

修士学位論文

題目

ソースファイル群の類似性を用いたソフトウェア再利用元の推定

指導教員

井上 克郎 教授

報告者

坂口 雄亮

平成 29 年 2 月 7 日

大阪大学 大学院情報科学研究科
コンピュータサイエンス専攻 ソフトウェア工学講座

内容梗概

ソフトウェア開発において、既存ソフトウェアの再利用が盛んに行われている。再利用の方法の1つとして、既存ソフトウェアのソースコードを開発中のソフトウェアにコピーして取り込み、開発中ソフトウェアでの用途に合わせて変更を加える場合がある。既存ソフトウェアに不具合が見つかり修正された場合、再利用したソフトウェアのコピーを更新しなければいけない。しかし、既存ソフトウェアのソースコードを独自に変更している場合、更新を取り込むと変更箇所が競合する可能性があるだけでなく、再利用したソフトウェアの記録がなされていないことが多くある。そのため、既存ソフトウェアを安全に更新するためには、再利用されたソフトウェアとそのバージョンを確認する必要がある。

本研究では、開発中ソフトウェアが再利用するソフトウェアとそのバージョンを特定する手法を提案する。ソフトウェアはソースファイルの集合とみなし、その各ソースファイルについて類似するソースファイルを既存の高速検索手法を用いて、大量のソフトウェアのソースファイルから検索する。そして、与えられたソフトウェアがもつソースファイルの集合と類似度がより高いものが上位になるように、検索結果から得られたソフトウェア群の順位付けを行う。この操作によって、再利用したソフトウェアとそのバージョンを知ることが可能になる。

実験として、Mozilla Firefox と Android が再利用するソフトウェア 72 件に対して手法を適用した。その結果、手法が出力するソフトウェア集合の中に再利用元が含まれるものが 71 件、その集合を絞り込んだ集合の中に再利用元が含まれるものが 66 件であった。その 66 件について、再利用元の順位が 1 位であるものは 50 件であった。

主な用語

ソフトウェア再利用
オープンソースソフトウェア
ソースコード比較

目次

1	はじめに	4
2	背景	6
2.1	ソースコードの再利用	6
2.2	オープンソースソフトウェアの管理状況	6
2.3	起源分析	7
2.4	LSH アルゴリズムを利用した高速検索手法	7
2.4.1	ソースファイル間の類似度	8
2.4.2	Locality-Sensitive Hashing	8
2.4.3	実装	9
2.4.4	性能	9
3	提案手法	10
3.1	類似ソースファイルの検索	11
3.2	候補ソフトウェア集合の取得	11
3.3	候補ソフトウェアの順位付け	12
3.3.1	順序関係	13
3.3.2	有力ソフトウェアの順位付け	13
3.3.3	トポロジカルソート	14
4	実験	17
4.1	データベース	17
4.2	評価方法	18
4.3	対象プロジェクト	18
4.3.1	Mozilla Firefox	18
4.3.2	Android	20
4.4	手法出力と正解ソフトウェアに対する評価	20
4.5	距離関数と順位の評価	23
4.6	順序関係の評価	23
5	妥当性への脅威	30
6	まとめと今後の課題	31
	謝辞	32

1 はじめに

ソフトウェア開発において、既存のソースコードを再利用することが盛んに行われており、ソフトウェアに必要な機能を開発するコストの削減に役立っている。再利用可能なソースコードとしてはオープンソースソフトウェアが活用されている。コードを再利用することによって、機能を開発するコストを大幅に削減することが出来る。オープンソースソフトウェアは、多くのユーザによって使用されているため、安全性が比較的高いというメリットがあることから、企業の製品開発にも利用されるようになってきている [2], [3].

ソースコードの再利用の行われ方は様々あるが、C言語でのソフトウェア開発では、既存のソフトウェアのソースコードをコピーして取り込むという方法で再利用が行われる。コピーして取り込むと、開発中のソフトウェアに合わせた変更が容易になるというメリットがある一方、再利用したソースコードに含まれる脆弱性を取り込んでしまうおそれがある。そのような場合、ソフトウェアの更新を速やかに行わなければいけないが、更新された部分と独自に変更した部分が競合する可能性があるため、再利用しているソフトウェアとバージョンの情報が必要になる。しかし、多くのプロジェクトにおいてそのような情報が記録されていないことが判明している [14].

ソースファイルが一致するかどうかの単純な判定は容易であるが、実際には再利用後にソースコードの変更が加えられることがあるため不十分であり、ソースファイルの内容によって再利用元を検索する必要がある。ソースファイルに基づいて再利用元を検索する既存手法として、Kawamitsuら [9] は、リポジトリ内の最も類似したファイルのリビジョンを再利用元のバージョンとして識別するコード比較手法を提案した。この手法は、あらかじめ再利用元ライブラリのリポジトリを用意しなければならない。また彼らは、locality-sensitive hashing[4]を用いた類似ファイルの高速検索手法を提案した [15]。これは、検索のクエリとして与えられたソースコードのファイルに対し、大量のソースコードのファイルの中からソースコードの内容がクエリに類似するものを高速に検索する。1つのファイルについての正確性や有効性について述べられているが、入力として複数ファイルを与え、それらの情報を合わせて利用することが課題とされている。

本研究では、複数ファイルの検索結果であるソースファイル群を利用して、開発中ソフトウェアが再利用しているソフトウェアの再利用元を自動的に推定する手法を提案する。提案手法では、開発中ソフトウェアのソースファイルを入力とする。まず、各ソースファイルについて、locality-sensitive hashing を用いた検索手法を適用し、大量のソフトウェアのソースファイルから成るデータベースから類似するファイルを検索する。データベースは、フリーなオペレーションシステムである Debian GNU/Linux が採用してきたほぼすべてのパッケージで構築されている。検索後、各入力ソースファイルの類似ファイル群をデータベース

に登録されているソフトウェア毎にまとめ、ソフトウェア集合を得る。そして、入力ソースファイル集合により近いものが上位となるようなソフトウェア間の半順序関係を用いて順位付けられた再利用元ソフトウェアリストを得る。

手法の評価のために、Mozilla Firefox と Android のソースコードに再利用されているソフトウェアを対象として手法の妥当性を評価した。

以降、第2章では、本研究の背景について述べる。第3章で提案手法について説明し、第4章で手法の評価を述べる。第5章では妥当性への脅威について述べ、最後に第6章で、まとめと今後の課題を記述する。

2 背景

2.1 ソースコードの再利用

ソフトウェア開発者は、既存のソフトウェアのソースコードを再利用してソフトウェアを開発する [2]。ソースコードの再利用は、機能を開発するコストを大幅に抑えるため、ソフトウェア開発の効率化につながる。Rubin ら [12] は、企業内での新規ソフトウェア開発において、既存のソフトウェアの再利用が行われていることを報告している。Mohagheghi ら [11] は、再利用されたコンポーネントはそうでないものと比べて安全性が高いと報告している。

再利用のされ方として、ソフトウェアの機能を開発する際に、新しく開発するのではなく、その機能を持つオープンソースソフトウェア (以下 OSS とも表記) が一般的に利用される。その際に、再利用したコードを開発中のソフトウェアに合うように独自の変更が加えられることがある。Java プロジェクトでの OSS の利用は、バイナリファイルの再利用が一般的である。また、新たにソフトウェアを開発する際には、既存ソフトウェアを全て再利用し、細かな改変が加えられ類似したソフトウェアが開発する派生開発、ソフトウェアプロダクトで全体に共通するコア資産と独自性を持つ機能部品に分割し、これらを組み合わせてソフトウェアが開発するソフトウェアプロダクトラインエンジニアリング [8] もソースコードの再利用と言える。

2.2 オープンソースソフトウェアの管理状況

OSS は安全性が高いとは言え、再利用したソースコードにバグやセキュリティに関する脆弱性が発見されることもある。脆弱性が公開されると開発者は、脆弱性の含まれたコードを再利用しているのかを確認する必要がある。再利用している場合、開発者のソフトウェアにも脆弱性が含まれることになり、それに対する修正を手元のソースコードにも適用すべきである。脆弱性が含まれたコードを再利用しているのかを確認するためには、再利用しているソフトウェアを管理する必要がある。

しかし、多くのプロジェクトにおいて、再利用しているオープンソースソフトウェアの管理がなされていないことがわかっている。Xia ら [14] は、オープンソースライブラリを再利用しているプロジェクトについて、そのライブラリが脆弱性の含まれるソースコードを保有しているバージョンであるかどうかを調査した。その結果、zlib について、45 プロジェクト中 14 プロジェクト (31.1%) が脆弱性のあるバージョンを使用していることがわかった。同様に、libcurl については、28 プロジェクト中 24 プロジェクト (85.7%)、libpng については、50 プロジェクト中 46 プロジェクト (92%) であった。また、計 123 プロジェクト中 27 プロジェクト (22%) で、ライブラリのソースコードを再利用後に編集しており、23 プロジェクト (18.7%) では、再利用しているライブラリのバージョンに関する情報が残っていなかつ

