

Mining API Mapping for Language Migration

Hao Zhong¹, Suresh Thummalapenta⁴, Tao Xie⁴, Lu Zhang^{2,3}, Qing Wang¹

¹Laboratory for Internet Software Technologies, Institute of Software, Chinese Academy of Sciences, Beijing, 100190, China

²Key Laboratory of High Confidence Software Technologies (Peking University), Ministry of Education, China

³Institute of Software, School of Electronics Engineering and Computer Science, Peking University, China

⁴Department of Computer Science, North Carolina State University, Raleigh, NC 27695-8206, USA

zhonghao@itechs.iscas.ac.cn, {sthumma,txie}@ncsu.edu, zhanglu@sei.pku.edu.cn, wq@itechs.iscas.ac.cn

ABSTRACT

Since the inception of programming languages, researchers and practitioners developed various languages such as Java and C#. To address business requirements and to survive in competing markets, companies often have to release different versions of their projects in different languages. Migrating projects from one language to another language (such as from Java to C#) manually is a tedious and error-prone task. To reduce manual effort or human errors, tools can be developed for automatic translation of projects from one language to another, but these tools require the knowledge of how Application Programming Interfaces (APIs) of one language are mapped to the APIs of the other language, referred to as API mapping relations. In this paper, we propose a novel approach that mines API mapping relations from one language to another using API client code, referred to as MAM (Mining API Mapping). MAM accepts a set of projects with versions in two languages and mines API mapping relations between those two languages based on how APIs are used by the two versions. These mined API mapping relations assist in translation of projects from one language to another. For MAM, we implemented a tool and conducted two evaluations to show the effectiveness of MAM. The results show that our tool mines 25,805 unique mapping relations of APIs between Java and C# with more than 80% accuracy. The results also show that mined API mapping relations reduce 54.4% compilation errors and 43.0% defects during translation of projects with an existing translation tool, called Java2CSharp. The reduction in compilation errors and defects is due to our new mined mapping relations that are not available with the existing translation tools.

1. INTRODUCTION

A programming language serves as a means for instructing computers to achieve a programming task at hand. Since their inception, various programming languages came into existence due to several factors such as existence of many platforms or requirements for different programming styles. The HOPL¹ website lists 8,512 different programming languages. To address business require-

ments and to survive in competing markets, companies or open source organizations often have to release different versions of their projects in different languages. For example, many well-known projects such as Lucene² and WordNet³ provide multiple versions in different languages. For some open source projects, although their project teams do not officially provide multiple versions, external programmers often create their versions in different languages. For example, the WordNet team does not provide a C# version, but Simpson and Crowe developed a C# version of WordNet.Net⁴. Totally, as described by Jones [6], about one-third of the existing projects have multiple versions in different languages.

Translating projects from one language to another language (e.g., from Java to C#) manually is a tedious and error-prone task. Some companies have suffered from huge losses because of failures in language translation. For example, Terekhov and Verhoef [10] stated that at least three companies went bankrupt and another company lost 50 million dollars due to failed language translation projects. A natural way to address this issue is to develop a translation tool that can automatically translate projects from one language to another. However, it is challenging to develop such a translation tool as the translation tool should have knowledge of how Application Programming Interfaces (APIs) in one programming language are mapped to APIs in the other language. In the literature, there exist approaches [3, 7, 16] that address the problem of language translation partially. These approaches expect programmers to describe how APIs of one language is mapped to APIs of another language. As programming languages provide a large number of APIs, writing mappings manually for all APIs is tedious and error-prone. As a result, existing approaches [3, 7, 16] support only a subset of APIs for translation. Such a limitation results in many compilation errors in translated projects and limits these approaches' usage in practice.

In this paper, we propose a novel approach that automatically mines how APIs of one language are mapped to APIs of another language. We refer to this mapping as *mapping relations of APIs* and our approach as MAM. MAM mines mapping relations based on API usages in the client code rather than based on API implementations with three major reasons. (1) API libraries often do not come with source files especially for those C# libraries. (2) Relations mined based on API implementations often have lower confidence than relations mined based on API usages. The reason is that API implementations have only one data point for analysis, whereas API usages can have many data points (i.e., call sites) for mining. (3) Mapping relations of APIs are often complex and cannot be mined based on the information available in the API implementations. First, mapping parameters of an API method in one language

¹<http://hop1.murdoch.edu.au>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICSE '2010 Cape Town, South Africa

Copyright 2010 ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

²<http://lucene.apache.org/>

³<http://wordnet.princeton.edu/>

⁴<http://opensource.ebswift.com/WordNet.Net/>

with an API method in the other language can be complex. For example, consider the following two API methods in Java and C#:

m_1 in Java: `BigDecimal java.math.BigDecimal.multiply (BigDecimal p_1^1)`

m_2 in C#: `Decimal System.Decimal.Multiply (Decimal p_1^2 , Decimal p_2^2)`

Here, m_1 has a receiver, say v_1^1 , of type `BigDecimal` and has one parameter p_1^1 , whereas m_2 has two parameters p_1^2 and p_2^2 . For these two API methods, v_1^1 is mapped to p_1^2 , and p_1^1 is mapped to p_2^2 . Second, an API method of one language can be mapped to more than one API method in the other language. For example, consider the following two API methods:

m_3 in Java: `E java.util.LinkedList.removeLast()`

m_4 in C#: `void System.Collections.Generic.LinkedList.RemoveLast()`

Although the method names of m_3 and m_4 are the same, m_3 in Java cannot be directly mapped with m_4 in C#. The reason is that m_3 in Java returns the last element removed from the list, whereas m_4 does not return any element. Therefore, m_3 is mapped to two API methods m_4 and m_5 (shown below) in C#. The API method m_5 returns the last element and should be called before calling m_4 .

m_5 in C#: `void System.Collections.Generic.LinkedList.Last()`

To deal with the complexity of mining API mapping, we construct a graph, referred to as *API transformation graph* (ATG), for aligned methods of the client code in both languages. These ATGs precisely capture inputs and outputs of API methods, and help mine complex mapping relations of API methods.

This paper makes the following major contributions:

- The first approach that mines API mapping between different languages using API client code, referred to as MAM. MAM addresses an important and yet challenging problem that is not addressed by previous work on language translation.
- API transformation graphs (ATGs) proposed to capture inputs and outputs of API methods, and a technique for building ATGs and comparing built ATGs for mining API mapping. As ATGs describe data dependencies among inputs and outputs of API methods, MAM is able to mine complex mapping relations between API methods from two languages.
- A tool implemented for MAM and two evaluations on 15 projects that include 18,568 classes and 109,850 methods with both Java and C# versions. The results show that our tool mines 6,695 mapping relations of API classes with accuracy of 86.7% and 19,110 mapping relations of API methods with accuracy of 90.0%. The results also show that the mined API mapping relations reduce 55.4% of compilation errors and 43.0% defects during translation of projects from Java to C# using Java2CSharp.

The rest of this paper is organized as follows. Section 2 presents definitions. Section 3 illustrates our approach using an example. Section 4 presents our approach. Section 5 presents our evaluations. Section 6 discusses issues of our approach. Section 7 presents related work. Finally, Section 8 concludes.

2. DEFINITIONS

We next present definitions of terms used in the rest of the paper.

API. An Application Programming Interface (API) [8] is a set of classes and methods provided by frameworks or libraries.

API library. An API library is a framework or library that provides reusable API classes and methods.

Client code. Client code is application code that reuses or extends API classes and methods provided by API libraries.

