

# Repository-Level Compositional Code Translation and Validation

ALI REZA IBRAHIMZADA\*, University of Illinois Urbana-Champaign, USA

KAIYAO KE, University of Illinois Urbana-Champaign, USA

MRIGANK PAWAGI, Indian Institute of Science, India

MUHAMMAD SALMAN ABID, Cornell University, USA

RANGEET PAN, IBM Research, USA

SAURABH SINHA, IBM Research, USA

REYHANEH JABBARVAND, University of Illinois Urbana-Champaign, USA

Code translation transforms programs from one programming language (PL) to another. One prominent use case is application modernization to enhance maintainability and reliability. Several rule-based transpilers have been designed to automate code translation between different pairs of PLs. However, the rules can become obsolete as the PLs evolve and cannot generalize to other PLs. Recent studies have explored the automation of code translation using Large Language Models (LLMs). One key observation is that such techniques may work well for crafted benchmarks but fail to generalize to the scale and complexity of real-world projects with inter- and intra-class dependencies, custom types, PL-specific features, etc. We propose ALPHATrans, a neuro-symbolic approach to automate *repository-level* code translation. ALPHATrans translates both source and test code, and employs multiple levels of validation to ensure the translation *preserves* the functionality of the source program. To break down the problem for LLMs, ALPHATrans leverages program analysis to decompose the program into fragments and translates them in the *reverse call order*.

We leveraged ALPHATrans to translate *ten* real-world open-source projects consisting of (836, 8575, 2719) classes, methods, and tests. ALPHATrans translated the entire repository of these projects consisting of 6899 source code fragments. 99.1% of the translated code fragments are syntactically correct, and ALPHATrans validates the translations' runtime behavior and functional correctness for 25.8%. On average, the integrated translation and validation take 36 hours (min=4, max=122) to translate a project, showing its scalability in practice. For the syntactically or semantically incorrect translations, ALPHATrans generates a report including existing translation, stack trace, test errors, or assertion failures. We provided these artifacts to two developers to fix the translation bugs in four projects. They were able to fix the issues in 20.1 hours on average (5.5 hours for the smallest and 34 hours for the largest project) and achieve all passing tests. Without ALPHATrans, translating and validating such big projects could take weeks, if not months.

Additional Key Words and Phrases: Automated Code Translation, Code Translation and Validation, Program Analysis, Large Language Models

## 1 Introduction

Application modernization offers numerous benefits to developers, including better performance, maintainability, productivity, reliability, and security [26, 27, 29, 30]. Manual migration or modernization of real-world projects can be time-consuming and error-prone. Code translation can help automatically convert programs from one programming language (PL) to another.

Transpilers solely rely on program analysis and perform rule-based translation, failing to translate code between languages that greatly differ in syntax or semantics [3]. This also makes them very

---

\*Author was an intern at IBM Research at the time of this work.

---

Authors' Contact Information: Ali Reza Ibrahimzada, alirezai@illinois.edu, University of Illinois Urbana-Champaign, Urbana, Illinois, USA; Kaiyao Ke, kaiyaok2@illinois.edu, University of Illinois Urbana-Champaign, Urbana, Illinois, USA; Mrigank Pawagi, mrigankp@iisc.ac.in, Indian Institute of Science, Bengaluru, Karnataka, India; Muhammad Salman Abid, ma2422@cornell.edu, Cornell University, Ithaca, NY, USA; Rangeet Pan, rangeet.pan@ibm.com, IBM Research, Yorktown Heights, NY, USA; Saurabh Sinha, sinhas@us.ibm.com, IBM Research, Yorktown Heights, NY, USA; Reyhaneh Jabbarvand, reyhaneh@illinois.edu, University of Illinois Urbana-Champaign, Urbana, Illinois, USA.

PL-specific; they cannot generalize to newer features of the same PL pairs easily, let alone other PLs. Finally, the translations lack readability, requiring much effort to understand and validate them, and naturalness, failing to create idiomatic code in the target PL [42]. State-of-the-art code translation techniques attempt to harvest the emerging abilities of the Large Language Model (LLM) in code synthesis to overcome the limitations of transpilers [39, 42, 60]. However, these techniques are still limited to translating simple programs in crafted benchmarks or selected slices of real-world projects due to the following challenges:

- (1) *Problem complexity*. The source and target PLs can be fundamentally different in programming paradigms, typing, and memory management. Some PLs have specific properties that may not exist in others, e.g., constructor overloading in Java. Such complexities are beyond the abilities of existing LLMs to handle, causing them to hallucinate when translating types, code constructs, or even method names [42], making translations non-compilable or useless.
- (2) *Validation*. The translation should preserve the functionality of the source project. Most existing techniques follow a “translation first and validation next” approach, which can postpone the validation and not benefit from the potential use of validation as feedback to correct the translation [42]. A few techniques use formal methods [39, 60] to verify translations on the go. However, these techniques cannot scale to real-world projects. One possible solution for validation is reusing the tests in the source language. However, due to (1) multiple invocations of different methods in unit tests and (2) inherent long call chains in real-world projects, testing a translated method *in isolation* is impossible.
- (3) *Limited context window*. Concerning repository-level translation, the entire project and, in many cases, even the entire class cannot fit into the context window of current LLMs [23]. Even assuming an unlimited context window, LLMs are shown to have a short attention span [35], preventing them from properly capturing the intra- and inter-procedural dependencies in real-world projects.

This paper presents ALPHATRANS, a neuro-symbolic<sup>1</sup> approach for automated repository-level code translation and validation. ALPHATRANS leverages static analysis to resolve PL-specific features of the source language (§4.1), decompose the source project into smaller fragments (§4.2), and create a compilable project skeleton in the target language (§5). It then starts translating fragments in the reverse call order and validates them using existing tests when possible (§6). After translating each fragment, ALPHATRANS updates the project skeleton and ensures the whole project compiles, gradually translating and validating the source project into the target PL. To improve translation quality, static analysis again comes to the rescue: ALPHATRANS collects relevant context to each fragment, including translated callee methods and surrounding contexts, e.g., class declaration, global variables/fields, etc. It also uses relevant in-context examples based on the specific properties of the fragment to be translated. The current version of ALPHATRANS implements two levels of dynamic validation: (1) running the source tests on the translated fragment in isolation using language interoperability (§6.1) and (2) decomposing, translating, and executing the unit tests on the translated project (§6.3). Finally, ALPHATRANS recomposes the translated fragments to create the program in the target PL (§6.2).

The idea of compositional translation and validation proposed by us is PL-agnostic; however, implementing the program transformation component is PL-specific. For the first version of ALPHATRANS, the implementation supports translating from Java to Python. Our motivations for choosing this PL pair are: (1) Java offers many features that are not supported or common in other PLs by default (e.g., method/constructor overloading, complex types, circular dependencies, local

<sup>1</sup>The keyword symbolic here refers to a general term of symbolic learning in contrast to machine learning and should not be confused with symbolic execution. We refer to combining LLMs and program analysis as a neuro-symbolic approach.

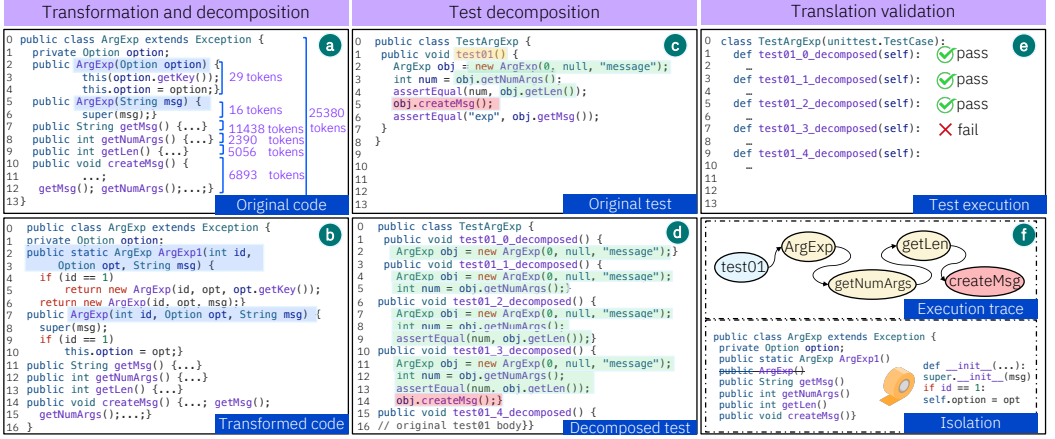


Fig. 1. Illustration of key challenges in repository-level code translation and ALPHATrans addressing them.

or anonymous inner classes, interfaces, etc.); (2) Python programs are not compiled but interpreted, which makes many translation issues that can be caught at the compile time stay there until test execution and challenge the validation; and (3) both PLs are popular (top-5 on the TIOBE index [54]).

Using ALPHATrans to translate ten real-world Java projects to Python corroborates its *effectiveness* (translating 6,899 code fragments, with 99.1% syntactically correct translations and 25.8% validated translations), *scalability* (completing translations in 36 hours, on average), and *practicality* (human subjects were able to fix the issues in ALPHATrans translations and achieve a green test suite within 20.1 hours, on average). It is worth noting that these results were achieved using an open-access LLM (DeepSeek-Coder-33b-Instruct [22]) with a moderate size, and using bigger/stronger models such as GPT-4 [40] or Claude-3 [2] will improve the results.

