

# ARIST: An Effective API Argument Recommendation Approach

Son Nguyen, Cuong Tran Manh, Kien T. Tran, Tan M. Nguyen, Thu-Trang Nguyen, Kien-Tuan Ngo and Hieu Dinh Vo\*

Faculty of Information Technology, VNU University of Engineering and Technology, Hanoi, Vietnam

## ARTICLE INFO

### Keywords:

Effective argument recommendation, code completion, program analysis, statistical language model

## ABSTRACT

Learning and remembering to use APIs are difficult. Several techniques have been proposed to assist developers in using APIs. Most existing techniques focus on recommending the right API methods to call, but very few techniques focus on recommending API arguments. In this paper, we propose ARIST, a novel automated argument recommendation approach which suggests arguments by predicting developers' *expectations* when they define and use API methods. To implement this idea in the recommendation process, ARIST combines program analysis (PA), language models (LMs), and several features specialized for the recommendation task which consider the functionality of formal parameters and the positional information of code elements (e.g., variables or method calls) in the given context. In ARIST, the LMs and the recommending features are used to suggest the promising candidates identified by PA. Meanwhile, PA navigates the LMs and the features working on the set of the valid candidates which satisfy syntax, accessibility, and type-compatibility constraints defined by the programming language in use. Our evaluation on a large dataset of real-world projects shows that ARIST improves the state-of-the-art approach by 19% and 18% in top-1 precision and recall for recommending arguments of frequently-used libraries. For general argument recommendation task, i.e., recommending arguments for every method call, ARIST outperforms the baseline approaches by up to 125% top-1 accuracy. Moreover, for newly-encountered projects, ARIST achieves more than 60% top-3 accuracy when evaluating on a larger dataset. For working/maintaining projects, with a personalized LM to capture developers' coding practice, ARIST can productively rank the expected arguments at the top-1 position in 7/10 requests.

## 1. Introduction

Application programming interfaces (APIs) are extensively used in software development to reuse defined components in source code such as libraries and frameworks. Learning and remembering to use APIs are difficult. Thus, researchers have introduced the techniques to help developers in using APIs [1, 2, 3, 4, 5, 6, 7]. Most existing approaches focus on recommending the right methods to call and leave the arguments for developers to complete [1, 2, 3, 4]. According to the empirical studies in the existing work [6], more than half of method calls in practice require arguments (aka., actual parameters), and unfamiliarity with argument usage could cause bugs [8, 6, 9]. Moreover, effectively determining the correct arguments passed to method calls is a non-trivial task [5, 6, 7]. Indeed, the arguments usually contain one or more identifiers, while recommending identifiers remains one of the most challenging tasks in automated code completion [10, 11, 12].

Recognizing the importance of argument recommendation, researchers have introduced automated tools to suggest arguments to complete method calls. Asaduzzaman *et al.* [5] and Zhang *et al.* [6] follow an information retrieval (IR) direction to recommend arguments for the APIs that have been frequently used in the past. These approaches construct

a database of the seen arguments and their corresponding context. For a request from developers to recommend an argument, hereinafter *argument recommendation request*, the key idea is that two similar contexts should have similar arguments. The IR direction requires an appropriate design of predefined static features to represent contexts and methods determining their similarity. For example, the coding context could be represented by containing (enclosing) method name/declaration, calling methods, or the code tokens prior to the calling methods [5, 6]. More importantly, since the argument candidates are suggested by searching for the similar contexts, these IR approaches cannot recommend a new argument that they have not seen before.

Our empirical study performed on 1,000 top-ranked Java projects on Github (Section 2) shows that arguments and their usages (an argument with an API method call) are very unique. Specifically, 63% of the argument usages are unique, while the proportion of unique arguments is 47%. Hence, searching for the previously seen arguments to complete API method calls may not be effective. Moreover, the recommended arguments must conform to the syntactic, type-comparability, and accessibility constraints defined by programming languages. Thus, using the seen arguments without considering those constraints might cause incorrect recommendations.

In this work, we introduce ARIST, a novel approach which combines program analysis (PA), language models (LMs), and several argument-recommending features considering the functionality of formal parameters and positional information to predict developers' expected arguments. For an argument recommendation request of an API

\*Corresponding author

✉ sonnguyen@vnu.edu.vn (S. Nguyen); tranmanhcuong@vnu.edu.vn (C.T. Manh); 18020026@vnu.edu.vn (K.T. Tran); 18020050@vnu.edu.vn (T.M. Nguyen); trang.nguyen@vnu.edu.vn (T. Nguyen); tuangokien@vnu.edu.vn (K. Ngo); hieuvd@vnu.edu.vn (H.D. Vo)

ORCID(s): 0000-0002-8970-9870 (S. Nguyen); 0000-0002-3596-2352 (T. Nguyen); 0000-0001-7136-7529 (K. Ngo); 0000-0002-9407-1971 (H.D. Vo)

arXiv:2306.06620v1 [cs.SE] 11 Jun 2023

method call, our idea is that the argument candidates must conform to the syntactic and semantic constraints defined by the programming language in use. Also, the candidates should meet the expectations of the API developers and API users. We aim to benefit from the strengths of both PA and LMs as well as the recommending features in which the LMs and the features are applied to suggest the arguments that satisfy the stakeholders' expectations, while PA enforces the syntactic, accessibility, and type-compatibility constraints in the recommendation process.

Particularly, ARIST applies PA to analyze the given partially completed code (*context*) to identify the syntax-valid, accessible, and type-compatible argument candidates (so-called *valid candidates*) rather than retrieving candidates from a set of the seen arguments. The identified valid candidates are ranked according to their likelihood to be the next argument of the API method call. In ARIST, the likelihood of candidates is assigned by the LMs and the designed features. However, the identified set of valid candidates could be huge. The reason is that the accessible code elements (e.g., variables, fields, methods, or literals) in source code could be combined in many ways to form type-compatible candidates (i.e., expressions) for the request. This problem of a huge valid candidate set could cause the ineffectiveness and inefficiency of the ranking stage. To tackle this problem, before ranking the identified valid candidates, we filter out the less promising candidates by applying a set of heuristic rules and a light-ranking stage.

To evaluate ARIST, we conducted several experiments on two large datasets used in the existing studies with +1.6M argument recommendation requests in two very large projects (Eclipse and Netbeans) [5] and 9.2K projects [10], respectively. Our results show that ARIST outperforms the state-of-the-art approach [5] in both precision and recall up to 19% and 18% in recommending arguments for *frequently-used libraries* of Eclipse and Netbeans. For *general argument recommendation* task (i.e., recommending arguments for every method call) in these projects, the performance of ARIST is also significantly better than that of the state-of-the-art *n*-gram LM with nested cache (SLP) [13], GPT-2 [14], and CodeT5 [15] which are selected as the representative baseline approaches in general argument recommendation. Notably, GPT-2 is the core engine of Tabnine [16] and IntelliCode [17], two of the most popular AI-assisted code completer. Particularly, ARIST improves up to 125% in top-1 accuracy for CodeT5, GPT-2, and SLP.

Moreover, after training with the much larger dataset [10], ARIST achieves more than 60% top-3 accuracy in recommending arguments for *newly-encountered projects*. For *working/maintaining projects*, with a personalized LM to capture developers' coding practice, ARIST can productively rank the expected arguments at the top-1 position in up to +7/10 recommendation requests.

In brief, this paper makes the following contributions:

1. ARIST, a novel argument recommendation tool combining program analysis and language models in an effective and efficient recommendation process.
2. An extensive empirical evaluation showing that ARIST significantly outperforms the baseline approaches, and ARIST is effective in assisting developers in coding tasks.

## 2. Empirical Study

For automatically recommending arguments, one would face the following challenges. First, the recommended candidates must conform to the syntactic, type-compatibility, and accessibility constraints defined by the programming language in use. Second, the tool must correctly predict the argument which a developer intends to fill the method call at hand. To better understand the characteristics of arguments which might affect the automation of argument recommending, we conducted an empirical study to answer the following research questions:

**RQ1:** *What are the expression types [18] of arguments?*

For this RQ, we explore how common expression types (e.g., Simple Name, Method Invocation or String Literal) of arguments are in practice.

**RQ2:** *What are the expected data types of arguments?*

For this RQ, we investigate how common the data types (e.g., String or Object) of arguments are expected by the calling methods. In attempt to estimate the numbers of valid candidates for each expected data type of arguments, we design a simple method to generate syntax-valid, accessible, and type-valid candidates for the calling methods.

**RQ3:** *How unique are the arguments and their usages (i.e., an argument with a calling method)?* For this RQ, we investigate how often an argument of a calling method is likely to be not seen before (i.e., unique).

**RQ4:** *Where are the arguments from?* For this RQ, we investigate how often an argument is in the same projects as their calling methods.

**Definition 1.** An *argument usage* is a 3-tuple  $\langle arg, m\_call, pos \rangle$  where *arg* is an argument name, *m\_call* is a calling method name in which *arg* is the *pos*<sup>th</sup> argument.

**Definition 2.** An *Argument Recommendation (AR) request* *r* is a 3-tuple  $r = \langle P, m\_call, pos \rangle$  where *P* is a partially completed code (context), *m\_call* is a calling method which is partially completed in *P* where the *pos*<sup>th</sup> argument of *m\_call* is unfilled and requested for recommendations.

### 2.1. Data collection and processing

In this empirical study, we collected a dataset of 1,000 top-ranked, high-quality, long-history Java projects on GitHub. The projects are ranked based on the numbers of stars, folks, and commits of the repositories. The dataset could be found on our website [19]. In this dataset, all duplicated Java files, migrated projects, and the forks of the same projects are filtered out. This dataset includes about 460K files, +4.2M methods, and +42.4M LOC. For each method call having arguments, we collected and analyzed the name and origin of the calling method as well as the names, expression types, expected data types, and origin of arguments.

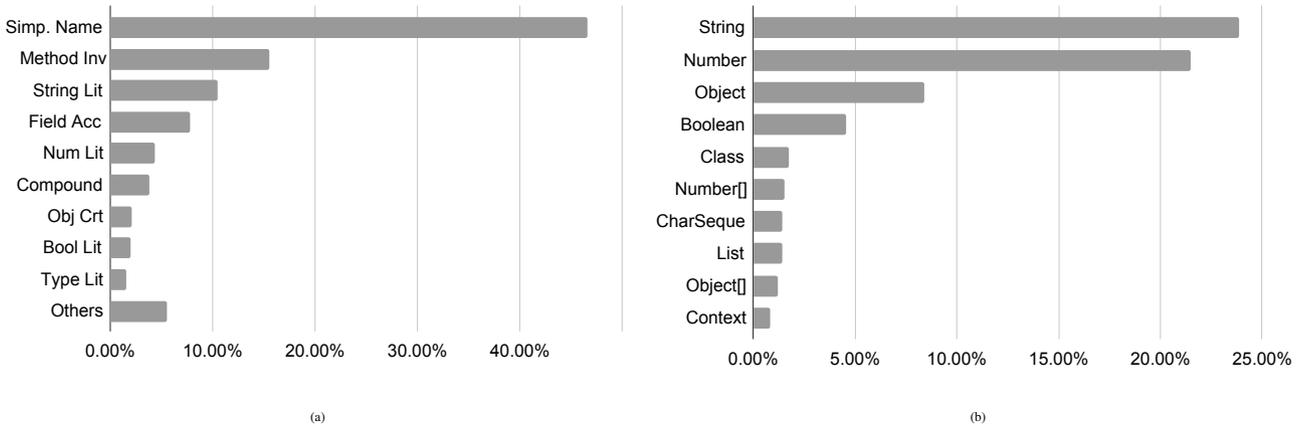


Figure 1: (a): Argument usage distribution by expr. types [18]; (b): Top-10 most popular argument data types

## 2.2. Results and Analysis

*RQ1: What are the expression types [18] of arguments?*

Figure 1a shows the portions of argument usages by expression types [20, 18, 6]. As seen, 70% of arguments are *Simple Name*, *Method Invocation*, and *Field Access*. In other words, *at least seven out of ten arguments are/contain identifiers*. Meanwhile, existing studies [10, 11] have confirmed that recommending identifiers remains a much harder task compared to recommending other token kinds such as keywords, separators, and operators of the syntactic grammar [21]. We also found that 60% of method calls have at least one argument in our dataset. Thus, *argument recommendation is a very frequent task, still, AR is very hard*.