The definitions of API library and client code are relative to each other. For example, Lucene uses J2SE⁵ as an API library, whereas Nutch⁶ uses Lucene as an API library. Therefore, we consider

⁵<http://java.sun.com/j2se/1.5.0/>

⁶<http://lucene.apache.org/nutch/>

Lucene as client code and API library for the J2SE API library and Nutch, respectively. In general, for programmers of client code, source files of API libraries are often not available.

Mapping relation. For a set of entities E_1 defined by a language L_1 and another set of entities E_2 defined by another language L_2 , a mapping relation is a triple $\langle E_1, E_2, b_i \rangle$ where E_1 and E_2 have the same behavior b_i .

We use mapping relations of API classes for translating data such as variable, parameters, and constants, so we require that two mapped API classes have the same behavior of storing data, referred to as the s behavior.

1-to-1 mapping relation of API classes. For an API class c_1 defined by L_1 and an API class c_2 defined by L_2 , an 1-to-1 mapping relation of API classes is a triple $\langle c_1, c_2, s \rangle$, where s denotes the s behavior.

One API class defined by L_1 can have more than one 1-to-1 mapping relations with API classes defined by L_2 . For example, data in `java.util.ArrayList` of Java can be stored in either `System.Collections.ArrayList` or `System.Collections.Generic.List` of C#, so the Java class has two 1-to-1 mapping relations with the two C# classes.

1-to-many mapping relation of API classes. For an API class c_1 defined by L_1 and a set of API classes C_2 defined by L_2 , an 1-to-many mapping relation of API classes is a triple $\langle c_1, C_2, s \rangle$, where s is the s behavior.

For example, the current time in `java.lang.System` of Java is stored in `System.Environment` of C#, whereas the environment settings in `java.lang.System` is stored in `System.Environment` of C#, so the Java class has an 1-to-many mapping relations with the two C# classes.

We use mapping relations of API methods for translating API methods that use input to produce desirable outputs, so we require two mapped API methods have the same behavior of inputs, outputs, and functionalities, referred as the “ r ” behavior.

Merged API method. A merged API method cm is a set of API methods M combined by inputs and outputs, so the t behavior of cm is the combination of M . Consider two API methods m_1 and m_2 defined in classes C_1 and C_2 of L_1 , respectively, with the following signatures:

m_1 signature: $o_1 \ C_1.m_1(inp_1^1, inp_2^1, \dots, inp_m^1)$

m_2 signature: $o_2 \ C_2.m_2(inp_1^2, inp_2^2, \dots, inp_n^2)$

We merge methods m_1 and m_2 to create a new merged API method m_{new} if the output o_1 of m_1 is used either as a receiver or as a parameter for m_2 (i.e., $o_1 == C_2$ or $o_1 == inp_i^2$) in client code. The signature of the new merged API method m_{new} is shown below:

m_{new} signature: $o_2 \ m_{new}(inp_1^1, inp_2^1, \dots, inp_m^1, inp_1^2, inp_2^2, \dots, inp_n^2)$

For the Java code example shown in Figure 1, consider the `file` variable, which is a return variable for the constructor and a receiver object for the `exists` method. As the output of one API method is passed as receiver object of another API method, we can combine these two methods to create a new merged API method `boolean File.exists(string)` as shown in Figure 3 (b). The merged API method accepts a `string` parameter that represents a file name and returns a boolean value that describes whether a file exists or not. For simplicity, we consider each API method as a merged API method when we define mapping relations of API methods.

Mapping relation of API methods. For a merged API method cm_1 defined by L_1 and a merged API method cm_2 defined by L_2 , a mapping relation of API methods is a triple $\langle c_1, C_2, t \rangle$, where t is the t behavior.

Java code:

```
1 File file = new File("test");
2 Boolean b = file.exists();
```

Translated C# code:

```
3 FileInfo file = new FileInfo("test");
4 Boolean b = System.IO.File.Exists(file.FullName) ||
  System.IO.Directory.Exists(file.FullName);
```

Figure 1: Java code and its translated C# code.

```
IndexFiles.java:
5 public class IndexFiles {
6   static final File INDEX_DIR = new File("index");
7   public static void main(String[] args) {
8     ...
9     if (INDEX_DIR.exists()) {...}
10    ...
11    INDEX_DIR.delete();
12  }
13 }

IndexFiles.cs:
10 class IndexFiles{
11   internal static readonly System.IO.FileInfo INDEX_DIR
12   = new System.IO.FileInfo("index");
13   public static void Main(String[] args){
14     ...
15     bool tmpBool;
16     if (System.IO.File.Exists(INDEX_DIR.FullName))
17       tmpBool = true;
18   else
19     tmpBool = System.IO.Directory
20       .Exists(INDEX_DIR.FullName);
21   ...
22 } }
```

Figure 2: Two versions (Java and C#) of client code.

3. EXAMPLE

We next use an example to illustrate challenges in mining API mapping relations. Figure 1 shows a Java code example and its translated C# code. This Java code example accepts a string input that represents the name of a file or directory and returns a boolean value that describes whether the file or directory exists. To achieve this functionality, the code example declares a local variable, called `file`, of type `java.io.File` and invokes the `exists` method. The method takes the string input and `file` as its inputs and produces the desirable boolean value. Here, we consider `file` (a receiver) as a special input for the `exists` method.

To translate this code example into C#, a translation tool needs to know mapping relations of API classes, so it can translate inputs, outputs, and variables into C#. For example, the translation tool needs to know the mapped API class for `java.io.File` in C# to translate the variable `file` to C#. In addition, the translation tool needs to know the mapped API methods, so it can add code for invoking proper API methods that use translated inputs and variables to produce desirable outputs. For this example, the translation tool adds code for invoking the `Exists` method and the `FullName` method to achieve the functionality. Here, we consider field accesses as special type of method calls.

Some projects such as Lucene have both Java and C# versions. Our approach has three major steps to mine the preceding two types of mapping relations of APIs from these projects.

Aligning client code. First, our approach aligns classes and methods (between the two versions) that implement similar functionality. As these aligned classes and methods implement similar functionality, our approach mines mapping relations of API classes and methods from these aligned classes and methods. To align classes and methods, our approach uses a mapping algorithm based on similarities in the names of classes and methods.

Aligning client code based on names of classes and methods is

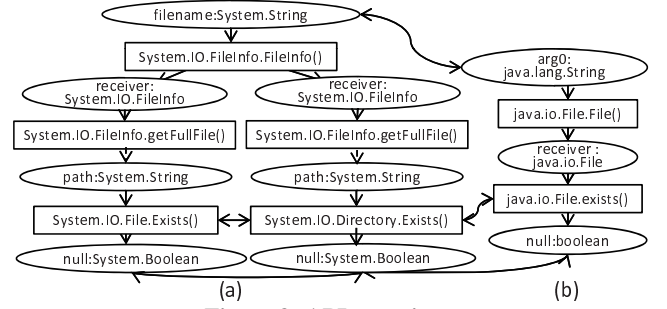


Figure 3: API mapping

based on our observation on how many existing projects such as rasp⁷ are migrated from one language to another. We observed that while migrating the rasp project from C# to Java, programmers first renamed source files from C# to Java and systematically addressed the compilation errors by replacing C# APIs with Java APIs. During this procedure, names of classes, methods, fields of classes, or local variables in methods often remain the same or similar between the two versions. Therefore, we use name similarities for aligning client code of the two versions. For example, our approach aligns `IndexFiles.java` with the `IndexFiles.cs` (shown in Figure 2) since the names of their classes and methods are similar.