To the best of our knowledge, ALPHATrans is the *first technique to translate an entire repository*, including tests, and generates validated translations (considering existing tests). Compared to the only repository-level translation attempt using GPT-4 [42] (translating Apache Commons CLI from Java to Python) that resulted in non-compilable code, let alone the translation being validated, ALPHATrans translated the entire repository, producing 99% syntactically correct and 70.7% functionally validated code. The effort by human subjects in fixing the translation bugs by ALPHATrans and achieving green tests creates pragmatic bug datasets for testing, bug localization, and program repair research. We want to show the artifacts we generate can be used by other people. Our code and artifacts are publicly available [31].

## 2 Challenges in Repository-Level Code Translation

To illustrate the most notable challenges in repository-level code translation and validation, we use the hypothetical example in Figure 1, inspired by the complexities in real-world Java projects.

**Challenge 1: Class Size.** The class consists of 25,380 tokens (a). Instructions for translating the code, in-context examples, and the model’s response can also significantly increase the number of input tokens. While some commercial LLMs support tens of thousands of tokens, many open-access LLMs do not. For example, DeepSeek-Coder-33b-Instruct [22] used in this paper has a context window of 16,384 tokens, of which only 4,096 tokens can be used for generation. To address this challenge, ALPHATrans decomposes Java application classes into smaller *field* or *method fragments* and translates each separately in a reverse call order (§4.2.1, §4.2.2).

**Challenge 2: PL-specific Properties.** Java programs frequently use method and constructor overloading, which are not supported by default in Python (a). This example shows instances of

constructor overloading (lines 2 and 5). In Python, declaring two constructors is allowed, however at runtime, the last declaration always overrides all previous constructors, which can result in unexpected behavior. To address this issue, ALPHATrans employs program analysis to refactor the original code while preserving the functionality (through test execution). The transformation includes changing the constructor’s name, updating the references, and, in many cases, changing the constructor’s implementation. The transformed code (b) makes the source program amenable to translation to Python.

**Challenge 3: Validation.** To illustrate the challenges with validation, consider test01 (c) that invokes *four* methods in its body (ArgExp, getNumArgs, getLen, and createMsg) to test the functionality of method getMsg in the assert statement. Suppose we can successfully translate all methods except createMsg. If we choose test translation, the most natural way of validating code translation, the execution of the translated test results in a runtime error when invoking createMsg. As a result, a translation issue in one method casts a shadow in validating the translation of the other methods. We refer to this issue as the *test translation coupling effect*. To overcome this challenge, ALPHATrans executes source language tests as-is (i.e., without translation) by leveraging a language-interoperability framework called GraalVM [41] (f). In this setting, a test in the source language is executed every time one of its invoked application methods (method fragments) is translated. This approach validates *functional equivalence* of each method *in isolation* since other methods invoked in the test or the body of the translated method can remain in the source language.

**Challenge 4: Test Translation.** GraalVM has certain limitations (§6.1), which prevents ALPHATrans from validating all the code fragments in isolation. Furthermore, we need to translate tests regardless of whether they are used for validation to maintain the translated projects in the target language. Test errors as a result of *test translation coupling effect* under-approximates the quality of translation: failing to validate the translation of four methods because of one incorrect translation. Root causing the translation bugs also requires additional efforts from developers, i.e., looking at the stack trace and coverage. To overcome this challenge, ALPHATrans decomposes the original test suite into *test fragments* (d). Executing the translated decomposed test suite results in three test passes (e), validating the *runtime behavior* of three methods that the original test suite could not promptly provide.

An alternative approach is parsing the stack trace and code coverage results for each runtime error during translation. However, test decomposition is a cleaner way to see the results per test execution promptly. It is also done once before translation. In translation to interpreted languages such as Python, specifically, the execution of test fragments can validate the runtime behavior of methods before waiting for functional validation. For fragments that GraalVM cannot validate, if ALPHATrans can successfully translate all the methods invoked during test execution and test passes, such test will also be used for validating *functional correctness*.

### 3 Overview of Approach

ALPHATrans consists of three main phases, shown in Figure 2: program transformation and decomposition (§4), type translation and skeleton construction (§5), and compositional translation and validation (§6). The first two phases aim to decompose and simplify the repository-level code translation problem for LLMs, helping the third phase yield high-quality validated translations.

The program transformation and decomposition phase first refactors the PL-specific properties of the source program into programming paradigms common among many PLs (§4.1). Next, it decomposes the source project into smaller units, i.e., *fragments*, and stores fragment dependencies in a data structure called *schema* (§4.2).

The type translation and skeleton construction phase takes the schema as input and produces *target project skeleton*, i.e., a compilable project in the target language with method signatures but

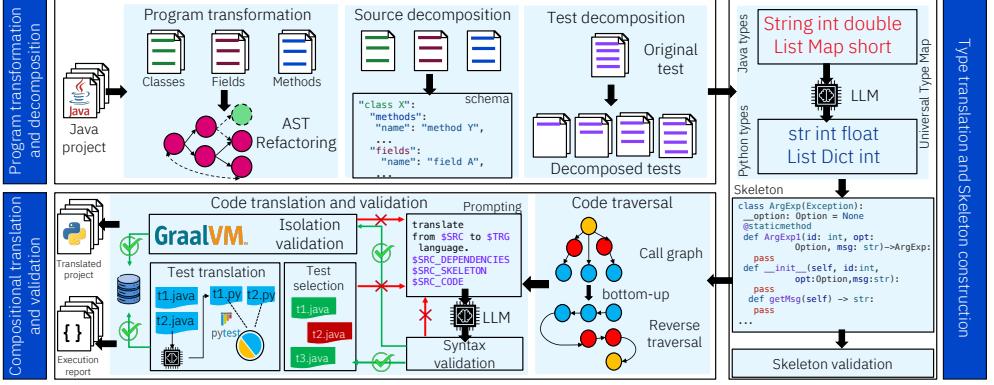


Fig. 2. Overview of ALPHATRANS.

no method implementation (§5.2). The first translation step is also performed by this phase, where it translates the source PL types to the target PL to ensure that class skeletons are compilable (§5.1). The outcome of type translation is a type mapping from the source to the target PL, which ALPHATRANS can reuse in translating other projects.

The compositional translation and validation phase takes schema and project skeleton as inputs and translates fragments, in the reverse call order, from the source to the target PL by prompting an LLM. After translating a fragment, this phase updates the class skeleton with its translation and checks whether the skeleton is still compilable. For method fragments, ALPHATRANS looks for corresponding tests and, if any exist, validates them. The first level of validation is performed through GraalVM’s language-interoperability capability to isolate the validation of the method using tests in the source language. In the second level, ALPHATRANS translates and executes the corresponding tests. In case of compilation errors or test failures, ALPHATRANS re-prompts the LLM with feedback (from the compilation and runtime errors) to improve the translation. If no improvement is achieved within a certain budget, ALPHATRANS continues to the next fragment until all are translated. For methods whose translations are not compilable or result in test errors/failures, ALPHATRANS generates reports consisting of existing translations and relevant artifacts, such as stack traces, test errors, assertion failures, and test coverage information.

## 4 Program Transformation and Decomposition

### 4.1 Program Transformation

This component performs semantics-preserving refactoring of method and constructor overloading in Java code to make it amenable to translation to Python. Other Java-specific features, namely, circular dependencies, inner classes, interfaces, and abstract classes, will be handled when constructing the project skeleton in Python (§5.2). The reason for resolving method and constructor overloading in the source language is that we have to change the implementation, i.e., call sites to methods and constructors. Therefore, such changes should be validated using source tests.

#### Algorithm 1: Constructor Overloading

**Inputs:** Overloaded Constructors  $OCs$   
**Output:** Code without Overloaded Constructors  $NOCs$

```

1 if !hasThisCall( $OCs$ ) then
2    $ids \leftarrow createConstructorIDS(OCs);$ 
3    $NOCs \leftarrow merge(OCs, ids);$ 
4 else
5   if hasOnlyThisCall( $OCs$ ) then
6      $refactor \leftarrow createRefactorMethod(OCs);$ 
7      $NOCs \leftarrow merge(OCs, refactor);$ 
8   else
9      $refactor \leftarrow createRefactorMethod(OCs);$ 
10     $ids \leftarrow createConstructorIDS(OCs);$ 
11     $NOCs \leftarrow merge(OCs, refactor, ids);$ 
12 refactorCallSites();
13 return  $NOCs$ ;

```

<pre> 1 public ArgExp(Option opt) { 2   this.option = opt; 3   this.threshold = 1; 4 } 5 public ArgExp(String msg) { 6   this.msg = msg; 7 } </pre> <p>Original code</p>	<pre> 1 public ArgExp(Option opt) { 2   this(opt.getKey()); 3 } 4 5 public ArgExp(String msg) { 6   super(msg); 7 } </pre> <p>Original code</p>	<pre> 1 public ArgExp(Option opt) { 2   this(opt.getKey()); 3   this.option = opt; 4 } 5 public ArgExp(String msg) { 6   super(msg); 7 } </pre> <p>Original code</p>
<pre> 1 public ArgExp(int id, Option opt, String msg){ 2   if (id == 0) 3   { 4     this.option = opt; 5     this.threshold = 1; 6   } 7   else { 8     this.msg = msg 9   } </pre> <p>Transformed code</p>	<pre> 1 public static ArgExp ArgExp1(Option opt) { 2   return new ArgExp(opt.getKey()); 3 } 4 5 public ArgExp(String msg) { 6   super(msg); 7 } 8 9 </pre> <p>Transformed code</p>	<pre> 1 public static ArgExp ArgExp1( 2   int id, Option opt, String msg){ 3   if (id == 1) 4   { 5     return new ArgExp(id, opt, opt.getKey()); 6     return new ArgExp(id, opt, msg); 7   } 8   public ArgExp(int id, Option opt, String msg){ 9     super(msg); 10    if (id == 1) 11      this.option = opt; 12 } </pre> <p>Transformed code</p>
(a)	(b)	(c)

Fig. 3. Constructor overloading patterns and their corresponding transformations.