た。さらに、6プロジェクト(4.9%)で、ディレクトリ名や他のライブラリのソースコードが含まれており、管理が容易ではなかった。このように、再利用しているソフトウェアの管理について、そのソフトウェア名やバージョンなどの再利用元を推定することは有用であると考えられる。

2.3 起源分析

本研究では、再利用元の特定に関する研究として、起源分析が上げられる。起源分析とは、ソースコードの起源を複数のソフトウェアから特定する技術で、ソースコードの類似性を用いるため、コードクローン検出の一部と言える。コードクローンとは、ソースコード中に存在する同一または類似した部分を持つコード片のことであり、既存コードのコピーアンドペーストによる再利用等が原因で生じる。コードクローン検出技術の例として、Kamiyaら[7]は、ファイル間の類似したコード片を検出するために、CCFinderを提案した。また、Sasakiら[13]が提案したFCFinderは、プロジェクト間でのファイルクローンを高速に検出する。

起源分析では、Inoueら[5]は、Ichi trackerという名前のツールを提案した。これは、インターネット上の様々なリポジトリ間で、特定のソースファイルの類似ファイルを検索し、コード断片との類似度でクラスタリングを行う。そして、ソースファイルを時系列順に並べることによりソースコードの再利用の経緯を可視化する。このツールでは、1つのファイルについて注目している。Kawamitsuら[9]は、リポジトリに含まれているソースコードについて、ライブラリのバージョンを推定する手法を提案した。ソースコードの類似度として最長共通部分列に基づいた類似度[8]を用いて、最も類似度の高いソースコードが再利用元であるという仮定に基づきバージョンを提示した。この手法はあらかじめ再利用を行ったソフトウェアのリポジトリと再利用元のリポジトリが必要なため、再利用元が不明である場合は利用することが出来ない。Javaソフトウェアの場合、Ishioら[6]は、ソフトウェアに含まれるjarファイルを分析し再利用しているライブラリを自動的に特定した。

2.4 LSH アルゴリズムを利用した高速検索手法

本研究では、類似ソースファイルの検出に、川満ら[15]のlocality-sensitive hashing(以下LSHとも表記)アルゴリズムを利用した高速検索手法を用いる。この手法では、検索のクエリとして与えられたソースファイルに対し、データベース中のソースファイルからソースコードの内容がクエリに類似するものを検索する。LSHを用いることにより、実際の類似度ではなく、類似度の推定値を計算することで高速な検索を可能にしている。本節では、この手法におけるソースファイル間の類似度の定義について述べ、次に、それに対応するLSH

である MinHash について述べる。そして、性能について簡潔に説明する。

2.4.1 ソースファイル間の類似度

ソースファイル f_1, f_2 の類似度は、各ファイルの 3-gram 集合間の Jaccard 係数を利用し以下のように定義される。Jaccard 係数は、2つの集合間の類似度であり、文字列間の類似度を計算する手法としてよく利用される。

$$\text{sim}(f_1, f_2) = \frac{|3\text{-grams}(f_1) \cap 3\text{-grams}(f_2)|}{|3\text{-grams}(f_1) \cup 3\text{-grams}(f_2)|}$$

ただし、 $3\text{-gram}(f)$ は、ソースファイル f から抽出された 3-gram の多重集合である。多重集合の要素は、ソースコードからコメント、空白行を除いたものから得た字句列から長さ 3 の部分文字列である。たとえば、字句列 $ABABA$ からなるファイルから得られる 3-gram の多重集合は $\{\$A_1, \$AB_1, ABA_1, BAB_1, ABA_2, BA\$_1, A\$_1\}$ となる。ここで、各要素の添え字は、同じ 3-gram の複数回の出現を区別するためのもので、多重集合を通常の集合の表現で扱うためのものである。

本手法においても、ソースファイル間の類似度を同様に定義する。

2.4.2 Locality-Sensitive Hashing

クエリのソースファイルとデータベース中の大量のソースファイルの間の Jaccard 係数を求めるには、膨大な時間がかかってしまう。川満らの手法では、locality-Sensitive Hashing を用いている。LSH とは、近似最近傍検索などに利用されるアルゴリズムである [4]。

LSH は、ハッシュ値の値が等しくなる確率が類似度と等しくなるようなハッシュ関数族 F を利用し、以下の式で表される [1]。

$$\Pr_{h \in F}[h(x) = h(y)] = \text{sim}(x, y)$$

ただし、 $\text{sim}(x, y)$ は、 x, y 間の類似度であり、 $0 \leq \text{sim}(x, y) \leq 1$ を満たす。類似度 sim として、Jaccard 係数を用いる場合、関数族 F として、MinHash が利用できる。

MinHash とは、2 値ベクトルとして表現された集合に対する確率的なハッシュ関数である。集合の各要素をハッシュ関数を利用してハッシュ値を求める。そして、その最小値を記録し、別の集合に対しても同様にハッシュ値の最小値を求め、それぞれの最小値が一致する確率は Jaccard 係数に一致する。この性質より、Jaccard 係数は、ランダムなハッシュ関数による最小値が一致しているかどうかの確率から求めることができ、厳密性は失われるが、計算量を大きく減らすことが可能となる。

2.4.3 実装

川満らの手法では、文献 [10] の方法を用いて LSH を構成した。各ソースファイルについて、3-gram の多重集合に $r \times b$ 個のハッシュ関数を用いて、 r 次元の MinHash の値からなるベクトルを b 個求める。そして、クエリのソースファイルとデータベース中のソースファイル間の対応するベクトルを比較し、ベクトルが 1 つでも一致したものを出力する。この時、類似度を p とすると、 b 個のベクトルのうち 1 つ以上のベクトルが一致する確率は、 $1 - (1 - p^r)^b$ となる。パラメータ b, r を調整することで、出力するソースファイルを調整することが出来る。類似度 p に対する最尤推定量 \hat{p} は、ベクトルが一致した個数を x とすると、

$$\hat{p} = \sqrt[r]{\frac{x}{b}}$$

となる。また、類似度 p がある程度低いと、 $x = 0$ となり、類似度の推定値が 0 になる特性を持つ。川満らの手法では、64bit 長のハッシュ値を使用し、LSH のパラメータは $b = 120, r = 8$ を利用した。その時の $1 - (1 - p^r)^b$ のグラフを図 1 に示す。本手法においても、同様のパラメータを利用する。

2.4.4 性能

検索結果および類似度の推定値に関する評価では、類似度が高いものは検索結果に表れ、低いものは表れないことを示した。

ケーススタディとして、あるプロジェクトが再利用している OSS のうち使用しているバージョンが記録されている 197 ファイルについて手法を適用した。データベースには、そのライブラリを含む 200 プロジェクトから合計 567113 ファイル登録されている。その結果、全 197 件について、記録されている再利用元を検索することができ、そのうち 185 件については類似度の推定値が、検索結果中の推定値の最高値となった。さらに、検索時間は 1 件あたり 1 秒以内であった。

川満らの手法では、1 つのソースファイルについて、再利用元の特定についての有効性があるといえる。今後の課題として、複数ファイルの検索結果を利用することで、再利用元の特定の精度向上が述べられている。

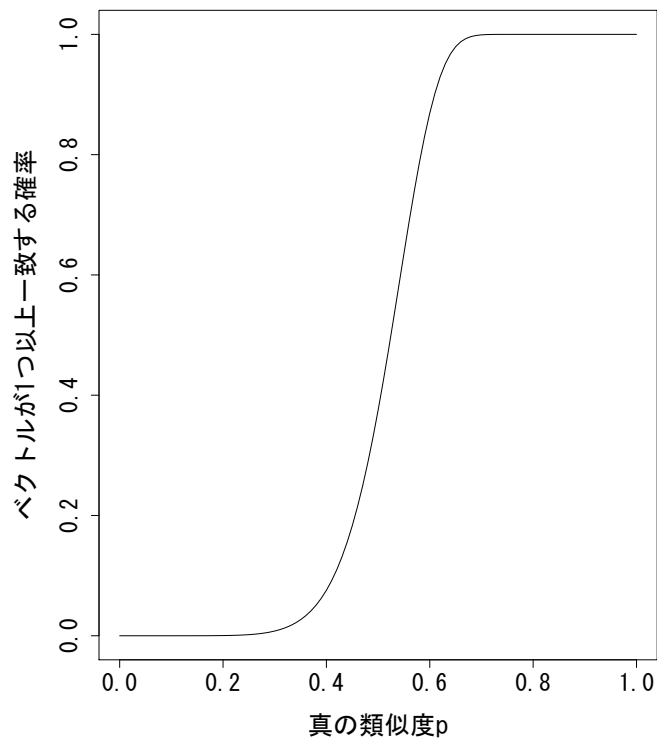


図 1: LSH のパラメータ

3 提案手法

本研究では、開発中ソフトウェアが再利用しているソフトウェアの再利用元を推定する手法を提案する。既存の様々なソフトウェアが蓄えられたデータベースが存在することを前提とし、入力として開発中のソフトウェアのソースファイル集合が与えられると、その入力ソースファイル集合と類似するソースファイルを持つソフトウェア群をデータベースから抽出し、類似度の順に並べ替えたリストを出力する。また、出力されたソフトウェアを利用者が分析するための情報として、ソフトウェア名やバージョン番号、類似ソースファイル集合もあわせて表示する。類似ソースファイル集合とは、類似ソースファイルのファイル名と類似度の集合である。リストのソフトウェアは類似する順に並んでいるため、開発中ソフトウェアが再利用しているソフトウェアの再利用元を推定することが可能となる。