Mining API mapping of classes. Next, our approach mines mapping relations of API classes by comparing entities such as names of fields in aligned classes, or variable names or constants in aligned methods. Our approach uses a text-based similarity metric for comparing these entities and considers the entities as similar if the metric is greater than a given threshold. These mapping relations of API classes help translate variables from one language to another. For example, our approach identifies the constant value “index” in Lines 6 (Java) and 11 (C#) (Figure 2) and maps the API classes associated with these constants. Based on this constant value, our approach maps the API class `java.io.File` of Java to `System.IO.FileInfo` of C#.

Mining API mapping of methods. After mapping API classes between the two languages, our approach maps API methods. Mapping API methods is challenging since often an API method of one language can be mapped to multiple API methods of the other language. Furthermore, mapping relations of API methods should also describe how parameters and returns are mapped between them. To address these challenges, our approach constructs a graph, referred to as *API Transformation Graph (ATG)*, for each aligned method of the client code in both languages. These ATGs precisely capture inputs and outputs of API methods, and help mine mapping relations of API methods. Figure 3 shows a mapping relation between API method `Exists` from one language to another. (@Hao, could you please add notations used in the figure here.) Section 4.2 presents more details on how we mine these mapping relations of API methods using ATGs. Our approach uses these mapping relations to assist translation tools such as `Java2CSharp` for conducting language translation.

4. APPROACH

Our approach accepts a set of projects as data sources and mines API mapping relations between two languages L_1 and L_2 . As mined API mapping relations describe mapping relations of APIs between the two languages, this mapping is useful for language translation between the two languages. For each project used as a data source, our approach requires at least two versions of the

⁷<http://sourceforge.net/projects/r-asp/>

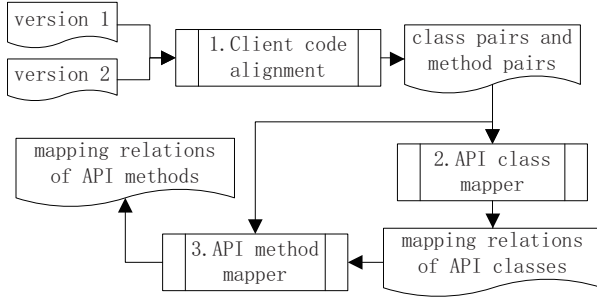


Figure 4: Overview of our approach

project (one version in L_1 and the other version in L_2). Figure 4 shows the overview of our approach.

First, our approach aligns client code in languages L_1 and L_2 so that the aligned source files implement similar functionalities (Section 4.1). Second, our approach mines mapping relations of API classes (Section 4.2). Finally, our approach mines mapping relations of API methods (Section 4.3) defined by the mapped API classes.

4.1 Aligning API Client Code

Initially, our approach accepts two versions of a project (one version in L_1 and the other version in L_2) and aligns classes and methods of the two versions. Aligned classes or methods between the two versions implement a similar functionality. As they implement a similar functionality, APIs used by these classes or methods can be replaceable.

To align classes and methods of the two versions, our approach uses name similarities between entities (such as class names or method names) defined by the two versions of the project. In our approach, we have two different kinds of entity names: entity names defined by the two versions of the project and entity names of third-party libraries used by the two versions of the project. The first kind often comes from the same programmer or the same team, or programmers may refer to existing versions for naming entities such as classes, methods, and variables. Therefore, name similarity of the first kind is often reliable to distinguish functionalities compared to the second kind. Our approach uses Simmetrics⁸ to calculate name similarities.

We next describe how our approach aligns client-code classes. The first step is to find candidate class pairs by names. For two sets of classes (C and C'), our approach returns candidate class pairs (M) with a similarity greater than a given threshold, referred to as $SIM_THRESHOLD$. As some projects may have more than one class with the same name, M may contain more than one mapping pair for a class in a version. To align those classes, our approach uses package names of these classes to refine M and returns only one mapping pair with the maximum similarity⁹.

In each aligned class pair, our approach further aligns methods within the class pair. The algorithm for methods is similar to the algorithm for classes and also may return more than one candidate method pair due to overloading. Here, the algorithm for methods relies on criteria such as the number of parameters and names of parameters to refine candidate method pairs. For the example shown in Section 3, our approach correctly aligns the class `IndexFiles` and the method `main` in Java to the class `IndexFiles` and the method `Main` in C#, respectively, as their names are quite similar.

4.2 Mapping API classes

⁸<http://sourceforge.net/projects/simmetrics/>

⁹For C#, we refer to namespace names for package names.

In this step, our approach mines mapping relations of API classes. As described in Section 2, mapping relations of API classes are used to translate variables. Consequently, our approach mines mapping relations of API classes based on how aligned client code declares variables such as fields of aligned classes, parameters, and local variables of aligned methods. In particular, for each aligned class pair $\langle c_1, c_2 \rangle$, our approach analyzes each field pair $\langle f_1, f_2 \rangle$ and considers $\langle f_1.type, f_2.type \rangle$ as a relation, if the similarity between $f_1.name$ and $f_2.name$ is greater than $SIM_THRESHOLD$. Similarly, for each aligned method pair $\langle m_1, m_2 \rangle$, our approach analyzes each local variable pair $\langle var_1, var_2 \rangle$ and considers $\langle var_1.type, var_2.type \rangle$ as a relation, if the similarity between $var_1.name$ and $var_2.name$ is greater than $SIM_THRESHOLD$. Also, our approach analyzes each parameter pair $\langle para_1, para_2 \rangle$ of m_1 and m_2 , and our approach considers $\langle para_1.type, para_2.type \rangle$ as a relation when the similarity between $para_1.name$ and $para_2.name$ is greater than $SIM_THRESHOLD$.

For the example shown in Section 3, our approach mines the mapping relation between `java.io.File` and `System.IO.FileInfo` based on the mapped fields of Lines 6 and 11 (Figure 2). The mapping relation of API classes helps translate the variable declared in Line 1 (Figure 1) to the variable declared in Line 3 (Figure 1).

4.3 Mapping API methods

In this step, our approach mines mapping relations of API methods. This step has two major sub-steps. First, our approach builds a graph, referred to as API transformation graph, for each client code method. Second, our approach compares the two graphs of each pair of client-code methods for mining mapping relations of API methods.

4.3.1 API Transformation Graph

We propose API Transformation Graphs (ATGs) to help deal with two major challenges. (1) Mapping parameters of an API method in one language with parameters of an API method in the other language can be complex. (2) An API method of one language can be mapped to more than one API method in the other language.

An ATG of a client-code method m is a directed graph $G\langle N_{data}, N_m, E \rangle$. N_{data} is a set of the fields F of m 's declaring class, local variables V of m , parameters P_1 of m , parameters P_2 of API methods invoked by m , and returns R of all methods. N_m is a set of methods invoked by m . E is a set of directed edges. An edge $d_1 \rightarrow d_2$ from a datum $d_1 \in N_{data}$ to a datum $d_2 \in N_{data}$ denotes that d_2 is data-dependent on d_1 , referred to data dependency from d_1 to d_2 . An edge $d_1 \rightarrow m_1$ from a datum $d_1 \in N_{data}$ to a method $m_1 \in N_m$ denotes that d_1 is a parameter or receiver of m_1 . An edge $m_1 \rightarrow d_1$ from a method $m_1 \in N_m$ to a datum $d_1 \in N_{data}$ denotes that d_1 is the return of m_1 .

