# CodeGRAG: Bridging the Gap between Natural Language and Programming Language via Graphical Retrieval Augmented Generation

Anonymous ACL submission

### Abstract

Utilizing large language models to generate codes has shown promising meaning in software development revolution. Despite the intelligence shown by the general large language models, their specificity in code generation can still be improved due to the syntactic gap and mismatched vocabulary existing among natural language and different programming languages. In this paper, we propose CodeGRAG, a Graphical Retrieval Augmented Code Generation framework to enhance the performance of LLMs. CodeGRAG builds the graphical view of code blocks based on the control flow and data flow of them to fill the gap between programming languages and natural language, which can facilitate natural 017 language based LLMs for better understanding of code syntax and serve as a bridge among different programming languages. To take the extracted structural knowledge into the foundation models, we propose 1) a hard meta-graph prompt template to transform the challenging graphical representation into informative knowledge for tuning-free models and 2) a soft prompting technique that injects the domain knowledge of programming languages into the model parameters via finetuning the models with the help of a pretrained GNN expert model. CodeGRAG significantly improves the code generation ability of LLMs and can even offer performance gain for cross-lingual code generation. Implementation is available at https://anonymous.4open.science/r/Code-5970/.

### 1 Introduction

042

In recent years, large language models (LLMs) (Achiam et al., 2023; Touvron et al., 2023a) have shown great impact in various domains. Automated code generation emerges as a captivating frontier (Zheng et al., 2023; Roziere et al., 2023; Shen et al., 2023), promising to revolutionize software development by enabling machines to write and optimize



Figure 1: Illustration of the gap between the programming language and the natural language.

043

044

045

046

047

052

056

058

060

061

062

063

064

065

067

068

code with minimal human intervention.

However, syntactic gap and mismatched vocabulary among natural language (NL) and programming languages (PL) exist, hindering LLM's performance on code generation. As illustrated in Figure 1, programming language (marked in blue) contains special tokens such as "int" or "++" that natural language (marked in yellow) doesn't possess, leading to vocabulary mismatch. Besides, the relations between tokens in programming languages are often structural, e.g., the complex branching and jumps, whereas natural language is arranged simply in sequential manner, leading to syntactic gap. For example, in the control flow graph of the raw code (marked in pink), two "if" blocks (marked in purple) are adjacent and are executed sequentially under certain condition, but they appear to be intervaled in raw textual code.

As discussed above, the innate structures of programming languages are different from that of the sequential-based natural language. The challenges of enhancing a general-purposed large language models for code-related tasks can be summarized into two folds.

(C1) How to solve the gap between different languages and better interpret the inherent logic

of code blocks. Code, unlike natural language, possesses a well-defined structure that governs its syntax and semantics. This structure provides valuable information about the relationships between different parts of the code, the flow of execution, and the overall organization of the functions (Jiang et al., 2021; Guo et al., 2020). General-purpose LLMs regard a code block as a sequence of tokens. By ignoring the inherent structure of codes, they miss out on essential cues that could help them better understand and generate code. In addition, the multi-lingual code generation abilities of LLMs is challenging due to the gap among different programming languages.

> (C2) How to inject the innate knowledge of programming languages into general purpose large language models for enhancement. Despite the well representation of the programming knowledge, the ways to inject the knowledge into the NL-based foundation models is also challenging. The structural representation of codes could be hard to understand, which poses a challenge to the capability of the foundation models.

087

880

095

100

101

102

103

104

106

107

108

110

111

112

113

114

115

116

117

To solve the above challenges, we propose Code-GRAG, a graphical retrieval augmented generation framework for code generation. For (C1), we propose to interpret the code blocks using the composed graph based on the data-flow and controlflow of the code block, which extracts both the semantic level and the logical level information of the code. The composed graphical view could 1) better capture the innate structural knowledge of codes for NL-based language models to understand and 2) model the innate function of code blocks that bridging different programming languages. For (C2), we propose a meta-graph prompting technique for tuning-free models and a softprompting technique for tuned models. The metagraph prompt summarizes the overall information of the extracted graphical view and transforms the challenging and noisy graphical representation into informative knowledge. The soft-prompting technique deals with the graphical view of codes with a pretrained GNN expert network and inject the processed knowledge embedding into the parameters of the general-purpose foundation models with the help of supervised finetuning.

