

コードクローン間の依存関係に基づくリファクタリング支援

吉田 則裕[†] 肥後 芳樹[†] 神谷 年洋^{††}
楠本 真二[†] 井上 克郎[†]

コードクローンとは、互いに一致または類似したコード片（ソースコードの断片）を持つコード片を意味し、ソフトウェア保守を困難にしている要因の 1 つとされている。例えば、あるコード片にバグが含まれていた場合、そのコード片のコードクローン全てに対して、修正の是非を検討する必要がある。コードクローンを取り除く方法として、リファクタリングの適用が考えられる。リファクタリングとは、ソフトウェアの外部的振舞いを変化させることなく、内部の構造を改善する作業のことである。しかし、クローンセット（互いに一致または類似したコード片の集合）に含まれるコード片と周辺のコード片間に依存関係が存在すると、リファクタリングが困難になる場合がある。本稿は、クローンセット間の依存関係を利用したリファクタリング支援手法を提案する。まず、異なるクローンセットに含まれるコード片間の依存関係に着目し、そのような依存関係を持つコード片の集合をチェンドクローンセットと定義する。そして、メトリクスを用いてチェンドクローンセットの特徴を判定し、適用可能なリファクタリングパターンを提示する。最後に、リファクタリング支援ツールとして実装することで、いくつかのオープンソースソフトウェアに適用し、有効性の評価を行う。

On Refactoring Support Based on Code Clone Dependency Relations

NORIHITO YOSHIDA,[†] YOSHIKI HIGO,[†] TOSHIHIRO KAMIYA,^{††}
SHINJI KUSUMOTO[†] and KATSURO INOUE[†]

Code clone is a set of code fragments identical or similar to each other. It is generally said that code clone is one of the factors that make software maintenance more difficult. If we modify one of them, it is necessary to determine whether or not we have to modify the others. Refactoring is a disciplined technique for restructuring an existing body of code, altering its internal structure without changing its external behavior. However, there are dependency relations between code fragments belonging to the different clone sets (a clone set is an equivalence class of code clones), and it is difficult to apply refactoring to such code clones. In this paper, we propose a refactoring support method by using dependency relations. At first, we focus on dependency relations between code fragments belonging to the different clone set, and we define “chained clone set” as such code fragments. Then, we define the metrics for providing an appropriate refactoring pattern to each “chained clone set”. Finally, we present the “chained clone set” refactoring support tool that we have developed, together with some case studies.

1. はじめに

ソフトウェア保守を困難にする要因の 1 つとして、コードクローンが指摘されている^{1)~5)}。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、または類似したコード片（ソースコードの断片）を意味する。例えば、あるコード片にバグが含まれていた場合、そのコード片のコードクローン全てに対して、修正の是非を検討する必要がある。このような作業に要する

コストが、特に大規模ソフトウェアの保守において問題となる。

ソフトウェア保守性を改善する技術の 1 つとして、リファクタリング⁶⁾がある。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変化させることなく、内部の構造を改善する技術のことである。Fowler は、リファクタリングを検討すべき箇所にあられる特徴を Bad Smell と呼び、その代表例としてコードクローン (Duplicated Code) を挙げている。また、コードクローンを単一のモジュールに集約する手法として、“Pull Up Method” や “Extract Method”, “Extract SuperClass” などのリファクタリングパターンを紹介している。

これまでに、我々はコードクローン検出ツール

[†] 大阪大学 大学院情報科学研究科
Graduate School of Information Science and Technology, Osaka
University
^{††} 産業技術総合研究所
National Institute of Advanced Industrial Science and Technology

CCFinder⁷⁾ およびリファクタリング支援環境 Aries⁸⁾を開発してきている。CCFinder は、ソースコードに字句解析と正規化処理を行うことで得られたトークン列の同値性に基づいてコードクローン検出を行う。CCFinder の特徴は、表現上の差異があるコードクローンを検出できること、および百万行単位のソースコードであっても実行時間で解析できることである。Aries は、CCFinder の出力情報を基に、リファクタリングに適した単位 (e.g. クラス, メソッド単位) でクローンセット (互いに一致または類似したコード片の集合) を検出し、更にメトリクスで特徴付けすることでリファクタリングパターンの提示を行う。

これまでに、Aries を用いて様々なソースコードを解析した結果、異なるクローンセットに含まれるコード片間に依存関係が存在する場合が確認されている。例えば、クローンセット S_a に 2 つのメソッド m_{a1} , m_{a2} が含まれ、同様にクローンセット S_b に 2 つのコード片 m_{b1} , m_{b2} が含まれるときに、メソッド m_{a1} がメソッド m_{b1} を呼び出し、メソッド m_{a2} がメソッド m_{b2} を呼び出しているという場合である。

Aries は、上述した呼出関係の解析は行っていないため、ユーザは自らクローンセット S_a と S_b 間の呼出関係を把握する必要がある。もし、ユーザが S_a に対して m_{a1} と m_{b1} , m_{a2} と m_{b2} の呼出関係を考慮せずに集約を試みると、呼出関係が保存されない可能性がある。有効なリファクタリング支援を行うためには、クローンセット S_a と S_b は、まとめてユーザに提示するべきであると考えられる。

本稿では、クローンセット間の依存関係を利用したリファクタリング支援手法を提案する。まず、異なるクローンセットに含まれるコード片間の依存関係に着目し、そのような依存関係を持つコード片の集合をチェーンクローンセットと定義する。そして、チェーンクローンセットの特徴に応じて、適用可能なリファクタリングパターンを提示するためのメトリクスを定義する。最後に、提案手法をリファクタリング支援ツールとして実装し、2 つのオープンソースソフトウェアに適用することで有効性の評価を行う。

2. コードクローン

2.1 コードクローンの定義

あるトークン列中に存在する 2 つの部分トークン列 α , β が等価であるとき、 α と β は互いにクローンであるという。またペア (α, β) をクローンペアと呼ぶ。 α , β それぞれを真に包含する如何なるトークン列も等価でないとき、 α , β を極大クローンと呼

ぶ。また、クローンの同値類をクローンセットと呼ぶ。ソースコード中でのクローンを特にコードクローンという⁹⁾。

2.2 コードクローン検出ツール: CCFinder

CCFinder はプログラムのソースコード中に存在する極大クローンを検出し、その位置をクローンペアのリストとして出力する。検出されるコードクローンの最小トークン数はユーザが前もって設定できる。

