

# Correctness of Classical Compiler Optimizations using CTL

Carl Christian Frederiksen

July 2, 2001

## Contents

<b>1</b>	<b>Abstract</b>	<b>3</b>
<b>2</b>	<b>Introduction</b>	<b>3</b>
2.1	Overview of the remainder of the paper . . . . .	4
<b>3</b>	<b>Language</b>	<b>4</b>
3.1	Semantics . . . . .	5
<b>4</b>	<b>Temporal logic</b>	<b>7</b>
4.1	CTL . . . . .	8
4.2	CTL-R . . . . .	9
4.3	CTL-R Models . . . . .	11
4.3.1	Fine structure on $AP$ . . . . .	11
4.3.2	Control flow model . . . . .	11
4.3.3	Computation prefix model . . . . .	12
4.3.4	Backwards until . . . . .	13
<b>5</b>	<b>Transformation</b>	<b>13</b>
5.1	Transformation rule syntax . . . . .	14
5.2	Concrete transformation rules . . . . .	15
5.3	Transformation . . . . .	16
<b>6</b>	<b>Correctness proofs</b>	<b>17</b>
6.1	A method for showing program equivalence . . . . .	17
6.2	Copy propagation . . . . .	17
6.3	Elimination of recomputation of available expressions . . . . .	20
<b>7</b>	<b>Skip removal and insertion</b>	<b>22</b>
7.1	Common subexpression elimination . . . . .	25

<b>8</b>	<b>Transformations with multiple program points</b>	<b>26</b>
8.1	Transformation syntax . . . . .	26
8.2	Concrete transformation . . . . .	27
8.3	Transformation . . . . .	27
8.4	Loop invariant hoisting . . . . .	28
<b>9</b>	<b>Summary</b>	<b>33</b>
<b>10</b>	<b>Acknowledgements</b>	<b>33</b>

# 1 Abstract

In this paper, global compiler optimizations are captured by conditional rewrite rules of the form  $\mathcal{I} \Rightarrow \mathcal{I}'$  if  $\phi$ , where  $\mathcal{I}$  and  $\mathcal{I}'$  are program instructions and  $\phi$  is a condition expressed in CTL, a formalism well suited to describe properties involving the control flow of a given program. The goal: to formally prove that if the condition  $\phi$  is satisfied, then the rewrite rule  $\mathcal{I} \Rightarrow \mathcal{I}'$  can be applied to the program without changing the semantics of the program.

This paper is based on the as yet unpublished ideas of David Lacey, Neil Jones and Eric Van Wyk[4]. The aim is to prove correct some conditional rewrite rules corresponding to classical compiler optimizations, beyond those considered by Lacey, et al.

# 2 Introduction

This paper presents a framework for specifying program transformations.

Programs are regarded as being effective specifications of partial functions; A program reads the input, performs some (possibly nonterminating) computation and outputs the result. Two programs are defined to be equivalent if and only if they have the same termination properties and if one of the programs terminate for some input  $v$ , then both programs terminate on input  $v$ , outputting the same value – i.e. the two programs compute the same partial function.

The main focus is on the *correctness* of program transformations, in the sense that the subject program and the transformed program are semantically equivalent. The framework in which the transformations are specified is computable, i.e. given a subject program and a transformation rule it is possible to automatically compute the transformed program resulting from the application of the transformation.

The transformation rules (conditional rewrite rules) described in the paper are transformations that a compiler can utilize in order to obtain a more efficient target program. Other issues such as *optimality*, *efficiency* and *generality* are not considered in detail. Thus the paper will not describe the transformation strategy (i.e. the sequence in which the program transformations are applied).

The approach to applying a program transformation to a subject program:

1. Instantiate an applicable rule schema to a concrete transformation rule.  
For example, from the constant propagation transformation, we might obtain the following rewrite rule.

$$x := y \Rightarrow x := 0$$

assuming a side condition:  $x$  is not assigned until the assignment  $y := 0$  is reached, following all possible backwards paths in the control flow graph of the program.

2. Verify the side condition. The side condition is specified using CTL-R – an extension of CTL – so verification amounts to model checking.

### 3. Construct the transformed program.

In an actual implementation, the set of substitutions  $\Theta$  such that the side condition holds can be found by model checking. Each CTL-R formula is model checked and a set of answer substitutions is found. By unification, a set of substitutions that satisfies all CTL-R formulae in the side condition can be obtained. Goal: All substitutions in the resulting set satisfy the entire side condition, and each one gives rise to a sound transformation, provided that the transformation rule was sound to begin with.

## 2.1 Overview of the remainder of the paper

Section 3 introduces a small imperative language without procedure calls or exceptions and the notion of a *computation prefix* is defined – a computation prefix is the sequence of program states that is computed when executing the program for a given input value. The computation prefixes will be a key factor in proving semantic program equivalence for the transformations. Section 4 introduces an extension of CTL (CTL-R) that will be used for specifying the side conditions in the transformation rules. In Section 5 the transformation rules are introduced. In order to keep the presentation simple, the framework handles only transformations involving one statement being transformed. The application of a transformation rule to a program is formally described as well. Section 6 investigates two transformation rules: *copy propagation* and *elimination of recomputation of available expressions*. Both transformations are proven correct. Next in Section 7 some auxiliary transformations are proven correct outside the framework. In particular the transformations *skip insertion*, *skip removal* and *dead assignment insertion* are proven correct. From these and the transformations from Section 6, the *elimination of common subexpressions* transformation is easily obtained and the correctness follows from the correctness of the other transformations. Section 8 extends the framework to the case with several statements being affected by the transformation. The development is analogous to the single statement case. As an example, the *loop invariant hoisting* transformation is specified and proven correct. Finally Section 9 gives a summary of the work presented.

## 3 Language

The language used in this paper is a small imperative language without procedures. A program has the following form.

$$\pi = \text{read } x; l_1 : \mathcal{I}_{l_1}; \dots; l_m : \mathcal{I}_{l_m}; l_{m+1} : \text{write } y;$$

Each statement, except the `read` statement, is labeled by a unique *label* ranging over the set of labels: *Label* which is not further specified. The set of labels occurring in the program  $\pi$  is denoted  $labels(\pi)$ . The label  $l_1$  at the first statement in the program after the `read` statement is denoted  $entry(\pi)$ , and the label  $l_{m+1}$  at the `write` statement is denoted  $exit(\pi)$ .

**Syntax.** The syntax of a program is given below.

$$\begin{array}{ll}
\pi \in \textit{pgm} & ::= \text{read } x; \textit{stmts}; p : \text{write } y; \\
\textit{stmts} & ::= p : \textit{stmt}; \textit{stmts} \mid \epsilon \\
\mathcal{I} \in \textit{stmt} & ::= x := \textit{exp} \mid \text{if } x \text{ then } p_1 \text{ else } p_2 \mid \text{skip} \\
e \in \textit{exp} & ::= c \mid x \mid \text{op}(\textit{exp}_1, \dots, \textit{exp}_n) \\
x & \in \textit{Variable} \\
p & \in \textit{Label} \\
c & \in \textit{Value}
\end{array}$$

There are three possible statements that can occur between the **read** statement and the **write** statement: assignments, conditionals and **skip** statements. Note that conditionals are not allowed to perform tests directly on expressions, but only on variables. This is not a serious restriction since we can compute the expression, assign it to a variable that does not occur in the program and use that variable in the test instead. All labels in the program are assumed to be unique, and any label referred in a conditional in the program must occur in the program. Note that the **read** statement always uses the variable  $x$  and the **write** statement always uses the variable  $y$ . The statement occurring at label  $p$  is referred to as  $\mathcal{I}_p$ . The set of variables occurring in the program is denoted  $\textit{vars}(\pi)$ . Binary operators are written in infix notation in the example programs, but should be interpreted as having the form given in the syntax definition.

### 3.1 Semantics

In this section the semantics is defined. The domain of values *Value*, the operations on it *Op* and the interpretation of the operators  $\llbracket \cdot \rrbracket_{op} : \textit{Value}^* \rightarrow \textit{Value}$  are left unspecified. However it is assumed that there exists a distinguished element *true* of *Value*.

**Next label function.** Suppose  $\pi \in \textit{pgm}$  is a program and suppose  $\prec$  is a total ordering on *Label* such that  $p \prec p'$  if and only if  $p$  appears syntactically before  $p'$  in the program. The *next-label function*  $Nx_\pi : \textit{Label} \rightarrow \textit{Label}$  is defined by:

$$Nx_\pi(p) = \min_{\prec} \{p' \in \textit{Label} \mid p \prec p'\}.$$

The next label function is a device for getting from a statement to the next label in the program. The next label function  $Nx_\pi$  takes a label occurring in the program  $\pi$  and returns the label occurring immediately syntactically after, i.e. the “smallest” label that is “greater than” the given label with respect to the ordering  $\prec$  on the labels. Since each statement is labeled by precisely one statement, the next label function is well-defined in the domain  $\textit{Label} \setminus \{\textit{exit}(\pi)\}$ .

**Store.** The program store is a finite map that maps variables to their values. The set of stores is denoted  $\textit{Sto} = \textit{Variable} \rightarrow \textit{Value}$ . Suppose  $\sigma : X \rightarrow \textit{Value}$ , where  $\textit{dom}(\sigma) = X \subseteq \textit{Variable}$ :

1.  $\sigma[x \mapsto v]$  is the extension of  $\sigma$  to  $X \cup \{x\}$  such that  $\sigma[x \mapsto v](x) = v$  and  $\forall y \in X \setminus \{x\} : \sigma[x \mapsto v](y) = \sigma(y)$ .
2.  $\sigma|_{X'}$  denotes the restriction of  $\sigma$  to the domain  $X' \subseteq X$ ,  $\sigma|_{X'} : X' \rightarrow \text{Value}$  where  $\forall x \in X' : \sigma(x) = \sigma|_{X'}(x)$ .
3.  $\sigma \setminus x = \sigma|_{X \setminus \{x\}}$  denotes  $\sigma$  restricted to its domain excluding  $x$ .

**Expression evaluation.** Expression evaluation is defined by the function

$$\llbracket \cdot \rrbracket : \text{Expr} \rightarrow \text{Sto} \rightarrow \text{Value}.$$

$$\begin{aligned} \llbracket c \rrbracket \sigma &= c \\ \llbracket x \rrbracket \sigma &= \sigma(x) \\ \llbracket \text{op}(exp_1, \dots, exp_n) \rrbracket \sigma &= \llbracket \text{op} \rrbracket_{op}(\llbracket exp_1 \rrbracket \sigma, \dots, \llbracket exp_n \rrbracket \sigma) \end{aligned}$$

The following lemma states that when evaluating an expression, the result depends only upon the variables occurring the expression, and their bindings.

**Lemma 1.** Suppose  $E \in \text{Exp}$  and  $\sigma \in \text{Sto}$  then  $\llbracket E \rrbracket \sigma = \llbracket E \rrbracket \sigma|_{\text{vars}(E)}$

**Program state.** A *program state* for a given program  $\pi$ , is a pair  $(p, \sigma)$  where  $p$  is a label in  $\text{labels}(\pi)$  and  $\sigma$  is a store such that  $\text{vars}(\pi) \subseteq \text{dom}(\sigma)$ . The set of all program states is denoted  $\text{PgmState}$ .

**Semantic transition relation.** Let the semantic transition relation  $\rightarrow \subseteq \text{PgmState} \times \text{PgmState}$  for a program  $\pi \in \text{pgm}$  be defined by:

1. If  $\mathcal{I}_p = \text{skip}$  then  $(p, \sigma) \rightarrow (N\pi(p), \sigma)$ .
2. If  $\mathcal{I}_p = (x := exp)$  then  $(p, \sigma) \rightarrow (N\pi(p), \sigma[x \mapsto \llbracket exp \rrbracket \sigma])$ .
3. If  $\mathcal{I}_p = (\text{if } x \text{ then } p_1 \text{ else } p_2)$  and  $(\llbracket x \rrbracket \sigma = \text{true})$  then  $(p, \sigma) \rightarrow (p_1, \sigma)$
4. If  $\mathcal{I}_p = (\text{if } x \text{ then } p_1 \text{ else } p_2)$  and  $(\llbracket x \rrbracket \sigma \neq \text{true})$  then  $(p, \sigma) \rightarrow (p_2, \sigma)$
5.  $(\text{exit}(\pi), \sigma) \rightarrow (\text{exit}(\pi), \sigma)$ .