提案手法は、以下の3つのステップから構成され、以降各ステップについての詳細を説明する。

1. 開発中ソフトウェアのソースファイル集合を入力し、各ソースファイルをクエリとし

て川満らの手法 [15] を用いて，類似するソースファイル集合を得る．

2. 類似ソースファイル集合を，それらを保有するソフトウェア毎にまとめ再利用元候補ソフトウェア集合を得る．
3. 再利用元候補ソフトウェア集合の各ソフトウェアに順位付けを行い，再利用元ソフトウェアリストを出力する．

3.1 類似ソースファイルの検索

このステップでは，入力されたソースファイルと類似するソースファイル集合を川満らの手法を用いて得る．

入力された開発中ソフトウェアのソースファイル集合の拡張子を調べて，検索対象となるソースファイル集合 $F = \{f_1, f_2, \dots, f_n\}$ を得る．そして，川満らの手法を用いて，入力ソースファイル集合 F 中の各ソースファイルをクエリとして検索を行う．入力ソースファイル集合 F 中のあるファイル f_i に対し，類似度の推定値 \hat{sim} が 0 より大きい，すなわち，真の類似度 sim がある程度大きい類似ソースファイル集合 $\tilde{F}(f_i)$ を抽出する．

本手法では，検索を行い類似するソースファイルを得たあとに，真の類似度 sim を計算する．まず，入力ソースファイルと類似ソースファイルの SHA-1 ハッシュ値を比較する．ハッシュ値が一致する場合は，ソースファイルの内容が完全一致しているので真の類似度 $sim = 1.0$ となる．一致しない場合，2.4.1 項の手順で真の類似度 sim を得る．本研究では，真の類似度 sim による閾値を 0.8 とし，それ未満のものは類似ソースファイルとみなさず集合 $\tilde{F}(f_i)$ に含めない．また，あるソースファイル f_i の検索結果において，同じソフトウェアが持つ類似ソースファイルが複数ある場合，類似度が最も高いソースファイルをそのソフトウェアから再利用したソースファイルとみなして $\tilde{F}(f_i)$ に含める．同一の類似度のソースファイルが複数ある場合は，ソースファイル名がアルファベット順で先頭のを $\tilde{F}(f_i)$ に含める．

3.2 候補ソフトウェア集合の取得

ステップ 2 では，ステップ 1 で得られた類似ソースファイルの集合から，(再利用元) 候補ソフトウェア集合を得る．

抽出された各類似ソースファイルは，それを保有するソフトウェアがただ 1 つ存在する．本手法ではそのようなソフトウェアを (再利用元) 候補ソフトウェアと呼ぶ．集合 $\tilde{F}(f_i)$ の各類似ソースファイルから候補ソフトウェア集合 $SW = \{S_1, S_2, \dots, S_m\}$ を得る．

集合 $\tilde{F}(f_i)$ 内のある類似ソースファイル f からは，入力ソースファイル f_i と類似度 $sim(f_i, f)$ で類似しており，候補ソフトウェア S_j が保有しているという 3 つの情報が得られる．そこ

表 1: 手法行列の例

入力 ソースファイル	候補 ソフトウェア	S_1	S_2	S_3	S_4	S_5	S_6
	f_1		0.8	0	0.95	0.85	0
f_2		0.8	0	0.9	0.85	0	0.95
f_3		0	0	1.0	0.9	0	0.9
f_4		0	1.0	0	0	1.0	0
f_5		0	0.9	0	0	0.8	0

で、入力ソースファイル f_i を行、候補ソフトウェア S_j を列、成分を類似度 $sim(f_i, f)$ とするような $n \times m$ 行列 M を作る。行番号 i は、入力ソースファイル f_i の番号、列番号 j は、候補ソフトウェア S_j の番号である。ただし、 $f \in \tilde{F}(f_i) \wedge f \in S_j$ となるような類似ソースファイル f がない場合、行列の成分となる類似度を 0 とする。この行列から各候補ソフトウェア S_j は、行列の j 列目の列ベクトルを得て、 $S_j = (M_{1j}, M_{2j}, \dots, M_{nj})$ と表す。表 1 に、入力ソースファイル 5 つ、候補ソフトウェア 6 つで構成される行列の例を示す。入力ソースファイル f_1 の類似ソースファイルのうち、候補ソフトウェア S_1 が保有するものの類似度は 0.8 と読むことができる。

3.3 候補ソフトウェアの順位付け

ステップ 3 では、ステップ 2 で得られる候補ソフトウェア集合について、再利用元ソフトウェアである可能性であるように順位付けを行い出力する。

類似するソースファイルを 1 つでも持てば候補ソフトウェアとなる。そのため、候補ソフトウェア集合 SW は、非常に大きくなる可能性があり、そこから再利用元ソフトウェアを見つけ出すことは困難である。そこで、本手法では以下の手順により、候補ソフトウェアを絞り込み集合 F に類似しているように順位付けを行う。

1. 候補ソフトウェア集合 SW に順序関係を定義し、半順序集合を得る
2. 有力ソフトウェアを抽出し、集合 F との距離を設定し順位付ける
3. その他の候補ソフトウェアをトポロジカルソートをして順位付けしリストを出力する

3.3.1 順序関係

候補ソフトウェア集合 SW に以下のような順序関係を定義し、候補ソフトウェア間に順序を付け、候補ソフトウェアの半順序集合 SW_{po} を得る。

定義. 候補ソフトウェア S_a, S_b において、すべての入力ソースファイルについて S_a よりも S_b のほうが類似度が高いソースファイルを保有するとき、すなわち任意の i について $S_a[i] \leq S_b[i]$ が成り立つとき、 S_b は S_a よりも、再利用しているソフトウェアの再利用元である可能性が高いとし、 $S_a \leq S_b$ と表す。

半順序集合 SW_{po} は、定義よりソースファイルの類似度が1つでも異なれば比較を行うことができず、2つの候補ソフトウェア間に順序はつかない。これにより、ある候補ソフトウェアが他の候補ソフトウェアよりも保有する類似ソースファイル数が少なくても、類似度が高いソースファイルをもてば、順序が下になることはない。入力ソースファイル集合 F に複数のソフトウェアのソースファイルが含まれており、各ソフトウェアのソースファイル数が異なる場合でも、それらに対応する候補ソフトウェア間には順序はつかない。

表1から得られる候補ソフトウェア集合 $S = \{S_1, S_2, \dots, S_6\}$ に、順序関係を定義すると、 $S_1 \leq S_1, S_1 \leq S_3, S_1 \leq S_4, S_1 \leq S_6, S_2 \leq S_2, S_3 \leq S_3, S_4 \leq S_3, S_4 \leq S_4, S_4 \leq S_6, S_5 \leq S_2, S_5 \leq S_5, S_6 \leq S_6$ という順序関係を持つ候補ソフトウェアの半順序集合 SW_{po} が得られる。

3.3.2 有力ソフトウェアの順位付け

候補ソフトウェアの半順序集合 SW_{po} 内のどの候補ソフトウェアよりも小さくない候補ソフトウェアは、他の候補ソフトウェアよりも再利用元ソフトウェアである可能性が高いと言える。本手法では、このような候補ソフトウェアを「有力ソフトウェア」と呼ぶ。定義した順序関係では、有力ソフトウェア間の順序はつかない。そのため、有力ソフトウェア S と入力ソースファイル集合 F 間の距離を距離関数を使って求め、その長さで順位を付ける。本研究では距離関数として、ユークリッド距離、マンハッタン距離、キャンベラ距離を採用する。実験によって、どの距離関数が良いかを調べる。各距離関数は、以下のように定義される。

