

# Anti-unification algorithms and their applications in program analysis

Peter E. Bulychev, Egor V. Kostylev, Vladimir A. Zakharov

Faculty of Computational Mathematics and Cybernetics,  
Moscow State University, Moscow, RU-119899, Russia  
([peter.bulychev@gmail.com](mailto:peter.bulychev@gmail.com), [jegor\\_kostylev@hotmail.com](mailto:jegor_kostylev@hotmail.com), [zakh@cs.msu.su](mailto:zakh@cs.msu.su))

**Abstract.** A term  $t$  is called a *template* of terms  $t_1$  and  $t_2$  iff  $t_1 = t\eta_1$  and  $t_2 = t\eta_2$ , for some substitutions  $\eta_1$  and  $\eta_2$ . A template  $t$  of  $t_1$  and  $t_2$  is called *the most specific* iff for any template  $t'$  of  $t_1$  and  $t_2$  there exists a substitution  $\xi$  such that  $t = t'\xi$ . The anti-unification problem is that of computing the most specific template of two given terms. This problem is dual to the well-known unification problem, which is the computing of the most general instance of terms. Unification is used extensively in automatic theorem proving and logic programming. We believe that anti-unification algorithms may have wide applications in program analysis. In this paper we present an efficient algorithm for computing the most specific templates of terms represented by labeled directed acyclic graphs and estimate the complexity of the anti-unification problem. We also describe techniques for invariant generation and software clone detection, which are based on the concepts of the most specific templates and anti-unification.

The anti-unification problem is that of finding the most specific template (pattern) of two terms. This problem has been first considered by G.D. Plotkin [15] and J. Reynolds [16]. It is dual to the well-known unification problem, which is the computing of the most general instance of terms. Unification is extensively used in automatic theorem proving, logic programming, typed lambda calculus, term rewriting, etc. The unification problem has been studied thoroughly in many papers (see [1] for survey) and a wide variety of efficient unification algorithms have been developed by many authors (see, for example, [10, 14, 17]).

The anti-unification problem attracted far less attention. The algebraic properties of anti-unification operation have been studied in [6, 13]. To the extent of our knowledge all anti-unification algorithms introduced so far (see [4, 9, 15, 16, 19]) deal with tree-like representation of terms. It is obvious that the anti-unification problem for terms represented by labeled trees can be solved in linear time. In [4] it has been shown that the anti-unification problem for labeled trees of size  $n$  can be solved in time  $O(\log^2 n)$  using  $O(n/\log^2 n)$  processors on EREW PRAM model of parallel computation, and in time  $O(\log n)$  using  $O(n/\log n)$  processors on CRCW PRAM model of computation. In [9] it has been proved that the anti-unification problem for labeled trees is in NC<sup>2</sup>. However, if one needs to deal with large sets of sizable terms then it is more suitable to represent such sets of terms by labeled directed acyclic graphs (dags). One of the aims of this paper is to develop an efficient algorithm for computing the most specific templates and estimate the complexity of anti-unification problem for terms represented by labeled dags.

Only few papers concern the application of anti-unification. R. Gluck and M.H. Sorensen used anti-unification (the most specific generalization) of terms to guarantee the termination of a positive supercompilation algorithm developed in their paper [19]. The utility of anti-unification in the setting of symbolic mathematical computing has been studied in [12, 21]. We believe that anti-unification algorithms may have wide applications in many areas of computer science and software engineering. In this paper we study the perspectives of using the concepts of the most specific templates and anti-unification for invariant generation and software clone detection.

The generation of invariants is the key technique in the analysis and verification of programs, since the effectiveness of automated program verification is highly sensitive to the ease with which

invariants, even trivial ones, can be automatically deduced. Much efforts (see [8, 20]) are directed toward the development of powerful invariant generating techniques for particular classes of programs. As opposed to these attempts, we present a light-weight technique for invariant generation. Our anti-unification based algorithm for invariant generation operates on the syntactic level. Therefore, it is of little sensitivity to program semantics and can reveal only trivial invariants. But, due to the efficiency of anti-unification algorithms, this technique provides a way for processing large pieces of code in short time.