*RQ2: What are the expected data types of arguments?*

As seen in Figure 1b, a significant portion of argument usages require an argument in the common data types [20, 18], such as *Object* or *String*. About 54% of the usages expect that the data type of arguments which is *Object*, *String*, or *Number* (e.g., *int* or *Double*). Consequently, when argument recommendations are requested at such positions, there are very large sets of (accessible) argument candidates that are type-compatible. Indeed, there are more than a half of argument usages (about 550K/1M usages) having at least 100 type-compatible (*valid*) candidates. Note that the valid candidates are limited by the set of expression types in Table 1. Without limiting expression types, the number of valid candidates could be infinite. For example, there are many argument usages in Eclipse project where the number of valid alternatives is up to 100K. Thus, given an AR request, *the set of argument candidates which are valid could be very large, and identifying the expected argument could be a challenging task*.

*RQ3: How unique are the arguments and their usages (i.e., an argument with a calling method)?* In our large corpus, argument usages are quite unique. Indeed, there are more than 6 out of 10 usages which appear **only once** in our dataset (Figure 2a). Even for sole arguments (the argument part in argument usages), they are still quite infrequent: 47% of the arguments are unique (Figure 2b). As a result, when being requested for recommendations, *a very large*

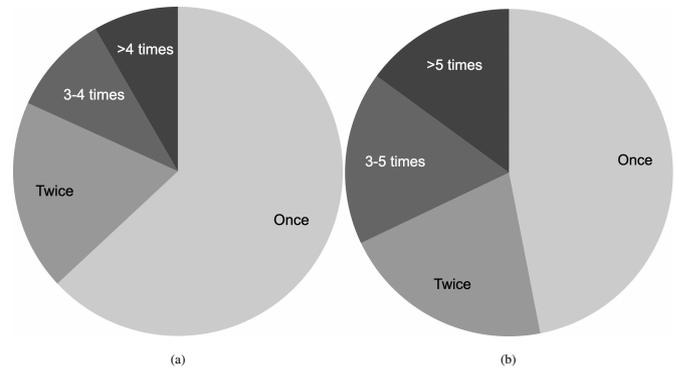


Figure 2: (a): Argument usage frequency; (b): Sole argument frequency

*number of the arguments might not be identified by searching for similar arguments/usages in the set of cases that are previously encountered.*

*RQ4: Where are the arguments from?* Regarding the origin of calling methods and arguments, nearly 50% of method calls are the invocations of the methods in the containing projects, while 84% of arguments can be only found in the containing projects. For these cases, those statistical approaches might fail to learn argument usages from a set of projects to apply to other projects. Moreover, the methods having the same method name can be declared with different sets of parameters. For example, there are several methods named *insert*, yet declared with different sets of parameters, such as `int insert(User u)`, `void insert(Record rec)`, or `void insert(int off, char[] str)`. In these cases, *one might not have a correct suggestion by blindly using the examples from the seen projects to recommend arguments in another project*.

## 3. Key Ideas and Approach Design

### 3.1. Key Ideas

To address the challenges discussed above, in completing argument for an API method, *m\_call*, our idea is

to recommend the argument candidates which satisfy the expectation of three stakeholders: the *creator* of the programming language in use, *API developers*, and *API users*.

For the *language creators*, the argument candidates are expected to be *valid* in terms of the syntax, type-compatibility, and accessibility. These expectations are encoded in the corresponding language compiler/interpreter as a set of rules to ensure that the written programs are valid/well-form. For the argument recommendation (AR) request in Figure 3, the argument candidates must be syntax-valid [18] and type-comparable expressions such as a variable (Simple Name) typed `Object` or a `String` literal. Additionally, the candidates must be accessible from the requesting location. In Figure 3, a variable candidate must be either fields of `AnnotationTypes` or static variables in `Nebeans`. Thus, in the argument recommendation process, the syntax, type-compatibility, and accessibility rules should be enforced to determine whether an argument candidate is *valid* or not.

For *API developers*, when defining *m\_call*, they might expect that the future arguments of *m\_call* must be *compatible* with the corresponding parameters in terms of their datatype and functionality. The type-compatibility expected by the language creators is *general* for all types developed in the language. Meanwhile, the type-compatibility expectation of the API developer is *specific* for the type of the corresponding parameter. Additionally, the functionality of parameters and arguments could be carried by their name. For example, the argument of `getProp` defined at lines 5–7 of Figure 3, when the method is called at line 10, is expected to be a **property**. The correct argument is `PROP_BACKGROUND_DRAWING` which indeed is a property. Liu *et al.* [7] show that arguments are very similar or even identical to their corresponding parameters. Inspired by this study, we estimate how the candidate fulfills the developers' expectation by measuring the similarity between the candidate's name and the corresponding parameter's name of *m\_call* (*parameter similarity*, Section 4.1).

The *API users* also expect that the candidates are valid regarding the programming language and reflect the users' intention. Indeed, the specific given context with each candidate must conform to the rules about the syntax, type-compatibility, and accessibility of the language in use. Additionally, to predict the intention of the developers when they use *m\_call*, we rely on the naturalness (repetition) of code because source code naturally written by developers does not occur randomly [22]. For example, the candidates of the AR request inside `isBackgroundDrawing` should intuitively be a property related to `BackgroundDrawing`. Thus, compared to `PROP_COMBINE_GLYPHS`, `PROP_BACKGROUND_DRAWING` is more likely to be the expected argument of the AR request in Figure 3, although both of the candidates are valid in terms of the syntax, type-compatibility, and accessibility. In fact, the naturalness of code can be captured by language models trained on the seen source code [22, 13].

Additionally, like natural language processing, besides the repetition of code, the positional information, which

```

1 package org.netbeans.editor;
2 public class AnnotationTypes {
3     public static String PROP_BACKGROUND_DRAWING = ...;
4     public static String PROP_COMBINE_GLYPHS = ...;
5     public Object getProp(Object prop){
6         //...
7     }
8     public Boolean isBackgroundDrawing() {
9         loadSettings();
10        bool b = (bool) getProp(/*AR_REQUEST*/);
11            //Expected: PROP_BACKGROUND_DRAWING
12    }
13 }

```

Figure 3: An example in Netbeans project

gives the model the relative position of the tokens in the context, also plays an important role in predicting the next tokens/argument [23]. Inspired by this idea, to capture the positional information of the candidates and apply it to predict the next arguments, we design two features specialized for the recommendation task, *creating-distance* and *accessing-recentness*, which are the “distances” to the positions where the candidate is declared and defined/used (Section 4.2).

### 3.2. Approach Design

To implement the above ideas in ARIST, we combine program analysis (PA), language models (LMs), and several features specialized for the argument recommendation task. For an AR request, the syntax, type-compatibility, and accessibility rules are enforced into the recommendation process to identify the *valid* argument candidates. The identified valid candidates are ranked according to their likelihood to be the next argument of the method call. In ARIST, the likelihood is assigned by the LMs and the selected features. In this novel combination, the LMs and selected features are applied to reduce the uncertainty in deciding the most likely candidates in the valid candidate set identified by PA. Meanwhile, PA can be used to navigate the LMs and selected features to apply on the right set of candidates. This can also boost the performance of the LMs in recommending the arguments where their usages are less frequent or have not been seen by the models. However, the LMs and selected features might still not effectively rank the expected arguments at the top positions because of the large identified candidate sets and the infrequency of arguments. Plus, these large sets could cause inefficiency in candidate ranking, especially for sophisticated LMs. To deal with the problem of large sets of valid candidates, ARIST reduces the identified candidate sets by filtering out the candidates which are less likely to be the expected arguments.

In ARIST, we design a set of filtering rules to eliminate the unpromising candidates. Additionally, a light-ranking stage is applied using the selected features and a light-weight LM, such as an *n*-gram LM, to efficiently isolate a relatively large number (e.g., 20–100) of promising candidates. Then, these promising candidates are passed to the heavy-ranking stage where a sophisticated model with more expensive

computation such as a deep neural network LM is performed to recommend the most appropriate arguments.

However, applying LMs in ARIST might face the problem of rare/unseen method calls/arguments [10] (out of vocabulary). To mitigate this problem, we tokenize source code into sequences of sub-tokens by camelCase and under\_score naming conventions (for the traditional count-based LMs such as  $n$ -gram) or employing byte pair encoding (BPE) [24] (for more sophisticated models such as GPT-2 or T5).

## 4. Features Specialized for API Argument Recommendation

In ARIST, we design a set of *static* features which capture the expectation of developers to estimate the likelihood to be the expected arguments of AR requests. In this version of ARIST, the features specialized for argument recommendation include the similarities between the candidates and the corresponding parameter of the calling method (*parameter similarity*) and how recently the candidates are created and accessed in the context (*creating-distance* and *accessing-recentness*). These features could always be applied to determine candidates' likelihood even when the candidates and calling methods are infrequent or unseen in the past. This could mitigate the negative impact of the out-of-vocabulary problem of LMs in ARIST and help it improve the recommendation performance.

### 4.1. Parameter Similarity

In a method declaration, a parameter's name usually conveys information about the semantics of the parameter in the method body as well as the meaning of the expected argument when the method is invoked. Liu *et al.* [7] have shown that parameter names and argument names are very similar. The similarity between an argument ( $c$ ) and its corresponding parameter ( $p$ ) is calculated as follows:

$$parasim(c, p) = \frac{|com\_terms(c, p)| + |com\_terms(p, c)|}{|terms(c)| + |terms(p)|} \quad (1)$$

In formula (1),  $terms(s)$  is the decomposition of  $s$ , based on camelCase and under\_score convention.  $com\_terms(n_1, n_2)$  returns the longest subsequence of  $terms(n_1)$ , where each term in the subsequence appears in  $terms(n_2)$ .

We have performed an empirical study to investigate the *parameter similarity* in Eclipse and Netbeans. Our results show that 68.24% of the arguments are very similar to their corresponding parameters' names, while 22.33% of the arguments are completely different (Figure 4). The very low similarity is often caused by short and generic names. This confirmed the characteristic of parameter similarity reported by Liu *et al.* [7]. Inspired by Liu *et al.* [7], we also applied parameter similarity to improve the precision of the AR process in ARIST. Note that, to avoid vanishing other features of the candidates which are different from the corresponding parameter, the parameter similarities of

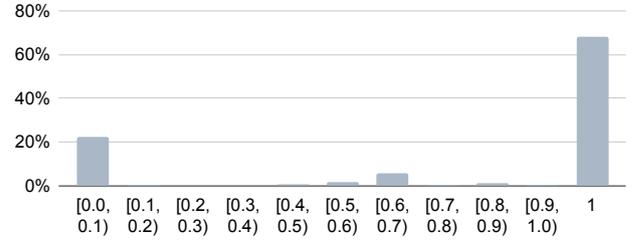


Figure 4: Distribution of parameter similarity

the candidates are normalized to the range of  $[x, 1.0]$  where  $x > 0$  to exclude the value of 0.