4.3.2 Building API Transformation Graphs

Our approach builds an ATG for each method m in the client code. ATG includes information such as inputs and outputs for each client-code method. In particular, for each method m , our approach first builds subgraphs for its variables, API methods, and field accesses. Our approach adds additional edges to the built ATG (and sub-graphs inside ATG) representing data dependencies among built sub-graphs. We use two notations for representing nodes in the ATG. A rectangle represents a method labeled with the method name, whereas an ellipse represents a datum such as fields, local variables, and parameters. An ellipse is labeled as " $n:t$ ", where n is the name of the variable and t is its type. We use the following rules for adding nodes and edges to the ATG.

1. $\forall f \in F \cup V \cup P_1$, our approach adds a node to the built ATG.

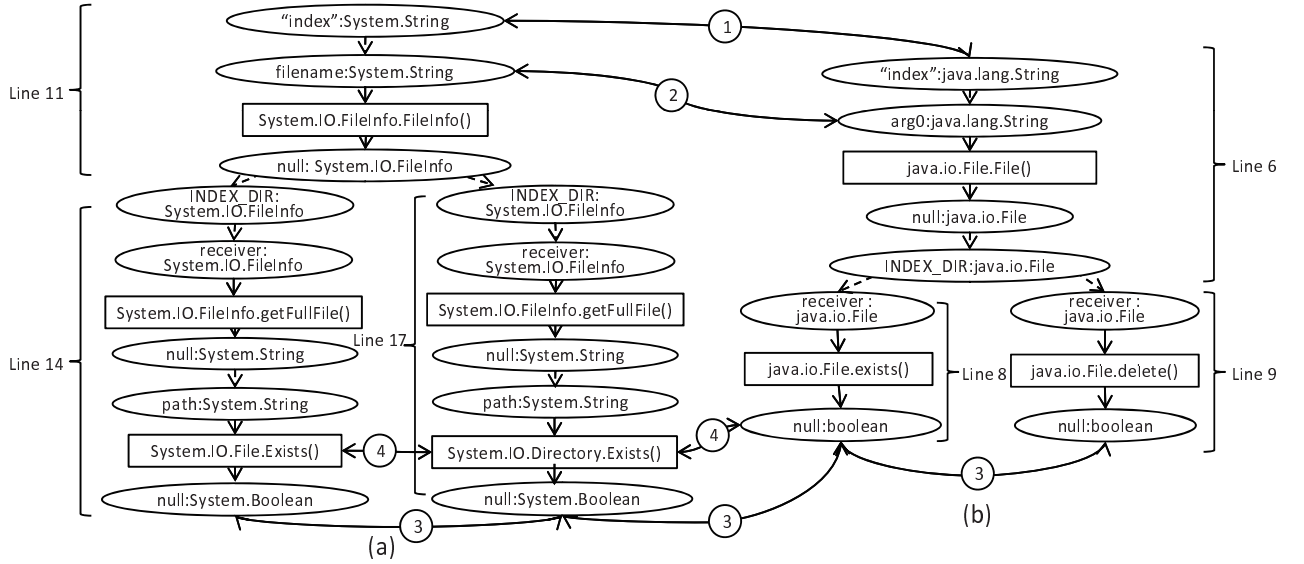
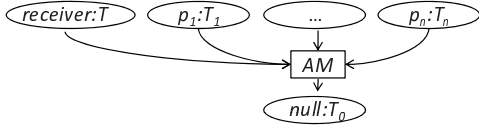


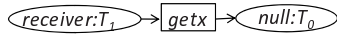
Figure 5: Built ATGs and the main steps of comparing ATGs

The reason for considering these variables such as fields in the declaring class or local variables in method m used in client code is that these variables are useful to analyze data dependencies among API methods.

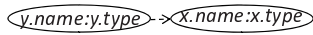
2. \forall API methods of the form " T_0 $T.AM(T_1p_1, \dots, T_n p_n)$ " invoked by method m , our approach adds a receiver node (of type T) and parameter nodes to the built ATG as shown below. Our approach does not add a receiver node for static API methods.



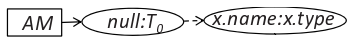
3. $\forall f \in F \cup V$, if f is a non-primitive variable of type T_1 and a field x of T_1 is accessed as $f.x$, our approach adds nodes to the built ATG as shown below. As Java often uses getters and setters whereas C# often use field accesses, our approach treats field accesses as a special type of method calls.



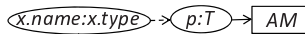
4. \forall statements of the form $x = y$, where $x \in F \cup V \wedge y \in F \cup V$, our approach adds an edge from y to x . This edge represents that x is data-dependent on y .



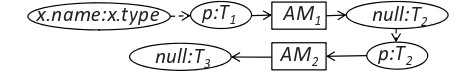
5. \forall statements of the form $x = AM()$, where $x \in F \cup V$, our approach adds an edge from AM to x if the return of AM is assigned to x . This edge represents that x is data-dependent on the return of AM .



6. \forall API methods $AM(x)$ invoked by method m , our approach adds an edge from x to the parameter node of AM . This edge represents that the parameter of AM is data-dependent on x .



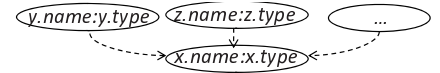
7. \forall statements of the form $m_2(m_1(x))$, our approach adds an edge from the return node of m_1 to the parameter node of m_2 . This edge represents that the parameter of m_2 is data-dependent on the return of m_1 .



8. \forall statements of the form $x.m()$, our approach adds an edge from x to m since x is the receiver of m . This edge represents that the receiver of m is data-dependent on x .



9. \forall statements of the form $x = y op z op \dots, op \in \{+, -, *, /\}$, our approach adds edges from y , z , and others to x , since these variables are connected by binary operations and the return is assigned to x . The edge denotes the data dependency from y , z , and other variables to x . For simplicity, our approach ignores op info. We discuss the issue in Section 6.



For each method m in the client code, our approach applies the preceding rules for each statement from the beginning to the end of m . Within each statement, our approach applies these rules based on their nesting depth in the abstract syntax tree. For example, for the statements of the form $m_2(m_1(x))$, our approach first applies these rules on m_1 and then on m_2 .

Figures 5a and 5b show partial ATGs for C# (IndexFiles.cs) and Java (IndexFiles.java) code examples shown in Figure 2, respectively. Figure 5 also shows corresponding line numbers of each sub-graph. Our approach applies Rules 2 and 8 for Lines 6 and 9 (Figure 2) to build corresponding sub-graphs in the ATG. Our approach applies Rules 2, 3, and 6 to build corresponding sub-graphs for Lines 11 and 14 (Figure 2).

Algorithm 2: ATG Comparison Algorithm

Input: G is the ATG of a method (m); G' is the ATG of m 's mapped method.