**Corollary 1.** The semantic transition relation is deterministic.

*Proof.* Immediate from the definition since only one of the clauses in the definition can hold for a state for a given program.  $\square$

**Initial state.** The initial state  $in_\pi(v) \in \text{PgmState}$  for a program  $\pi \in \text{pgm}$  and input  $v \in \text{Value}$  is defined by

$$in_\pi(v) = (\text{entry}(\pi), [x \mapsto v, y_1 \mapsto \text{true}, \dots, y_k \mapsto \text{true}])$$

where  $\text{vars}(\pi) \setminus \{x\} = \{y_1, \dots, y_k\}$ . In short the initial state is the pair consisting of the entry label and a store that maps the variable  $x$  to the input of the program, and all other variables to the value *true*.

**Sequence notation.** Let  $S$  be some arbitrary set. The set of sequences of finite length over  $S$  is denoted  $S^*$ . The set of infinite sequences over  $S$  is denoted  $S^\omega$  and the set of sequences of finite or infinite length is denoted  $S^{*\omega}$ .

**Computation prefix.** A *computation prefix* for a program  $\pi \in \text{pgm}$  and a value  $v \in \text{Value}$  is a finite or infinite sequence  $C \in \text{PgmState}^{*\omega}$

$$(p_0, \sigma_0) \rightarrow (p_1, \sigma_1) \rightarrow \dots$$

such that  $(p_0, \sigma_0) = \text{in}_\pi(v)$  and  $(p_i, \sigma_i) \rightarrow (p_{i+1}, \sigma_{i+1})$  for all  $i$ .

The above notation is used when the program and program input is clear from the context. To explicitly specify the program and input, computation prefixes are written:

$$\pi, v \vdash (p_0, \sigma_0) \rightarrow (p_1, \sigma_1) \rightarrow \dots$$

The set of computation prefixes, for a given program  $\pi$  and an input value  $v$ , is denoted  $\mathcal{T}_{\text{pfx}}(\pi, v)$ .

**Next prefix relation.** Suppose  $\pi \in \text{pgm}$  is a program,  $v \in \text{Value}$  is the input and  $C \in \mathcal{T}_{\text{pfx}}(\pi, v)$  is a finite computation prefix such that

$$C = (p_0, \sigma_0) \rightarrow \dots \rightarrow (p_n, \sigma_n).$$

The *next prefix relation* is a relation  $\rightarrow \subseteq \mathcal{T}_{\text{pfx}}(\pi, v) \times \mathcal{T}_{\text{pfx}}(\pi, v)$  where  $C \rightarrow C'$  for a finite computation prefix  $C' \in \mathcal{T}_{\text{pfx}}(\pi, v)$  if and only if  $C'$  has the form

$$C' = (p_0, \sigma_0) \rightarrow \dots \rightarrow (p_n, \sigma_n) \rightarrow (p_{n+1}, \sigma_{n+1}).$$

**Program semantics.** Suppose  $\pi \in \text{pgm}$  is a program. The partial meaning function  $\llbracket \cdot \rrbracket : \text{pgm} \rightarrow \text{Value} \rightarrow \text{Value}$  is defined as  $\llbracket \pi \rrbracket v = \sigma_k(\mathbf{y})$  if there exists a computation prefix  $C = \pi, v \vdash (p_1, \sigma_1) \rightarrow \dots \rightarrow (p_k, \sigma_k)$  such that  $p_k = \text{exit}(\pi)$ .

## 4 Temporal logic

The basic idea is to express a program transformation rule as a number of rewrite rules each with a side condition specified in a variant of CTL. The use of temporal logic allows expressing properties of computation prefixes of the subject program. The side condition corresponds to some analysis that a compiler might perform in order to ensure that the semantics of the subject program is preserved.

$\mathcal{M}, n \models \text{true}$	iff	$\text{true}$
$\mathcal{M}, n \models \text{false}$	iff	$\text{false}$
$\mathcal{M}, n \models \text{ap}$	iff	$\text{ap} \in V(n)$
$\mathcal{M}, n \models \neg\phi$	iff	$\neg\mathcal{M}, n \models \phi$
$\mathcal{M}, n \models \phi_1 \wedge \phi_2$	iff	$\mathcal{M}, n \models \phi_1 \wedge \mathcal{M}, n \models \phi_2$
$\mathcal{M}, n \models A\psi$	iff	$\forall (n = n_0 \rightarrow n_1 \rightarrow \dots) \in S^\omega : \mathcal{M}, (n_i)_{i \geq 0} \models \psi$
$\mathcal{M}, n \models E\psi$	iff	$\exists (n = n_0 \rightarrow n_1 \rightarrow \dots) \in S^\omega : \mathcal{M}, (n_i)_{i \geq 0} \models \psi$
$\mathcal{M}, (n_i)_{i \geq 0} \models X(\phi)$	iff	$n_1 \models \phi$
$\mathcal{M}, (n_i)_{i \geq 0} \models \phi_1 U \phi_2$	iff	$\exists i \geq 0 : \mathcal{M}, n_i \models \phi_2 \wedge \forall 0 \leq j < i : \mathcal{M}, n_j \models \phi_1$

Figure 4.1: Semantics of CTL.

## 4.1 CTL

First the ordinary definition of CTL is briefly presented.

**CTL Model.** A *CTL model* is a triple  $\mathcal{M} = (S, \rightarrow, V)$  where  $S$  is a set of states,  $\rightarrow \subseteq S \times S$  is a *total* relation between elements of  $S$ . The function  $V : S \rightarrow 2^{AP}$  is called the *valuation*, mapping states to atomic propositions that hold in the given state. The atomic propositions are “facts” that can hold for any given state in the model.

A *path* in a CTL model is an *infinite* sequence of states  $(n_i)_{i \geq 0} \in S^\omega$ , such that  $n_i \rightarrow n_{i+1}$  for all  $i \geq 0$ .

**CTL syntax.** The syntax of CTL is given below. A  $\phi$  formula is called a *state formula* since  $\phi$  formulae express properties over a single state in the model. Similarly a  $\psi$  formula is called a *path formula*.

$$\begin{aligned}
\phi &::= \text{true} \mid \text{false} \mid \text{ap} \mid \neg\phi \mid \phi_1 \wedge \phi_2 \\
&\quad \mid A\psi \mid E\psi \\
\psi &::= X\phi \mid \phi_1 U \phi_2 \\
\text{ap} &\in AP
\end{aligned}$$

**CTL satisfaction.** A CTL *clause* is written  $\mathcal{M}, n \models \phi$  where  $\mathcal{M}$  is a CTL model,  $n$  is a state in the model and  $\phi$  is a CTL formula. If the model is clear from the context, the  $\mathcal{M}$  can be omitted from the clause. A CTL clause is *satisfied* if and only if the clause evaluates to true by the semantics defined in Figure 4.1.

It is crucial to note that a given atomic proposition does not hold invariantly true or false for the entire model. The “truth” of a atomic proposition depends on the state for which it is evaluated. In the definition of CTL there are two path quantifiers forall ( $A$ ) that asserts some property over all paths in the CTL model starting with the current state, and exists ( $E$ ). The path operator next



( $X$ ) simply asserts a state property for the next state in the path. The until ( $U$ ) path operator, asserts some property  $\phi_1$  that must hold for all states in the path until eventually a state is reached for which the  $\phi_2$  property holds. The until operator does not assert anything for the remainder of the path.

It is possible to define a number of derived operators that can be used to write the CTL formulae more succinctly. Since the paper uses an extension of CTL, the derived operators will be introduced after the extension, as it subsumes ordinary CTL.

## 4.2 CTL-R

We present an extension to CTL called CTL-R, in order to express side conditions for the transformation rules. The extension is three fold.

**First** the ordinary definition of CTL only allows for reasoning about the “future” in the model, i.e. we can only reason about states that can be reached by the  $\rightarrow$  relation. Since we in essential ways will use the side condition to express some program analysis that must hold prior to the application of the transformation (e.g. expression  $x+y$  is available at label  $p$ ), it is vital that we allow reasoning about the “past” in the model as well. This is done by introducing two new versions of exists ( $E$ ) and forall ( $A$ ):  $\overleftarrow{E}$  and  $\overleftarrow{A}$ . For the interpretation of the “past” quantifiers the subformulae must hold over backwards paths, i.e. following the  $\leftarrow$  relation. Since the definition of a CTL model only requires the  $\rightarrow$  relation to total, any well-defined CTL model might not be a well-defined CTL-R model. We handle this by dropping the totality requirement on the  $\rightarrow$  relation, but at the same time we must specify how finite paths are handled.

**CTL-R Model.** A *CTL-R model* is a triple  $\mathcal{M} = (S, \rightarrow, V)$  where  $S$  is a set of states,  $\rightarrow \subseteq S \times S$  is a relation (not necessarily total) between elements of  $S$ . The function  $V : S \rightarrow 2^{AP}$  is called the *valuation*, mapping states to atomic propositions that hold in the given state.

**Second**, the definition of CTL will be extended to handle models with “dead states”, states with no successor following the direction of the path. Since the definition is extended to handle backwards paths, this applies in both directions. The key for defining a path in a such a model is the notion of a *maximal path*. We will not allow reasoning about a finite path that does not end in a dead state.

**Maximal path.** A path  $(n_i)_{i \geq 0}$  in  $S^{*\omega}$  is *maximal* if  $(n_i)_{i \geq 0}$  is infinite, or of the form  $(n_i)_{0 \leq i \leq k}$  and  $n_k$  is *dead* meaning:

$$\begin{aligned} \neg \exists n \in S : n \rightarrow n_k \text{ if } (n_i)_{0 \leq i \leq k} \text{ is a backwards path,} \\ \neg \exists n \in S : n_k \rightarrow n \text{ if } (n_i)_{0 \leq i \leq k} \text{ is a forwards path.} \end{aligned}$$

The set of all maximal paths over  $S$  will be denoted  $S^{\max}$ .

$\mathcal{M}, n \models true$	iff	$true$
$\mathcal{M}, n \models false$	iff	$false$
$\mathcal{M}, n \models ap(x_1, \dots, x_n)$	iff	$ap(x_1, \dots, x_n) \in V(n)$
$\mathcal{M}, n \models \neg\phi$	iff	not $\mathcal{M}, n \models \phi$
$\mathcal{M}, n \models \phi_1 \wedge \phi_2$	iff	$\mathcal{M}, n \models \phi_1$ and $\mathcal{M}, n \models \phi_2$
$\mathcal{M}, n \models A\psi$	iff	$\forall (n = n_0 \rightarrow n_1 \rightarrow \dots) \in S^{\max} : \mathcal{M}, (n_i)_{i \geq 0} \models \psi$
$\mathcal{M}, n \models \overleftarrow{A}\psi$	iff	$\forall (\dots \rightarrow n_1 \rightarrow n_0 = n) \in S^{\max} : \mathcal{M}, (n_i)_{i \geq 0} \models \psi$
$\mathcal{M}, n \models E\psi$	iff	$\exists (n = n_0 \rightarrow n_1 \rightarrow \dots) \in S^{\max} : \mathcal{M}, (n_i)_{i \geq 0} \models \psi$
$\mathcal{M}, n \models \overleftarrow{E}\psi$	iff	$\exists (\dots \rightarrow n_1 \rightarrow n_0 = n) \in S^{\max} : \mathcal{M}, (n_i)_{i \geq 0} \models \psi$
$\mathcal{M}, (n_i)_{i \geq 0} \models X(\phi)$	iff	$n_1 \models \phi$
$\mathcal{M}, (n_i)_{i \geq 0} \models \phi_1 U \phi_2$	iff	$\exists i \geq 0 : \mathcal{M}, n_i \models \phi_2 \wedge \forall 0 \leq j < i : \mathcal{M}, n_j \models \phi_1$

Figure 4.2: Semantics of CTL-R .

**Third**, in order to be able to express general transformation descriptions, a finer structure is introduced on the atomic propositions. In particular instead of considering the elements of  $AP$  as “tokens”, we will regard the set of atomic propositions as a set of terms.

**CTL-R syntax.** The syntax of CTL-R is given below.

$$\begin{aligned}
\phi &::= true \mid false \mid ap(x_1, \dots, x_n) \mid \neg\phi \mid \phi_1 \wedge \phi_2 \\
&\quad \mid A\psi \mid E\psi \mid \overleftarrow{A}\psi \mid \overleftarrow{E}\psi \\
\psi &::= X\phi \mid \phi_1 U \phi_2
\end{aligned}$$

**CTL-R operators.** Similarly to an ordinary CTL clause, a CTL-R clause is written  $\mathcal{M}, n \models \phi$  where  $\mathcal{M}$  is the CTL-R model,  $n$  is a state in the model and  $\phi$  is a CTL-R formula. The clause is *satisfied* if and only if the clause evaluates to *true* under the definition of CTL-R satisfaction given in Figure 4.2. As before, if the model is clear from the context then  $\mathcal{M}$  can be omitted in the clause.

