

Domain Adaptive Code Completion via Language Models and Decoupled Domain Databases

Ze Tang¹, Jidong Ge^{1*}, Shangqing Liu^{2*}, Tingwei Zhu¹, Tongtong Xu³, Liguang Huang⁴, Bin Luo¹

¹National Key Laboratory for Novel Software Technology, Nanjing University, China

²Nanyang Technological University, Singapore

³Huawei Software Engineering Application Technology, China

⁴Department of Computer Science, Southern Methodist University, USA

{zetang, tingweizhu33}@smail.nju.edu.cn, {gjd, luobin}@nju.edu.cn, liu.shangqing@ntu.edu.sg, xutongtong9@huawei.com, lghuang@lyle.smu.edu

Abstract—Large Language Models (LLMs) have demonstrated remarkable performance in code completion. However, due to the lack of domain-specific knowledge, they may not be optimal in completing code that requires intensive domain knowledge for example completing the library names. Although there are several works that have confirmed the effectiveness of fine-tuning techniques to adapt language models for code completion in specific domains. They are limited by the need for constant fine-tuning of the model when the project is in constant iteration.

To address this limitation, in this paper, we propose k NM-LM, a retrieval-augmented language model (R-LM), that integrates domain knowledge into language models without fine-tuning. Different from previous techniques, our approach is able to automatically adapt to different language models and domains. Specifically, it utilizes the in-domain code to build the retrieval-based database decoupled from LM, and then combines it with LM through Bayesian inference to complete the code. The extensive experiments on the completion of intra-project and intra-scenario have confirmed that k NM-LM brings about appreciable enhancements when compared to CodeGPT and UnixCoder. A deep analysis of our tool including the responding speed, storage usage, specific type code completion, and API invocation completion has confirmed that k NM-LM provides satisfactory performance, which renders it highly appropriate for domain adaptive code completion. Furthermore, our approach operates without the requirement for direct access to the language model’s parameters. As a result, it can seamlessly integrate with black-box code completion models, making it easy to integrate our approach as a plugin to further enhance the performance of these models.

Index Terms—domain adaptive code completion, retrieval-augment language model

I. INTRODUCTION

Large language models (LLM) [1]–[5], have achieved state-of-art performance in code completion, and some of them have been successfully used as the auto-completion plugin (e.g. GitHub Copilot [6] and ChatGPT [7]) in modern Integrated Development Environment (IDE). Nevertheless, code in IDE has distinctive domain-specific features, such as imported third-party libraries and intra-project references. These features can vary significantly across domains and are continuously updated. Thus, LLMs, which are developed for

completing general code, may not perform well when applied to new domains such as personal projects, as noted in prior works [8]–[10].

A standard practice to adapt a pre-trained model to a new domain (i.e., domain adaptation) is model fine-tuning [11]. By fine-tuning a language model on a new domain, the model’s performance on that domain can be improved without having to retrain the entire model from scratch. Nevertheless, fine-tuning can be infeasible for large language models. For instance, language models such as GPT-3 or GPT-4, are often deployed as black-box systems. The black-box limitation renders the **parameters inaccessible** to users, thereby preventing them from fine-tuning. More importantly, domain-specific features in code are subject to **frequent changes**. This is due to the continuous development of projects and the incessant updates of third-party libraries utilized in the code. Consequently, the employment of fine-tuning for domain-intensive code completion may not be a feasible option.

In recent years, a series of retrieval-augmented language models (R-LM) [12]–[15] have emerged to address this challenge. The key ingredient of R-LMs is their ability to utilize the domain database at test time without having to rely on the information encoded in the model’s weights only. In these models, the retrieval component first searches for nearest neighbor examples in an external datastore (e.g., code from the same project); then, the base model references these examples during the prediction. One prominent example of such a retrieval-based model is k -nearest neighbors language model (k NN-LM), which predicts a token by linearly interpolating the base LM’s output with a non-parametric nearest neighbor distribution. This distribution is constructed by searching for the k -nearest neighbors (k NN) in the datastore and weighing them according to their distance to the current test context. Notably, k NN-LM requires a large datastore that stores each token in the domain and is sensitive to the manual-selected interpolated weight for soft voting, making them challenging to apply to code completion. Developers may not have the resources or patience to find suitable hyper-parameters for their own code.

In this paper, we introduce a plug-and-play auto-completion

*Corresponding author.

solution that does not need careful configuration, the k -nearest mistakes language model (k NM-LM). Unlike prior R-LMs, k NM-LM retrieves information from a separate database that is decoupled from the language model. This is because the weights in the language model already contain some content knowledge about the program language, such as grammar and built-in methods. Thus, storing this information redundantly in the database would be both wasteful and meaningless. To overcome this challenge, we split the domain knowledge into two subsets: what the language model knows and what it doesn't. The latter is saved as a decoupled database, which is a subset of the domain database that stores only tokens that cannot be correctly predicted by the language model (i.e., mistakes collection). Decoupling offers two primary benefits. Firstly, the decoupled database requires less storage and search resources. Secondly, the decoupled database and language model are two distinct and non-intersecting systems. As a result, we do not need to manually select the interpolated weight but can leverage statistical approaches, such as Bayesian inference, to automatically combine the database and language model from a statistical perspective.

Specifically, we start by creating a datastore at the token level. This datastore only contains code tokens that the language model fails to predict correctly. We then retrieve tokens from the datastore that are most like the code being completed and normalize them into a distribution. Next, we adopt the error rate of the language model on the original domain database as the prior probability and assess whether the language model can predict the code tokens prior to the completion position as new evidence (likelihood probability). The posterior probability is then calculated and utilized to merge the distributions obtained from the datastore and the language model.

We conducted an evaluation of k NM-LM on intra-project and intra-scenario code completion, revealing the performance improvement in both CodeGPT and UnixCoder through our proposed approach. We further performed validation on three key aspects that could affect the user experience when using k NM-LM as a code completion plugin. These included completion speed and space utilization, the performance in completing specific code types and completing specific API calls from third-party libraries. Our experimental results demonstrate that k NM-LM outperforms several baseline models and is better suited for domain adaptive code completion than other R-LMs. Additionally, k NM-LM does not require access to the language model's weights, making it suitable for situations where only black-box access to the language model is available. In summary, the primary contributions of this paper are as follows:

- We propose k NM-LM, a retrieval-augmented language model (R-LM), that can be used for domain adaptive code completion without fine-tuning. All code, data, and results can be found at our anonymous repository¹.

- Different from previous R-LMs, k NM-LM retrieves from a decoupled database and uses Bayesian inference to interpolate between the database and language model.
- Experimental results on intra-project and intra-scenario code completion demonstrate that k NM-LM improves the performance of both CodeGPT and UniXcoder.
- We also investigate potential issues that may impact user experience when using k NM-LM as a code completion plugin. Our results indicate that k NM-LM performs satisfactorily across all considered factors.

II. BACKGROUND

A. k -Nearest Neighbors Language Model

k -Nearest Neighbors Language Model (k NN-LM) [15] is a retrieval-augmented language model (R-LM), in which uses a nearest neighbor retrieval mechanism to augment the pre-trained language model, without any additional fine-tuning. Given a sequence of tokens $c_t = (x_1, x_2, \dots, x_{t-1})$, autoregressive language models, such as GPT-3, estimate $p_{LM}(y|c_t)$, the probability distribution over the next token.

Datastore. Let $f_e(\cdot)$ be the encoding function that maps a context c_t to a fixed-length vector representation computed by the pre-trained language model. First, k NN-LM generates the key-value datastore by running the encoding function $f_e(\cdot)$ over a corpus \mathcal{D} . Typically, each key k_t is the vector representation $f_e(c_t)$ of the context c_t , and each value v_t is the next token x_t for the context c_t , as:

$$(K, V) = \{(f_e(c_t), x_t) | (c_t, x_t) \in \mathcal{D}\} \quad (1)$$

Inference. At test time, given the input context c_t , the pre-trained language model generates the probability distribution over the next token $p_{LM}(y|c_t)$ and the context representation $f_e(c_t)$. k NN-LM queries the datastore to retrieve k -nearest keys according to a vector distance function $d(\cdot, \cdot)$. Then, the model computes a softmax over the (negative) distances, which gives a distribution over the next token:

$$p_{kNN}(y|c_t) \propto \sum_{(k_t, v_t) \in \mathcal{D}} \mathbb{1}_{y=v_t} \exp(-d(k_t, f_e(c_t))) \quad (2)$$

The prediction is then interpolated with the prediction from the language model:

$$p(y|c_t) = \lambda p_{kNN}(y|c_t) + (1 - \lambda) p_{LM}(y|c_t) \quad (3)$$

where λ is a hyper-parameter that ranges from 0 to 1 and needs to be carefully selected.