For overloaded methods, ALPHATRANS makes each method name unique by adding an integer suffix (starting at 0) to the name and updates all call sites based on the new method names. Resolving overloaded constructors is not as straightforward, as they should have the same name as the enclosed declaring class. Furthermore, the invocation of constructors inside each other and the Java inheritance mechanism makes constructor overloading complex. Our algorithm (Algorithm 1) for resolving the constructor overloading handles three prominent<sup>2</sup> use cases shown in Figure 3.

The first pattern (Figure 3-a) shows multiple independent constructors. ALPHATRANS merges these constructors into one and uses an `id` parameter to differentiate between them. All the call sites of the constructors will be updated accordingly, based on the given `id`. The second use case (Figure 3-b) is more challenging to resolve, as one constructor calls the other using `this()`. ALPHATRANS transforms the first constructor into a factory method and invokes the second constructor inside it. Factory methods are static, and ALPHATRANS updates the call sites by directly invoking them on the class instance, e.g., `ArgExp.ArgExp1(id,opt,msg)`. The last use case (Figure 3-c) is similar to the second, except that both constructors implement some code. ALPHATRANS refactors the first constructor into a factory method and adds an extra `id` parameter to differentiate between behaviors implemented by different constructors. The call sites of the constructors will be updated like in previous cases. Real-world projects often combine these patterns, which ALPHATRANS handles using Algorithm 1.

## 4.2 Program Decomposition

Translating the entire repository of real-world projects is a very complex problem. As a result, ALPHATRANS breaks down projects into fragments, performs the translation and validation at the fragment level, and re-composes the translation as a repository in the target language.

**4.2.1 Source Decomposition.** Real-world projects can include hundreds of files with thousands of lines of code, which exceed the context window of state-of-the-art LLMs. ALPHATRANS employs static analysis to decompose code into smaller fragments, i.e., *field fragments* and *method fragments*. A field fragment includes modifiers, type, name, and field value. A field fragment can belong to an application or test class. A method fragment includes the full signature of an individual method in the source language. A method fragment can be an application or test method (e.g., helper methods or unit tests). During decomposition, ALPHATRANS extracts meta-information related to the fragments, such as their location in code (e.g., start and end line), code (e.g., implementation between start and end line), dependencies (e.g., callers and callees), types (of inputs and output), and other necessary information such as file paths, class inheritance, imports, and method annotations. ALPHATRANS stores fragments and their corresponding collected meta-data in a data structure

<sup>2</sup>Following the best practices for constructor overloading from *Stack Overflow* and analyzing the use of constructor overloading in open-source projects, we categorized the use cases into the mentioned rules.

called *schema*, which components in other phases will use. ALPHATrans also extracts the call graph to guide the translation, i.e., translate fragments in reverse call order.

**4.2.2 Test Decomposition.** The burden of validating *functional equivalence* in ALPHATrans is on GraalVM. Yet, we still need to translate and execute tests to validate the fragments that GraalVM cannot be used to validate (§6.1). Unit tests in real-world projects can invoke multiple methods and include multiple assert statements. Furthermore, long call chains are inevitable in real-world projects due to the high degree of intra- and inter-procedural dependencies. As we show later (§7.4), the average number of direct method invocations and method executions in tests for our studies subjects are 3 and 27, respectively. This can result in *test translation coupling effect* discussed in §2.

To enable test translation for validating runtime validation or functional equivalence, ALPHATrans decomposes each unit test into a series of *test fragments*, as shown in Figure 1-d. It uses each statement with a call to an application method as a cutting point. For statements enclosed by branches, loops, or exception-handling blocks, ALPHATrans includes the entire block. This process generates an ordering of executable test fragments for each unit test. Each test fragment includes all the statements of the lower-order fragments, along with additional statements that invoke *only one* method that was not invoked by previous fragments. ALPHATrans executes test fragments in increasing order until a test fails and skips running following fragments, as they will also fail.

## 5 Type Translation and Skeleton Construction

### 5.1 Type Translation

Automatically resolving types is a challenging problem [21, 53], and a large body of work has attempted to address this, mostly using symbolic rule-based approaches [5, 8, 9, 32, 46, 47]. ALPHATrans employs a Retrieval-Augment Generation (RAG) [34] technique for finding the equivalent types in target language. To that end, it first extracts all the types in the source language of a given project. Custom, application-level types will be resolved during the translation as ALPHATrans translates fragments and classes in the target language. For the remaining types, it crawls the online API documentation for the source language and retrieves the relevant description of each type.

To form the prompt,<sup>3</sup> ALPHATrans uses the retrieved description and instructs the model with an in-context example to return the most relevant type in the target language, given the use of types in the source language and the retrieved description. To account for potential hallucination in LLM’s response, i.e., returning a type that does not exist in Python, ALPHATrans employs a simple Python script, uses the translated type as an annotation, executes the script, and keeps those without any syntactic or runtime issue. The types in the source language and their corresponding in the target language form a data structure called *Universal Type Mapping*. In practice, ALPHATrans reuses or augments the mapping when translating new projects.

### 5.2 Skeleton Construction

ALPHATrans builds the project’s structure in the target language before translation. This step is necessary for compositional translation and validation, as ALPHATrans can insert the translated

```

0 class ArgExp(Exception):
1     __option: Option = None
2     @staticmethod
3     def ArgExp1(id: int, opt: Option, \
4                 msg: str) -> ArgExp:
5         pass
6     def __init__(self, id: int, opt: \
7                 Option, msg: str):
8         pass
9     def getMsg(self) -> str:
10        pass
11    def getNumArgs(self) -> int:
12        pass
13    def getLen(self) -> int:
14        pass
15    def createMsg(self) -> None:
16        pass

```

Fig. 4. Target PL skeleton for example in Figure 1-b.

<sup>3</sup>Due to space limit, we could not include the prompt. Please refer to the artifacts to see the prompts used for type resolution.



fragments into the project, compile it, or even execute the existing translated test suites, gradually completing the translation. At this step, ALPHATRANS also resolves Java-specific features in Python before starting the translation. Specifically, it resolves circular imports and dependencies, inner classes, interfaces, and abstract classes. Figure 4 shows the class skeleton corresponding to the illustrative example of Figure 1-b.

At the first step, ALPHATRANS creates classes corresponding to each application class in Java. The fields for Python classes are set to `None`, and ALPHATRANS uses the information in schema (§4.2.1) to ensure the naming corresponds to the type of their access modifier in the source language. In the example of Figure 4, the translation of field `private Option option;` in Java is `__option: Option = None`. The classes also include method signatures, with their body set to `pass`. ALPHATRANS uses the universal type mapping (§5.1) to use relevant types in the full method signature. Once the initial skeleton is created, ALPHATRANS leverages the extracted call graph during program decomposition to detect circular dependencies. If there exist any, ALPHATRANS resolves them with local imports. For inner classes, ALPHATRANS unfolds them in Python and uses *dot notation* to access certain methods and fields (e.g., `Class.methodName`). Finally, ALPHATRANS implements best practices in Python and sub-classes all abstract classes and interfaces from `abc.ABC class`. `ABC` is a class from the `abc` module in Python standard library, which is used for defining abstract base classes.

## 6 Compositional Translation and Validation

In this section, we discuss ALPHATRANS main algorithm for compositional translation and validation. ALPHATRANS translates fragments in the reverse call order. For each fragment, it employs the logic demonstrated by Algorithm 2. The algorithm takes a fragment  $F$ , LLM  $M$ , and a series of parameters as inputs, translates the fragment, and returns translation outcomes: (1) syntax check (“non-parseable”, “parseable”), (2) functional equivalence check (“graal-fail”, “graal-success”, “graal-error”), and (3) translated test execution check (“not-exercised”, “test-fail”, “test-success”).

ALPHATRANS employs iterative and feedback-based prompting. That is, if one of the mentioned checks fails, e.g., the translated fragment is not syntactically correct, it prompts the model for another translation attempt. To control the number of iterations, ALPHATRANS considers a reprompting budget, i.e., *repromptBudget*. The algorithm takes the minimum (*min<sub>budget</sub>*) and maximum (*max<sub>budget</sub>*) values for the budget and dynamically sets reprompting budget to a number between them based on the coverage information, e.g., the budget is close to *max<sub>budget</sub>* for a fragment if it is exercised multiple times (high hit rate based on coverage information) by different unit tests. The rationale is to give more importance to fragments covered by more tests to eventually increase translation validation success. The main translation and validation loop (lines 3–31) runs until the assigned budget is exhausted.

Inside the loop, ALPHATRANS first crafts a unique prompt based on the template shown in Figure 5 and then instructs the LLM to translate the fragment (lines 4–5). It then validates the generated translation in multiple steps. The first step checks for syntactical correctness and assigns proper labels to *TVO* (lines 6–10). Then, ALPHATRANS leverages GraalVM for isolation-based validation of fragment  $F$  (lines 12–20), if there exists a test in the source language covering the fragment during its execution. Finally, it translates and executes decomposed fragment tests: if there are no eligible tests (a test becomes eligible if all its dependencies are translated) for the fragment, ALPHATRANS simply assigns the “not-exercised” label to the fragment and moves on to the next one (lines 21–23). Otherwise, it translates the tests, executes them to validate the fragment, and assigns test outcome labels in *TVO* (lines 24–31). In case of a test failure, ALPHATRANS extracts all involved fragments and re-prompts them with feedback extracted from test execution.