CCFinder のコードクローン検出手順は以下の 4 つの STEP からなる。

STEP1 (字句解析): ソースファイルを字句解析することによりトークン列に変換する。入力ファイルが複数の場合には、個々のファイルから得られたトークン列を連結し、単一のトークン列を生成する。

STEP2 (変換処理): 実用上意味を持たないコードクローンを取り除くこと、及び、些細な表現上の違いを吸収することを目的とした変換ルールによりトークン列を変換する。例えば、この変換により変数名は同一のトークンに置換されるので、変数名が付け替えられたコード片もコードクローンであると判定することができる。

STEP3 (検出処理): トークン列の中から指定された長さ以上一致している部分をクローンペアとして全て検出する。

STEP4 (出力整形処理): 検出されたクローンペアについて、ソースコード上での位置情報を出力する。

2.3 リファクタリング支援環境: Aries

リファクタリング支援環境 Aries は、CCFinder の出力情報を基に、リファクタリングに適した単位 (e.g. クラス, メソッド単位) でクローンセットを検出し、更にメトリクスで特徴付けすることでリファクタリングパターンの提示を行う。

リファクタリングに適した単位とは、ソースコード上の構造的なまとまりのことである。クローンセットに含まれるコード片が構造的なまとまりを持っているなら、容易に集約することが出来る。現在、Aries は Java 言語を対象として実装されているため、用いる構造的なまとまりは以下の 12 種類である。

宣言 : class { }, interface { }
 メソッド : メソッド本体, コンストラクタ, スタティックイニシャライザ
 文 : if, for, while, do, switch, try, synchronized

クローンセットの特徴付けに用いるメトリクスの 1 つとして、分散度メトリクス $DCH(S)$ ⁸⁾ について説明

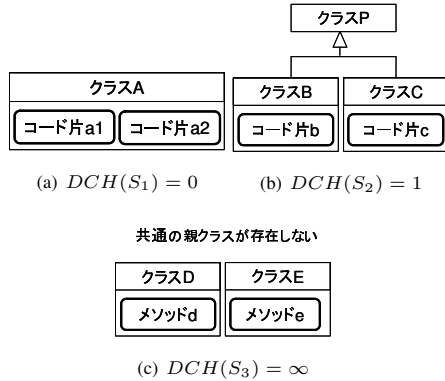


図 1 DCH メトリクスの算出例
Fig. 1 Example of DCH metric calculation

する。クローンセット S はコード片 f_1, f_2, \dots, f_n を含んでいるとする。クラス C_i はコード片 f_i を含んでいるクラスとする。もしクラス C_1, C_2, \dots, C_n が共通の親クラスを持つ場合は、その共通の親クラスの中で、クラス階層的に最も下位（最も深い階層）に位置するクラスを C_p で表すとする。また $D(C_k, C_h)$ はクラス C_k と C_h のクラス階層における距離を表すとする。この時、

$$DCH(S) = \max \{D(C_1, C_p), \dots, D(C_n, C_p)\}$$

と表される。直観的には、 $DCH(S)$ メトリクスはクローンセット S に含まれる各コード片間のクラス階層内における最大の距離を示す。図 1(a)~(c) は、それぞれクローンセットに含まれる 2 つのコード片に対して、 $DCH(S)$ メトリクスを算出した例である。 $DCH(S)$ の値は、全てのコード片が 1 つのクラス内に存在する場合は 0 (図 1(a))、あるクラスとその直接の子クラス内に存在する場合は 1 となる (図 1(b))。例外的に、コードクローンが存在するクラスが共通の親クラスを持たない場合は ∞ とする (図 1(c))。このメトリクスは、クラスライブラリ等の修正不可能なクラスを除外したクラスを対象として計算される。

$DCH(S)$ メトリクスにより、クローンセット S のコード片を集約したモジュールを置くことが出来るクラス（集約先）を特定することが出来る。例えば、 $DCH(S)$ メトリクスの値が 1 の場合は、そのクローンセットが存在するクラスの親クラスに集約できることがわかる。また、 $DCH(S)$ メトリクスの値が ∞ の場合は、分析対象内に集約先位置になるクラスが存在しないため、クラスの作成、もしくは継承関係のない

クラスへの集約を検討するべきであることがわかる。

3. 提案手法

3.1 本研究の動機

リファクタリングを行う際は、“保守性を悪化させるコード（Bad Smell）の特定”と“外部的振る舞いを変更させず保守性を向上させる修正作業”を行う必要がある。

既に述べたように、保守性を悪化させるコードとしてコードクローンが指摘されている。門田らは、コードクローンを含むモジュールの方が保守コストが大きくなる事例を紹介している¹⁰⁾。彼らが紹介した事例では、より行数の多いコードクローンを含むモジュールほど、保守コストが大きくなっていった。

検出したコードクローンを容易にリファクタリングするためには、次の 2 つの条件を満たしていることが望ましい。

(条件 1) クローンセットに含まれるコード片が構造的なまとまりを持っている。

(条件 2) クローンセットに含まれるコード片と周辺のコード間に結合が少ない。

(条件 1) については、2.3 節で述べている。(条件 2) の例として、メソッド呼出を含む場合 (図 2(a)) や、フィールド変数を使用 (参照, 代入) している場合 (図 2(b)) が挙げられる。図 2(a) では、メソッド a_1 と b_1 がコードクローンとなっており、点線矢印で表すようにメソッド a_1 が a_2 を呼び出し、メソッド b_1 が b_2 を呼び出している。図 2(b) では、メソッド a_1 とメソッド a_2 はフィールド変数 v_a を参照しており、メソッド b_1 とメソッド b_2 はフィールド変数 v_b を参照している。このような周辺のコードと結合度が高いコードクローンを集約することは困難である。なぜなら、集約作業を行う際に、メソッド間の呼出関係やメソッドとフィールド変数の関係に不整合が起きないように考慮する必要があるからである。

ここまで述べたように、コード量が多いコードクローンにかかる保守コストや、周辺のコードと結合が多いコードクローンに対するリファクタリングの難しさが問題になりやすい。我々はコードクローン分析を行う中で、そのようなコードクローンのパターンを見つけた。それは、異なるクローンセット間に依存関係が存在する場合である (図 3(a))。図 3(a) は、クラス A, B 間にまたがる 3 つのクローンセットを表している。これら 3 つのクローンセットは、他のクローンセットとの間にメソッド呼出関係やフィールド変数の参照関係が存在する。図 3(a) のクラス A, B のよう