Anti-unification algorithms can be of significant value in software refactoring. One of the major activity in this area is the detection and extraction of duplicate code. Code duplication can be a significant drawback, leading to bad design, and increased probability of bug occurrence and propagation. Consequently, duplicate code detectors are a useful class of program analysis tools (see [18]). We describe a simple anti-unification based algorithm for finding software clones. The algorithm checks fragments of code (sequences of program statements) and assign two pieces of code to the same clone if they are not too different from their most specific template.

## 1 Preliminaries

Given an alphabet of variables  $\mathcal{Y}$  and an alphabet  $\mathcal{F} = \{f_1^{(m_1)}, \dots, f_k^{(m_k)}\}$  of functional symbols of arities  $m_1, \dots, m_k$ , we define the set of terms  $Term[\mathcal{Y}, \mathcal{F}]$  as the smallest set of expressions that contains  $\mathcal{Y}$  and satisfies the following property: if  $f^{(m)} \in \mathcal{F}$  and  $t_1, \dots, t_m \in Term[\mathcal{Y}, \mathcal{F}]$  then  $f^{(m)}(t_1, \dots, t_m) \in Term[\mathcal{Y}, \mathcal{F}]$ .

Let  $\mathcal{X} = \{x_1, \dots, x_n\}$  and  $\mathcal{Y} = \{y_1, y_2, \dots\}$  be two sets of variables. Then a  $\mathcal{X}$ - $\mathcal{Y}$ -substitution is a mapping  $\theta : \mathcal{X} \rightarrow Term[\mathcal{Y}, \mathcal{F}]$ . Usually a substitution  $\theta$  is represented as the set of bindings  $\theta = \{x/\theta(x) : x \neq \theta(x)\}$ . We denote the set of all  $\mathcal{X}$ - $\mathcal{Y}$ -substitutions by  $Subst[\mathcal{X}, \mathcal{Y}, \mathcal{F}]$ . An *application* of a substitution  $\theta$  to a term  $t = t(x_1, \dots, x_n)$  yields the term  $t(\theta(x_1), \dots, \theta(x_n))$  obtained from  $t$  by replacing every variable  $x_i$ ,  $1 \leq i \leq n$ , with the term  $\theta(x_i)$ . The *composition*  $\theta\xi$  of substitutions  $\theta \in Subst[\mathcal{X}, \mathcal{Y}, \mathcal{F}]$  and  $\xi \in Subst[\mathcal{Y}, \mathcal{Y}, \mathcal{F}]$  is defined as follows:  $(\theta\xi)(x) = (\theta(x))\xi$  holds for every variable  $x$  from  $\mathcal{X}$ . A *renaming* is any bijection  $\rho$  from  $\mathcal{Y}$  onto  $\mathcal{Y}$ .

With the notion of composition of substitutions at hand, we can define a quasi-order  $\sqsubseteq$  and an equivalence relation  $\sim$  on the set of substitutions  $Subst[\mathcal{X}, \mathcal{Y}, \mathcal{F}]$ : a relation  $\theta_1 \sqsubseteq \theta_2$  holds iff there exists  $\xi \in Subst[\mathcal{Y}, \mathcal{Y}, \mathcal{F}]$  such that  $\theta_2 = \theta_1\xi$ , and  $\theta_1 \sim \theta_2$  holds iff  $\theta_2 = \theta_1\rho$  holds for some renaming  $\rho$ . Quasi-order  $\sqsubseteq$  induces the partial order  $\preceq$  on the quotient set  $Subst^{\sim}[\mathcal{X}, \mathcal{Y}]$ . Poset  $(Subst^{\sim}[\mathcal{X}, \mathcal{Y}], \preceq)$  has been studied in [6, 13]. This poset is a complete lattice, the least element of the lattice is the equivalence class of the *empty substitution*  $\varepsilon = \{x_1/y_1, \dots, x_n/y_n\}$ . The least upper bound  $\theta_1^{\sim} \uparrow \theta_2^{\sim}$  is called the *most general instance* of substitutions  $\theta_1$  and  $\theta_2$ , whereas the greatest lower bound  $\theta_1^{\sim} \downarrow \theta_2^{\sim}$  is called the *most specific template* of  $\theta_1$  and  $\theta_2$ . The operation  $\downarrow$  of computing the most specific template of substitution is called *anti-unification* (or *generalization*). For the sake of simplicity we will often skip the superscript  $\sim$  in our notation. It is easy to check (see [6, 13]) that composition of substitutions is left-distributive over anti-unification, i. e.  $\eta(\theta_1 \downarrow \theta_2) = \eta\theta_1 \downarrow \eta\theta_2$ .