Due to inherent intra- and inter-procedural dependencies in real-world projects, the number of fragments involved in re-prompting could be high, logarithmically increasing the translation



time. ALPHATrans filters out those with GraalVM label “*graal-success*”, ranks the remaining based on suspiciousness score, and re-prompts *topK* suspicious fragments. The suspiciousness score for fragments is calculated such that a fragment with more failing tests will get a higher score and, therefore, ranked higher among other fragments.

Our prompt template consists of *five* distinct parts as shown in Figure 5. The first part is the *persona* message used by DeepSeek-Coder-33b-Instruct during instruction fine-tuning and is required for producing the best outputs. The next part introduces the *In-Context Learning* (ICL) example, which reflects the complexities of code translation and instructs LLM on how to deal with them. The green part indicates the natural language instruction given to the model. After describing the objective, the prompt embeds the source Java code along with *partial translation* as a skeleton, which includes all dependencies and translations of the fragments invoked by the current one. The prompt concludes with *### Response: special keyword* to guide the model for code generation.

## 6.1 Language Interoperability

GraalVM [41, 58] is a Java Development Kit by Oracle. It offers the *Polyglot* API [20], which allows the integration of programs written in different guest languages within a Java-based host application. In the context of this paper, GraalVM allows developers to execute Python code from Java and vice versa. ALPHATrans leverages the Polyglot API to perform *in-isolation* validation of the fragments by replacing the Java implementation of a method with its translated Python version while keeping the rest of the project in Java. It then executes the Java tests covering the fragment to validate the functional equivalence of the translation.

The Polyglot API allows access to Python data objects from Java and vice-versa, as these objects reside in a shared memory space. However, objects must be cast to appropriate types for passing parameters to and processing returned values from polyglot calls. The Polyglot API can perform this casting implicitly for only a few simple data types. ALPHATrans builds on top of the Polyglot API to provide a framework to create a Python program state that is isomorphic to the Java program state. The Python translation is restricted to this isomorphic state, and the states are synchronized after method calls to preserve the isomorphism. ALPHATrans allows for the casting of user-defined types as well as several built-in and library types. Using both the static and the dynamic type information of Java objects, ALPHATrans can disambiguate the target types when casting Python

---

### Algorithm 2: Compositional Translation and Validation

---

**Inputs:** Fragment  $F$ , Model  $M$ ,  $min_{budget}$ ,  $max_{budget}$ , and  $topk$  suspicious methods  
**Output:** Translation and Validation Outcome  $TVO$

```

1  $feedback \leftarrow 0$ ;  $TVO \leftarrow \{\}$ ;
2  $repromptBudget \leftarrow$ 
    $getAdaptiveBudget(F, min_{budget}, max_{budget})$ ;
3 while  $repromptBudget > 0$  do
4    $prompt \leftarrow generatePrompt(F, feedback)$ ;
5    $translation \leftarrow translateFragment(prompt, M)$ ;
6   if  $!syntacticCheck(translation)$  then
7      $TVO["syntax\_outcome"] \leftarrow "non - parseable"$ ;
8      $feedback \leftarrow getFeedback(translation)$ ;
9      $repromptBudget \leftarrow repromptBudget - 1$ ;
10    continue;
11    $TVO["syntax\_outcome"] \leftarrow "parseable"$ ;
12   if  $!graalCheck(translation)$  then
13      $TVO["graal\_outcome"] \leftarrow "graal - fail"$ ;
14      $feedback \leftarrow getFeedback(translation)$ ;
15      $repromptBudget \leftarrow repromptBudget - 1$ ;
16     continue;
17   else if  $graalLimitation(translation)$  then
18      $TVO["graal\_outcome"] \leftarrow "graal - error"$ ;
19   else
20      $TVO["graal\_outcome"] \leftarrow "graal - success"$ ;
21   if  $!hasEligibleTests(F)$  then
22      $TVO["test\_outcome"] \leftarrow "not - exercised"$ ;
23     break;
24    $testTranslation \leftarrow getTestTranslation(F)$ ;
25   if  $!testCheck(translation, testTranslation)$  then
26      $TVO["test\_outcome"] \leftarrow "test - fail"$ ;
27      $repromptSuspiciousMethods(testTranslation, topK)$ ;
28      $repromptBudget \leftarrow repromptBudget - 1$ ;
29     continue;
30    $TVO["test\_outcome"] \leftarrow "test - success"$ ;
31   break;
32 return  $TVO$ ;
```

---

objects to Java types. It further preserves object identities and aliasing during such casting, and can also propagate exceptions across language boundaries.

To validate the translation of a method  $m$  in isolation, ALPHATRANS creates an instrumented version of the Java source code. We will refer to this instrumented Java project as the *primal project*,  $P_m$ . During instrumentation, ALPHATRANS replaces the original Java implementation,  $m_J$  of  $m$  with a polyglot call to its Python implementation,  $m_P$ .  $m_P$  resides inside a Python project, which we will refer to as the *dual project*,  $D_m$ . The structure of  $D_m$  is similar to that of the original Java project. All other methods in  $D_m$  wrap a call to the corresponding methods in  $P_m$ . Doing so provides an interface for  $m_P$  to execute with access to all other methods and fields, although these are defined only in  $P_m$ . Using the call graph for the Java project, ALPHATRANS determines all test methods that invoke  $m$  and executes them in  $P_m$  to validate the translation  $m_P$ .

This *in-isolation* validation approach is limited in the sense that it can handle only a limited number of built-in and library types. In certain cases, like reference cycles involving maps or objects with impure methods for hashing, the isomorphism between Java and Python states may not be maintained. Furthermore, it may sometimes not be possible to disambiguate target types when casting Python objects to Java types, for example, if the target object has type `List<Object>`.

## 6.2 Target Program Recomposition

ALPHATRANS updates project skeletons after each fragment translation, gradually constructing the project in target PL. Specifically, it combines class skeletons with the body of the translated fragment and creates stand-alone Python files. The recomposed Python files are then used for dynamic validation.

## 6.3 Test Translation

Similar to translating application method fragments, ALPHATRANS also translates test fragments. Using the dependency information captured during static analysis, it crafts prompts for unit tests along with their dependencies for the model to translate. The ICL examples used for prompting test fragments differ from prompts used for translating application method fragments. The focus of ICL examples here is to prevent LLM from hallucinating the usage of assert statements in the source to target PL. To construct ICL examples for test fragment translation, we created a pool of in-context examples, where each example shows the Python assert statements equivalent to Java assert statements in the context of a test. When prompting a test fragment, the ALPHATRANS detects the assert statement in the fragment and retrieves the corresponding examples from the pool. For translated tests, only syntactic validation is performed as there is no other means of validating their translations.

## 7 Evaluation

To evaluate different aspects of ALPHATRANS, we investigate the following research questions:

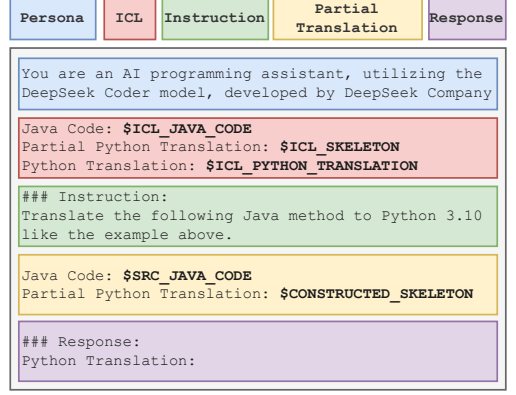


Fig. 5. Prompt structure utilized in ALPHATRANS.

- RQ1: Effectiveness.** To what extent ALPHATrans can automatically resolve types from source to target PL? Can ALPHATrans effectively translate real-world projects?
- RQ2: Translation Bugs and Fixes.** How much effort do developers spend completing the partial translations created by ALPHATrans and achieving all passing tests? What is the nature of translation bugs?
- RQ3: Impact of Test Decomposition.** To what extent does test decomposition impact the validation results?
- RQ4: Impact of Test Coverage.** To what extent does a test suite with higher coverage impact the validation results?

## 7.1 Experiment Setup

We followed three steps for selecting subjects:

**1- Mining.** We mined GitHub and retrieved a list of repositories that use Java as the primary language, are self-contained (include build files, etc.), and have more than 30 stars with at least one commit pushed within the last 12 months.

**2- Filtering.** We filtered out projects based on the number of call edges in their call graphs: removed those with less than 2,000 call edges to ensure the subject projects are big enough to challenge ALPHATrans. We also removed those with more than 30,000 call edges to reduce the computation and manual effort for further steps. Per GraalVM requirements, we only selected projects we could successfully build (compile and achieve green tests) using Java at 21.

**3- Reduction.** ALPHATrans currently supports the following Java APIs: core Java API (`java.util`, `java.text`, `java.lang`, `java.io`, `java.nio`, `java.net`, `java.time`, and `java.math`) and third-party libraries (`org.opentest4j`, `org.slf4j.Logger`, and `org.junit`). We automatically removed all other third-party library dependencies and their usage in the source code in the selected projects. We chose a project if at least 50% of its total methods were preserved after such process. Table 1 shows the list of ten projects used in the evaluation of ALPHATrans and details about their size (classes, methods, tests, and fragments).

ALPHATrans uses CodeQL [19] and tree-sitter [57] for static analysis. For running tests, validating translation, and computing coverage, ALPHATrans uses GraalVM 21.0.3 + 7.1 [41], JUnit 4 and 5 [52], Pytest 8.2.1 [45], JaCoCo [51], and Python’s coverage [4]. ALPHATrans can work with API-access and open-access LLMs. We considered the following criteria for selecting the LLM: (1) for better reproducibility, we prioritized open-access models; (2) due to computing constraints, we wanted an LLM with moderate size ( $> 10B$  but  $< 70B$  parameters); (3) the model should perform reasonably well in code-related tasks; and (4) the model should have fast inference time due to the huge number of prompts. Per the mentioned criteria, we selected DeepSeek-Coder-33b-Instruct for our experiments. We prompted DeepSeek-Coder-33b-Instruct under temperature 0 setting to make translations reproducible. We used the LLM’s default setting for other parameters. For the base prompting (Algorithm 2), we set the minimum and maximum values of the reprompting budget to 3 and 5. For the feedback prompting, we set the reprompting budget to 1, i.e., ALPHATrans attempts to fix issues with feedback for only once.