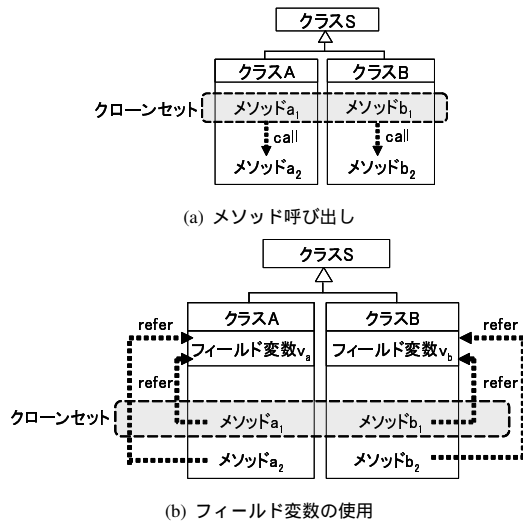


図 2 周辺のコードと結合している例
Fig. 2 Examples of codes are coupled each other

に、類似したクラスには、“Extract SuperClass” パターンの適用を検討すべきと指摘されている⁶⁾。クラス A, B に “Extract SuperClass” パターンの適用した結果が図 3(b) である。

図 3(a) から図 3(b) へのリファクタリングを支援する方法として、図 3(a) 中に含まれる 3 つのクローンセットを同時に提示することが考えられる。

その理由の 1 つ目は、1 つのクローンセットを対象としたリファクタリングは、呼出関係が原因で困難な場合があるからである。例えば、図 3(a) のクローンセットの 1 つ (メソッド a_1 とメソッド b_1) に対して集約を試みると、親クラスに新たに作成したメソッド s_1 から子クラスのメソッド a_2 と b_2 を呼び出すことが出来なくなる (図 3(c))。この問題を解決するためには、子クラスのメソッドに対応する抽象メソッドを親クラスに追加する必要がある。3 つのクローンセットを同時にリファクタリングする際にはこのような工夫は必要ない。

2 つ目の理由は、図 3(a) のようなクローンセットの組み合わせを手作業で見つけ出すことは難しいからである。なぜなら、大規模ソフトウェアは大量のクローンセット含んでいることが多いため、それらの全てに対して呼出関係の有無を確認することは困難だからである。また、呼出関係の理解支援にはコールグラフ (Call Graph) がよく用いられるが、頂点数 (メソッド数) が膨大になりやすい大規模ソフトウェアへの適用は現実的でない。

本稿では、図 3(a) のようなクローンセットの組み合

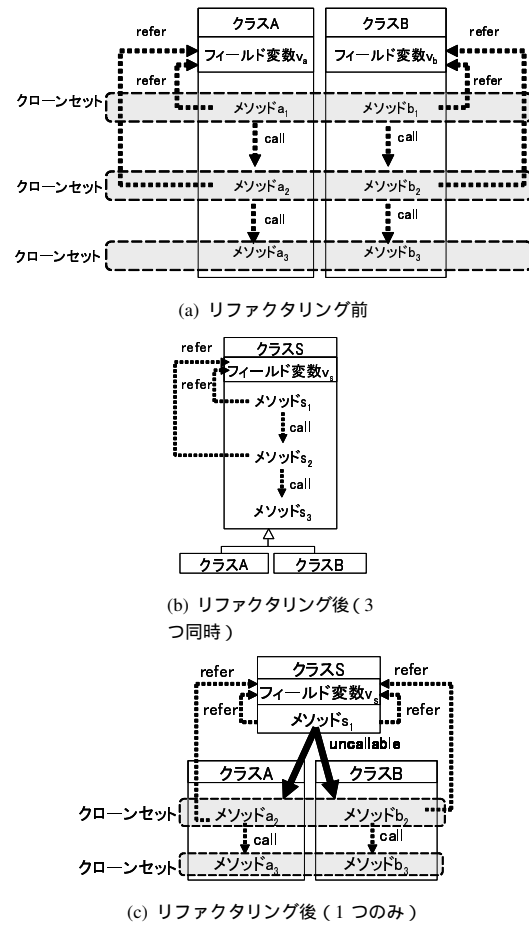


図 3 異なるクローンセット間に依存関係が存在する例
Fig. 3 Example of clone sets that have dependency relations

わせを“チェンドクローンセット”と呼び、“チェンドクローンセット”に対するリファクタリング支援手法を提案する。

3.2 チェンドクローン

チェンドクローンを定義するための準備として、メソッドチェーンを定義する。メソッドの集合が与えられたとき、それらメソッド間の依存関係を表す有向グラフが連結グラフになるなら、そのメソッドの集合をメソッドチェーンと定義する。

ここで扱う依存関係は、以下の 2 種類である。

- (1) メソッドの呼出関係
- (2) 同一フィールド変数の共有 (参照または代入)

図 4(a) は呼出関係を含むメソッドチェーンの例である。この例では、メソッド a が b を、メソッド b が c を呼び出している。また、図 4(b) は、図 4(a) のメソッドチェーンの依存関係をラベル付き有向グラフで表したものである。有向辺に付随しているラベル “call” は、

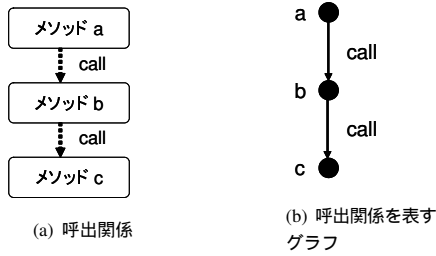


図 4 呼出関係と呼出関係を表すグラフ

Fig. 4 Method chain whose dependency relations are method invocations

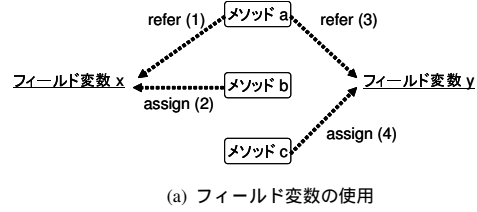
依存関係の種類がメソッドの呼出関係であることを表している。例えば、メソッド a が b を呼び出しているとき、有向辺 (a, b) を引き、ラベル “call” を付ける。なお、1つのメソッドが同一のメソッドを2回呼び出している場合は、それらメソッド間に呼出関係を表す有向辺を2本追加する。