**Derived CTL-R operators.** Let  $\gamma \in \{A, E, \overleftarrow{A}, \overleftarrow{E}\}$ . Then the following operators can be defined.

$$\begin{aligned}
\mathcal{M}, n \models \gamma F(\phi) &\quad \text{iff} \quad \mathcal{M}, n \models \gamma(true U \phi) \\
\mathcal{M}, n \models \gamma G(\phi) &\quad \text{iff} \quad \mathcal{M}, n \models \neg(\gamma^{-1}F(\neg\phi)) \\
\mathcal{M}, n \models \gamma(\phi_1 W \phi_2) &\quad \text{iff} \quad \mathcal{M}, n \models \gamma(\phi_1 U \phi_2) \vee \gamma G(\phi_1) \\
A^{-1} &= E \\
E^{-1} &= A \\
\overleftarrow{A}^{-1} &= \overleftarrow{E} \\
\overleftarrow{E}^{-1} &= \overleftarrow{A}
\end{aligned}$$

Similarly we can derive the ordinary boolean operators  $\vee, \Rightarrow$  etc.

### 4.3 CTL-R Models

In this section the models needed for the correctness proofs are presented.

#### 4.3.1 Fine structure on $AP$

Since the CTL-R side condition is obtained by instantiating the free variables in the transformation, we define a finer structure on the atomic propositions. This admits reasoning about the instantiated parts of the program.

**Structure on  $AP$ .** The set of atomic propositions used in the models in this paper collects syntactic information about the subject program  $\pi$ . The set is easily computed in linear time by traversing the subject program.

$$\begin{aligned}
 AP \quad ::= \quad & node(Label) \\
 & | \quad entry \\
 & | \quad stmt(stmt) \\
 & | \quad def(Variable) \\
 & | \quad use(Variable) \\
 & | \quad inexp(Variable, exp) \\
 & | \quad isvar(Variable)
 \end{aligned}$$

Each label, statement, variable and expression occurring in an atomic proposition must occur in the subject program. When the control flow model is constructed, the valuation is generated by traversing the program syntax and inserting those atomic propositions that hold at a particular label  $p$  in the valuation  $V(p)$ . As an example consider the following program.

```

      read x;
1 :   y := 5;
2 :   x := x + y;
3 :   write y;

```

The valuation for this program would be the following.

$$\begin{aligned}
 V(1) &\supseteq \{ node(1), \quad entry, \quad stmt(y := 5), \quad def(y) \} \\
 V(2) &\supseteq \{ node(2), \quad stmt(x := x + y), \quad def(x), \quad use(x), use(y) \} \\
 V(3) &\supseteq \{ node(3), \quad stmt(write y), \quad use(y) \}
 \end{aligned}$$

The remaining atomic propositions are independent of the label  $p$ :

$$V(p) \supseteq \{ isvar(x), isvar(y), inexp(x, x + y), inexp(y, x + y) \}.$$

#### 4.3.2 Control flow model

The model in which the satisfaction of the side condition will be decided, is a model of the control flow of the subject program that includes syntactic information about the program. The model is sufficiently detailed to express many of the classical compiler optimizations.

**Control flow model.** The *control flow model*  $\mathcal{M}_{cf}(\pi) = (S, \rightarrow_{cf}, V)$  defines a model based on the program  $\pi \in pgm$ :

1. The set of states  $S$  is defined by  $S = labels(\pi)$ .
2. The state transition relation,  $\rightarrow_{cf}: S \times S$ , is defined as  $p \rightarrow_{cf} p'$  if and only if

$$\begin{aligned} \mathcal{I}_p = (\text{if } x \text{ then } p_1 \text{ else } p_2) & \wedge p' \in \{p_1, p_2\} & \vee \\ (\mathcal{I}_p = (\text{skip}) \vee \mathcal{I}_p = (x:=e)) & \wedge p' = Nx_\pi(p) & \vee \\ \mathcal{I}_p = (\text{write } y) & \wedge p' = p. \end{aligned}$$

3. The *valuation*  $V : S \rightarrow 2^{AP}$  is defined as:

$$\begin{aligned} node(q) \in V(p) & \Leftrightarrow p = q \\ entry \in V(p) & \Leftrightarrow p = entry(\pi) \\ stmt(\mathcal{I}) \in V(p) & \Leftrightarrow \mathcal{I}_p = \mathcal{I} \\ def(x) \in V(p) & \Leftrightarrow \exists e \in expr(\pi) : \mathcal{I}_p = (x:=e) \\ use(y) \in V(p) & \Leftrightarrow \exists x \in vars(\pi), e \in expr(\pi), \exists p_1, p_2 \in labels(\pi) : \\ & \mathcal{I}_p = (x:=e) \wedge y \in vars(e) \vee \\ & \mathcal{I}_p = (\text{if } y \text{ then } p_1 \text{ else } p_2) \vee \\ & \mathcal{I}_p = (\text{write } y) \\ inexp(x, e) \in V(p) & \Leftrightarrow x \in vars(\pi) \wedge e \in expr(\pi) \wedge x \in vars(e) \\ isvar(x) \in V(p) & \Leftrightarrow x \in vars(\pi) \end{aligned}$$

The control flow model can be extracted from the subject program by a simple pass over the program. The last two types of atomic propositions express some global program properties that are independent of the label. The global properties do not belong in a CTL-R formula, but in order to avoid formalizing the machinery to deal with them, they have been included in the set of atomic propositions.

**Lemma 2.** Suppose  $(p_0, \sigma_0) \rightarrow \dots \rightarrow (p_k, \sigma_k)$  and  $\forall 0 \leq i < k : \forall x \in X : \mathcal{M}_{cf}(\pi), p_i \models \neg def(x)$  then  $\sigma_0|_X = \sigma_k|_X$ .

**Lemma 3.** Suppose  $\pi \in pgm$ ,  $\sigma_1 \in Sto$ ,  $(p_1, \sigma_1) \rightarrow (p_2, \sigma_2)$  and  $\mathcal{M}_{cf}(\pi), p_1 \models def(x)$  then  $\sigma_1 \setminus x = \sigma_2 \setminus x$ .

*Proof.* The statement at  $p_1$  must have the form  $x:=e$  for some  $e \in exp$ . By the semantics  $\sigma_2 = \sigma_1[x \mapsto v]$  where  $v = \llbracket e \rrbracket \sigma_1$ . Clearly this implies  $\sigma_1 \setminus x = \sigma_2 \setminus x$ .  $\square$

**Lemma 4.** Suppose  $(p_1, \sigma_1) \rightarrow (p_2, \sigma_2)$ ,  $(p_1, \sigma'_1) \rightarrow (p_2, \sigma'_2)$ ,  $\sigma_1 \setminus x = \sigma'_1 \setminus x$  and  $p_1 \models \neg use(x)$  then  $\sigma_2 \setminus x = \sigma'_2 \setminus x$ .

### 4.3.3 Computation prefix model

In order to prove correctness we define a model that directly corresponds to a given computation prefix of the subject program. The model relation for the computation prefix model is defined as the semantic transition relation. Thus no computation prefix model contains any branching. Intuitively a computation prefix model is simply a computation prefix interpreted as a CTL-R model.

**Computation prefix model.** Suppose  $\pi \in pgm$  and  $C \in \mathcal{T}_{pfx}(\pi, v)$  such that  $C = \pi, v \vdash (p_0, \sigma_1) \rightarrow \dots \rightarrow (p_t, \sigma_t)$  is a finite computation prefix. The computation prefix model for the prefix  $C$  is defined as

$$\mathcal{M}_{pfx}(C) = (\{(p_0, \sigma_1), \dots, (p_t, \sigma_t)\}, \rightarrow, V)$$

where the valuation  $V$  is identical to the valuation for the control flow model.

Next we establish a lemma stating that any atomic proposition that holds for the control flow model of a given program will also hold for any computation trace.

**Lemma 5.** *Suppose  $\pi \in pgm$  and  $C \in \mathcal{T}_{pfx}(\pi, v)$  is a finite computation prefix for some value  $v \in Value$ ,*

$$C = \pi, v \vdash (p_0, \sigma_0) \rightarrow \dots \rightarrow (p_t, \sigma_t) \text{ and } p_t = exit(\pi).$$

*For all  $0 \leq i \leq t$ , the CTL-R clause  $\mathcal{M}_{cf}(\pi), p_i \models ap(x_1, \dots, x_n)$  is satisfied if and only if  $\mathcal{M}_{pfx}(C), (p_i, \sigma_i) \models ap(x_1, \dots, x_n)$  is satisfied.*

#### 4.3.4 Backwards until

In the control flow model, the until operator asserts a rather strong statement; The CTL-R formula  $\overleftarrow{A}(\phi_1 U \phi_2)$  states that, for all backwards paths  $\phi_1$  should hold until  $\phi_2$  holds. If  $\phi_1$  occurs in a loop then there exists a backwards path that loops infinitely. Thus the state where  $\phi_2$  holds must also occur within the loop. For the sake of handling more powerful transformations it would be desirable to be able to express that  $\phi_1$  holds until  $\phi_2$  holds, while allowing the path to “go out of loops”. Thus a more relaxed version is presented:  $\overleftarrow{A}(\phi_1 \wedge \neg entry W \phi_2)$ . Infinite backward paths are admitted since they are unrealizable for any actual computation prefix. At the same time, any finite backwards path must eventually satisfy  $\phi_2$  since any backwards path in the computation prefix model ends in a state where *entry* holds.

**Lemma 6.** *Suppose  $\pi \in pgm$  is a program and  $C \in \mathcal{T}_{pfx}(\pi)$  is a computation prefix of  $\pi$  then  $\mathcal{M}_{cf}, p_t \models \overleftarrow{A}(\phi_1 \wedge \neg entry W \phi_2)$  implies*

$$\begin{aligned} &\forall \text{finite paths} : (\rightarrow n_1 \rightarrow n_0 = p_t) \in Label^* \\ &\exists i \geq 0 : n_i \models \phi_2 \wedge \forall 0 \leq j < i : n_j \models \phi_1. \end{aligned}$$

*That is:  $\overleftarrow{A}(\phi_1 U \phi_2)$  is satisfied for all finite paths in the control flow model.*

## 5 Transformation

This section defines framework for describing a schema of concrete program transformations. An important point is that the applicability of a transformation rule to a program at a given label is computable. It is thus possible to write a compiler that admits transformation rules that the compiler can use to transform (optimize) compiled programs.

## 5.1 Transformation rule syntax

Free variables are allowed in the transformation rules. This allows for specification of transformation templates that intuitively can be instantiated (by binding the free variables) to give concrete a transformation rule.

**Transformation rules.** A *transformation rule* has the syntax:

$$\begin{aligned}
rw &::= vstmt_1 \Rightarrow vstmt_2 \text{ if } \phi \\
vstmt &::= v := vexp \mid \text{if } v \text{ then } \mathbf{v}_1 \text{ else } \mathbf{v}_2 \mid \text{skip} \\
vexp &::= c \mid x \mid \mathbf{v} \mid \text{op}(vexp_1, \dots, vexp_n) \\
\phi &::= true \mid false \mid ap(\mathbf{v}_1, \dots, \mathbf{v}_n) \mid \neg\phi \mid \phi_1 \wedge \phi_2 \\
&\quad \mid \forall y \in Y(\mathbf{v}_1, \dots, \mathbf{v}_n) : \phi \mid \exists y \in Y(\mathbf{v}_1, \dots, \mathbf{v}_n) : \phi \\
&\quad \mid A\psi \mid E\psi \mid \overleftarrow{A}\psi \mid \overleftarrow{E}\psi \\
\psi &::= X\phi \mid \phi_1 U \phi_2 \\
Y &::= vars \\
\mathbf{v} &\in FormulaVariable \\
c &\in Value \\
x &\in Variable \\
y &\in QuantifierVariable
\end{aligned}$$

The rewrite rule “*rw*” may contain free variables ranging over substructures of statements (to be bound in the instantiation phase). The side conditions are expressed in CTL-R, extended with two constructs forall ( $\forall$ ) and exists ( $\exists$ ). The purpose of the two constructs is two allow reasoning about substructures of statements. The constructs will be removed by “unrolling” before any model checking is performed, thus the formula to be model checked is an ordinary CTL-R formula.

The functions  $Y$  are names for functions that, given an expression in the program generates a set of expressions that can be quantified over. An example could be:  $\forall y \in vars(E) : use(y)$ . Once  $E$  is instantiated to an expression, say  $a + b$ , the set  $vars(E) = \{a, b\}$  can be computed and the construct can be “unrolled” to  $use(a) \wedge use(b)$ . The functionality admits expressing transformations such as the elimination of recomputation of available expressions transformation, as will be demonstrated subsequently.

Applying a transformation rule to a given program is done in three steps.

