API Reranking for Automatic Code Completion: Leveraging Explicit Intent and Implicit Cues from Code Context

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have significantly advanced software development, particularly in automatic code completion, where selecting suitable APIs from vast third-party libraries plays a critical role. However, current solutions either focus on recommending APIs based on user queries or code context, without considering both aspects simultaneously. To bridge this gap, we propose a novel framework APIRANKER to rerank candidate API documents based on both the explicit developer intent and implicit cues in the incomplete code context. To generate training data for this task, we introduce a self-supervised ranking framework that automatically constructs data by assessing the relevance of API documents to code context with a perplexity-driven approach 017 via comments. To enhance API relevance detection, we propose a novel reranking model that predicts relevance scores by capturing a hidden reasoning state to estimate relevance. The experimental results show the effective-022 ness of our approach, in both recommending more accurate APIs and enhancing automatic code completion. The code is available¹ and the dataset will be released.

1 Introduction

034

The introduction of LLMs has led to advancements in automatic code completion (Husein et al., 2024), with recent models adopting the *retrievethen-generate* paradigm (Nashid et al., 2024). This approach enables LLMs to dynamically retrieve upto-date Application Programming Interface (API) information from documents, rather than relying solely on static training data.

A crucial aspect of the code generation process is selecting suitable API docs from massive amounts of third-party libraries. The choice of API not only determines the functionality of the generated code

¹https://anonymous.4open.science/r/ APIRanker-C442



Figure 1: Example of retrieval-augmented code completion with different retrieved top-k API documents.

but also affects its efficiency, maintainability, and overall integration with the existing software system (Wang et al., 2024b). Current research predominantly focuses on two approaches: 1) retrieving API docs based on user queries (query-based API recommendation) (Wu et al., 2023), which neglects the code context, and 2) completing user-written code based on the preceding code context (Peng et al., 2022), which fails to capture the developer's underlying intention behind the API usage. However, automating code completion with user preferable APIs requires a more comprehensive model that considers both the explicit intent conveyed by user queries and the implicit cues embedded in the code context.

Consider the following practical scenario shown in Fig. 1: Alice is a developer, she encounters a problem during her daily work, i.e., "convert the 'datatime' column to UTC timezone and sort it from earliest to latest for further analysis". She uses the LLM to complete her code. However, without knowing Alice's intention or the current code context, a large number of relevant

but unsuitable APIs may be recommended (*e.g.*, datetime.astimezone, numpy.sort, etc.). If the target APIs are not ranked in the top-k retrieval results, they will not be used for code completion, causing the auto-completed code to misalign with her intended behavior. Now consider Alice inputs both her intent and the current code context, the target APIs (*e.g.*, df.dt.tz_localize, pandas.DataFrame.sort_values) can be successfully retrieved, the LLM is likely to complete her code by smoothly integrating with the recommended API docs, effectively solving her problem.

063

064

065

077

094

097

101

103

104

106

107

108

110

111

112

113

114

Based on our previous observations, there is a need for completing code using the correct APIs. However, recommending appropriate APIs based on both the developer's intent and the incomplete code context is a challenging task: (i) Lack of code completion dataset annotated with relevant APIs hinders learning-based approaches. Curating training datasets for API recommendation requires manually evaluating the relevance of API documentation to the developer's requirements. This process depends on domain expertise, additionally, developers rarely express intent in their code, which makes large-scale data collection impractical. (ii) The relevance of APIs to the developer's requirements is hard to learn and capture. API documentation varies significantly in format, writing style, and level of detail across different third party libraries, making it difficult to directly link to developers' requirements. Furthermore, developers' requirements are expressed through implicit intent and subtle semantics in the incomplete code, recommending APIs that align with both aspects effectively is challenging.

To tackle the above challenges, we propose a novel framework named APIRANKER, which is designed to rerank candidate APIs by jointly considering the developer's intent and the incomplete code context. To address the challenge of lacking training data, we propose a self-supervised ranking framework to automatically construct ranking data. Specifically, we leverage a perplexity-driven relevance ranking approach, which uses LLM as an evaluator to automatically estimate the relevance of API documents to developer requirements by measuring the perplexity of completed code. To bridge the gap between perplexity and true semantic relevance, we employ a perplexity alignment strategy using code comments as semantic anchors. To better learn and capture the relevance of API documentation to the developer's requirements, we design a novel reranking model consisting of two key components. First, a hidden reasoning state extractor is designed to capture the relevance of the API documentation to implicit cues within code context, by extracting reasoning states from LLMs during inference. Second, a relevance estimator uses the reasoning states to explicitly predict a relevance score, learning to differentiate the impact of various API docs on code completion effectiveness. 115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

164

In summary, our paper makes the following contributions: (1) Joint consideration of developers' intent and code context for API recommendation. Prior research typically focuses on either the developer's intent or the code context in isolation. To the best of our knowledge, our work first thoroughly investigated API recommendations based on both aspects simultaneously. (2) A self-supervised ranking framework for ranking data construction. We introduce a novel self-supervised approach that generates ranking data in the form of *(incomplete code, target code, targe* API docs, relevance scores automatically, eliminating the need for manual annotation. (3) An API reranking model for better code completion. We design a novel reranking model that leverages LLMs' semantic understanding to rerank APIs based on the relevance of API documents to a developer's requirements. The experimental results show the effectiveness of our model over a set of baselines, showing its potential to enhance automatic code completion by reranking candidate API documents. We hope our study can lay the foundations for this research topic.

2 Related Work

API Recommendations. API recommendation methods typically rely on two main sources: natural language queries and contextual code information. Some studies focus on query intent, such as BIKER (Huang et al., 2018) and CLEAR (Wei et al., 2022), while others emphasize code context, like GAPI (Ling et al., 2021) and MEGA (Chen et al., 2023). Deep learning models like Deep-API (Gu et al., 2016) and CodeBERT (Feng et al., 2020) enhance recommendations through embedding-based methods, using pretrained models to calculate similarities between queries and API docs. However, limited labeled data hampers model performance (Ma et al., 2024). In contrast, our approach leverages LLMs and automatically generated data to reduce reliance on QA data, im165 proving API recommendation performance.

166 **Retrieval-augmented** Code Generation. Retrieval-augmented generation (Gao et al., 2023) 167 has proven valuable in code generation (Parvez 168 et al., 2021), especially as code libraries are frequently updated (Lu et al., 2022). For instance, 170 CodeGen4Libs (Liu et al., 2023) recommends 171 class-level API docs through a two-stage process 172 of retrieval and fine-tuning. DocPrompting (Zhou 173 et al., 2022) enables continuous updates to the 174 documentation pool, ensuring that the most 175 current code libraries are used for generation. 176 ToolCoder (Zhang et al., 2023b) integrates API 177 search tools and uses automated data annotation to 178 179 teach the model how to use tool usage information, thereby enhancing code generation.