図 5(a) はフィールド変数を使用しているメソッドチェーンの例である。図 5(a) の有向辺に付属しているラベル “refer” はフィールド変数の参照を表しており、“assign” はフィールド変数への代入を表している。図 5(a) のメソッド a と b は、フィールド変数 x を共有しており、メソッド a, c はフィールド変数 y を共有している。また、図 5(b) は、図 5(a) のメソッドチェーンに含まれるメソッド間の共有関係をラベル付き有向グラフで表したものである。有向辺に付属しているラベル “ $R_x(refer)$ ” はフィールド変数 x への参照による共有関係、ラベル “ $A_y(assign)$ ” は、フィールド変数 y への代入による共有関係を表している。例えば、メソッド a が変数 x を参照し、かつメソッド b が変数 x を使用（参照または代入）しているとき、有向辺 (a, b) を引き、ラベル “ R_x ” を付ける。また、メソッド c が変数 y に代入し、かつメソッド a が変数 y を使用しているとき、有向辺 (c, a) を引き、ラベル “ A_y ” を付ける。

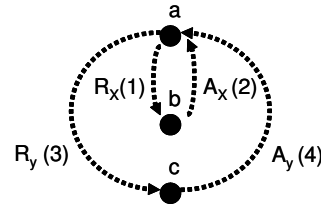
次に、メソッドチェーンを用いてチェンドクローンを定義する。2つのメソッドチェーンが互いにチェンドクローンとなるのは、各メソッドチェーンが持つ依存関係のグラフが同形であり、対応する頂点（メソッド）が同一クローンセットに含まれ、対応する辺（依存関係）のラベルは等しいときである。また、互いにチェンドクローンであるメソッドチェーンの同値類を、チェンドクローンセットと呼ぶ。

3.3 チェンドクローンセットに対するリファクタリング

ここでは、チェンドクローンセットに対して考え



(a) フィールド変数の使用



(b) フィールド変数の共有関係を表すグラフ

図 5 フィールド変数の使用と共有関係を表すグラフ（括弧内の数字は対応する関係を表す）

Fig. 5 Chained method whose dependency relations are sharing variables

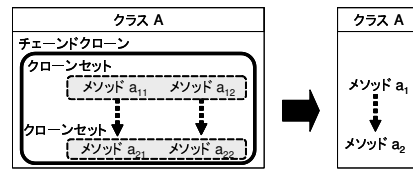


図 6 ケース 1

Fig. 6 Case 1

られるリファクタリングについて説明するため、適用可能なリファクタリングパターンが異なる4つのケースを紹介する。

ケース 1 は、チェンドクローンセットが1つのクラスに含まれている場合である。ケース 1 では、その1つのクラス内にクローンセットを集約可能である。図 6 は、ケース 1 のチェンドクローンセットに対するリファクタリングの例である。互いにクローンであるメソッド a_{11} とメソッド a_{12} を集約しメソッド a_1 とし、同様に互いにクローンであるメソッド a_{21} とメソッド a_{22} を集約しメソッド a_2 としている。

ケース 2 は、チェンドクローンセットが以下の2つの条件を満たす場合である。

- チェンドクローンセットに含まれるメソッドは、全て兄弟クラスに属する。
- 各メソッドチェーンは、それぞれ1つのクラスに含まれている。

ケース 2 は、“Pull Up Method” パターンを適用することで、リファクタリングできる。つまり、兄弟クラス

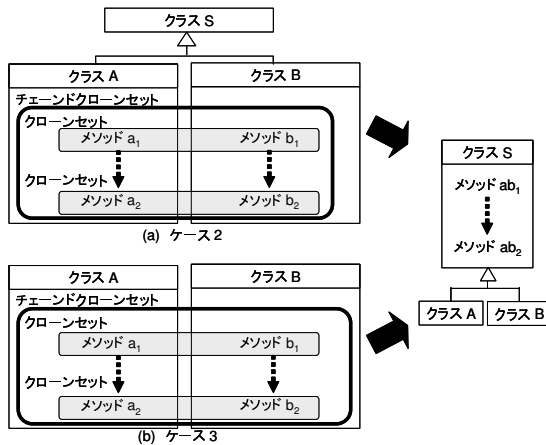


図7 ケース2とケース3
Fig.7 Case 2 and Case 3

にまたがって存在するクローンセットを、親クラスに集約することでリファクタリングできる。図7(a)は、ケース2に対するリファクタリングの例である。この例では、兄弟クラスであるクラスA、Bにまたがって存在する2つのクローンセットを、親クラスSに作成したメソッドに集約している。具体的には、クラスAのメソッド a_1 とクラスBのメソッド b_1 を集約し、親クラスSのメソッド ab_1 とし、同様にクラスAのメソッド a_2 とクラスBのメソッド b_2 を集約し、親クラスSのメソッド ab_2 としている。

ケース3は、チェンドクローンセットが以下の2つの条件を満たす場合である。

- チェンドクローンセットに含まれるメソッドを持つクラスは、いずれも共通の親クラスを持たない。
- 各メソッドチェーンは、それぞれ1つのクラスに包含されている。

ケース3は、“Extract SuperClass”パターンを適用することで、リファクタリングできる。つまり、チェンドクローンセットに含まれるメソッドを持つクラスに対して、共通の親クラスを作成し、クラス間をまたがって存在するクローンセットを、新たに作成した親クラスに集約することでリファクタリングできる。図7(b)は、ケース3に対するリファクタリングの例である。この例では、まず共通の親クラスを持たない2つのクラスA、Bに、共通の親クラスSを作成している。その後、ケース2と同様に兄弟クラスとなったクラスA、Bにまたがって存在する2つのクローンセットを、親クラスSに作成したメソッドに集約している。

ケース4は、チェンドクローンセットが以下の条件を満たす場合である。

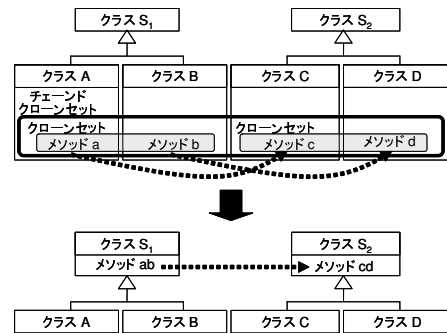


図8 ケース4
Fig.8 Case 4

- 各メソッドチェーンは、複数のクラスにまたがって存在する。つまり依存関係が複数のクラス間にまたがっている。

ケース4は、チェンドクローンセット単位でリファクタリングできない場合である。だが、チェンドクローンセットを複数のクローンセットとして扱い、それぞれの親クラスに集約することでリファクタリングできるため、クローンセット単位でのリファクタリングを検討すべきである。図8は、ケース4にクローンセット単位でのリファクタリングを適用した例である。この例では、兄弟クラスであるクラスA、Bにまたがって存在する2つのクローンセットを、それぞれの親クラスに作成したメソッドに集約している。具体的には、クラスAのメソッド a とクラスBのメソッド b を集約し、親クラス S_1 のメソッド ab とし、同様にクラスCのメソッド c とクラスDのメソッド d を集約し、親クラス S_2 のメソッド cd としている。

