times) of any two people were different, and person i $(1 \leq i \leq N)$ had the P_i -th fastest round-trip time. • The person (or persons) with the fastest return time was awarded the fastest return award. Here, P_1, P_2, \dots, P_N is a permutation of $1, 2, \dots, N$. How many people could have received the fastest return award? There are T test cases. Answer each of them. Constraints • $1 \le T \le 500$ • $2 < N < 10^3$ • P_1, P_2, \ldots, P_N is a permutation of $1, 2, \ldots, N$ • All input values are integers ullet The sum of N over all test cases does not exceed 10^3 The input is given from standard input in the following format:

Listing 1: Solution Code 1(C++)

Print T lines. The i-th line $(1 \leq i \leq T)$ should contain the answer for the i-th

Each test case i $(1 \le i \le T)$ is given in the following format:

case_1
...
case_T

P_1 P_2 ... P_N
Output Format

test case.

```
#include <bits/stdc++.h>
#include <cstdlib>
#include <math.h>
using namespace std;
int main()
    int t;
    cin >> t;
    while(t--){
        int n;
        cin >> n;
        vector < int > a(n);
        vector < int > b(n);
        for(int i=0;i<n;i++){</pre>
             cin >> a[i];
             a[i]--;
             b[a[i]]=i;
        }
         int ans=0;
        for(int i=0;i<n;i++){</pre>
```

Listing 2: Solution Code 2(Ruby)

```
t = read_line.to_i

def solve
  n = read_line.to_i
  p = read_line.split.map { |x| x.to_i - 1 }

  n.times.count do |i|
    n.times.all? do |j|
    # i < j => p[i] < p[j]
    !(i < j) || p[i] < p[j]
    end
  end
end

t.times do
  puts solve
end</pre>
```

A.8.2 Problem-to-Duplicate

The only difference between the two problems is that one includes multiple test cases, while the other does not. Aside from this, the problems are nearly identical, and are therefore classified as a Near Match.

```
https://codeforces.com/contest/1702/problem/G2
Problem: G2. Passable Paths (hard version)
You are given a tree with n vertices numbered from 1 to n. A tree is an
undirected connected graph with n-1 edges and no cycles.
You are also given q queries. In each query, you are given a set of k vertices.
For each query, determine whether there exists a simple path in the tree that
passes through all the given \boldsymbol{k} vertices.
Input
The first line contains an integer n (2 \le n \le 2 \times 10^5) - the number of vertices in
Each of the next n-1 lines contains two integers u and v (1 \leq u,v \leq n) - the
endpoints of an edge in the tree.
The next line contains an integer q (1 \le q \le 2 \times 10^5) - the number of queries.
Each of the next q lines describes a query:
      • The first integer k (2 \le k \le n) - the number of vertices in the query.
      • Followed by k integers v_1, v_2, \dots, v_k (1 \leq v_i \leq n) - the vertices in the
        query.
Output
For each query, print YES if there exists a simple path that passes through all
the given k vertices. Otherwise, print NO.
```

https://www.codechef.com/JULY21A/problems/KPATHQRY

Problem: KPATHQRY - Path Queries on Trees

You're given a tree with N vertices numbered from 1 to N. Your goal is to handle Q queries. For each query, you are given K nodes v_1, v_2, \ldots, v_K . Find whether there exists a simple path in the tree that covers all the given vertices.

Input

- ullet The first line contains a single integer T the number of test cases.
- For each test case:
 - The first line contains a single integer N the number of vertices.
 - Each of the following $N\,-\,1$ lines contains two integers u and v an edge in the tree.
 - The next line contains a single integer ${\it Q}$ the number of queries.
 - Each of the following Q lines describes a query:
 - st Starts with K_i the number of vertices in the query.
 - * Followed by K_i integers: $v_1, v_2, \ldots, v_{K_i}$.

Output

For each query, print "YES" if a simple path covering all the given nodes exists, otherwise print "NO".

You may print each character of the string in uppercase or lowercase (for example, "yEs", "yes", "Yes", and "YES" will all be treated as identical). Constraints

- $1 \le T \le 10$
- $1 < N < 10^5$
- $1 \leq u, v, v_j \leq N$
- $1 \le Q \le 10^5$
- $1 < K_i < N$
- The edges form a valid tree.
- All vertices in a single query are distinct.
- The sum of N over all test cases does not exceed $2\cdot 10^5$.
- The sum of K_i over all queries in a single test case does not exceed $10^5\,.$

Subtasks

• Subtask #1 (100 points): Original constraints.