3 Methodology

181

182

183

184

185

187

190

191

192

193

194

195

196

197

198

199

203

3.1 Task Definition

Given a query q, which contains natural language (NL) intent x and the corresponding incomplete code snippet c, the objective is to successfully complete code snippet c via recommending correct API documents D for code completion from a large corpus of candidates.

3.2 Self-supervised Ranking Framework

The absence of a code completion dataset annotated with relevant APIs makes training models a challenging task. This is primarily due to the high cost of manual annotation and the difficulties of verifying the correctness of generated code based on the APIs. To address this challenge, we propose a perplexity-driven relevance ranking approach, leveraging the perplexity of LLM-generated code to construct training data. However, perplexity can be influenced by factors such as code formatting and syntactic variations, introducing noise that distorts accurate relevance estimation. To mitigate this issue, we propose a perplexity alignment strategy that enriches the code context with comments, reducing the perplexity shifts and aligning code semantics.

Perplexity-driven Relevance Ranking. Evaluating the relevance of API documentation to preceding code context is a time-consuming process, requiring complex execution environments (Wei et al., 2023). These obstacles lead to a scarcity of training data, further limiting the learning-based approach's progress. To address this challenge,



Figure 2: Overview of the Self-supervised Ranking Framework.

we propose a perplexity-driven relevance ranking method, which assesses the relevance by measuring the perplexity of LLM-generated code.

Specifically, as illustrated in Fig. 2, we construct data from GitHub repositories², API documents, and an LLM as a perplexity evaluator. For a specific code file that has cross-file dependencies, we randomly select a middle position to split it into incomplete code c and the target code y, ensuring ample context for retrieval and completion. We directly use the incomplete code c as query q and retrieve the top n API documents $D = \{d_1, d_2, ..., d_n\}$ as candidates via the retrieval model. We use dependency analysis tool³ to identify direct and indirect dependency files of the code file, then each dependency file is regarded as a potential API document for user-defined functions, as it contains both the detailed implementations and definitions of those functions. Including these dependency files ensures that query q is paired with relevant API documents.

For each API document $d \in D$, the perplexity (PPL) of the target code y is defined as:

$$PPL(y|d,q) = e^{-\frac{1}{N}\sum_{i=1}^{N} \log P(y_i|d,q,y_{$$

where P represents the probability distribution over the LLM's vocabulary, and N is the number of tokens in the target code y. The relevance score r between API document d and query q is then defined by the perplexity of the target code y as:

$$\mathbf{r}(d,q) = \frac{1}{\mathbf{PPL}(y|d,q)}.$$
 (2)

Using the relevance score r, we can compare the relevance of different API docs for the same query, since lower perplexity indicates that LLM has less difficulty in correctly completing the code with the API docs. A higher value of r indicates a higher relevance of the document to the query.

³https://github.com/IBM/import-tracker, https: //maven.apache.org

242

243

245

246

²https://github.com



Figure 3: Overview of APIRANKER. (a) is the training process of APIRANKER based on the collection of different API documents D_r and the same incomplete code c as the query. (b) illustrates the structure of the hidden reasoning state extractor. (c) illustrates the structure of the relevance estimator.

Perplexity Alignment via Anchor Comments. The perplexity score used for relevance estimation is intended to better reflect code semantics, while minimizing the influence of non-semantic factors, such as formatting variations (e.g., line breaks, indentation) and code syntactic variations (e.g., bracket placement, variable declaration). These variations can cause significant shifts in perplexity, making it unreliable to reflect the true relevance of the API document and the code context. To address this issue, we incorporate a strategy of perplexity alignment via adding anchor comments into the target code. Anchor comments provide semantic information that enhances the logic representation during perplexity calculation, thereby reducing the sensitivity of perplexity shifts and enhancing the alignment of perplexity and logic relevance.

248

249

250

251

255

257

261

262

263

265

267

269

271

272

Specifically, given an incomplete code c and the corresponding target code y, LLM is asked to add comments to each line of code y, as described by the following equations:

$$\hat{y} = \text{LLM}(c, y), \tag{3}$$

where \hat{y} is the generated code with comments. The prompt of perplexity alignment guided by anchor comments is then constructed as Fig. 4.

273Based on the target code with comments \hat{y} , the274relevance score r is calculated by Equation 2, of-275fering a more accurate measure of the relevance276between the API document and the query.

• Instruction: Based on the following two consecutive parts of the same code, Part A (the first half) and Part B (the second half), both enclosed within <code> and </code> tags, you should add comments to each line of code in Part B as much as you can.

- Part A (the first half): c
- \bullet Part B (the second half): y

Figure 4: Prompt of perplexity alignment.

278

279

281

282

287

290

291

294

297

299

3.3 API Reranking Model Architecture

Inspired by LLMs' strong abilities (Naveed et al., 2023) in natural language comprehension, we incorporate LLM into our model architecture to capture user intents. This design allows us to bypass the need for modeling or learning from natural language intent. However, LLMs still struggle to capture implicit cues in code context and differentiate between different API documents for the same query. To tackle this issue, we propose a novel reranking model architecture, APIRANKER, for selecting the most suitable APIs from a set of candidates. APIRANKER includes a hidden reasoning state extractor that leverages the reasoning state to capture the relevance of the API document to implicit cues from the code context. Additionally, a relevance estimator is employed to detect the reasoning state and explicitly predict a relevance score between the API document and the query.

Hidden Reasoning State Extractor. A large number of tokens in natural language generation are produced solely for fluency, contributing little to the underlying reasoning process. Inspired by

349 350 351

352

353

355

356

357

358

359

360

361

363

364

365

366

368

370

371

372

373

374

375

376

378

379

381

383

347

348

the previous studies on hidden reasoning state (Hao et al., 2024; Ouyang et al., 2022), we extract the representation of the reasoning state from the last hidden state of the LLM, allowing us to capture semantic relevance rather than relying on general tokens for linguistic coherence.