3.4 チェンドクローンセットの分類

前節の4つのケースのように、チェンドクローンセットを分類する。前節の4つのケースには、それぞれ適合するための条件があった。それらは、次の2つである。

- C1 チェンドクローンセットに含まれるメソッドが所属するクラス間の関係についての条件
- C2 メソッドチェーンに含まれるメソッドが所属するクラス間の関係についての条件

ここでのクラス間の関係とは、クラス階層上の関係のことである。クラス間の関係は、次に3つに分類できる。

- R1 全ての同一クラス
- R2 共通の祖先クラスを持つ
- R3 共通の祖先クラスを持たないクラス

条件の種類とクラス間の関係を組み合わせることにより、チェンドクローンセットを表1のように分類

できる。

次の4つの分類については、以下に示すリファクタリングを行うことができると考えられる。

分類1 前節のケース1である。図6の例のように、チェンドクローンセットを包含しているクラスに、全てのクローンセットを集約することができる。

分類2 前節のケース2である。図7(a)の例のように、“Pull Up Method”パターンを適用できる。

分類3 前節のケース3である。図7(b)の例のように、“Extract SuperClass”パターンを適用できる。

分類4 前節のケース4である。図8からわかるように、チェンドクローンセット単位でのリファクタリングを行うことが出来ないが、クローンセット単位でのリファクタリングを検討すべきである。

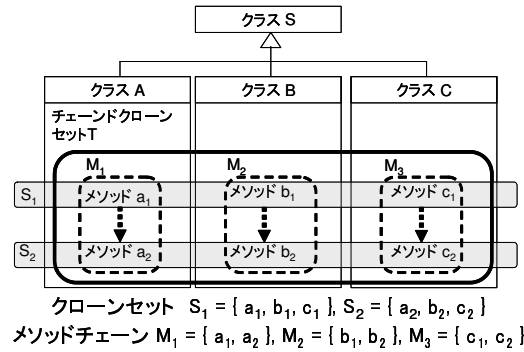
3.5 チェンドクローンセットの分類を目的としたメトリクス

ここでは、チェンドクローンセットの分類を行うためのメトリクスを2つ提案する。1つはC1を評価するメトリクス、もう1つはC2を評価するメトリクスである。これらメトリクスは、メソッド間のクラス階層上における関係を表す。この関係は、2.3節で述べた $DCH(S)$ メトリクス(クローンセット S に含まれる各コード片間のクラス階層内における最大の距離)によって表すことができる。

まず、 $DCH(S)$ メトリクスを用いて、C1を表す $DCHS(T)$ メトリクスを定義する。チェンドクローンセット T を n 個のクローンセット S_1, S_2, \dots, S_n に分割する(すなわち、 $S_1 \cup S_2 \cup \dots \cup S_n = T, S_i \cap S_j = \emptyset, 1 \leq i \leq n, 1 \leq j \leq n, i \neq j$)。更に、クローンセット S_i には、複数のメソッドが含まれるとする。このとき、 $DCHS(T)$ メトリクスの定義は以下のようになる。

$$DCHS(T) = \max\{DCH(S_1), \dots, DCH(S_n)\}$$

同様に $DCH(S)$ メトリクスを用いて、C2を表す $DCHD(T)$ メトリクスを定義する。チェンドクローンセット T 中には、 n 個のメソッドチェーン



$$DCHS(T) = \max\{DCH(S_1), DCH(S_2)\} = 1$$

$$DCHD(T) = \max\{DCH(M_1), DCH(M_2), DCH(M_3)\} = 0$$

図9 提案するメトリクスの算出例
 Fig. 9 Example of proposed metrics calculation

M_1, M_2, \dots, M_n が含まれるとする。更に、メソッドチェーン M_i には、複数のメソッドが含まれるとする。このとき、 $DCHD(T)$ メトリクスの定義は以下のようになる。

$$DCHD(T) = \max\{DCH(M_1), \dots, DCH(M_n)\}$$

図9は、提案する2つのメトリクスの算出例である。ここでは、分類2のチェンドクローンセットを例として用いる。このチェンドクローンセットには、クラスA, B, Cにまたがって2つのクローンセット S_1, S_2 が存在する。各クローンセットについてそれぞれ $DCH(S_1), DCH(S_2)$ を求めると、クラスA, B, Cは共通の直接の親クラスSを持っているため両者とも1になる。 $DCHS(T)$ の値は、これらの最大値の1である。また、このチェンドクローンセットには、3つのメソッドチェーン M_1, M_2, M_3 が含まれている。各メソッドチェーンについてはそれぞれ $DCH(M_1), DCH(M_2), DCH(M_3)$ を求めると、各メソッドチェーンはそれぞれ1つのクラスに包含されているため全て0になる。よって、 $DCHD(T)$ の値は、これらの最大値の0である。

3.6 実装

提案手法をAriesのコンポーネントの1つとして実装した。具体的には、Ariesに対して、以下の3つの機能を追加した。

- (F1) チェンドクローンセットの検出機能
 - (F2) 提案したメトリクスの算出機能
 - (F3) チェンドクローンセットおよびメトリクス値の表示機能
- (F1)は、まずCCFinderおよびAriesを用いて各ク

表1 チェンドクローンセットの分類
 Table 1 Categorization of chained clone sets

C1 \ C2	R1	R2	R3
	R1	分類1	分類4
R2	分類2		
R3	分類3		

ローンセットが含むコード片を検出する。次に、それらコード片を対象に、メソッド呼出関係と変数の共有関係を表すグラフを構築する。その後、構築したプログラム依存グラフに含まれる部分グラフから同形グラフを検出することにより、チェンドクローンセットを検出する。(F2)は、AriesのDCH(S)メトリクスを計算する機能を拡張した。(F3)を実現するために、チェンドクローンセットを閲覧するビューを追加した。

4. ケーススタディ

4.1 概要

提案手法の有効性を確かめるため、ケーススタディを行った。具体的には、以下の2つを確認した。

- クローンセット単位の検出と比較して、検出できたチェンドクローンセットの規模が大きいか
- クローンセット単位でのリファクタリングと比較して、容易にリファクタリングできているか

