

# メソッド呼び出しに基づくJavaプログラムの 類似度計算と可視化

東京情報大学 総合情報学部  
宇田川 佳久

1

## 目 次

1. 研究の背景と概要
2. 開発した機能と処理流れ
3. メソッドの構文解析と類似度の計算
4. 類似したメソッド集合の視覚的表示
5. 頻出するメソッド集合の抽出
6. おわりに

2

## 1. 研究の背景と概要

- ICTによる社会・生産の改革が進められている **(DX)**。
- プログラミング技法は複雑さを増し、ソフトウェア開発の効率化やエンジニアの育成にも改革が求められている **(EX)**。
- インターネットに公開されている**サンプルプログラム**を参考にする方法がある。
- サンプルプログラムの違いや類似性は、プログラムの中身を**目視で確認**しなければならず、この作業がプログラミング学習の制約となっていた。
- 類似するプログラムを検索することで、プログラミング**開発や学習を支援**することを研究テーマとした。
- 本研究の目的は、**類似するプログラムの検索方法と表示方法を実証**することである。

3

## 従来の研究との違い

- (1)メソッド単位の類似度を**APIの種類と引用回数**に基づいて算出する。
- (2)メソッド単位の類似度を**グラフで表示**し、関連性を直感的に把握できる機能を提供している。
- (3)頻出するAPIの組み合わせを算出する機能により、**サンプルプログラムが提示しているAPIの組み合わせを頻度順に提示**できる。

