

平成 22 年 5 月 14 日現在

研究種目：若手研究(B)
研究期間：2008 ~ 2009
課題番号：20700024
研究課題名（和文） 超大規模ソフトウェアを対象としたコードクローン分析基盤環境
研究課題名（英文） Code-Clone Analysis Environment for Huge Software Repositories
研究代表者
松下 誠 (Makoto Matsushita)
大阪大学・大学院情報科学研究科・准教授
研究者番号：60304028

研究成果の概要（和文）: 本研究では、超大規模なソフトウェアを対象としたコードクローン分析を行うことを目的として、コードクローンを作成した開発者の情報を収集することによるクローン分析手法や、コードクローン分析手法を UML で記述された設計文書に対して適用することによる設計の再利用手法などを実現する、コードクローン分析のための基盤環境について研究を行い、ツールの実装を通じて手法の有効性を確認した。

研究成果の概要（英文）: In this research, we propose several methodologies on detecting huge software repositories, including a method that employs software developer analysis, and a method for detecting reusable design documents. Through experiments with prototype tools, we confirmed on its effectiveness.

交付決定額

（金額単位：円）

	直接経費	間接経費	合計
2008 年度	1,600,000	480,000	2,080,000
2009 年度	1,700,000	510,000	2,210,000
年度			
年度			
年度			
総計	3,300,000	990,000	4,290,000

研究分野：ソフトウェア工学

科研費の分科・細目：情報学・ソフトウェア

キーワード：コードクローン、リポジトリ分析

1. 研究開始当初の背景

ソフトウェアに対する分析手法として、近年、ソースコード中の重複した部分を発見するコードクローン検出技術が注目されてきている。コードクローンの分析を行うことにより、単にソースコードを共有している部分が発見できるだけでなく、ソフト

ウェアの再利用の程度や、あるいはソースコードの盗用といったソフトウェアのさまざまな側面を見ることが出来る。しかしながら、コードクローン分析は、特に対象となるソースコードが超大規模になった場合、アルゴリズムを工夫しても単一計算機上で分析することは非常に困難であり、そのた

めこれまで単一の大規模システム程度にし
か分析されてこなかった。

研究代表者はこれまで、コードクローン
や大規模ソフトウェアに対する検索技術等、
ソフトウェア工学の各領域において研究を
行ってきた。たとえば、Java 言語を対象と
して、複数のソフトウェアの集合である十
数万クラス規模の大規模ソースコードに
対し、実用レベルで利用可能なソースコード
検索システムに関する研究を行っており、
Java クラス間の依存関係解析の結果を用
いて利用頻度の高いクラスを優先した検索
結果の表示を行うための手法を確立してい
る。また、膨大なソフトウェアの中からコ
ードクローンを発見するためのツールや、
ツールが発見した結果を用いてコードク
ロンの性格等によって分類するための
環境、および分析するための手法を確立し
てきている。

一方、超分散型のプログラム実行環境と
して、グリッドコンピューティング等の環
境が注目を集めている。この種の環境は、
単一の処理を細かく分割し、大量の計算機
等を用いて演算を行うことで、大量の計算
を低いコストで実行することを目指してい
る。そこで、超分散型のプログラム実行
環境をソフトウェア分析を行うための実行
環境として用いることにより、さらに大規
模なソースコードを対象としたコードク
ローン分析環境を構築することが可能では
ないかと考えた。

2. 研究の目的

本研究は、これまでの成果を踏まえ、よ
り大規模なソースコード、具体的には、世
界中で広く用いられている全てのソース
コードを扱うことができる程の規模を対象
としたソースコード分析基盤を構築するこ
とを目指している。また、構築したコード
クローン分析基盤を運用することにより、
世界中のソフトウェアがどのような性格を
持っているのか、全体像を明らかにするこ
とを目指す。また、コードクローン分析基
盤を他のソースコード解析技術へ応用する
可能性についても検討を行う。

そのために、超大規模なソフトウェアを
対象として、コードクローンの分析を行う
ための分析基盤を構築し、オープンソース

ソフトウェア等のインターネット上に存在
する多種多様なソフトウェアを対象として
分析を行う。これにより、世の中に存在す
る全てのソフトウェアが、どの程度再利用
を行い、あるいはソースコードの単純な流
用を行っており、またそれらの結果と、一
般的に言われているソフトウェアの品質が
どのような関係となっているか、を明らか
にすることを目標とする。

3. 研究の方法

本研究ではまず、ソフトウェア分析を行う
ためのモデルについて検討を行う。大規模
かつ複数のソフトウェアを同時に取り扱う
ことを前提として、コードクローン分析を
分散型で行うためのモデルの構築を行う。
本モデルは、以下のような点に考慮して設計する。

- 入力となるソフトウェアの大きさ
には依存しないこと。

本システムでは、数億行あるいは数
GB といった非常に大規模で、かつ複
数のソフトウェアを分析対象とす
ることを目指す。

- 分析結果が容易に統合できること。

本システムでは、分析する対象およ
び内容を、細かな問題に分解し、分
析を行った後その結果を結合する、
といった一般的な方法を用いる。そ
のため、個別の分析時における入出
力情報、および分析内容が分割可能
なようにすることはもちろん、その
結果を集めて1つの分析結果とす
ることが可能でなければならない。

- 可搬性が高いこと。

分析システムが動作する環境や、分
析内容は状況に応じて変化するこ
とが考えられる。そのため、本分析
モデルでは、動作環境や分析内容に
依存しない形でソースコードおよ
び分析内容を扱えるようにしなけ
ればならない。