なお、この章におけるチェンドクローンセットは、極大チェンドクローンセットを指す。

適用対象は、次の2つのオープンソースソフトウェアである。

- ANTLR 2.7.4¹¹⁾(4.7万行, 285クラス)
- JBoss 3.2.6¹²⁾(64万行, 3364クラス)

ANTLRは、3つのプログラミング言語(Java, C#, C++)に対応したコンパイラ・コンパイラである。JBossは、J2EEアプリケーションサーバである。

4.2 チェンドクローンセットの検出

前述の2つのソフトウェアに対し、提案手法に基づくチェンドクローンセットの検出、および従来手法に基づくクローンセットの検出を行った。提案手法ではメソッド単位のコードクローンのみを扱うため、従来手法に基づく検出でもメソッド単位のクローンセットのみを対象とした。また、CCFinderが検出するコードクローンの最小トークン数は30に設定した。

検出結果の比較を行うために、2つの評価基準として、メソッド数と、メソッド行数を用いる。ここで、メソッド数は、検出単位毎(クローンセット毎やチェンドクローンセット毎)に含まれるメソッド数を求め、総クローンセット数や各分類に属する全てのチェンドクローンセット数で除算した値とした。メソッド行数は、検出単位毎に最長メソッドの行数を求め、総クローンセット数や各分類に属する全てのチェンドク

与えられたチェンドクローンセットを真に包含する如何なるチェンドクローンセットも存在しないとき、そのチェンドクローンセットを極大チェンドクローンセットと呼ぶ。

ローンセット数で除算した値とした。

設定した評価基準に基づいて、クローンセットの検出結果(表2)と分類1, 2, 3に属したチェンドクローンセットの検出結果(表3(a))を比較する。

まず、ANTLRでは分類2に属したチェンドクローンセットのメソッド数やメソッド行数が極めて大きかった。分類2のメソッド数はクローンセットの8.3(19/2.3)倍、メソッド行数は5.0(54.0/10.8)倍であった。一方、分類1, 3に属したチェンドクローンセットのメソッド数は両者ともクローンセットの1.7(4.0/2.3)倍、メソッド行数はそれぞれ2.6(27.7/10.8)倍, 3.2(35.0/10.8)倍であった。分類2に属するチェンドクローンセットの多くは、図10のような、Java, C#, C++に対応した出力を行う箇所から検出された。これら言語に対応した出力処理は類似しており、大量のコードクローンを含んでいた。

次に、JBossの結果を見てみると、分類1, 2, 3に属したチェンドクローンセットのメソッド数はそれぞれクローンセットの1.8倍~2.8倍、メソッド行数がそれぞれ2.4倍~3.2倍であった。これらの結果から、チェンドクローンセットの規模がクローンセットに比べて大きいことがわかる。続いて、分類1, 2, 3に属したチェンドクローンセットの検出結果(表3(a))と、分類4に属したチェンドクローンセットの検出結果(表3(b))を比較する。3.4節で述べたように、表3(a)で示した分類1, 2, 3は、チェンドクローンセット単位でのリファクタリング可能であるが、表3(b)で示した分類4はチェンドクローンセット単位でのリファクタリングが出来ない。分類1, 2, 3と比較して分類4に属するチェンドクローンセッ

表2 クローンセットの検出結果

Table 2 Detection result of clone sets

検出数		メソッド数		メソッド行数	
ANTLR	JBoss	ANTLR	JBoss	ANTLR	JBoss
152	377	2.3	2.4	10.8	6.63

表3 チェンドクローンセットの検出結果

Table 3 Detection result of chained clone set

(a) 分類1, 2, 3

分類	検出数		メソッド数		メソッド行数	
	ANTLR	JBoss	ANTLR	JBoss	ANTLR	JBoss
1	3	16	4.0	5.8	27.7	16.2
2	6	17	19	4.5	54.0	17.1
3	1	13	4.0	6.8	35.0	21.5

(b) 分類4

検出数		メソッド数		メソッド行数	
ANTLR	JBoss	ANTLR	JBoss	ANTLR	JBoss
0	4	-	19	-	54.5

トは少ないことがわかる．特に，ANTLR からは検出されなかった．一方，JBoss では 4 種類のチェンドクローンセットが検出された．これらのメソッド数、メソッド行数は、分類 1~3 に比べて数倍の大きさになっている．

4.3 チェンドクローンセットに対するリファクタリングの例

ANTLR から検出された全てのチェンドクローンセットと JBoss から検出された 8 つのチェンドクローンセットを対象として、提案手法に基づくリファクタリングを行った．ここでは、それらの中から 2 つの例を紹介する．

まず、図 10、図 12 で示すチェンドクローンセットに対し、提案手法により提示されたリファクタリングパターンを適用できることを確認した．図 10 は分類 2 であるから、“Pull Up Method” パターンを適用した．その結果、図 11 のようになった．ANTLR パッケージ中の examples ディレクトリ以下にある全てのテストケース（計 86 ファイル、文法ファイル）を用いて回帰テストを行い、外部的振る舞い（出力結果）が変化していないことを確認した．また、図 12 は分類 3 であるから、“Extract SuperClass” パターンを適用した．その結果、図 13 のようになった．JBoss パッケージの testsuite ディレクトリ以下にある全テストケース（計 65 ファイル、JUnit フレームワーク¹³）を用いて回帰テストを行い、外部的振る舞いが変化していないことを確認した．

更に、従来手法に基づいて、チェンドクローンセットを構成するクローンセットに対し集約を試みると、工夫が必要となる場合があることを確認した．具体的には、図 10、図 12 からそれぞれクローンセットを 1 つ選び、従来手法により提示されたリファクタリングパターンの適用をした．まず、図 10 のクローンセット 1 に対し集約を試みると、提示されたリファクタリングパターンの適用に加えて、クローンセット 2 のメソッドに対応する抽象メソッドを CodeGenerator クラスに追加する必要があった．なお、クローンセット 2 に対して集約を試みた場合も同様であることが確認できた．また、図 12 のクローンセット 2 を新たに作成した親クラスに集約を試みると、提示されたリファクタリングパターンの適用に加えて、クローンセット 2 のメソッドに対応する抽象メソッドを親クラスに追加する必要があった．