A.8.3 Simplified-to-Full

The first paragraph below is a simplified version of the second paragraph. The simplification removes irrelevant background information and input/output format details, preserving only the core problem statement.

https://codeforces.com/problemset/problem/1183/F Concise problem description

You have n numbers from a_1 to a_n , and you want to select at most 3 numbers such that no number is a multiple of another. You aim to maximize the sum of the selected numbers. Output this maximum sum.

https://codeforces.com/problemset/problem/1183/F Full problem description

Topforces Strikes Back Problem Description

- One important contest will take place on the most famous programming platform (Topforces) very soon!
- The authors have a pool of n problems and should choose at most three of them into this contest. The prettiness of the i-th problem is a_i . The authors have to compose the most pretty contest (in other words, the cumulative prettinesses of chosen problems should be maximum possible).
- But there is one important thing in the contest preparation: because of some superstitions of authors, the prettinesses of problems cannot divide each other. In other words, if the prettinesses of chosen problems are x,y,z, then x should be divisible by neither y, nor z, y should be divisible by neither x, nor z and z should be divisible by neither x, nor y. If the prettinesses of chosen problems are x and y then neither x should be divisible by y nor y should be divisible by x. Any contest composed from one problem is considered good.
- Your task is to find out the maximum possible total prettiness of the contest composed of at most three problems from the given pool.
- ullet You have to answer q independent queries.
- If you are Python programmer, consider using PyPy instead of Python when you submit your code.

Input Format

- The first line of the input contains one integer q ($1 \leq q \leq 2 \cdot 10^5$) the number of queries.
- The first line of the query contains one integer n ($1 \leq n \leq 2 \cdot 10^5$) the number of problems.
- The second line of the query contains n integers a_1,a_2,\ldots,a_n ($2\leq a_i\leq 2\cdot 10^5$), where a_i is the prettiness of the i -th problem.
- It is guaranteed that the sum of n over all queries does not exceed $2\cdot 10^5$.

Output Format

• For each query print one integer - the maximum possible cumulative prettiness of the contest composed of at most three problems from the given pool of problems in the query.

A.9 Online Demo and Test Cases of CPRetriever

We deployed an open-source competitive programming problem retrieval platform, **CPRetriever**, on May 21, 2025 (https://cpret.online/). Within the first week, the platform processed nearly 2,000 search queries, and the announcement post on Codeforces received over 250 upvotes, indicating strong community interest.

The platform supports two primary retrieval functionalities: (i) **similar problem retrieval**, which assists users in expanding problem-solving perspectives and identifying knowledge gaps, and (ii) **duplicate problem detection**, which aids problem setters in identifying previously seen ideas or solutions.

A.9.1 Duplicate Problem Retrieval in a Recent Contest

We evaluated CPRetriever on the 2025 CCPC National Invitational Contest (Northeast), which featured 12 problems (https://codeforces.com/gym/105924). The system successfully identified six problems with highly similar or identical historical counterparts. Manual inspection of the top three retrievals per query suggests a **minimum duplicate rate of 50%**, highlighting potential fairness concerns in contest scoring.

Table 9: Detected duplicates in the 2025 CCPC Northeast Contest.

Contest Problem	Matched Historical Problem	Similarity Level	Rank
A. GD Ultimate Rhythm Lab	Nowcoder - Xiao Rui Rui's Sequence	Same approach	1
D. Defend the Carrot	SPOJ - UOFTBB	Almost identical	1
E. Tree Edge Removal	Luogu - [JRKSJ R7] Stem	Almost identical	1
F. Youthful Oath II	Codeforces - 80B Depression	Almost identical	1
J. Kingdom: Memories	AtCoder - R Walk	Almost identical	3
L. Bathhouse	Codeforces - 219E Parking Lot	Same approach	2

A.9.2 Similar Problem Retrieval: MEX Variants

We further evaluated the system using the classic "interval MEX" problem to retrieve its variants across multiple contests, demonstrating CPRetriever's utility for **idea exploration and knowledge transfer**.

Table 10: Top retrieval results for an interval MEX query.