> The main contributions of the paper can be summarized as follows:

 Novel GraphRAG framework for code generation. We propose CodeGRAG that bridges the gap among natural language and programming languages, transfers knowledge among different programming languages, and enhances the ability of LLMs for code generation. CodeGRAG requires only one calling of LLMs and can offer multi-lingual enhancement. 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

163

164

165

- Effective graphical view to inform and stimulate the structural programming knowledge of LLMs. We propose an effective graphical view to purify the semantic and logic knowledge from the code space, which offers more useful information than the raw code block and can summarize the cross-lingual knowledge.
- Effective soft prompting technique to preserve the programming domain knowledge and inject it into LLMs parameters. We propose an effective soft prompting technique, which injects the domain knowledge of programming languages into the model parameters via finetuning LLMs with the assistance of a pretrained GNN expert model.

# 2 Methodology

## 2.1 Overview

In this paper, we leverage both generative models and retrieval models to produce results that are both coherent and informed by the expert graphical knowledge of programming language. The overall process of CodeGRAG is illustrated in Figure 2, which mainly consists of three stages: graphical knowledge base preparation, knowledge querying, and graphical knowledge augmented generation.

# 2.2 Graphical Knowledge Base Preparation

In this section, we discuss how to extract informative graphical views for code blocks. We analyze the syntax and control information of code blocks and extract their graphical views to better represent the codes. This process can be formulated as,  $\forall c_i \in D_{pool}$ :

$$g_i \leftarrow \text{GraphExtractor}(c_i),$$
 (1)

$$\mathsf{KB.append}(\langle c_i, g_i \rangle), \tag{2}$$

where  $c_i$  is the raw code block and  $g_i$  is the corresponding extracted graphical view.

To capture both the semantic and the logical information, we propose to combine the data flow graph (Aho et al., 2006) and the control flow graph (Allen, 1970) with the read-write signals (Long



Figure 2: Overview of CodeGRAG. (**Top**) **Knowledge Preparation.** We extract the control flow and data flow of each external code block and compose them using the read-write signal to obtain the semantic and logical expression of each code block, which is then abstracted into graphical view as hard knowledge document and embedded into GraphEmb as soft knowledge document. The GraphEmb is encoded by a pretrained GNN expert model constrained by the alignment and structure preserving objectives. (**Bottem**) **Retrieval Augmented Generation.** We extract query from the task input and retrieve from the external corpus. For tuning free models, we use the hard graphical view to stimulate the structural programming knowledge of LLMs for enhanced generation. For tunable models, we use the soft GraphEmb and inject the programming domain knowledge into LLMs parameters via finetuning them with the GNN expert signals. The expert signals informed LLMs can then produce enhanced generation.

et al., 2022) to represent the code blocks, both of them are constructed on the base of the abstract syntax tree.

166

167

168

169

171

172

174

175

Abstract Syntax Tree (AST). An abstract syntax tree (AST) is a tree data structure that represents the abstract syntactic structure of source code. An AST is constructed by a parser, which reads the source code and creates a tree of nodes. Each node in the tree represents a syntactic construct in the source code, such as a statement, an expression, or a declaration. ASTs have good compactness and can represent the structure of the source code in a clear and concise way.

179Data Flow Graph (DFG). The data flow graph180(DFG) is a graphical representation of the flow of181data dependencies within a program. It is a directed182graph that models how data is transformed and183propagated through different parts of a program. In184DFG, nodes are operands and edges indicate data185flows. Two types of edges are considered: 1) opera-186tion edges that connect the nodes to be operated and187the nodes that receive the operation results; 2) func-188tion edges that indicate data flows for function calls

and returns. These edges connect nodes, including non-temporary operands and temporary operands, which refer to variables and constants that explicitly exist in the source code, and variables existing only in execution, respectively. 189

190

191

192

194

195

196

197

198

199

200

201

202

203

204

205

206

207

209

210

211

**Control Flow Graph (CFG).** The control flow graph (CFG) is a graphical representation of the flow of control or the sequence of execution within a program. It is a directed graph that models the control relationships between different parts of a program. Based on compiler principles, we slightly adjust the design of CFG to better capture the key information of the program. Nodes in CFG are operations in the source code, including standard operations, function calls and returns. Edges indicate the execution order of operations.

**Composed Syntax Graph.** A composed syntax graph composes the data flow graph and the control flow graph with the read-write flow existing in the code blocks. An illustration of the extracted composed syntax graph is displayed in Figure 3. Different edge types along with their concrete names are given in colors. As for the node names, the



Figure 3: Illustration of the extracted composed syntax graph from the code block. The arrows in the bottom part indicate the names of different edges, which are extracted based on the ASTs.

middle figure displays the concrete types of nodes (operands) and the right figure displays the properties of nodes.