Anti-unification operation can be naturally extended to the terms. The term  $t$  is called the most specific template of terms  $t_1$  and  $t_2$  ( $t = t_1 \downarrow t_2$  in symbols) if  $\{x/t\} = \{x/t_1\} \downarrow \{x/t_2\}$ .

## 2 Anti-unification algorithms

In this section we study the complexity of anti-unification of substitutions represented as labeled directed acyclic graphs (dags). Dags are the most suitable structures for succinct representation of substitutions. A node  $V$  that has no incoming arcs is called a *root* of a dag. A labeled single-rooted dag  $\mathcal{G}(t)$  associated with a term  $t$ ,  $t \in Term[\mathcal{Y}, \mathcal{F}]$ , is arranged as follows. If  $t$  is a constant or a variable then  $\mathcal{G}(t)$  consists of single node labeled with  $t$  (this node is the root of  $\mathcal{G}(t)$ ). If  $t = f^{(m)}(t_1, \dots, t_m)$  then the root of  $\mathcal{G}(t)$  is labeled with  $f^{(m)}$  and has  $m$  outgoing arcs that are

---

```

procedure MST ( $\mathcal{G}(\theta'), \mathcal{G}(\theta'')$ )
  set  $\mathcal{G} := \mathcal{G}(\{x_1/y_1, \dots, x_n/y_n\})$ ,  $\mathcal{Q} := \emptyset$ ,  $\widehat{\mathcal{Q}} := \emptyset$ ;
  for  $i := 1$  to  $n$  do
    if exists  $y$ , such that  $(y, V_{\mathcal{G}(\theta')}(x_i), V_{\mathcal{G}(\theta'')}(x_i)) \in \mathcal{Q}$ 
      then remove  $V_{\mathcal{G}_\eta}(x_i)$  from  $\mathcal{G}_\eta$ ; set  $mark(mem(y)) := x_i$ 
      else set  $\mathcal{Q} := \mathcal{Q} \cup \{(y_i, V_{\mathcal{G}(\theta')}(x_i), V_{\mathcal{G}(\theta'')}(x_i))\}$ ,  $mem(y_i) := V_{\mathcal{G}}(x_i)$ 
    fi
  od;
  while exists  $(y, V', V'') \in \mathcal{Q}$  such that  $mark(V') = mark(V'')$  do
    set  $\mathcal{Q} := \mathcal{Q} \setminus \{(y, V', V'')\}$ ,  $\widehat{\mathcal{Q}} := \widehat{\mathcal{Q}} \cup \{(y, V', V'')\}$ ;
    if  $mark(V') = f^{(m)} \in \mathcal{F}$  then
      set  $mark(mem(y)) := f^{(m)}$ ;
      for  $i := 1$  to  $m$  do
        if exists  $z$  such that  $(z, desc(V', i), desc(V'', i)) \in \mathcal{Q} \cup \widehat{\mathcal{Q}}$  then
          let  $W_i = mem(z)$ 
        else
          let  $W_i$  be a new node in  $\mathcal{G}$  and  $z$  be a new variable in  $\mathcal{Y}$ ;
          set  $mark(W_i) := z$ ,  $mem(z) := W_i$ ,  $\mathcal{Q} := \mathcal{Q} \cup \{(z, desc(V', i), desc(V'', i))\}$ 
        fi;
        set  $desc(mem(y), i) := W_i$ 
      od
    fi
  od;
  return  $\mathcal{G}(\theta' \downarrow \theta'') := \mathcal{G}$ 
end of MST.