### 4.2. Creating-distance and Accessing-recentness

From a method call's position, the local variables/parameters, instance variables, or class variables (so-called *variables*) used as the arguments of this method call are often *created* or *accessed* very recently. In Figure 6, variable provider is created in the same scope as where it is used as an argument (line 3). For the first argument of ignoreProvider, parameter user is also recently *used* at line 2. In ARIST, for a request  $r = \langle P, m\_call, pos \rangle$ , we utilize the features which reflect how recent the variable is to the first time (*creating-distance*) and the last time (*accessing-recentness*) it appearing in the context to improve the AR performance. In practice, these features of a variable might require recording the real-time changes by developers, and hence be hard to collect. In ARIST, the *creating-distance* and *accessing-recentness* are statically approximated as follows.

**Creating-distance.** An accessible variable  $c$  could be created/declared outside the containing method, even  $c$  could be a field of parent classes. Thus, the *creating-distance* of  $c$  is computed as the *scope-distance* between the code closest block/region containing the calling method, and the code block/region<sup>1</sup> where  $c$  is created.

**Definition 3. (creating-distance).** Given an AR request  $r = \langle P, m\_call, pos \rangle$  and a candidate  $c$  which is an accessible variable, the creating-distance,  $create\_dis(c, r)$ , is the scope-distance of  $B$  and  $B'$ ,  $scope\_dis(B, B')$ , where  $B$  is the closest block where  $c$  is created, and  $B'$  is the closest block containing  $m\_call$ :

- $scope\_dis(B, B') = 0$  if  $B = B'$ .
- $scope\_dis(B, B') = 1 + scope\_dis(B, B'')$ , where  $B''$  is the closest block containing  $B'$ .

For AR request  $r$  in Figure 6,  $B'$  contains only line 4, while variable provider is created in block  $B$  from lines 2–5, and  $create\_dis(provider, r) = 1$ . In general,  $scope\_dis(B_1, B_2) = null$  if  $B_1$  does not (in)directly contain  $B_2$ . In Definition 3, since  $c$  is accessible for the given context,  $B'$  must be contained by  $B$ , and  $create\_dis(c, r) \neq$

<sup>1</sup>A block is a sequence of statements, local class declarations, and local variable declaration statements within braces [25].

*null*. For the candidates which are accessible fields,  $B$  is considered as the block containing all the blocks in the containing method. The closest block  $B$  where a global variable is created can be considered as the outermost block.

**Accessing-recentness.** The access variable  $c$  could be used in the code blocks/regions which have no ordering or containing relation, e.g different methods. Thus, *accessing-recentness* should be estimated within methods and not by using scope-distance. In ARIST, the *accessing-recentness* of candidate  $c$  is simply calculated as the distance in lines of code (LOC) as follows:

**Definition 4. (accessing-recentness).** For an AR request  $r = \langle P, m\_call, pos \rangle$  and a candidate  $c$ , the *accessing-recentness*  $access\_rec(r, c)$  is the distance in LOC between the line (indexed  $L$ ) where  $m\_call$  is called and the line (indexed  $L'$ ) where  $c$  is most recently used/defined in the containing method,  $access\_rec(c, r) = |L - L'|$ .

For AR request  $r$  in Figure 6,  $access\_rec(user, r) = 0$  and  $access\_rec(provider, r) = 1$ . Note that if  $c$  is never defined and used in the containing method, then  $access\_rec(c, r) = null$ . For the candidates which are fields or global variables, the methods, which define/use them, could be invoked in the containing method. An inter-procedural analysis can be performed to more accurately measure  $access\_rec$ . However, those methods can call other methods. Thus, a (recursive) inter-procedural analysis with a certain depth is required. For an efficient estimation of  $access\_rec$ , we set the depth as 0 in this version of ARIST.

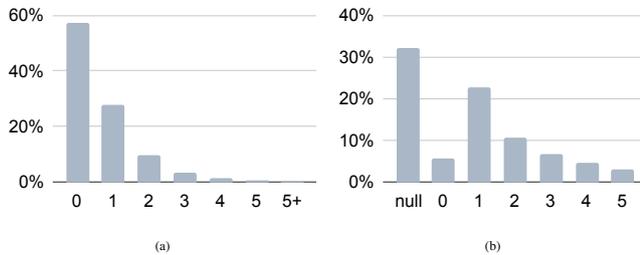


Figure 5: (a) *creating-distance*; (b) *accessing-recentness*

In this version of ARIST, the *creating-distance* and *accessing-recentness* of candidates are computed based on heuristics. In fact, they can be computed more accurately by applying data-flow analyses. However, as the program to be completed is often not compilable, partial analysis techniques such as PPA [26] could produce unreliable results, especially for complex code. Additionally, an accurate analysis might be costly. Meanwhile, the *creating-distance* and *accessing-recentness* of candidates, which are estimated by the current simple method, can still capture the recentness of the candidates which are variables. Indeed, Figure 5 shows the distribution of the *creating-distance/accessing-recentness* of the arguments which are a variable (*variable arguments*) in our large corpus. As seen, the variable arguments tend to be “closer” to the requesting position. For *creating-distance*, nearly 3/5 variable arguments are

```

1 public int insert(User user){
2     String provider = user.getProvider();
3     if(StringUtils.equals("`standard`", provider)){
4         ignoreProvider(user, /* AR_REQUEST */);
5     }
6 }

```

Figure 6: An example in *xdiamond* project

declared in the same scope as their corresponding method calls (*creating-distance* is 0). Meanwhile, the variable arguments are often either used for the first time in the containing method (*accessing-recentness* is *null*), or reused very quickly (within a few LOCs). However, reusing variables at the same line is less likely. From these results, we conclude that for variable arguments, developers usually use the variables recently declared or defined/used. Thus, *creating-distance* and *accessing-recentness* should be used as recommending features in ARIST’s AR process.

## 5. ARIST: An Effective API Argument Recommendation Approach

Figure 7 illustrates the argument recommendation (AR) process of ARIST. For an AR request, ARIST identifies the set of valid argument candidates which are syntax-valid, type-compatible, and accessible (*valid*) for the given context. Next, the set of valid candidates is reduced by filtering out unpromising candidates. This could improve the candidate ranking performance of the final step in not only accuracy but also recommending time.

### 5.1. Identifying Valid Candidates

Given an AR request  $r = \langle P, m\_call, pos \rangle$ , ARIST analyzes the partially completed code  $P$  to generate a set of code expressions based on the syntax, type-compatibility, and accessibility constraints defined by the programming language in use such as Java. The main idea is that from accessible expressions, ARIST constructs more complex accessible expressions whose data types are compatible with the expected data type(s)  $t$ . The generation starts with the accessible variables/methods, constants, static methods, and classes in the given partially completed code  $P$ .

First, the expected data type  $t$  of the  $pos^{th}$  parameter in the method declaration corresponding to  $m\_call$  is identified. In ARIST, to avoid misleadingly omitting valid candidates, all the overloading method declarations are also considered. For example, both `Pos getPosition(Character)` and `Pos getPosition(String)`, which are the members of the same class in the context, are considered for the request  $r = \langle P, getPosition, 1 \rangle$ . In this case,  $t$  could be either `Character` or `String`.

Next, the initial set of accessible expressions is constructed by identifying all the accessible variables, methods, constants, static methods, and classes in the given context  $P$ . These expressions are used to generate the more complex expressions based on the grammar and accessibility

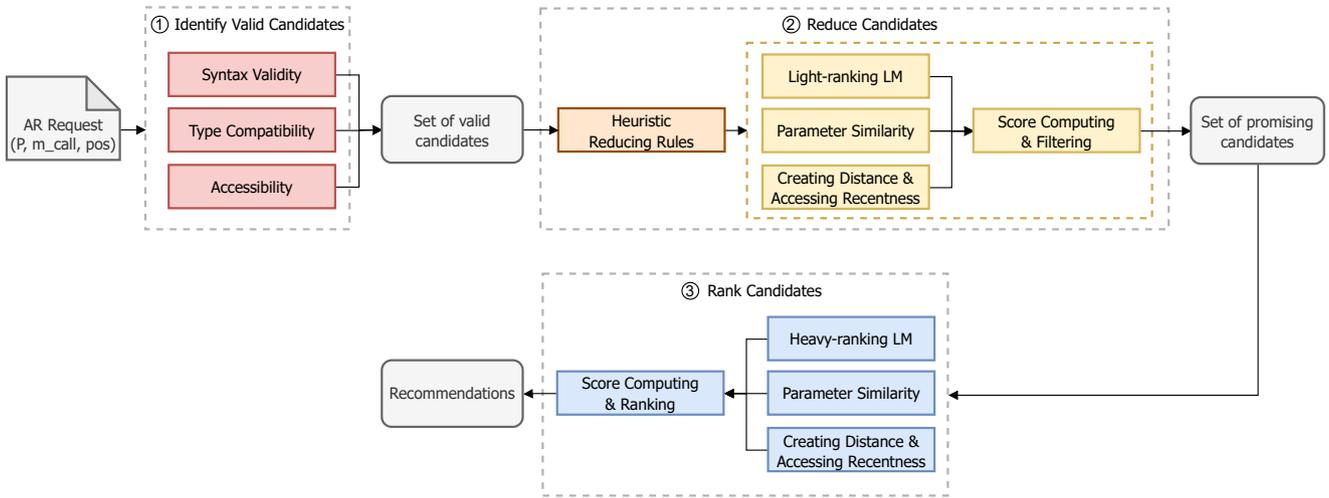


Figure 7: ARIST's API Argument Recommendation Process Overview

constraints. Essentially, valid expressions could be infinitely complex. However, most of the arguments belong to a small number of expression types in practice [6, 18]. As shown in Table 1, 99.43% of the arguments in our dataset in Section 2 are categorized into 17 types of expressions. Among these types, 15/17 types are supported by ARIST, we leave out two very complex expression types *lambda expression* and *compound expression*. These complex expressions are too costly to be generated and account for only about 4% of the arguments. As a result, 95% of all the arguments in our dataset are covered by ARIST. Finally, the generated expressions are checked if their data type is compatible with at least one expected data type of  $r$ . Particularly, the data type  $t'$  of a valid candidate must be *compatible* with  $t$  without casting. In other words,  $t = t'$  or  $t'$  inherits from  $t$ .

Without loss of generality, the *arguments* of *method invocation* and *object creation*, the *dim* of *array creation*, and the *index* of *array access* can be treated as new arguments. Thus, these elements are left for the next AR requests of users in ARIST to avoid infinite argument generation. For the example shown in Figure 3, the request expects an `Object`, and after apply our candidate identification method, there are about 40K valid candidates. To improve AR performance, it's essential to narrow down such very large sets of valid candidates. In the next section, we introduce our method to reduce the identified candidate sets.

## 5.2. Reducing Candidates

For a request  $r = \langle P, m\_call, pos \rangle$ , the set of valid candidates could be very large, especially when  $r$  expects an argument of a common data type such as `Object` or `String`. Meanwhile, as only one candidate among them must be recommended efficiently to complete the method call, the very large candidate set could cause incorrect recommendations. Moreover, it could become too expensive for heavy-ranking models (e.g., DNN models) to evaluate all the valid argument candidates. The reason is that these LMs could be generative models which are not designed to rank a large set

of candidates. In ARIST, we develop a set of heuristic rules and apply a light-ranking step to efficiently reduce the valid candidate set by filtering out unpromising candidates before passing them to the heavy-ranking stage.

