

# コードクローンに対する集約結果に基づいた 削減可能なソースコード行数の測定手法

中川 将<sup>†</sup> 肥後 芳樹<sup>†</sup> 楠本 真二<sup>†</sup>

<sup>†</sup> 大阪大学大学院情報科学研究科

〒 565-0871 大阪府吹田市山田丘 1-5

E-mail: †{t-nakagw,higo,kusumoto}@ist.osaka-u.ac.jp

あらまし コードクローンとは、ソースコード中に存在する互いに一致または類似しているコード片のことである。コードクローンを一つのメソッドやクラスに集約することにより、ソースコード行数の削減が可能である。削減可能なソースコード行数を推定する手法として、コードクローンの情報を解析する手法が既存研究では提案されている。しかし、コードクローンの中には集約するとコンパイルやテストに失敗するものが存在すると著者らは考えた。そのため、コードクローンの集約により削減可能なソースコード行数を推定するとき、コードクローンの情報のみでは正確な推定を行うことはできない。そこで本研究では、より正確な削減可能なソースコード行数を推定する手法を提案する。提案手法では、それぞれのコードクローンの集合に対して集約、コンパイル、テストを自動で行い、それらが実際に集約可能なかを判定したうえで、削減行数を測定する。さらに提案手法を用いてオープンソースソフトウェアの削減可能なソースコード行数の測定を行った。その結果、コードクローンの情報のみで推定された削減可能な行数とは異なる数値を得た。

キーワード コードクローン, ソフトウェア保守, リファクタリング支援

## 1. ま え が き

一度完成したソフトウェアであっても機能追加やバグ修正などの保守が必要であり、保守にはコストが必要となる [1]。ソフトウェアの保守に必要なコストはそのソフトウェアの規模や複雑さから算出されることが多い [2]。ソフトウェアの規模はソースコードの行数から測定できる。しかし、冗長なソースコードが多く含まれているとソフトウェアの規模を正確に測定できない。したがって、ソフトウェアの規模を測定し保守に必要なコストを見積もるためには、冗長なソースコードを削減した場合の行数を推定する必要がある。

冗長なソースコードの一つとしてコードクローン (以下、クローン) が挙げられる。クローンとは、ソースコード中に存在する互いに一致または類似しているコード片のことである。クローンの主な発生要因としてコピーアンドペーストが挙げられる [3]。クローンとなっているコード片を一つのモジュールに集約することで、冗長なソースコードを削減できる。

クローンを集約する方法の一つとしてリファクタリングが挙げられる。リファクタリングとはソフトウェアの外部的振る舞いを保ちつつ内部の構造を改善することと定義されている [4]。リファクタリングの一つにメソッド抽出という手法がある。メソッド抽出とは、既存のソースコードから一部のコード片を切り出し新たなメソッドにすることである。メソッド

抽出を行うことで、クローンの集合 (以下、クローンセット) を一つのメソッドに集約できる [5]。クローンの集約を行うことでソースコードの行数が削減可能である。そのため、クローンの情報から削減可能なソースコード行数を推定できると考えられている [6]。しかし、検出されたクローンの中には、集約を行うことでコンパイルやテストに失敗するクローンがある。また、異なるクローンが部分的に重なることで、どのクローンを集約するべきか選択する必要がある。このような問題点から、リファクタリングを行うことで削減可能なソースコード行数を推定することは容易ではない。

既存研究では、削減可能なソースコード行数を推定する手法としてクローンの情報を解析するという手法を提案している [7]。しかし、実際にクローンの集約を行っていない。そのため、算出された推定値に対して、その行数が実際に削減可能であるかどうか分からない。また、実際にクローンの集約を行うと、コンパイルとテストに失敗して集約を行えない可能性も考えられる。例えば、検出されたクローン間の対応する変数の型が異なる場合、集約を行うとコンパイルに失敗するので集約を行うことはできない。また、テストを行わなければ外部的振る舞いが保たれていることを保証できない。しかし、検出された大量のクローン全てに対して、手作業で実際に集約が可能かを確認することには多大な労力が必要となる。