```

---

**Fig. 1.** Anti-unification algorithm MST.

labeled with integers from 1 to  $m$ ; the arc labeled with  $i$ ,  $1 \leq i \leq m$ , leads to the root of  $\mathcal{G}(t_i)$  associated with the subterm  $t_i$ . Given a node  $V$  of such a dag, we denote by  $mark(V)$  the label of  $V$  and by  $desc(V, i)$  the  $i$ -th descendant of  $V$  (i. e., the node that is at the end of the arc outgoing from  $V$  and labeled with  $i$ ). A labeled dag  $\mathcal{G}$  represents a substitution  $\theta = \{x_1/t_1, \dots, x_n/t_n\}$  if it contains dags  $\mathcal{G}(t_1), \dots, \mathcal{G}(t_n)$  associated with the terms  $t_1, \dots, t_n$  as subgraphs, and the root of each subgraph  $\mathcal{G}(t_i)$  has an extra label  $x_i$ . We denote these roots by  $V_{\mathcal{G}}(x_1), \dots, V_{\mathcal{G}}(x_n)$ , respectively. Two nodes  $V$  and  $U$  of a dag  $\mathcal{G}$  associated with a substitution  $\theta$  are said to be *equivalent* if the subgraphs rooted at  $V$  and  $U$  are associated with the same term. A dag  $\mathcal{G}$  is called *reduced* if

- every node is reachable from a root extra labeled with a variable  $x$ ,  $x \in \mathcal{X}$ , and
- all nodes of  $\mathcal{G}(\theta)$  are pairwise non-equivalent.

Every substitution  $\theta$  is associated with the unique reduced dag, which is denoted by  $\mathcal{G}(\theta)$ . The *size*  $N(\theta)$  of a substitution  $\theta$  is the number of nodes in the reduced dag  $\mathcal{G}(\theta)$  associated with  $\theta$ .

A sequential anti-unification algorithm MST, which computes the most specific templates of substitutions represented by reduced labeled dags, is depicted in Fig. 1. The algorithm gets as input a pair of reduced dags  $\mathcal{G}(\theta')$  and  $\mathcal{G}(\theta'')$  associated with substitutions  $\theta'$  and  $\theta''$  and outputs the reduced dag  $\mathcal{G}$  associated with the most specific template  $\eta$  of  $\theta'$  and  $\theta''$ . Every node  $W$  in  $\mathcal{G}$  matches some pair of nodes (subterms)  $V'$  and  $V''$  in the dags  $\mathcal{G}(\theta')$  and  $\mathcal{G}(\theta'')$ , respectively. Such a node  $W$  is specified by a 3-tuple  $(y, V', V'')$ , where  $y$  is a variable from  $\mathcal{Y}$  used as a unique identifier of  $W$ . A node named  $y$  is denoted by  $mem(y)$ . The algorithm MST operates with two sets  $\mathcal{Q}$  and  $\widehat{\mathcal{Q}}$ . The set  $\mathcal{Q}$  is a worklist of nodes to be handled. When handling a node  $W$  specified by  $(y, V', V'')$  the algorithm assigns the corresponding label to  $W$ , checks all its descendants, and moves  $W$  from the worklist  $\mathcal{Q}$  to the list of processed nodes  $\widehat{\mathcal{Q}}$ .

**Theorem 1.** *Let  $\theta'$  and  $\theta''$  be a pair of substitutions. Then the algorithm *MST* correctly computes a reduced dag associated with the most specific template  $\eta = \theta' \downarrow \theta''$  of  $\theta'$  and  $\theta''$ ; it requires time  $O(n \log n)$ , where  $n = N(\eta)$ .*

Since every node in  $\mathcal{G}$  matches subterms corresponding to exactly one pair of nodes in  $\mathcal{G}(\theta')$  and  $\mathcal{G}(\theta'')$ , the size  $N(\eta)$  of  $\eta$  does not exceed  $N(\theta') \times N(\theta'')$ . Hence, we arrive at the following corollary.

**Corollary 1.** *The most specific template of substitutions  $\theta'$  and  $\theta''$  represented by reduced dags can be computed in time  $O(n^2 \log n)$ , where  $n = \max(N(\theta'), N(\theta''))$ .*

As can be seen from the assertions below, the upper bound can not be substantially improved.