An illustration of the composed graphical view is in Figure 3. After obtaining the composed syntax graphs, we use them to inform the generalpurpose LLMs to bridge the gap between NL and PLs, where both the semantic level and the logic level information are preserved.

# 2.3 Knowledge Querying

Given a target problem to be completed, we generate informative query of it and use it to retrieve graphical knowledge from the constructed knowledge base.

We extract the problem description of each task to reduce the ambiguity and then concatenate it with the function declaration to serve as the query content, where the functionality and input format of the expected code block are contained. The query of the retrieval includes problem description  $Q_p$ and function description  $Q_c$ , while each content of the retrieval pool includes raw code block  $V_c$  and its graphical view  $V_g$ .

To expressively represent the components, we use the encoder  $\phi(\cdot)$  of the pretrained NL2Code model to represent the problem description and code snippets. The retrieval function is:

$$\mathbf{h}^{\mathbf{V}} = \phi(V_c \| V_g), \tag{3}$$

$$\mathbf{h}^{\mathbf{Q}} = \phi(Q_p \| Q_c), \tag{4}$$

Distance = 
$$1 - \frac{\mathbf{h}^{\mathbf{Q}} \cdot \mathbf{h}^{\mathbf{v}}}{\|\mathbf{h}^{\mathbf{Q}}\| \cdot \|\mathbf{h}^{\mathbf{V}}\|}.$$
 (5)

### 2.4 Graphical Knowledge Augmented Generation

After we obtain the returned graphical view, we inject it to the foundation LLMs for graphical knowledge augmented generation. Since the graphical view is hard to understand, we propose 1) a metagraph template to transform the graphical view into informative knowledge for tuning-free model and 2) a soft prompting technique to tune the foundation models for their better understanding of the graphical views with the assistance of an expert GNN model. 246

247

248

249

250

252

253

254

255

257

258

259

260

261

263

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

281

### 2.4.1 Hard Meta-Graph Prompt

The original graphical view of a code block could contain hundreds of nodes and edges. A full description of it could cost an overly long context, along with the understanding challenge posed by the long edge lists. Therefore, we propose to use a meta-graph template to abstract the information of the graphical view. The abstracted meta-graph consists of the canonical edge types and node types, which describes the basic topology of the graphical view (Sun and Han, 2013), with the textual features obtained from the ASTs contained in the node and edge features.

Then we use the meta-graph template to transform the retrieved graphical view into digestable knowledge and insert it into the final prompt for generation. As illustrated in Figure 4 in the Appendix, the final prompt consists of three components: the system prompt illustrated in the blue part, the retrieved knowledge and hints illustrated in the green part, and the problem (including task description, function declaration, etc.) illustrated in the yellow part. The three parts are concatenated to be fed into LLMs for knowledge augmented generation.

### 2.4.2 Soft Prompting with the Expert

Directly hard prompt to the LLMs poses a challenge to the digesting capability of the backbone

4

212 213

- 23
- 239 240

241

242

243

245

327

328

329

330

LLMs, which could fail under the case where the backbone LLMs cannot well understand the graph components. To compress the graphical knowledge into model parameters and help the backbone LLMs to better understand the programming language, we propose a soft prompting technique. The overall procedure can summarized into expert encoding of graphical views, finetuning with the expert signal, and inference.

**Expert Encoding of Graphical Views.** We design a graph neural network to preserve the semantic and logical information of code blocks. The representation of each node  $\mathbf{n}_i^{(0)}$  and edge  $\mathbf{e}_i^{(0)}$  are first initialized with vectors corresponding to the node text and edge text encoded by  $\phi_1$ . A message passing process is first conducted to fuse the semantic and structural information into each node representation.

291

296

297

302

310

311

312

313

314

315

316

317

318

319

321

322

323

$$\mathbf{m}_{ij}^{(l)} = \mathbf{W}^{(l)}(\mathbf{n}_i^{(l-1)} \| \mathbf{e}_{ij}^{(l-1)}), \tag{6}$$

$$Q_j^{(l)} = \mathbf{W}_{\mathbf{Q}}^{(l)} \mathbf{n}_j^{(l-1)}, \tag{7}$$

$$\mathbf{K}_{ij}^{(l)} = \mathbf{W}_{\mathbf{K}}^{(l)} m_{ij}^{(l)}, \quad \mathbf{V}_{ij}^{(l)} = \mathbf{W}_{\mathbf{V}}^{(l)} \mathbf{m}_{ij}^{(l)}, \quad (8)$$

$$a_{ij}^{(l)} = \operatorname{softmax}_{i \in N(j)}(\mathbf{Q}_j^{(l)}\mathbf{K}_{ij}^{(l)}), \qquad (9)$$

$$\mathbf{n}_{j}^{(l)} = \sum_{i \in N(j)} a_{ij}^{(l)} \mathbf{V}_{ij}^{(l)}.$$
 (10)