## 7.2 RQ1: Effectiveness of ALPHATrans

In this RQ, we evaluate ALPHATrans in (1) type translation and skeleton construction (§7.2.1) and (2) compositional translation and validation (§7.2.2).

**7.2.1 Effectiveness in type resolution and skeleton validation:** ALPHATrans extracted 1,797 distinct types from the source projects and attempted to translate them to equivalent Python types. Of these, 915 are application-level types (e.g., classes defined within the Java projects) and were directly resolved during skeleton construction. ALPHATrans prompts DeepSeek-Coder-33b-Instruct to

Table 1. Effectiveness of ALPHATRANS in program transformation, automated type translation, and skeleton validation. **ATR**: Automated Types Resolution, **SV**: Skeleton Validation.

Subjects	# Classes	# Methods	# Original Tests	Method Coverage (%)	ATR (%)	SV (%)	# Fragments			
							Fields		Application	Test
							Application	Test	Methods	Methods
cli [10]	58	664	437	93.4	96.6	100	104	57	273	2,180
codec [11]	156	1780	992	86.5	96.0	100	425	140	680	2,849
csv [12]	41	694	309	87.7	92.3	100	146	35	235	1,272
exec [13]	56	407	70	53.2	78.9	100	104	27	248	327
fast-pfor [33]	82	971	82	50.5	87.4	100	127	14	748	302
fileupload [14]	49	381	39	20.3	98.3	100	121	39	192	194
graph [15]	118	879	146	58.4	97.1	100	216	29	541	977
jansi [49]	48	474	107	25.4	84.8	100	378	0	409	123
pool [16]	98	1097	73	24.6	91.6	100	203	91	682	649
validator [17]	130	1,228	464	61.9	95.5	100	421	209	646	1,463
<b>Total/Average</b>	836	8,575	2,719	56.2	91.9	100	2,245	641	4,654	10,336

resolve the remaining 882, out of which it successfully translated 738  $((915 + 738)/1,797 = 91.99\%)$  of them: generated types passed the syntactic and runtime check. The column *ATR* in Table 1 shows the results of automated type resolution. Since type resolution is essential to project skeleton construction, we manually checked the type mappings generated by DeepSeek-Coder-33b-Instruct, and also attempted to translate the remaining unresolved 144 types.

Through manual investigation of the automatically resolved types, we observed that DeepSeek-Coder-33b-Instruct’s type resolution for 182 cases, while *correct*, can be *improved*. For example, ALPHATRANS translated `java.io.File`, a class concerning file manipulation functionality to `str`. The resolved type can represent file paths in Python but lacks features for file manipulation. We suspect this translation is impacted by the Java use case provided in the prompt. While this translation is correct concerning the use case, we replaced it with `pathlib.Path` to have a more generic type mapping. We also augmented the type mapping with additional types in the target language for 38 types. For example, ALPHATRANS translated `java.nio.Buffer` to `bytearray`, which is correct as they both provide a mutable sequence of bytes with efficient in-place modifications. However, `array.array` and `memoryview` also provide similar functionality with efficient and low-level data manipulation capabilities. Consequently, we augmented type mapping to `typing.Union[bytearray, array.array, memoryview]`. Given that type mapping can be reused, this one-time manual effort increases the chance of ALPHATRANS’s success on unseen projects.

Using the universal type mapping, ALPHATRANS successfully creates and validates project skeletons in target PL, achieving 100% syntax and runtime validation (column *SV* in Table 1). The skeleton validation step ensures all module imports, class structures, method signatures, and type annotations are done properly, making the subsequent steps easier. Applying ALPHATRANS to unseen projects, if a class skeleton cannot be validated, ALPHATRANS removes it from the target project, updates the skeleton based on the class dependencies, and proceeds to the next phase.

**Summary.** ALPHATRANS can successfully transform projects to remove method and constructor overloading. Moreover, it can automatically translate 91.99% of the source PL types and use that to create and validate project skeletons in the target PL.

**7.2.2 Effectiveness in compositional translation and validation:** Table 2 shows the detailed compositional translation and validation results. The *AMF* column indicates the total number of application method fragments. The numbers in subsequent columns demonstrate the effectiveness of ALPHATRANS in translation and validation of *AMFs* only<sup>4</sup>. The *Syntax Check* column indicates the

<sup>4</sup>Due to space limit, please refer to our artifact [31] for details about all translation and validation of other fragments shown in Table 1.

Table 2. Effectiveness of ALPHATrans in compositional translation and validation when using test translation and language interoperability as validation techniques. **AMF**: #Application Method Fragments, **SNEF**: Source Non-Exercised Fragments, **GS**: Graal Success, **GF**: Graal Fail, **GE**: Graal Error, **TNEF**: Target Non-Exercised Fragments, **ATP**: Fragments All Test Pass, **OTF**: Fragments One Test Fail, **MTF**: Fragments Many Test Fail, **ATF**: Fragments All Test Fail, **TPR**: Test Pass Rate, **RE**: Runtime Error, **AF**: Assertion Failure.

Subjects	AMF	Syntax Check (%)	SNEF (%)	GraalVM			Test Translation													M1	
				GS (%)	GF (%)	GE (%)	TNEF (%)	ATP (%)	OTF (%)			MTF (%)			ATF (%)			TPR (%)			
									Overall	RE	AF	Overall	RE	AF	Overall	RE	AF		All	Some	
cli	273	100	6.6	70.7	11.7	11.0	25.6	10.3	12.8	51.4	48.6	35.9	91.0	9.0	8.8	100	0.0	10.1	0	14	
codec	680	98.5	13.5	25.4	49.9	11.2	61.0	4.1	4.4	60.0	40.0	12.8	55.1	44.9	4.1	75.2	24.8	9.4	11	24	
csv	235	98.7	12.3	38.7	26.8	22.1	83.0	0.9	1.3	33.3	66.7	0.0	0.0	0.0	2.6	100	0.0	0.2	0	0	
exec	248	100	46.8	33.5	2.0	17.7	31.5	4.4	2.0	40.0	60.0	8.1	93.6	6.4	7.3	100	0.0	19.3	6	6	
fast-pfor	748	95.3	49.5	11.9	24.6	14.0	35.4	4.8	2.0	86.7	13.3	4.5	79.3	20.7	3.7	87.4	12.6	20.1	6	14	
fileupload	192	100	79.7	8.9	1.0	10.4	7.8	3.6	7.3	28.6	71.4	1.6	87.5	12.5	0.0	0.0	0.0	63.4	2	3	
graph	541	99.6	41.6	24.4	23.7	10.4	53.0	0.4	1.1	100	0.0	1.7	100	0.0	2.2	100	0.0	11.0	0	1	
jansi	409	99.8	74.6	8.1	11.5	5.9	23.5	0.2	1.7	100	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.1	0	1	
pool	682	100	75.4	5.6	2.8	16.3	21.7	1.6	1.0	100	0.0	0.3	100.0	0.0	0.0	0.0	0.0	6.6	4	2	
validator	646	99.2	38.1	30.5	11.1	20.3	44.0	3.1	6.3	92.7	7.3	5.3	78.0	22.0	3.3	97.5	2.5	8.7	2	3	
Total/ Average	4,654	99.1	43.8	25.8	16.5	13.9	38.7	3.3	4.0	69.3	30.7	7.0	68.5	11.5	3.2	66.0	4.0	15.0	31	68	

percentage of *AMFs* that pass syntactic validation. Column *SNEF* shows the percentage of *AMFs* not covered by source PL tests. These results show that ALPHATrans successfully creates syntactically correct code (99.1%, on average). We also observe that 43.8% of *AMFs* are not covered during the execution of any test, i.e., we cannot go beyond syntactic check and validate their runtime behavior or functional equivalence.

For 56.2% of *AMFs* that can be covered by source project tests, ALPHATrans attempts to validate the functional equivalence using GraalVM. The super column *GraalVM* shows the percentage of *AMFs* that GraalVM executes and successfully validates (*GS*), executes but there is a test assertion failure (*GF*), and cannot execute due to its limitation (*GE*) mentioned in §6.1. On average, 25.8% (min=5.6% and max=70.7%), 16.5% (min=1% and max=49.9%), and 13.9% (min=5.9% and max=22.1%) of *AMFs* resulted in *Graal Success*, *Graal Fail*, and *Graal Error*, respectively. Note that these numbers add up to 56.2% of *AMFs* that were covered by source program tests. With respect to only covered methods, GraalVM Success will be 45.91%. Furthermore, our analysis shows that a high portion of methods that are not covered by tests are either abstract methods or getter/setter methods. If their translations are syntactically correct, they are also likely functionally equivalent, which can ramp up the success rate (§7.5). Spearman Rank Order Correlation [50] indicates a strong positive correlation between the method coverage (Table 1) and *GS* numbers ( $\rho = 0.88$ ), confirming that with better method coverage, the validated *AMFs* are very likely to be higher.

Regardless of GraalVM’s validation for functional equivalence, ALPHATrans translates and executes the test fragments on the recomposed translated project. Super column *Test Translation* shows the results of test translation and execution. Column *TNEF* indicates the ratio of *AMFs* where execution of translated tests never reached them. Columns *ATP* through *ATF* show the number of *AMFs* that ALPHATrans executed using translated test fragments, categorized per the test execution results. On average, for 3.3% of *AMFs*, all the test fragments that covered them were marked as pass (*ATP*). For 4.0% and 7.0%, at least one (*OTF*) or more than one test (*MTF*) failed. For 3.2% of them, all the test fragments failed.