ユークリッド距離

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

マンハッタン距離

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

キャンベラ距離

$$d(x, y) = \sum_{i=1}^n \frac{|x_i - y_i|}{|x_i| + |y_i|}$$

入力ソースファイル集合 F の各要素はソースファイルであるため、類似度ベクトル F_{sim} を求める。入力ソースファイル f_i 自身との類似度であるため、類似度 $sim(f_i, f_i) = 1$ である。よって、 F_{sim} の全てのベクトルの大きさは1である。有力ソフトウェア S と入力ソースファイル集合 F 間の距離 $d(S_j, F_{sim})$ は、上の式に当てはめると、

ユークリッド距離

$$d(S_j, F_{sim}) = \sqrt{\sum_{i=1}^n (S_j(i) - 1)^2}$$

マンハッタン距離

$$d(S_j, F_{sim}) = \sum_{i=1}^n |S_j(i) - 1|$$

キャンベラ距離

$$d(S_j, F_{sim}) = \sum_{i=1}^n \frac{|S_j(i) - 1|}{S_j(i) + 1}$$

となる。距離関数からわかるように、入力ソースファイルとの類似度が高いほど、また、類似ソースファイルを多く持つほど、距離 $d(S_j, F)$ は短くなる。よって、有力ソフトウェアにおいては、距離が短い順に順位を付ける。ただし距離が同じ場合は、入力によるものとする。

例の半順序集合 SW_{po} における有力ソフトウェアは、 S_2, S_3, S_6 である。各有力ソフトウェアと集合 F との距離は、距離関数としてユークリッド距離を用いた場合、それぞれ約 1.73, 約 1.41, 約 1.43 となる。有力ソフトウェアを距離が短い S_3, S_6, S_2 の順に 1 位, 2 位, 3 位と順位を付ける。 k 位である候補ソフトウェアが k 列目となるように列を入れ替えた行列を表 2 に示す。順位が確定した有力ソフトウェア名には○印を付けており、残りの候補ソフトウェアは順位がついていない。

3.3.3 トポロジカルソート

本手法では、ソフトウェアの再利用元を順位付けて出力することを目的としている。再利用元が有力ソフトウェアでない可能性もあるため、残りの候補ソフトウェアも順位をつける。残りの候補ソフトウェアについては、半順序集合 SW_{po} の各要素を半順序関係に矛盾しないように並べ替えて順位付ける。そこで本手法では、トポロジカルソートを利用する。トポロジカルソートとは、グラフ理論において、有向非巡回グラフ (DAG) の有向辺の情報を満たすようにグラフを一行に並べることである。順序関係のない 2 つのデータはどちらを先にし

表 2: 有力ソフトウェアの順位付け後の行列

入力 ソースファイル	候補 ソフトウェア	S_3	S_6	S_2	S_1	S_4	S_5
	f_1		0.95	0.85	0	0.8	0.85
f_2		0.9	0.95	0	0.8	0.85	0
f_3		1.0	0.9	0	0	0.9	0
f_4		0	0	1.0	0	0	1.0
f_5		0	0	0.9	0	0	0.8

てもよく，データの並べ方に自由度があるため，トポロジカルソートには一般にたくさんの解がありうる．

半順序集合は， $S_a \leq S_b$ のとき， S_a から S_b に向かう辺をつくることで有向グラフとみなすことができる．ただし，反射律 $S_a \leq S_a$ と，反対称律 $S_a \leq S_b$ かつ $S_b \leq S_a$ を満たすような辺はつukらない．また， $S_a \leq S_b$ かつ $S_b \leq S_c$ のとき推移律である $S_a \leq S_c$ を満たすような辺もつukらない．反射律と反対称律を満たす辺が存在しないということは，グラフ上では閉路が存在しない，すなわちグラフが DAG であることに対応する．

半順序集合を表現した DAG に対してトポロジカルソートを適用する．半順序関係にトポロジカルソートを使用し得られる解は，半順序関係の順序を満たすような全順序の 1 つに相当する．有力ソフトウェアは DAG において，出力辺が 0 であるノードで極大元にあたる．有力ソフトウェアを除く候補ソフトウェアの半順序集合 SW_{po} に以下の手順でトポロジカルソートを行い順位付ける．

1. DAG の極大元であるノードをすべて選択
2. 選択したノードを選択した順に順位付け
3. DAG から選択したノードとそのノードへの入力辺を削除
4. DAG にノードが残っていれば 1. へ戻る

例の半順序集合 SW_{po} から生成された有力ソフトウェアを含む DAG を図 2 に示す．有力ソフトウェアであるノードは \odot で囲っており極大元であることがわかる．有力ソフトウェアを除く DAG の極大元は， S_4 と S_5 であり，この順に選択しそれぞれ 4 位，5 位と順位付ける．そして DAG からノード S_4, S_5 とそのノードへの入力辺を削除し，再び極大元を選ぶ．残ったノード S_1 が極大元であるので選択し，6 位と順位付ける．有力ソフトウェアの順位

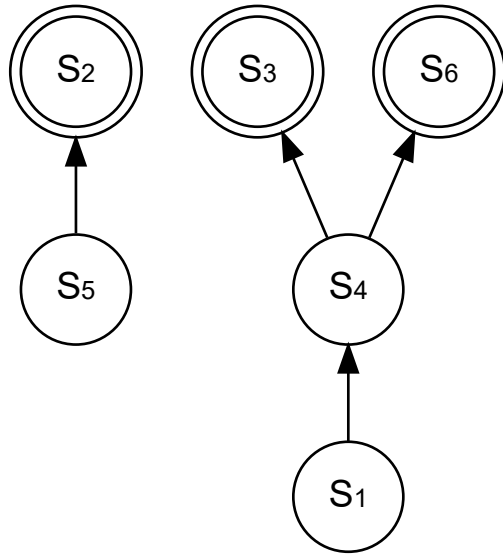


図 2: 有向非巡回グラフの例

付けと同様に k 位である候補ソフトウェアが k 列目となるように列を入れ替え, すべてのソフトウェアの順位が確定した行列を表 3 に示す.

手法では最後に例のような候補ソフトウェアを順位付け並べ替えた行列をリストとして出力する. リストを先頭から見ることでソフトウェアの再利用元を推定することが出来る. この例では, 入力ソースファイル f_1, f_2, f_3 の再利用元はソフトウェア S_3 , 入力ソースファイル f_4, f_5 の再利用元はソフトウェア S_2 であると推定することが出来る.

表 3: トポロジカルソート後の行列

入力 ソースファイル	候補 ソフトウェア	S_3	S_6	S_2	S_4	S_5	S_1
	f_1		0.95	0.85	0	0.85	0
f_2		0.9	0.95	0	0.85	0	0.8
f_3		1.0	0.9	0	0.9	0	0
f_4		0	0	1.0	0	1.0	0
f_5		0	0	0.9	0	0.8	0

4 実験

手法が出力する候補ソフトウェアリストと有力ソフトウェアから再利用元ソフトウェアを推定できるか確認するために、提案手法を Java で実装して実験を行った。対象言語は C/C++、Java とした。本節では、はじめにデータベースに登録するソフトウェアおよびファイルについて説明する。そして、Mozilla Firefox および Android が再利用しているソフトウェアをデータセットとして実験を行い評価した。

4.1 データベース

データベースに登録するデータセットとして、フリーなオペレーションシステムである Debian GNU/Linux が使用しているアーカイブのスナップショット¹を利用する。このアーカイブには 2005 年以来的の使用してきたほぼすべてのパッケージのソースコードが含まれている。本研究では、1 パッケージを再利用元ソフトウェアの 1 つとみなす。マシンインタフェース²を利用し、自動的にパッケージのダウンロードを行った。アクセスした日付は、2016 年 8 月 19 日である。ただし、“404 Not Found”，“503 Service Temporarily Unavailable”のエラーによりダウンロードが行えなかったファイル、ダウンロードができたが破損しているファイルもある。OS に合わせてパッチを適用しているパッケージもあるが、本手法では独自の変更が加えられていない、“*.orig.*”というパターンにマッチしたもののみをダウンロードした。

ダウンロードしたデータセットは、33,496 種類のパッケージを含んでいる。バージョン違いを含めたパッケージの総数は 357,075 であり、パッケージの内容が完全一致している重複

¹<http://snapshot.debian.org/>

²<https://anonscm.debian.org/cgit/mirror/snapshot.debian.org.git/plain/API>

したものを除く総数は188,212である。パッケージの重複は、パッケージのSHA-1ハッシュ値によって確認した。1パッケージあたりのバージョン数は、中央値3で、最大は308である。また、重複を除くパッケージの合計サイズは圧縮形式で約840GBである。