そこで本研究では、クローンの検出、ソースコードの変更、

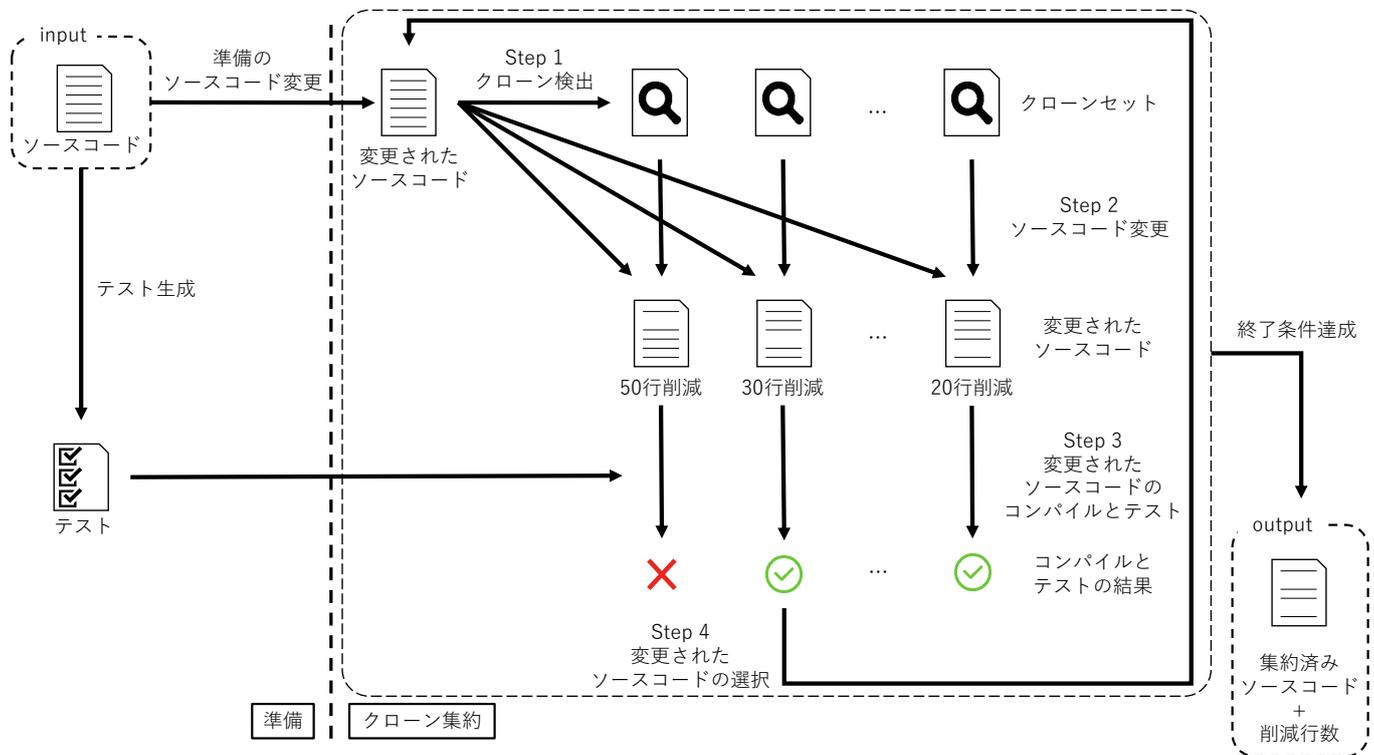


図1 提案手法の概要

コンパイル及びテストを全て自動で行い、実際に削減行数を測定することで、より正確なソースコードの削減可能行数を算出する手法を提案する。また、提案手法を用いて Java 言語で記述されたプロジェクトに対して実験を行い、削減可能行数を算出した。その結果、既存研究とは異なる数値を得た。

## 2. 提案手法

本研究では、クローンの検出、ソースコードの変更、コンパイル、テストを繰り返し行いソースコードの削減可能行数を算出する手法を提案する。提案手法の概要を図1に示す。本研究における提案手法の入力はソースコードである。出力は集約可能なクローンを全て集約したソースコードと削減されたソースコード行数である。

本研究の提案手法は以下の二つの工程で構成される。

- 準備
- クローン集約

### 2.1 準備

準備ではソースコードの変更前後で外部的振る舞いが変化していないかを確認するために、対象とするソースコードの単体テストを用意する。単体テストはソースコードから自動生成する。