### 5.2.1. Candidate Reducing by Heuristic Rules

In large projects such as Eclipse or Netbeans, their sets of static methods or static variables (public static fields) are usually very large. For example, Eclipse has about +74K static methods/variables, while this figure for Netbeans is +52K. Since static methods/variables and static methods can be used everywhere (i.e., always accessible), these large sets of static methods/variables could cause huge sets of valid candidates (e.g., Figure 3). Thus, we design the heuristic rules in ARIST to reduce the subset of candidates which are generated from static methods/variables. The valid candidate  $c$  generated from static methods/variables **must satisfy at least one** (i.e., *OR* relation) of the reducing rules to be kept in the reduced candidate set.

**Rule of shared sub-tokens.** Variable  $c$  shares at least a subtoken (case-insensitive) with the names of either the calling method, the name or data type of the calling object, or the containing method. For a method  $m$  called by object  $o$  in the containing method  $M$ , the arguments of  $m$  usually contain certain subtokens in the names of  $m$ ,  $o$ , and  $M$ , and the data type of  $o$ . In Figure 3, both the expected argument's name, `PROP_BACKGROUND_DRAWING` and the called method's name `getProp` also contains sub-token `prop`. Both the containing method's name, `isBackgroundDrawing`, and the argument's name also contain `background`. In our dataset, we found that 31.3% of the arguments share at least a subtoken with the names of either the called method, the name or data type of calling object, or the containing method. For Figure 3, this rule effectively and correctly reduces the candidate set to only 3.5K candidates from more than 52K valid candidates.

**Rule of recently-used classes.** Variable  $c$  is a member of a class which is recently used to access its members

**Table 1**

Argument expression types and their grammar. This grammar just shows the types' structures, and the accurate grammar can be found in Java Specification [18].

Expr. type	Grammar	Distr. (%)
<simple name>	<identifier>	49.31
<method invoc>	[<expr>"."]<identifier>"("<arguments>")"	13.01
<field access>	<expr>".<identifier>	11.07
<string lit>	-	4.91
<number lit>	-	3.27
<bool lit>	-	3.69
<null lit>	-	1.99
<char lit>	-	0.88
<type lit>	-	0.71
<lit>	<string lit> <bool lit> <null lit>  <number lit> <char lit> <type lit>	-
<cast expr>	"("<type>")"<expr>	0.95
<obj creation>	"new "<class type> "["<arguments>"]"	2.77
<array creation>	"new "<type> "["<dim>"]"	0.71
<compound expr>	<unary expr> <binary expr>	3.68
<this expr>	-	0.93
<lambda expr>	((<lambda para list>) <lambda body>	0.75
<array access>	<expr> "["<index>"]"	0.71
<static method ref>	<class type>::<identifier>	0.09
<expr>	<simple name> <method invoc> <cast expr>   <field access> <lit> <array access>	-
Total		99.43

(fields/methods) in the coding context. The rule ensures that ARIST never misses the static methods/variables or static methods in the class it has been encountered.

Besides these rules, we also apply other rules, such as:

- $c$  is a boundary variable (e.g., max/min).
- $c$  is a member of the class/parent classes of the context.
- $c$  is declared in the same package scope of the context.

We have preliminarily evaluated those rules on +35K AR requests in 10% randomly selected files in Netbeans. The results show that 91.5% of (unpromising) candidates have been eliminated without sacrificing much accuracy, i.e., there are only about 2% of the requests where the expected arguments are incorrectly eliminated.

### 5.2.2. Candidate Reducing by Light-ranking Model

Although the reducing rules filter out a large portion of valid yet unpromising candidates, the number of remaining candidates might still be large. For AR, a time-sensitive task, this could create difficulties in applying powerful (yet computationally expensive) models to quickly evaluate and rank these candidates. To further narrow down the search space, we select top- $RT$  candidates ranked based on their *parameter similarity*, *creating-distance*, *accessing-recentness* and the evaluation of a light-ranking language model. Specifically, a light-ranking model should be computationally efficient language model. In this work, we use a count-based method such as  $n$ -gram LM [27] or nested-cached LM [13] for this light-ranking stage.

$RT$  is the predefined *reducing threshold* and should be able to balance between candidate reduction performance and accuracy sacrifice (incorrectly eliminating expected arguments). Specially, given an AR request  $r = \langle P, m\_call, pos \rangle$ , the ranking score of a candidate  $c$  is computed as:

$$score_{lr}(c, r) = \sqrt[1+v]{\mathcal{P}_{lr}(c, P) \times parasim(c, p) \times recent(c, r)^v} \quad (2)$$

where  $p$  is the parameter's name, and  $\mathcal{P}_{lr}(c, P)$  is the likelihood that  $c$  is the next sequence following  $P$ , which is evaluated by a lightweight LM trained on the lexical form of a corpus. In formula (2),  $parasim(c, p)$  is the *parameter similarity* of  $c$ . Since the features of recentness are only applied for variables, if  $c$  is not a variable then  $v = 0$ , and only  $\mathcal{P}_{lr}$  and  $parasim$  are considered in  $score_{lr}$ . Otherwise,  $v = 1$  and  $recent(c, r) = Prob(D = create\_dis(c, r)) \times Prob(U = access\_rec(c, r))$ , such that  $Prob(D = d)$  and  $Prob(U = u)$  are the probability in a dataset that *creating-distance* equals to  $d$  and *accessing-recentness* equals to  $u$ . Setting  $v = 1$  when  $c$  is variable and  $v = 0$  otherwise also makes  $score_{lr}$  comparable in these two cases. Candidate  $c$  is ranked higher if  $c$  is more similar to the recommending parameter's name, "closer" to the requesting position, and more likely to be the next sequence following  $P$ .

In summary, the set of the valid candidates is reduced by the heuristic reducing rules. The reduced candidate set is further narrowed down efficiently up to a reducing threshold,  $RT$ , which remains the most promising candidates. To

isolate top- $RT$  promising ones, the candidates are evaluated based on their similarity with the recommending parameter, recentness, and likelihood to be the next sequence following the coding context.

### 5.3. Ranking Candidates

To produce recommendations, ARIST ranks the reduced candidate set based on their probability to be the argument of the calling method. In ARIST, the probability of a candidate is evaluated by a heavy-ranking model combined with the parameter similarity and recentness of the candidate. Compared to the count-based models (e.g.,  $n$ -gram LM) in the light-ranking stage, the heavy-ranking models applied in this stage should be capable of leveraging longer dependencies. The suitable options for heavy-ranking models could be neural network LMs such as LSTM, GPT-2, or CodeT5.

**Selective Prediction.** In fact, an LM could be more effective for arguments in certain expression types. To verify this assumption, we have performed an experiment to preliminarily evaluate the performance of the trained  $n$ -gram LM with nested-cached [13] and GPT-2 [14] in recommending arguments of method calls of Lucene project<sup>2</sup>. The results show that GPT-2 is better than  $n$ -gram LM with nested-cached in the cases whose the argument expression types are *Simple Name*, *Array Access*, *Type Literal*, *Object Creation*, *Array Creation*, *Lambda Expr*, and *Method Reference*. For the cases with the remaining expression types,  $n$ -gram with nested-cached performed better. Thus, given an AR request  $r = \langle P, m\_call, pos \rangle$ , we adopt ensemble methods [28] in ARIST to compute the likelihood of a candidate  $c$  to be the next sequence following  $P$ :

$$\mathcal{P}(c, P) = \begin{cases} \mathcal{P}_{hr}(c, P) & expr\_type(c) \in E \\ \mathcal{P}_{lr}(c, P) & \text{Otherwise} \end{cases} \quad (3)$$

where  $\mathcal{P}_{lr}(c, P)$  and  $\mathcal{P}_{hr}(c, P)$  are evaluated by a lightweight LM and a heavy-weight LM, respectively. In formula (3),  $E$  is the set of expression types where the heavy-ranking model performs better than the light-ranking model.  $E$  is empirically computed beforehand by evaluating the performance of the models, or  $E$  can be gradually updated while being applied. As only the candidates of certain expression types are evaluated by the heavy-ranking model, this method reduces the workload for the heavy-ranking model in ARIST while potentially improving the overall performance.

**Combining score.** Given a request  $r = \langle P, m\_call, pos \rangle$ , the final ranking score of a candidate  $c$  to be the argument for  $r$  is computed as:

$$score(c, r) = \sqrt[1+v]{\mathcal{P}(c, P) \times parasim(c, p) \times recent(c, r)^v} \quad (4)$$

where  $p$  is the parameter's name. Similar to  $score_{lr}$ , in formula (4), if  $c$  is a variable,  $v = 1$  and  $recent(c, r) = Prob(D = create\_dis(c, r)) \times Prob(U = access\_rec(c, r))$ , such that  $Prob(D = d)$  and  $Prob(U = u)$  are the probability

<sup>2</sup><https://github.com/apache/lucene>

**Table 2**  
Statistics of the datasets

	Small corpus	Large corpus
#Projects	Eclipse & Netbeans	9,271
#Files	53,787	961,493
#LOCs	7,218,637	84,236,829
#AR requests	700,696	913,175

in a dataset that *creating-distance* equals to  $d$  and *accessing-recentness* equals to  $u$ . If  $c$  is not a variable,  $v = 0$ . Finally, the candidates are ranked according to their  $score(c, r)$  score.

## 6. Empirical Methodology

We seek to answer the following research questions:

- **RQ5: Accuracy Comparison.** How accurate is ARIST in argument recommendation (AR)? How is it compared with the state-of-the-art AR approaches [5, 13, 16]?
- **RQ6: Sensitivity Analysis.** How do various factors affect ARIST, e.g. argument recommendation scenarios, training data's sizes, context lengths?
- **RQ7: Intrinsic Analysis.** How do ARIST's components, e.g., candidate identification, candidate reduction, heavy-ranking model, and reducing threshold, contribute to its performance?
- **RQ8: Time Complexity.** What is our running time?

### 6.1. Datasets

In this study, we used two datasets for evaluation, *small corpus* and *large corpus* (Table 2). For comparison, we used the same dataset (*small corpus*) as in the state-of-the-art technique for AR [5]. The *small corpus* contains two large projects: Eclipse and Netbeans. This corpus contains a relatively large code base, about 54K files and 7.2M LOCs. For further evaluation, we used the *large corpus* by Allamanis *et al.* [29] containing +9.2K popular Java projects from Github. The large corpus is widely used to evaluate the existing code completion approaches [10, 11, 13]. For a non-bias evaluation, these datasets in this evaluation are different from the dataset in our empirical study (Section 2).

### 6.2. Evaluation Setup, Procedure, and Metrics

#### RQ5. Accuracy Comparison

**Baselines.** We compare ARIST with the state-of-the-art argument recommendation (AR) techniques: (1) **PARC** [5], an IR-based AR approach and specialized for AR, represents argument usages by several static features, constructs a usage database, and searches for similar usages when requested; (2) **GPT-2**, the core engine of Tabnine [16] and IntelliCode [17], two of the most popular AI-assisted code completer; (3) **CodeT5** [15], an encoder-decoder model specialized for source code; (4) **SLP** [13], a representative code completer is carefully adapted from  $n$ -gram models

for source code and shown that it can even better than RNN and LSTM based models [13]. For each baseline, we applied a post-processing step to exclude the candidates which are type-incompatible in the recommendation sets resulted by the code model. Note that there are several others powerful models designed for general SE tasks such as CodeBERT [30] or CuBERT [31] which could be applied for the AR task. However, for a practical evaluation, we select CodeT5 [15], which is one of the most recent and powerful models specialized for code, as a baseline.

**Procedure.** Since the source of PARC is not publicly available, for a fair comparison with PARC, we use the original results of PARC [5], and apply the same dataset (the *small corpus*) and evaluation setting for the other approaches. Specially, we randomly shuffle and divide the source files of each project into 10 equal folds. We perform 10-fold cross-validation: each fold was chosen for testing, and the remaining folds were used for training. Since PARC focuses on AR for *frequently-used API libraries* [5], it fits well with recommending certain API methods in the common libraries such as SWT, AWT, or Java Swing which are frequently used in developing applications.