B. Bayesian Inference

Bayesian inference is a fundamental concept in probability theory that allows us to update our beliefs about an event considering new evidence or information. It states that the probability of a hypothesis or event (A) given some observed evidence (C) can be calculated using Bayes' rule:

$$p(A|C) \propto p(A) \times p(C|A) \quad (4)$$

where $p(A)$ is the prior probability of A, $p(C|A)$ is the probability of observing C given that A is true (the likelihood), and

¹https://github.com/zetang94/ASE2023_kNM-LM

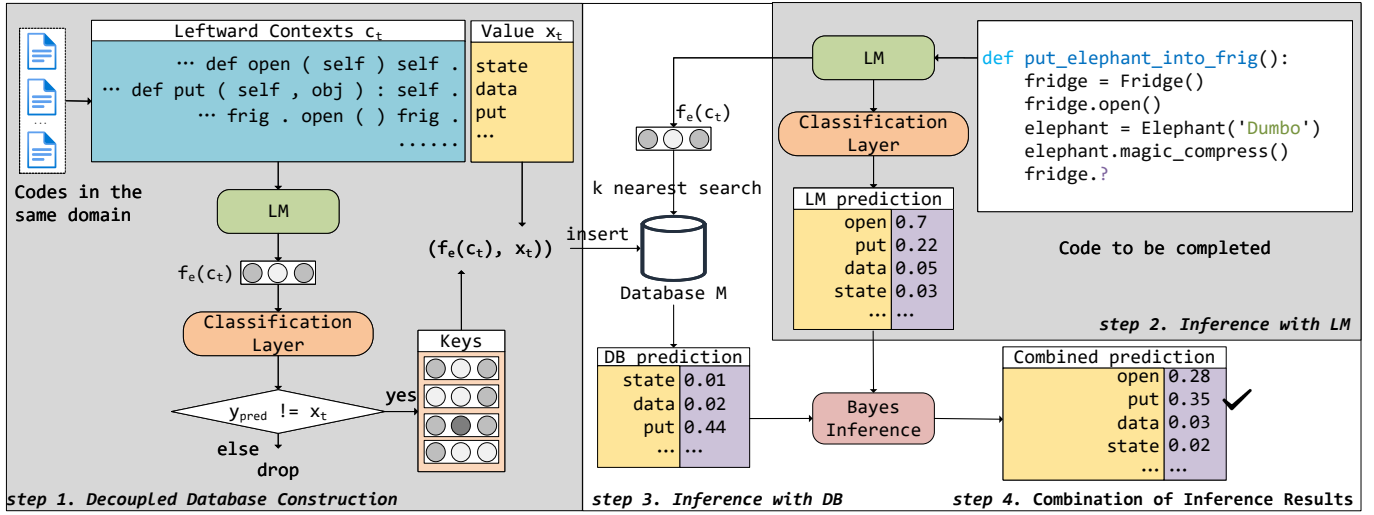


Fig. 1. The primary workflow of our proposed k NM-LM. A grey box indicates that this part needs to access the (black-box) language model.

$p(A|C)$ is the posterior probability of A given C . Bayesians consider $\theta = p(A)$ and use $p(\theta)$ to describe the distribution of $p(A)$. Specifically, if $p(\theta) = \text{Beta}(a, b)$, and event C indicates that A occurs α times in N new observations. Then we have the following conclusions:

$$\begin{aligned}
 \textbf{Likelihood:} \quad & p(C|\theta) = \theta^{(\alpha-1)}(1-\theta)^{(N-\alpha-1)} \\
 \textbf{Posterior:} \quad & p(\theta|C) = \text{Beta}(a+\alpha, b+N-\alpha) \\
 \textbf{Expectation:} \quad & p(A|C) = E(p(\theta|C)) = \frac{(a+\alpha)}{(a+b+N)}
 \end{aligned} \tag{5}$$

III. APPROACH

The code completion task aims to predict the next token or the next line of code given a known code sequence $c_t = (x_1, x_2, \dots, x_{t-1})$. For domain-adaptive code completion, we assume that some in-domain code can be obtained (e.g., code already written in current project or code from other projects that are developed for the same application scenario). Retrieval-augmented language models (R-LMs) can be employed to enhance the performance of general (black-box) LMs for completing domain-intensive code. The primary workflow of our proposed method, k -nearest mistakes language model (k NM-LM), is illustrated in Fig. 1.

Initially, we analyze the R-LMs from a statistical perspective (section III-A), discovering that manual selection of combination coefficient in previous R-LMs can be avoided by decoupling the retrieval module and LM. Next, we discuss the construction of the decoupled retrieval database (section III-B). During the inference phase, we separately predict the next token probability using the LM (section III-C) and the retrieval module (section III-D). Finally, we employ Bayesian inference to combine the inferences (section III-E).

A. Analysis of R-LM from a Statistical Perspective

We initiate our discussion by defining two distinct events associated with predicting the next token using LM:

Event E : The LM correctly predicts the next token of c_t . $\text{leftmargin} = *$

- Set notation: $E = \{x_t | \arg \max(p_{LM}(y|c_t)) = x_t\}$
- Description: The predicted token y with the highest probability matches the ground truth x_t .

Event E' : The LM incorrectly predicts the next token of c_t . $\text{leftmargin} = *$

- Set notation: $E' = \{x_t | \arg \max(p_{LM}(y|c_t)) \neq x_t\}$
- Description: The predicted token y with the highest probability does not match the ground truth x_t .

The primary goal of retrieval-augmented approaches is to enhance the base LM's performance through error correction. Ideally, for tokens belonging to Event E , the retrieval module should not impact the final prediction, allowing the predictions from the LM to be utilized directly. Conversely, for tokens belonging to Event E' , the retrieval module aims to “correct” the erroneous prediction by searching for previously similar mistakes made by the LM and using the corresponding ground truths as predictions. The idea is similar to the way that humans use a “collection of mistakes” to remind themselves not to make the same mistake twice. We can use the total probability theorem to formulate the retrieval-argument process, as:

$$\begin{aligned}
 p(y|c_t) &= p(y|c_t, E' \cup E) \\
 &= p(y|c_t, E')p(E'|c_t) + p(y|c_t, E)p(E|c_t) \\
 &= p_{\text{Retrieval}}(y|c_t)p(E'|c_t) + p_{LM}(y|c_t)p(E|c_t)
 \end{aligned} \tag{6}$$

From Equation 6, it is essential to decouple the retrieval module's abilities from the LM, as they predict tokens that belong to distinct events. Previous R-LMs do not decouple the retrieval module from the LM, necessitating the use of an ensemble strategy [16] to combine their strengths. However, the ensemble strategy requires manual and empirical determination of the combination coefficient. Inappropriate coefficient may even lead to a decline in performance. By decoupling

the retrieval module, the combination coefficient $p(E'|c_t)$ can be computed mathematically. Moreover, $p(E'|c_t)$ is associated with both the domain-specific context and the LM, facilitating the adaptation of code completion tasks with different LMs across various domains. In the subsequent section, we will discuss the construction of the decoupled database.

B. Step 1. Decoupled Database Construction

Given the known in-domain code \mathcal{D} , we build the decoupled database \mathcal{M} as a key-value store. The values consist of tokens associated with event E' , while the corresponding keys are the embedded vectors of their preceding contexts:

$$\mathcal{M} = \{(f_e(c_t), x_t) | (c_t, x_t) \in \mathcal{D}, x_t \in E'\} \quad (7)$$

Here, $f_e(\cdot)$ denotes a function that maps a code sequence to a fixed-length representation with dimension d . In CodeGPT, for example, $f_e(\cdot)$ might be the output of the last self-attention layer. It is worth noting that the decoupled database is considerably smaller than the one constructed in k NN-LM, which saves storage space and search time. The reduction ratio is inversely proportional to the LM's error rate err in \mathcal{D} , as:

$$p(E') = err = \frac{||\mathcal{M}||}{||\mathcal{D}||} \quad (8)$$

$||\cdot||$ signifies the set size. Importantly, the error rate also represents the probability of event E' occurring in \mathcal{D} . Assuming that the code to be completed are independent and identically distributed (i.i.d.) with code in \mathcal{D} , we can use the error rate as the prior probability $p(E')$ at inference time.

C. Step 2. Inference with LM

Language models are treated as a black-box system in our approach, as we do not need to access their parameters. We input the code context c_t into the LM and obtain the probability of the next token $p_{LM}(y|c_t)$. Additionally, we need to acquire the embedded context vector $f_e(c_t)$. All of this can be achieved by utilizing the corresponding APIs provided by the LM. For example, GPT-3 offers APIs that support embedding context [17] and obtaining the probabilities of next tokens based on the context [18]. After obtaining the embedded code context, we use it as the key to search the database.