1. All free variables in the transformation rule are *instantiated*, by matching a statement in the program with the left-hand side of the rewrite rule in the transformation rule. All remaining free variable are instantiated. Next the program dependent subformulae in the side condition are unrolled to give an ordinary CTL-R formula.
2. The side condition is *verified* by model checking the side condition at the label where the transformation rule was instantiated. If the clause holds, the instantiated transformation can be applied to the program, provided that the transformation rule is sound.

3. The new program is *constructed*, simply by replacing the statement in the program with the (instantiated) right-hand side of the rewrite rule.

## 5.2 Concrete transformation rules

Given a program and a transformation rule, we can instantiate the transformation rule at a given program point to obtain a *concrete* transformation rule. Intuitively a concrete transformation rule is obtained by instantiating the free variables in the rule and unrolling the program dependent subformulae in the side condition.

**Closed transformation rule.** Suppose  $\pi \in \text{pgm}$  is a program,  $p \in \text{labels}(\pi)$  is a label in  $\pi$ ,  $T$  is a transformation rule  $lhs \Rightarrow rhs$  if  $\phi$  and  $\theta$  is a substitution binding all variables not occurring in  $lhs$ . The purpose of  $\theta$  is to define unreferenced variables, that are referenced in the side condition, but not in left hand side of the rewrite rule. Define  $\theta'$  as the substitution obtained by unifying the left hand side of the transformation rule with the statement at label  $p$  the  $\pi$ , i.e.  $\theta' = \text{Unify}(lhs, \mathcal{I}_p)$ . Provided that there exists a  $\theta'$  unifying  $lhs$  and  $\mathcal{I}_p$ , the transformation rule  $T'$  is obtained from  $T$  by closure with respect to  $\theta'' = \theta \circ \theta'$  if  $T' = \theta''(T)$  and  $T'$  contains no free variables. The transformation rule  $T'$  is called a *closed transformation rule*, since it contains no free variables. The type of  $\theta''$  is given by  $\theta'' : \text{FormulaVariables} \rightarrow \text{Label} \cup \text{stmt} \cup \text{exp} \cup \text{Variable}$ .

From a closed transformation rule, the *concrete* transformation rules can be obtained by “unrolling” the program quantifiers in the side condition.

**Concrete transformation rule.** If  $lhs \Rightarrow rhs$  if  $\phi$  is a closed transformation rule, then the *concrete transformation rule* is defined as  $lhs \Rightarrow rhs$  if  $\text{unroll}(\phi)$  where  $\text{unroll}$  is a function that recursively traverses a CTL-R formula and expands the forall and exists constructs. When a forall or exists formula is reached, the set that the formula quantifies over is computed. The resulting is set required to be finite for any input. The function  $\llbracket \cdot \rrbracket$  maps the names of set extraction functions to their definitions, the syntactic elements  $Y$  are mapped via. function the  $\llbracket \cdot \rrbracket : Y \rightarrow \text{exp} \rightarrow \mathcal{P}(\text{exp})$  to set extraction functions.

As an example consider  $\forall y \in \text{vars}(a + b) : \text{use}(y)$ . Syntactic element  $\text{vars}$  is associated via. the  $\llbracket \cdot \rrbracket$  with a function that extracts the variables in an expression, resulting in the set  $\{a, b\}$ . The function  $\llbracket \text{vars} \rrbracket : \text{exp} \rightarrow \mathcal{P}(\text{exp})$  recursively traverses an expression and extracts all variables occurring the the expression. The formula would then be unrolled to the CTL-R formula  $\text{use}(a) \wedge \text{use}(b)$ .

$$\begin{array}{ll}
\text{unroll}(\text{true}) & = \text{true} \\
\text{unroll}(\text{false}) & = \text{false} \\
\text{unroll}(\text{ap}(x_1, \dots, x_n)) & = \text{ap}(x_1, \dots, x_n) \\
\text{unroll}(\neg\phi) & = \neg\text{unroll}(\phi) \\
\text{unroll}(\phi_1 \wedge \phi_2) & = \text{unroll}(\phi_1) \wedge \text{unroll}(\phi_2) \\
\text{unroll}(\forall v \in Y(x_1, \dots, x_n) : \phi) & = \bigwedge_{y \in \llbracket Y \rrbracket} \text{unroll}([v \mapsto y]\phi) \\
\text{unroll}(\exists v \in Y(x_1, \dots, x_n) : \phi) & = \bigvee_{y \in \llbracket Y \rrbracket} \text{unroll}([v \mapsto y]\phi) \\
\text{unroll}(\gamma X(\phi)) & = \gamma X(\text{unroll}(\phi)) \\
\text{unroll}(\gamma(\phi_1 \text{ } U \text{ } \phi_2)) & = \gamma(\text{unroll}(\phi_1) \text{ } U \text{ } \text{unroll}(\phi_2))
\end{array}$$

Suppose given a transformation rule  $T$ , program  $\pi \in \text{pgm}$  and a label  $p$ . The concrete transformation rule  $T'$  is obtained by *instantiating* the transformation rule  $T$  by the substitution  $\theta$  at label  $p$  in program  $\pi$  if  $T'$  is obtained by first computing the closed transformation rule  $T''$  from  $T$  with respect to  $\theta$  and then computing  $T'$  by unrolling.

### 5.3 Transformation

We now define what it means to apply a rewrite rule to a given program at a specified program point.

**Apply.** Suppose  $\pi, \pi' \in \text{pgm}$  are programs,  $p$  is a label in  $\pi$ ,  $T$  is a transformation rule and  $\theta$  is a substitution. The relation  $\text{Apply}(\pi, p, T, \theta, \pi')$  holds if and only of all of the following hold.

1. (Instantiation): Given  $\theta$  there exists an instantiation  $T'$  by the substitution  $\theta'$  of  $T$  at label  $p$  on the form

$$\mathcal{I} \Rightarrow \mathcal{I}' \quad \text{if} \quad \phi.$$

2. (Verification): The CTL-R clause  $\mathcal{M}_{cf}(\pi), p \models \phi$  is satisfied.
3. (Construction): The transformed program  $\pi' \in \text{pgm}$  has the form

$$\pi' = \text{read } x; \ l_1 : \mathcal{I}'_{l_1}; \dots; \ l_m : \mathcal{I}'_{l_m}; \ l_{m+1} : \text{write } y,$$

where the statements of  $\pi'$  are defined by

$$\mathcal{I}'_{l_i} = \begin{cases} \theta'(\mathcal{I}') & \text{if } l_i = p \\ \mathcal{I}_{l_i} & \text{otherwise} \end{cases}$$

The program  $\pi'$  is obtained by *applying* the transformation  $T$  to to program  $\pi$  at label  $p$  given the substitution  $\theta$  if  $\text{Apply}(\pi, p, T, \theta, \pi')$  holds.

By the definition of *Apply*, the labels in a transformed program appear in the same order as the labels in the subject program.

**Corollary 2.** Suppose  $\pi, \pi' \in \text{pgm}$  and  $\text{Apply}(\pi, p, T, \theta, \pi')$  then  $Nx_\pi = Nx_{\pi'}$ .



## 6 Correctness proofs

**Lemma 7.** *Suppose  $\pi \in \text{pgm}$  and  $\text{Apply}(\pi, p, RW, \theta, \pi')$  for some  $p \in \text{label}(\pi)$  and some transformation rule  $RW$ . If  $\text{in}_\pi(v) = (p_0, \sigma_0)$  and  $\text{in}_{\pi'}(v) = (p'_0, \sigma'_0)$  for some  $v \in \text{Value}$  then  $p_0 = p'_0$  and  $\sigma_0|_{\text{vars}(\pi) \cap \text{vars}(\pi')} = \sigma'_0|_{\text{vars}(\pi) \cap \text{vars}(\pi')}$ .*

### 6.1 A method for showing program equivalence

**Corollary 3.** *Suppose  $\pi \in \text{pgm}$  and  $\text{Apply}(\pi, p, T, \theta, \pi')$  holds for some label  $p$  and some transformation  $T$ . Let  $C \in \mathcal{T}_{\text{pfx}}(\pi)$  and  $C' \in \mathcal{T}_{\text{pfx}}(\pi')$  be computation prefixes of the same length of  $\pi$  and  $\pi'$ :*

$$\begin{aligned} C &= \pi, v \vdash (p_0, \sigma_0) \rightarrow \dots \rightarrow (p_n, \sigma_n) \\ C' &= \pi', v \vdash (p'_0, \sigma'_0) \rightarrow \dots \rightarrow (p'_n, \sigma'_n), \end{aligned}$$

*such that  $p_i = p'_i$  for all  $0 \leq i \leq n$ .*

*If there exists a relation  $\mathcal{R}$  between computation prefixes such that:  $CR C'$  implies  $p_i = p'_i$  for all  $0 \leq i \leq n$  and  $p_n = \text{exit}(\pi)$  implies that  $\sigma_n(y) = \sigma'_n(y)$  then  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .*

*Proof.* Suppose  $\pi$  does not terminate on input  $v$ , then  $p_i \neq \text{exit}(\pi)$  for all  $i > 0$ , so  $p'_i \neq \text{exit}(\pi) = \text{exit}(\pi')$  by Lemma 2 implying that  $\pi'$  does not terminate on input  $v$ . Conversely nontermination of  $\pi'$  implies nontermination of  $\pi$ . Otherwise  $\pi$  terminates on input  $v$ . Let  $n$  be given such that  $p_n = \text{exit}(\pi) = \text{exit}(\pi')$ . By assumption  $\sigma_n(y) = \sigma'_n(y)$  so  $\llbracket \pi \rrbracket(v) = \llbracket \pi' \rrbracket(v)$ . Since  $v$  was an arbitrary value  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .  $\square$

### 6.2 Copy propagation

The copy propagation transformation propagates assignments of the form  $x := y$  to other assignments of the same type. The purpose is to render intermediate variables in chains of copying assignments dead.

**Example.** As an example of copy propagation, consider the following program fragment.

```
1 : y := z;
2 : a := 5 + z;
3 : b := a * 2;
4 : x := y;
```

Performing copy propagation on the statement at label 4 results in the fragment:

```
1 : y := z;
2 : a := 5 + z;
3 : b := a * 2;
4 : x := z;
```

Since the program variables  $y$  and  $z$  are not defined (assigned) between 2 and 4, the two variables evaluate to the same value at label 4. Thus either variable

can be used in the assignment. If the variable  $y$  is not used in subsequent computations then it can be removed by a transformation eliminating definitions of dead variables.

**Copy propagation.** The following transformation rule defines the copy propagation transformation (CP).

$$x:=y \Rightarrow x:=z \quad \text{if} \quad \begin{array}{l} \text{isvar}(z) \wedge \\ \overleftarrow{A}(\neg \text{def}(z) \wedge \neg \text{def}(y) \wedge \neg \text{entry } W \text{ stmt}(y:=z)) \end{array}$$

Informally the side condition states:

1.  $z$  must be bound to a program variable <sup>1</sup>.
2. Along all backwards paths in the control flow graph starting from statement to be transformed, the program variables that  $y$  and  $z$  are bound to must not be assigned to until the assignment  $y:=z$  is reached (which *must* eventually happen). Note that  $y$  and  $z$  are implicitly forced to be bound to program variables, since they appear on the left hand side of an assignment (where only program variables can appear).

**Theorem 1.** Suppose given a program  $\pi, \pi' \in \text{pgm}$ , a label  $p \in \text{labels}(\pi)$  and suppose  $\text{Apply}(\pi, p, \text{CP}, \theta, \pi')$  then  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .

*Proof.* Let  $\pi, \pi' \in \text{pgm}$ ,  $\text{Apply}(\pi, p, \text{CP}, \theta, \pi')$ . variables. Suppose the instantiated transformation rule is given by

$$x:=y \Rightarrow x:=z \quad \text{if} \quad \begin{array}{l} \text{isvar}(z) \wedge \\ \overleftarrow{A}(\neg \text{def}(z) \wedge \neg \text{def}(y) \wedge \neg \text{entry } W \text{ stmt}(y:=z)). \end{array}$$

Note that the variables  $x, y, z$  are bound to actual program variables, since  $x$  and  $y$  appears on the left hand side of an assignment, and the satisfaction of the side condition implies that  $z$  is variable.

Let  $v \in \text{Value}$  be the input to the programs and let  $C \in \mathcal{T}_{\text{pfx}}(\pi, v)$  and  $C' \in \mathcal{T}_{\text{pfx}}(\pi', v)$  be two finite prefixes of the same length of  $\pi$  and  $\pi'$  respectively:

$$\begin{array}{lcl} C & = & \pi, v \vdash (p_0, \sigma_0) \rightarrow \dots \rightarrow (p_t, \sigma_t) \\ C' & = & \pi', v \vdash (p'_0, \sigma'_0) \rightarrow \dots \rightarrow (p'_t, \sigma'_t). \end{array}$$

The idea of the proof is to show that the transformation preserves the set of computation prefixes, i.e. that any state transition sequence occurring in the computation trace of the original program  $\pi$  on input  $v$ , will also occur in the computation trace of the transformed program  $\pi'$  on input  $v$  and vice versa. We can then use Lemma 3 to conclude that the program semantics is preserved. First we prove that  $(p_i, \sigma_i) = (p'_i, \sigma'_i)$  for all  $i \leq t$  by induction over the length of  $C$  and  $C'$ .