次に、前述のモデルを用いることを前提と
し、コードクローン分析を分散型ソフトウ
ェア分析モデル上で行うためのアルゴリ
ズムを決定する。既存のアルゴリズムは、単一
のシステムが全てローカルに処理を行うこ
とを前提としていたため、これを分散環境
でも用いることが出来るように拡張する。

ソースコードの大きさが小さいときには、既存のコードクローン検出アルゴリズムを用いても十分高速であるため、まず入力となる大規模なソースコードを小さなソースコードの集合として捉え、それぞれを既存のアルゴリズムで分析し、その結果を統合して全体の結果とすることを考える。

以上の結果を踏まえ、特に分散型ソフトウェア分析モデルがどの程度の性能を持つかを確認するために、システムの試作を行う。分析システムは超大規模なソースコードを入力とし、コードクローンの分布状況を視覚的に表現するものとする。

4. 研究成果

2年間の研究期間の間、コードクローン分析基盤環境の構築を目指し研究を行った結果、以下のような成果を得た。

(1) 2008年度の成果

ソフトウェアの保守性を悪化させる問題の1つとして、コードクローンの存在があげられる。しかし近年の研究では、コードクローンはすぐに消滅するものが多く、また意図してコードクローンとなるように作られているものもあり、ある時点で発見されたすべてのコードクローンが不適切とは限らないことがわかってきた。そのため、ある時点で存在するコードクローンを適切に管理することによって、問題のあるコードクローンを確実に発見することが重要であるとされている。しかしながら、問題を抱えているコードクローンを与えられたソースコードのみから発見することは困難である。

そこで、コードクローン分析と同時に、そのコードクローンを作成した開発者の情報を収集することによって、各開発者がどのようにコードクローンを作成、あるいは解消しているか、分析するための手法を提案した。また、オープンソースソフトウェアの1つである PostgreSQL のソフトウェア開発を対象として、提案手法を用いた分析を行った。その結果、開発者の中には、有意にコードクローン数を増やす、あるいは減らす者がいることを、Kruskal-Wallis 検定を用いて定量的に示すことができた。本研究により、どの開発者が作成したかどうかを管理することにより、コードクローンの発生を予想することができると期待される。

(2) 2009年度の成果

ソフトウェアの再利用を適切に行うことによって、ソフトウェア開発の生産性と成果物の信頼性の双方が向上するといわれている。再利用を行う対象は、大きく設計情報とソフトウェア部品に分けることが出来る。再利用コンポーネントやクラスライブラリなどにより、ソフトウェア部品の再利用は積極的に行われているが、設計情報を再利用する試みはあまり行われていない。設計情報の再利用を積極的に行うことを目的として、プロダクトライン開発と呼ばれる手法があるが、これは現在行われている開発の設計情報を、派生して行う開発に適用するものであり、過去すでに行われた設計情報を再利用することは想定されていない。一方、ソフトウェア部品の再利用を促進することを目指し、再利用可能な部品を開発中に適宜推薦することができる、部品の自動推薦手法に関する研究が行われてきている。

そこで設計情報の中でもソフトウェア設計の際によく用いられる UML 図に着目し、UML 図を再利用する際に用いることができる自動推薦ツールを作成した。本ツールは、既存 UML 図の作成ツールに UML 図の自動推薦機能を組み込むことによって実現しており、現在作成中の UML 図中に記述された文字情報を元に、それと類似した情報を持つ UML 図を既存の UML 図群を対象として検索を行い、検索結果を適宜別ウィンドウへ表示する。また、検索結果を選択することによって、現在編集中の UML 図へ、選択した UML 図あるいはその一部を取り込むことができる。本ツールを用いた実験を行った結果、全体の1割程度の図を記述した状態で検索を行うと、再利用可能な類似 UML 図を高い確率で推薦できることがわかった。これにより、UML 図を作成する際に既存の UML 図の再利用を行いやすくなると期待される。

5. 主な発表論文等

(研究代表者、研究分担者及び連携研究者には下線)

[学会発表](計4件)

[1] 松下 誠、早瀬 康裕、松下 誠、井上 克郎: "状況に応じた設計情報の再利用を支援する UML 図の自動推薦ツール"、信学技報、vol. 109、No.456、SS2009-54、pp.37-42 (2010年3月8日、鹿児島大学、鹿児島、日本)。

[2] 松下 誠: "ソフトウェア工学の新潮流 (1)リポジトリマイニング"、ソフトウェアエンジニアリング最前線 2009、pp.21-24 (2009年9月9日、東京女子大学、東京、日

本).

[3] Yue Jia, David Binkley, Mark Harman, Jens Krinke, Makoto Matsushita: "KClone: A Proposed Approach to Fast Precise Code Clone Detection", Proc. of Third International Conference on Software Clones (2009年3月24日, Fraunhofer IESE, Kaiserslautern, Germany).

[4] 東 誠, 肥後 芳樹, 早瀬 康裕, 松下 誠, 井上 克郎: "コードクローンの複雑度メトリクスを用いた開発者の特徴分析", ソフトウェアエンジニアリング最前線 2008, pp.103-106 (2008年9月2日, 東洋大学, 東京都, 日本).

6 . 研究組織

(1)研究代表者

松下 誠 (Makoto Matsushita)
大阪大学・大学院情報科学研究科・准教授
研究者番号：60304028