**Theorem 2.** *Suppose that  $\mathcal{F}$  contains a functional symbol of arity  $m > 1$ . Then there exists an infinite sequence of pairs of substitutions  $(\theta'_i, \theta''_i)$ ,  $i \geq 1$ , such that*

$$\frac{1}{6} N(\theta'_i) \times N(\theta''_i) \leq N(\theta'_i \downarrow \theta''_i).$$

**Corollary 2.** *If  $\mathcal{F}$  contains a functional symbol of arity  $m > 1$  then time complexity of the anti-unification problem for two substitutions represented by reduced dags is  $\Omega(n^2)$ .*

Following the ideas suggested in [4, 5, 9] we proved that the anti-unification problem for terms represented by labeled reduced dags is in  $\text{NC}^2$ . We also designed a parallel anti-unification algorithm which computes the most specific template of substitutions  $\theta_1$  and  $\theta_2$  in time  $O(\log^2 n)$  using  $O(n^5)$  processors, where  $n = \max(N(\theta_1), N(\theta_2))$ .

### 3 Generating invariants with the help of anti-unification

In this section we demonstrate how anti-unification algorithms can be applied to the generation of program invariants. We introduce a simple nondeterministic formal model of sequential programs and show that the most specific invariants of the form  $x_1 = t_1 \wedge x_2 = t_2 \wedge \dots \wedge x_n = t_n$  can be computed purely automatically by conventional static analysis techniques (see [11]) adapted to the lattice of finite substitutions.

Let  $\mathcal{V} = \{v_1, \dots, v_n\}$  be a finite set of variables and  $\text{Term}[\mathcal{V}, \mathcal{F}]$  be the set of terms over  $\mathcal{V}$  and a set of functional symbols  $\mathcal{F}$ . Then a *program* is a pair  $\Pi = \langle L, E \rangle$ , where  $L$  is a finite set of program points and  $E$  is a finite set of assignment statements. Every program point is a non-negative integer. We will assume that  $L$  includes 0 which is the *entry point* of  $\Pi$ . Every *assignment statement*  $e$  of a program  $\Pi$  is an expression of the form  $l_{in} : v \leftarrow t : l_{out}$ , where  $v \in \mathcal{V}$ ,  $t \in \text{Term}[\mathcal{V}, \mathcal{F}]$ , and  $l_{in}, l_{out}$  are program points. The integer  $l_{in}$  is called the *entry point* of  $e$  (denoted  $in(e)$ ) and the integer  $l_{out}$  is called the *exit point* of  $e$  (denoted  $out(e)$ ).

Any finite sequence of statements  $tr = e_1, e_2, \dots, e_m$  is called a *trace* of a program  $\Pi$  if  $in(e_1) = 0$  and  $out(e_i) = in(e_{i+1})$  for every  $i$ ,  $1 \leq i < m$ . We say that such a trace  $tr$  leads to the point  $out(e_m)$ . The set of all traces of a program  $\Pi$  leading to a point  $l$  will be denoted by  $Tr_{\Pi}(l)$ .

The semantics of our programs is defined as follows. Let  $M = \{D_M, \bar{f}_1, \dots, \bar{f}_k, =\}$  be a first-order model, where  $D_M$  is a semantic domain with equality relation  $=$ , and functions  $\bar{f}_1, \dots, \bar{f}_k$  are interpretations of the functional symbols from  $\mathcal{F}$ . A *data state*  $\sigma$  is a mapping  $\mathcal{V} \rightarrow D_M$ . Given a term  $t$  and a data state  $\sigma$  we write  $t[\sigma]$  to denote the value of  $t$  in the data state  $\sigma$ .

Let  $M$  be a model,  $\sigma_0$  be a data state, and  $tr = e_1, e_2, \dots, e_m$  be a trace of  $\Pi$  such that  $e_i = l_i : v_i = t_i : l_{i+1}$ . Then the *run* of a program  $\Pi$  for the data state  $\sigma_0$  and the trace  $tr$  is the finite sequence  $(e_1, \sigma_1), (e_2, \sigma_2), \dots, (e_m, \sigma_m)$ , such that the data state  $\sigma_i$  agrees with  $\sigma_{i-1}$  except for the variable  $v_i$ , where the value  $v_i[\sigma_{i-1}]$  is changed to  $t_i[\sigma_{i-1}]$ , for every  $i$ ,  $1 \leq i \leq m$ . The final state  $\sigma_m$  of this run is called the result of the run and denoted by  $r(\sigma_0, tr)$ .

A first-order formula  $\Phi(v_1, \dots, v_n)$  is called an *M-invariant* of a program  $\Pi$  at a point  $l$  iff  $M, r(\sigma_0, tr) \models \Phi(v_1, \dots, v_n)$  holds for every data state  $\sigma_0$  and trace  $tr \in Tr_{\Pi}(l)$ . If  $\Phi(v_1, \dots, v_n)$  is an *M-invariant* for every model  $M$  then it is called a *strong invariant*.