Output: S is a set of mapping relations for API methods

```
begin
   $P \leftarrow \text{findVarPairs}(m, m')$ 
  for Pair  $p$  in  $P$  do
     $SM \leftarrow G.\text{nextMethods}(p.\text{sharp})$ 
     $JM \leftarrow G.\text{nextMethods}(p.\text{java})$ 
     $\Delta S = \text{mapping}(SM, JM)$ 
    while  $\Delta S \neq \phi \mid \Delta SM \neq \phi \mid \Delta JM \neq \phi$  do
       $S.\text{addAll}(\Delta S)$ 
      for Method  $sm$  in  $SM$  do
        if  $sm.\text{isMapped}$  then
           $SM.\text{replace}(sm, sm.\text{nextMethod}())$ 
        else
           $SM.\text{replace}(sm, sm.\text{mergeNextMethod}())$ 
      for Method  $jm$  in  $JM$  do
        if  $jm.\text{isMapped}$  then
           $JM.\text{replace}(jm, jm.\text{nextMethod}())$ 
        else
           $JM.\text{replace}(jm, jm.\text{mergeNextMethod}())$ 
       $\Delta S = \text{mapping}(SM, JM)$ 
  end
```

4.3.3 Comparing API transformation graphs

The second sub-step compares each pair of built ATGs for mining mapping relations of API methods. Our mapped API methods should satisfy three criteria. (1) The mapped API methods implement the same functionality. (2) The mapping relation describes the relation between parameters and receivers of mapped API methods. (3) The mapping relation describes the relation between return of mapped API methods. The two mapped API methods in two different languages satisfying the preceding three criteria are replaceable in the client code. Therefore, these mapped API methods assist for translating client code from one language to another.

(@Hao: Could you please change the image to Algorithm 1 from Algorithm 2) Algorithm 1 presents major steps of comparing ATGs for mining mapping relations of API methods. Consider two methods m and m' of two different languages L and L' , respectively, in the client code. Consider that the associated ATGs of m and m' are compared to mine mapping relations of API methods. First, the `findVarPairs` function finds mapped variables $\in F, V$, and P_1 in m and m' . `findVarPairs` maps two variables v and v' of methods m and m' , respectively, if the similarity metric value on their names is greater than `SIM_THRESHOLD`. For constants in m and m' , `findVarPairs` maps those two constants, if they have exactly the same value. Our algorithm uses these variable and constant mappings to compute mappings between API methods that use these variables and constants.

Mapping inputs: The first criterion is based on the inputs of jm and sm . We map jm with sm , if there is a one-to-one mapping between inputs of jm and sm . Here, our approach considers both the receiver and the parameters as the inputs of an API method.

Mapping functionalities: The second criterion is based on functionalities of jm and sm . We consider that jm and sm implement the same functionality, if the similarity metric value between the name of jm and the name of sm is greater than `SIM_THRESHOLD`.

Mapping outputs: The third criterion is based on the returns of jm and sm . Consider the returns of jm and sm as r_1 and r_2 , respectively. We map jm with sm , if the type of r_1 is mapped with the type of r_2 in mapping relations of API classes.

Our algorithm first attempts to map the first API method jm in

m with the first API method sm in m' . Our algorithm uses the `nextMethods` function to get these jm and sm API methods. If our algorithm successfully maps jm with sm , our algorithm moves to the next available API methods in m and m' of the client code. If our algorithm is not able to map jm with sm , our algorithm merges sm and jm with their next available API methods in the corresponding ATGs, respectively, and attempts to map merged API methods. Our algorithm uses the `replace` function for merging an API method with its next available API method. For two merged API methods, our algorithm uses the maximum similarity of method names between jm and sm as a similarity metric value for mapping their functionalities. With each iteration, sm or jm or the mapping relation (represented as S) in the algorithm may change. Therefore, we repeat our algorithm till S , sm , and jm do not change anymore.

We next explain our algorithm using the illustrative example shown in Figure 5. The numbers shown in circles represent the major steps in our algorithm for mining mapping relations of API methods. We next explain each step in detail.

(@Hao, Can you please change the numbers in Figure 5 to show more explicitly as steps such as Step 1 ...) **Step 1: mapping parameters, fields, local variables, and constants.** Given two ATGs of each method pair $\langle m, m' \rangle$, this step maps variables such as parameters, fields, and local variables by comparing their names, and maps constants by comparing their values. As shown in Figure 5, Step 1 maps two constants since both the constants have the same value “index”.

Step 2: mapping inputs of API methods. Step 2 mines mapping relations of API methods using variable and constant mapping relations. Initially, this step identifies the first API methods in the two ATGs and tries to map their receiver and parameters of the two API methods. In our current example, this step maps the parameter `filename` to the parameter `arg0` as these parameters are of the same type and their associated constants are mapped.

Step 3: mapping outputs of API methods. In contrast to Step 2 that maps parameters, Step 3 maps returns of API methods. In this step, if our approach is not able to map returns, our approach merges the next API method (Section 2) and then attempts to map returns of merged API methods. In our current example shown in Figure 5, the return of `System.IO.FileInfo.FileInfo()` cannot be mapped to the return of `java.io.File.File()`. Therefore, our approach merges subsequent API methods in the ATG till the `Exists` API method, since the returns (shown as `Boolean`) can be mapped only after the `Exists` API method. Figure 5 shows Step 3 along with the mapped returns.

Step 4: mapping functionalities. After our approach maps parameters and returns, this step further maps functionalities of those merged API methods. Given two merged API methods with mapped parameters and returns, this step uses the similarity metric value based on their method names as a criterion for mapping their functionalities. In the preceding example, this step maps the two merged API methods shown in Figure 5a to the merged API methods of the `java.io.File.exists()` as all three merged API methods include the method named `exists`.

Our approach applies the preceding steps on ATGs (as shown in Figures 5a and 5b). After finding out the mapped pair of API methods as shown in Figure 5, our approach merges all variables and outputs to corresponding parameters and receivers and produces the mapping relation of API method as shown in Figure 3.

5. EVALUATIONS

We implemented a tool named MAM based on our approach and conducted two evaluations using our tool to show the effectiveness

Project	Source	Java version		C# version	
		#C	#M	#C	#M
neodatis	SourceForge	1298	9040	464	3983
db4o	SourceForge	3047	17449	3051	15430
numerics4j	SourceForge	145	973	87	515
fpml	SourceForge	143	879	144	1103
PDFClown	SourceForge	297	2239	290	1393
OpenFSM	SourceForge	35	179	36	140
binaryNotes	SourceForge	178	1590	197	1047
lucene	Apache	1298	9040	464	3015
logging	Apache	196	1572	308	1474
hibernate	hibernate	3211	25798	856	2538
rasp	SourceForge	320	1819	557	1893
llrp	SourceForge	257	3833	222	978
simmetrics	SourceForge	107	581	63	325
aligner	SourceForge	41	232	18	50
fit	SourceForge	95	461	43	281
Total		11668	75685	6900	34165

Table 1: Subjects

of our approach. In our evaluations, we address the following two research questions.

1. How effectively can our approach mine various mapping relations of APIs (Section 5.1)?
2. How much benefit can the mined mapping relations of APIs offer in aiding language translation (Section 5.2)?

Table 1 shows 15 open source projects with both Java and C# versions used as subjects in our evaluations. Column “Project” lists names of subjects. Column “Source” lists sources of these subjects. These subjects are collected from famous open source societies such as SourceForge¹⁰, Apache¹¹, and hibernate¹². Columns “Java version” and “C# version” list versions in Java and C#, respectively. In these two columns, sub-columns “#C” and “#M” list number of classes and methods, respectively. As shown in the table, Java versions are much larger than C# versions for all subjects. We found two major factors for such a difference. First, Java versions of some of the projects are more up-to-date. For example, the latest Java version of *numerics4j* is 1.3, whereas the latest C# version is 1.2. Second, for some projects, translation from Java to C# is still in progress. For example, the website¹³ of *neodatis* states that *neodatis* is a project in Java and is being ported to C#. This observation further confirms the usefulness of our approach as our approach aids translating projects from one language to other languages. In total, all these projects include 18,568 classes and 109,850 methods.