A global attention-based readout is then applied to obtain the graph representation:

$$\mathbf{g} = \sum_{i} \operatorname{softmax}(f_{\text{gate}}(\mathbf{n}_{i}^{L})) f_{\text{feat}}(\mathbf{n}_{i}^{L}).$$
(11)

The expert encoding network is optimized via the contrastive learning based self-supervised training, which includes the intra-modality contrastive learning and inter-modality contrastive learning. The intra-modality constrastive learning serves for preserving the modality information, while the inter-modality contrastive learning serves for modality alignment.

• Alignment Contrastive Learning. There are two types of alignment to be ensured: 1) NL-Code (NC) alignment and 2) Code-Graph (CG) alignment. We define the positive pairs for NC alignment purpose as  $\mathcal{I}_{NC}^+ = \{\langle \mathbf{h}_i^V, \mathbf{h}_i^Q \rangle | i \in D_{\text{train}}\}$  and define the negative pairs for NC alignment purpose as  $\mathcal{I}_{NC}^- = \{\langle \mathbf{h}_i^V, \mathbf{h}_j^Q \rangle | i \neq j, i \in D_{\text{train}}\}$ .

And we define the positive pairs for CG alignment purpose as  $\mathcal{I}_{CG}^+ = \{ \langle \phi_1(c_i), \phi_2(g_i) \rangle | i \in \}$   $D_{\text{train}}$  and define the negative pairs for CG alignment purpose as  $\mathcal{I}_{CG}^- = \{\langle \phi_1(c_i), \phi_2(g_j) \rangle | i \neq j, i \in D_{\text{train}}, j \in D_{\text{train}} \}.$ 

• Structure Preserving Contrastive Learning. To preserve the structural information of the graphical views, we perform intra-modality contrastive learning among the graphical views and their corrupted views. Concretely, we corrupt each of the graphical view  $g_i$  with the edge dropping operation to obtain its corrupted view  $g'_i$ . The positive pairs for structure-preserving purpose are then designed as  $\mathcal{I}^+_{\text{preserve}} = \{\langle \phi_2(g_i), \phi_2(g'_i) \rangle | i \in D_{\text{train}} \}$ . The negative pairs for structure preserving designed as  $\mathcal{I}^-_{\text{preserve}} = \{\langle \phi_2(g_i), \phi_2(g'_i) \rangle | i \notin D_{\text{train}} \}$ .

**Finetuning with the Expert Soft Signal.** To help the backbone LLMs to digest the graphical views, we tune the LLMs with the expert soft signal using supervised finetuning. The prompt for finetuning consists of the system prompt, retrieved knowledge where the expert encoded graphical view is contained using a token embedding, and task prompt, which is illustrated in Figure 5 in the Appendix. **Inference.** After the finetuning stage, we used

**Inference.** After the finetuning stage, we used the tuned models to generate codes using the soft prompting template as illustrated in Figure 5 in the Appendix.

# **3** Experiments

- **RQ1** Does the proposed CodeGRAG offer performance gain against the base model?
- **RQ2** Does the proposed graph view abstract more informative knowledge compared with the raw code block?
- **RQ3** Can soft prompting enhance the capability of the backbone LLMs? Does finetuning with the soft prompting outperforms the simple supervised finetuning?
- **RQ4** Are the proposed pretraining objectives for the GNN expert effective?
- **RQ5** What is the impact of each of the components of the graphical view?
- RQ6 How is the compatibility of the graphical 368 view? 369

Model	Retrieved Knowledge		Python
GPT-3.5-Turbo	N/A	57.93	71.95
	Code Block (Nashid et al., 2023; Lu et al., 2022)		72.56
	Meta-Graph		72.56
	(Multi-Lingual) Code-Block (Nashid et al., 2023; Lu et al., 2022)	62.20	70.12
	(Multi-Lingual) Meta-Graph	64.02	77.44
GPT-4omini	N/A	63.41	78.66
	Code Block (Nashid et al., 2023; Lu et al., 2022)	65.24	78.66
	Meta-Graph		79.88
	(Multi-Lingual) Code-Block (Nashid et al., 2023; Lu et al., 2022)	65.85	79.27
	(Multi-Lingual) Meta-Graph	67.07	80.49