4.4 考 察

ここでは、今回のケーススタディに基づいて、提案手法の妥当性・制限等について考察を行う．

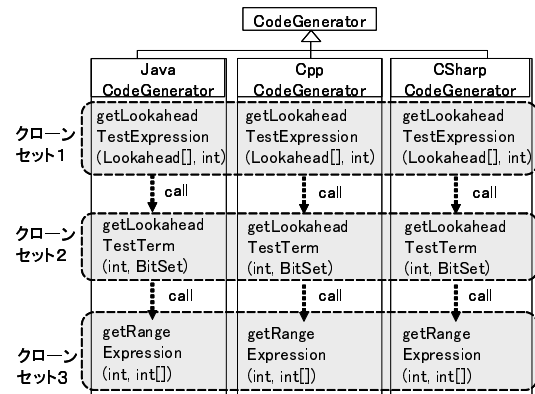


図 10 ANTLR から検出されたチェンドクローンセットの例
Fig. 10 Example of a chained clone set in ANTLR

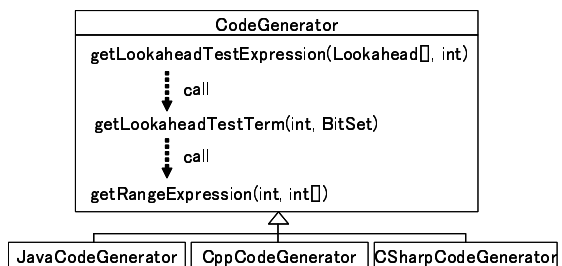


図 11 図 10 のチェンドクローンセットをリファクタリングした例
Fig. 11 Example of refactoring of a chained clone set of Fig. 10

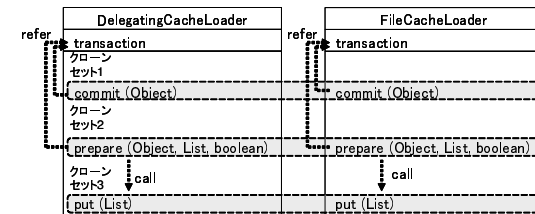


図 12 JBoss から検出されたチェンドクローンセットの例
Fig. 12 Example of a chained clone set in JBoss

- (1) 対象ソフトウェア ケーススタディでは、Java 言語で開発された 2 つのオープンソースソフトウェアを対象として有効性の確認を行った．オブジェクト指向型言語であれば、Java 言語以外で開発されたソフトウェアであっても適用可能であると考えられる．なぜなら、Java 言語以外のオブジェクト指向型言語で開発されたソフトウェアであっても、リファクタリングや依存関係解析を適用可能だからである．今後、Java 言語以外で開発されたソフトウェアや商用ソフトウェアなど、様々な種類のソフトウェアを対象に有効性の評価を行う必要がある．
- (2) 被支援者の知識、経験 今回のケーススタディに

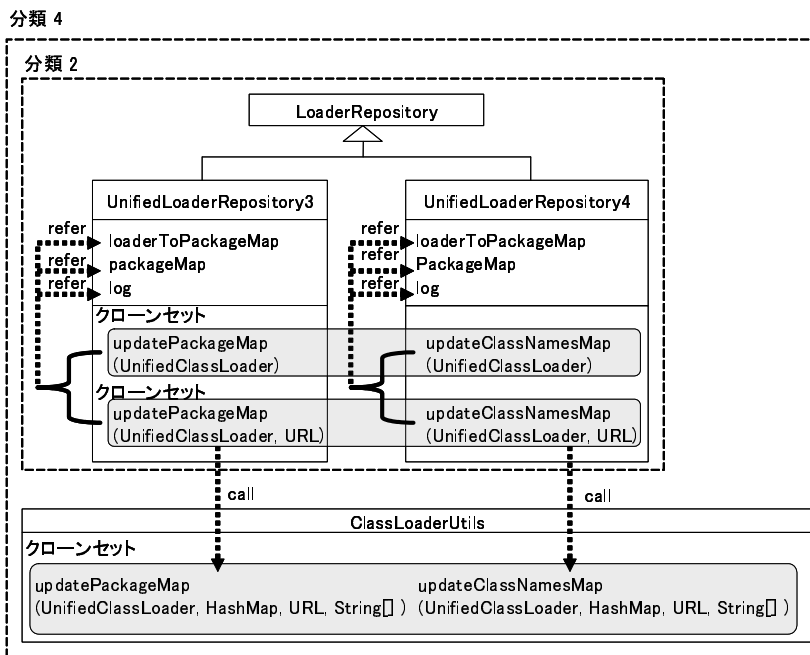


図 14 分類 4 のチェンドクローンセットの例

Fig. 14 Example of Category 4

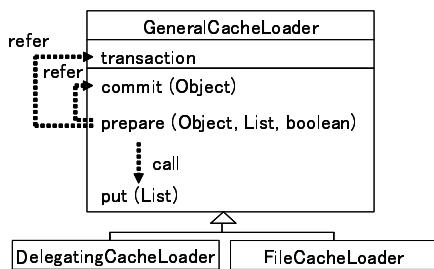


図 13 図 12 のチェンドクローンセットをリファクタリングした例

Fig. 13 Example of refactoring of a chained clone set of Fig. 12

おけるリファクタリング作業は、全て著者が行った。そのため、手法の詳細を知らない人やリファクタリング経験が少ない人に対しても、有効な支援が行えるかどうか評価する必要がある。例えば、一般の開発者に本稿のツールを使用してもらい、リファクタリングにかかった時間を計測することで効率を評価するということが考えられる。

- (3) 対象とするコード片や依存関係の種類 提案手法は、メソッドより小さい単位のコード片からなるクローンセットは対象としていない。また、メソッドの呼出関係および変数の共用による依存関係のみを扱っているため、その他のデータ依存関係や制御依存関係を対象としていない。今後、対象とするコード片や依存関係を増やすことで検出可能なチェンドクローンセットの規模を大きく

し、有効性の評価を行う必要がある。

- (4) 分類 4 のチェンドクローンセットへの支援

JBoss から分類 2 のチェンドクローンセットを内包する分類 4 のチェンドクローンセットが検出された (図 14)。提案手法は、これに対しクローンセット単位でのリファクタリングを提示するが、内包された分類 2 のチェンドクローンセットに“Pull Up Method”リファクタリングパターンを適用可能である。このような場合は、内包されたチェンドクローンセットに対してリファクタリングパターンの提示を行うべきであると考えられる。