2021/5/21

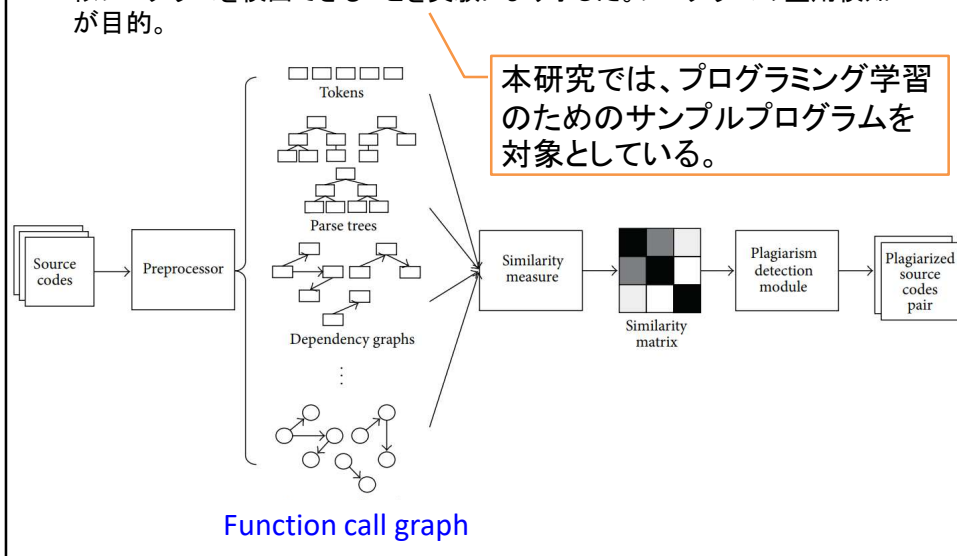
4

## 関連研究

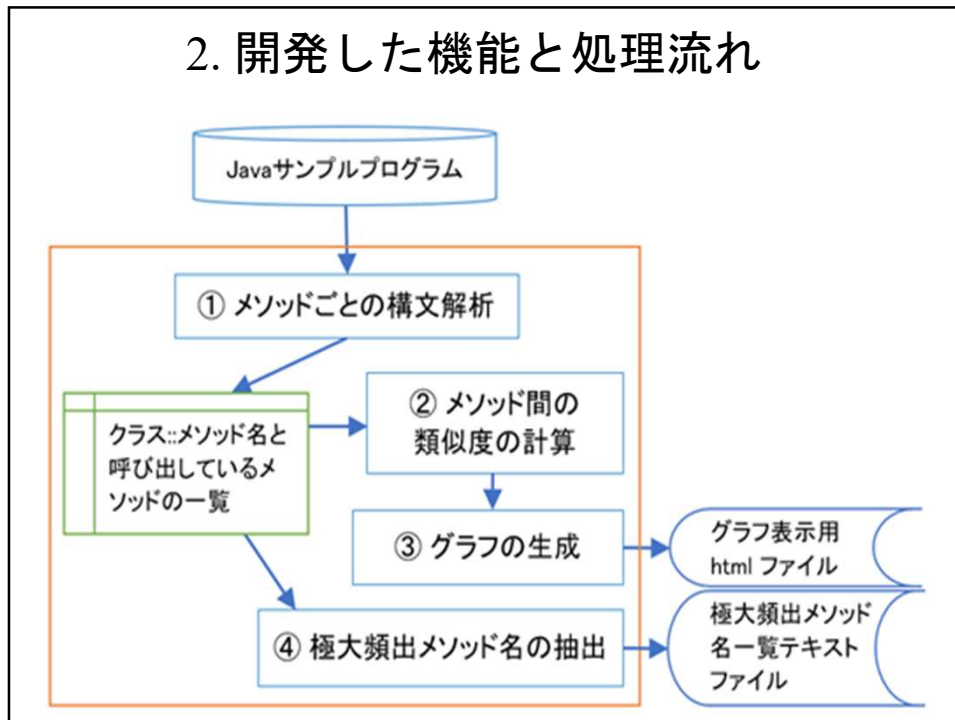
1. F. Zhang, Guofan Li, Cong Liu, and Qian Song, "Flowchart-Based Cross-Language Source Code Similarity Detection," Hindawi Scientific Programming Vol. 2020, Article ID 8835310, pp.1-15.
2. R. Pelánek, T. Effenberger, M. Vanek, V. Sassmann, and D. Gmiterko, "Measuring Item Similarity in Introductory Programming," Proceedings of the Fifth Annual ACM Conference on Learning at Scale, June 2018, pp.1-4.
3. C. Ragkhitwetsagul, "Measuring Code Similarity in Large-scaled Code Corpora," IEEE International Conference on Software Maintenance and Evolution (ICSME), Oct. 2016.
4. H-J. Song, S-B. Park, and S. Y. Park, "Computation of Program Source Code Similarity by Composition of Parse Tree and Call Graph," Mathematical Problems in Engineering Vol. 2015, Article ID 429807, pp.1-12.
5. I. Keivanloo, J. Rilling, and Y. Zou, "Spotting working code examples," Proceedings of the 36th International Conference on Software Engineering, May 2014, pp.664-675.
6. F. A. Shamsi and A. Elnagar, "An Intelligent Assessment Tool for Students' Java Submissions in Introductory Programming Courses," Journal of Intelligent Learning Systems and Applications, Apr. 2012, pp.59-69.

5

2015年 H-Je Song氏らは[4], プログラムのローカルな制御構造を循環複雑度で指標化し, プログラムの全体構造をメソッド呼び出しで表すことを提案した。これらの2種類の指標を使うことで, 従来手法よりも高い精度で類似プログラムを検出できることを実験により示した。プログラムの盗用検知が目的。



## 2. 開発した機能と処理流れ



## 3. メソッドの構文解析と類似度の計算

- Javaプログラムの構文解析では、JDT(Java Development Tools)-Core の `Scanner` クラスを使用した。

```

public class LongestCommonSubsequence_2 {
    int lcs(char[] X, char[] Y, int m, int n) {
        int L[][] = new int[m+1][n+1];
        for (int i=0; i<=m; i++) {
            for (int j=0; j<=n; j++) {
                if (i == 0 || j == 0)
                    L[i][j] = 0;
                else if (X[i-1] == Y[j-1])
                    L[i][j] = L[i-1][j-1] + 1;
                else
                    L[i][j] = max(L[i-1][j], L[i][j-1]);
            }
        }
        return L[m][n];
    }
    int max(int a, int b) {
        return (a > b)? a : b;
    }
}
  