**Theorem 3.** *Let a model  $H$  be an Herbrand model. Then a formula  $\Phi(v_1, \dots, v_n)$  is a strong invariant of  $\Pi$  at a point  $l$  iff  $\Phi(v_1, \dots, v_n)$  is an  $H$ -invariant of  $\Pi$  at the same point.*

An invariant  $\Phi$  is called *the most specific strong invariant* if, for every strong invariant  $\Psi$ , the formula  $\Phi \rightarrow \Psi$  is valid. An invariant  $\Phi(v_1, \dots, v_n)$  of the form  $\exists y_1 \dots \exists y_k (v_1 = t_1 \wedge v_2 = t_2 \wedge \dots \wedge v_n = t_n)$  is called an *equality invariant*.

We are in a position to show how anti-unification can be used in generating the most specific strong equality invariants.

Every statement  $e = l_{in} : v \Leftarrow t : l_{out}$  gives rise to a  $\mathcal{V}$ - $\mathcal{V}$ -substitution  $\theta_e = \{v/t\}$ . The substitutions introduced thus provide a way of characterizing the equality invariants of programs. With every trace  $tr = e_1, e_2, \dots, e_m$  of a program  $\Pi$  we associate a substitution  $\eta_{tr} = \theta_{e_m} \dots \theta_{e_2} \theta_{e_1} \varepsilon$  which is a composition of substitutions associated with the statements  $e_m, \dots, e_2, e_1$  and the empty substitution  $\varepsilon$ . Then, given a program  $\Pi$  and a point  $l$  of  $\Pi$ , we denote by  $\theta_{\Pi, l}$  the substitution  $\downarrow_{tr \in Tr_{\Pi}(l)} \eta_{tr}$  which is the most specific template of all substitutions associated with the traces of  $\Pi$  leading to the point  $l$ .

**Theorem 4.** *Let  $\Pi = \langle L, E \rangle$  be a program and  $l$  be a point of  $\Pi$ . Suppose that  $\theta_{\Pi, l} = \{v_1/t_1, v_2/t_2, \dots, v_n/t_n\}$ . Then the formula*

$$\Phi_{\Pi, l} = \exists y_1 \dots \exists y_k (v_1 = t_1 \wedge v_2 = t_2 \wedge \dots \wedge v_n = t_n),$$

where  $\{y_1, \dots, y_k\}$  is the set of all variables occurred in the terms  $t_1, t_2, \dots, t_n$ , is the most specific strong equality invariant of the program  $\Pi$  at the point  $l$ .

To effectively compute the substitutions  $\theta_{\Pi, l}$  consider the system of equations

$$\Omega(\Pi) : \begin{cases} \Theta_l = \downarrow_{e \in E, out(e)=l} \theta_e \Theta_{in(e)}, & l \in L, l \neq 0, \\ \Theta_0 = \varepsilon, \end{cases}$$

where the  $\Theta_l$ ,  $l \in L$ , are the unknown substitutions.

**Theorem 5.** *For every program  $\Pi = \langle L, E \rangle$ , the set of substitutions  $\{\theta_{\Pi, l} : l \in L\}$  is the least solution to the system  $\Omega(\Pi)$ .*

This theorem leans upon left distributivity of composition of substitutions over anti-unification. To solve the system  $\Omega(\Pi)$  one can involve anti-unification algorithms and any iterative technique used in program static analysis for computing the least fixed points of monotonic operators on lattices (see [11]).

## 4 Duplicate code detection using anti-unification

Detecting duplicate pieces of program code is another task that can be effectively solved with the help of anti-unification. Two sequences of program statements form duplicate code if they are similar enough according to a selected measure of similarity. A set of duplicated pieces of code is called a *clone*. Although there is a huge amount of papers dealing with the duplicate code detection problem (see [18] for survey), so far as we know, no generally recognized definitions of code cloning is developed yet. Pieces of code can be viewed as similar based on syntactic criteria or at the semantic level.