5. 関連研究

CCFinder や Balazinska²⁾ らの手法を用いることにより、クローンセットの検出を行うことはできる。本稿の手法では、チェンドクローンセット単位でのリファクタリングを提示することにより、大規模なリファクタリングを実現することが出来た。また、CCFinder が検出するクローンセットには容易にリファクタリングできないものが含まれており、Balazinska らの手法も同様と考えられる。本稿では、それら容易にリファクタリングできないクローンセットを組み合わせることによって容易にリファクタリングできる場合があることを示し、それらクローンセットに対するリファクタリング

手法を提案した。

Komondoor⁴⁾らの手法は、プログラムスライシング技術を用いて、ソースコードからプログラム依存グラフを構築し、そのグラフ上で同形である箇所をコードクローンとして検出している。よって、本稿で用いているCCFinderが検出できないコードクローン（一部の文の出現順序が異なっている reodered clone 等）を検出することができる。しかし、プログラム依存グラフの構築にかかる計算コストは非常に大きいため、大規模ソフトウェアへの適用は現実的でない。本稿の手法は、CCFinderが検出したコードクローンに対して、メソッド呼出関係と変数の利用関係のみを解析しているため、ケーススタディで示した規模のソフトウェアに適用可能である。実際に、ANTLRを対象としたチェンドクローンセットの検出を約1分4秒で行うことができた。

6. ま と め

本稿では、クローンセットに含まれるメソッド間の依存関係に着目し、チェンドクローンセットを定義した。そして、チェンドクローンセットに対しリファクタリングパターンを提示するためのメトリクスを提案した。最後に、提案手法をリファクタリング支援ツールとして実装し、2つのオープンソースソフトウェアに適用することで、有効性の評価を行った。有効性の評価として、チェンドクローンセットの規模がクローンセットと比べて大きいこと、およびチェンドクローンセットのリファクタリングがクローンセット単位のリファクタリングと比べて容易であることを確認した。

今後の課題としては、様々なソフトウェアを対象とした有効性の評価、チェンドクローンセットの中に異なる分類のチェンドクローンセットが包含されている場合への対処、対象とする依存関係やコード片の拡大が挙げられる。

謝辞 本研究は一部、文部科学省リーディングプロジェクト「eSociety 基盤ソフトウェアの総合開発」、日本学術振興会の科研費（課題番号：17200001）の支援を受けている。

参 考 文 献

- 1) Baker, B. S.: A Program for Identifying Duplicated Code, *Proc. Computing Science and Statistics*, Vol.6, pp.49–57 (1992).

- 2) Balazinska, M., Merlo, E., Dagenais, M., Lague, B. and Kontogiannis, K.: Advanced clone-analysis to support object-oriented system refactoring, *Proc. Working Conference on Reverse Engineering (WCRE2000)*, pp.98–107 (2000).
- 3) Baxter, I., Yahin, A., Moura, L., Anna, M. and Bier, L.: Clone Detection Using Abstract Syntax Trees, *Proc. International Conference on Software Maintenance (ICSM98)*, pp.368–377 (1998).
- 4) Komondoor, R. and Horwitz, S.: Using Slicing to Identify Duplication in Source Code, *Proc. International Static Analysis Symposium (SAS2001)*, pp.40–56 (2001).
- 5) Krinke, J.: Identifying Similar Code with Program Dependence Graphs, *Proc. Working Conference on Reverse Engineering (WCRE2001)*, pp.301–309 (2001).
- 6) Fowler, M.: *Refactoring: improving the design of existing code*, Addison Wesley (1999).
- 7) Kamiya, T., Kusumoto, S. and Inoue, K.: CCFinder: A multi-linguistic token-based code clone detection system for large scale source code, *IEEE Transactions on Software Engineering*, Vol. 28, No. 7, pp. 654–670 (2002).
- 8) Higo, Y., Kamiya, T., Kusumoto, S. and Inoue, K.: ARIES: Refactoring Support Environment Based on Code Clone Analysis, *Proc. IASTED International Conference on Software Engineering and Applications (SEA2004)*, pp.222–229 (2004).
- 9) 井上克郎, 神谷年洋, 楠本真二: コードクローン検出法, *コンピュータソフトウェア*, Vol.18, No.5, pp.47–54 (2001).
- 10) 門田暁人, 佐藤慎一, 神谷年洋, 松本健一: コードクローンに基づくレガシーソフトウェアの品質の分析, *情報処理学会論文誌*, Vol.44, No.8, pp. 2178–2188 (2003).
- 11) ANTLR: <http://www.antlr.org>.
- 12) JBoss: <http://www.jboss.org>.
- 13) JUnit: <http://www.junit.org>.

(平成 xx 年 x 月 x 日受付)

(平成 xx 年 x 月 x 日採録)

吉田 則裕 (学生会員)

平 16 九工大・情報工・知能情報
卒・平 18 阪大大学院博士前期課程
修了。現在同大学院博士後期課程
1 年。コードクローン分析の研究に
従事。人工知能学会会員。



肥後 芳樹 (正会員)

平 14 阪大・基礎工・情報中退。平 18 同大学院博士後期課程修了。現在日本学術振興会特別研究員。コードクローン分析・リファクタリング支援の研究に従事。



楠本 真二 (正会員)

昭 63 阪大・基礎工・情報卒。平 3 同大学院博士課程中退。同年同大・基礎工・情報・助手。平 8 同講師。平 11 同助教授。平 14 阪大・情報・コンピュータサイエンス・助教授。平 17 同教授。博士 (工学)。ソフトウェアの生産性や品質の定量的評価、プロジェクト管理に関する研究に従事。電子情報通信学会, IEEE, JFPUG, PM 各会員。



神谷 年洋

平 8 阪大・基礎工・情報中退。平 13 同大学院博士課程了。同年科学技術振興事業団研究者。平 17 産業技術総合研究所研究員。博士 (工学)。オブジェクト指向関連技術, ソフトウェア保守 (メトリクス, コードクローン), 認知科学に関する研究に従事。電子情報通信学会, 電気学会, IEEE, 各会員。



井上 克郎 (正会員)

昭 54 阪大・基礎工・情報卒。昭 59 同大学院博士課程了。同年同大・基礎工・情報・助手。昭 59 ~ 昭 61 ハワイ大マノア校・情報工学科・助教授。平 1 阪大・基礎工・情報・講師。平 3 同学科・助教授。平 7 同学科・教授。工学博士。平 14 阪大・情報・コンピュータサイエンス・教授。ソフトウェア工学の研究に従事。電子情報通信学会, 日本ソフトウェア科学会, IEEE, ACM 各会員。