300

301

306

310

311

313

314

315

317

319

321

325

326

328

329

Specifically, as illustrated in Fig. 3, given a query q (*i.e.*, the incomplete code c) and an API document d, we prompt the LLM to perform code completion based on d and extract the sequence of hidden states through the decoder layer of LLMs as:

$$h = \text{DecoderLayer}(d, q).$$
 (4)

where $h = \{h_1, h_2, ..., h_m\}$ represents the sequence of hidden states, m is the number of hidden states. During this process, the LLM's parameters are kept frozen. To align the dimensions between the LLM and the state extractor, we introduce a linear layer as:

$$h' = W_s * h + b_s, \tag{5}$$

where h' is the hidden states after aligning, W_* and b_* denote the trainable parameters in this section. Each layer of the state extractor consists of self-attention, cross-attention, and a feed-forward network (FFN) followed by layer normalization as:

$$p' =$$
SelfAttention $(p, p, p),$ (6)

$$p'' = \text{CrossAttention}(p', h', h'),$$
 (7)

s = LayerNorm(FFN(p'') + p''), (8)

where p denotes a set of learnable vectors used to capture the reasoning states and s represents the sequence of reasoning states, which serves as input for the next layer of the state extractor. We initialize the extractor with transformer weights pre-trained on code data, whereas the cross-attention layers are randomly initialized.

334**Relevance Estimator.** To assess relevance from335the reasoning states, we propose a relevance esti-336mator that aggregates semantic information from337the reasoning states and predicts relevance scores.338Specifically, as illustrated in Fig. 3(c), we use a339learnable query vector with multi-head attention,340where reasoning state s serves as both the key and341value of attention. The final relevance score is then342predicted by a neural network, as described by the343following equations:

$$g' = \text{LayerNorm}(\text{MHA}(g, s, s)),$$
 (9)

345
$$g'' = \text{LayerNorm}(g' + \text{FFN}(g')), \quad (10)$$

346
$$\hat{r} = W_r * g'' + b_r,$$
 (11)

where g is a learnable vector, representing the relevance of states from the last reasoning states s of state extractor, MHA denotes multi-head attention, and \hat{r} represents the predicted relevance score.

3.4 Training and Inference

Training Objective. APIRANKER is trained on a dataset consisting of pairwise comparisons between different API candidates for the same query. As illustrated in Fig. 3, we use a cross-entropy loss, where each training sample is labeled by performing comparisons between API document pairs. The difference in rewards represents the log odds of one document being preferred over the other, with this preference determined by the relevance function r(Section 3.2). To speed up comparison training, we construct pairs from a set of K documents selected evenly based on the difference in r values, chosen from the top n candidate documents, and train on all comparisons for each query as a single batch. Formally, the training objective of the reranking model is defined as:

$$\mathcal{L} = -\frac{1}{\binom{K}{2}} \mathbb{E}_{(q,\hat{y},d_w,d_l)\sim\mathcal{D}_r}$$

$$[\log \sigma(\hat{r}(q,\hat{y},d_w) - \hat{r}(q,\hat{y},d_l))],$$
(12)

where σ denotes the logistic function, $\hat{r}(q, \hat{y}, d)$ is the scalar output of the reranking model for query q, target code with comments \hat{y} and API document $d. d_w$ is the preferred document out of the pair of d_w and d_l , and D_r is the training data based on score of relevance function r.

Inference. During the inference stage, given a set of candidate documents D retrieved by the retrieval model based on query q (*i.e.*, NL intent and incomplete code), each document $d \in D$ is evaluated by APIRANKER, which produces a new ordering of the candidate documents based on the relevance between document and the query.

4 **Experiments**

4.1 Experimental Setup

Dataset.To study how API recommendation384benefits the automatic code completion task, we385construct a dataset APIRAC (API Retrieval-386Augmented Completion) for this task.We col-lect 110,646 API docs from the dataset CodeRAG-388bench (Wang et al., 2024c) as retrieval sources.389Additionally, we gather 4,400 large-scale reposito-390ries from GitHub, based on the dataset presented391

in the RLCoder (Wang et al., 2024a), with an equal number of Python and Java repositories, and split 393 them into training and validation sets with a 10:1 ratio. For a given code file, we add its associated dependency files to the retrieval sources as API candidates. Finally, we construct (incomplete code, target code, API docs, relevance scores training data using our self-supervised ranking framework. For the test data, we select DS-1000 (Lai et al., 2023) 400 as the automatic code completion dataset, which 401 includes general open-domain coding completion 402 tasks. We use the human-annotated API documen-403 tation for DS-1000, provided by CodeRAG-bench, 404 as a dataset for API recommendation. There is no 405 overlap between the ground truth documents from 406 test dataset and the APIs used in training. Overall 407 statistics of the dataset are given in Table 1. Further 408 details can be found in Appendix A.1. 409

Baselines. We consider the following retrieval baselines: (1) Unixcoder: A cross-modal pretrained model for programming language (Guo et al., 2022). (2) GIST-large: A model that enhances text embedding fine-tuning by selecting negative samples (Solatorio, 2024). (3) Arctic-Embed 2.0: A multilingual text embedding model for retrieval (Yu et al., 2024). (4) NV-Embed-v2: A embedding model ranked No.1 in the retrieval subcategory of the Massive Text Embedding leader-board (as of Aug 30, 2024) (Lee et al., 2024).

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437 438

439

440

441

442

Then we consider the following reranking methods: (1) **Unsupervised Passage Re-ranker (UPR)**: A pointwise approach based on query generation (Sachan et al., 2022). (2) **Relevance Generation (RG)**: A pointwise approach based on relevance generation (Liang et al., 2022). (3) **Pairwise Ranking Prompting- Sorting (PRP-Sorting)**: A pairwise method based on the log-likelihood of document generation, and it optimizes time complexity through heap sort (Qin et al., 2023). (4) **Pairwise Ranking Prompting-Sliding (PRP-Sliding)**: A variant of PRP based on the sliding window approach. (5)**RepoCoder**: A reranking method through iterative retrieval of code snippets based on code generation results (Zhang et al., 2023a).

In addition, we adopt two widely adopted API recommendations approaches in software engineering, **BIKER** (Huang et al., 2018) and **GAPI** (Ling et al., 2021), which consider either user intent or the code context, respectively. For automatic code completion, we consider the following code LLM: **Starcoder2-7B** (Lozhkov et al., 2024) and

Sets	Avg. query	Number canonical	Source	Avg. intent	Code Lines/ incomplete	Words target
train	4,000	-	Github	-	38.2	41.4
val	400	-	Github	-	37.9	40.5
test	513	1.4	Stackflow	84	10.7	5

Table 1: Dataset Statistics.

CodeLlama-Instruct-7B (Roziere et al., 2023), which are both trained and optimized for coderelated tasks. Further details on the above baselines can be found in Appendix A.2. 443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

Implementation Details. For API recommendation, we rerank the top 50 docs retrieved by different retrieval models. In our model, we chose CodeLlama-Instruct-7B as the LLM and Unixcoder as the initial weight of the hidden reasoning state extractor. Further details can be found in Appendix A.3.