```

1	1 .	17	36 new
2	2 ++	18	40 static
3	4 +	19	42 1
4	5 -	20	47
5	6 [	21	49 {
6	12 <=	22	57 public
7	15 >	23	63 :
8	19 ==	24	67 ]
9	22 n	25	68 class
10	23 (	26	73 =
11	25 ;	27	82 for
12	26 )	28	83 if
13	29 ?	29	84 return
14	31	30	105 char
15	32	31	109 int
16	33 }	32	112 void
		33	114 else

2021/5/21

## 3. メソッドの構文解析と類似度の計算

- Javaプログラムを構成するトークンを100種類以上に分類する機能を提供している。また、コメントを除外するため、実行文の解析を効率的に行える。
- 一行に複数のメソッドがある場合にも対応している。
- `println()` や `printStackTrace()` など、共通的に使用されるメソッドは、除外する機能を作り込んだ。

```
private static void timeMeasure(String MID) {
    try {
        System.out.println(MID);
        Method method = Class.forName("Sample1607").getMethod(MID);
        long startTime = System.currentTimeMillis();
        method.invoke(null); // メソッドの実行。
        long endTime = System.currentTimeMillis();
        System.out.println(MID+" ReadTime: "+(endTime - startTime)+" msec");
    } catch (Exception e) {
        e.printStackTrace();
    }
}
```

2021/

- 任意個のメソッドとメソッド呼び出しを管理するために、MapのArrayListを使っている。
- グラフ表示には、Vis.jsを使っている。

```
int cnt= 0;
ArrayList< Map<String, Integer> > MCall;
for( Javaソースファイル) {
    if(クラス名::メソッド宣言){
        cnt++;
        MCall.add(cnt, Map(呼び出されているメソッド名, 呼出し数));
    }
    Num= Mcall.size(); //メソッド宣言の数;
    GR= ∅;
    for( j=0; j < Num-1; j++) { // メソッド間の類似度を計算する
        for( k=j+1; k < Num; k++) {
            // MCallのメソッド名と呼出し数を使ってCos類似度を計算する。
            Sim[j][k]= CosSim( j, k );
            if( Sim[j][k] > 閾値 )
                // ノードj, kをアークSim[j][k]で結合したネットワークグラフの生成
                GR= GR U Gen_Graph(j, k);
        }
    }
}
```

2021/

10

## メソッド抽出処理の概要 (事例)

```

1 package Sample_1;
2 import java.io.*;
3 import java.lang.reflect.Method;
4
5 class Sample1607 {
6     static String FName = "c:\\temp\\Aprion\\VMD_Struct.txt";
7     public static void main(String args[]) {
8         timeMeasure("ReadByte");
9         timeMeasure("ReadByteBuffered");
10        timeMeasure("ReadText");
11        timeMeasure("ReadTextBuffered");
12    }
13    public static void ReadByte() throws IOException {
14        FileInputStream fis = new FileInputStream(FName);
15        int code;
16        while ((code = fis.read()) != -1) {
17            Integer.toHexString(code);
18        }
19        fis.close();
20    }
21    中略
22    public static void ReadTextBuffered() throws IOException {
23        BufferedReader br = new BufferedReader(new FileReader(FName));
24        int data;
25        while ((data = br.read()) != -1) {
26            Character.toChars(data);
27        }
28        br.close();
29    }
30    private static void timeMeasure(String MID) {
31        try {
32            System.out.println(MID);
33            Method method = Class.forName("Sample1607").getMethod(MID);
34            long startTime = System.currentTimeMillis();
35            method.invoke(null); // メソッドの実行。
36            long endTime = System.currentTimeMillis();
37            System.out.println(MID+" ReadTime : "+(endTime - startTime)+"
38            2021/5/21 catch (Exception e) {
39                e.printStackTrace();
40            }
41        }
42    }
43    }
44    }
45    }
46    }
47    }
48    }
49    }
50    }
51    }
52    }
53    }
54    }
55    }
56    }

```

```

Sample1607::main(String[])
timeMeasure,4
Sample1607::ReadByte()
FileInputStream,1
close,1
read,1
toHexString,1
Sample1607::ReadByteBuffered()
BufferedReader,1
FileInputStream,1
close,1
read,1
toHexString,1
Sample1607::ReadText()
FileInputStream,1
InputStreamReader,1
close,1
read,1
toChars,1
Sample1607::ReadTextBuffered()
BufferedReader,1
FileReader,1
close,1
read,1
toChars,1
Sample1607::timeMeasure(String)
currentTimeMillis,2
forName,1
getMethod,1
invoke,1