For cases with test failure, columns *RE* and *AF* show the breakdown of whether test failure was due to assertion failure or runtime error. From the breakdown of results, we can observe that most translated test executions terminated with a runtime error due to translation bugs and never reached the assert statement. Our manual investigation confirms that a high rate of runtime errors is due to a relatively small number of fragments with translation bugs. Although test decomposition

helps with *test translation coupling effect* (§2), there is still a high degree of runtime errors due to the long call chains in these projects (the average number of methods executed per test in the original and decomposed test suites are 27.4 and 21.8). As a result, the overall pass rate, i.e., the percentage of recomposed test fragments for the translated projects that pass (column *TPR*), is low.

Finally, we calculated the number of *AMFs* that GraalVM could not execute (numbers under *GE* column) but translated test fragments exercised (column *M1*). *All* indicates the number of *AMFs* with all passing tests, including test fragment with assert statement, indicating the validation of functional equivalence. *Some* corresponds to the number of *AMFs* with at least one passing test, which indicates runtime validation. Overall, test translation can validate the functional correctness and runtime behavior of 31 and 68 fragments that GraalVM could not exercise.

**Summary.** ALPHATRANS effectively performs compositional translation and validation of 6,899 source code fragments, achieving overall 99.1% syntactical correctness, 43.7% *runtime behavior validation*, and 25.8% *functional equivalence*.

### 7.3 RQ2: Translation Bugs and Fixes

This research question presents the manual effort for fixing translation bugs in a subset of studied subjects, namely *Commons-FileUpload*, *Commons-CLI*, *Commons-CSV*, and *Commons-Validator*. We also discuss some of the translation bugs and fixes for them to better illustrate the challenges in code translation.

**7.3.1 Human Study.** Our two human subjects were selected due to their relative familiarity with the selected projects. Their effort indicates an upper bound for the amount of time required to fix translation bugs since developers of the source projects are likely to fix the bugs better and faster. We shared with them the source program in Java, the translations by ALPHATRANS, and all the reports and artifacts generated by ALPHATRANS during translation.

For *Commons-FileUpload*, achieving green tests took roughly 5.5 hours and required 120 line additions and 114 line deletions from partial translation generated by ALPHATRANS. For *Commons-CLI*, the manual fix took roughly 11 hours to fully achieve all passing tests, making 614 and 1,253 line additions and deletions, respectively. For *Commons-CSV*, the project was very dense, with a lot of method calls, making it harder to manually fix bugs. Nevertheless, a developer achieved all green tests in 30 hours with 2,676 and 999 line additions and deletions, respectively. Finally, for *Commons-Validator*, the developer spent 34 hours to fix translation bugs, with 3,585 and 2,416 line additions and deletions, respectively. One of the major feedback from developers was that test decomposition greatly helped locate and fix translation bugs: in case of a test failure, developers only need to investigate the last call statement in the failed test fragment instead of looking at the stack trace and other prior calls (more quantified details about the impact of test decomposition on validation in §7.4).

**7.3.2 Translation Bugs.** Our publicly available artifacts [31] contain partial translations and fixed versions as separate commits. These commits can serve as useful benchmarks for evaluating techniques such as fault localization, program repair, and test generation. This section shows *four* instances of such translation bugs. The two most prevalent sources of translation bugs are mismatches between APIs and behavioral differences in the PLs. The code snippet below demonstrates a bug that happened due to a mismatch in the logic of *Calendar* (Java) and *datetime* (Python). Line 3 in Java sets the *MONTH* field to 0, which corresponds to the first month of the year (*January*). Similarly, the Python translation sets the *month* attribute to 0; however, in the Python library, *January* is the first month, i.e., the correct translation should use index 1.

```

1 ----- JAVA SOURCE CODE -----
2 Calendar calendar = Calendar.getInstance()
3 calendar.set(Calendar.MONTH, 0);

```

```

1 ----- PYTHON TRANSLATION -----
2 calendar = datetime.datetime.now()
3 - calendar = calendar.replace(month=0)
4 + calendar = calendar.replace(month=1)

```

The next example shows the difference in implicit type casting between the two Pls. Line 5 in Java source code concatenates a String with `nullStr = null`. During execution, Java runtime silently casts `null` to a String and then performs the binary operation on it. In Python, concatenating an `str` with `None` results in a `TypeError` as the operands of the binary operation has different types. A correct Python translation requires explicit casting of `None` to `str` as shown in Line 5.

```

1 ----- JAVA SOURCE CODE -----
2 qChar = " ";
3 nullStr = null;
4
5 this.qNullStr = qChar + nullStr + qChar;

```

```

1 ----- PYTHON TRANSLATION -----
2 qChar = " "
3 nullStr = None
4 - self.qNullStr = qChar + nullStr + qChar
5 + self.qNullStr = qChar + str(nullStr) + qChar

```

The third example shows an instance of `write(int b)` method from `ByteArrayOutputStream` class, where the least significant 8 bits of the integer  $(b2 \ll 4) \mid (b3 \gg 2)$  are directly written to the stream. The incorrect Python translation attempts to construct a bytes object using a singleton list with the input integer before writing it to an object of `io.BytesIO`. However, this neglects that the `bytes()` constructor requires the integers in the input iterable to be strictly in the range of  $[0, 255]$ . Thereby, a `ValueError` is thrown when `b2` is large. The correct translation requires `0xF`, a mask that maintains only the 4 lowest bits of `b2` before left-shifting by 4 as shown in Line 3 under Python translation. Given that `b3` and `b4` each contain no more than 8 bits, this change ensures the least significant 8 bits of  $(b2 \ll 4) \mid (b3 \gg 2)$  are correctly written to the `BytesIO` object.

```

1 ----- JAVA SOURCE CODE -----
2
3 out.write((b2 << 4) | (b3 >> 2));

```

```

1 ----- PYTHON TRANSLATION -----
2 - out.write(bytes([(b2 << 4) | (b3 >> 2)]))
3 + out.write(bytes([(b2 & 0xF) << 4 | (b3 >> 2)]))

```

The last code snippet shows an example where the Java behavior of an iterator is unavailable in Python. Specifically, the incorrect Python translation uses `next()` to implement both `.next()` and `.hasNext()` calls of Java iterator. The issue is that calling `next()` increments the iterator in Python. The correct translation should implement `PeekableIterator` interface in Python with a method `hasNext() -> bool`.

```

1 ----- JAVA SOURCE CODE -----
2 Iterator<String> headers = ls.keySet().iterator();
3
4 assertEquals("content", headers.next());
5
6 assertFalse(headers.hasNext());

```

```

1 ----- PYTHON TRANSLATION -----
2 - headers = iter(ls.keys())
3 + headers = PeekableIterator(ls.keys())
4 self.assertEqual("content", next(headers))
5 - self.assertFalse(next(headers, None) is not None)
6 + self.assertFalse(headers.hasNext())

```

**Implications.** In order to obtain correct translation, especially for translating APIs, models need to generate test cases as well, which will validate the translated fragment in isolation. This could be an interesting direction for applying the agentic approach, where the orchestrating agent can decide when to generate test cases, and the test case generator agent gets all the information by running static analysis tools, gathering context from previous runs, collecting API documentation by crawling internet, and can finally generating the translation based on all the information.

**Summary.** Although ALPHATRANS cannot validate all the translations, it provides partial translations and artifacts that developers can use to complete the translation and achieve green tests in a reasonable time (20.1 hours, on average).

## 7.4 RQ3: Impact of Test Decomposition

We previously showed the effectiveness of test translation in validating the runtime behavior or even the functional correctness of application method fragments (§7.2.2). To better understand the



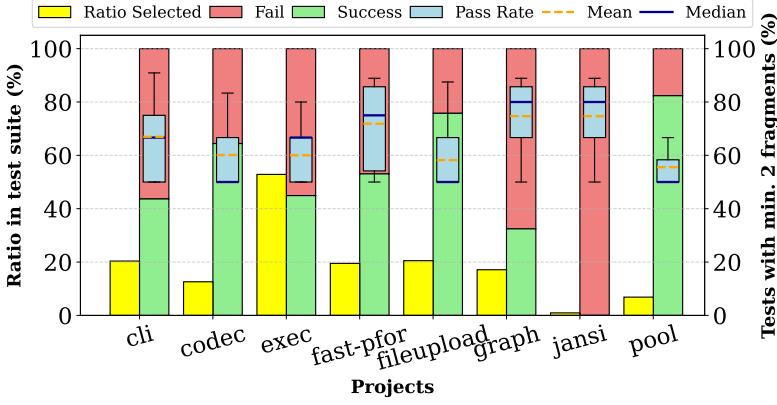


Fig. 6. Effectiveness of decomposing test suites in ALPHATRANS for validating earlier fragments in failing tests.

effect of test decomposition and how it helps with the *test translation coupling effect*, we collected translated unit tests with at least two corresponding decomposed test fragments. We further filtered out these unit tests and kept those that *all* their decomposed test fragments were executed, regardless of passing or failing. The yellow bars in Figure 6 show the percentage of our selected unit tests from the translated test suites. For *Commons-CSV* and *Commons-Validator*, none of the unit tests met the criteria, so we excluded them from further investigation.