Evaluation Metrics To evaluate the performance of API recommendation, we report the common evaluation metrics (Zhang et al., 2017; Wei et al., 2023): Recall@k, NDCG@k, MRR@k, and MAP, with k set to 10. We use Recall@k as the primary metric since retrieval-augmented generation primarily relies on key information that appears in the context. To evaluate the performance of code completion based on API recommendation, we adopt Pass@k and Improve@k metrics to measure the execution correctness of programs, with k set to 1. Further details can be found in Appendix A.4.

4.2 Experimental Results

API Recommendation Evaluation. Table 2 shows the experimental results of our approach and reranking baselines on API recommendation task. It is obvious that: (1) Regarding the Recall and overall ranking performance of recommending correct APIs, our approach APIRANKER outperforms all other reranking baselines by a large margin across different retrieval models, demonstrating substantial improvements in both the coverage and ranking quality of relevant documents. For example, APIRANKER achieves a Recall rate of 32.36% on Arctic-Embed 2.0, surpassing the next best method (i.e., UPR) by a significant margin of 11.72%. Similarly, in terms of NDCG@10, APIRANKER outperforms the next best method (*i.e.*, RG) with a value of 18.39%, surpassing it by a significant margin of 5.33%, which clearly indicates its superior reranking capability. (2) Our approach APIRANKER demonstrates consistent

Retrieval	Size	dim	Reranking Method	Recall@10	NDCG@10	MRR@10	MAP
Unixcoder BIKER GAPI	126M - -	768 - -		2.44 11.65 7.18	1.35 8.57 4.67	0.98 6.65 3.55	1.15 6.27 2.99
GIST-large	335M	1024	- RG UPR PRP-Sorting PRP-Sliding RepoCoder APIRANKER	15.25 15.69 <u>16.65</u> 7.00 12.65 11.09 25.50	6.88 <u>10.93</u> 9.56 2.35 10.15 4.15 14.52	3.87 9.10 7.01 0.87 <u>9.15</u> 5.91 10.33	4.79 9.42 7.66 2.42 <u>10.04</u> 7.21 10.83
Arctic-Embed 2.0	568M	1024	- RG UPR PRP-Sorting PRP-Sliding RepoCoder APIRANKER	18.86 19.98 <u>20.64</u> 8.45 12.51 17.67 32.36	10.83 <u>13.06</u> 12.00 2.80 11.52 8.18 18.39	7.72 10.30 8.38 1.12 <u>11.20</u> 7.27 12.48	8.71 10.84 9.22 3.11 <u>12.73</u> 7.66 13.09
NV-Embed-v2	7.9B	4096	RG UPR PRP-Sorting PRP-Sliding RepoCoder APIRANKER	27.12 22.53 25.53 14.98 23.59 <u>28.41</u> 30.17	13.65 <u>14.37</u> 14.34 4.76 13.33 13.72 15.49	8.76 11.30 <u>10.27</u> <u>1.82</u> 9.46 8.00 9.53	9.80 12.29 <u>11.37</u> 4.25 10.93 9.09 11.08

Table 2: Evaluation results on the APIRAC dataset. All results in the table are reported in percentage (%). The best method is in boldface, and the second best method is underlined for each metric.

Model	Complexity	Parameters	Train	Method Inference	principle
RG	O(N)	6.7B	-	Pointwise	Perplexity
UPR	O(N)	6.7B	-	Pointwise	Perplexity
PRP-Sorting	O(logN*N)	6.7B	-	Pairwise	Perplexity
PRP-Sliding	O(K*N)	6.7B	-	Pairwise	Perplexity
APIRANKER	O(N)	6.7B+160M	Pairwise	Pointwise	Semantics

Table 3: Comparison of the reranking method. N is the number of documents retrieved for reranking. K is the number of documents to be returned after reranking.

improvements across all evaluation metrics postreranking, irrespective of the underlying retrieval model. Even in the case of the strong retrieval baseline (e.g., NV-Embed-v2), where other methods all show degraded performance compared to the original retrieval results, APIRANKER still demonstrates stable improvements, outperforming retrieval baseline across all metrics, which highlights the effectiveness of our method in enhancing ranking quality and coverage in diverse retrieval scenarios. Overall, our method shows substantial and stable improvements in reranking different API candidates compared to other models, validating the effectiveness of our method for API recommendation. Further details on the above baselines can be found in Appendix A.5.

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

503

504

Methods Comparison and Analysis. As illustrated in Table 3, APIRanker demonstrates several key advantages over other reranking methods: (1) Its linear complexity O(N) ensures scalability, making it suitable for large-scale applications. (2) The pairwise training method improves its ability to learn relevance preference, while the pointwise inference ensures efficient inference. It achieves higher accuracy with just 160M additional parameters. (3) By leveraging semantic understanding of relevance estimator, APIRanker excels in tasks that require a deep comprehension of API documents in comparison to models that rely solely on perplexity. In the Appendix A.6, we provide further discussion on PPL and log-probability based uncertainty. 505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

Retrieval-augmented Completion Evaluation. Using a retrieval-augmented generation approach, we evaluated the performance of different reranking models in improving code completion performance on Arctic-Embed 2.0. The experimental results showed that: (1) As illustrated in Fig. 5(a), APIRANKER consistently outperforms the noretrieval baseline and leads to stable improvements in passing across both code LLMs, whereas other reranking models (*i.e.*, UPR, RG) cause performance degradation. This decline can be attributed to the noise of wrongly recommended APIs, which disrupts the code LLM's ability to generate code that was previously correct. In contrast, API-RANKER offers stable and reliable improvements



Figure 5: Effect of API Recommendations on Code Completion: Pass@1 and Improve@1 Evaluation for CodeLlama-Instruct-7B and StarCoder2-7B. The top 10 documents were used as context, with the number of documents incrementing from 1 to 10 in each trial. The best result from 10 runs was reported.

Model	Recall@10			
Wodel	GIST-large	Arctic-Embed		
APIRANKER	25.50	32.36		
w/o Perplexity Alignment	19.90	31.57		
w/o Reasoning State Extractor	24.66	27.91		
w/o Relevance Estimator	0.49	0.88		

Table 4: Ablation study.

in retrieval-augmented generation, demonstrating practical usability in real-world applications. (2) As illustrated in Fig. 5(b), we analyze the percentage of cases where recommended APIs enable code LLMs to correct outputs that it initially failed to generate (without recommended APIs). API-**RANKER** consistently outperforms other reranking models, achieving higher improvements in both code LLMs, highlighting its superior performance in scenarios where the code LLM's capabilities fall short and external API knowledge is needed.