---

<sup>1</sup>Otherwise  $z$  could be bound to an expression

**Base case.** Since  $\text{vars}(\pi) = \text{vars}(\pi')$  then by Lemma 7:  $(p_0, \sigma_0) = (p'_0, \sigma'_0)$ .

**Induction step.** Since the semantics is deterministic there exists exactly one  $C_2 \in \mathcal{T}_{pfx}(\pi, v)$  and one  $C'_2 \in \mathcal{T}_{pfx}(\pi', v)$  such that  $C \rightarrow C_2$  and  $C' \rightarrow C'_2$ . Assuming that the two computation prefixes are equal up to index  $t$ :

$$\forall 0 \leq i \leq t : (p_i, \sigma_i) = (p'_i, \sigma'_i)$$

then we need to show equivalence at index  $t + 1$ :  $(p_{t+1}, \sigma_{t+1}) = (p'_{t+1}, \sigma'_{t+1})$ .

If  $p_t \neq p$  then  $\mathcal{I}_{p_t} = \mathcal{I}'_{p_t}$ . By the induction assumption  $\sigma_t = \sigma'_t$ . Lemma 2 states that the next label functions for  $\pi$  and  $\pi'$  are identical, so  $(p_{t+1}, \sigma_{t+1}) = (p'_{t+1}, \sigma'_{t+1})$  since the transition relation is deterministic (Corollary 1).

Otherwise  $p_t = p$ , implying that  $\mathcal{I}_{p_t} = (x := y)$  and  $\mathcal{I}'_{p_t} = (x := z)$ . Clearly the next label for either computation prefix is given by the next label function, since the statement at  $p_t$  is an assignment in both programs. Lemma 2 states that the next label function is identical for the two programs  $\pi$  and  $\pi'$ , so  $p_{t+1} = p'_{t+1}$ .

Both statements at  $p_t$  defines (assigns) variable  $x$ , so  $p_t \models \text{def}(x)$  is satisfied for both  $\mathcal{M}_{cf}(\pi)$  and  $\mathcal{M}_{cf}(\pi')$ . By Lemma 3  $\sigma_t \setminus x = \sigma_{t+1} \setminus x$  and  $\sigma'_t \setminus x = \sigma'_{t+1} \setminus x$ . The induction assumption implies  $\sigma_t = \sigma'_t$  so by transitivity  $\sigma_{t+1} \setminus x = \sigma'_{t+1} \setminus x$ . Thus the stores  $\sigma_{t+1}$  and  $\sigma'_{t+1}$  are equivalent if  $\sigma_{t+1}(x) = \sigma'_{t+1}(x)$ .

The side condition must by assumption hold at  $p_t$  in the model for the control flow of program  $\pi$ :

$$\mathcal{M}_{cf}(\pi), p_t \models \overleftarrow{A}(\neg \text{def}(z) \wedge \neg \text{def}(y) \wedge \neg \text{entry } W \text{ stmt}(y := z)).$$

Using Lemma 6 a strong until ( $U$ ) version of the above statement is satisfied in the model  $\mathcal{M}_{cf}(\pi)$  for all finite paths:

$$\begin{aligned} & \forall \text{finite paths} : (\rightarrow n_1 \rightarrow n_0 = p_t) \in \text{Label}^* \\ & \exists i \geq 0 : n_i \models \text{stmt}(y := z) \wedge \forall 0 \leq j < i : n_j \models (\neg \text{def}(z) \wedge \neg \text{def}(y)). \end{aligned}$$

Substituting the definition of the valuation for  $\text{stmt}$  valuations, we get:

$$\begin{aligned} & \forall \text{finite paths} : (\rightarrow n_1 \rightarrow n_0 = p :) \in \text{Label}^* \\ & \exists i \geq 0 : \mathcal{I}_{n_i} = (y := z) \wedge \forall 0 \leq j < i : n_j \models \neg \text{def}(z) \wedge n_j \models \neg \text{def}(y). \end{aligned}$$

Due to Lemma 5 any computation prefix must be described by the control flow model, thus the following holds for the computation prefix

$$\begin{aligned} \exists i \leq t : & \mathcal{I}_{p_i} = (y := z) \wedge \\ & \forall t \geq j > i : \mathcal{M}_{pfx}(C), p_j \models \neg \text{def}(z) \wedge \mathcal{M}_{pfx}(C), p_j \models \neg \text{def}(y). \end{aligned}$$

Let  $i$  be given such that the above holds. The program variables  $z$  and  $y$  are not defined at statements  $\mathcal{I}_{p_t}, \dots, \mathcal{I}_{p_{i+1}}$ , so Lemma 2 implies that  $y$  and  $z$  does not change their value from  $\sigma_{i+1}$  to  $\sigma_t$ :  $\sigma_{i+1}(y) = \sigma_t(y)$  and  $\sigma_{i+1}(z) = \sigma_t(z)$ .

Furthermore  $\mathcal{I}_{p_i} = (y := z)$  so by the semantics:  $\sigma_{i+1}(y) = \sigma_i(z)$ . Clearly  $z$  is not defined at  $\mathcal{I}_{p_i}$ :  $p_i \models \neg \text{def}(z)$ , so by Lemma 2:  $\sigma_i(z) = \sigma_{i+1}(z)$ . Thus by transitivity  $\sigma_t(y) = \sigma_t(z)$ .

$$\begin{aligned}
\sigma_{t+1}(x) &= \sigma_t(y) && \text{(semantics of } \mathcal{I}_{p_t} = (x:=y)) \\
&= \sigma_t(z) && \text{(argument above)} \\
&= \sigma'_t(z) && \text{(induction assumption)} \\
&= \sigma'_{t+1}(x) && \text{(semantics of } \mathcal{I}_{p_{t+1}} = (x:=z))
\end{aligned}$$

This proves that  $\sigma_{t+1} = \sigma'_{t+1}$ , implying that  $(p_{t+1}, \sigma_{t+1}) = (p'_{t+1}, \sigma'_{t+1})$ .

Per induction for any finite computation prefixes  $C \in \mathcal{T}_{pfx}(\pi, v)$  and  $C' \in \mathcal{T}_{pfx}(\pi', v)$  of the same length  $t$ , the two computation prefixes are identical.

It then follows by Lemma 3 that  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .  $\square$

### 6.3 Elimination of recomputation of available expressions

The elimination of recomputation of available expressions transformation detects assignments of expressions that are already available in some variable, provided that the changes in the store does not affect the variables which the expression depends on.

**Example.** As an example, consider the following program fragment.

```

1 : x := a + b
2 : if ... then 2 else 3
3 : z := 5
4 : v := a + b

```

Applying the transformation at label 4 results in the following fragment.

```

1 : x := a + b
2 : if ... then 2 else 3
3 : z := 5
4 : v := x

```

The transformation is semantics preserving since the expression  $a+b$  evaluates to the same value at label 1, and 4 because neither variable changes its value from 1 to 4 regardless of the path taken. The expression is *available* in the variable  $x$  at label 4 since  $x$  is not modified either.

**Elimination of recomputation of available expressions.** The following rewrite rule defines the “elimination of recomputation of available expressions” transformation (*ER*).

$$\mathbf{x} := \mathbf{e} \Rightarrow \mathbf{x} := \mathbf{z} \quad \text{if} \quad \neg \text{inexp}(\mathbf{z}, \mathbf{e}) \wedge \overleftarrow{A} X \overleftarrow{A} ( \neg \text{def}(\mathbf{z}) \wedge \forall y \in \text{vars}(\mathbf{e}) : \neg \text{def}(y) \wedge \neg \text{entry } W \text{stmt}(\mathbf{z} := \mathbf{e}) )$$

Informally the side condition states:

1.  $\mathbf{z}$  may not occur in the expression  $\mathbf{e}$ . If it did, the expression  $\mathbf{e}$  would (potentially) evaluate to a different value in the store after the assignment.
2. The second part of the side condition states that for all backwards paths from the current label back to the statement  $\mathbf{z} := \mathbf{e}$  (which must eventually be reached), the value of  $\mathbf{z}$  is unchanged. Furthermore any program variable occurring in  $\mathbf{e}$  does not change its value on the path. This implies that the expression  $\mathbf{e}$  evaluates to the same expression at label where the transformation takes place and at the  $\mathbf{z} := \mathbf{e}$  statement.
3. The side condition implicitly states that the CTL-R variables  $\mathbf{x}$  and  $\mathbf{z}$  are bound to program variables, since they appear on the left hand side of assignments. The CTL-R variable  $\mathbf{e}$  must be bound to some program expression, since it appears on the right hand side of assignments.

**Theorem 2.** Suppose  $\pi, \pi' \in \text{pgm}$ ,  $\text{Apply}(\pi, p, ER, \theta, \pi')$  then  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .

*Proof.* Let  $\pi, \pi' \in \text{pgm}$ ,  $\text{Apply}(\pi, p, ER, \theta, \pi')$ . Let  $v \in \text{Value}$  be the input to the programs  $\pi$  and  $\pi'$ , and let  $C \in \mathcal{T}_{\text{pfx}}(\pi, v)$  and  $C' \in \mathcal{T}_{\text{pfx}}(\pi', v)$  be program prefixes of the same length:

$$\begin{aligned} C &= \pi, v \vdash (p_0, \sigma_0) \rightarrow \dots \rightarrow (p_t, \sigma_t) \\ C' &= \pi', v \vdash (p'_0, \sigma'_0) \rightarrow \dots \rightarrow (p'_t, \sigma'_t). \end{aligned}$$

Suppose the instantiated transformation rule is given by:

$$x := e \Rightarrow x := z \quad \text{if} \quad \neg \text{inexp}(z, e) \wedge \overleftarrow{A} X \overleftarrow{A} (\neg \text{def}(z) \wedge \forall y \in \text{vars}(e) : \neg \text{def}(y) \wedge \neg \text{entry } W \text{ stmt}(z := e))$$

First we prove by induction over the length of  $C, C'$  that the two computation prefixes are identical:  $(p_i, \sigma_i) = (p'_i, \sigma'_i)$  for all  $i \leq t$ .

**Base case.** Since  $\text{vars}(\pi) = \text{vars}(\pi')$  then by lemma 7:  $(p_0, \sigma_0) = (p'_0, \sigma'_0)$ .

**Induction step.** Since the semantics is deterministic there exists exactly one  $C_2 \in \mathcal{T}_{\text{pfx}}(\pi, v)$ ,  $C'_2 \in \mathcal{T}_{\text{pfx}}(\pi', v)$  such that  $C \rightarrow C_2$  and  $C' \rightarrow C'_2$ . Assuming that  $(p_i, \sigma_i) = (p'_i, \sigma'_i)$  for all  $i \leq t$ , we need to show that  $(p_{t+1}, \sigma_{t+1}) = (p'_{t+1}, \sigma'_{t+1})$ .

If  $p_t \neq p$  then  $\mathcal{I}_{p_t} = \mathcal{I}'_{p_t}$ , so the claim follows by Lemma 2 and the induction assumption.