D. Step 3. Inference with DB

We follow Khandelwal et al. [15] to generate the probability of the next token from the datastore. With the embedded context vector $f_e(c_t)$ generated by the LM, the retrieval module searches for keys similar to $f_e(c_t)$ in the datastore using k -nearest neighbors search. Similarity is defined by the vector distance $d(\cdot, \cdot)$ (we employ Euclidean distance in our experiments). The model then computes the distribution over the retrieved results as follows:

$$p_{kNM}(y|c_t) \propto \sum_{(k_i, v_i) \in \mathcal{M}} \mathbb{1}_{y=v_i} \exp(-d(k_i, f_e(c_t))) \quad (9)$$

$\mathbb{1}$ is an indicator function, which aggregates probability mass for each vocabulary item across all its occurrences in the

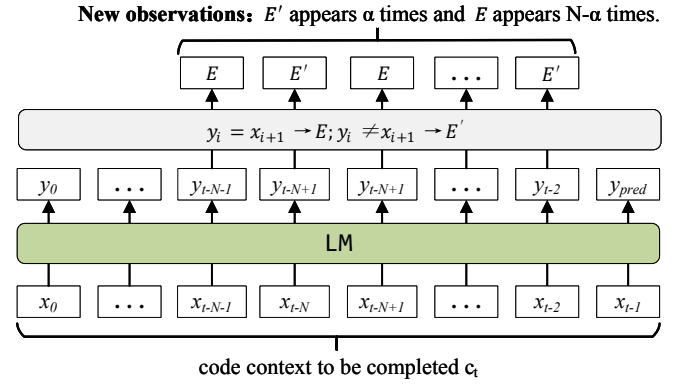


Fig. 2. Calculation of $p(c_t|E')$. As the code context to be completed is known, we leverage the LM to make predictions for each token within the context and determine if the predictions are accurate. These new observations can be used to calculate likelihood probability and to update the combined beliefs of the LM and retrieval module. Noted that for decoder-only models like GPT-3, this process does not incur additional computational costs.

retrieved targets, and probabilities for items not appearing in the retrieved targets are set to zero.

E. Step 4. Combination of Inference Results

Combination of inference results from the LM and database is achievable by utilizing Equation 6. Nevertheless, calculating $p(E'|c_t)$ directly poses challenges, as it represents the probability about the subsequent token of c_t , which remains unknown. We employ Bayes' rule to solve this, as:

$$p(E'|c_t) \propto p(E') \times p(c_t|E') \quad (10)$$

The prior probability, denoted as $p(E')$, has been computed using Equation 8. The likelihood probability $p(c_t|E')$ represents the probability of observing c_t given the event E' . This likelihood probability indicates that the combination coefficient is also influenced by the code context being completed. Given that we already know c_t , for each token $x_i \in c_t$, we can predict the token y_i by inputting the context $x_{0:i-1}$ into the language model. If $y_i = x_{i+1}$, we label it as event E , otherwise, it is classified as event E' . Bayesian analysis treats these newly occurring events as new observations, and they are utilized to update the prior probability.

In the context of parallel computing, our focus is on the latest $N - 1$ observations within c_t (can also be viewed as N -gram assumption). For illustration, as depicted in Fig. 2, we utilize $y_{t-N-1:t-2}$ as the newly considered observations. Assuming that event E' occurs α times, then event E occurs $N - \alpha$ times. The computation of the likelihood probability $p(c_t|E')$ is expressed as follows:

$$p(c_t|E') \approx p(x_{t-N-1:t-2}|E') = p(E')^\alpha (1 - p(E'))^{N-\alpha} \quad (11)$$

From the perspective of Bayesian inference, $p(E')$ is not a fixed number but a distribution represented by $\theta = p(E')$. Given that the distribution's expectation equals to err (Equation 8), it is reasonable to consider $p(\theta) \sim \text{Beta}(err \cdot N, (1 -$

$err) \cdot N)^2$. With the new evidence on c_t , we can deduce that the conditional distribution $p(\theta|c_t) \sim \text{Beta}(err * N + \alpha, (1 - err) * N + N - \alpha)$. As a result, the calculation of $p(y|c_t)$ can be formulated as:

$$p(y|c_t) = \lambda p_{kNM}(y|c_t) + (1 - \lambda) p_{LM}(y|c_t) \quad (12)$$

$$\lambda = p(E'|c_t) = \frac{1}{2} \left(\frac{\alpha}{N} + \frac{||\mathcal{M}||}{||\mathcal{D}||} \right)$$

Equation 12 incorporates two hyper-parameters: k and N . In this context, k is employed for the k -nearest neighbor search to compute $p_{kNM}(y|c_t)$, while N represents the number of observations utilized for calculating the likelihood probability.

The values of α , $||\mathcal{D}||$, and $||\mathcal{M}||$ in Equation 12 are determined by three key components pertaining to the completion process: the code context to be completed, the specific language model utilized, and the chosen database. Specifically, when the error rate of the LM on the database and code context c_t is elevated, λ is correspondingly large, leading the combined approach to favor the prediction outcome of the retrieval module. In contrast, when the model’s error rate on the database and code context c_t is low, λ is reduced, and the combination strategy tends to rely more on the prediction results generated by the LM.

IV. EXPERIMENTAL SETUP

To assess the efficacy of k NM-LM, we have formulated two research questions: leftmargin=*

- RQ1: What is the performance of k NM-LM for intra-project code completion?
- RQ2: What is the performance of k NM-LM for intra-scenario code completion?

Furthermore, in relation to RQ1, we have explored the ability of k NM-LM to complete specific types of code. Regarding RQ2, we have investigated the effectiveness of k NM-LM in completing lines that involve Android API calls, along with an ablation study.

A. Baselines

As our proposed k NM-LM is a retrieval-augmented framework, for code completion, we choose two state-of-art pre-trained models as the base model: leftmargin=*

- CodeGPT [1] is a GPT-based pre-trained model trained on the CodeSearchNet dataset [20] for code completion and generation tasks. Specifically, we use CodeGPT-small-java-adaptedGPT2 for Java code and CodeGPT-small-py-adaptedGPT2 for Python code.
- UniXcoder [2] is also pre-trained on the CodeSearchNet dataset and utilizes cross-modal content from code. It applies mask attention matrices with prefix adapters to control the model behavior. As UniXcoder is a cross-language model, we use UniXcoder-base for both Java and Python code.

²The choice of the Beta distribution is informed by its flexibility and suitability for modeling probabilities between 0 and 1 [19]. Additionally, we assume equal significance of prior beliefs and new observations, leading us to utilize the same N for the Beta distribution.

We also compare k NM-LM with other competitive retrieval-augmented frameworks, including: leftmargin=*

- **BM25** [14] is a block R-LM. It saves domain code as a document datastore, then searches the datastore to find code like the code to be completed based on the BM25 [21] algorithm. Finally, the search results are concatenated with the unfinished code and jointly fed into the language model.
- **ReACC** [13] is also a block R-LM. Different from BM25, it uses GraphCodeBERT [22] to train a retrieve model, and then encodes the domain code into a fixed length vector to store. ReACC only published their fine-tuned retrieval model on Python language, we only use it for Python language.
- **Hybrid** (BM25+ReACC) [13] is an ensemble approach based on BM25 and ReACC. It combines search results from both BM25 and ReACC, then uses the re-ranked results to generate the next code token. Due to the need of using ReACC, we only use it for Python language.
- **k NN-LM** [15] is a token R-LM, which saves each token in the domains code into a key-value datastore, and use k NN algorithm to retrieve similar code tokens.

B. Implementation Details

In order to enable comparison with model fine-tuning (section VII-B), we employ the pre-trained Code-GPT and UniXcoder models as the base models, rather than relying on the black-box GPT-3 model. We utilize Elasticsearch for BM25 algorithm. ReACC is implemented using the released “microsoft/reacc-py-retriever” as the retriever model. Hybrid uses the combination weight 0.9 from their original paper. In the case of k NN-LM and k NM-LM, we set k (the number of retrieved nearest neighbors) to 8 for RQ1 and 1024 for RQ2 and employ L2 distance as the distance function for neighbor retrieval. We set λ to 0.1 and window size N to 8. In k NM-LM, we let the retrieval module to participate in computing the likelihood probability, if y_i can be both correctly predicted by the LM and the retrieval module, we just ignore this token. Faiss-gpu [23] is used for accelerated k nearest neighbors search, same as Khandelwal et al. [24].

Metrics Accuracy, representing the ratio of correctly predicting the next token, serves as the evaluation metric for token completion, while the metric for line completion involves Exact Match accuracy (EM) and Levenshtein Edit Similarity (ES) [25]. To ensure fairness in the analysis of time and space complexity, all experiments are performed on a machine equipped with 2 NVIDIA 3090Ti-24GB GPU cards.