Ablation Study. As illustrated in Table 4, we conduct an ablation study to assess the contribution of different techniques by removing key components (i.e., Perplexity Alignment via Comments, Hidden Reasoning State Extractor and Relevance Estimator) of our approach separately. The experimental results show that: (1) No matter which component we drop, it hurts the overall performance of our model, which signals the importance and effectiveness of all three components. (2) The recall rate shows a significant drop in reranking performance on the candidate documents retrieved by GIST-large and Arctic-Embed 2.0 when the Hidden Reasoning State Extractor and Relevance Estimator are removed separately. Notably, the removal of the Relevance Estimator causes an enormous decrease, which makes the model fail to work properly. This 559



Figure 6: The example of code completion based on API recommendation.

justifies the importance and necessity of these two components in our reranking model architecture.

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

594

Case Study and Manual Evaluation. As illustrated in Fig. 6, we present an example of generation using CodeLlama, based on API documents retrieved (e.g., colored in red) by Arctic-Embed 2.0 and reranked by our model. The retrieved APIs can't properly handle matrix exponentiation for Numpy, causing the completed code fails to pass test cases. APIRANKER reranks the retrieved APIs (e.g., colored in blue), successfully moving the correct API (e.g., colored in yellow) to the top, thus enhancing the LLM to produce the correct solution. This highlights that, APIRANKER can benefit automatic code completion task by providing more accurate and effective API recommendations. In addition, we conduct a manual evaluation to verify the validity of the relevance score predicted by our framework. Further details can be found in Appendix A.7.

5 Conclusions.

This research aims to recommend correct APIs to enhance automatic code completion, by jointly considering both natural language intent and incomplete code. To perform this task, we propose an approach APIRANKER that utilizes a self-learning ranking framework to automatically construct training data. Then we propose a novel reranking model to predict the relevance score between the API documents and the query, based on the LLM's reasoning capabilities. The experimental results show the effectiveness of our approach in both API recommendation and automatic code completion. We hope our study lays the foundations for this research and provides valuable insights.

532

533

6 Limitations.

595

Several limitations are concerned with our work. Firstly, due to the limited availability of code com-597 pletion test sets that support code evaluation in 598 other languages, and the difficulty in constructing queries that simultaneously include both intent and incomplete code, our test is based on Python, one of the most popular programming languages used by developers. However, during the training of our method, we used data from two programming languages Java and Python, and we believe that our approach can easily adapt to other programming languages. Secondly, our approach does not explicitly create intent but rather leverages the language comprehension ability of LLMs to reduce the need for learning natural language intent. Exploring how 610 to automatically generate high-quality intent from 611 code is an interesting research topic for our future work. 613

614 References

615

616

617

618

619

620

621

622

625

628

635

641

642

643

- Yujia Chen, Cuiyun Gao, Xiaoxue Ren, Yun Peng, Xin Xia, and Michael R Lyu. 2023. Api usage recommendation via multi-view heterogeneous graph representation learning. *IEEE Transactions on Software Engineering*, 49(5):3289–3304.
- Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155*.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:2312.10997.
- Xiaodong Gu, Hongyu Zhang, Dongmei Zhang, and Sunghun Kim. 2016. Deep api learning. In Proceedings of the 2016 24th ACM SIGSOFT international symposium on foundations of software engineering, pages 631–642.
- Michael Günther, Jackmin Ong, Isabelle Mohr, Alaeddine Abdessalem, Tanguy Abel, Mohammad Kalim Akram, Susana Guzman, Georgios Mastrapas, Saba Sturua, Bo Wang, et al. 2023. Jina embeddings 2: 8192-token general-purpose text embeddings for long documents. arXiv preprint arXiv:2310.19923.
- Daya Guo, Shuai Lu, Nan Duan, Yanlin Wang, Ming Zhou, and Jian Yin. 2022. Unixcoder: Unified crossmodal pre-training for code representation. *arXiv preprint arXiv:2203.03850*.

Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. 2024. Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.06769*. 645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

- Qiao Huang, Xin Xia, Zhenchang Xing, David Lo, and Xinyu Wang. 2018. Api method recommendation without worrying about the task-api knowledge gap. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, pages 293–304.
- Rasha Ahmad Husein, Hala Aburajouh, and Cagatay Catal. 2024. Large language models for code completion: A systematic literature review. *Computer Standards & Interfaces*, page 103917.
- Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. 2024. A survey on large language models for code generation. *arXiv preprint arXiv:2406.00515*.
- Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2023. Ds-1000: A natural and reliable benchmark for data science code generation. In *International Conference on Machine Learning*, pages 18319–18345. PMLR.
- Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024. Nv-embed: Improved techniques for training llms as generalist embedding models. *arXiv* preprint arXiv:2405.17428.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*.
- Chunyang Ling, Yanzhen Zou, and Bing Xie. 2021. Graph neural network based collaborative filtering for api usage recommendation. In 2021 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER), pages 36–47. IEEE.
- Mingwei Liu, Tianyong Yang, Yiling Lou, Xueying Du, Ying Wang, and Xin Peng. 2023. Codegen4libs: A two-stage approach for library-oriented code generation. In 2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE), pages 434–445. IEEE.
- Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. 2024. Starcoder 2 and the stack v2: The next generation. *arXiv preprint arXiv:2402.19173*.
- Shuai Lu, Nan Duan, Hojae Han, Daya Guo, Seungwon Hwang, and Alexey Svyatkovskiy. 2022. Reacc: A retrieval-augmented code completion framework. *arXiv preprint arXiv:2203.07722.*

Zexiong Ma, Shengnan An, Bing Xie, and Zeqi Lin. 2024. Compositional api recommendation for libraryoriented code generation. In *Proceedings of the 32nd IEEE/ACM International Conference on Program Comprehension*, pages 87–98.

704

706

710

711

712

713

714

715

716

717

718

719

720

721

722

724

725

726

727

728

729

730

731

733

734

737

741

743

744

745

747

748

751

- Marcellino Marcellino, Davin William Pratama, Steven Santoso Suntiarko, and Kristien Margi. 2021.
 Comparative of advanced sorting algorithms (quick sort, heap sort, merge sort, intro sort, radix sort) based on time and memory usage. In 2021 1st International Conference on Computer Science and Artificial Intelligence (ICCSAI), volume 1, pages 154–160. IEEE.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*.
- Noor Nashid, Taha Shabani, Parsa Alian, and Ali Mesbah. 2024. Contextual api completion for unseen repositories using llms. *arXiv preprint arXiv:2405.04600*.
- Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. 2023. A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Arnold Overwijk, Chenyan Xiong, Xiao Liu, Cameron VandenBerg, and Jamie Callan. 2022. Clueweb22: 10 billion web documents with visual and semantic information. arXiv preprint arXiv:2211.15848.
- Md Rizwan Parvez, Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Retrieval augmented code generation and summarization. *arXiv preprint arXiv:2108.11601*.
- Yun Peng, Shuqing Li, Wenwei Gu, Yichen Li, Wenxuan Wang, Cuiyun Gao, and Michael R Lyu. 2022.
 Revisiting, benchmarking and exploring api recommendation: How far are we? *IEEE Transactions on Software Engineering*, 49(4):1876–1897.
- Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, et al. 2023. Large language models are effective text rankers with pairwise ranking prompting. *arXiv preprint arXiv:2306.17563*.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. 2023. Code llama: Open foundation models for code. arXiv preprint arXiv:2308.12950.