```

## 類似度の計算

- 呼び出しているメソッド名と個数で、Wordのベクトルを表現する。
- Wordのベクトル間の Cos類似度を計算する。
- 呼び出している個数を反映した類似度を計算できる。
- なお、Wordの発生頻度は考慮していない(tf-idfによる補正はしていない)。

$$\mathbf{a}=(a_1, a_2, \dots, a_n), \mathbf{b}=(b_1, b_2, \dots, b_n)$$

$$\cos(a, b) = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \sqrt{\sum_{i=1}^n b_i^2}}$$





## 4.2 Vis.js ネットワークグラフの生成

```

1 <!doctype html>
2 <html>
3 <head>
4 <title>Java Program Network</title>
5 <style type="text/css">
6 #mynetwork {
7   width: 1200px;
8   height: 800px;
9   border: 1px solid lightgray;
10 }
11 </style>
12 <script type="text/javascript" src="../vis.js"></script>
13 <link href="../vis-network.css" type="text/css" />
14 <script type="text/javascript">
15 function draw() {
16   <nodes の定義>
17   <edges の定義>
18   var container = document.getElementById
19   var data = {
20     nodes: nodes,
21     edges: edges
22   };
23   var options = {};
24   var network = new vis.Network(container, data, options);
25 }
26 </script>
27 </head>
28 <body onload="draw()">
29 <div id="mynetwork"></div>
30 </body>
31 </html>

```

Vis.js をインストールしたパスに変更する。

表示したい Node edge を定められたフォーマットで書き込む。

15

## 変数nodesとedgesの仕様

```

var nodes = [
  {id:<番号>, label: '<ラベル名>', shape: '<形状>'},
  ... ];
var edges = [
  {from: <番号>, to: <番号>, color: '<色>',
    width:<幅>, length:<長さ>, label: '<類似度>' },
  ... ];

```

2021/5/21

16



## 変数nodesとedgesの生成例

```

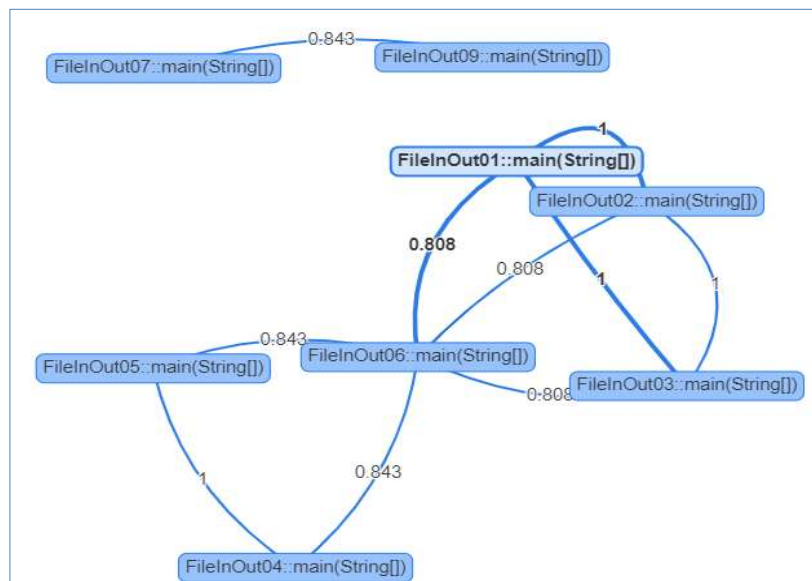
var nodes = [
  {id:1, label: 'FileInOut01::main(String[])', shape: 'box'},
  {id:2, label: 'FileInOut02::main(String[])', shape: 'box'},
  {id:3, label: 'FileInOut03::main(String[])', shape: 'box'},
  {id:4, label: 'FileInOut06::main(String[])', shape: 'box'},
  {id:5, label: 'FileInOut04::main(String[])', shape: 'box'},
  {id:6, label: 'FileInOut05::main(String[])', shape: 'box'},
  {id:7, label: 'FileInOut07::main(String[])', shape: 'box'},
  {id:8, label: 'FileInOut09::main(String[])', shape: 'box'}
];
var edges = [
  {from: 1, to: 2, color: 'red', width: 2, length: 200, label: '1'},
  {from: 1, to: 3, color: 'red', width: 2, length: 200, label: '1'},
  {from: 1, to: 4, color: 'red', width: 2, length: 200, label: '0.808'},
  {from: 2, to: 3, color: 'red', width: 2, length: 200, label: '1'},
  {from: 2, to: 4, color: 'red', width: 2, length: 200, label: '0.808'},
  {from: 3, to: 4, color: 'red', width: 2, length: 200, label: '0.808'},
  {from: 5, to: 6, color: 'red', width: 2, length: 200, label: '1'},
  {from: 5, to: 4, color: 'red', width: 2, length: 200, label: '0.843'},
  {from: 6, to: 4, color: 'red', width: 2, length: 200, label: '0.843'},
  {from: 7, to: 8, color: 'red', width: 2, length: 200, label: '0.843'}
];