データベースに登録するファイルは、C/C++のソースファイルは拡張子が、.c, .h, .cpp, .hppであるもの、Javaは.javaであるものを対象とした。ただし、拡張子の大文字と小文字は区別しない。登録対象となるソースファイルの総数は、50,903,100ファイル(18,569,351,349LOC)であり、重複を除くと8,926,509ファイル(4,511,358,044LOC)である。LOCは、コメントや空白のみの行も含み、ソースファイルの重複は、ファイルのSHA-1ハッシュ値を使用し除いた。C/C++のソースファイルは、46,699,686ファイル(17,722,308,883LOC)であり、重複を除くと7,658,951ファイル(4,198,847,423LOC)である。また、Javaのソースファイルは、4,203,414ファイル(847,042,466LOC)であり、重複を除くと1,267,558ファイル(312,510,621LOC)である。登録したファイルの合計ファイルサイズは、約566GBである。

4.2 評価方法

実験は、対象プロジェクトが再利用しているソフトウェアのうち、バージョン番号がわかっているものについて行う。そのようなソフトウェアを入力ソフトウェアと呼ぶ。入力ソフトウェアに対応するデータベース内の同名同バージョンのソフトウェアを正解ソフトウェアとする。ただし、正解ソフトウェアとなるものがデータベースにないような入力ソフトウェアは用いない。入力ソフトウェアに手法を適用し、以下の項目について調査し、評価と考察を行う。

- 正解ソフトウェアが出力された候補ソフトウェアリストに含まれているか。含まれているなら、正解ソフトウェアが有力ソフトウェアとして識別されているか。
- 正解ソフトウェアが有力ソフトウェアであるものについて、3つの距離関数がそれぞれ正解ソフトウェアを第何位に出力したか。
- 順序関係を用いることで再利用元候補を絞り込むことができるか。

4.3 対象プロジェクト

4.3.1 Mozilla Firefox

Mozilla Firefox³(以下Firefoxとも表記)は、Mozilla Foundationが開発するオープンソースなウェブブラウザである。開発言語は、主にC++とJavaScriptである。実験で使用した

³<https://www.mozilla.org/ja/firefox/new/>

表 4: Mozilla Firefox が再利用しているソフトウェア

ソフトウェア名	バージョン	入力ソフトウェア の対象ファイル数	正解ソフトウェア の対象ファイル数	入力ソフトウェアのロケーション
cairo	1.10	226	710	/gfx/cairo
graphite2	1.3.6	89	105	/gfx
gtest	1.6.0	22	31	/testing/gtest
hunspell	1.3.3	13	61	/extensions/spellcheck
libav	11.3	76	1907	/media
libevent	2.0.21	130	126	/ipc/chromium/src/third_party
libffi	3.1	209	205	/js/src/ctypes
libjpeg-turbo	1.4.2	79	132	/media
libopus	1.1	235	258	/media
libpng	1.6.19	25	77	/media
libsoundtouch	1.9.0	28	39	/media
libvorbis	1.3.5	65	93	/media
libvpx	1.4.0	487	569	/media
nspr	4.12	462	462	/
nss	3.21.1	1079	1079	/security
snappy	1.0.4	7	8	/other-licenses
sqlite	3.9.1	2	271	/db
srtp	1.4.4	71	79	/nertwerk
stlport	5.2.1	303	553	/build
zlib	1.2.8	27	51	/modules

firefox バージョンは 45.0 である。入力ソフトウェアのバージョン番号は、コミットログやバージョンが記入されたファイルを確認した。入力ソフトウェアと正解ソフトウェアに関する情報を表 4 に示す。

データベースには、Firefox を始めとする Mozilla Foundation が開発するソフトウェアやそれから派生したソフトウェアが登録されている。Firefox に合わせて固有の変更がなされたソフトウェアのソースファイルが、データベース内の Firefox を始めとするソフトウェアにも含まれると考えられる。その場合、そのようなソフトウェアが半順序関係により上位にくると考えられる。そのため、Firefox に対する実験では、表 5 に示すソフトウェアを除外して実験を行った。具体的には、データベース内に登録しているパッケージ名が表 5 のパッケージ名列を含んでいるものを除外した。

表 5: 実験で除外するソフトウェア

ソフトウェア名	概要	パッケージ名
Mozilla Firefox	Mozilla Foundation が開発するウェブブラウザ	firefox
Iceweasel	商標の関係により Firefox から独立したウェブブラウザ	iceweasel
Mozilla Thunderbird	Mozilla Foundation が開発する電子メールクライアント	thunderbird
Icedove	商標の関係により Thunderbird から独立した電子メールクライアント	icedove
Mozilla Sunbird	Mozilla Foundation が開発するスケジュール管理ソフト	sunbird
Iceowl	商標の関係により Sunbird から独立したスケジュール管理ソフト	iceowl
SeaMonkey	Mozilla Foundation から独立した SeaMonkey Council が開発するインターネット統合アプリケーション	seamonkey
Iceape	商標の関係により SeaMonkey から独立したインターネット統合アプリケーション	iceape
Gecko	Mozilla Foundation が開発する HTML レンダリングエンジン群	gecko

4.3.2 Android

Android⁴ は、Google が開発する携帯情報端末向けのオープンソースなオペレーションシステムである。開発言語は、主に C/C++ 及び Java である。実験で使用したバージョンは、4.4.2 であり、“/external” 以下のソフトウェアを入力ソフトウェアとして用いる。入力ソフトウェアと正解ソフトウェアに関する情報を表 6 に示す。

4.4 手法出力と正解ソフトウェアに対する評価

入力ソフトウェアに対する正解ソフトウェアが、出力された候補ソフトウェアリストに含まれているかどうか、さらに有力ソフトウェアであるかどうかについて調べる。結果より再利用しているソフトウェアの再利用元の推定を行えるかどうかを評価することができる。手法は入力ソフトウェア単位に適用している。まず、Firefox の結果を表 7 に示す。正解ソフトウェアが候補ソフトウェアリストに含まれる結果となるものは、20 件中 19 件 (95%) であった。また、正解ソフトウェアが有力ソフトウェアであるのは 17 件 (85%) であった。続いて、Android の結果を表 8 に示す。Android では全ての正解ソフトウェアが候補ソフトウェアリストに含まれた。また、正解ソフトウェアが有力ソフトウェアである結果となったものは 52 件中 49 件 (94.2%) であった。

結果より、Firefox, Android 共に正解ソフトウェアが候補ソフトウェアリストに含まれていた。これより、利用しているソフトウェアの情報がわからなくとも手法を用いることで、

⁴<https://source.android.com/>

表 6: Android が再利用しているソフトウェア

ソフトウェア名	バージョン	入力ソフトウェア の対象ファイル数	正解ソフトウェア の対象ファイル数
arduino	0022	30	200
blktrace	1.0.1	47	46
bouncycastle	1.49	644	3268
bsdif	4.3	2	2
bzip2	1.0.6	15	15
checkpolicy	2.1.11	13	13
dhcpcd	5.5.6	53	49
dnsmasq	2.51	21	20
dropbear	0.49	572	570
e2fsprogs	1.41.14	377	366
easymock	2.5.2	68	140
emma	2.0	198	197
expat	2.1.0	27	55
flac	1.2.1	55	192
genext2fs	1.4.1	2	1
grub	0.97	116	114
guava	11.0.2	972	366
hamcrest	1.1	21	139
ipsec-tools	0.7.3	135	145
iptables	1.4.11.1	208	222
iputils	20121221	26	26
jhead	2.87	8	9
jpeg	6b	92	85
jsr305	0.1	43	43
junit	4.10	164	320
libgsm	1.0.13	23	41
libmtp	1.0.1	48	45
libnl-headers	2.0	42	219
libogg	1.2.0	7	6
libpcap	0.9.8	64	83
libpng	1.2.46	46	45
libsepol	2.2	122	122
libusb	1.0.8	14	13
libvorbis	1.3.1	93	93
libvpx	1.2.0	486	454
libxml2	2.7.8	103	157
libxslt	1.1.26	44	71
mksh	43	20	20
netperf	2.4.4	33	38
openssh	5.9	298	297
openssl	1.0.1e	896	1111
pixman	0.30.0	105	104
safe-iop	0.3.1	2	3
scrypt	1.1.6	9	21
speex	1.2rc1	99	120
sqlite	3.7.11	14	198
srtplib	1.4.4	82	81
stlport	5.2.1	555	553
tagsoup	1.2	19	19
tcpdump	3.9.8	218	225
valgrind	3.8.1	1149	1100
zlib	1.2.8	75	51

データベースに再利用元であるソフトウェアがあれば高い確率でそれが含まれるソフトウェアのリストを得ることが可能であると言える。さらに、ほとんどの正解ソフトウェアが有力ソフトウェアであるという結果から、有力ソフトウェアから再利用元であるソフトウェアを推定することが可能であると言える。候補ソフトウェア数と有力ソフトウェア数の評価については、4.6節で行う。