Devendra Singh Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. Improving passage retrieval with zero-shot question generation. *arXiv preprint arXiv*:2204.07496. 752

753

754

755

756

757

759

760

761

762

763

764

765

766

767

768

769

770

772

773

774

775

776

778

779

780

781

782

783

784

785

786

787

788

789

790

791

793

794

796

797

798

799

800

801

802

803

804

805

- Aivin V Solatorio. 2024. Gistembed: Guided in-sample selection of training negatives for text embedding fine-tuning. *arXiv preprint arXiv:2402.16829*.
- Yanlin Wang, Yanli Wang, Daya Guo, Jiachi Chen, Ruikai Zhang, Yuchi Ma, and Zibin Zheng. 2024a. Rlcoder: Reinforcement learning for repository-level code completion. arXiv preprint arXiv:2407.19487.
- Yong Wang, Yingtao Fang, Cuiyun Gao, and Linjun Chen. 2024b. Api recommendation for novice programmers: Build a bridge of query-task knowledge gap. *IEEE Transactions on Reliability*.
- Zora Z. Wang, Akari Asai, Xinyan V. Yu, Frank F. Xu, Yiqing Xie, Graham Neubig, and Daniel Fried. 2024c. Coderag-bench: Can retrieval augment code generation? *arXiv preprint arXiv:2406.14497*.
- Moshi Wei, Nima Shiri Harzevili, Alvine Boaye Belle, Junjie Wang, Lin Shi, Song Wang, and Zhen Ming Jiang. 2023. A survey on query-based api recommendation. *arXiv preprint arXiv:2312.10623*.
- Moshi Wei, Nima Shiri Harzevili, Yuchao Huang, Junjie Wang, and Song Wang. 2022. Clear: contrastive learning for api recommendation. In *Proceedings* of the 44th International Conference on Software Engineering, pages 376–387.
- Di Wu, Xiao-Yuan Jing, Hongyu Zhang, Yang Feng, Haowen Chen, Yuming Zhou, and Baowen Xu. 2023. Retrieving api knowledge from tutorials and stack overflow based on natural language queries. *ACM Transactions on Software Engineering and Methodology*, 32(5):1–36.
- Puxuan Yu, Luke Merrick, Gaurav Nuti, and Daniel Campos. 2024. Arctic-embed 2.0: Multilingual retrieval without compromise. *arXiv preprint arXiv:2412.04506*.
- Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and Weizhu Chen. 2023a. Repocoder: Repository-level code completion through iterative retrieval and generation. *arXiv preprint arXiv:2303.12570*.
- Jingxuan Zhang, He Jiang, Zhilei Ren, and Xin Chen. 2017. Recommending apis for api related questions in stack overflow. *IEEE Access*, 6:6205–6219.
- Kechi Zhang, Huangzhao Zhang, Ge Li, Jia Li, Zhuo Li, and Zhi Jin. 2023b. Toolcoder: Teach code generation models to use api search tools. *arXiv preprint arXiv:2305.04032*.
- Shuyan Zhou, Uri Alon, Frank F Xu, Zhiruo Wang, Zhengbao Jiang, and Graham Neubig. 2022. Docprompting: Generating code by retrieving the docs. *arXiv preprint arXiv:2207.05987*.

A Appendix

807

810

811

812

813

814

815

816

817

820

821

822

825

832

833

839

841

844

849

853

A.1 Dataset Construction Details

We collect 110,646 API documentations from the dataset CodeRAG-bench (Wang et al., 2024c) as retrieval sources. These documents come from two main sources: official Python library documentation provided by devdocs.io⁴ and content obtained from ClueWeb22 (Overwijk et al., 2022), a largescale web corpus, covering a wide range of topics, from basic programming techniques to advanced library usage. Each page in ClueWeb22 includes code snippets and textual explanations. An API documentation typically includes the API's purpose, function description, input parameters, output values, and example requests and responses. Howerever, we take into account the diversity of API descriptions in real-world scenarios and do not strictly restrict the structure or form of API documents, as long as they can be converted into text or code formats.

> To support efficient vector queries based on cosine similarity, we create vector search libraries using Milvus⁵, a high-performance vector database designed for scalability and providing fast, scalable similarity search and retrieval.

> In the training and evaluation data, for a specific code file that has cross-file dependencies, we treat all the dependency files (i.e., both direct and indirect dependencies) of the code file as the candidate documentation for the code file. For the test data, we select DS-1000 (Lai et al., 2023) as the query (*i.e.*, NL intent and incomplete code) and code completion dataset, which includes general open-domain coding completion tasks (*e.g.*, Matplotlib, Numpy, Pandas, Sklearn, Tensorflow).

A.2 Baselines Setup detail

Retrieval Baselines. Since the performance of code retrieval models (*e.g.*, Unixcoder, Code-Bert (Feng et al., 2020), jina-base-v2-code (Günther et al., 2023)) is not ideal (with poor retrieval performance), we do not conduct reranking experiments on it. Additionally, CodeBert and jina-base-v2-code are unable to recall any relevant API documents in the top 50, we do not report retrieval results for these models.

We consider the following retrieval baselines, which are dense retrievers that encode both the

query and code documentation into vector spaces for retrieving semantically relevant documentation based on vector similarity: (1) Unixcoder: Unixcoder (Guo et al., 2022) is a unified crossmodal pre-trained model for programming language. (2) GIST-large: GIST-large (Solatorio, 2024) is a method that improves text embedding fine-tuning by selectively choosing negative samples. (3) Arctic-Embed 2.0: Arctic-Embed 2.0 (Yu et al., 2024) is an open-source text embedding model built for accurate and efficient multilingual retrieval. (4) NV-Embed-v2: NV-Embed-v2 (Lee et al., 2024) is a generalist embedding model that ranks No. 1 in the retrieval sub-category of the Massive Text Embedding (MTEB) leaderboard (Muennighoff et al., 2022). GIST-large and Arctic-Embed 2.0 are also ranked highly on the MTEB leaderboard.