We categorized the selected unit tests into two groups: Those with all their test fragments pass (green bars in Figure 6), and those with at least one test fragment fail (red bars in Figure 6). For the *unit tests* in the latter group, we calculate the pass rate of the *decomposed test fragments*. The blue box chart in Figure 6 shows the distribution of the measured pass rate per unit test. These unit tests would have been marked as *fail* without test decomposition. However, we can observe that these tests can be decomposed into test fragments that 63.43% of them pass, on average. These results confirm how decomposed test fragments were useful in helping developers localize translation bugs more easily and resolve the translation bugs faster.

**Summary.** Test decomposition unburdens validation of application method fragments from incorrect translations. On average, 63.43% of test fragments for unit tests that would have been marked as failed achieved a test pass.

## 7.5 RQ4: Impact of Test Coverage

While existing developer-written tests are useful for validating functional equivalence, they can pose two major issues for automated code translation and validation. First, the coverage for these tests can be extremely low (e.g., 20.3% for Commons-Pool [16]), preventing most of the code from being validated (as shown in RQ1, the translation validation rate strongly correlates with test suites' (method) coverage). Second, executing a developer-written test can have a long call sequence. To show the positive impact of better, more focused tests with higher coverage on translation validation, we automatically generated more tests using EvoSuite [18]. We generated EvoSuite tests with the test generator's default settings: used DynaMOSA [43] as optimization algorithm and set a timeout of 120 seconds as stopping criteria.

Table 3 compares the properties of the developer-written tests and EvoSuite tests. Since we configure EvoSuite to generate several tests per each Java class, the average *Method Coverage* of its tests is higher than developer-written tests (79.3% compared to 56.2%). Furthermore, the average number of methods executed per single test is almost half that of decomposed test suites (11.4

Table 3. Effectiveness of test augmentation in exercising and validating more application method fragments.

Subjects	Developer-Written Test				EvoSuite Test			
	Method Coverage (%)	# Decomposed Tests	Avg. Methods Executed / Test	TPR (%)	Method Coverage (%)	# Tests	Avg. Methods Executed / Test	TPR (%)
cli	93.4	3036	34.3	10.1	83.5	569	12.2	13.9
codec	86.5	3522	10.6	9.4	86.9	1141	8.0	48.6
csv	87.7	1219	52.6	0.2	68.5	220	39.2	8.6
exec	25.4	311	19.0	19.3	76.8	245	6.3	28.2
fast-pfor	61.9	249	41.6	20.1	82.6	1843	4.3	30.9
fileupload	50.5	93	3.5	63.4	84.4	231	5.3	41.1
graph	58.4	933	25.0	11.0	86.3	800	9.0	9.1
jansi	24.6	187	13.6	1.1	70.3	332	9.1	14.2
pool	20.3	287	6.5	6.6	69.1	394	7.4	5.6
validator	53.2	1479	11.7	8.7	84.5	1305	13.4	3.3
<b>Total/Average</b>	56.2	11316	21.8	15.0	79.3	7080	11.4	20.3

compared to 21.8 methods). We translated EvoSuite tests and executed the translated projects using them. As corroborated by the numbers under columns *TPR*, we can see that we achieve a higher test pass rate (20.3% compared to 15.0%) on the same translated code. Note that not all the EvoSuite tests have assertions, and even if they do, the quality of the assertions could be lower compared to developer-written tests (e.g., checking trivial properties that result in a test pass). Thereby, we only claim that higher *TPR* of such tests enhance runtime behavior validation, which is still promising in code translation. Unfortunately, EvoSuite is not compatible with Java 21, which is required for the GraalVM component of ALPHATrans. Otherwise, we could use the generated tests in ALPHATrans. We anticipate including it in the loops could have also improved the overall quality of translations.

**Summary.** Augmenting existing test suite increases code coverage for exercising and validating more application method fragments. Moreover, the generated tests are more focused and on average invokes 50% less methods than developer-written tests.

## 8 Related Work

There are generally two main approaches to translating code from one programming language to another: using transpilers and statistical machine translation and leveraging language models.

*Code translation using non-LLM-based approaches.* In this domain, tools like C2Rust [24], CxGo [56], Sharpen [44], and Java2CSharp [25] have been developed to translate code from C to Rust, C to Go, and Java to C# respectively. However, a recent study [42] revealed that, apart from C2Rust and CxGo, other tools lack proper maintenance. For CxGo, language models outperform traditional approaches, whereas for C2Rust, language models generate safer but less effective code, aligning with the primary goal of translating C code to Rust. In terms of statistical machine translation, works by Nguyen *et al.* [36–38], and Chen *et al.* [6] focus on translating Java to C#. Additionally, deep learning approaches have been utilized for code translation [47, 48]. However, none of these efforts have tackled the translation of real-world Java projects to Python.

*Code translation using LLMs.* Recently, large language models have been employed for code translation [7, 28, 42, 55, 59, 62], demonstrating high success rates on crafted examples but poor performance on real-world projects. Other studies [1, 63] have also utilized language models for code translation, mainly focusing on crafted benchmarks. Recently, there have been works that use transpiler output to guide the code translation [60]. However, the limitation of such work is the availability of robust and well-maintained transpilers, which, in many cases, may not be a feasible solution. Nitin *et al.* [39] introduced a specification-based translation, where natural language specification has been captured from the source code, which helps the translation process. Whereas

Yang et al. [61] used tests to assist the translation. However, compared to previous works, the major differences are (a) the first attempt to translate a real-world project, (b) modular translation, and (c) a validation-guided translation approach.

## 9 Threats to Validity

Like most approaches, ALPHATrans possesses some limitations and comes with a list of threats to the validity. In this section, we will discuss how we mitigated various threats.

*External Validity.* One of the key external threats is the generalizability of our approach. The translation pipeline is very generic and can be extended for more language pairs. Also, the majority of the tools that we used support a large set of programming languages. However, to expedite the research in repository-level translation, we build the first version of ALPHATrans to be specific towards Java to Python.

*Internal Validity.* One major internal threat can be that the results are calculated based on a single run. To mitigate that threat, we ran ALPHATrans with greedy decoding and temperature 0, which reduces the model’s creativity but makes the output consistent with several runs. Another threat in this group can be the manual validation of the translated types. To address that, several authors have verified the types individually and consulted API documents when necessary.

*Construct Validity.* In order to mitigate construct validity, ALPHATrans is built and validated with well-vetted tools like, GraalVM [41], JaCoCo [51], coverage for Python [4], CodeQL [19], etc.

## 10 Concluding Remarks

In this paper, we introduced ALPHATrans, a neuro-symbolic approach that combines the power of static analysis and emerging abilities of Code LLMs in code synthesis to automate repository-level code translation and validation. ALPHATrans decomposes the program into smaller fragments and translates the fragments in the reverse call order, originally building the project in the target language. In addition to syntax check, ALPHATrans implements two levels of validation through GraalVM and test translation. Our results demonstrate the effectiveness of ALPHATrans in translating ten real-world Java projects to Python, achieving 99.1% syntactical correctness, 43.7% runtime behavior validation, and 25.8% functional equivalence. ALPHATrans is the first approach to translate and validate the entire repository, and we envision several research directions to advance repository-level code translation and validation as follows:

One of the major challenges of repository-level code translation is identifying suitable library APIs in the target PL. Often, equivalent Python APIs may not exist, requiring new code generation or translation of the library API itself. Even if similar libraries exist, the logic of libraries might be different in two PLs. ALPHATrans supports translating frequently used APIs and aims to build a generic pipeline. Supporting all the libraries in the pipeline remains an open challenge. Furthermore, while the idea of compositional translation and validation is PL-agnostic, the static analysis makes the extension of ALPHATrans to translating from other source projects challenging. Devising LLM-enabled or PL-agnostic static analysis approaches can benefit code translation approaches such as ALPHATrans. We also showed that the quality of the source project test suite can significantly impact the translation validation results. As part of our future work, we plan to integrate an LLM-based test generator into the ALPHATrans pipeline to advance the validation component.

## 11 Data Availability

The implementation of ALPHATrans and all the artifacts required for reproducing the results presented in this paper are publicly available [31].