Table 1: Results of Hard Meta-Graph Prompt on Humaneval-X. (Pass@1)

Table 2: Results of Soft Prompting. (Pass@1)	l I

Model	Finetune	CodeForce (C++)	APPS (Python)
	N/A	12.83	5.09
Gemma 7b	SFT	14.76	21.09
	Soft Prompting	19.13	26.15
Llama2 13b	N/A	9.61	7.29
	SFT	11.88	12.06
	Soft Prompting	13.62	12.74
CodeLlama 7b	N/A	5.20	24.41
	SFT	9.87	26.15
	Soft Prompting	11.09	30.26

### 3.1 Setup

In this paper, we evaluate CodeGRAG with the widely used HumanEval-X (Zheng et al., 2023) dataset, which is a multi-lingual code benchmark and CodeForce dataset in which we collect real-world programming problems from codeforces<sup>1</sup> website. For CodeForce dataset we include problems categorized by different difficulty levels corresponding to the website and select 469 problems of difficulty level A for testing. We use greedy decoding strategy to do the generation. The evaluation metric is Pass@1. More details of the retrieval pool and the finetuning setting can be found in Section A in the Appendix.

### 3.2 Main Results

The main results are summarized in Table 1 and Table 2. From the results, we can draw the following conclusions.

**RQ1.** The proposed CodeGRAG could offer performance gain against the base model, which validates the effectiveness of the proposed graphical retrieval augmented generation for code generation framework.

**RQ2.** The model informed by the meta-graph (CodeGRAG) could beat model informed by the raw code block. From the results, we can see that the proposed graph view could summarize the use-ful structural syntax information and filter out the noises, which could offer more informative knowl-edge hints than the raw code block. In addition, inserting the intermediate representations of codes into the prompt can stimulate the corresponding programming knowledge of LLMs.

**RQ3.** From Table 2, we can see that finetuning with the expert soft prompting could offer more performance gain than that brought by simple supervised finetuning. This validates the effectiveness of the designed pretraining expert network and the technique of finetuning with soft prompting, which injects the programming domain knowledge into the LLMs parameters and informs the models with the structural information for gap filling.

# **3.3** Impacts of the pretraining objectives for the expert GNN (RQ4)

To study the effectiveness of the proposed pretraining objectives for the expert GNN, we remove each

<sup>&</sup>lt;sup>1</sup>https://codeforces.com/

Model	Finetune	CodeForce (C++)	APPS (Python)
Gemma 7b	Soft Prompting	19.13	26.15
	w/o Alignment	7.88	28.58
	w/o Structure-Preserving	11.70	21.50
Llama2 13b	Soft Prompting	13.62	12.74
	w/o Alignment	11.79	10.76
	w/o Structure-Preserving	5.50	11.09
CodeLlama 7b	Soft Prompting	11.09	30.26
	w/o Alignment	10.92	29.45
	w/o Structure-Preserving	10.66	26.59

Table 3: Ablation studies on the GNN pretraining losses.

objective to yield different expert GNNs. The results are in Table 3.

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

From the results, we could see that both the Alignment and the Structure Preserving contribute to the expressiveness of the expert GNN model. The alignment pretraining objective helps to promote the alignment among natural language, programming language, and their graphical views. The structure preserving objective helps to preserve the innate data-flows and control-flows information of code blocks. The two objectives collaborate with each other to yield expressive programming domain knowledge GNN expert model, which encodes external programming knowledge and injects the knowledge into LLMs parameters.

# 3.4 Impacts of the Components of the Graphical View (RQ5)

In this section, we adjust the inputs of the graphical components to the LLMs. Concretely, we study the information contained in node names, edge names, and the topological structure. The results are presented in Table 4.

Table 4:	The impacts o	f the grap	h components.
----------	---------------	------------	---------------

Datasets	Python	C++
Edge Type Only	73.78	61.59
Edge Type + Node Type	75.61	59.76 59.15
Edge Type + Topological	77.44	64.02

The edge type refers to the type of flows between 438 operands (child, read, write, etc.), the node type 439 refers to the type of operands (DeclStmt, temp, 440 441 etc.), the node name refers to the name of the intermediate variables, and the topological information 442 refers to the statistics of the concrete numbers of 443 different types of edges. From the results, we can 444 observe that 1) the edge features matter the most 445

in constructing the structural view of code blocks for enhancement, 2) the type of nodes expresses the most in representing operands information, and 3) the overall structure of the graphical view also gives additional information.