正解ソフトウェアが候補ソフトウェアリストには含まれていたが、有力ソフトウェアではなかったものに対して考察を行う。考察対象となるソフトウェアは、Firefox の libvpx, srtp, Android の expat, sqlite, zlib である。Firefox の libvpx, srtp, Android の expat, zlib では、入力ソフトウェアには含まれているが正解ソフトウェアには含まれていないソースファイルがあった。また、有力ソフトウェアとなったソフトウェアはそのようなソースファイルを含んでおり、その他のソースファイルの類似度は正解ソフトウェアのものと同等であった。このため、順序関係により正解ソフトウェアが有力ソフトウェアとならなかった。ソフトウェア毎にソースファイルの差異について簡単に記述する。libvpx は、入力ソフトウェア内のファイル README_MOZILLA によると、ビルドを行うためにソースファイルを追加すると記述されていた。srtp は、入力ソフトウェアには含まれるが正解ソフトウェアには含まれないソースファイルが、正解ソフトウェアであるバージョン 1.4.4 以降の srtp には含まれていた。srtp のバージョン 1.4.5 が有力ソフトウェアとしてリストの上位に位置していたため、実際の再利用元はバージョン 1.4.4 と次のバージョン 1.4.5 の間のリビジョンであると考えられる。expat は、入力ソフトウェアではソースファイル expat_config.h がスクリプト configure により生成されていた。zlib は、正解ソフトウェアのディレクトリ contrib 内のファイルがライセンスの関係上含まれていなかった。Android の sqlite は、sqlite の形式の差異によって正解ソフトウェアが有力ソフトウェアとならなかった。公式の sqlite⁵ によると、高速化のために 100 以上のソースファイルを 1 つのソースファイルに集約した形式の sqlite と集約していない形式の sqlite を配布している。Android が再利用する sqlite は集約形式であり、Debian から取得した sqlite は集約していない形式であった。集約される各ソースファイルと集約された 1 つのソースファイル間の類似度が低くなり、類似ファイルとして検索できなかった。そのため、集約形式の正解ソフトウェアは有力ソフトウェアとならなかった。

次に、正解ソフトウェアが候補ソフトウェアリストに入らなかった Firefox の sqlite について考察する。Firefox が再利用する sqlite は、Android と同様に集約形式であった。Android と異なり候補ソフトウェアリストに入らなかった理由について述べる。入力ソフトウェアは、集約ソースファイルとソースファイル sqlite.h の 2 つのソースファイルを持ち、それ以外は削除されていた。正解ソフトウェアは、集約元のソースファイル群を持つが、sqlite.h はコンパイル時に生成されるため含まれない。そのため、ソースファイルが 1 つも類似しなかつ

⁵<https://www.sqlite.org/>

たため候補ソフトウェアリストに入らなかった。

4.5 距離関数と順位の評価

提案手法で有力ソフトウェアと入力ソースファイル集合間の距離を求めるために用いた3つの距離関数について評価を行う。正解ソフトウェアが有力ソフトウェアであったものに対して、各距離関数によって定まった順位を調べる。実装では距離 $d(S_j, F_{sim})$ が同じ場合、データベースに登録しているソフトウェアに割り当てられたIDの順に並ぶようになっている。表9に、Firefoxにおける正解ソフトウェアの各距離関数による順位と有力ソフトウェア、候補ソフトウェアの数を示す。表10に、Androidの結果を示す。各順位内の括弧内の数字は、距離が同じ場合、同順位としてみなしたときの順位である。結果より、jheadを除きどの距離関数を用いても正解ソフトウェアの順位の変動は無かった。jheadはマンハッタン距離、キャンベラ距離による正解ソフトウェアの順位は1位であり、ユークリッド距離のみ正解ソフトウェアが2位であるため、今回用いた3つの距離関数は大きく差がないことがわかった。

次に、正解ソフトウェアの順位の評価を行う。距離関数はマンハッタン距離を用いる。FirefoxとAndroidにおける正解ソフトウェアの順位の数値を図3に示す。図3より、ほとんどの正解ソフトウェアは5位以内である。よって、手法を用いることで有力ソフトウェアの上位にあるソフトウェアが再利用元と推定することが可能であると言える。50位よりも順位が低かったものは、Firefoxのgtestの126位とzlibの59位であった。これらを含めた多くの1位ではないソフトウェアは、表9, 10より、同距離を同順位としてみなすと1位であった。このように、距離が同じソフトウェアが多くある（多くの場合、入力ソフトウェアが持つファイルと同内容のファイルをもつソフトウェアが多数ある）と、同順位としてみなさない順位が大きく下がるため、再利用元を推定することは難しい可能性がある。これはソフトウェア再利用時、ソースコードに独自の変更を加えない場合に多く起こると考えられる。

4.6 順序関係の評価

提案手法で候補ソフトウェアを絞り込み有力ソフトウェアを得るために用いた順序関係について評価を行う。順序関係を用いると絞り込まれた有力ソフトウェアから再利用元を推定することが容易になると考えられる。また複数のソフトウェアを入力し順序関係を用いず距離のみを用いて順序付けた場合、ソースファイル数の少ないソフトウェアの再利用元はリストの下位になるが、順序関係を用いることで有力ソフトウェアとしてリストの上位になると考えられる。そこで候補ソフトウェア数と有力ソフトウェア数の比較を行い推定することが容易になるか調査する。また複数のソフトウェアを入力し、順序関係の有無による正解ソフ

表 7: 出力結果 (Firefox)

ソフトウェア名	候補ソフトウェアリスト に含まれる	有力ソフトウェアである	候補ソフト ウェア数	有力ソフト ウェア数
cairo	✓	✓	127	23
graphite2	✓	✓	30	6
gtest	✓	✓	688	140
hunspell	✓	✓	158	8
libav	✓	✓	592	59
libevent	✓	✓	494	13
libffi	✓	✓	731	4
libjpeg-turbo	✓	✓	667	3
libopus	✓	✓	98	35
libpng	✓	✓	154	7
libsoundtouch	✓	✓	57	1
libvorbis	✓	✓	727	1
libvpx	✓		159	21
nspr	✓	✓	581	2
nss	✓	✓	3962	2
snappy	✓	✓	120	11
sqlite			53	2
srtp	✓		126	9
stlport	✓	✓	20	1
zlib	✓	✓	2691	88

表 8: 出力結果 (Android)

ソフトウェア名	候補ソフトウェアリスト に含まれる	有力ソフトウェアである	候補ソフト ウェア数	有力ソフト ウェア数
arduino	✓	✓	16	1
blktrace	✓	✓	470	6
bouncycastle	✓	✓	65	3
bsdif	✓	✓	5	1
bzip2	✓	✓	986	6
checkpolicy	✓	✓	46	1
dhcpcd	✓	✓	10	2
dnsmasq	✓	✓	44	1
dropbear	✓	✓	490	1
e2fsprogs	✓	✓	5227	49
easymock	✓	✓	6	1
emma	✓	✓	1	1
expat	✓		1181	2
flac	✓	✓	182	1
genext2fs	✓	✓	2	1
grub	✓	✓	3864	1
guava	✓	✓	41	11
hamcrest	✓	✓	3	1
ipsec-tools	✓	✓	315	3
iptables	✓	✓	470	22
iputils	✓	✓	16	2
jhead	✓	✓	18	3
jpeg	✓	✓	1491	70
jsr305	✓	✓	220	1
junit	✓	✓	21	1
libgsm	✓	✓	316	2
libmtp	✓	✓	318	1
libnl-headers	✓	✓	14	2
libogg	✓	✓	460	15
libpcap	✓	✓	575	304
libpng	✓	✓	899	43
libsepol	✓	✓	62	1
libusb	✓	✓	61	1
libvorbis	✓	✓	1095	1
libvpx	✓	✓	1246	23
libxml2	✓	✓	512	2
libxslt	✓	✓	487	24
mksh	✓	✓	53	1
netperf	✓	✓	11	1
openssh	✓	✓	2236	15
openssl	✓	✓	542	35
pixman	✓	✓	392	1
safe-iop	✓	✓	1	1
scrypt	✓	✓	8	1
speex	✓	✓	883	44
sqlite	✓		963	9
srtp	✓	✓	288	40
stlport	✓	✓	862	55
tagsoup	✓	✓	52	1
tcpdump	✓	✓	550	1
valgrind	✓	✓	1489	11
zlib	✓		3553	1

表 9: 各距離関数による正解ソフトウェアの順位 (Firefox)