また、メソッド抽出で切り出されたメソッドの宣言を行うクラスを用意する。それに加えて、メソッドのアクセス修飾子及びローカル変数の未初期化によるコンパイルエラーとコーディングスタイルによる削減行数の変化を防ぐために、ソースコードに変更を加える。合計で以下の四つの変更を行う。それぞれの変更内容については3.1で詳しく述べる。

- 抽出されたメソッドの宣言を行うクラスの用意
- メソッドのアクセス修飾子の変更

- ローカル変数の初期化
- フォーマッタの適用

この準備はクローン集約を行う度に実行する必要はない。

### 2.2 クローン集約

クローン集約は次の4ステップから構成される。これらはすべて自動で行われる。

Step 1: クローン検出

Step 2: ソースコード変更

Step 3: 変更されたソースコードのコンパイルとテスト

Step 4: 変更されたソースコードの選択

Step 1では、対象とするソースコード中に存在するクローンセットを検出する。Step 2では、ソースコードに対して Step 1で検出したクローンセットの一つを集約する変更を加える。集約の際に抽出されたメソッドは、2.1で説明した準備で用意したクラス内に宣言する。Step 3では、変更されたソースコードに対してコンパイルとテストを行い、外部的振る舞いが変化していないかどうかを判定する。以降、Step 2で変更されたソースコードの中で、行数が削減されコンパイルとテストを通過するようなソースコードを選択可能なソースコードと定義する。Step 1で検出したすべてのクローンセットに対して Step 2と Step 3を行い、Step 4で選択可能なソースコードの中で削減行数が最も多いソースコードを選択する。そして選択したソースコードに対して繰り返し Step 1以降を実行する。終了条件に到達した時点でクローン集約を終了し、削減行数を測定し出力する。クローン集約の終了条件は以下とする。

- クローンが検出されなくなる
- 選択可能なソースコードが存在しなくなる



る必要がある。本来ならプロジェクト毎のコーディング規約を再現したフォーマッタを用意するべきだが、それは困難なため、本研究では適用するフォーマッタは Eclipse の基本設定のフォーマッタで統一した。

### 3.2 クローン集約

#### Step 1: クローン検出

クローン検出の例を図 4 に示す。クローンはブロック単位で検出する。本研究では、Eclipse JDT<sup>(注1)</sup> で Statement として定義されているものの内、以下をブロックとして検出する。

- Block
- DoStatement
- EnhancedForStatement
- ForStatement
- IfStatement
- SwitchStatement
- SynchronizedStatement
- TryStatement
- WhileStatement

ブロック単位で検出されたクローンは字句単位で検出されたクローンと比べると粗粒度であり検出されるクローンの数は少なくなるが、文の集合として検出されメソッドとして抽出することが容易であるという特徴がある [9]。本研究では Eclipse JDT を用いて抽象構文木解析を行いソースコード中のブロックを検出する。この時、return 文を含むコード片をメソッドとして抽出することは難しいとされているため、そのようなブロックは検出を行わない [10]。

識別子の正規化を行うことでより多くのクローンを検出できるため、検出したブロックに対して以下のルールで正規化を行う。

- 変数名は "\$" + "数字" で正規化する
- 同一の変数名は同一の名前で正規化する
- リテラルは全て "\$" で正規化する
- 修飾された変数名は一つの変数として正規化する
- クラス名は正規化しない
- メソッド名は正規化しない

同一の変数名を同一の名前で正規化することで誤検出を防ぐ。クラス名とメソッド名を正規化しない理由は、クラスとメソッドはメソッドの引数として与えることができないので、クローンとして検出されるのを防ぐためである。

正規化を行った後にブロックのハッシュ化を行う。ハッシュ化には MD5 [11] を利用する。MD5 は 128 ビットのハッシュ値を出力するため、ハッシュ値の衝突の可能性は十分に低いと考えられる。