Meanwhile, GPT-2, CodeT5, and SLP can be applied to recommend arguments for a wider range of APIs/method calls. However, GPT-2, CodeT5, and SLP are designed to predict next tokens for the given context. In this work, we adapt GPT-2, CodeT5, and SLP for AR and evaluate the AR performance of ARIST and the two baselines in the *general AR task*, i.e., recommending arguments for any method call. To adapt these approaches for AR task, we apply the same mechanism used by Karampatsis *et al.* [10] to combine multiple predicted tokens to form a whole argument. Particularly, while predicting next tokens, the beam search algorithm is applied to find the  $k$  highest probability sequences of tokens, which each of them corresponds to a recommended argument.

**Metrics.** To compare with PARC [5] in recommending arguments for method in *frequently-used libraries*, we adopt the same metrics used in PARC, *Precision* and *Recall* at top  $k$ . Specially,  $Precision(k) = \frac{R(k)}{S}$  and  $Recall(k) = \frac{R(k)}{A}$ , where  $A$  is the total number of AR requests in the test set,  $S$  is the number of requests where the expression type of the expected argument is supported by the tools, and  $R(k)$  is the number of requests where the recommended argument in top- $k$  of the ranked list matches the expected argument.

Note that an expected argument could be a literal or a lambda expression. AR tools might not be able to correctly recommend the whole expressions [32]. Thus, if the approaches can correctly suggest the expression type, we consider it as a correct recommendation. Instead of suggesting incorrect ones, recommending default values (e.g., "`<EMPTY_STRING>`" for string literals and or  $\emptyset$  for number literals) or a template of the expected expression types would be more appropriate for these cases. This is still useful for developers and can be done by a post-processing step.

For *general AR task* in the other experiments, we apply top- $k$  accuracy (aka., *Recall* at top- $k$ ), and *Mean Reciprocal*

*Rank (MRR)*, which are widely used in code completion [10, 13]. Top- $k$  accuracy represents whether the target arguments are presented in top- $k$  ranked lists or not. *MRR* is the average of reciprocal ranks for the test set,  $MRR = \frac{1}{A} \sum_{i=1}^A \frac{1}{rank_i}$ , where  $A$  is the total number of AR requests in test set, and  $rank_i$  is the position of expected argument in the ranked list of request  $i^{th}$  in the test set. *MRR* focuses more on the rank of target argument in each suggestion list. For example, an *MRR* value of 0.5 shows that AR technique could suggest correct arguments at position 2 of ranked lists on average.

**Evaluation Setup.** For ARIST,  $n$ -gram LM with nested-cache is used for the light-ranking model, the reducing threshold  $RT = 20$ , GPT-2 is used for heavy-ranking model, the sets of expression types selective prediction and recentness probabilities are derived by analyzing the dataset used in Section 2 which can be found in our website [19].  $n$ -gram LM uses the Jelinek-Mercer smoothing with a fixed confidence of 0.5 and  $n = 6$ . For detailed implementation, see our website [19].

For *SLP*, we use their default settings in training and predicting for the evaluation<sup>3</sup>. For *GPT-2*, we use the pre-trained model with 345M parameters by OpenAI<sup>4</sup>. For *CodeT5*, we use the pre-trained *CodeT5-base* model with 220M parameters released by Salesforce<sup>5</sup>. In our experiments, all the pre-trained models are fine-tuned on the training data.

In this paper, all our experiments were run on a server with Intel(R) Xeon(R) CPU @ 2.30GHz, 32GB RAM, and Tesla P100 GPU.

## RQ6. Sensitivity Analysis.

In ARIST, the heavy-ranking model performs complex computations and should be deployed on a powerful server as a central/public resource for all users. Thus, this model should *arguably* not learn from a user's local source files and use the learned knowledge for others. Meanwhile, the light-ranking model which is very light-weight can be deployed in developers' (much less powerful) computer. This enables ARIST to *personalize* the argument recommendations by learning developers' coding practice from their local source files. As a result, ARIST can better support developers in several scenarios. In this work, we aim to investigate the impact of 3 different real-world AR scenarios on ARIST's performance on the large corpus:

- *New-project.* This setting reflects the scenario when ARIST encounters completely new projects in practice. The tool learns from a fixed training set and then is tested on the separate testing set for evaluating performance. This usage is cross-project: 400 projects (about 5% of the projects) for testing and the remaining for training.
- *Working-project.* This scenario comes from the idea that when developers need argument recommendations from the tool for the projects that they are working on. In

<sup>3</sup><https://github.com/SLP-team/SLP-Core>

<sup>4</sup><https://openai.com/blog/tags/gpt-2>

<sup>5</sup><https://blog.salesforceairesearch.com/codet5/>

this scenario, some source files available at that time are used for *personalizing* the argument recommendation task. Specially, instead of fixing the training set like *new-project* setting, ARIST could iteratively learn after each recommendation during testing. In this setting, the light-ranking model, which could be deployed on developers' machine, can learn the their coding practice from the available source files.

- *Maintaining-project*. This setting simulates the scenario that developers make AR requests when they maintain their projects. In this scenario, developers usually make small changes (e.g., several files/methods) to the existing code. The tool is tested on one file at a time in the test set. ARIST could learn from all the remaining files in the corpus.

Besides, we also study the impacts of the following factors on the performance of ARIST: *training data size*, *context length*, and *origin of called methods*.

### RQ7. Intrinsic Analysis.

We study the impacts of the following components: valid candidate identification, candidate reduction, candidate ranking, reducing size, and static features. We create different variants of ARIST with different combinations and measured their performance. Particularly, to evaluate the impact of a component, we create a variant where the component is disabled and the others are on.

## 7. Empirical Results

### 7.1. Accuracy Comparison (RQ5)

Table 3 shows the comparison results of ARIST against the state-of-the-art techniques in argument recommendation (AR), PARC [5]. For the results in Table 3, we applied the same evaluation procedure used in PARC [5]. Particularly, the techniques are tested on recommending arguments for method calls of *frequently-used libraries*: SWT for Eclipse and Java Swing and AWT for Netbeans. As seen, ARIST *outperforms PARC in both Precision and Recall for both Eclipse and Netbeans*. For Netbeans, ARIST achieves 16% and 15% better in the *Precision* and *Recall* than PARC at the top position. With Eclipse, both the *Precision* and *Recall* at the top position of ARIST are also much better than PARC, 19% and 18%. With higher *Recall* and *Precision*, ARIST is able to recommend more diverse argument expression types as well as recommends arguments more precisely.

Analyzing the results, we found the main reason for the higher accuracy of ARIST compared to PARC. Particularly, arguments are quite unique in practice (Section 2), and ARIST analyzes the code context to identify accessible and type-compatible argument candidates and learns to recommend arguments. Meanwhile, PARC retrieves the candidates stored in a database. Thus, ARIST can give the correct recommendations for the requests, even when their expected arguments have not been seen, while PARC cannot.

For *general AR task* (i.e., recommending argument for every method call), Table 4 shows the performances

**Table 3**

Performance of ARIST and PARC in AR task for the methods in the *frequently-used libraries*

Project		ARIST		PARC	
		<i>Precision</i>	<i>Recall</i>	<i>Precision</i>	<i>Recall</i>
Netbeans	Top-1	52.92%	51.67%	46.46%	44.86%
	Top-3	70.18%	68.28%	66.20%	66.75%
	Top-10	78.36%	76.15%	72.06%	69.57%
Eclipse	Top-1	56.66%	55.04%	47.65%	46.65%
	Top-3	67.88%	65.63%	65.05%	63.68%
	Top-10	73.14%	70.76%	72.26%	70.73%

```

1 // Package: org.netbeans.editor.ext
2 public class FormatSupport {
3     public TokenItem insertToken(TokenItem beforeToken,
4         TokenID tokenID, TokenContextPath
5         tokenContextPath, String tokenImage) {...}
6 }
7 // Package: org.netbeans.editor
8 public interface ImageTokenID extends TokenID {
9     public String getImage();
10 }
11 // Package: org.netbeans.editor.ext
12 public class ExtFormatSupport extends FormatSupport {
13     public TokenItem insertImageToken(TokenItem
14         beforeToken, ImageTokenID tokenID,
15         TokenContextPath tokenContextPath) {
16         return super.insertToken(beforeToken, tokenID,
17             tokenContextPath, /* AR_REQUEST */);
18         //Expected a String: tokenID.getImage()
19     }
20 }

```

**Figure 8:** An argument recommendation request in Netbeans

of ARIST, GPT-2, CodeT5 [15], and SLP [13]. As seen, ARIST's *performance is significantly better than those of GPT-2, CodeT5, and SLP*. Particularly, compared to GPT-2, the top-1 accuracy of ARIST is 24% and 14% better for Netbeans and Eclipse. The top-3 accuracy of ARIST is 16.4% and 13.5% better for those of CodeT5. For SLP, the improvements are much more significant, 87% and 125% for Netbeans and Eclipse. Especially, the top-5 accuracy of ARIST is about 80% on average. This means for 4/5 requests, the expected arguments can be found at top-5 positions of the recommended lists of ARIST. Plus, *MRR* of ARIST is about 0.71. In other words, ARIST correctly recommended the expected arguments from the first to the second positions on average. This is about 12%, 24%, and 73% relatively better than GPT-2, CodeT5, and SLP.

Analyzing the cases where ARIST performed better than GPT-2, CodeT5, and SLP, we found that even having high accuracy in predicting next tokens, these techniques failed to construct the expected arguments which contain multiple tokens. This causes their lower accuracy compared to ARIST's. For example, for the AR request at line 12 in Figure 8, GPT-2, CodeT5, and SLP can produce the recommendations which are partially correct. The closest recommendation of CodeT5 and GPT-2 is `tokenContextPath`. However, this

```

1 public final class DebuggerEngine implements
  ContextProvider {
2   public class Destructor {
3     //...
4   }
5 }
6 public final class DebuggerManager implements
  ContextProvider {
7   public DebuggerEngine[] startDebugging (DebuggerInfo
  info) {
8     //...
9     ep.setDestructor (/* AR_REQUEST */);
10    //Expected: engine.new Destructor ()
11  }
12 }
13 }

```

**Figure 9:** An example which argument recommendation tools often fail to suggest correct argument

recommendation is type-incompatible (TokenContextPath vs String). Thus, this candidate is filtered out from the candidate list, and the remaining ones such as text are even less relevant. SLP suggested variable tokenImage which is unavailable in the context, i.e., the variable is not declared in insertImageToken, and neither ExtFormatSupport nor FormatSupport has any field named tokenImage. Meanwhile, ARIST generated the valid candidate set including tokenID.getImage() in the context before feeding them to the LMs. Because of the likelihood and the *parasim* of tokenID.getImage(), the candidate is ranked first by ARIST.

By analyzing the cases where ARIST did not work well (i.e., cannot rank the expected arguments at the top-3 positions), we found that one of the main reasons is, that the expected expression types are unpopular and not supported by ARIST. Consequently, ARIST cannot construct and recommend the correct arguments. Figure 9 shows an example in which ARIST failed to suggest the correct argument. In this case, setDestructor requires an object type Destructor (line 9). This is a challenging case, because to fill a correct argument, a new object type Destructor should be created and passed to the method. Meanwhile, Destructor is a non-static inner class of DebuggerEngine, and to initialize a Destructor object, an object of the outer class (engine) must be used. This argument expression type is very unpopular and not supported by ARIST. Consequently, ARIST fails to generate a correct candidate for this challenging case. Also, for this complicated case, GPT-2, CodeT5, and SLP cannot suggest the expected argument.