Considering the context limitations of retrieval and code generation, as well as the excessive length of some API documentation, the retrieval model's maximum token encoding length is uniformly set to 512. An API documentation typically includes the API's purpose, function description, input parameters, output values, and example requests and responses. Most API documentation includes essential information, and although some longer API docs may be truncated at the "example" section, the necessary details, including the description of the API's role and function, are typically present within the first 512 tokens.

Reranking Baselines. We consider the following reranking baselines, which are based on LLMs:

(1) Unsupervised Passage Re-ranker (UPR): UPR (Sachan et al., 2022) is a pointwise approach based on query generation. The prompt template for UPR is shown in Fig. 7. In this approach, the relevance score of an API document d to the query q is measured by the probability of generating the query.

• Instruction:	Please	write a	question	based	on this	pas-
sage.						
• Passage: d						
• Question: q						

Figure 7: The prompt template for UPR. d is the API document, q is the query.

(2) **Relevance Generation (RG)**: RG (Liang et al., 2022) is a pointwise approach based on relevance generation. The prompt template for RG is shown in Fig. 8. In this approach, the relevance of

898

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

⁴https://devdocs.io

⁵https://github.com/milvus-io/milvus

899

900

901 902

respectively.

• Passage: d • Query: q

document, q is the query.

"Passage A" or "Passage B".

913

914

915

931

933

(3) **Pairwise Ranking Prompting- Sorting** (**PRP-Sorting**): PRP-Sorting (Qin et al., 2023) is a pairwise method based on the log-likelihood of document generation, and it optimizes time complexity through heap sort algorithm (Marcellino

et al., 2021). The prompt template for PRP-Sorting

is shown in Fig. 9. In this approach, to compare

two API documents d_A and d_B , the one that is

more relevant to the query q is determined based

on which has a higher probability of generating

an API document d to the query q is defined as:

 $s_i = \begin{cases} 1 + p(\text{Yes}), & \text{if output Yes} \\ 1 - p(\text{No}), & \text{if output No} \end{cases}$

• Instruction: Does the passage answer the query?

where p(Yes) and p(No) denote the probabilities

of LLMs generating the tokens of "Yes" or "No"

Figure 8: The prompt template for UPR. d is the API

(13)

Instruction: Given a query "q", which of the following two passages is more relevant to the query?
Passage A: d_A
Passage B: d_B

Figure 9: The prompt template for PRP-Sorting and

Figure 9: The prompt template for PRP-Sorting and PRP-Sliding. d is the API document, q is the query.

(4) **Pairwise Ranking Prompting-Sliding** (**PRP-Sliding**): PRP-Sliding is a variant of PRP, which is based on the sliding window approach. The prompt template and comparison function for PRP-Sliding are the same as those for PRP-Sorting.

(5) **RepoCoder**: RepoCoder (Zhang et al., 2023a) is a reranking method through iterative retrieval of code snippets based on the result of code generation. The API documentation is provided as retrieval source for the RepoCoder. We conducted the experiment using a 2-iteration approach, following the method described in the RepoCoder paper.

In order to comparing the performance of different reranking methods, we uniformly use CodeLlama-Instruct-7B as the base LLMs. The comparison between methods is made using the same retrieval source. The maximum token length of an API document is set to 512. The baseline of API recommendations. We adopt two common approaches of API recommendations: (1) BIKER (Huang et al., 2018): BIKER is an API recommendation approach that bridges lexical and knowledge gaps by using word embeddings for similarity and similar questions retrieval for supplementary information. Here, we take the queries in the test set as the source of similar questions. (2) GAPI (Ling et al., 2021): GAPI uses the code context as the query for API usage recommendation and employs graph neural networks to capture high-order collaborative signals. However, due to differences in task setup and dataset, the project structure information is not available in our dataset. We only used text attributes to nodes as input for API prediction since lacking project structural information in our datasets.

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

Code Completion Baselines. For automatic code completion, we consider the following code LLMs: (1) Starcoder2-7B (Lozhkov et al., 2024), which is trained on a vast programming dataset and achieves superior performance on code-related tasks. (2) CodeLlama-Instruct-7B (Roziere et al., 2023), which is a fine-tuned version of Code Llama, optimized to follow natural language instructions for code generation.

A.3 Implementation Details

In our approach, we chose CodeLlama-Instruct-7B as the perplexity evaluator in the self-supervised learning ranking framework and as the base LLM of the reranking model. Additionally, UnixCoder is chosen as the retriever in the self-supervised learning ranking framework and as the initial weight of the hidden reasoning state extractor. All experiments were conducted on two A800 GPUs.

In our self-supervised learning ranking framework, we set the total length of the incomplete code and the target code to be no more than 1024 tokens, ensuring that the ratio of 0.4 to 0.5 of the total length is considered as the incomplete code. The prompt template for the perplexity evaluator is shown in 10.

In the design of the reranking model, we set the number of learnable vectors in the hidden reasoning state extractor to 32. We employed the AdamW optimizer with a learning rate of 1e-4. The learning rate schedule was managed using the WarmupCosineLR scheduler, where the learning rate linearly warms up for the first 75 steps and then follows a cosine decay towards a minimum ratio

of 0.0001 over a total of 750 steps. The batch size 984 was set to 384, and the number of gradient accumulation steps was 4. The input length was capped at a maximum of 1152 tokens. We constructed pairs from a set of 4 documents, selected evenly based on the difference in values from the perplexity evaluator, chosen from the top 20 candidate documents 990 retrieved by the retriever. The prompt template for 991 training is shown in Fig. 10. During the inference stage, we reranked the top 50 documents retrieved 993 by different retrieval models. The input length was 994 capped at a maximum of 1600 tokens. The prompt 995 template for inference was the same as for training.

Instruction: Refer to the following documentation (between "— Documentation —" and "— End Documentation —") to complete the code.
API document: d
query: q

Figure 10: The prompt template for APIRANKER. d is the API document, q is the query.

For retrieval-augmented code completion, we use top-k API documents as a context for automatic code completion, keeping only the first 512 tokens in each document. The prompt template of retrieval-augmented code completion is shown in Fig. 11. During decoding, code is generated using greedy decoding. The length of the output to a maximum of 2048 tokens.

Instruction: Refer to the following documentation (between "— Documentation —" and "— End Documentation —") to complete the code. Based on the following problem description and existing code, please write the code to achieve the desired output. Place the executable code between <code> and </code> tags, without any other non-executable things.
the top-k API documents: d1, ..., dk
query: q

Figure 11: The prompt template for PRP-Sorting and PRP-Sliding. d_i is the *i*-th API document, q is the query.