最後にそれぞれのブロックのハッシュ値を比較し、同じ値のブロックをクローンとして検出する。

#### Step 2: ソースコード変更

Step 1 で検出されたクローンセットの一つに対してメソッド抽出を用いてソースコードに変更を加える。ソースコード

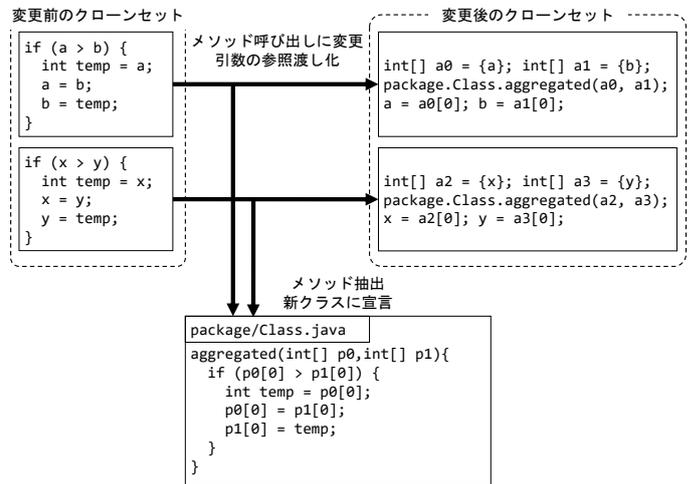


図 5 ソースコード変更の例

変更の例を図 5 に示す。クローンをメソッドに抽出することで、それぞれのクローンをメソッド呼び出しに変更でき、ソースコード行数を減らすことができる。抽出されたメソッドは 2.1 で作成したクラスに宣言する。このメソッドは完全修飾名で呼び出される。

また、抽出されたメソッド内で引数として与えられた変数に変更がある場合、メソッド呼び出し以降の外部的振る舞いに影響が出る可能性がある。そこで、抽出されたメソッドに参照渡しで引数を与えることで外部的振る舞いを保つようにする。引数を参照渡しにするために配列を利用する。メソッド呼び出し前にそのメソッド内で利用される変数の型の配列を定義する。これらの配列を各変数で初期化した後にそれらの配列をメソッドに引数として与える。抽出されたメソッド内では、引数として与えられた配列の添字 0 の要素を参照する。メソッド呼び出し後に元の変数に配列の添字 0 の要素を代入する。

#### Step 3: 変更したソースコードのコンパイルとテスト

Step 2 で変更した結果、行数が削減したソースコードに対してコンパイルとテストを行う。テストには 3.1 で生成したテストを用いる。初めにコンパイルを行い、コンパイルに成功すると次にテストを行う。コンパイルとテストのどちらにも成功した場合、変更されたソースコードを選択可能なソースコードとして記録し、ソースコードを変更を加える前の状態に戻し、次のクローンセットに対して Step 2 から繰り返す。

コンパイルあるいはテストのどちらかに失敗した場合、そのソースコードは記録せず、成功した場合と同様にソースコードを変更を加える前の状態に戻し、次のクローンセットに対して Step 2 から繰り返す。

#### Step 4: 変更したソースコードの選択

Step 1 で検出されたすべてのクローンセットに対して Step 2, Step 3 の処理を行った後に Step 4 に移る。Step 4 では、選択可能なソースコードの中から最も削減行数が大きいソースコードを選択する。選択されたソースコードに対して、Step

(注1) Java Development Tools <https://www.eclipse.org/jdt/>

1 から繰り返す。

## 4. 実験

本章では、提案手法を用いて行った実験と、その実験結果について述べる。

### 4.1 実験概要

オープンソースソフトウェア（以下、OSS）に対して提案手法を適用することで、削減可能なソースコード行数を算出する。今回は、3.2 で述べたクローン検出の際に正規化を行った場合と、正規化を行わなかった場合の二通りの実験を行った。また、算出した値に対して既存研究と削減行数の比較を行った。

### 4.2 実験対象

実験対象は Java 言語で記述された OSS である。本研究では、既存研究で対象となったプロジェクトの中から EvoSuite でのテスト生成に成功した jEdit, JFreeChart, JRuby の三つと、GitHub 上で公開されている fastjson, Checkstyle の二つを加えた合計五つのプロジェクトを対象として実験を行った。fastjson と Checkstyle は、GitHub 上で公開されており以下の三つの条件を満たすプロジェクトの中で、スター数が上位二つのプロジェクトである。

- プロジェクトが Apache Maven で管理されている
- 対象とするソースコードの行数が 50,000 行以上である
- EvoSuite でのテストの生成に成功する

Maven で管理されているプロジェクトを選んだ理由は、EvoSuite は Maven のプラグインとして提供されているため、テストの生成が容易であるからである。それに加えて、Maven を用いることでコンパイルとテストを容易に行うことができる。