The authors of the paper [7] came up with a proposal for detecting code clones by analyzing the patterns of program expressions. We think that this is one of the most simple yet effective approach to checking the similarity code pieces. Following this line of research we have developed a new anti-unification based duplicate code detection algorithm. Its key idea is as follows. Given two expressions  $E_1$  and  $E_2$ , one need to compute their most specific template  $E = E_1 \downarrow E_2$  and estimate how much  $E_1$  and  $E_2$  differs from  $E$ . The latter can be done based on anti-unification

distance which is defined as follows. Let  $E$  be the most specific template of expressions (terms, statements, sequences of statements, etc.)  $E_1$  and  $E_2$ , such that  $E_1 = E\eta_1$  and  $E_2 = E\eta_2$ . Then *anti-unification distance*  $\rho(E_1, E_2)$  is the sum of sizes  $N(\eta_1)$  and  $N(\eta_2)$  of substitutions  $\eta_1, \eta_2$ . Anti-unification distance  $\rho(E_1, E_2)$  can be seen as a variant of tree edit distance introduced in [2]. If  $\rho(E_1, E_2)$  is minor then  $E_1$  and  $E_2$  belong to the same clone. The efficiency of anti-unification algorithms guarantees that this approach is applicable to large programs independent of the source language.

In order to find clones which consist of sequences of statements (not just a single statement) we developed a compound algorithm that consists of three phases. In the beginning anti-unification is used to partition all statements of a program under analysis into clusters. Such clusterization makes it possible to view the code as a sequence of cluster identifiers. At the second stage the algorithm finds all pairs of identical sequences of cluster identifiers. Finally, the matching pairs of sequences having similar statements in corresponding positions are checked once again for global similarity. This checking also involves the computation of anti-unification distance. Two sequences of program statements are assigned to the same clone if the distance between them is below some certain threshold.

Now we discuss the stages of our duplication detection algorithm in some detail. In the very beginning of the whole algorithm an abstract syntax tree for the analyzed program is built. For the sake of efficiency and memory saving, this tree can be transformed into a reduced dag. Every statement is associated with a subgraph (a tree or a dag) in an abstract syntax graph, and anti-unification algorithm is used to compute a distance between any pair of program statements. Clusterization of program statements is performed in two passes. During the first pass the most frequent templates of program statements in the source code are discovered and a preliminary clusterization is performed. Every cluster  $C$  is characterized by the template  $E_C = \downarrow_{E \in C} E$ . Each new statement  $E'$  is compared with the templates  $E_C$  of all existing clusters. If the distance  $\rho(E', E_C)$  is below some threshold  $d_1$  then the updated template of  $C$  becomes equal to  $E_C \downarrow E'$ . If no such clusters are found then  $E'$  forms a new cluster. During the second pass all statements are processed again. For every statement  $E'$  the algorithm searches the cluster  $C$  from the set produced at the first pass whose template  $E_C$  is the most similar to  $E'$ . When such cluster  $C$  is found,  $E'$  is assigned to  $C$ .

After the first stage of our algorithm all statements are assigned to clusters and marked with corresponding clusters identifiers. At the second stage the algorithm searches for long enough pairs of sequences of statements which are labeled identically, i.e. the statements at the same position in both sequences are marked with the same ID. Detected pairs are considered as clone candidates and their similarity have to be checked at the next stage. This checking is performed at the third stage as follows. Every sequences  $B = E_1, E_2, \dots, E_n$  is treated as a whole expression. If anti-unification distance  $\rho(B', B'')$  between sequences  $B'$  and  $B''$  is below a certain threshold then this pair is reported as a clone.

The algorithm described above has been implemented in a software tool CLONE DIGGER aimed at detecting similar code in Python and Java programs (see [3]).

## 5 Conclusion

The main contribution of our work is twofold.

1. We introduced both sequential and parallel anti-unification algorithms for computing the most specific templates (patterns) of expressions represented as labeled directed acyclic graphs. All previously known anti-unification algorithms operate only with tree-like structures. Since the size of tree representation of some expressions is exponent of the size of their dag representation, our algorithms extend the field of application of anti-unification techniques. We also proved that time complexity of our anti-unification algorithms is close to the optimal one. This provides a firm foundation for the development of various anti-unification based techniques for program analysis.