ソフトウェア名	ユークリッド距離	マンハッタン距離	キャンベラ距離	有力ソフトウェア数
cairo	1(1)	1(1)	1(1)	23
graphite2	2(2)	2(2)	2(2)	6
gtest	126(1)	126(1)	126(1)	140
hunspell	8(1)	8(1)	8(1)	8
libav	3(1)	3(1)	3(1)	59
libevent	1(1)	1(1)	1(1)	13
libffi	1(1)	1(1)	1(1)	4
libjpeg-turbo	1(1)	1(1)	1(1)	3
libopus	3(1)	3(1)	3(1)	35
libpng	1(1)	1(1)	1(1)	7
libsoundtouch	1(1)	1(1)	1(1)	1
libvorbis	1(1)	1(1)	1(1)	1
nspr	1(1)	1(1)	1(1)	2
nss	1(1)	1(1)	1(1)	2
snappy	1(1)	1(1)	1(1)	11
stlport	1(1)	1(1)	1(1)	1
zlib	59(1)	59(1)	59(1)	88

ソフトウェアの順位の変動を調べる。

Firefox の入力ソフトウェアに対する正解ソフトウェアが有力ソフトウェアとなった 17 件に対して評価を行い、距離関数はマンハッタン距離を用いる。入力ソフトウェア単位で手法を適用する場合と入力ソフトウェア 17 件全てをまとめて手法を適用する場合を考える。前者については、候補ソフトウェア数と有力ソフトウェア数の比較のみ行う。入力ソフトウェア単位の場合の候補ソフトウェア数と有力ソフトウェア数は、表 7 に示すとおりであり、それぞれを箱ひげ図で示すと図 4 のようになる。入力ソフトウェア 17 件をまとめた場合の候補ソフトウェア数と有力ソフトウェア数は、表 11 に示す。図 4 より、順序関係を用いることで候補数を大幅に削減できることがわかる。また表 11 より、全てのソフトウェアをまとめて手法を適用した場合、有力ソフトウェア数が 598 と少し多い結果となったが候補ソフトウェア数の約 10% であり、比較的少数の候補だけを抽出していると評価することができる。入力ソフトウェア単位の有力ソフトウェア数の合計よりも多くなった理由としては、2 つ以上の入力ソフトウェアのソースファイルを持つソフトウェアが有力ソフトウェアとなったと推定できる。

次に、順序関係の有無による正解ソフトウェアの順位について調べる。全てのソフトウェアをまとめた入力に対して、順序関係の有無によって得られる 2 つのリスト中の 17 つの正

表 10: 各距離関数による正解ソフトウェアの順位 (Android)

ソフトウェア名	ユークリッド距離	マンハッタン距離	キャンベラ距離	有力ソフトウェア数	候補ソフトウェア数
arduino	1(1)	1(1)	1(1)	1	16
blktrace	1(1)	1(1)	1(1)	6	470
bouncycastle	1(1)	1(1)	1(1)	3	65
bsdif	1(1)	1(1)	1(1)	1	5
bzip2	1(1)	1(1)	1(1)	6	986
checkpolicy	1(1)	1(1)	1(1)	1	46
dhcpcd	1(1)	1(1)	1(1)	2	10
dnsmasq	1(1)	1(1)	1(1)	1	44
dropbear	1(1)	1(1)	1(1)	1	490
e2fsprogs	1(1)	1(1)	1(1)	49	5227
easymock	1(1)	1(1)	1(1)	1	6
emma	1(1)	1(1)	1(1)	1	1
flac	1(1)	1(1)	1(1)	1	182
genext2fs	1(1)	1(1)	1(1)	1	2
grub	1(1)	1(1)	1(1)	1	3864
guava	6(6)	6(6)	6(6)	11	41
humcrest	1(1)	1(1)	1(1)	1	3
ipsec-tools	1(1)	1(1)	1(1)	3	315
iptables	1(1)	1(1)	1(1)	22	470
iputils	1(1)	1(1)	1(1)	2	16
jhead	2(2)	1(1)	1(1)	3	18
jpeg	6(1)	6(1)	6(1)	70	1491
jsr305	1(1)	1(1)	1(1)	1	220
junit	1(1)	1(1)	1(1)	1	21
libgsm	2(1)	2(1)	2(1)	2	316
libmtp	1(1)	1(1)	1(1)	1	318
libnl-headers	2(2)	2(2)	2(2)	2	14
libogg	1(1)	1(1)	1(1)	15	460
libpcap	6(6)	6(6)	6(6)	304	575
libpng	2(1)	2(1)	2(1)	43	899
libsepol	1(1)	1(1)	1(1)	1	62
libusb	1(1)	1(1)	1(1)	1	61
libvorbis	1(1)	1(1)	1(1)	1	1095
libvpx	1(1)	1(1)	1(1)	23	1246
libxml2	1(1)	1(1)	1(1)	2	512
libxslt	24(24)	24(24)	24(24)	24	487
mksh	1(1)	1(1)	1(1)	1	53
netperf	1(1)	1(1)	1(1)	1	11
openssh	1(1)	1(1)	1(1)	15	2236
openssl	3(3)	3(3)	3(3)	35	542
pixman	1(1)	1(1)	1(1)	1	392
safe-iop	1(1)	1(1)	1(1)	1	1
scrypt	1(1)	1(1)	1(1)	1	8
speex	3(1)	3(1)	3(1)	44	883
srtp	1(1)	1(1)	1(1)	40	288
stlport	31(12)	31(12)	31(12)	55	862
tagsoup	1(1)	1(1)	1(1)	1	52
tcpdump	1(1)	1(1)	1(1)	1	550
valgrind	1(1)	1(1)	1(1)	1	1489

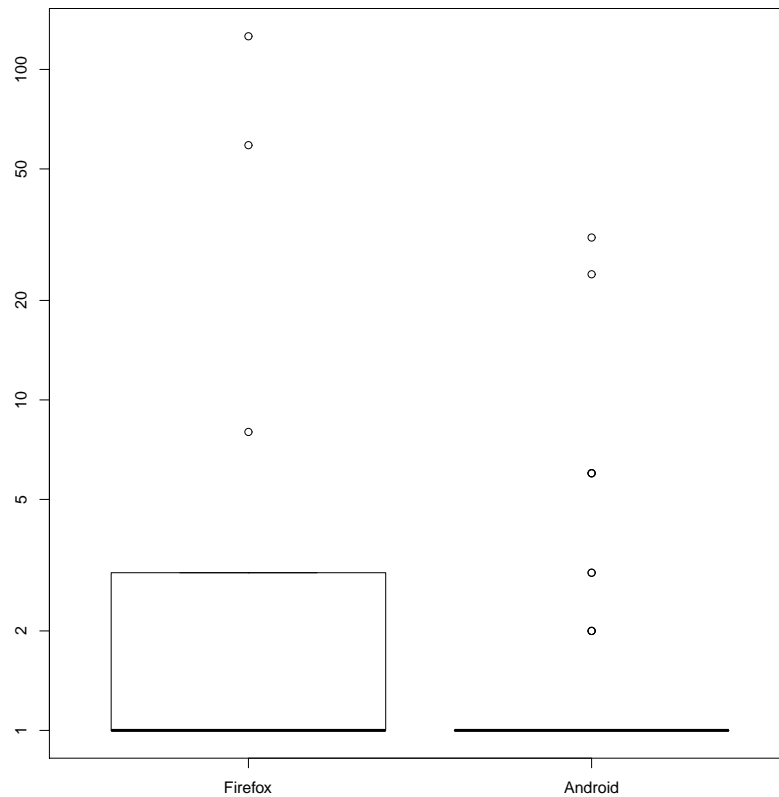


図 3: 正解ソフトウェアの順位

解ソフトウェアの順位を表 12 に示す。ソフトウェア名に “*” がついているものは、入力ソフトウェアをまとめて入力した場合の出力では正解ソフトウェアが有力ソフトウェアとならなかったものである。入力ソフトウェアをまとめた場合の出力でも正解ソフトウェアが有力ソフトウェアものは、14 件であった。その 14 件に対し、正解ソフトウェアの順位は全て上がっていることが確認できる。しかし、順位が 50 位以内に入るものは 2 件という結果であった。これは距離関数の定義より、多くの類似ソースファイルを持てば距離が短くなり、より上位になるという性質から複数ソフトウェアのソースファイルをもつソフトウェアが上位に位置しているためである。正解ソフトウェアが有力ソフトウェアにならなかった gtest, libopus, zlib 全ては、表 9 より、同距離のソフトウェアがあることがわかる。そのソフトウェアが他の入力ソフトウェアのソースファイルを含んでおり、順序関係で順序が付き正解ソフトウェアが有力ソフトウェアにならなかった。これより順序関係を用いることで、絞り込んだ再利用元候補にソフトウェアの再利用元が含まれ、その順位が上げることが可能と言える。ただ