対象としたプロジェクトの名前、バージョン及びソースコードの総行数を表 1 に示す。バージョンは、既存研究と同じプロジェクトに関しては同じバージョン、異なるプロジェクトは実験時の最新リリースのバージョンである。なお、プロジェクトのテストやチュートリアルなどのソースコードは実験対象

表 1 実験対象の OSS (\* 単位は K LoC)

プロジェクト名	バージョン	* 行数
jEdit	5.4.0	162
JFreeChart	1.0.19	236
JRuby	1.7.27	332
fastjson	1.2.54	50
Checkstyle	8.15	81

に含まない。ソースコードの行数は 3. で説明したフォーマッタを適用した後の行数である。

### 4.3 実験結果

実験結果を表 2 に示す。実験結果から分かるように、検出されたクローンセットの中で、集約することで行数が減り、コンパイルとテストを通過するクローンセットは全体の 5 - 10 % ほどとそれほど多くない。また、集約可能なクローンセット数と削減行数から、一つのクローンセットあたり 10 行程度行数が削減されていることが分かる。

次に既存研究との比較を行う。比較対象は jEdit, JFreeChart, JRuby の三つのプロジェクトである。既存研究と本研究の提案手法で算出した削減行数を表 3 に示す。なお、既存研究は対象のプロジェクトに対してフォーマッタを適用していないため、フォーマッタを適用していないプロジェクトに対して提案手法で削減行数を測定した結果を記載する。

表から分かるように、既存研究と本研究とでは得られた削減行数は大きく異なる。jEdit では既存研究と比べると提案手法がより多い削減行数を算出したが、jEdit 以外では既存研究と比較して非常に少ない削減行数を算出した。

## 5. 考察

本章では、4. で述べた実験の考察を行う。

### 5.1 実験結果に対する考察

実験の結果から、集約することで行数が減り、なおかつコンパイルとテストを通過するクローンセットは検出されたクローンセット全体の内 5 - 10 % ほどであることが分かった。コンパイルに失敗する主な原因として以下が挙げられる。

- クローン間で対応する変数の型が異なる
- super を使用して親クラスを参照している
- ブロック内で例外を throw しているがその例外の対応はブロック外で行っている
- ブロック内で break や continue しているが繰り返し処理はブロック外で行っている

これらの中には、提案手法の実装方法によってはコンパイ

表 3 算出された削減可能行数の比較結果 (\* 単位は LoC)

プロジェクト名	* 既存研究	* 提案手法
jEdit	136	362
JFreeChart	9,700	757
JRuby	2,161	222

表 2 実験結果 (\* 単位は LoC)

プロジェクト名	正規化を行った場合			正規化を行わなかった場合		
	検出されたクローンセット数	集約可能なクローンセット数	* 削減行数	検出されたクローンセット数	集約可能なクローンセット数	* 削減行数
jEdit	397	19	213	273	17	139
JFreeChart	872	43	377	607	24	217
JRuby	802	40	326	515	44	238
fastjson	491	44	600	189	6	400
Checkstyle	100	5	36	51	2	1

ルが可能であるものも存在する。例えば、super を使用して親クラスを参照しているコードを含むクローンが同じクラス内のみ存在した場合、抽出されたメソッドをそのクラスに宣言することでコンパイルを通過する可能性がある。

## 5.2 既存研究との比較に対する考察

既存研究との比較を行った結果、算出された削減可能行数が大きく異なるという結果になった。このような結果となった理由として以下の要素が考えられる。

- クローン検出法の違い
- コンパイル、テストの有無
- 削減行数の計測法

### クローン検出法の違い

本研究ではブロック単位でのクローン検出を行っているが、既存研究では字句単位でのクローン検出を行っている。字句単位でのクローン検出はブロック単位でのクローン検出より多くのクローンの検出が可能である。クローン検出法の違いから、集約対象となるクローンの数に差が出るのが考えられる。

### コンパイル、テストの有無

本研究では実際にクローンを集約したソースコードに対してコンパイルとテストを行い集約可能かを判定している。しかし、既存研究では行っていない。したがって、既存研究で集約可能と判断されたクローンセットは、実際は集約できない可能性があり、それが原因で削減行数に差が出ているのが考えられる。

### 削減行数の計測法

本研究では実際にクローンの集約を行い、集約の前後でのソースコード行数の変化から削減行数を測定している。しかし、既存研究では jEdit 以外のプロジェクトに対して実際の削減行数の測定は行っていない。したがって、実際の削減行数が既存研究の推定値とは異なっている可能性がある。