We conducted all evaluations on a PC with Intel Qual CPU @ 2.83GHz and 1.98M memory running Windows XP. More details of our evaluation results are available at <https://sites.google.com/site/asergpr/projects/mam>.

5.1 Mining API mapping

To evaluate the first research question, we use the top 10 projects from Table 1 as subjects for mining API mapping.

Aligning client code. We first present the results of aligning client code. We use the *SIM_THRESHOLD* value as 0.6, which is based on our initial empirical experience. We choose a relatively low threshold as it helps our approach to take into account as much client code as possible.

Table 2 shows our evaluation results. In column “Aligned”, sub-columns “#C” and “#M” list number of aligned classes and aligned

Project	Java version		C# version		Aligned	
	%C	%M	%C	%M	#C	#M
db4o	87.8%	65.5%	87.6%	74.1%	2674	11433
fpml	93.7%	70.5%	93.5%	56.2%	134	620
PDFClown	86.5%	51.0%	88.6%	82.1%	257	1143
OpenFSM	97.1%	72.1%	94.4%	92.1%	34	129
binaryNotes	98.9%	61.1%	89.3%	92.7%	176	971
neodatis	44.7%	54.8%	100.0%	93.6%	408	3728
numerics4j	57.2%	48.6%	95.4%	89.9%	75	174
lucene	34.9%	26.6%	97.6%	79.8%	453	2406
logging	91.8%	18.1%	58.4%	19.3%	180	285
hibernate	26.4%	1.2%	99.1%	12.6%	848	319
Average	53.2%	30.8%	88.8%	69.2%	524	2121

Table 2: Results of Aligning client code

methods. For each project of Column “C# version” and Column “Java version”, sub-column “%C” lists the percentage of the aligned classes among total classes of corresponding versions. Sub-column “%M” lists the percentage of the aligned methods among total methods of corresponding versions. Row “Total” of the two sub-columns lists the percentage of aligned methods/classes among the total methods/classes as shown in Table 1. We find that the results of Table 2 fall into three categories. This first category includes *db4o*, *fpml*, *PDFClown*, *OpenFSM*, and *binaryNotes*. In this category, our approach achieves relatively high percentages for both Java and C# versions. For all these five projects, “%M” is relatively smaller than “%C” because methods of those unaligned classes cannot be aligned and hence are counted as unaligned¹⁴. The second category includes *neodatis*, *numerics4j*, and *lucene*. In this category, our approach aligns C# versions well but does not align Java versions so well. We find that the translation of *neodatis* and *lucene* from Java to C# is in progress, and the Java version of *numerics4j* is more up to date than its C# version. As a result, some Java classes or methods do not have corresponding implementations in C# versions in these projects and hence are not mapped. The third category includes *logging* and *hibernate*. In this category, our approach does not align classes and methods of the two projects well. Although both of the two projects seem to be translated from existing Java versions, the programmers of the two projects often do not refer to names of existing Java versions for naming entities. For these two projects, the percentage of aligned classes is relatively high, and the percentage of aligned methods is relatively low. We find that even if our approach aligns a wrong class pair, our approach does not align methods within the wrong pair as the method names of a wrong pair are quite different. These results suggest that we can take method names into account when aligning classes in future work. For all these projects, our approach does not align all classes and all methods. We discuss these issues in Section 6.

In summary, as shown by Row “Average”, our approach aligns most classes and methods on average. The result confirms that many programmers refer to existing versions of another language to name entities of a version under development.

Mining API mapping. Table 3 shows the results of mined mapping relations of API classes and methods. Columns “Class” and “Method” list results of mining API classes and API methods, respectively. Sub-column “Num.” lists number of mined mapping relations. The number of mined API mapping relations are highly proportional to the sizes of projects as shown in Table 1, except for *logging* and *hibernate*. As classes and methods of these two projects are not well aligned, our approach does not mine many mapping relations of APIs from these two projects. For the remain-

¹⁰<http://www.sf.net>

¹¹<http://www.apache.org/>

¹²<http://www.hibernate.org/>

¹³<http://wiki.neodatis.org/>

¹⁴Another factor lies in that Java versions usually have many getters and setters and these getters and setters often do not have corresponding methods in C# versions.

Project	Class			Method		
	Num.	Acc.	J2SE	Num.	Acc.	J2SE
db4o	3155	83.3%	117	10787	90.0%	297
fpml	199	83.3%	41	508	83.3%	216
PDFCrown	539	96.7%	36	514	100.0%	111
OpenFSM	64	86.7%	16	139	73.3%	12
binaryNotes	287	90.0%	31	671	90.0%	55
neodatis	526	96.7%	41	3517	100.0%	539
numerics4j	97	83.3%	2	429	83.3%	29
lucene	718	90.0%	83	2725	90.0%	522
logging	305	73.3%	45	56	90.0%	19
hibernate	1126	66.7%	87	7	13.3%	5
Total	6695	86.7%	344	19110	90.0%	1768

Table 3: Results of mining API mapping

ing projects, our approach mines many mapping relations of API classes and API methods. Sub-column “Acc.” lists accuracies of the top 30 mined API mapping (*i.e.*, percentages of correct mapping relations). For mined API mapping from each project, we manually inspect top 30 mined mapping relations of APIs and classify them as correct or incorrect based on programming experiences. We find that our approach achieves high accuracies, except for *hibernate*. Although our approach does not align *logging* quite well either, the accuracies of API mapping from *logging* are still relatively high. To mine API mapping of classes, our approach requires that names of classes, methods, and variables are similar. To mine API mapping of methods, our approach requires that two built ATGs are similar. These two requirements are relatively strict. As a result, if the first step does not align client code well, our approach may miss some mapping relations of APIs but does not introduce many false mapping relations. In other words, our approach is robust to mine accurate API mapping. Sub-column “J2SE” lists mined API mapping between J2SE and .NET. We next compare these mapping relations with manually built mapping relations.

Row “Total” lists the total result after we merge all duplicated mapping relations. In summary, our approach mines a large number of mapping relations of APIs. These mined mapping relations are accurate and cover various libraries.

Comparing with manually built API mapping. Some translation tools such as Java2CSharp¹⁵ include manually written mapping relations of APIs. For example, one item from the mapping files of Java2CSharp is as follows:

```
package java.math :: System {
    class java.math.BigDecimal :: System.Decimal {
        method multiply(BigDecimal)
        { pattern = Decimal.Multiply(@0, @1); }
    }
}
```

This item describes mapping relations between the `multiply` method of Java and the `Multiply` method of C#. The pattern string describes mapping relations of inputs. In particular, “@0” denotes the receiver of the `multiply` method, and “@1” denotes the first parameter of the `multiply` method. Based on this item, Java2CSharp translates the following code snippet from Java to C# as follows:

```
BigDecimal m = new BigDecimal(1);
BigDecimal n = new BigDecimal(2);
BigDecimal result = m.multiply(n);
->
Decimal m = new Decimal(1);
Decimal n = new Decimal(2);
Decimal result = Decimal.Multiply(m, n);
```

To compare with manually built mapping files of Java2CSharp, we translate our mined API mapping relations using the following strategy. First, for each Java class, we translate its mapping relations of classes with the highest support values into mapping files