2. We also showed that anti-unification machinery can be successfully applied to the solution of two important problems in program analysis — generation of program invariants and duplicate code detection. Since anti-unification deals with program expression on syntactic level only, our techniques for invariant generation and clone detection are insensitive to any semantical properties of functions and predicates involved in programs. Thus, program invariants computed with the help of anti-unification capture only primitive relationships between data structures, and even a small modification of a program (say, a transposition of program statements) makes similar pieces of code unrecognizable by our duplicated code detection algorithm. This is the principal drawback of any anti-unification based technique for program analysis. On the other hand, anti-unification algorithms are very efficient and simple, they provide a way for processing large pieces of code in reasonable time. Non-trivial relationships and structures (program invariants and clones) revealed by these means can be used as a raw material for more advanced program analysis procedures. Therefore, anti-unification based techniques for program analysis can find practical use in the front end of many tools for program optimization and verification.

## References

1. Baader F., Snyder W. Unification theory. In J.A. Robinson and A. Voronkov, editors, *Handbook of Automated Reasoning*, 2001, volume 1, p. 447-533.
2. Bille P. A survey on tree distance and related problems. *Theoretical Computer Science*, 2005, v. 337, N 1-3, p. 217-239.
3. Bulychev P. Duplicate code detection using Clone Digger. *PythonMagazine*, 2008. v. 9, p. 18-24.
4. Delcher A.L., Kasif S. Efficient parallel term matching and anti-unification. *Journal of Automated Reasoning*, v. 9, N 3, 1992, p. 391-406.
5. Dwork C., Kanellakis P.C., Stockmeyer L. Parallel algorithms for term matching. *SIAM Journal of Computing*, v. 17, N 4, 1988, p. 711-731.
6. Eder E. Properties of substitutions and unifications. *Journal of Symbolic Computations*, v. 1, 1985, p. 31-46.
7. Evans W., Fraser C., Ma F. Clone detection via structural abstraction. *Proc. of 14th Working Conference on Reverse Engineering*, 2007, p. 150-159.
8. Kovac L.I., Jebelean T. An algorithm for automated generation of invariants for loops with conditionals. *Proc. of the 7-th Int. Symp. on Symbolic and Numeric Algorithms for Scientific Computing*, 2005, p. 245-250.
9. Kuper G.M., McAloon K.W., Palem K.V., Perry K.J. A note on the parallel complexity of anti-unification. *Journal of Automated Reasoning*, v. 9, N 3, 1992, p. 381-389.
10. Martelli A., Montanari U. An efficient unification algorithm. *ACM Transactions on Program, Languages and Systems*, v. 4, N 2, 1982, p. 258-282.
11. Nielson F., Nielson H.R., Hankin C. *Principles of program analysis*. Springer, 1999, 446 p.
12. Oancea C.E., So C., Watt S.M. Generalization in Maple. *Maple Conference*, 2005, p. 277-382.
13. Palamidessi C. Algebraic properties of idempotent substitutions. *Lecture Notes in Computer Science*, v. 443, 1990, p. 386-399.
14. Paterson M.S., Wegman M.N. Linear unification. *The Journal of Computer and System Science*, v. 16, N 2, 1978, p. 158-167.
15. Plotkin G.D. A note on inductive generalization. *Machine Intelligence*, 1970, v. 5, N 1, 1970, p. 153-163.
16. Reynolds J.C. Transformational systems and the algebraic structure of atomic formulas. *Machine Intelligence*, v.5, N 1, 1970, p. 135-151.
17. Robinson J.A. A machine-oriented logic based on the resolution principle. *Journal of the ACM*, v. 12, N 1, 1965, p. 23-41.
18. Roy C.K., Cordy J.R. A survey on software clone detection research. Technical Report N 2007-541, School of Computing Queen's University at Kingston Ontario, Canada.
19. Sorensen M.H., Gluck. R. An algorithm of generalization in positive supercompilation. *Proc. of the 1995 Int. Symposium on Logic Programming*, MIT Press, 1995, p. 465-479.
20. Tiwari A., Rueb H., Saidi H., Shankar N. A technique for invariant generation. *Proc. of the 7-th Int Conf. on Tools and algorithms for the Construction and Analysis of Systems*, 2001, p. 113-127.
21. Watt S.M. Algebraic generalization. *ACM SIGSAM Bulletin*, v. 39, N 3, 2005, p. 93-94.