Otherwise  $p_t = p$ , so  $\mathcal{I}_{p_t} = (x := e)$ ,  $\mathcal{I}'_{p_t} = (x := z)$ . By Lemma 4  $\sigma_{t+1} \setminus x = \sigma'_{t+1} \setminus x$ , so the claim  $\sigma_{t+1} = \sigma'_{t+1}$  holds if we can prove  $\sigma_{t+1}(x) = \sigma'_{t+1}(x)$ . The side condition must hold at  $p$ :

$$p \models \overleftarrow{A} X \overleftarrow{A} (\neg \text{def}(z) \wedge \forall y \in \text{vars}(e) : \neg \text{def}(y) \wedge \neg \text{entry } W \text{ stmt}(z := e)).$$

Using Lemma 6 this implies

$$\begin{aligned} &\forall \text{finite paths } : (\dots \rightarrow n_1 \rightarrow n_0 = p) \in \text{Label}^* : \\ &\exists i \geq 1 : n_i \models \text{stmt}(z := e) \wedge \\ &\quad \forall 1 \leq j < i : (n_j \models \neg \text{def}(z)) \wedge \forall y \in \text{vars}(e) : n_j \models \neg \text{def}(y). \end{aligned}$$

Since the control flow model describes the control flow for any possible computation – Lemma 5 – the above applies for computation prefix  $C$ . Applying the definition of the valuation of the atomic propositions, the following must hold for the concrete computation prefix  $C$ :

$$\begin{aligned} \exists i < t : \mathcal{I}_{p_i} &= (z := e) \wedge \\ \forall i < j \leq t : p_j &\models \neg \text{def}(z) \wedge \\ &\forall y \in \text{vars}(e) : p_j \models \neg \text{def}(y). \end{aligned}$$

The rewrite rule does not change  $\text{def}$  atomic propositions, so

$$\forall y \in \text{vars}(e) : p_j \models \neg \text{def}(y)$$

must also hold for the control flow model for  $\pi'$ :  $\mathcal{M}_{cf}(\pi')$  for  $i < j \leq t$ . Similarly  $p_i \models \text{def}(z)$  holds in  $\mathcal{M}_{cf}(\pi')$ . The side condition states that  $z \notin \text{vars}(e)$ , so

$$\forall y \in \text{vars}(e) : p_j \models \neg \text{def}(y)$$

holds for the program point with the assignment  $z := e$ , i.e. the property holds for  $j = i$  as well. By Lemma 2 we conclude  $\sigma'_i|_{\text{vars}(e)} = \sigma'_t|_{\text{vars}(e)}$ , that is: the expression  $e$  evaluates to the same value in the stores  $\sigma'_t$  and  $\sigma'_i$ .

Next we prove that  $\sigma'_{i+1}(z) = \llbracket E \rrbracket \sigma'_i|_{\text{vars}(e)}$ :

Suppose  $p_i \neq p$  then  $\mathcal{I}'_{p_i} = (z := e)$  so by the semantics  $\llbracket e \rrbracket \sigma'_i|_{\text{vars}(e)} = \sigma'_{i+1}(z)$ .

Otherwise  $p_i = p$ . Since the transformation rule was instantiated  $x = z$  and  $\mathcal{I}'_{p_i} = (x := x)$ , so  $\sigma'_i = \sigma'_{i+1}$  by the semantics. By the induction assumption  $\sigma_i = \sigma'_i$  and  $\sigma_{i+1} = \sigma'_{i+1}$ , so by transitivity  $\sigma_i = \sigma_{i+1}$ . The statement at  $p_i$  in  $\pi$  is  $\mathcal{I}_p = (x := e)$ , thus  $\sigma_{i+1}(x) = \llbracket e \rrbracket \sigma_i|_{\text{vars}(e)}$ . Due to equivalence  $\sigma'_{i+1}(z) = \llbracket e \rrbracket \sigma'_i|_{\text{vars}(e)}$  holds.

Putting the pieces together:

$$\begin{aligned} \sigma_{t+1}(x) &= \llbracket e \rrbracket \sigma_t && \text{(semantics of } \mathcal{I}_{p_t} = (x := e)) \\ &= \llbracket e \rrbracket \sigma'_t|_{\text{vars}(e)} && \text{(Lemma 1)} \\ &= \llbracket e \rrbracket \sigma'_t|_{\text{vars}(e)} && \text{(induction assumption } \sigma_t = \sigma'_t) \\ &= \llbracket e \rrbracket \sigma'_i|_{\text{vars}(e)} && \text{(1}^{st} \text{ argument above)} \\ &= \sigma'_{i+1}(z) && \text{(2}^{nd} \text{ argument above)} \\ &= \sigma'_t(z) && \text{(Lemma 2)} \\ &= \sigma'_{t+1}(x) && \text{(semantics of } \mathcal{I}'_{p_t} = (x := z)) \end{aligned}$$

Thus  $\sigma_{t+1}(x) = \sigma'_{t+1}(x)$  so  $\sigma_{t+1} = \sigma'_{t+1}$  which proves the induction step.

Thus any two computation prefixes of the same length  $C \in \mathcal{T}_{pfx}(\pi, v)$  and  $C' \in \mathcal{T}_{pfx}(\pi', v)$  are equivalent. By Lemma 3 the transformation is semantics preserving:  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ . □

## 7 Skip removal and insertion

As Lemma 2 states, the framework is not strong enough to handle transformations that insert or remove statements. For completeness three transformations are proven correct outside the framework.

1. *Statement insertion transformation*: inserting a **skip** statement between two statements in the subject program.
2. *skip removal transformation*: deleting a **skip** statement from the program while changing all references to the **skip** statement to the statement immediately syntactically after.
3. *dead assignment insertion*: rewriting a **skip** statement to an assignment to a variable that does not occur in the subject program.

These transformations and their proofs are not particularly interesting, they simply prove intuitive facts. Armed with these transformations, it is possible to perform the elimination of recomputation of available expressions transformation, for instance, without leaving the transformed program cluttered with useless **skip** statements.

As will be demonstrated subsequently, some transformations “insert” statements in the program and thus (when specified in this framework) will have to rely on the presence of **skip** statements in the “right” places. The **skip** insertion transformation is targeted to handle these issues.

**Statement insertion.** Suppose  $\pi \in \text{pgm}$  has the form

$$\text{read } x; l_1 : \mathcal{I}_{l_1}; \dots; l_m : \mathcal{I}_{l_m}; l_{m+1} : \text{write } y;$$

and  $p \in \text{Label} \setminus \text{labels}(\pi)$  is a unique label. For any label  $p' \in \text{labels}(\pi)$  the *statement insertion transformation* is defined as  $\pi' = \text{SkipIns}(\pi, p', p)$  where  $p' = l_i$  for some  $i$  and  $\pi' \in \text{pgm}$  is a program of the form

$$\text{read } x; \dots l_{i-1} : \mathcal{I}_{l_{i-1}}; p : \text{skip}; l_i : \mathcal{I}_{l_i}; \dots l_{m+1} : \text{write } y;.$$

**Theorem 3.** Suppose  $\pi \in \text{pgm}$ ,  $l_i \in \text{labels}(\pi)$ ,  $p \in \text{Label} \setminus \text{labels}(\pi)$  and  $\pi' = \text{SkipIns}(\pi, l_i, p)$  then  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .

*Proof.* Proof idea: show by induction for any two computation prefixes  $C \in \mathcal{T}_{\text{pfx}}(\pi)$  and  $C' \in \mathcal{T}_{\text{pfx}}(\pi')$ :

$$\begin{aligned} C &= \pi, v \vdash (p_0, \sigma_0) \rightarrow \dots \rightarrow (p_m, \sigma_m) \\ C' &= \pi', v \vdash (p'_0, \sigma'_0) \rightarrow \dots \rightarrow (p'_n, \sigma'_n) \end{aligned}$$

that for all  $C_2 \in \mathcal{T}_{\text{pfx}}(\pi)$  and  $C'_2 \in \mathcal{T}_{\text{pfx}}(\pi')$  such that  $C \rightarrow C_2$ ,  $C' \rightarrow C'_2$  and  $CRC'$  implies  $C_2RC'_2$  where  $CRC'$  if and only if  $p \neq p_m = p'_n$  and  $\sigma_m = \sigma'_n$ .

**Base case.** Suppose  $p'_0 \neq p$  then clearly the initial states are equivalent. Otherwise  $\sigma_0 = \sigma'_0$  and  $p_0 = Nx_{\pi'}(p'_0)$ , so it follows by the semantics of **skip** that  $(p'_1, \sigma'_1) = (p_0, \sigma_0)$ .

**Induction step.** Suppose  $p'_{n+1} \neq p$  implying  $\mathcal{I}_{p_m} = \mathcal{I}'_{p'_n}$ . The label  $p$  does not occur in the program  $\pi$ , so  $p_m \neq p$ . The induction assumption  $(p_m, \sigma_m) = (p'_n, \sigma'_n)$  then implies that  $(p_{m+1}, \sigma_{m+1}) = (p'_{n+1}, \sigma'_{n+1})$ .

Otherwise  $p'_{n+1} = p$ . The semantics of **skip** and the definition of  $\pi'$  implies that  $p'_{n+2} = p_{m+1}$ . By the induction assumption it follows that  $\sigma_{m+1} = \sigma'_{n+1} = \sigma'_{n+2}$ . Thus the induction assumption holds for  $m+1$  and  $n+2$ .

By induction for any computation prefix  $C$  of  $\pi$  there exists a computation prefix  $C'$  of  $\pi'$  such that  $(p_m, \sigma_m) = (p'_n, \sigma'_n)$  and conversely for any computation prefix  $C'$  of  $\pi'$  such that  $p'_n \neq p$  there exists a computation prefix  $C$  of  $\pi$  such that the final states are equivalent.

Clearly  $\pi$  terminates if and only if  $\pi'$  terminates and if they terminate, they compute the same store. Thus the equivalence  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$  follows by the definition of program semantics.  $\square$

**Skip removal.** Suppose  $\pi \in pgm$  has the form

$$\text{read } x; l_1 : \mathcal{I}_{l_1}; \dots; l_m : \mathcal{I}_{l_m}; l_{m+1} : \text{write } y; .$$

The *skip removal transformation* is defined by:  $\pi' = \text{SkipRm}(\pi, p)$  if  $p = l_i$  for some  $i$  and

$$\pi' = [p \mapsto Nx_\pi(p)](\text{read } x; \dots; l_{i-1} : \mathcal{I}_{l_{i-1}}; l_{i+1} : \mathcal{I}_{l_{i+1}}; \dots; l_{m+1} : \text{write } y; ).$$

**Theorem 4.** Suppose  $\pi \in pgm$  is a program,  $p \in \text{labels}(\pi)$  is a label in  $\pi$ , and  $\pi' = \text{SkipRm}(\pi, p)$  then  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .

*Proof.* Proof idea: show by induction for any two computation prefixes  $C \in \mathcal{T}_{pfx}(\pi)$  and  $C' \in \mathcal{T}_{pfx}(\pi')$ :

$$\begin{aligned} C &= \pi, v \vdash (p_0, \sigma_0) \rightarrow \dots \rightarrow (p_m, \sigma_m) \\ C' &= \pi', v \vdash (p'_0, \sigma'_0) \rightarrow \dots \rightarrow (p'_n, \sigma'_n) \end{aligned}$$

that for all  $C_2 \in \mathcal{T}_{pfx}(\pi)$  and  $C'_2 \in \mathcal{T}_{pfx}(\pi')$  such that  $C \rightarrow C_2$ ,  $C' \rightarrow C'_2$  and  $C\mathcal{R}C'$  implies  $C_2\mathcal{R}C'_2$  where  $C\mathcal{R}C'$  if and only if  $p \neq p_m = p'_n$  and  $\sigma_m = \sigma'_n$ .

**Base case.** Suppose  $p'_0 \neq p$  then clearly the initial states are equivalent. Otherwise  $\sigma_0 = \sigma'_0$  and  $p'_0 = Nx_\pi(p_0)$ , so it follows by the semantics of **skip** that  $(p_1, \sigma_1) = (p'_0, \sigma'_0)$ .

**Induction step.** Suppose  $p_{m+1} \neq p$ . The label  $p$  does not occur in the program  $\pi'$ , so  $p'_n \neq p$ . The induction assumption  $(p_m, \sigma_m) = (p'_n, \sigma'_n)$  then implies that  $(p_{m+1}, \sigma_{m+1}) = (p'_{n+1}, \sigma'_{n+1})$ .

Otherwise  $p'_{n+1} = p$ . The semantics of **skip** and the definition of  $\pi'$  implies that  $p'_{n+2} = p_{m+1}$ . By the induction assumption it follows that  $\sigma_{m+1} = \sigma'_{n+1} = \sigma'_{n+2}$ . Thus the induction assumption holds for  $m+1$  and  $n+2$ .

By induction for any computation prefix  $C$  of  $\pi$  there exists a computation prefix  $C'$  of  $\pi'$  such that  $(p_m, \sigma_m) = (p'_n, \sigma'_n)$  and conversely for any computation



prefix  $C'$  of  $\pi'$  such that  $p'_n \neq p$  there exists a computation prefix  $C$  of  $\pi$  such that the final states are equivalent.

Clearly  $\pi$  terminates if and only if  $\pi'$  terminates and if they terminate, they compute the same store. Thus the equivalence  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$  follows by the definition of program semantics.  $\square$

**Dead assignment insertion.** Suppose  $\pi \in \text{pgm}$  is a program,  $x \in \text{Variable} \setminus \text{vars}(\pi)$  is a variable not occurring in  $\pi$  and  $e \in \text{expr}$  is an expression such that  $y \in \text{vars}(\pi)$  for all  $y \in \text{vars}(e)$ . Further suppose that for some label  $p$ :  $\mathcal{I}_p = \text{skip}$ . The *dead assignment insertion transformation* is defined as

$$\text{DeadIns}(\pi, p, x, e) = \text{read } x; l_1 : \mathcal{I}'_{l_1}; \dots; l_m : \mathcal{I}'_{l_m}; l_{m+1} : \text{write } y;$$

where the statements in the transformed program are given by

$$\mathcal{I}'_{l_i} = \begin{cases} x := e & \text{if } l_i = p \\ \mathcal{I}_{l_i} & \text{otherwise} \end{cases}.$$

**Theorem 5.** Suppose  $\pi \in \text{pgm}$ ,  $p \in \text{labels}(\pi)$ ,  $x \in \text{Variable} \setminus \text{vars}(\pi)$  and  $e \in \text{expr}$  such that  $y \in \text{vars}(\pi)$  for all  $y \in \text{vars}(e)$ . If  $\pi' = \text{DeadIns}(\pi, p, x, e)$  then  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .

*Proof.* Using Lemma 3 it can be shown by induction that the stores in two computation prefixes of the same length over  $\pi$  and  $\pi'$ , given the same input  $v \in \text{Value}$ , are equivalent when restricted to  $\text{vars}(\pi)$ . This implies semantic program equivalence, since  $y \in \text{vars}(\pi)$ .  $\square$

## 7.1 Common subexpression elimination

At this point it is possible to prove correct the elimination of common subexpressions transformation. In fact the correctness follows directly from the correctness of the transformations already treated since they can be combined to obtain the common subexpression elimination transformation. The transformation detects expressions that are *very busy*, i.e. expressions that are computed eventually on all paths leading from a given program point, where the variables in the expression have not changed their value. An assignment of the expression to a fresh variable is inserted at the particular program point, rendering the subsequent recomputations redundant.

This can be achieved by using some of the transformations already proven formally correct.

1. Decide on a suitable label  $p$ , and fresh variable  $x$  and some expression  $e$ . Usually a compiler would select a label where some expression  $e$  is very busy.
2. Insert a `skip` statement labeled  $p'$  at label  $p$ .
3. Perform the dead assignment insertion transformation at label  $p'$  with variable  $x$  and expression  $e$ .

4. Eliminate as many redundant assignments by repeatedly applying the elimination of recomputation of available expressions transformation at various labels in the program.
5. Delete all **skip** statements from the transformed program.

Since the purpose of the paper is to prove correctness of some typically employed compiler optimizations, we will be satisfied with claiming that it is possible to perform elimination of common subexpressions by the above procedure.

## 8 Transformations with multiple program points

This section presents an extension of the framework that handles transformations in which more than one program statement is affected.

### 8.1 Transformation syntax

**Transformation rules.** A *multi-label transformation rule* has the syntax:

$$\begin{aligned}
\textit{rule} &::= \textit{rw}_1 \dots \textit{rw}_k \text{ if } \textit{cond}_1 \dots \textit{cond}_k \\
\textit{rw} &::= p : \textit{vstmt}_1 \Rightarrow \textit{vstmt}_2 \\
\textit{cond} &::= p \models \phi \\
\textit{vstmt} &::= v := \textit{vexp} \mid \text{if } v \text{ then } \mathbf{v}_1 \text{ else } \mathbf{v}_2 \mid \text{skip} \\
\textit{vexp} &::= c \mid x \mid \mathbf{v} \mid \text{op}(\textit{vexp}_1, \dots, \textit{vexp}_n) \\
\phi &::= \textit{true} \mid \textit{false} \mid \textit{ap}(\mathbf{v}_1, \dots, \mathbf{v}_n) \mid \neg \phi \mid \phi_1 \wedge \phi_2 \\
&\quad \mid \forall y \in Y(\mathbf{v}_1, \dots, \mathbf{v}_n) : \phi \mid \exists y \in Y(\mathbf{v}_1, \dots, \mathbf{v}_n) : \phi \\
&\quad \mid A\psi \mid E\psi \mid \overleftarrow{A}\psi \mid \overleftarrow{E}\psi \\
\psi &::= X\phi \mid \phi_1 \text{ U } \phi_2 \\
Y &::= \textit{vars} \\
\mathbf{v} &\in \textit{FormulaVariable} \\
c &\in \textit{Value} \\
x &\in \textit{Variable} \\
y &\in \textit{QuantifierVariable} \\
p &\in \textit{LabelVariable}
\end{aligned}$$

The main difference from the single rewrite case is that since the transformation now matches against and changes statements at several labels, a device for referring to a specific label is needed. Thus each rewrite rule “*rw*” is now annotated with a label variable. For each rewrite rule labeled by some label, there must exist exactly one side condition annotated with the same label.

As previously noted, the label invariant atomic propositions do not really belong in the CTL-R clauses and in the situation where several clauses are present this is even more pronounced, since there is no natural clause to insert them in. However this approach does avoid introducing inessential concepts.

## 8.2 Concrete transformation

The concrete transformation is computed in a manner similar to the single rewrite rule case. When closing the transformation, a substitution  $\theta'$  is computed such that  $\theta'$  unifies all the left hand sides of the rewrite rules with their corresponding side conditions, for some label tuple  $\vec{p}$ .

**Closed transformation rule.** Suppose  $\pi \in pgm$  is a program,  $\vec{p} = (p_1 \dots p_k) \in Label^*$  is a tuple of labels occurring in  $\pi$ ,  $T$  is a multi-label transformation rule

$$\begin{array}{l} p_1 : lhs_1 \Rightarrow rhs_1 \dots p_k : lhs_k \Rightarrow rhs_k \\ \text{if} \\ p_1 \models \phi_1 \dots p_k \models \phi_k \end{array}$$

and  $\theta$  is a substitution binding all variables not occurring in  $lhs_i$  for any  $i$ . Define  $\theta'$  as the substitution that unifies the left hand side  $lhs_i$  with the statement at label  $p_i$ , for all  $1 \leq i \leq k$ . Provided that such a substitution exists, the multi-label transformation rule  $T'$  can be obtained from  $T$  by closure with respect to  $\theta'' = \theta \circ \theta'$  if  $T' = \theta''(T)$  and  $T'$  contains no free variables. The multi-label transformation rule  $T'$  is called a *closed transformation rule*, since it contains no free variables. The type of  $\theta''$  is given by  $\theta'' : FormulaVariables \rightarrow Label \cup stmt \cup exp \cup Variable$ .

From a closed multi-label transformation rule, the *concrete* transformation rules can be obtained by “unrolling” the program quantifiers in the all the CTL-R clauses in the the side condition.

**Concrete transformation rule.** Suppose  $T$  is the closed multi-label transformation rule as given above, then the *concrete multi-label transformation rule* is computed by unrolling all the side conditions

$$\begin{array}{l} p_1 : lhs_1 \Rightarrow rhs_1 \dots p_k : lhs_k \Rightarrow rhs_k \\ \text{if} \\ p_1 \models unroll(\phi_1) \dots p_k \models unroll(\phi_k). \end{array}$$

Suppose given a multi-label transformation rule  $T$ , program  $\pi \in pgm$  and a label tuple  $\vec{p}$ . The concrete multi-label transformation rule  $T'$  is obtained by *instantiating* the multi-label transformation rule  $T$  by the substitution  $\theta$  at label tuple  $\vec{p}$  in program  $\pi$  if  $T'$  is obtained by first computing the closed transformation rule  $T''$  from  $T$  with respect to  $\theta$  and then computing  $T'$  by unrolling.

## 8.3 Transformation

The definition of application of a multi-label transformation at some specific labels in the subject program is analogous to the single rewrite rule case.

**Apply.** Suppose  $\pi, \pi' \in \text{pgm}$  are programs,  $\vec{p} = (p_1, \dots, p_k) \in \text{Label}^*$  is a tuple of labels in  $\pi$ ,  $T$  is a multi-label transformation rule and  $\theta$  is a substitution. The program  $\pi'$  is obtained by *applying* the multi-label transformation  $T$  at the labels  $\vec{p}$  with substitution  $\theta$  if and only if  $\text{Apply}^*(\pi, \vec{p}, T, \theta, \pi')$ , where the relation  $\text{Apply}^*$  holds if and only if all of the following cases hold.

1. (Instantiation): Given the substitution  $\theta$  there exists an instantiation  $T'$  of  $T$  with the substitution  $\theta'$ , where  $T'$  has the form

$$\begin{aligned} & p_1 : \mathcal{I}_1 \Rightarrow \mathcal{I}'_1 \dots p_k : \mathcal{I}_k \Rightarrow \mathcal{I}'_k \\ \text{if} \\ & p_1 \models \phi_1 \dots p_k \models \phi_k. \end{aligned}$$

2. (Verification): For all  $i = 1 \dots k$  the CTL-R clause  $\mathcal{M}_{cf}(\pi), p_i \models \phi_i$  is satisfied.
3. (Construction): The transformed program  $\pi$  has the form

$$\pi' = \text{read } \mathbf{x}; l_1 : \mathcal{I}'_{l_1}; \dots; l_m : \mathcal{I}'_{l_m}; l_{m+1} : \text{write } \mathbf{y},$$

where the statements of  $\pi'$  are defined by

$$\mathcal{I}'_{l_i} = \begin{cases} \theta(\mathcal{I}_i) & \text{if } \exists j : l_i = p_j \\ \mathcal{I}_i & \text{otherwise} \end{cases}$$

**Corollary 4.** Suppose  $\pi, \pi' \in \text{pgm}$  and  $\text{Apply}^*(\pi, \vec{p}, T, \theta, \pi')$  then  $Nx_\pi = Nx_{\pi'}$ .

## 8.4 Loop invariant hoisting

**Code motion/loop invariant hoisting.** A restricted version of a “code motion” transformation ( $CM$ ) that covers the “Loop invariant hoisting” transformation is defined as

$$\begin{aligned} & p : \text{skip} \Rightarrow \mathbf{x} := \mathbf{e} \\ & q : \mathbf{x} := \mathbf{e} \Rightarrow \text{skip} \\ \text{if} \\ & p \models A(\neg use(\mathbf{x}) \ U \ node(q)) \\ & q \models \overleftarrow{A} X \overleftarrow{A} (\neg def(\mathbf{x}) \wedge \forall y \in vars(\mathbf{e}) : \neg def(y) \wedge \neg entry \ W \ node(p)). \end{aligned}$$

This transformation involves two (different) statements in the subject program. The transformation moves an assignment at label  $q$  to label  $p$  provided that two conditions are met:

1. The assigned variable  $\mathbf{x}$  is dead after  $p$  until  $q$  is reached. If this requirement holds, then introducing the assignment  $\mathbf{x} := \mathbf{e}$  at label  $p$  will not change the semantics of the program.
2. The second requirement (in combination with the first rewrite rule) states that the expression  $\mathbf{e}$  should be available at  $q$  *after* the transformation.

This transformation could be obtained by applying two transformations: One that inserts the statement  $\mathbf{x} := \mathbf{e}$  provided that  $\mathbf{x}$  is dead at  $p$ , followed by the elimination of available expressions transformation (ER). By formulating the transformation as a single transformation, the two labels  $p$  and  $q$  are explicitly linked. With the two transformations one would need some mechanism of controlling where to insert which assignments.

Since all paths from  $p$  *must* end in  $q$ , it is not possible move assignments to labels such that  $\mathbf{e}$  is still available in  $q$  and  $x$  is dead in all paths not leading to  $q$ , which would still be a semantics preserving transformation.

```

1: skip;
2: if ... then 3 else 6;
3: x := a + b;
4: y := y - 1;
5: if y then 3 else 6;
6: x := 0;

```

In this program is is possible to lift the statement  $\mathbf{x} := \mathbf{a} + \mathbf{b}$  from label 3 to label 1. In general the transformation by itself could slow down the computation, since there is no need to compute the expression if the expression is not needed.

**CM relation.** Suppose  $C \in \mathcal{T}_{pfx}(\pi, v)$  and  $C' \in \mathcal{T}_{pfx}(\pi', v)$  for some  $\pi, \pi'$  and  $v$  then define the CM relation on computation prefixes as:  $CR C'$  if and only if  $t = t' \wedge \forall 0 \leq i \leq t : p_i = p'_i$  and one of the following cases holds:

1.  $\sigma_t = \sigma'_t \wedge \forall 0 \leq i \leq t : p_i \notin \{p, q\}$
2.  $\sigma_t = \sigma'_t \wedge \exists 0 \leq i \leq t : p_i = q \wedge \forall i \leq j \leq t : p_j \neq p$
3.  $\exists 0 \leq i \leq t : p_i = p \wedge$   
 $(\sigma_i \setminus x = \sigma'_i \setminus x) \wedge (x \in \text{vars}(e) \Rightarrow \sigma_i(x) = \sigma'_i(x)) \wedge$   
 $\forall i \leq j \leq t : p_j \notin \{p, q\}$

The notation  $CR_i C'$  will be used to indicate that  $CR C'$  holds by case  $i$ , as defined above.