# 3.5 Compatibility Discussion of the Graphical Views(RQ6)

Despite the effectiveness of the proposed graphical views to represent the code blocks, the flexibility and convenience of applying the graphical views extraction process is important for wider application of the proposed method. In this section, we discuss the compatibility of CodeGRAG.

First of all, the extraction process of all the graphical views are front-end. Therefore, this extraction process applies to a wide range of code, even error code. One could also use convenient tools to reformulate the code and improve the pass rate of the extraction process.

In addition, we give the ratio of generated results that can pass the graphical views extraction process, which is denoted by Extraction Rate. The Pass@1 and the Extraction Rate of the generated results passing the graphical extraction process are given in Table 5.

Table 5: The extraction rate of the generated results passing the graphical extraction process.

Generated Codes Pass@1 Extraction Ra	ate
(C++) Code-RAG         62.20         92.07           (C++) CodeGRAG         64.02         92.68           (Python) Code-RAG         71.95         91.46           (Python) CodeGRAG         77.44         96.95	

From the results, we could see that the extraction rates are high for codes to pass the graphical views extraction process, even under the situation where the Pass@1 ratios of the generated results are low.

474

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

- 477

- 478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

504

505

This indicates that the application range of the proposed method is wide. In addition, as the code RAG also offers performance gains, one could use multiple views as the retrieval knowledge.

#### 4 **Related Work**

LLMs for NL2Code. The evolution of the Natural Language to Code translation (NL2Code) task has been significantly influenced by the development of large language models (LLMs). Initially, general LLMs like GPT-J (Radford et al., 2023), GPT-NeoX (Black et al., 2022), and LLaMA (Touvron et al., 2023a), despite not being specifically tailored for code generation, showed notable NL2Code capabilities due to their training on datasets containing extensive code data like the Pile (Gao et al., 2020) and ROOTS (Laurençon et al., 2022). To further enhance these capabilities, additional pretraining specifically focused on code has been employed. PaLM-Coder, an adaptation of the PaLM model (Chowdhery et al., 2023), underwent further training on an extra 7.8 billion code tokens, significantly improving its performance in code-related tasks. Similarly, Code LLaMA (Roziere et al., 2023) represents an advancement of LLaMA2 (Touvron et al., 2023b), benefiting from extended training on over 500 billion code tokens, leading to marked improvements over previous models in both code generation and understanding. These developments underscore the potential of adapting generalist LLMs to specific domains like NL2Code through targeted training, leading to more effective and efficient code translation solutions.

507 Code Search. The code search methods can be summarized into three folds. Early methods uti-508 lize sparse search to match the query and codes 509 (Hill et al., 2011; Yang and Huang, 2017), which suffers from mismatched vocabulary due to the 511 gap between natural language and codes. Neural 512 methods (Cambronero et al., 2019; Gu et al., 2021) 513 then focus on mapping the query and codes into 514 a joint representation space for more accurate retrieval. With the success of pretrained language 516 models, many methods propose to use pretraining 517 tasks to improve the code understanding abilities 518 and align different language spaces. For example, 520 CodeBERT (Feng et al., 2020) is pretrained on NL-PL pairs of 6 programming languages with the 521 masked language modeling and replaced token detection task. CodeT5 (Wang et al., 2021) supports both code-related understanding and generation 524

tasks through bimodal dual generation. UniXcoder (Guo et al., 2022) integrates the aforementioned pretraining tasks, which is a unified cross-modal pre-trained model. As retrieval augmented generation (RAG) shows its significance in promoting the quality of LLMs generation, works in code RAG start to accumulate. (Nashid et al., 2023; Lu et al., 2022) utilize the code blocks as the retrieved knowledge to inform the LLMs with similar code blocks for enhancement. (Zhou et al., 2022) uses the programming related document to serve as the retrieval content, injecting auxiliary external programming knowledge into the LLMs generation.

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

Code Representation. Early methods regard code snippets as sequences of tokens, assuming the adjacent tokens will have strong correlations. This line of methods (Harer et al., 2018; Ben-Nun et al., 2018; Feng et al., 2020; Ciniselli et al., 2021) take programming languages as the same with the natural language, using language models to encode the code snippets too. However, this ignoring of the inherent structure of codes leads to a loss of expressiveness. Methods that take the structural information of codes into consideration then emerge. Mou et al. (2016) used convolution networks over the abstract syntax tree (AST) extracted from codes. Alon et al. (2019) encoded paths sampled from the AST to represent codes. Further exploration into the graphical representation of codes (Allamanis et al., 2017) is conducted to better encode the structures of codes, where more intermediate states of the codes are considered.