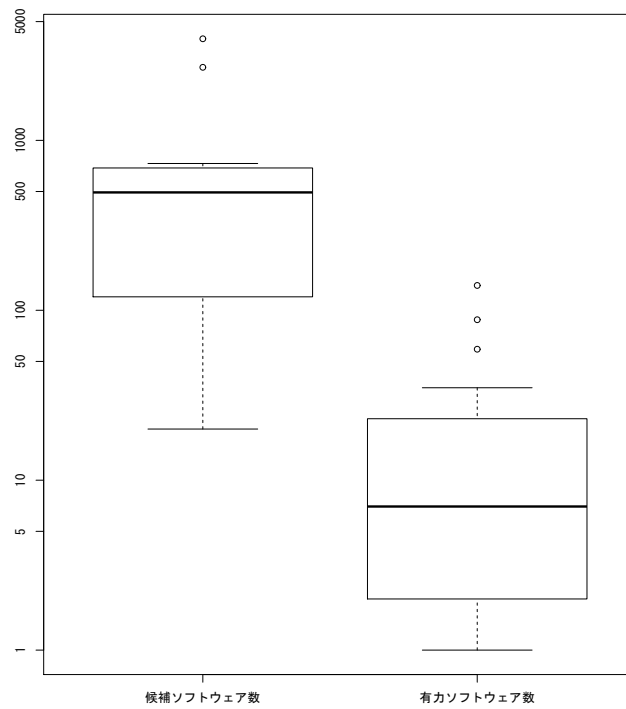


図 4: 入力ソフトウェア単位の場合の各ソフトウェア数

表 11: 入力ソフトウェアをまとめた場合の各ソフトウェア数

候補ソフトウェア数	有力ソフトウェア数
5969	598

し，入力に含まれるソフトウェアの数が増えるほど有力ソフトウェアとなる精度が落ち，順位の上昇率が減る可能性がある．

表 12: 順序関係の有無による正解ソフトウェアの順位

ソフトウェア名	順序関係なし	順序関係あり
cairo	359	130
graphite2	572	222
gtest*	1177	1855
hunspell	3963	558
libav	643	264
libevent	461	160
libffi	341	123
libjpeg-turbo	579	223
libopus*	269	1004
libpng	2642	481
libsoundtouch	2620	479
libvorbis	657	268
nspr	99	35
nss	17	8
snappy	3495	521
stlport	228	78
zlib*	1042	1653

5 妥当性への脅威

本研究の実験において、プロジェクトが再利用しているソフトウェアの正解として Debian GNU/Linux のパッケージとして配布されているソフトウェアを利用した。各ソフトウェアの公式なリリースの配布ではなく、Debian GNU/Linux プロジェクトによる再配布という形であり、ソフトウェアの名前、バージョン番号が情報が正しくない可能性がある。

本研究の実験において、入力ソフトウェアのバージョン番号は各プロジェクトのバージョン番号が記述されているコミットログやファイルの情報を利用した。背景でも述べているように、再利用しているソフトウェアの管理は正しく行われていないという報告があるため、実際のバージョン番号と記述されているバージョン番号が異なる可能性がある。

6 まとめと今後の課題

本研究では、開発中ソフトウェアが再利用しているソフトウェアの再利用元を推定する手法を提案した。開発中ソフトウェアのソースファイル集合を与えると、入力ソースファイル集合と類似するソースファイルを持つソフトウェアを類似する順に並べ替わりリストを表示する。この類似順に並んだソフトウェアリストを見ることで、どのソフトウェアを再利用したのか推定することが可能となる。提案手法を用いて実験を行い、プロジェクトが再利用するソフトウェア 72 件の内、71 件で再利用元を含むリストを出力した。また、手法では再利用元である可能性が高いソフトウェアを絞りこんでおり、絞り込んだソフトウェアの中に再利用元が含まれるものは 66 件であった。そのうち順位が第 1 位のもの 50 件であり、再利用元を推定することが可能であることを確認した。

今後の課題として、再利用元である可能性が高いソフトウェアの絞り込みを類似度による順序関係を用いたが、ソフトウェア全体が持つファイルの総数などによる絞り込みを用いることで精度の向上が期待できる。さらに、絞り込み得られた有力ソフトウェアのよりよい順位付けを提案することができれば、複数ソフトウェアを入力した場合の順位の向上が期待できる。また、再利用元の特定を行った後、どの程度ソースファイルを再利用し変更したのか、独自に新しいソースファイルをどの程度追加したかなどの調査を行うことが容易になると考えられる。

謝辞

本研究および研究生活において、常に適切な御指導及び御助言を賜りました大阪大学大学院情報科学研究科コンピュータサイエンス専攻 井上 克郎 教授に心より深く感謝いたします。

本研究において、随時適切な御指導及び御助言を賜りました大阪大学大学院情報科学研究科コンピュータサイエンス専攻 松下 誠 准教授に深く感謝いたします。

本研究において、逐次適切な御指導及び御助言を頂きました大阪大学大学院情報科学研究科コンピュータサイエンス専攻 石尾 隆 助教に深く感謝いたします。

本研究において、さまざまな御協力を頂きました大阪大学大学院情報科学研究科コンピュータサイエンス専攻 伊藤 薫 氏に深く感謝いたします。

最後に、その他様々な御指導、御助言等を頂いた大阪大学大学院情報科学研究科コンピュータサイエンス専攻井上研究室の皆様に深く感謝いたします。

参考文献

- [1] Moses S. Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of the Thirty-fourth Annual ACM Symposium on Theory of Computing*, STOC '02, pp. 380–388, New York, NY, USA, 2002. ACM.
- [2] Yael Dubinsky, Julia Rubin, Thorsten Berger, Slawomir Duszynski, Martin Becker, and Krzysztof Czarnecki. An Exploratory Study of Cloning in Industrial Software Product Lines. In *Proceedings of the 17th European Conference on Software Maintenance and Reengineering*, pp. 25–34, 2013.
- [3] Christof Ebert. Open source software in industry. *IEEE Software*, Vol. 25, pp. 52–53, 2008.
- [4] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*, STOC '98, pp. 604–613, New York, NY, USA, 1998. ACM.
- [5] K. Inoue, Y. Sasaki, Pei Xia, and Y. Manabe. Where does this code come from and where does it go? – integrated code history tracker for open source systems –. In *Proceedings of the 34th IEEE/ACM International Conference on Software Engineering*, pp. 331–341, 2012.
- [6] Takashi Ishio, Raula Gaikovina Kula, Tetsuya Kanda, Daniel M. German, and Katsuro Inoue. Software ingredients: Detection of third-party component reuse in java software release. In *Proceedings of the 13th International Conference on Mining Software Repositories*, pp. 339–350, 2016.
- [7] Toshihiro Kamiya, Shinji Kusumoto, and Katsuro Inoue. Ccfinder: A multilinguistic token-based code clone detection system for large scale source code. *IEEE Trans. Softw. Eng.*, Vol. 28, No. 7, pp. 654–670, July 2002.
- [8] Tetsuya Kanda, Takashi Ishio, and Katsuro Inoue. Extraction of product evolution tree from source code of product variants. In *Proceedings of the 17th International Software Product Line Conference*, pp. 141–150, Tokyo, Japan, 2013. ACM.
- [9] Naohiro Kawamitsu, Takashi Ishio, Tetsuya Kanda, Raula Gaikovina Kula, Coen De Roover, and Katsuro Inoue. Identifying source code reuse across repositories using lcs-based source code similarity. In *Proceedings of the 2014 IEEE 14th International*

- Working Conference on Source Code Analysis and Manipulation*, SCAM '14, pp. 305–314, Washington, DC, USA, 2014. IEEE Computer Society.
- [10] Jure Leskovec, Anand Rajaraman, and Jeff Ullman. *Mining of Massive Datasets*. Cambridge University Press, 2011.
 - [11] P. Mohagheghi, R. Conradi, O.M. Killi, and H. Schwarz. An empirical study of software reuse vs. defect-density and stability. In *Proceedings of the 26th International Conference on Software Engineering*, pp. 282–291, May 2004.
 - [12] Julia Rubin, Krzysztof Czarnecki, and Marsha Chechik. Managing cloned variants: A framework and experience. In *Proceedings of the 17th International Software Product Line Conference*, pp. 101–110, August 2013.
 - [13] Yusuke Sasaki, Tetsuo Yamamoto, Yasuhiro Hayase, and Katsuro Inoue. Finding file clones in freebsd ports collection. In *MSR*, 2010.
 - [14] Pei Xia, Makoto Matsushita, Norihiro Yoshida, and Katsuro Inoue. Studying reuse of out-dated third-party code in open source projects. *Computer Software*, Vol. 30, No. 4, pp. 98–104, 2013.
 - [15] 川満直弘, 石尾隆, 井上克郎. LSH アルゴリズムを利用した類似ソースコードの検索. 第 2016-SE-191 巻, 2016.