## 7.2. Sensitivity Results (RQ6)

### 7.2.1. Impact of argument recommendation scenarios

Table 5 shows the performance of ARIST in different AR scenarios with the *large corpus*. For newly-encountered projects (*new-project* scenario), ARIST achieved an impressive performance. By ARIST, developers can find their expected arguments at top-3 positions in 6/10 requests. Meanwhile, with GPT-2 which is the engine of Tabnine [16], the top-10 accuracy is only 56.1%. For *working* projects (*working-project* scenario), with *personalizing only the*

**Table 4**  
Comparison in *general AR task*

Project		ARIST	GPT-2	CodeT5	SLP
Netbeans	Top-1	65.15%	52.63%	59.97%	34.91%
	Top-3	78.16%	57.69%	67.16%	48.10%
	Top-5	81.10%	57.87%	67.57%	55.02%
	Top-10	83.53%	57.88%	67.60%	67.20%
	<i>MRR</i>	0.72	0.55	0.63	0.44
Eclipse	Top-1	64.19%	56.53%	61.20%	28.52%
	Top-3	76.29%	61.89%	67.21%	41.60%
	Top-5	79.23%	62.09%	67.53%	49.46%
	Top-10	81.65%	62.10%	67.54%	62.67%
	<i>MRR</i>	0.70	0.59	0.64	0.38

**Table 5**  
Top-*k* accuracy of ARIST in different scenarios

	New project	Working project	Maintain. project
Top-1	53.42%	69.96%	74.49%
Top-3	61.50%	81.14%	83.23%
Top-5	64.21%	83.74%	85.38%
Top-10	67.96%	85.88%	87.38%
<i>MRR</i>	0.58	0.76	0.79

*light-ranking model* which can be deployed on developers' machines, ARIST achieved nearly 70% top-1 accuracy and 0.76 *MRR*. Especially, for *maintaining* projects (*maintaining-project* scenario), ARIST can correctly recommend arguments at the top position in +7/10 requests. ARIST with personalization can adjust the parameters in its light-ranking model for the working projects to capture developers' coding practice and reduce the cases where the expected arguments are misleadingly eliminated during candidate reduction. With this high accuracy, ARIST *can effectively assist developers in completing method calls while preserving their privacy*.

**Unseen expected arguments.** We also studied the performance of ARIST in recommending the expected arguments which have not been seen in the training corpus (*unseen expected arguments*) for the *new-project* scenario. We found that among the AR requests whose expected arguments is unseen, about 70% of them are correctly suggested by ARIST at the top-1 positions. This confirms our strategy of combining PA and LMs in the argument recommendation task. Specifically, PA is first used to identify correct arguments included in the valid candidate sets even when they are unseen, and then the LMs and the features specialized for AR are applied to learn recommending arguments rather than retrieving what has been stored in the training corpus.

**Accuracy by argument expression types.** Table 6 shows the performance of ARIST by argument expression types. As seen, ARIST is very effective for arguments which are *Simple Name* (e.g., variables or parameters), literals, or *This Expr* expressions. Especially, for *Simple Name* which

**Table 6**  
ARIST’s performance by the expected expression types

Expression type	Distribution (%)	Top-1 (%)
Simple Name	48.14	83.66
Method Invocation	15.19	45.51
Field Access	6.09	31.01
Array Access	0.74	53.26
Cast Expr	0.99	18.46
String Literal	10.03	98.14
Number Literal	5.06	95.66
Character Literal	0.47	87.93
Type Literal	0.90	81.92
Bool Literal	1.50	78.43
Null Literal	0.79	84.45
Object Creation	2.09	51.96
Array Creation	0.29	43.14
This Expr	1.06	91.05
Super Expr	0.00	0.00
Compound Expr	5.65	3.69
Lambda Expr	0.73	78.83
Method Reference	0.28	0.56
Total	100.00	69.96

accounts for a haft of the expected arguments, ARIST can correctly rank the expected arguments at the top position in more than 8/10 requests. The reason is that the major portion of expected *Simple Name* arguments are local variables or parameters of the containing methods. These arguments are correctly recommended because features *parameter similarity*, *creating-distance*, and *accessing-recentness* are encoded in ARIST. Meanwhile, ARIST does not work very well for *Method Invocation* and *Field Access*, only 45.5% and 31.0%, respectively. However, this is still much better than the performance of GPT-2, with about 29% and 28% top-1 accuracy for the same expression types. We found that among these AR requests, GPT-2 suggested invalid candidates in many cases. Also, GPT-2 prefers the candidates which are close to the requesting positions, even when they are inaccessible. Meanwhile, the expected arguments which are valid field access/method calls might not be very close to the positions or even defined in different files.

**Accuracy by argument variable types.** We also investigated the performance of ARIST in recommending arguments which are variables (parameters, local variables, class/instance variables). This kind of arguments accounts for 54.23% of all the arguments. We found that for local variables and parameters<sup>6</sup>, ARIST’s top-1 accuracy is 88.54%. Meanwhile, for others such as class variables and instance variables, the corresponding figure is 64.28%. Although local variables and parameters usually are method-specific and less frequent, they are frequently used as arguments. This property of argument usage is captured by ARIST thanks to features *parameter similarity*, *creating-distance*, and *accessing-recentness*.

<sup>6</sup>The arguments which are local variables and parameters are parts of *Simple Name* in Table 6

**Table 7**  
Impact of Context Length on Performance

	$l_1$	$l_2$	$l_3$	$l_4$	$l_5$
Top-1 (%)	62.05	65.86	66.14	67.00	67.83
<i>MRR</i>	0.70	0.72	0.72	0.73	0.74
Run. time (s)	0.33	0.39	0.42	0.51	0.56

### 7.2.2. Impact of testing context length

Because the LMs in ARIST recommend based on the given code sequences (the containing classes), the length of the code sequence could impact ARIST’s accuracy. To measure that impact, we randomly picked a fold of Netbeans for testing and the remaining folds for training in this experiment. For each AR request  $r = \langle P, m\_call, pos \rangle$  whose containing class in context  $P$  has  $n$  code tokens,  $C = \langle t_1, t_2, \dots, t_n \rangle$ , we evaluated the performance of ARIST in five different context lengths. To do this, we varied the context-beginning point in five locations,  $l_1 - l_5$ : the  $i^{th}$  point  $l_i$  contains  $i \times \lfloor n/5 \rfloor$  last tokens of  $C$ . Formally, the code sequence fed to LMs is  $C_i = \langle t_k, t_{k+1}, \dots, t_n \rangle$  where  $k = n - (i \times \lfloor n/5 \rfloor)$ . In Table 7, both top-1 accuracy and *MRR* increase by 9.3% and 5.7%, respectively, when we move the point from  $l_1$  to  $l_5$ , while running time increases 1.5 times. This is expected because ARIST has more information to decide the likelihood of candidates then performs better. However, ARIST needs more time to process longer contexts.

### 7.2.3. Impact of origin of calling method

For *intra-project* method calls (49% of the AR requests) where calling methods are declared in the requesting projects, calling methods might not be frequently used and are less likely to be seen by the models in the training set. Meanwhile, the *inter-project* method calls (51% of the AR requests) are usually the methods in common libraries such as JDK, Apache libraries, and Google Guava, and appear frequently in other projects. In this work, we also studied the impact of calling method’s origin on ARIST’s performance. Particularly, ARIST performed quite stably for *intra-project* and *inter-project* method calls, 72.9% and 63.8% top-1 accuracy, respectively. For *intra-project* calls, the impact of less frequent calls is mitigated thanks to valid candidate identification and static features. Meanwhile, as the methods, which are designed to be used cross-project, they often accept quite general argument data types such as String or Object. The identified sets of valid candidates for *inter-project* calls could be very large and contain many unexpected candidates (much *noise*). On the other hand, the sets of valid candidates for *intra-project* calls are much smaller because they require project-specific data types. This causes the higher accuracy for *intra-project* calls compared to ARIST’s accuracy for *inter-project* calls.

### 7.2.4. Impact of training data size

To study the impact of training data size on ARIST’s performance, we randomly picked a fold of Netbean for testing and varied the training data size by gradually adding

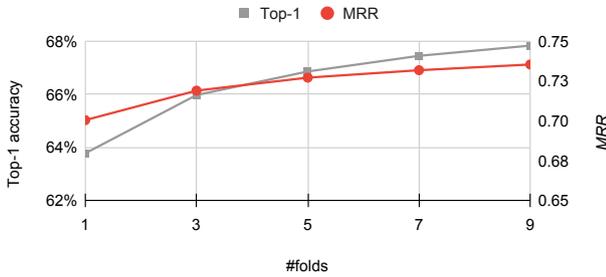


Figure 10: Impact of training data size

Table 8  
Impact of Valid Candidate Identification

	Top-1 (%)	MRR	Run. time (s)
ON	69.96	0.76	0.444
OFF	47.50	0.51	0.809

more data by fold to the training data. As expected, the accuracy slightly improves when we increase the training data size (Figure 10) because the models have observed more and thus performed better. In particular, the top-1 accuracy and *MRR* increase by more than 3.4% and 2.6%, when data increases from 1 fold to 3 folds. However, the increasing trends slow down when data size increases from 3 folds to 9 folds. The reason is that the added data does not contain much more new knowledge as in the smaller sizes.

### 7.3. Intrinsic Evaluation Results (RQ7)

#### 7.3.1. Impact of Valid Candidate Identification

To study the contribution of the valid candidate identification to the performance, we built a variant of ARIST where the candidate identification is disabled and evaluated it on the large corpus and *working-project* scenario. In this variant, we applied a beam search algorithm to construct the whole arguments from tokens recommended by the light-ranking model. As seen in Table 8, the valid candidate identification component significantly contributes to both the recommendation accuracy and running time of ARIST. Particularly, without the candidate identification, the performance decreases by 30% for top-1 accuracy and 50% for *MRR*, while the recommendation process is two times slower than the full version of ARIST. As expected, we found that in most of the cases where the full version worked better, the *n*-gram LM with nested cache constructed many invalid candidates. This confirms the effectiveness of our strategy to identify valid candidates before ranking.

Furthermore, we measured the performance of the candidate identification in ARIST. There are 91.0% of the requests where the expected arguments are correctly identified. The average number of identified candidates for a request (before candidate reduction) is about 7.2K.

Table 9  
Impact of Candidate Reduction

	Top-1 (%)	MRR	Run. time (s)
ON	69.96	0.76	0.444
OFF	61.98	0.69	2.424

#### 7.3.2. Impact of Candidate Reduction

We constructed a variant of ARIST without the candidate reduction component to measure the component’s contribution to the performance of ARIST with our large corpus and *working-project* scenario (Table 9). With much larger candidate sets (+7K candidates on average), ARIST’s performance is greatly reduced in both accuracy and efficiency when disabling the candidate reduction component. Specially, both top-1 accuracy and *MRR* significantly decrease, by about 11% and 8%, respectively. Moreover, the variant is also much slower, 6 times (2.4 sec) slower than that when enabling this component. Thus, candidate reduction should be enabled to maintain both the performance of ARIST.

#### 7.3.3. Impact of Light-ranking Model

To evaluate the impact of selecting the light-ranking model, we built two variants of ARIST with different count-based LMs applied for its light-ranking component: (plain) *n*-gram LM [27] and *n*-gram LM with nested-cache [13]. In both variants,  $n = 6$  and Jelinek-Mercer smoothing are applied, the reduction threshold is 20, and GPT-2 is applied in the heavy-ranking stage. For the large corpus, our results show that ARIST with plain *n*-gram LM achieved top-1 accuracy of only about 53%, while the corresponding figure of ARIST when applying *n*-gram LM with nested-cache was about 70% with a marginal running time increase. These results are expected and consistent with the findings of Helleendoorn *et al.* [13] where the *n*-gram LM carefully adapted for source code can yield performance that surpasses the plain *n*-gram LM.