#### 5 Conclusion

Despite the expanding role of LLMs in code generation, there are inherent challenges pertaining to their understanding of code syntax. General large language models trained mainly on sequentialbased natural language cannot well understand the structural-based programming language, e.g., the branching and jumping in codes. This paper proposes an effective way to build a graphical view of codes to better inform LLMs for code generation. To take the challenging structural graphical knowledge into LLMs, a meta-graph prompt is proposed for tuning-free models and a soft-prompting technique is proposed to inject the structural programming domain knowledge into the parameters of LLMs. By integrating external structural knowledge, CodeGRAG enhances LLMs' comprehension of code syntax and empowers them to generate code with improved accuracy and fluency.

674

675

676

677

678

625

626

# 576 Limitations

577In this paper, we propose a graphical retrieval aug-578mented generation method that can offer enhanced579code generation. Despite the efficiency and effec-580tiveness, there are also limitations within this work.581For example, dependency on the quality of the ex-582ternal knowledge base could be a potential concern.583The quality of the external knowledge base could584be improved with regular expression extraction on585the noisy texts and codes.

### Acknowledgments

This document has been adapted by Emily All-588 away from the instructions for earlier ACL and 589 NAACL proceedings, including those for NAACL 2024 by Steven Bethard, Ryan Cotterell and Rui Yan, ACL 2019 by Douwe Kiela and Ivan Vulić, NAACL 2019 by Stephanie Lukin and Alla Roskovskaya, ACL 2018 by Shay Cohen, Kevin Gimpel, and Wei Lu, NAACL 2018 by Margaret Mitchell and Stephanie Lukin, BibT<sub>F</sub>X suggestions for (NA)ACL 2017/2018 from Jason Eisner, ACL 2017 by Dan Gildea and Min-Yen Kan, NAACL 2017 by Margaret Mitchell, ACL 2012 by Maggie Li and Michael White, ACL 2010 by Jing-Shin Chang and Philipp Koehn, ACL 2008 by Johanna D. Moore, Simone Teufel, James Allan, and Sadaoki Furui, ACL 2005 by Hwee Tou Ng and Kemal Oflazer, ACL 2002 by Eugene Charniak and Dekang Lin, and earlier ACL and EACL formats written by several people, including John Chen, Henry S. Thompson and Donald Walker. Addi-606 tional elements were taken from the formatting instructions of the International Joint Conference on Artificial Intelligence and the Conference on Computer Vision and Pattern Recognition. 610

# References

611

612

613

615

617

618

619

624

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Alfred V Aho, Monica S Lam, Ravi Sethi, and Jeffrey D Ullman. 2006. Compilers: Principles techniques and tools. 2007. *Google Scholar Google Scholar Digital Library Digital Library*.
- Miltiadis Allamanis, Marc Brockschmidt, and Mahmoud Khademi. 2017. Learning to represent programs with graphs. *arXiv preprint arXiv:1711.00740*.

- Frances E Allen. 1970. Control flow analysis. *ACM Sigplan Notices*, 5(7):1–19.
- Uri Alon, Meital Zilberstein, Omer Levy, and Eran Yahav. 2019. code2vec: Learning distributed representations of code. *Proceedings of the ACM on Programming Languages*, 3(POPL):1–29.
- Tal Ben-Nun, Alice Shoshana Jakobovits, and Torsten Hoefler. 2018. Neural code comprehension: A learnable representation of code semantics. *Advances in Neural Information Processing Systems*, 31.
- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, et al. 2022. Gpt-neox-20b: An open-source autoregressive language model. *arXiv preprint arXiv:2204.06745*.
- Jose Cambronero, Hongyu Li, Seohyun Kim, Koushik Sen, and Satish Chandra. 2019. When deep learning met code search. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 964–974.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Matteo Ciniselli, Nathan Cooper, Luca Pascarella, Denys Poshyvanyk, Massimiliano Di Penta, and Gabriele Bavota. 2021. An empirical study on the usage of bert models for code completion. In 2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR), pages 108–119. IEEE.
- 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*.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.
- Jian Gu, Zimin Chen, and Martin Monperrus. 2021. Multimodal representation for neural code search. In 2021 IEEE International Conference on Software Maintenance and Evolution (ICSME), pages 483– 494. IEEE.
- 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*.
- Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey

788

Svyatkovskiy, Shengyu Fu, et al. 2020. Graphcodebert: Pre-training code representations with data flow. *arXiv preprint arXiv:2009.08366*.