```

2021/5/21

17

## Visjsによるグラフの表示例



20

18

## 5. 頻出するメソッド集合の抽出

### 5.1 Aprioriアルゴリズムと極大頻出集合

- 頻出するデータを見つけ出す問題は、大量のデータから有用なパターンを見つけるための基本手法。
- **Aprioriアルゴリズム**は、典型的なもの。
- 頻出データは大量に検出されることが知られており、その中から有意な頻出系列を検出する技法として、**極大頻出集合**がある。
- 本研究では、呼び出されているメソッド名の極大頻出集合を計算することで、**プログラミングのテーマを見つける**ことを試みた。

2021/5/21

19

## 極大頻出集合

```
timeMeasure
FileInputStream close read toHexString
BufferedInputStream FileInputStream close read toHexString
FileInputStream InputStreamReader close read toChars
BufferedReader FileReader close read toChars
currentTimeMillis forName getMethod invoke
```

```
1 close (66.6667) ←
2 read close (66.6667) ←
3 read (66.6667) ←
4 FileInputStream read close (50) ←
5 FileInputStream read (50) ←
6 FileInputStream close (50) ←
7 FileInputStream (50) ←
8 toHexString FileInputStream read close (33.3333) ←
9 toHexString FileInputStream read (33.3333) ←
10 toHexString FileInputStream close (33.3333) ←
11 toHexString FileInputStream (33.3333) ←
12 toHexString read close (33.3333) ←
13 toHexString read (33.3333) ←
14 toHexString close (33.3333) ←
15 toHexString (33.3333) ←
16 toChars close read (33.3333) ←
17 toChars close (33.3333) ←
18 toChars read (33.3333) ←
19 toChars (33.3333) ←
```

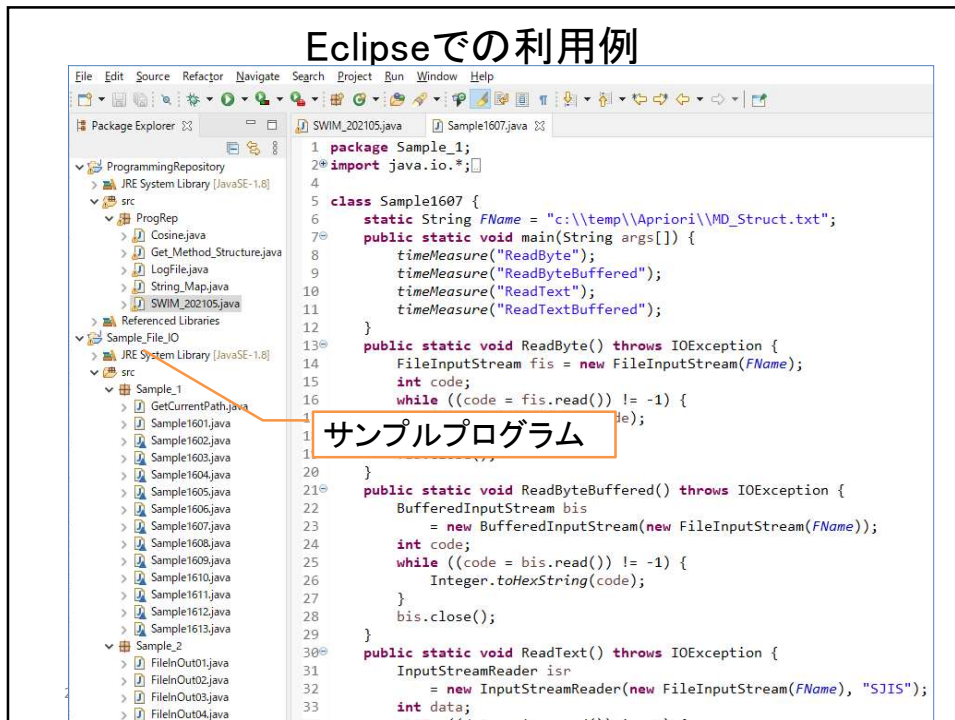
支持度30% 以上で抽出した  
極大頻出アイテム集合(データ)

```
1 toHexString FileInputStream read close (33.3333) ←
2 toChars close read (33.3333) ←
```