#### 7.3.4. Impact of Reduction Threshold

Additionally, we evaluated the reduction threshold, *RT*, on the accuracy and running time of ARIST. We varied *RT* from 10 to 50 to study its impact on ARIST’s performance (Table 10). As seen, when  $RT = 10$ , ARIST ran faster because of small reduced sets of valid candidates. However, with a small reduction threshold, in many requests, the expected arguments were incorrectly eliminated before feeding the reduced candidate set to the heavy-ranking stage. The accuracy of ARIST is slightly improved when increasing *RT* from 10 to 50. This is because with a larger threshold, there are fewer requests where the identified candidate sets are over-reduced. Meanwhile, the running time is significantly increased because the heavy-ranking stage accounts for about 86.3% of the total running time (Section 7.4). Thus, in this version, we kept  $RT = 20$  to balance between the accuracy and recommending time.

**Table 10**  
Impact of reducing threshold,  $RT$

$RT$	10	20	30	40	50
Top-1 (%)	63.77	64.67	65.10	65.34	65.49
Run. time (s)	0.342	0.406	0.418	0.464	0.508

**Table 11**  
Impact of heavy-ranking stage

$P_{hr}$	Top-1 (%)	$MRR$	Run. Time (s)
OFF	65.37	0.72	0.090
GPT-2	70.71	0.76	0.732
CodeT5	68.59	0.74	0.186
LSTM	49.26	0.61	0.198
$n$ -gram	36.89	0.51	0.137

### 7.3.5. Impact of Candidate Ranking Component

We also studied the impact of selecting heavy-ranking model on the overall performance of ARIST. In ARIST's process, any approach that can evaluate the likelihood of candidates could be applied as a heavy-ranking model. In this work, we selected several representative approaches for the evaluation. In this experiment, we randomly selected a fold of Netbeans for testing and the remaining folds for training to evaluate all variants of ARIST. In Table 11, ARIST without heavy-ranking stage (OFF) can run very fast (0.09s), yet achieved about 65.37% of top-1 accuracy. The top-1 accuracy of ARIST can be improved by 5% and 8.2% when applying CodeT5 and GPT-2, respectively. This is expected because GPT-2 and CodeT5 can capture dependencies which are longer than that by the light-ranking model. However, the running time when using these powerful models greatly increases by about 8.6 times when GPT-2 is applied. With the less powerful LMs such as vanilla LSTM and  $n$ -gram, the performance of ARIST might not be improved, or even be significantly worsened.

### 7.3.6. Impact of Static Features

To investigate the impact of the static features used in this version of ARIST: *parameter similarity*, *creating-distance*, and *accessing-recentness*, we built a variant of ARIST where the application of only these static features is disabled. In this experiment, Eclipse is used for both ARIST and the variant. We found that the full version of ARIST's top-1 accuracy is 0.4% better than that of the variant where the static features are not applied. The reason for this marginal improvement is that the heavy-ranking model applied in the current version of ARIST is GPT-2. This powerful model can dynamically learn these features to rank candidates. Still, for LSTM and  $n$ -gram which are less powerful, these features can improve 45.4% and 20% respectively in top-1 accuracy. Thus, these static features should be enabled to ensure ARIST's best performance for any models used in ARIST.

## 7.4. Time Complexity (RQ8)

All experiments were run on a server with Intel(R) Xeon(R) CPU @ 2.30GHz, 32GB RAM, and Tesla P100 GPU. ARIST took 38 hours for training. Our average recommendation time is 0.444s/request. Valid candidate identifying took 9%, while candidate reducing and ranking took 4.7% and 86.3%.

Analyzing the impact of ARIST's components on its efficiency, we found that the valid candidate identification (VCI) step significantly affects ARIST's efficiency. As seen in Table 8, ARIST with VCI recommended almost 2 times faster than when the beam search technique is applied. Meanwhile, Candidate Reduction reduces ARIST's running time by about 6 times compared to that when that component is disabled (Table 9). This is because with only less than 5% of overall running time of ARIST, the candidate reducing step significantly narrows down the candidate sets, from 7K to less than 50. With those much smaller sets of candidates, the ranking stage, which accounts for 86.3% the overall ARIST's recommending time, saves much time. As a result, ARIST's overall running time of is greatly reduced. Additionally, the ARIST's running time is also proportional to the reducing threshold (Table 10). Table 11 shows that selecting heavy-ranking models could greatly impact ARIST's efficiency. For example, ARIST with GPT-2 is 8 times slower than ARIST without the heavy-ranking stage.

## 7.5. Threats to Validity

The main threats to the validity of our work consist of internal, construct, and external threats.

**Threats to internal validity** include the influence of the hyper-parameters used in ARIST. ARIST's performance would be affected by different hyper-parameter settings, which are tuned empirically in our empirical study. Thus, there is a threat to the hyper-parameter choosing, and further improvement is possible. However, current settings have achieved acceptable performance. Another threat to validity relates to using the beam search algorithm in our adaptation of the baseline methods for recommending arguments. To reduce this threat and balance the cost of a practical evaluation, we used a relatively large beam size for the beam search algorithm.

**Threats to construct validity** relate to the suitability of our evaluation procedure. We used *Precision*, *Recall*, *top-k accuracy*, and *MRR*. They are the classical evaluation measures for argument recommendation and code completion and are used in almost all the previous AR and code completion work [5, 6, 10, 32, 11]. In fact, a user study could be the best evaluation method, which investigates how developers interact with the proposed approach and whether the proposed approach can indeed help developers find their expected arguments. In future work, we are planning to perform a human evaluation to measure actual performance of ARIST in a real-world setting.

**Threats to external validity** mainly lie in the procedure used in our experiments. Our data has only Java code. Thus, we cannot claim that similar results would have been

observed in other programming languages, especially for dynamically typed languages (untyped languages). Further studies are needed to validate and generalize our findings to other languages. Another threat relates to the quality of the datasets. To reduce this threat, we used the same very large datasets which are widely used in existing studies. A threat relates to our selections for baseline approaches. To reduce this threat, we selected the state-of-the-art  $n$ -gram with nested cache, CodeT5 [15], and the combination of GPT-2 [14] and a type-based post-processing step. Moreover,  $n$ -gram with nested cache is used to represent for the classical direction which is widely used in existing studies [10, 33, 12]. Meanwhile, CodeT5 and GPT-2 are used as a representative advanced DNN approach instead of vanilla LSTM or RNN. Moreover, GPT-2 is used as the core engine of Tabnine [16] and IntelliCode [17] which are two of the most popular AI-assisted code completers.

## 8. Related Work

**Argument recommendation.** ARIST is related to PARC by Asaduzzaman *et al.* [5] and Precise by Zhang *et al.* [6]. In comparison, there are fundamental differences between ARIST and both PARC and Precise. First, PARC and Precise follow IR direction which searches for similar arguments that have been seen in the corpus. The seen arguments which are represented by a set of static features are stored in a database. Similar arguments are found by *simhash* algorithm in PARC [5] and  $k$ -nearest neighbor algorithm in Precise [6]. Meanwhile, ARIST is a learning-based approach applying statistical language models learning the dynamic features to recommend arguments. Secondly, in the recommendation process, the syntactical and semantic constraints are enforced very differently in PARC and Precise as ARIST. In PARC and Precise, the type-compatibility constraint is used to validate the similar argument candidates at the last step. The validation becomes not very useful when the expected arguments are unseen or not sufficiently similar to the seen arguments. For ARIST, these constraints are applied to identify the set of valid candidates for the given request at the first step. Our results show that ARIST is very effective in identifying the valid argument candidates, i.e., the expected arguments are identified in +90% of the requests. This ensures that ranking models work in the valid set of candidates even for unseen expected arguments. The empirical results also show ARIST's ability to recommend unseen arguments.

**API recommendation.** There exists a rich literature of approaches on API recommendation. These studies mainly focus on suggesting methods to call. In modern IDEs such as Eclipse [34] or IntelliJ IDEA [35], APIs are suggested primarily based on program analysis, which utilizes type-compatibility and accessibility information to support completion of method calls/arguments in the current context. To improve the API recommendation performance, instead of just listing all valid candidates alphabetically, several approaches also leverage previous API usages to rank, filter, or

group candidates for better suggestions [36, 37, 38, 39, 40, 1, 3, 2]. Huang *et al.* [3] discover task-API knowledge gap and the necessity of incorporating API documentation and Stack Overflow posts to search for relevant API usages. Recently, Chen *et al.* [2] combine the textual and dependency information of source code, which analyzes control and data flow graphs to capture correlated API usages and overcome the limitation on long dependencies of token-sequence-based API recommendation. Kang *et al.* [41] apply pre-trained LM to recommend Cross-Library APIs. However, while those approaches focus on recommending method calls, ARIST is designed to predict actual parameters for method calls. Thus, ARIST and those approaches can be applied together to recommend longer code sequences.

**Code completion and suggestion.** Traditional work on code completion mainly focused on defining heuristic rules based on language-specific grammar and type-checking information to complete current code fragments [37, 38]. Hindle *et al.* [22] show that source code is highly repetitive (*natural*). Several studies [22, 42, 43, 13, 32, 33] have been proposed, exploiting this naturalness property to predict next tokens. Tu *et al.* [42] demonstrate that code also has localness property, which could be exploited by the combination of cache mechanism and traditional  $n$ -gram model. Hellendoorn and Devanbu [13] introduce nested-cache, an  $n$ -gram model integrated with nested scopes, locality, and dynamism, which showed promising performance on code completion task.

In recent years, deep neural network (DNN) language models have made great progress while applying to learn source code [10, 44, 11, 45, 4, 46]. Neural networks such as RNN and vanilla LSTM could capture longer dependencies compared to  $n$ -gram model, which could benefit code completer in large context circumstances [45]. In fact, training a DNN language model requires large datasets, which might produce a huge vocabulary set and cause Out-of-Vocabulary (OOV) problem [44, 10]. Li *et al.* [44] introduce a pointer mixture model to mitigate the OOV issue. Karampatsis *et al.* [10] propose a solution for large-scale open-vocabulary DNN language models, which employ the BPE algorithm, beam search algorithm, and cache mechanism.

Although these code completion approaches are originally designed to predict any next tokens and not whole arguments, they can be adapted for argument recommendation (e.g., applying beam search algorithm). We showed that ARIST outperforms the representative code completion approaches. In fact, any existing model for code completion can be applied as a ranking model in our argument recommendation approach.

**Learning-based approaches for SE tasks.** Several studies have been proposed for specific SE tasks including code suggestion [47, 22, 48], program synthesis [49, 50], pull request description generation [51, 52], code summarization [53, 54, 55], code clones [56], fuzz testing [57], fault localization [58], and program repair [59, 60].

## 9. Conclusion

We introduce ARIST, a novel automated argument recommendation approach combining program analysis (PA) and statistical language models (LMs) in the recommendation process. In ARIST, the statistical LMs are used to suggest the frequent/natural candidates identified by applying PA. Meanwhile, PA is applied to navigate LMs working on the set of valid candidates. Our empirical evaluation on the large datasets of real-world projects shows that ARIST improves the state-of-the-art approach up to 19% and 18% in both precision and recall for recommending arguments of *frequently-used libraries*. For *general argument recommendation task*, ARIST also outperforms the baseline approaches by up to 31% and 125% top-1 accuracy. Moreover, for new projects, ARIST achieves more than 60% top-3 accuracy after training with a larger dataset. Especially, with solely personalizing the light-ranking stage which is deployed on developers' machines, ARIST can correctly rank the expected arguments at the top-1 position in up to +7/10 recommendation requests in less than 0.5 seconds while preserving the privacy of developers.