## References

- [1] Wasi Uddin Ahmad, Md Golam Rahman Tushar, Saikat Chakraborty, and Kai-Wei Chang. 2023. AVATAR: A Parallel Corpus for Java-Python Program Translation. In *ACL*. ACL, Toronto, Canada, 2268–2281.
- [2] Anthropic AI. 2024. Claude 3. <https://www.anthropic.com/news/claude-3-family>.
- [3] Andrés Bastidas Fuertes, María Pérez, and Jaime Meza Hormaza. 2023. Transpilers: A Systematic Mapping Review of Their Usage in Research and Industry. *Applied Sciences* 13 (2023), 3667.
- [4] Ned Batchelder. 2024. Coverage.py. <https://pypi.org/project/coverage>.
- [5] Dante Broggi and Yi Liu. 2023. On the Interoperability of Programming Languages via Translation. In *CSCE*. IEEE, Las Vegas, NV, USA, 2579–2585.
- [6] Xinyun Chen, Chang Liu, and Dawn Song. 2018. Tree-to-tree Neural Networks for Program Translation. In *NIPS*. Curran Associates Inc., Red Hook, NY, USA, 2552 – 2562.
- [7] Peng Di, Jianguo Li, Hang Yu, Wei Jiang, Wenting Cai, Yang Cao, Chaoyu Chen, Dajun Chen, Hongwei Chen, Liang Chen, et al. 2024. CodeFuse-13B: A Pretrained Multi-lingual Code Large Language Model. In *ICSE-SIEP*. ACM, New York, NY, USA, 418–429.
- [8] George Dony, Girase Priyanka, Gupta Mahesh, Gupta Prachi, and Sharma Aakanksha. 2010. Programming Language Inter-Conversion. *International Journal of Computer Applications* 1, 20 (2010), 63–69.
- [9] Hadeel A. Osman Eman J. Coco and Niemah I. Osman. 2018. JPT : A Simple Java-Python Translator. *CAIJ* 5, 2 (2018), 1–18.
- [10] The Apache Software Foundation. 2024. Apache Commons CLI. <https://github.com/apache/commons-cli>
- [11] The Apache Software Foundation. 2024. Apache Commons Codec. <https://github.com/apache/commons-codec>
- [12] The Apache Software Foundation. 2024. Apache Commons CSV. <https://github.com/apache/commons-csv>
- [13] The Apache Software Foundation. 2024. Apache Commons Exec. <https://github.com/apache/commons-exec>
- [14] The Apache Software Foundation. 2024. Apache Commons FileUpload. <https://github.com/apache/commons-fileupload>
- [15] The Apache Software Foundation. 2024. Apache Commons Graph. <https://github.com/apache/commons-graph>
- [16] The Apache Software Foundation. 2024. Apache Commons Pool. <https://github.com/apache/commons-pool>
- [17] The Apache Software Foundation. 2024. Apache Commons Validator. <https://github.com/apache/commons-validator>
- [18] Gordon Fraser and Andrea Arcuri. 2011. EvoSuite: automatic test suite generation for object-oriented software. In *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering (Szeged, Hungary) (ESEC/FSE '11)*. Association for Computing Machinery, New York, NY, USA, 416–419. <https://doi.org/10.1145/2025113.2025179>
- [19] GitHub. 2024. CodeQL. <https://codeql.github.com>
- [20] GraalVM. 2024. Polyglot API. <https://www.graalvm.org/latest/reference-manual/polyglot-programming>.
- [21] Giovanni Guizzo, Jie M. Zhang, Federica Sarro, Christoph Treude, and Mark Harman. 2023. Mutation analysis for evaluating code translation. *Empirical Software Engineering* 29 (2023), 23 pages.
- [22] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming – The Rise of Code Intelligence. arXiv:2401.14196
- [23] Ali Reza Ibrahimzada. 2024. Program Decomposition and Translation with Static Analysis. In *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings*. 453–455.
- [24] Immunant. 2024. C2Rust Transpiler. <https://github.com/immunant/c2rust>.
- [25] Paul Irwin. 2024. Java to CSharp Converter. <https://github.com/paulirwin/JavaToCSharp>.
- [26] Suman Jain and Indrveer Chana. 2015. Modernization of Legacy Systems: A Generalised Roadmap. In *ICCCT*. ACM, New York, NY, USA, 62–67.
- [27] Pooyan Jamshidi, Aakash Ahmad, and Claus Pahl. 2013. Cloud Migration Research: A Systematic Review. *IEEE Transactions on Cloud Computing* 1 (2013), 142–157.
- [28] Mingsheng Jiao, Tingrui Yu, Xuan Li, Guanjie Qiu, Xiaodong Gu, and Beijun Shen. 2023. On the evaluation of neural code translation: Taxonomy and benchmark. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, Las Vegas, NV, USA, 1529–1541.
- [29] Ravi Khadka, Belfrit V Batlajery, Amir M Saeidi, Slinger Jansen, and Jurriaan Hage. 2014. How Do Professionals Perceive Legacy Systems and Software Modernization? In *ICSE*. ACM, New York, NY, USA, 36–47.
- [30] Musawwer Khan, Islam Ali, Wasif Nisar, Muhammad Qaiser Saleem, Ali S Ahmed, Haysam E Elamin, Waqar Mehmood, and Muhammad Shafiq. 2022. Modernization Framework to Enhance the Security of Legacy Information Systems. *Intelligent Automation & Soft Computing* 32 (2022), 543–555.
- [31] Intelligent CAT Lab. 2024. Artifact Website. <https://github.com/Intelligent-CAT-Lab/AlphaTrans>.
- [32] Kevin Lano and Hanan Siala. 2024. Using model-driven engineering to automate software language translation. *Automated Software Engineering* 31 (2024), 59 pages.

- [33] Daniel Lemire. 2024. JavaFastPFOR. <https://github.com/lemire/JavaFastPFOR>
- [34] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.
- [35] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranajpe, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the Middle: How Language Models Use Long Contexts. *TACL* 12 (2024), 157–173.
- [36] Anh Tuan Nguyen, Tung Thanh Nguyen, and Tien N Nguyen. 2013. Lexical Statistical Machine Translation for Language Migration. In *FSE*. ACM, New York, NY, USA, 651–654.
- [37] Anh Tuan Nguyen, Tung Thanh Nguyen, and Tien N Nguyen. 2014. Migrating Code with Statistical Machine Translation. In *ICSE Companion*. ACM, New York, NY, USA, 544–547.
- [38] Anh Tuan Nguyen, Tung Thanh Nguyen, and Tien N Nguyen. 2015. Divide-and-Conquer Approach for Multi-phase Statistical Migration for Source Code. In *ASE*. IEEE, Las Vegas, NV, USA, 585–596.
- [39] Vikram Nitin, Rahul Krishna, and Baishakhi Ray. 2024. SpecTra: Enhancing the Code Translation Ability of Language Models by Generating Multi-Modal Specifications. *arXiv:2405.18574*
- [40] OpenAI. 2023. GPT-4 Technical Report. *arXiv:2303.08774*
- [41] Oracle. 2024. GraalVM. <https://www.graalvm.org>.
- [42] Rangeet Pan, Ali Reza Ibrahimzada, Rahul Krishna, Divya Sankar, Lambert Pouguem Wassi, Michele Merler, Boris Sobolev, Raju Pavuluri, Saurabh Sinha, and Reyhaneh Jabbarvand. 2024. Lost in Translation: A Study of Bugs Introduced by Large Language Models while Translating Code. In *ICSE*. ACM, New York, NY, USA, 866–866.
- [43] Annibale Panichella, Fitsum Meshesha Kifetew, and Paolo Tonella. 2018. Automated Test Case Generation as a Many-Objective Optimisation Problem with Dynamic Selection of the Targets. *IEEE Transactions on Software Engineering* 44, 2 (2018), 122–158. <https://doi.org/10.1109/TSE.2017.2663435>
- [44] Mono Project. 2023. Sharpen - Automated Java->C# coversion. <https://github.com/mono/sharpen>.
- [45] pytest dev. 2024. Pytest. <https://www.pytest.org>.
- [46] Lili Qiu. 1999. *Programming Language Translation*. Technical Report. Cornell University, USA.
- [47] Baptiste Roziere, Marie-Anne Lachaux, Lowik Chanasot, and Guillaume Lample. 2020. Unsupervised Translation of Programming Languages. In *NIPS*. Curran Associates Inc., Red Hook, NY, USA, 20601 – 20611.
- [48] Baptiste Roziere, Jie M Zhang, Francois Charton, Mark Harman, Gabriel Synnaeve, and Guillaume Lample. 2021. Leveraging Automated Unit Tests for Unsupervised Code Translation. *arXiv:2110.06773*
- [49] Fuse Source. 2024. Jansi. <https://github.com/fusesource/jansi>
- [50] Charles Spearman. 1961. The Proof and Measurement of Association between Two Things. *The American Journal of Psychology* 15 (1961), 72 – 101.
- [51] The JaCoCo Team. 2024. Java Code Coverage. <https://www.eclemma.org/jacoco/>
- [52] The JUnit Team. 2024. JUnit. <https://junit.org/junit5/>
- [53] Andrey A Terekhov and Chris Verhoef. 2000. The Realities of Language Conversions. *IEEE Software* 17 (2000), 111–124.
- [54] TIOBE. 2023. TIOBE Index. <https://www.tiobe.com/tiobe-index>.
- [55] Sindhu Tipirneni, Ming Zhu, and Chandan K Reddy. 2024. StructCoder: Structure-Aware Transformer for Code Generation. *Transactions on Knowledge Discovery from Data* 18 (2024), 1–20.
- [56] Go Transpile. 2024. C to Go Translator. <https://github.com/gotranspile/cxgo>.
- [57] Tree-Sitter. 2024. Tree-Sitter Library. <https://tree-sitter.github.io/tree-sitter/>
- [58] Thomas Würthinger, Christian Wimmer, Andreas Wöß, Lukas Stadler, Gilles Duboscq, Christian Humer, Gregor Richards, Doug Simon, and Mario Wolczko. 2013. One VM to Rule Them All. In *Onward!* ACM, New York, NY, USA, 187–204.
- [59] Weixiang Yan, Yuchen Tian, Yunzhe Li, Qian Chen, and Wen Wang. 2023. CodeTransOcean: A Comprehensive Multilingual Benchmark for Code Translation. In *EMNLP*. ACL, Singapore, 5067–5089.
- [60] Aidan ZH Yang, Yoshiaki Takashima, Brandon Paulsen, Josiah Dodds, and Daniel Kroening. 2024. VERT: Verified Equivalent Rust Transpilation with Large Language Models as Few-Shot Learners. *arXiv:2404.18852*
- [61] Zhen Yang, Fang Liu, Zhongxing Yu, Jacky Wai Keung, Jia Li, Shuo Liu, Yifan Hong, Xiaoxue Ma, Zhi Jin, and Ge Li. 2024. Exploring and Unleashing the Power of Large Language Models in Automated Code Translation. *FSE* 1 (2024), 1585–1608.
- [62] Xin Yin, Chao Ni, Tien N Nguyen, Shaohua Wang, and Xiaohu Yang. 2024. Rectifier: Code Translation with Corrector via LLMs. *arXiv:2407.07472*
- [63] Ming Zhu, Karthik Suresh, and Chandan K Reddy. 2022. Multilingual Code Snippets Training for Program Translation. *AAAI* 36 (2022), 11783–11790.