¹⁵<http://j2cstranslator.wiki.sourceforge.net>

Package	Class			Method		
	P	R	F	P	R	F
java.io	78.6%	73.3%	76.0%	93.1%	66.3%	79.7%
java.lang	82.6%	86.4%	84.5%	93.8%	81.5%	87.6%
java.math	50.0%	50.0%	50.0%	66.7%	66.7%	66.7%
java.net	100.0%	50.0%	75.0%	100.0%	50.0%	75.0%
java.sql	100.0%	66.7%	83.3%	100.0%	66.7%	83.3%
java.text	50.0%	50.0%	50.0%	50.0%	50.0%	50.0%
java.util	56.0%	87.5%	71.8%	65.8%	67.6%	66.7%
junit	100.0%	50.0%	75.0%	92.3%	88.9%	90.6%
orw.w3c	42.9%	75.0%	58.9%	41.2%	77.8%	59.5%
Total	68.8%	77.9%	73.4%	84.6%	73.9%	79.3%

Table 4: Results of comparing results

as relations of packages and classes. Second, for each Java method, we translate its mapping relations of methods with the highest support values into mapping files as relations of methods with pattern strings. For one-to-one mapping relations of methods, this step is automatic as mined mapping relations describe mapping relations of corresponding methods and inputs. For many-to-many mapping relations of methods, this step is manual as mined mapping relations do not include adequate details such as how to deal with multiple outputs. We further discuss this issue in Section 6.

The mapping files of Java2CSharp cover 13 packages defined by J2SE and 2 packages defined by JUnit¹⁶, and we treat these mapping files as a golden standard. We find 9 packages overlapping between the mined mapping files and the mapping files of Java2CSharp. We compare mapping relations of APIs within these mapping packages, and Table 4 shows the results. Columns “Class” and “Method” list results of comparing API classes and methods, respectively. Sub-columns “P”, “R”, and “F” denote precision, recall, and F-score. *Precision*, *Recall*, and *F-score* are defined as follows:

$$Precision = \frac{true\ positives}{true\ positives + false\ positives} \quad (1)$$

$$Recall = \frac{true\ positives}{true\ positives + false\ negatives} \quad (2)$$

$$F-score = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (3)$$

In these preceding formula, true positives represent those mapping relations that exist in both the mined API mapping and the golden standard¹⁷; false positives represent those transitions that exist in the mined API mapping but not in the golden standard; false negatives represent those mapping relations that exist in the golden standard but not in the mined API mapping. Row “Total” show the results when we compare mapping relations of all the packages listed in Table 4.

From sub-columns “P” of Table 4, we find that our approach achieves relatively high precisions, but the precisions are lower than the accuracies shown in Table 3. After inspecting those differences, and we find 25 new correct mapping relations of API classes from our mined mapping files. For example, these mined mapping files contain a mapping relation between `org.w3c.dom.Attr` and `System.Xml.XmlAttribute`, and the mapping relation does not exist in the mapping files of Java2CSharp. As these mapping relations are considered as false positives, these relations reduced our precisions. However, these mapping relations are useful and complements the existing mapping files of the Java2CSharp tool.

¹⁶<http://www.junit.org/>

¹⁷We ignore those mapping relations of APIs that do not have call sites in the projects listed in Table 3, since our approach cannot mine mapping relations for these APIs from the given projects.

Projects	No MF		MF		Ext. MF			
	<i>E</i>	<i>D</i>	<i>E</i>	<i>D</i>	<i>E</i>	% <i>E</i>	<i>D</i>	% <i>D</i>
rasp	973	159	708	123	627	11.4%	93	24.4%
llrp	2328	122	1540	114	269	82.5%	42	63.2%
simmetrics	217	13	12	0	6	50.0%	0	0%
aligner	368	34	289	0	262	9.3%	0	0%
fit	177	29	27	0	20	25.9%	0	0%
Total	4063	491	2576	237	1174	54.4%	135	43.0%

Table 5: Compilation errors and defects

From sub-columns “*R*” of Table 4, we find that our approach achieves relatively high recalls, but the recalls still have spaces for improvements. For example, our approach does not mine the mapping relation between `java.util.ResourceBundle` and `System.Resources.ResourceManager` as described in the mapping files of Java2CSharp. Although it exists in *hibernate*, our approach fails to mine the mapping relation since *hibernate* uses the two API classes in two classes with quite different names. Our approach also fails to mine the mapping relations between `java.util.getTime()` and `System.DateTime.Ticks` correctly as our approach ignores the operation of dividing 1,000 from ticks to milliseconds. We further discuss these issues in Section 6.

In summary, compared with the mapping files of Java2CSharp, our mined mapping files show a reasonable high precisions and recalls. The precisions are relatively high since our mined mapping relations are accurate and include new mapping relations that are not covered by Java2CSharp. The recalls are also relatively high as our approach mines many mapping relations although we still have spaces for further improvements.

5.2 Aiding Language Translation

To evaluate the second research question, we feed the mined API mapping relations to the Java2CSharp tool and investigate whether these relations can improve the tool’s effectiveness. We choose this tool because this tool is a relatively mature project at ILOG¹⁸ (now part of IBM) and supports the extension of user-defined mapping relations of APIs.

We use Java2CSharp to translate five projects listed in Table 1 from Java to C#, and Table 5 shows the results. For each translated C# project, Column “No MF” lists results without mapping files. Column “MF” lists results with the mapping files of Java2CSharp. Column “Ext. MF” lists results with mapping files that combine our mined API mapping relations with the existing mapping files of Java2CSharp. Sub-columns “*E*” and “*D*” list number of compilation errors and found defects, respectively. For each project, we find out those overlapping files between translated files and existing C# files and manually inspect the top 5 largest files by comparing existing C# files for API related defects. Sub-columns “%*E*” and “%*D*” list percentages of improvements over the results of “MF”. On average, mined API mapping relations help further reduce 54.4% compilation errors and 43.0% found defects. As the five projects use different libraries, the number of translated projects are different. In particular, *simmetrics* and *fit* use API classes of J2SE that are covered by mapping files. Consequently, the translated projects of *simmetrics* and *fit* have a few errors and defects. The *aligner* project also mainly uses J2SE, but it uses many API classes and methods from `java.awt` for its GUI. The mapping files of Java2CSharp do not cover any classes of `java.awt`, so the translated project has many errors. As the existing C# version of *aligner* does not have GUI, we do not compare those buggy translated GUI files and we do not find any defects. The mined files map `java.awt` to `System.Windows.Forms` and thus reduce compilation errors. However, the result is not signifi-

¹⁸<http://www.ilog.com/>

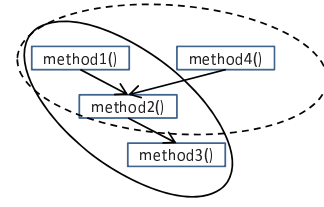


Figure 6: Merging technique

cant as many classes of the two packages are still not mapped. For *rasp* and *llrp*, they both use various libraries besides J2SE. Consequently, the translated projects have both many errors and defects. In particular, *llrp* uses `log4j`¹⁹ and `jdom`²⁰, and the mined mapping files contain mapping relations of the two libraries. As a result, the mined API mapping helps reduce compilation errors and defects significantly. For *rasp*, it uses some libraries such as `Neethi`²¹ and `WSS4J`²². Since the used subjects for mining and thus our mined API mapping do not cover the two libraries, the translated project of *rasp* contains many complex errors and defects.

In summary, the mined API mapping improves existing language translation tools such as Java2CSharp. In particular, the mined API mapping helps effectively reduce compilation errors and defects in the translated projects.