## References

- [1] A. T. Nguyen, M. Hilton, M. Codoban, H. A. Nguyen, L. Mast, E. Rademacher, T. N. Nguyen, D. Dig, Api code recommendation using statistical learning from fine-grained changes, in: Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering, 2016, pp. 511–522.
- [2] C. Chen, X. Peng, Z. Xing, J. Sun, X. Wang, Y. Zhao, W. Zhao, Holistic combination of structural and textual code information for context based api recommendation, IEEE Transactions on Software Engineering (2021).
- [3] Q. Huang, X. Xia, Z. Xing, D. Lo, X. Wang, Api method recommendation without worrying about the task-api knowledge gap, in: 2018 33rd IEEE/ACM International Conference on Automated Software Engineering (ASE), IEEE, 2018, pp. 293–304.
- [4] X. He, L. Xu, X. Zhang, R. Hao, Y. Feng, B. Xu, Pyart: Python api recommendation in real-time, in: 2021 IEEE/ACM 43rd International Conference on Software Engineering, IEEE, 2021, pp. 1634–1645.
- [5] M. Asaduzzaman, C. K. Roy, K. A. Schneider, Parc: Recommending api methods parameters, in: 2015 IEEE International Conference on Software Maintenance and Evolution, IEEE, 2015, pp. 330–332.
- [6] C. Zhang, J. Yang, Y. Zhang, J. Fan, X. Zhang, J. Zhao, P. Ou, Automatic parameter recommendation for practical api usage, in: 2012 34th International Conference on Software Engineering (ICSE), IEEE, 2012, pp. 826–836.
- [7] H. Liu, Q. Liu, C.-A. Staicu, M. Pradel, Y. Luo, Nomen est omen: Exploring and exploiting similarities between argument and parameter names, in: Proceedings of the 38th International Conference on Software Engineering, 2016, pp. 1063–1073.
- [8] A. Rice, E. Aftandilian, C. Jaspán, E. Johnston, M. Pradel, Y. Arroyo-Paredes, Detecting argument selection defects, Proceedings of the ACM on Programming Languages 1 (OOPSLA) (2017) 1–22.
- [9] M. Pradel, S. Heiniger, T. R. Gross, Static detection of brittle parameter typing, in: Proceedings of the 2012 International Symposium on Software Testing and Analysis, 2012, pp. 265–275.
- [10] R.-M. Karampatsis, H. Babii, R. Robbes, C. Sutton, A. Janes, Big code!= big vocabulary: Open-vocabulary models for source code, in: 2020 IEEE/ACM 42nd International Conference on Software Engineering (ICSE), IEEE, 2020, pp. 1073–1085.
- [11] F. Liu, G. Li, Y. Zhao, Z. Jin, Multi-task learning based pre-trained language model for code completion, in: Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering, 2020, pp. 473–485.
- [12] V. J. Hellendoorn, S. Proksch, H. C. Gall, 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.
- [13] V. J. Hellendoorn, P. Devanbu, Are deep neural networks the best choice for modeling source code?, in: Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering, Association for Computing Machinery, New York, NY, USA, 2017, p. 763–773.
- [14] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al., Language models are unsupervised multitask learners, OpenAI blog 1 (8) (2019) 9.
- [15] Y. Wang, W. Wang, S. Joty, S. C. Hoi, Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation, in: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, 2021, pp. 8696–8708.
- [16] Deep TabNine, <https://www.tabnine.com/blog/deep/>, accessed: January 2022 (2022).
- [17] A. Svyatkovskiy, S. K. Deng, S. Fu, 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.
- [18] Oracle, Java SE Specification, <https://docs.oracle.com/javase/specs/index.html>, [Online; accessed 19-Aug-2021] (2018).
- [19] S. N. et al., ARist, <https://ttrangnguyen.github.io/ARist/>, [Online; 05-Oct-2022] (2022).
- [20] E. documentation, Eclipse JDT API Specification, <https://help.eclipse.org/latest/index.jsp>, accessed: December 2021 (2021).
- [21] Oracle, Chapter 18. Syntax, <https://docs.oracle.com/javase/specs/jls/se7/html/jls-18.html>, [Online; accessed 19-Aug-2021] (2018).
- [22] A. Hindle, E. T. Barr, M. Gabel, Z. Su, P. Devanbu, On the naturalness of software, Communications of the ACM 59 (5) (2016) 122–131.
- [23] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, Advances in neural information processing systems 30 (2017).
- [24] R. Sennrich, B. Haddow, A. Birch, Neural machine translation of rare words with subword units, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Berlin, Germany, 2016, pp. 1715–1725.
- [25] J. Spec, Blocks and Statements, <https://docs.oracle.com/javase/specs/jls/se7/html/jls-14.html>, [Online; 05-Mar-2022] (2022).
- [26] B. Dagenais, L. Hendren, Enabling static analysis for partial java programs, in: Proceedings of the 23rd ACM SIGPLAN conference on Object-oriented programming systems languages and applications, 2008, pp. 313–328.
- [27] C. Parsing, Speech and language processing (2009).
- [28] T. G. Dietterich, Ensemble methods in machine learning, in: International workshop on multiple classifier systems, Springer, 2000.
- [29] M. Allamanis, C. Sutton, Mining source code repositories at massive scale using language modeling, in: 2013 10th Working Conference on Mining Software Repositories (MSR), IEEE, 2013, pp. 207–216.
- [30] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang, et al., Codebert: A pre-trained model for programming and natural languages, in: Findings of the Association for Computational Linguistics: EMNLP 2020, 2020, pp. 1536–1547.
- [31] A. Kanade, P. Maniatis, G. Balakrishnan, K. Shi, Learning and evaluating contextual embedding of source code, in: International Conference on Machine Learning, PMLR, 2020, pp. 5110–5121.
- [32] S. Nguyen, T. Nguyen, Y. Li, S. Wang, Combining program analysis and statistical language model for code statement completion, in: 2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE), IEEE, 2019, pp. 710–721.
- [33] M. Allamanis, E. T. Barr, P. Devanbu, C. Sutton, A survey of machine learning for big code and naturalness, ACM Computing Surveys (CSUR) 51 (4) (2018) 1–37.

- [34] E. Foundation, Eclipse Code Recommenders, <http://www.eclipse.org/recommenders>, [Online; accessed 19-Aug-2021] (2018).
- [35] JetBrains, IntelliJ Idea, <https://www.jetbrains.com/help/idea/auto-completing-code.html>, [Online; accessed 21-Aug-2021] (2021).
- [36] M. Bruch, M. Monperrus, M. Mezini, Learning from examples to improve code completion systems, in: Proceedings of the 7th joint meeting of the European software engineering conference and the ACM SIGSOFT symposium on the foundations of software engineering, 2009, pp. 213–222.
- [37] D. Hou, D. M. Pletcher, Towards a better code completion system by api grouping, filtering, and popularity-based ranking, in: Proceedings of the 2nd International Workshop on Recommendation Systems for Software Engineering, 2010, pp. 26–30.
- [38] D. M. Pletcher, D. Hou, An evaluation of the strategies of sorting, filtering, and grouping api methods for code completion, in: 2013 IEEE International Conference on Software Maintenance, IEEE Computer Society, Los Alamitos, CA, USA, 2011, pp. 233–242.
- [39] M. M. Rahman, C. K. Roy, D. Lo, Rack: Automatic api recommendation using crowdsourced knowledge, in: 2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER), Vol. 1, IEEE, 2016, pp. 349–359.
- [40] X. Gu, H. Zhang, D. Zhang, S. Kim, Deep api learning, in: Proceedings of the 2016 24th ACM SIGSOFT international symposium on foundations of software engineering, 2016, pp. 631–642.
- [41] Y. Kang, Z. Wang, H. Zhang, J. Chen, H. You, Apirecx: Cross-library api recommendation via pre-trained language model, in: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, 2021, pp. 3425–3436.
- [42] Z. Tu, Z. Su, P. Devanbu, On the localness of software, in: Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, 2014, pp. 269–280.
- [43] T. T. Nguyen, A. T. Nguyen, H. A. Nguyen, T. N. Nguyen, A statistical semantic language model for source code, in: Proceedings of the 9th Joint Meeting on Foundations of Software Engineering, 2013, pp. 532–542.
- [44] J. Li, Y. Wang, M. R. Lyu, 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.
- [45] C. Liu, X. Wang, R. Shin, J. E. Gonzalez, D. Song, Neural code completion (2016).
- [46] M. Izadi, R. Gismondi, G. Gousios, Codefill: Multi-token code completion by jointly learning from structure and naming sequences, in: Proceedings of the 30th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, 2022.
- [47] S. Nguyen, H. Phan, T. Le, T. N. Nguyen, Suggesting natural method names to check name consistencies, in: 2020 IEEE 42nd International Conference on Software Engineering, IEEE, 2020, pp. 1372–1384.
- [48] M. Allamanis, H. Peng, C. Sutton, A convolutional attention network for extreme summarization of source code, in: International conference on machine learning, PMLR, 2016, pp. 2091–2100.
- [49] M. Amodio, S. Chaudhuri, T. W. Reps, Neural attribute machines for program generation, CoRR (2017).
- [50] T. Gvero, V. Kuncak, Synthesizing java expressions from free-form queries, in: Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, 2015, pp. 416–432.
- [51] X. Hu, G. Li, X. Xia, D. Lo, Z. Jin, Deep code comment generation, in: 2018 IEEE/ACM 26th International Conference on Program Comprehension (ICPC), IEEE, 2018, pp. 200–20010.
- [52] Z. Liu, X. Xia, C. Treude, D. Lo, S. Li, Automatic generation of pull request descriptions, in: 34th IEEE/ACM International Conference on Automated Software Engineering, IEEE, 2019, pp. 176–188.
- [53] S. Iyer, I. Konstas, A. Cheung, L. Zettlemoyer, Summarizing source code using a neural attention model, in: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2016, pp. 2073–2083.
- [54] A. Mastropaolo, S. Scalabrino, N. Cooper, D. N. Palacio, D. Poshvanyk, R. Oliveto, G. Bavota, Studying the usage of text-to-text transfer transformer to support code-related tasks, in: 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), IEEE, 2021, pp. 336–347.
- [55] Y. Wan, Z. Zhao, M. Yang, G. Xu, H. Ying, J. Wu, P. S. Yu, Improving automatic source code summarization via deep reinforcement learning, in: Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, 2018, pp. 397–407.
- [56] L. Li, H. Feng, W. Zhuang, N. Meng, B. Ryder, Cclerner: A deep learning-based clone detection approach, in: International Conference on Software Maintenance and Evolution, IEEE, 2017, pp. 249–260.
- [57] P. Godefroid, H. Peleg, R. Singh, Learn&fuzz: Machine learning for input fuzzing, in: 2017 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE), IEEE, 2017, pp. 50–59.
- [58] Y. Li, S. Wang, T. N. Nguyen, Fault localization with code coverage representation learning, in: 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), IEEE, 2021, pp. 661–673.
- [59] N. Jiang, T. Lutellier, L. Tan, Cure: Code-aware neural machine translation for automatic program repair, in: 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), IEEE, 2021, pp. 1161–1173.
- [60] Y. Ding, B. Ray, P. Devanbu, V. J. Hellendoorn, Patching as translation: the data and the metaphor, in: 2020 35th IEEE/ACM International Conference on Automated Software Engineering, IEEE, 2020, pp. 275–286.