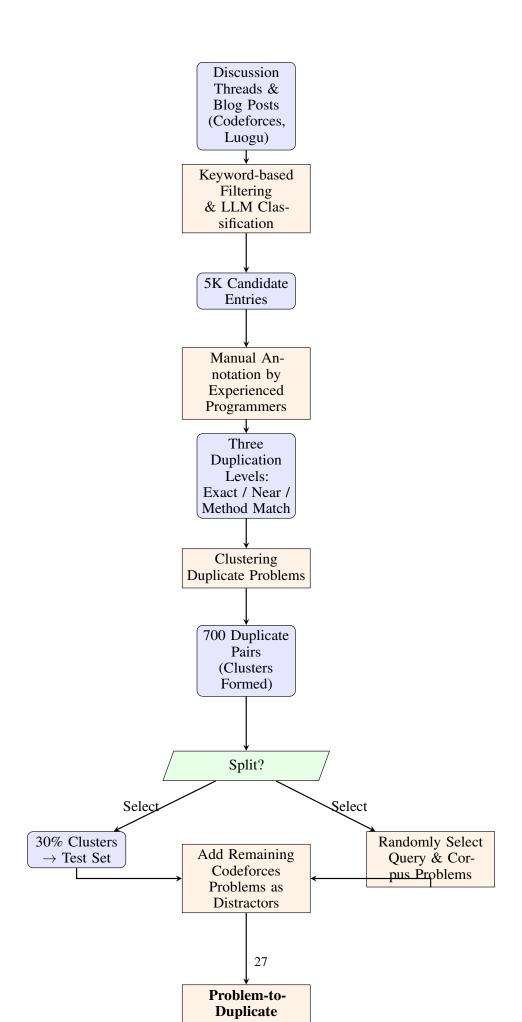
Rank	Problem	Description
1	Luogu P4137: RMQ Problem / MEX	Original problem
2	LOJ 6908: THUPC 2024 Prelim - "Matryoshka"	MEX of all subarrays of length k , then take MEX
5	AtCoder ABC194E: Mex Min	MEX of all subarrays of length k , then take minimum
6	Luogu P10032: Mex of Sequence	Repeated operations: $a'[i] = \max(a \setminus a[i])$
11	Nowcoder 237670: Classic Problem	MEX queries on permutations, optimized to $O(n+m)$
14	Luogu P8087: JROI-5 Interval	MEX of all subarrays of length k , then take maximum
15	AtCoder ABC290C: Max MEX	MEX of all subsequences of length k , then take minimum
16	Codeforces 1436E: Complicated Computations	MEX of all subarrays, then take MEX again
23	AtCoder ABC330E: Mex and Update	Supports element modification or full-array MEX queries
24	Luogu P11837: Making Mexes B	Minimum edits to ensure $mex(a) = i$

These case studies demonstrate that CPRetriever effectively identifies both **duplicates** and **semantically similar problems**, supporting practical contest preparation, knowledge expansion, and fairness monitoring in competitive programming environments.

A.10 Figures of data collection and processing pipeline



Figure 8: Overall construction pipeline of the **CPRet-PCPCD** dataset. Problems and accepted solutions are collected from multiple online judges, paired and enriched with metadata, filtered for quality, and annotated with timestamps to support temporally-aware retrieval research.



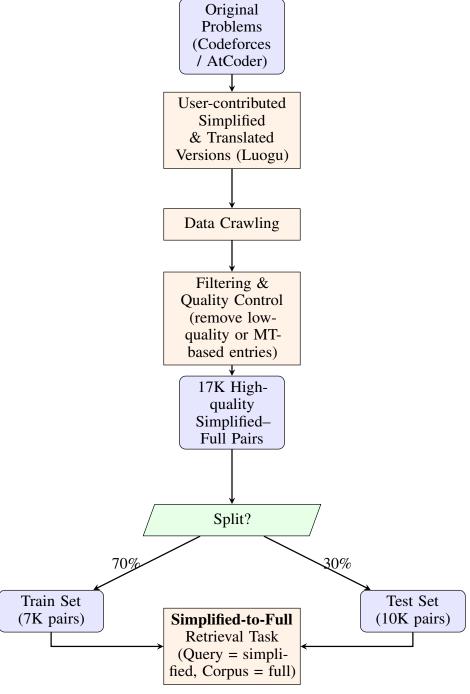


Figure 10: Data collection and processing pipeline for the **Simplified-to-Full** retrieval task. Problems are simplified or translated by Luogu users, then crawled, filtered for quality, and split into training and test sets.