5.3 Threats to Validity

The threats to external validity includes the representativeness of the subjects in true practice and the existing translation tool being compared. Although we applied our approach on 10 projects for mining API mapping and on additional 5 projects for helping language translation, our approach is evaluated only on these limited projects. Although Java2CSharp is the best translation tool as far as we know, other tools may perform better than the tool. The threat could be reduced by more evaluations on more subjects and more translation tools in future work. The threats to internal validity includes human factors for determining correct mined API mapping and for determining defects in translated projects. To reduce the threats, we inspected mined mapping relations and defects carefully, and we referred to existing mapping relations and existing C# files for determining correct API mapping and defects, respectively. The former threat could be further reduced by find more existing mapping relations of APIs as we did to J2SE. The latter threat could be reduced by running test cases when we improve our approach and produce translated projects without complication errors.

6. DISCUSSION AND FUTURE WORK

We next discuss issues in our approach and describe how we address these issues in our future work.

Aligning client code. Table 2 shows that our approach could not align client code in a few cases. The primary reason is that the functionality associated with a class or a method in one language version is split among multiple classes or methods in the other language version. To address this issue, we plan to align classes and methods of client code based on their functionalities through developing or adapting dynamic approaches such as Jiang and Su [5] in future work.

Mining richer API mapping. Table 4 shows that our approach does not achieve high recall for J2SE. Although we use ten large

¹⁹<http://logging.apache.org/log4j/>

²⁰<http://www.jdom.org/>

²¹<http://ws.apache.org/commons/neethi/>

²²<http://ws.apache.org/wss4j/>

projects as subjects, these projects still do not provide sufficient code examples for mining mapping relations of all APIs in J2SE. Our previous work [11, 12] show that it is feasible to use large-scale repositories available on the web as subjects with the help of code search engines such as Google code search²³. In future work, we plan to leverage these code search engines to mine richer API mapping.

Mining many-to-many mapping relations of API methods. Among our mined mapping relations of API methods, many relations are one-to-one relations. The reason is that our Algorithm 2 uses only forward analysis for merging API methods. We explain this issue using an illustrative example shown in Figure 6. After merging API methods `method1` and `method2`, if our algorithm still not able to map the merged API methods with an API method in the other language, our algorithm attempts to merge `method3` rather than `method4`, which could be a possible candidate for mining the mapping relation. In future work, we plan to incorporate backward analysis to enhance our existing algorithm. We expect that our enhanced algorithm can mine more many-to-many mapping relations.

Migrating many-to-many mapping relations of API methods. A mined many-to-many mapping relation of API methods can have multiple outputs and complex internal data processes. Although, our ATGs help identify all API methods, our implementation of ATGs is not complete for supporting an automatic translation. For example, we need to manually add an *or* operator for the two outputs of the API mapping shown in Figure 3. In future work, we plan to enhance our implementation to help automate migration with many-to-many mapping relations.

Migrating unmapped APIs. Our approach mines API mapping of methods along with the mappings of their inputs and outputs. These mappings are useful for translating API methods of one language to another. Sometimes, our approach may not be able to map inputs and outputs of mapped API methods. If our approach is not able to map outputs, our approach simply ignores those outputs that are not used in the client code. However, as inputs cannot be ignored, the translated code results in compilation errors. In future work, we plan to address this issue by analyzing how two versions of a project deal with a similar unmapped API problem for some other code examples.

7. RELATED WORK

Our approach is related to previous work on language migration and library migration.

Language migration. To reduce effort of language migration [9], researchers propose various approaches to automate the process [3, 7, 13, 14, 16]. Most of these approaches focus the syntax differences between languages. For example, Deursen *et al.* [13] propose an approach to identify objects in legacy code, and the results are useful to deal with differences between object-oriented and procedural languages. As shown by El-Ramly *et al.* [2]’s experience report, existing approaches and tools support only a subset of APIs, and consequently it becomes an important and yet challenging task to automate API transformation. Our approach mines API mapping between languages to aid language migration, addressing a significant problem unaddressed by the previous approaches and complementing these approaches.

Library migration. With evolution of libraries, some APIs may become incompatible across library versions. To deal with the problem, some approaches have been proposed. In particular, Henkel and Diwan [4] propose an approach that captures and replays API refactoring actions to keep client code updated. Xing and Strou-

lia [15] propose an approach that recognizes the changes of APIs by comparing the differences of two versions of libraries. Balaban *et al.* [1] propose an approach to help translate client code when mapping relations of libraries are available. Different from these approaches, our approach focuses on mapping relations of APIs among different languages. In addition, as our approach uses API transformation graphs to mine mapping relations of APIs, our approach helps mine mapping relations for those API methods whose input orders are changed or whose functionalities are split into several methods if our approach is applied in library migration.

8. CONCLUSION

Mapping relations of APIs are quite useful to language migration but are difficult to obtain due to various factors. In this paper, we propose a novel approach to mine mapping relations of APIs from existing different versions of a project automatically. We conducted two evaluations on our approach. The results show that our approach mines various API mapping between Java and C#, and API mapping improves existing language translators such as Java2CSharp.

9. REFERENCES

- [1] I. Balaban, F. Tip, and R. Fuhrer. Refactoring support for class library migration. In *Proc. 20th OOPSLA*, pages 265–279, 2005.
- [2] M. El-Ramly, R. Eltayeb, and H. Alla. An experiment in automatic conversion of legacy Java programs to C#. In *Proc. AICCSA*, pages 1037–1045, 2006.
- [3] J. Hainaut, A. Cleve, J. Henrard, and J. Hick. *Software Evolution*. Springer, 2008.
- [4] J. Henkel and A. Diwan. CatchUp!: capturing and replaying refactorings to support API evolution. In *Proc. 27th ICSE*, pages 274–283, 2005.
- [5] L. Jiang and Z. Su. Automatic mining of functionally equivalent code fragments via random testing. In *Proc. 18th ISSA*, pages 81–92, 2009.
- [6] T. Jones. *Estimating software costs*. McGraw-Hill, Inc. Hightstown, NJ, USA, 1998.
- [7] M. Mossienko. Automated COBOL to Java recycling. In *Proc. 7th CSMR*, pages 40–50, 2003.
- [8] D. Orenstein. QuickStudy: Application Programming Interface (API). *Computerworld*, 10, 2000.
- [9] H. Samet. Experience with software conversion. *Software: Practice and Experience*, 11(10), 1981.
- [10] A. Terekhov and C. Verhoef. The realities of language conversions. *IEEE Software*, pages 111–124, 2000.
- [11] S. Thummalapenta and T. Xie. PARSEWeb: A programmer assistant for reusing open source code on the web. In *Proc. 22nd ASE*, pages 204–213, November 2007.
- [12] S. Thummalapenta and T. Xie. SpotWeb: Detecting framework hotspots and coldspots via mining open source code on the web. In *Proc. 23rd ASE*, pages 109–112, 2008.
- [13] A. Van Deursen, T. Kuipers, and A. CWI. Identifying objects using cluster and concept analysis. In *Proc. 21st ICSE*, pages 246–255, 1999.
- [14] R. Waters. Program translation via abstraction and reimplement. *IEEE Transactions on Software Engineering*, 14(8):1207–1228, 1988.
- [15] Z. Xing and E. Stroulia. API-evolution support with Diff-CatchUp. *IEEE Transactions on Software Engineering*, 33(12):818–836, 2007.

²³<http://www.google.com/codesearch>

- [16] K. Yasumatsu and N. Doi. SPiCE: a system for translating Smalltalk programs into a C environment. *IEEE Transactions on Software Engineering*, 21(11):902–912, 1995.