679

684

692

693

701

703

704

705

707

709

710

711

714

715

716

718

719

720

721

722

723

724

725

726

727

728

729

730

731

733

734

- Jacob A Harer, Louis Y Kim, Rebecca L Russell, Onur Ozdemir, Leonard R Kosta, Akshay Rangamani, Lei H Hamilton, Gabriel I Centeno, Jonathan R Key, Paul M Ellingwood, et al. 2018. Automated software vulnerability detection with machine learning. *arXiv preprint arXiv:1803.04497*.
- Emily Hill, Lori Pollock, and K Vijay-Shanker. 2011. Improving source code search with natural language phrasal representations of method signatures. In 2011 26th IEEE/ACM International Conference on Automated Software Engineering (ASE 2011), pages 524– 527. IEEE.
- Xue Jiang, Zhuoran Zheng, Chen Lyu, Liang Li, and Lei Lyu. 2021. Treebert: A tree-based pre-trained model for programming language. In *Uncertainty in Artificial Intelligence*, pages 54–63. PMLR.
- Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, et al. 2022. The bigscience roots corpus: A 1.6 tb composite multilingual dataset. Advances in Neural Information Processing Systems, 35:31809–31826.
- Ting Long, Yutong Xie, Xianyu Chen, Weinan Zhang, Qinxiang Cao, and Yong Yu. 2022. Multi-view graph representation for programming language processing: An investigation into algorithm detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 5792–5799.
- 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.*
- Lili Mou, Ge Li, Lu Zhang, Tao Wang, and Zhi Jin. 2016. Convolutional neural networks over tree structures for programming language processing. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30.
- Noor Nashid, Mifta Sintaha, and Ali Mesbah. 2023. Retrieval-based prompt selection for code-related few-shot learning. In 2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE), pages 2450–2462. IEEE.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023.
   Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. arXiv preprint arXiv:2308.12950.

- Bo Shen, Jiaxin Zhang, Taihong Chen, Daoguang Zan, Bing Geng, An Fu, Muhan Zeng, Ailun Yu, Jichuan Ji, Jingyang Zhao, et al. 2023. Pangu-coder2: Boosting large language models for code with ranking feedback. *arXiv preprint arXiv:2307.14936*.
- Yizhou Sun and Jiawei Han. 2013. Mining heterogeneous information networks: a structural analysis approach. ACM SIGKDD explorations newsletter, 14(2):20–28.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. 2021. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. *arXiv preprint arXiv:2109.00859*.
- Yangrui Yang and Qing Huang. 2017. Iecs: Intentenforced code search via extended boolean model. *Journal of Intelligent & Fuzzy Systems*, 33(4):2565– 2576.
- Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, et al. 2023. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. *arXiv preprint arXiv:2303.17568*.
- 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** Implementation Details

For the size of retrieval pool, we use 11,913 C++ code snippets and 2,359 python code snippets. Due to the limited access, we do not use a large retrieval corpus for our experiment, which can be enlarged by other people for better performance. We also attach the graph extraction codes for both languages and all other expeirment codes here: https://anonymous.4open.science/r/Code-5970/

For the fintuning details, the learning rate and weight decay for the expert GNN training is 0.001 and 1e-5, repectively. We apply 8-bit quantization and use LoRA for parameter-efficient fine-tuning. The rank of the low-rank matrices in LoRA is uniformly set to 8, alpha set to 16, and dropout is set to 0.05. The LoRA modules are uniformly applied
to the Q and V parameter matrices of the attention
modules in each layer of the LLM. All the three
models are optimized using the AdamW optimizer.
For the CodeContest dataset, totally 10609 datapoints are used, and for APPS dataset, 8691 data
samples are used to train the model.

# **B Prompt Template**

```
      System Prompt

      Please continue to complete the [lang] function according to the requirements and function declarations. You are not allowed to modify the given code and do the completion only.\n

      Retrieved Knowledge

      The syntax graph of a similar code might be:\n [composed syntax graph desciption]

      You can refer to the above knowledge to do the completion. \n

      Problem

      The problem:\n [problem prompt]
```

Figure 4: Hard meta-graph prompt.



Figure 5: Soft prompting.