V. RQ1: INTRA-PROJECT CODE COMPLETION

A. Study Design

Auto code completion plugins are widely used in projects to streamline the coding process, with other files in the project often containing helpful information for code completion. Consequently, we have developed the intra-project code completion task. In this research question, we utilize Git commit history to determine the creation time of various methods in

TABLE I

STATISTICS OF INTRA-PROJECT COMPLETION DATASET. (·) INDICATES THE PROJECT’S COMMIT ID. THE TEST SET CONTAINS METHODS CREATED FROM THE OLD COMMIT TO THE LATEST COMMIT.

Repo Name	Database	Test
Froyo_Email	2.4M (56e48)	388K (56e48 → 508c9)
dropwizard	2.2M (04c45)	388K (04c45 → ad30b)
AmazeFileManager	1.9M (d01c5)	292K (d01c5 → 5c3cd)
rest-assured	1.8M (0c8f0)	274K (0c8f0 → 246ba)
logging-log4j1	1.4M (0d9e1)	353K (0d9e1 → e5c56)
feign	1.3M (6989b)	140K (6989b → d42fc)
reqquery	1.3M (a3f71)	164K (a3f71 → 1d6fa)
eureka	1.2M (22528)	184K (22528 → fc4e0)
galaxy	1.2M (0ed70)	166K (0ed70 → 8fcc6)
interview	1.1M (ed0d5)	164K (ed0d5 → 94be5)
android-priority-jobqueue	592K (0287d)	69K (0287d → f807c)
lottie-android	541K (c32f1)	88K (c32f1 → 42b48)
xUtils3	472K (0d10d)	35K (0d10d → 4f3c3)
Fragmentation	408K (0e972)	59K (0e972 → d5e79)
VirtualAPK	308K (e6174)	15K (e6174 → 01b73)
spring-boot-starter	181K (f0951)	16K (f0951 → a634e)
DiscreteScrollView	147K (49671)	15K (49671 → 9d979)
flow	108K (489b2)	18K (489b2 → 06eb0)
hover	259K (7e95c)	14K (7e95c → 03e44)
StickyHeaderListView	117K (00265)	6K (00265 → b871c)

the project. We then select the project snapshot at a specific commit as the basis for the written code used to build the database. Methods created after that commit serve as a test set to evaluate the performance of different approaches.

B. Dataset

We compile an intra-project code completion dataset by collecting 10 large projects (over 1M in size) and 10 small projects (under 1M in size) from the test set of Java-mid [26]. The dataset’s statistics are displayed in Table IV-B. Following the approach of Bogomolov et al. [9], we gather the Git history for each project and extract all method creation commits using Miner [27]. These methods are sorted by commit time, and we select the commit at which the first 80% of methods have been created. The project snapshot at this commit is used to build the database, while the remaining 20% of methods serve as the test set. Then, we separately evaluate the performance on each project.

C. Results

Token Completion Accuracy. As illustrated in Fig. 3, our proposed approach, *k*NM-LM, exhibits a significant improvement for all projects in the dataset. *k*NM-LM achieves the most substantial improvement on the “hover” project, with a 26.67% and 35.99% increase with CodeGPT and UniXcoder, respectively. The trend of the curve indicates that larger projects benefit more from using *k*NM-LM. This may be because large projects contain more code in the database, making it more likely to find useful information from the database. We also observe that the performance boost gained by using the prompt method (BM25) is inconsistent across different projects. In some cases, it may be worse than the base model, such as with CodeGPT on the “android-priority-jobqueue” project, where it decreases from 52.03% to 51.84%. This discrepancy may arise due to the absence of code in the database that closely

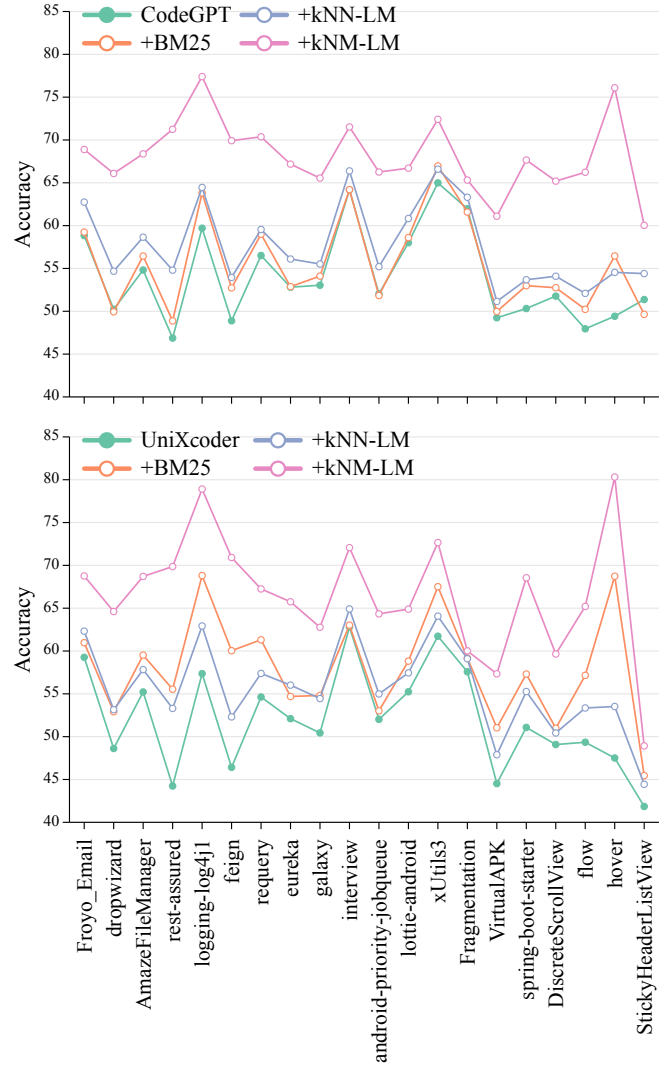


Fig. 3. Token completion accuracy on intra-project dataset. The abscissa is the project name, ordered from large to small based on the project size.

resembles the code requiring completion. When we merge the retrieved result with the input and feed them into the language model, there is a potential to introduce irrelevant noise that does not contribute to the accurate completion.

Analysis of Completing Different Code Token Types. Language models often face greater difficulties when dealing with certain code token types. For instance, predicting identifiers (such as variable or function names) is more challenging than predicting keywords (like “if” or “for”). However, developers primarily utilize code completion plugins to finalize intra-project APIs [28]. Hence, we follow previous works [29] to investigate the performance of models when predicting different types of tokens. Results are presented in Table II. From the results, we observe that all retrieval-augmented models can enhance the code completion performance of the base model in intra-project contexts, with *k*NM-LM showing the most significant improvement. Specifically, for completing punctuation, which is related to programming language syntax, *k*NM-LM improves the base models with increases

TABLE II
TOKEN COMPLETION ACCURACY OF DIFFERENT CODE TYPES ON INTRA-PROJECT DATASET. NUMBERS ARE SHOWN IN PERCENTAGE(%).

	Types	Percent	CodeGPT	+BM25	+ k NN-LM	+ k NM-LM	UniXcoder	+BM25	+ k NN-LM	+ k NM-LM
Large Projects	All	100.00	54.59	56.12	58.69	69.66	53.11	59.16	57.46	68.96
	Punctuation	34.56	64.97	66.74	68.10	86.07	67.24	73.40	70.70	84.90
	Identifier	30.63	38.37	40.72	41.48	47.68	38.43	46.69	41.42	47.33
	Operator	19.51	70.02	70.66	71.39	77.05	64.84	69.74	66.29	81.08
	Keyword	8.02	49.96	51.17	53.82	56.97	28.96	34.46	32.21	46.51
	Literals	7.28	39.87	38.67	59.08	76.54	41.36	42.01	63.87	75.86
Small Projects	All	100.00	53.71	55.10	56.59	66.70	51.00	56.92	54.05	64.19
	Punctuation	37.30	63.86	65.39	67.38	85.57	64.77	71.70	68.01	85.05
	Identifier	32.06	37.76	39.11	40.29	44.48	35.96	42.49	38.47	40.10
	Operator	17.54	69.83	70.86	70.46	73.48	63.05	66.79	64.41	76.12
	Keyword	10.55	48.92	50.12	50.00	55.62	32.01	37.19	35.09	43.19
	Literals	2.56	27.44	29.84	48.26	74.59	30.09	30.96	45.87	72.29

TABLE III
STATISTICS OF INTRA-SCENARIO COMPLETION DATASET (API-BENCH).