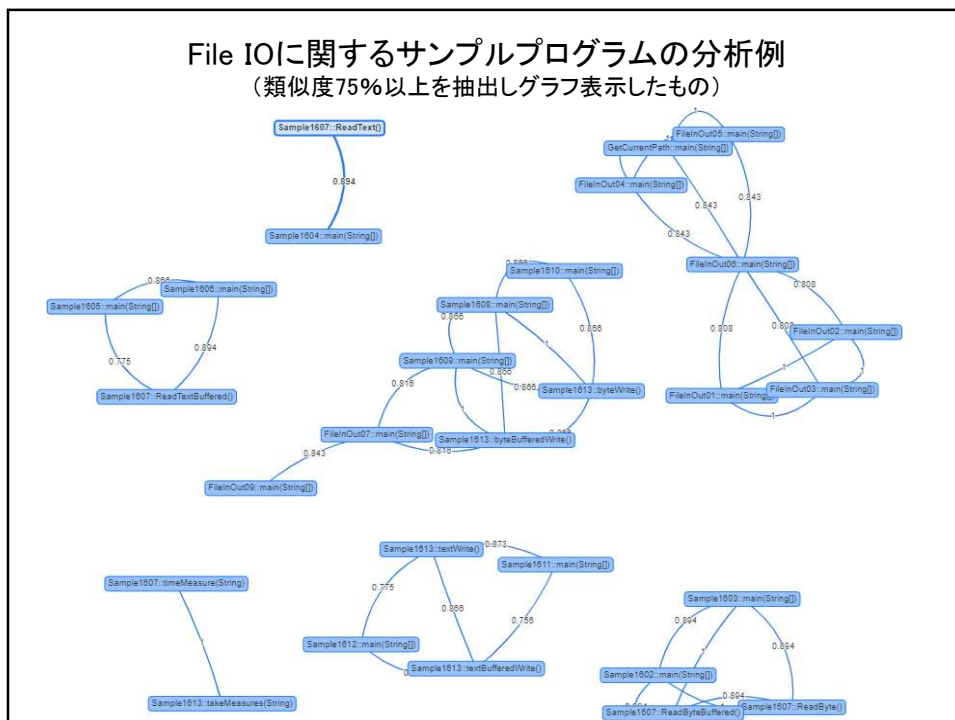
支持度30% 以上で抽出した  
頻出アイテム集合(データ)

20

## Eclipseでの利用例



## File IOに関するサンプルプログラムの分析例 (類似度75%以上を抽出しグラフ表示したもの)



## 6. おわりに

1. サンプルプログラムの利用を推進するための類似検索機能とグラフ表示方法を提案した。
  2. 加えて、頻出するメソッドを抽出することで、優先的に学習すべきメソッド集合を導出することを試みた。
  3. プログラムの構造解析により、呼び出しているメソッド名と回数を抽出し、Cos類似度でメソッド同士の類似度を計算した。
- 今後は、実用的なサンプルプログラムを使った実験を行う。
  - また、教育現場への適用を計画している。