A.4 Evaluation Metrics

To evaluate the performance of API recommendation, we report the common evaluation metrics (Zhang et al., 2017; Wei et al., 2023): (1) **Recall@k**, measures the proportion of correct API documents in the the top-k recommendation results. It is defined as follows:

$$\operatorname{Recall}@\mathbf{k} = \frac{\mathbf{R}}{\mathbf{N}},\tag{14}$$

1013 where N is the total number of relevant documents, 1014 and R is the number of relevant documents in top-k recommended results. (2) NDCG@k, evaluates1015the ranking of correct documents in the top-k rec-
ommendation results. As a normalized Discounted1016Cumulative Gaine, NDCG is calculated by dividing
by a special ideal DCG, where all relevant docu-
ments are ranked higher than irrelevant ones. It is
defined as:1015

$$NDCG@k = \frac{DCG@k}{ideal \ DCG@k}, \qquad (15) \qquad 102$$

$$DCG@k = \sum_{i=1}^{k} \frac{2^{rel(i)} - 1}{\log_2(i+1)},$$
 (16) 10

1024

1025

1026

1027

1028

1029

1030

1032

1034

1035

1036

1038

1039

1041

1042

1043

1044

1045

1046

1047

1049

1050

1052

where *i* represents the rank. rel(i) is a binary function to check whether the API in rank *i* is correct or not. If the API at rank *i* is a correct API, then the value rel(i) is 1; otherwise, the value is 0. (3) **MRR@k**, represents the reciprocal of the position where the first correct API appears in the top-k recommendation results. It is defined as:

$$MRR@k = \frac{1}{|Q|} \sum_{j=1}^{Q} \frac{1}{k_{Rank_i}},$$
 (17) 1031

where |Q| is the number of queries Q, and k_Rank_i means the rank position of the first correct answer in the top k recommended list for the *i*-th query. (4) **Mean Average Precision (MAP)**, evaluates the overall performance by taking into account the ranking of correct API documents. It is defined as:

$$MAP = \frac{1}{|Q|} \sum_{j=1}^{Q} \frac{\sum_{i=1}^{n} (P(i) \times rel(i))}{\#correct \text{ answers}}, \quad (18)$$

$$P(i) = \frac{\#\text{correct answers in top } i}{i}, \qquad (19)$$

where p(i) is the precision at a given cut-off rank *i*. The value of *k* is set to 10, and *n* is set to 50. We use Recall@k as the primary metric since retrievalaugmented generation primarily relies on key information that appears in the context.

To evaluate the performance of code completion based on API recommendation, we adopt Pass@k and Improve@k metrics to measure the execution correctness of programs: (1) Pass@k, is an evaluation metric that has been widely used in previous work (Jiang et al., 2024), computing the fraction of problems having at least one correct prediction within k samples. It is defined as:

$$pass@k := \mathbb{E}_{task}\left[1 - \frac{\binom{n-c}{k}}{\binom{n}{k}}\right], \quad (20) \quad 1054$$

1005

1006

1007

1008

1009

where n is the total number of sampled candidate 1055 code solutions, k is the number of randomly se-1056 lected code solutions from these candidates for 1057 each programming problem, with $n \ge k$, and c is 1058 the count of correct samples within the k selected. 1059 (2) Improve@k, is the proportion of cases in which 1060 the code LLM generates the correct output with 1061 the recommended API documentation, compared 1062 to when it initially failed without the recommended 1063 API documentation. It is defined as: 1064

1065

1066

1067

1068

1069

1070

1071

1074

1075

1076

1078

1079

1080

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1094

1095

1096

1097

1098

1099

1100

1101

1102

Improve@k =
$$\frac{\sum_{i=1}^{m} \text{correct}(i)}{\#\text{failures in } k \text{ samples}}$$
, (21)

where m is the number of problems that initially failed to generate the code in the k samples, and correct(i) is 1 if the *i*-th problem passes in the ksamples, and 0 if it fails. The value of k is set to 1 in our experiment. Given the differences in the capabilities of code LLMs, there are instances where a model, initially capable of generating correct outputs, may fail when code completion is based on API documents. Therefore, we use Improve@k to explore the potential for improvement.

A.5 Experimental Comparison of API Recommendation.

The results show that APIRanker significantly outperforms both BIKER and GAPI, demonstrating its superior effectiveness over query-based methods and code context-based method, verifying the effectiveness of APIRanker for combining user intent and code context for API recommendation.

Based on the experimental results, our model APIRanker is bertter than RepoCoder. (1) API-Ranker significantly outperforms RepoCoder in terms of different retrieval methods. This advantage is likely due to the larger scale of API document retrieval, which recalls a much larger number of similar documents compared to repositorylevel code retrieval. (2) APIRanker achieves consistent performance improvements, while Repocoder is more dependent on the quality of the retrieval model. For example, RepoCoder experiences a notable decline in terms of using GIST-large and Arctic-Embed, but shows an improvement when paired with the strong retrieval model NV-Embedv2.

A.6 Discussing PPL and Log-Probability Based Uncertainty.

PPL is the exponentiated average of the logprobability based uncertainty, which is more suitable for assessing the overall model performance, 1103 as it aggregates token-level uncertainties into a 1104 comprehensive score. While the Log-Probability 1105 based uncertainty is used to calculate the genera-1106 tion probability of each token, which offers more 1107 granular insight on individual token-level. In terms 1108 of our research of code completion, we care more 1109 about the model's ability to generate complete code 1110 sequence, thus PPL is more appropriate. Regarding 1111 tasks for more detailed token-level analysis (e.g., 1112 keyword analysis, code style analysis), examining 1113 log-probability based uncertainty could be more 1114 informative. 1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

A.7 Manual Evaluation.

To further validate the effectiveness of our reference, we conduct a user study. In particular, we randomly select 100 training API-query pairs that are scored based on our framework and asked 3 users (each with over 4 years of programming experience) to assess them. Users are asked to answer the question: "Which of the two API documents is more helpful for the query?", and every user is provided with three options (i.e., A is Better, B is Better, Cannot Determine/Both Equally). We calculate the agreement ratio between the manual evaluations and the automated scores. The results of the user study are as follows:

Consistency	Inconsistency	Indeterminate		
85%	11%	4%		

Table 5: The result of the user study.

We calculate the Pearson correlation between1130the manual evaluations and the automated scores.1131The high consistency ratio of 94.8% indicates that1132our method aligns well with human evaluations,1133demonstrating its effectiveness in generating train-1134ing data for API reranking.1135