PL	Domain	#num of projects		#num of functions	
		Database	Test	Database	Test
Java	Android	292	85	53,646	29,211
	ML	38	14	27,962	10,435
	Security	42	16	14,202	4,698
	Test	41	14	16,439	5,341
	DL	191	116	24,546	12,778
Python	ML	165	159	23,990	20,533
	Security	95	47	10,823	4,527
	Test	57	22	4,761	3,303

ranging from 17.66% to 21.71%. For completing identifiers and literals, which are project-specific components, the base models are subpar. This supports our assumption that general code completion language models lack domain-specific knowledge. Using k NM-LM, the accuracy of completing identifiers exhibited improvements of 9.31% and 8.90% in large projects, and 6.72% and 4.14% in small projects, respectively. This notable enhancement can be attributed to k NM-LM’s use of project code as the retrieval database, allowing domain-specific knowledge, such as custom functions and objects within the project, to be incorporated into the LM.

►RQ1◀ Results indicate that utilizing historical project code with R-LMs improves intra-project code completion, with k NM-LM exhibiting the most substantial improvement across all projects compared with baselines.

VI. RQ2: INTRA-SCENARIO CODE COMPLETION

A. Study Design

Completing code across different scenarios necessitates specialized knowledge, primarily manifested in understanding and utilizing the specific third-party libraries. For instance, mobile development code uses the Android library, while deep learning code depends on the PyTorch library. However, large language models may lack insight into the usage of these third-party libraries, due to their proprietary or frequently updated nature. In this research, we leverage projects that use the same third-party libraries to build the code scenario database. These projects contain scenario-specific knowledge, including

distinct API usage patterns. Then, we explore the potential advantages of using this scenario database for intra-scenario code completion. Furthermore, we evaluate the effectiveness of completing lines of code that involve Android API calls.

B. Dataset

We utilize the API-Bench dataset [30] to facilitate code completion within specific scenarios. This dataset encompasses four distinct code scenarios, involving both Java and Python. In each of these scenarios, we employ the train set to construct a datastore and then evaluate token completion using the corresponding test set. For a comprehensive overview of the dataset’s statistics, refer to Table III.

In addition, we curate an extra test set designed to evaluate the effectiveness of completing lines containing Android API calls. To achieve this, we leverage JavaParser³ to extract the object type associated with the method call, filtering out API calls from the Android library. Subsequently, we utilize the code context preceding the API call to facilitate completion of the entire code line, encompassing the API name and input parameters. This dedicated test set comprises 49,344 lines featuring Android API usage. Notably, the database employed for this task is identical to the one used for token completion within the Android dataset.

C. Results

Token Completion Accuracy. Results in Table IV suggest that R-LMs do not consistently outperform the base LMs in intra-scenario completion task. For example, in Java, BM25 has a negative impact on Android, ML, and Security. Similarly, in Python, both BM25 and ReACC perform worse than the base LMs across all four scenarios. While the Hybrid and k NN-LM models show some improvement, their performance gains are smaller than those observed in the intra-project completion. This is mainly because code within the same project is more similar than those in the same scenario [8], making it harder for the retrieval-augmented approaches to find useful code. On the other hand, k NM-LM shows a significant performance improvement of approximately 6-12% compared to the base LMs. This improvement can be attributed to the

³<https://github.com/javaparser/javaparser>

TABLE IV

TOKEN COMPLETION ACCURACY ON INTRA-SCENARIO DATASET (API-BENCH). NUMBERS ARE SHOWN IN PERCENTAGE(%). k NM-LM W/O BAYESIAN USES THE MANUALLY SET λ LIKE k NN-LM [15], AND k NM-LM W/O NEW OBSERVATIONS USES $p(E')$ AS λ .

Models	Java				Python			
	Android	ML	Security	Test	DL	ML	Security	Test
CodeGPT	66.77	64.98	65.4	64.74	51.71	51.81	52.77	50.09
+BM25	66.75	64.28	63.92	64.08	49.96	49.94	50.27	47.14
+ReACC	-	-	-	-	49.92	49.98	50.45	47.71
+Hybrid	-	-	-	-	52.16	52.32	52.94	50.37
+kNN	69.88	67.68	68.14	67.12	53.65	53.76	56.83	53.98
+kNM	72.34	71.00	71.40	70.54	59.50	59.27	64.85	60.87
+kNM w/o Bayesian	69.38	67.45	67.86	66.90	53.46	53.58	56.59	53.66
+kNM w/o new observations	72.17	70.90	71.30	70.36	59.49	59.26	64.68	60.73
UniXcoder	67.75	67.65	66.80	65.65	59.53	59.77	63.01	60.30
+BM25	67.59	66.51	65.60	64.88	56.99	56.66	61.49	56.90
+ReACC	-	-	-	-	56.30	56.29	59.32	56.53
+Hybrid	-	-	-	-	56.10	56.41	59.76	57.20
+kNN	69.83	71.28	69.20	67.63	62.30	62.59	69.86	65.88
+kNM	74.86	76.08	74.84	73.12	67.03	67.02	75.56	70.09
+kNM w/o Bayesian	69.76	71.22	69.15	67.55	62.24	62.55	69.88	65.84
+kNM w/o new observations	74.83	76.10	74.70	73.07	67.07	67.07	75.56	70.09

decoupled database construction. Traditional R-LMs retain all data when constructing the database, leading to the code that use the specific API call not being retrieved due to the few samples problem. In contrast, k NM-LM constructs a decoupled database that only retains code that the LM cannot predict accurately. This strategy increases the ratio of those specific API calls that cannot be correctly predicted by LM in the datastore, thus improving performance.

Ablation Study of k NM-LM. Table IV also includes an ablation study on the use of Bayesian inference. Specifically, two variants of the k NM-LM are evaluated: k NM-LM w/o Bayesian, which relies on a manually set λ value, and k NM-LM w/o new observations, which sets the prior probability of $p(E')$ as λ . The results indicate that manually setting the combination coefficient has the least favorable effect, while employing Bayesian inference has the most favorable effect in most cases. This is primarily because determining the optimal λ value for different scenarios and LMs can be challenging, while Bayesian inference can take into account the impact of both the scenarios and LMs and dynamically adjust λ based on the code text being completed.

Performance of Completing Lines that involve Android APIs. Table V presents an experimental evaluation of completing lines containing Android API. The results show that the performance of the basic LM is not satisfactory, as indicated by the EM metrics of 5.99% and 2.74% for completion of lines that contain Android API calls. However, the EM metrics of k NM-LM are significantly improved, increasing by 7.75% and 8.94%, respectively. Moreover, we utilize the Venn diagram [31] to examine the precision of each model’s completion on a per-method basis. The analysis demonstrates that k NM-LM outperforms other models in terms of accuracy, with 4,245 more accurately predicted code lines than all other models combined.

TABLE V

PERFORMANCE OF COMPLETING LINES THAT CONTAIN ANDROID APIS.

	CodeGPT		UniXcoder	
	ES	EM	ES	EM
Base	44.38	5.99	42.87	2.74
+BM25	44.65	6.02	42.77	2.71
+kNN-LM	48.48	7.91	44.79	3.40
+kNM-LM	54.10	13.74	49.59	11.68

TABLE VI

COMPARISON OF RESOURCE USAGE IN k NN-LM AND k NM-LM.

Model	Large projects		Small Projects	
	Speed (tokens/s)	DB size (MB)	Speed (tokens/s)	DB size (MB)
CodeGPT	3063.40	-	1773.92	-
+kNN-LM	1242.25	780.48	990.71	156.78
+kNM-LM	1282.68	391.12	953.27	79.17
UniXcoder	2730.87	-	1625.89	-
+kNN-LM	1175.03	713.76	937.06	142.42
+kNM-LM	1252.48	428.96	948.49	86.85

►RQ2◀ Utilizing code from the same scenario can improve intra-scenario code completion, with k NM-LM exhibits maximum enhancement. Also, Bayesian inference is important in determining the optimal combination coefficient.

VII. DISCUSSION

In this chapter, our primary focus is on assessing the effectiveness of k NM-LM from three key angles. As k NM-LM is an improvement over k NN-LM, we begin by discussing the distinctions between these two methods. Following this, since the aim of k NM-LM is to enable domain-specific code completion without the need of fine-tuning, we proceed to compare it with fine-tuned models. Lastly, we conduct case studies to provide concrete examples that validate the efficacy of our proposed approach.