## 6. 妥当性への脅威

本章では、本研究に対する妥当性への脅威について述べる。本研究では5つのOSSを対象に実験を行った。しかし、他のプロジェクトに対して実験を行った場合、本研究で得られた結果とは異なる可能性がある。

本研究では、ハッシュ値の比較を行うことでクローンを検出している。ハッシュ値の衝突が発生した場合、誤ったクローンが検出される可能性がある。しかし、本研究では128ビットのハッシュ値を出力するMD5を用いており、ハッシュ値の衝突の可能性は十分に低い。

本研究では、ソースコードの変更前後で外部的振る舞いに変化していないことを確認するために EvoSuite によって自動生成された単体テストを用いた。しかし、自動生成されたテストが不十分で、外部的振る舞いの変化を正確に確認できていない可能性がある。

## 7. あとがき

本研究では、コードクローンに対して実際に集約を行い、そ

の結果からより正確な削減可能なソースコード行数を測定する手法を提案した。また、実験の結果、既存研究とは大きく異なる削減可能行数が算出されるという結果となった。

今後の課題としては次のようなものが考えられる。

### 複数のリファクタリング手法の採用

本研究ではクローンの集約手法としてメソッドの抽出のみを使用した。しかし、クローンに対するリファクタリング手法はメソッドの抽出だけでなく、多くの手法が提案されている。例えば、メソッドの引き上げやテンプレートメソッドの形成などが挙げられる。これらの手法を採用することで、より正確な削減行数の算出が可能になることが期待できる。

### 回避可能なコンパイルエラーを回避するための実装

5.1 で述べたように、本研究ではコンパイルエラーとなったが、実装を工夫することでコンパイルが可能になるようなクローンセットが存在すると考えられる。このようなクローンセットに合わせて実装を改善することで、より正確な削減行数の算出が可能になることが期待できる。

### 謝辞

本研究は、日本学術振興会科学研究費補助金基盤研究(B)(課題番号:17H01725)の助成を得て行われた。

### 文 献

- [1] B.P. Lientz and E.B. Swanson, Software Maintenance Management, Addison-Wesley Longman Publishing Co., Inc., 1980.
- [2] T.M. Pigoski, Practical Software Maintenance: Best Practices for Managing Your Software Investment, 1st edition, Wiley Publishing, 1996.
- [3] 肥後芳樹, 楠本真二, 井上克郎, “コードクローン検出とその関連技術,” 電子情報通信学会論文誌. D, 情報・システム = The IEICE transactions on information and systems (Japanese edition), vol.91, no.6, pp.1465–1481, June 2008.
- [4] M. Fowler, Refactoring: Improving the Design of Existing Code, Addison-Wesley Signature Series, Pearson Education, 1999.
- [5] 肥後芳樹, 吉田則裕, “コードクローンを対象としたリファクタリング,” コンピュータソフトウェア, vol.28, no.4, pp.4\_43–4\_56, 2011.
- [6] N. Yoshida, T. Ishizu, B. Edwards, III, and K. Inoue, “How slim will my system be?: Estimating refactored code size by merging clones,” Proceedings of the 26th Conference on Program Comprehension, pp.352–360, ICPC '18, 2018.
- [7] 石津卓也, 吉田則裕, 崔 恩瀾, 井上克郎, “コードクローンに対するリファクタリング可能性に基づいた削減可能ソースコード量の調査,” Technical Report 7, 大阪大学, 名古屋大学, 奈良先端科学技術大学院大学, 大阪大学, Nov. 2017.
- [8] G. Fraser and A. Arcuri, “Evosuite: Automatic test suite generation for object-oriented software,” Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering, pp.416–419, ESEC/FSE '11, 2011.
- [9] 堀田圭佑, 楊 嘉晨, 肥後芳樹, 楠本真二, “粗粒度なコードクローン検出手法の精度に関する調査,” 情報処理学会論文誌, vol.56, no.2, pp.580–592, Feb. 2015.
- [10] R. Komondoor and S. Horwitz, “Effective, automatic procedure extraction,” 11th IEEE International Workshop on Program Comprehension, pp.33–42, May 2003.
- [11] R. Rivest, “The md5 message-digest algorithm,” April 1992.