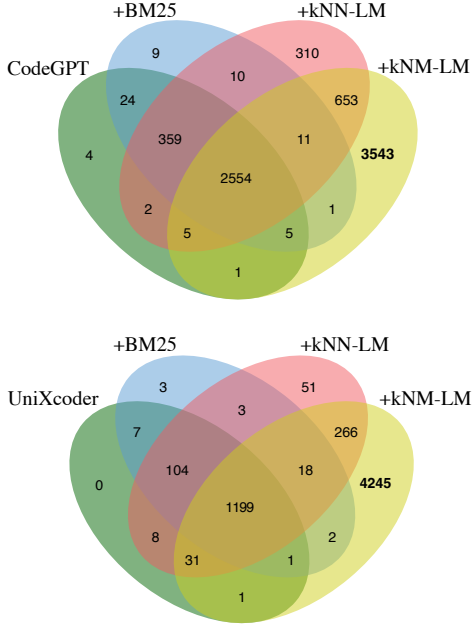


Fig. 4. Venn diagram of completing lines that contain Android APIs. It shows the number of samples that are completed correctly.

A. Compared with kNN -LM.

Speed and Space Usage Analysis. Table VI demonstrates that both kNN -LM and kNM -LM have a negative impact on completion speed. However, even with the worst-case scenario, the completion of 937.06 words can be achieved within 1 second. As developers typically anticipate receiving the completion results within 200 milliseconds [32], we think that the impact on completion speed is not significant. Storage space consumption, on the other hand, is a more significant consideration. Both models require database storage, but kNM -LM uses considerably less space than kNN -LM. In large projects, kNM -LM saves 284.83MB and 389.36MB compared to kNN -LM, while in small projects, it saves 55.57MB and 67.61MB. Given the significant performance gain from kNM -LM, the additional storage space required is acceptable.

Influence of Hyper-parameters. We compare the performance of kNM -LM and kNN -LM with varying hyper-parameter settings for intra-project code completion. Results, shown in Fig. 5, indicate that hyper-parameters have a significant impact on the performance of kNN -LM. In contrast, kNM -LM is less sensitive to hyper-parameters and achieves good performance for most settings. Therefore, kNM -LM is more suitable for domain adaptive code completion task, as it can provide satisfactory results without the carefully hyper-parameters configuration.

B. Compared with Fine-tuning.

We compared kNM -LM and fine-tuning on intra-scenario code completion task, as presented in Fig. 6. We found that fine-tuning is more effective than kNM -LM, but using kNM -LM after fine-tuning can still lead to some improvement. It is

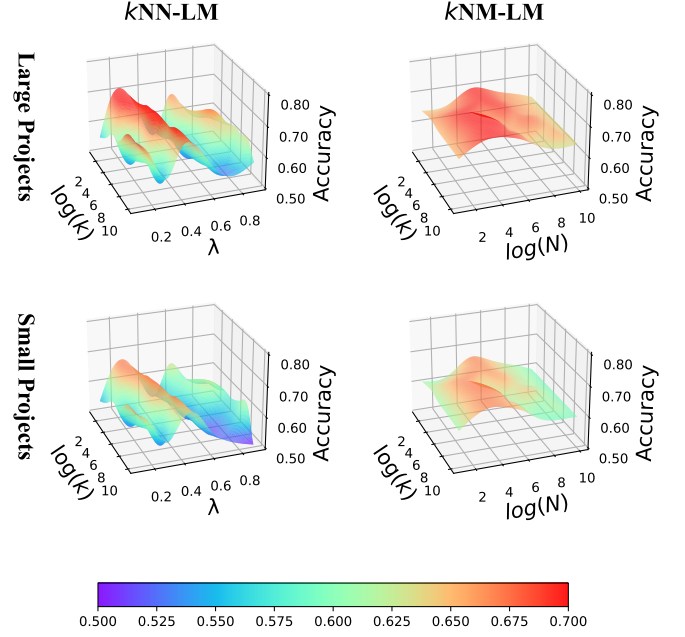


Fig. 5. Influence of hyper-parameters settings in kNN -LM and kNM -LM. The base LM is CodeGPT and the task is intra-project code completion. The number of searched neighbors k is used in both. kNN -LM needs to set the combination coefficient λ , and kNM -LM needs to set the window size N . The above two sub-figures are the results on large projects, and the below two are results on small projects.

worth noting that fine-tuning may not be a practical solution for domain adaptive code completion due to the continuously updated domain features and the black-box constraint of LM. In contrast, kNM -LM may be more suitable for individual users to obtain personalized code completion service.

C. Qualitative Analysis.

In this study, we conducted a qualitative analysis on completing lines that contain Android API calls, presented in Fig. 7. The study also included three examples that demonstrated the capabilities of the models used. For instance, Example ID=3877 showed that kNM -LM accurately completed complex API names, which a general model without Android-specific knowledge would fail to predict accurately. Example ID=40542 demonstrated that kNM -LM accurately predicted the API name and correctly filled in the corresponding input parameters. On the other hand, the BM25 and kNN -LM models failed to correct the API name and also incorrectly completed additional input parameters. Lastly, example ID=45427 verified that kNM -LM can ignore irrelevant code and accurately complete the intended code. In contrast, other models were affected by irrelevant code, including file writing semantics in the completed code.

VIII. THREATS TO VALIDITY

The first threat is related to the base LMs selected for the experiment. To mitigate this threat, we chose two state-of-the-art pretrained LMs, CodeGPT and UniXcoder, for the code

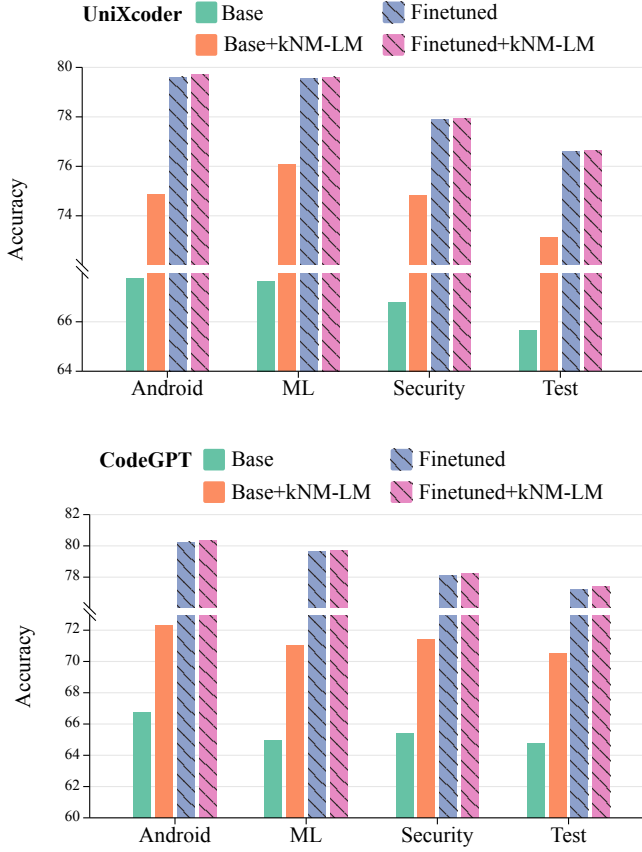


Fig. 6. Finetune VS k NM-LM.

completion task. We want to compare our proposed method with fine-tuning, hence we do not use black-box LM systems, such as Codex or ChatGPT. However, in section III-C, we introduced how to use our approach in GPT-3 through API calling. The second threat is related to the setting of the domain database and test set on intra-project code completion, as both use code from the same project. To mitigate this threat, we use the Git commit history to split the code for building the database and the code for testing, ensuring that the code in the test set were developed after those in the database. The third threat is related to the fact that we did not verify the effect of ReACC and Hybrid on the Java language, as the authors only released their pretrained models on Python. To mitigate this threat, we chose to compare our approach with BM25, which has similar performance compared to ReACC and Hybrid as reported in their paper. Lastly, in the context of intra-project code completion, our choice of the commit point at which 80% of the code is already written introduces a potential concern. This choice does not always reflect practical scenarios, we intend to investigate the impact of varying the size of the intra-project database in our future work. To ensure transparency and enable verification by other researchers, we open-sourced all our code, data, and experimental results.

ID: 3877

```
// the ability of completing the complex API name.
public void setUpTheme ( ) {
    Utils.setUpPopupTheme(context, settings);
    actionBar=getActionBar();
    actionBar.
```

UniXcoder: setUpTheme(context, settings);

+BM25: setUpTheme(context, settings);

+ k NN-LM:;

+ k NM-LM: setDisplayHomeAsUpEnabled(**true**);

Ground Truth: setDisplayHomeAsUpEnabled(**true**);

ID:40542

```
//the ability of filling arguments.
private void drawBackground (Canvas canvas, float
↳ fromX, float toX){
    int count = canvas.save();
    canvas.clipRect(
        fromX, (canvas.getHeight()-mStrokeWidth)/2,
        toX, (canvas.getHeight()+mStrokeWidth)/2);
    mBackgroundDrawable.
```

UniXcoder: drawBackground(canvas);

+BM25: drawBackground(canvas, fromX, toX);

+ k NN-LM: drawBackground(canvas, count);

+ k NM-LM: draw(canvas);

Ground Truth: draw(canvas);

ID: 45427

```
//the ability of ignoring long useless codes.
void createVideoThumbnail(String
↳ thumbnailSaveDir,
                                String originalPath,
                                ↳ int scale){
    Bitmap bitmap =
        ↳ ThumbnailUtils.createVideoThumbnail(
            originalPath,
            ↳ MediaStore.Video.Thumbnails.MINI_KIND);
    ...
    fileOutputStream = new
        ↳ FileOutputStream(targetFile);
    bitmap.
```

UniXcoder: encode(fileOutputStream);

+BM25: write(fileOutputStream);

+ k NN-LM: writePixels(0, 0, originalImageWidth);

+ k NM-LM: recycle();

Ground Truth: recycle();

Fig. 7. Qualitative examples.

IX. RELATED WORK

A. Code Completion

Code completion can be divided into three categories based on the code type to be completed: API name completion [33], [34], variable name completion [35], [36] and arbitrary token completion [29], [37]–[39]. API name completion aims to complete API calls from a specific third-party library by extracting API call sequence [40]. Variable name completion recommends variable references from declared variables, using data-flow graphs [36] or pointer network [35]. Arbitrary token

completion aims to complete arbitrary tokens in the code sequence or the abstract syntax tree (AST) of code, using Tree-based NN [41], [42], Graph NN [36], [43] or Transformer [44]. For example, CodeFill [29] improves GPT-2 [25] by predicting the code type when completing the next token, while Grammarformer [37] generates code completions with “holes” inserted in places where the model is uncertain. Our approach belongs to completing arbitrary token in code sequence.

Nowadays, many works begin to focus on code completion in specific domains, such as test code [45], [46] or repository level code completion [9], [47]. Our approach focuses on domain adaptive code completion, which is a generic approach that can adapt to different domains by switching the database. It does not require access to the parameters of the LM, hence can be used with black-box LLMs like GPT-4.

B. Retrieval-augment Language Model

Retrieval-augmented language models (R-LMs) utilize retrieval-based techniques to improve the performance of LMs and can be divided into two main categories: block R-LMs and token R-LMs. Block R-LMs [12], [13], [48] are similar to one-shot or few-shot learning [49], where one or a few examples are retrieved from a database instead of being randomly selected. Token R-LMs [15], [24], [50] retrieve tokens from database and then combine the retrieved results into the LM. Compared with block R-LMs, token R-LMs can update retrieval results at the same time of generating new tokens, hence our approach uses the architecture of token R-LM. However, token R-LMs suffer from high storage costs and require hyper-parameters selection to combine the inference results from the database and language model.

Various approaches have been proposed to address the limitations of token R-LMs. For example, k NN-Adapter [51] uses a trained network to determine the combination weights. To reduce the search cost, RetoMaton [52] uses the automaton states to save search time, while AdaptRet [53] uses a trained network to decide whether to use the retrieval module. GNN-LM [50] selects similar texts and builds a contextual graph to incorporate into the language model. In contrast, our proposed approach does not require training or the addition of an additional module. By decoupling the datastore and language model, we can save the storage cost and utilize Bayesian inference to select suitable hyper-parameters at the same time.

X. CONCLUSION

In this paper, we propose k NM-LM for domain adaptive code completion. By utilizing the in-domain code to construct the retrieval database, the language model can adapt to complete code in target domain without fine-tuning. Specifically, k NM-LM builds a decoupled domain database (saving only tokens that the language model cannot predict correctly) and employs Bayesian inference to combine the results from the language model and the database. Experiments have shown that k NM-LM achieves the best results on intra-project and intra-scenario code completion tasks. Notably, our approach does not require access to the weights in the language model,

nor does it require adding any additional neural network modules, making it feasible for leveraging the black-box language models on other domain-intensive code-related tasks, such as code summarization [43], [54], bug identification [55], localization [56] and repair [57]. To these ends, we make all our code, data, and models publicly available.

XI. ACKNOWLEDGMENTS

We are grateful to the anonymous reviewers for their useful comments and suggestions. This work was supported by National Key Research and Development Program of China (2022YFF0711404), Natural Science Foundation of Jiangsu Province, China (BK20201250), Cooperation Fund of Huawei-NJU Creative Laboratory for the Next Programming and CCF-Huawei Populus Grove Fund. Jidong Ge and Shangqing Liu are the corresponding authors.

REFERENCES

- [1] S. Lu, D. Guo, S. Ren, J. Huang, A. Svyatkovskiy, A. Blanco, C. Clement, D. Drain, D. Jiang, D. Tang *et al.*, “CodeXGLUE: A machine learning benchmark dataset for code understanding and generation,” *arXiv preprint arXiv:2102.04664*, 2021.
- [2] D. Guo, S. Lu, N. Duan, Y. Wang, M. Zhou, and J. Yin, “UniXcoder: Unified cross-modal pre-training for code representation,” in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, pp. 7212–7225.
- [3] W. Ahmad, S. Chakraborty, B. Ray, and K.-W. Chang, “Unified pre-training for program understanding and generation,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Online: Association for Computational Linguistics, Jun. 2021, pp. 2655–2668. [Online]. Available: <https://www.aclweb.org/anthology/2021.naacl-main.211>
- [4] C. Niu, C. Li, V. Ng, J. Ge, L. Huang, and B. Luo, “SPT-code: sequence-to-sequence pre-training for learning source code representations,” in *Proceedings of the 44th International Conference on Software Engineering ICSE*, 2022, pp. 2006–2018.
- [5] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 5485–5551, 2020.
- [6] “GitHub Copilot website.” [Online]. Available: <https://github.com/features/copilot>
- [7] “ChatGPT plugin in VSCode.” [Online]. Available: <https://marketplace.visualstudio.com/items?itemName=timkmecl.chatgpt>
- [8] F. F. Xu, J. He, G. Neubig, and V. J. Hellendoorn, “Capturing structural locality in non-parametric language models,” in *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. [Online]. Available: <https://openreview.net/forum?id=nnU3IUMJmN>
- [9] E. Bogomolov, S. Zhuravlev, E. Spirin, and T. Bryksin, “Assessing project-level fine-tuning of ML4SE models,” *CoRR*, vol. abs/2206.03333, 2022. [Online]. Available: <https://doi.org/10.48550/arXiv.2206.03333>
- [10] D. Zan, B. Chen, D. Yang, Z. Lin, M. Kim, B. Guan, Y. Wang, W. Chen, and J.-G. Lou, “CERT: Continual pre-training on sketches for library-oriented code generation,” in *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022*, 2022, pp. 2369–2375.
- [11] S. Gururangan, A. Marasovic, S. Swayamdipta, K. Lo, I. Beltagy, D. Downey, and N. A. Smith, “Don’t stop pretraining: Adapt language models to domains and tasks,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pp. 8342–8360.
- [12] K. Guu, K. Lee, Z. Tung, P. Pasupat, and M. Chang, “Retrieval augmented language model pre-training,” in *International conference on machine learning*. PMLR, 2020, pp. 3929–3938.

- [13] S. Lu, N. Duan, H. Han, D. Guo, S.-w. Hwang, and A. Svyatkovskiy, "ReACC: A retrieval-augmented code completion framework," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, pp. 6227–6240.
- [14] M. R. Parvez, W. Ahmad, S. Chakraborty, B. Ray, and K.-W. Chang, "Retrieval augmented code generation and summarization," in *Findings of the Association for Computational Linguistics: EMNLP 2021*, 2021, pp. 2719–2734.
- [15] U. Khandelwal, O. Levy, D. Jurafsky, L. Zettlemoyer, and M. Lewis, "Generalization through memorization: Nearest neighbor language models," in *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. [Online]. Available: <https://openreview.net/forum?id=HkIBjCEkvH>
- [16] T. G. Dietterich, "Ensemble methods in machine learning," in *Multiple Classifier Systems, First International Workshop, MCS 2000, Cagliari, Italy, June 21-23, 2000, Proceedings*, ser. Lecture Notes in Computer Science, J. Kittler and F. Roli, Eds., vol. 1857. Springer, 2000, pp. 1–15. [Online]. Available: https://doi.org/10.1007/3-540-45014-9_1
- [17] "Openai's text embedding api." [Online]. Available: <https://platform.openai.com/docs/guides/embeddings>
- [18] "Openai's completion api." [Online]. Available: <https://platform.openai.com/docs/api-reference/completions/create>
- [19] J. JOYCE, "Bayes' theorem," *Stanford Encyclopedia of Philosophy*, 2003.
- [20] H. Husain, H. Wu, T. Gazit, M. Allamanis, and M. Brockschmidt, "CodeSearchNet Challenge: Evaluating the state of semantic code search," *CoRR*, vol. abs/1909.09436, 2019. [Online]. Available: <http://arxiv.org/abs/1909.09436>
- [21] S. Robertson, H. Zaragoza *et al.*, "The probabilistic relevance framework: Bm25 and beyond," *Foundations and Trends® in Information Retrieval*, vol. 3, no. 4, pp. 333–389, 2009.
- [22] D. Guo, S. Ren, S. Lu, Z. Feng, D. Tang, L. Shujie, L. Zhou, N. Duan, A. Svyatkovskiy, S. Fu *et al.*, "GraphCodeBERT: Pre-training code representations with data flow," in *International Conference on Learning Representations ICLR, 2021*.
- [23] J. Johnson, M. Douze, and H. Jégou, "Billion-scale similarity search with gpus," *IEEE Transactions on Big Data*, vol. 7, no. 3, pp. 535–547, 2019.
- [24] U. Khandelwal, A. Fan, D. Jurafsky, L. Zettlemoyer, and M. Lewis, "Nearest neighbor machine translation," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. [Online]. Available: <https://openreview.net/forum?id=7wCBOFJ8hJM>
- [25] A. Svyatkovskiy, S. K. Deng, S. Fu, and N. Sundaresan, "IntelliCode Compose: Code generation using transformer," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2020, pp. 1433–1443.
- [26] U. Alon, S. Brody, O. Levy, and E. Yahav, "code2seq: Generating sequences from structured representations of code," in *International Conference on Learning Representations*.
- [27] E. Spirin, E. Bogomolov, V. Kovalenko, and T. Bryksin, "Psiminer: A tool for mining rich abstract syntax trees from code," in *2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR)*, 2021, pp. 13–17.
- [28] V. J. Hellendoorn, S. Proksch, H. C. Gall, and A. Bacchelli, "When code completion fails: A case study on real-world completions," in *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 2019, pp. 960–970.
- [29] M. Izadi, R. Gismondi, and G. Gousios, "Codefill: Multi-token code completion by jointly learning from structure and naming sequences," in *Proceedings of the 44th International Conference on Software Engineering*, 2022, pp. 401–412.
- [30] Y. Peng, S. Li, W. Gu, Y. Li, W. Wang, C. Gao, and M. Lyu, "Revisiting, benchmarking and exploring api recommendation: How far are we?" *IEEE Transactions on Software Engineering*, 2022.
- [31] P. Bardou, J. Mariette, F. Escudé, C. Djemiel, and C. Klopp, "jvenn: an interactive venn diagram viewer," *BMC bioinformatics*, vol. 15, no. 1, pp. 1–7, 2014.
- [32] C. Wang, J. Hu, C. Gao, Y. Jin, T. Xie, H. Huang, Z. Lei, and Y. Deng, "Practitioners' expectations on code completion," *arXiv preprint arXiv:2301.03846*, 2023.
- [33] X. He, L. Xu, X. Zhang, R. Hao, Y. Feng, and B. Xu, "Pyart: Python api recommendation in real-time," in *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 2021, pp. 1634–1645.
- [34] A. Svyatkovskiy, Y. Zhao, S. Fu, and N. Sundaresan, "Pythia: Ai-assisted code completion system," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 2727–2735.
- [35] J. Li, Y. Wang, M. R. Lyu, and I. King, "Code completion with neural attention and pointer networks," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 2018, pp. 4159–25.
- [36] M. Allamanis, M. Brockschmidt, and M. Khademi, "Learning to represent programs with graphs," in *International Conference on Learning Representations*.
- [37] D. Guo, A. Svyatkovskiy, J. Yin, N. Duan, M. Brockschmidt, and M. Allamanis, "Learning to complete code with sketches," in *International Conference on Learning Representations*, 2021.
- [38] S. A. Hayati, R. Olivier, P. Avvaru, P. Yin, A. Tomasic, and G. Neubig, "Retrieval-based neural code generation," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, Eds. Association for Computational Linguistics, 2018, pp. 925–930. [Online]. Available: <https://doi.org/10.18653/v1/d18-1111>
- [39] T. B. Hashimoto, K. Guu, Y. Oren, and P. Liang, "A retrieve-and-edit framework for predicting structured outputs," in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, S. Bengio, H. M. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., 2018, pp. 10073–10083. [Online]. Available: <https://proceedings.neurips.cc/paper/2018/hash/cd17d3ce3b64f227987cd92cd701cc58-Abstract.html>
- [40] H. Zhong, T. Xie, L. Zhang, J. Pei, and H. Mei, "Mapo: Mining and recommending api usage patterns," in *ECOOP 2009—Object-Oriented Programming: 23rd European Conference, Genoa, Italy, July 6-10, 2009. Proceedings 23*. Springer, 2009, pp. 318–343.
- [41] K. S. Tai, R. Socher, and C. D. Manning, "Improved semantic representations from tree-structured long short-term memory networks," in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2015, pp. 1556–1566.
- [42] W. Wang, K. Zhang, G. Li, S. Liu, Z. Jin, and Y. Liu, "A tree-structured transformer for program representation learning," *arXiv preprint arXiv:2208.08643*, 2022.
- [43] S. Liu, Y. Chen, X. Xie, J. K. Siow, and Y. Liu, "Retrieval-augmented generation for code summarization via hybrid gnn," in *International Conference on Learning Representations*.
- [44] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman *et al.*, "Evaluating large language models trained on code," *arXiv preprint arXiv:2107.03374*, 2021.
- [45] E. Dinella, G. Ryan, T. Mytkowicz, and S. K. Lahiri, "Toga: a neural method for test oracle generation," in *Proceedings of the 44th International Conference on Software Engineering*, 2022, pp. 2130–2141.
- [46] P. Nie, R. Banerjee, J. J. Li, R. J. Mooney, and M. Gligoric, "Learning deep semantics for test completion," 2023.
- [47] F. Zhang, B. Chen, Y. Zhang, J. Liu, D. Zan, Y. Mao, J.-G. Lou, and W. Chen, "Repocoder: Repository-level code completion through iterative retrieval and generation," *arXiv preprint arXiv:2303.12570*, 2023.
- [48] G. Izacard, P. Lewis, M. Lomeli, L. Hosseini, F. Petroni, T. Schick, J. Dwivedi-Yu, A. Joulin, S. Riedel, and E. Grave, "Few-shot learning with retrieval augmented language models," *arXiv preprint arXiv:2208.03299*, 2022.
- [49] Y. Wang, Q. Yao, J. T. Kwok, and L. M. Ni, "Generalizing from a few examples: A survey on few-shot learning," *ACM computing surveys (csur)*, vol. 53, no. 3, pp. 1–34, 2020.
- [50] Y. Meng, S. Zong, X. Li, X. Sun, T. Zhang, F. Wu, and J. Li, "Gnn-lm: Language modeling based on global contexts via gnn," in *ICLR 2022 Workshop on Deep Learning on Graphs for Natural Language Processing*.

- [51] Y. Huang, D. Liu, Z. Zhong, W. Shi, and Y. T. Lee, “ k nn-adapter: Efficient domain adaptation for black-box language models,” *arXiv preprint arXiv:2302.10879*, 2023.
- [52] U. Alon, F. Xu, J. He, S. Sengupta, D. Roth, and G. Neubig, “Neuro-symbolic language modeling with automaton-augmented retrieval,” in *International Conference on Machine Learning*. PMLR, 2022, pp. 468–485.
- [53] J. He, G. Neubig, and T. Berg-Kirkpatrick, “Efficient nearest neighbor language models,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 5703–5714.
- [54] X. Hu, G. Li, X. Xia, D. Lo, and Z. Jin, “Deep code comment generation,” in *Proceedings of the 26th conference on program comprehension*, 2018, pp. 200–210.
- [55] Y. Zhou, S. Liu, J. Siow, X. Du, and Y. Liu, “Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks,” *Advances in neural information processing systems*, vol. 32, 2019.
- [56] F. Niu, W. K. G. Assunção, L. Huang, C. Mayr-Dorn, J. Ge, B. Luo, and A. Egyed, “Rat: A refactoring-aware traceability model for bug localization,” in *2023 45th International Conference on Software Engineering (ICSE)*. IEEE, 2023, accepted.
- [57] W. Zhong, H. Ge, H. Ai, C. Li, K. Liu, J. Ge, and B. Luo, “Standup4npr: Standardizing setup for empirically comparing neural program repair systems,” in *37th IEEE/ACM International Conference on Automated Software Engineering*, 2022, pp. 1–13.