**Theorem 6.** Suppose  $\pi, \pi' \in \text{pgm}$ ,  $\text{Apply}^*(\pi, (p, q), CM, \theta, \pi')$  then  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .

*Proof.* Suppose  $\pi, \pi' \in \text{pgm}$  and  $\text{Apply}^*(\pi, (p, q), CM, \theta, \pi')$ . Let  $v \in \text{Value}$  be the input to program  $\pi, \pi'$  and let  $C \in \mathcal{T}_{pfx}(\pi, v)$ ,  $C' \in \mathcal{T}_{pfx}(\pi', v)$  be two finite computation prefixes of the same length of  $\pi$  and  $\pi'$  respectively.

$$\begin{aligned}
C_1 &= \pi, v \vdash (p_0, \sigma_1) \rightarrow \dots \rightarrow (p_t, \sigma_t) \\
C'_1 &= \pi', v \vdash (p'_0, \sigma'_1) \rightarrow \dots \rightarrow (p'_t, \sigma'_t)
\end{aligned}$$

Suppose the instantiated multi-label transformation rule is given by

```

p : skip  $\Rightarrow$  x := e
q : x := e  $\Rightarrow$  skip
if
p  $\models A(\neg \text{use}(x) \cup \text{node}(q))$ 
q  $\models \overleftarrow{A} X \overleftarrow{A} (\neg \text{def}(x) \wedge \forall y \in \text{vars}(e) : \neg \text{def}(y) \wedge \neg \text{entry } W \text{ node}(p))$ .

```

First we show that  $\mathcal{C}\mathcal{R}\mathcal{C}'$  by induction over the length of  $C, C'$ .

**Base case.** Since  $\text{vars}(\pi) = \text{vars}(\pi')$  then by Lemma 7:  $(p_0, \sigma_0) = (p'_0, \sigma'_0)$ .

**Induction step.** The semantics for the language is deterministic, so there exists exactly one  $C_2$  and one  $C'_2$  such that:  $C \rightarrow C_2$  and  $C' \rightarrow C'_2$ . Assuming that  $\mathcal{C}\mathcal{R}\mathcal{C}'$  we need to show that  $C_2\mathcal{R}C'_2$ .

The proof is spilt up into cases depending on the label  $p_t$  and each of these is further spilt up, depending on the case that the  $CM$  relation  $\mathcal{C}\mathcal{R}\mathcal{C}'$  holds for.

Suppose  $p_t \notin \{p, q\}$ , then  $\mathcal{I}_{p_t} = \mathcal{I}'_{p_t}$ .

1. Suppose  $\mathcal{C}\mathcal{R}_i C'$  where  $i = 1, 2$ . By assumption  $\sigma_t = \sigma'_t$  so  $\pi$  and  $\pi'$  evaluate the same statement in the same store with the same next label function. It then follows from Lemma 2  $(p_{t+1}, \sigma_{t+1}) = (p'_{t+1}, \sigma'_{t+1})$ , so  $C_2\mathcal{R}_i C'_2$ .
2. Otherwise  $\mathcal{C}\mathcal{R}_3 C'$ , so let  $0 \leq i \leq t$  be given such that:  
 $p_i = p \wedge (\sigma_i \setminus x = \sigma'_i \setminus x)$   
 The side condition for  $p$  must be satisfied at  $p_i$ :

$$p_i \models A(\neg \text{use}(x) \cup \text{node}(q)).$$

By the definition of CTL-R this is equivalent to:

$$\begin{aligned} & \forall (p_i = n_0 \rightarrow n_1 \rightarrow n_2 \rightarrow \dots) \in \text{Label}^* \\ & \exists \alpha \geq 0 : n_\alpha \models \text{node}(q) \wedge \forall 0 \leq \beta < \alpha : n_\beta \models \neg \text{use}(x) \end{aligned}$$

Since by Lemma 5 the control flow model describes all possible computation prefixes, the same must hold for the computation prefix  $C$ . By assumption we know that

$$\forall i < j \leq t : p_j \neq q.$$

Thus  $p_j \models \neg \text{use}(x)$  is satisfied for all  $i < j \leq t$ . In particular  $p_t \models \neg \text{use}(x)$  is satisfied and by assumption  $(\sigma_i \setminus x = \sigma'_i \setminus x)$  so using Lemma 4 we conclude:  $\sigma_{t+1} \setminus x = \sigma'_{t+1} \setminus x$ .

Suppose  $x \in \text{vars}(e)$  then  $\sigma_t = \sigma'_t$ , so clearly  $\sigma_{t+1} = \sigma'_{t+1}$ , fulfilling the requirement on the stores in  $C$  and  $C'$ .

It remains to be shown that the labels are equivalent at  $t + 1$ . Suppose  $p_{t+1}$  is given by the next label function. By the induction assumption  $p_t = p'_t$ , so clearly  $p_{t+1} = p'_{t+1}$  holds. Otherwise the statement at  $p_t$  is a conditional and the next state depends on the variable being tested. By the above reasoning  $p_t \models \neg \text{use}(x)$  implying that the conditional does not depend on the  $x$ . By induction assumption  $\sigma_t \setminus x = \sigma'_t \setminus x$ , so  $p_{t+1} = p'_{t+1}$ . Thus  $C_2\mathcal{R}_3 C'_2$  holds for  $i$  unchanged.

Suppose  $p_t = p$ , then  $p_t \models A(\neg use(x) \cup node(q))$ ,  $\mathcal{I}_{p_t} = (\text{skip})$  and  $\mathcal{I}'_{p_t} = (x:=e)$ . By the semantics  $p_{t+1}$  and  $p'_{t+1}$  is given by the next label functions  $Nx_\pi$ ,  $Nx_{\pi'}$ , which by Lemma 2 are identical. Since  $p_t = p'_t$  by the induction assumption, it follows that  $p_{t+1} = p'_{t+1}$ .

1. Suppose  $C\mathcal{R}C'$  by case 1 or 2. The induction assumption implies that  $\sigma_t = \sigma'_t$ , so by the semantics:  $\sigma_{t+1} \setminus x = \sigma'_{t+1} \setminus x$ . Clearly  $C_2\mathcal{R}_2C'_2$  holds for  $i = t$ .
2. Otherwise  $C\mathcal{R}_3C'$ : As previously shown this implies the satisfaction of  $p_t \models \neg use(x)$ , so by Lemma 4 the statement at  $p_t$  does not depend on the value on  $x$ , implying:  $\sigma_{t+1} \setminus x = \sigma'_{t+1} \setminus x$ .  
Suppose  $x \in vars(e)$ , which by the definition of the valuation implies the satisfaction of  $p_t \models use(x)$ , contradicting the above conclusion.  
Thus  $C\mathcal{R}_3C'$  holds for  $i$  unchanged.

Otherwise  $p_t = q$  so  $\mathcal{I}_{p_t} = (x:=e)$  and  $\mathcal{I}'_{p_t} = (\text{skip})$ . By the semantics  $p_{t+1}$  and  $p'_{t+1}$  is given by the next label function, which by Lemma 2 is the identical for  $\pi$  and  $\pi'$ . Since  $p_t = p'_t$  by the induction assumption then  $p_{t+1} = p'_{t+1}$ .

1. Suppose  $C\mathcal{R}_1C'$ : By Lemma 6 the side condition in the rewrite rule

$$p_t \models \overleftarrow{A}X\overleftarrow{A}(\neg def(x) \wedge \forall y \in vars(e) : \neg def(y) \wedge \neg entry\ W\ node(p)).$$

implies that any backwards path eventually ends up in a state  $p_k$  where  $p_k \models node(p)$  is satisfied, i.e.  $p_k = p$ . This contradicts the assumption  $C_1\mathcal{R}_1C'_1$ .

2. Suppose  $C\mathcal{R}_2C'$ : Per assumption there exists an  $i$  such that  $p_i = q$ . Looking at the side condition for  $p_t$  we get:

$$p_t \models \overleftarrow{A}X\overleftarrow{A}(\neg def(x) \wedge \dots \wedge \neg entry\ W\ node(p)).$$

This implies by Lemma 6

$$\begin{aligned} & \forall \text{finite paths} : (\dots \rightarrow p_{t-2} \rightarrow p_{t-1}) \in Label^* : \\ & \exists \alpha < t : p_\alpha \models node(p) \wedge \forall \alpha < \beta \leq t-1 : p_\beta \models \neg def(x) \end{aligned}$$

Since by assumption  $p_j \neq p$  for all  $i \leq j \leq t$  it follows from the definition of the valuation that the clause  $p_j \models \neg node(p)$  is satisfied, so  $\alpha < i < t$ . In particular  $p_i \models \neg def(x)$  is satisfied. But statement  $\mathcal{I}_q = (x:=e)$  clearly defines  $x$ , so  $p_i = q \models def(x)$  contradicting the above conclusion.

3. Otherwise  $C\mathcal{R}_3C'$ : By assumption there exists an  $i$  such that  $p_i = p$ ,  $(\sigma_i \setminus x = \sigma'_i \setminus x)$ ,  $(x \in vars(e) \Rightarrow \sigma_i(x) = \sigma'_i(x))$  and

$$\forall i \leq j \leq t : p_j \neq q \wedge \sigma_j \setminus x = \sigma'_j \setminus x.$$

We wish to show that  $C_2\mathcal{R}_2C'_2$  holds by case 2 for  $i = t$ , i.e.  $\sigma_{t+1} = \sigma'_{t+1}$ .

First we show that  $\sigma_{t+1} \setminus x = \sigma'_{t+1} \setminus x$ . We know from the assumption  $p_t = p$  that the statement at  $p_t$  defines  $x$ :  $\mathcal{M}_{cf}(\pi), p_t \models \text{def}(x)$ , so by Lemma 3:  $\sigma_t \setminus x = \sigma_{t+1} \setminus x$ .

$$\begin{aligned} \sigma_{t+1} \setminus x &= \sigma_t \setminus x && \text{(Lemma 3)} \\ &= \sigma'_t \setminus x && \text{(induction assumption)} \\ &= \sigma'_{t+1} \setminus x && \text{(semantics of skip)} \end{aligned}$$

Next we need to show that  $\sigma_{t+1}(x) = \sigma'_{t+1}(x)$ . First we make an observation:

The assumption  $C\mathcal{R}_3C'$  implies that  $\forall i < k < t: p_k \neq p$ , so the side condition (from  $p_t$ )

$$p_t \models \overleftarrow{A} X \overleftarrow{A} (\neg \text{def}(x) \wedge \forall y \in \text{vars}(e) : \neg \text{def}(y) \wedge \neg \text{entry } W \text{ node}(p))$$

implies that  $\mathcal{M}_{cf}(\pi), p_k \models \neg \text{def}(x)$  and  $\mathcal{M}_{cf}(\pi), p_k \models \neg \text{def}(y)$  for all  $y \in \text{vars}(e)$  and  $i < k < t$ . Clearly the same holds for  $\mathcal{M}_{cf}(\pi')$ , since the statements at  $p_k$  in  $\pi$  will be identical to the statement in  $\pi'$  at  $p_k$ , as  $p_k \notin \{p, q\}$  for  $i < k < t$ . By Lemma 2 variable  $x$  and  $\text{vars}(e)$  does not change from  $i + 1$  to  $t$ :  $\sigma'_{i+1}(x) = \sigma'_k(x)$  and  $\sigma_{i+1}|_{\text{vars}(e)} = \sigma_k|_{\text{vars}(e)}$  for all  $i < k \leq t$ .

Using the above observations it can be shown that  $x$  maps to the same value in  $\sigma_{t+1}$  and  $\sigma'_{t+1}$ :

$$\begin{aligned} \sigma_{t+1}(x) &= \llbracket e \rrbracket \sigma_t && \text{(semantics of } \mathcal{I}_{p_k} = (x := e)) \\ &= \llbracket e \rrbracket \sigma_t|_{\text{vars}(e)} && \text{(Lemma 1)} \\ &= \llbracket e \rrbracket \sigma_{i+1}|_{\text{vars}(e)} && \text{(1<sup>st</sup> argument above)} \\ &= \llbracket e \rrbracket \sigma_i|_{\text{vars}(e)} && \text{(semantics of } \mathcal{I}_{p_i} = (\text{skip})) \\ &= \llbracket e \rrbracket \sigma'_i|_{\text{vars}(e)} && \text{(induction assumption)} \\ &= \llbracket e \rrbracket \sigma'_i && \text{(Lemma 1)} \\ &= \sigma'_{i+1}(x) && \text{(semantics of } \mathcal{I}'_{p_i} = (x := e)) \\ &= \sigma'_t(x) && \text{(2<sup>nd</sup> argument above)} \\ &= \sigma'_{t+1}(x) && \text{(semantics of } \mathcal{I}'_{p_t} = (\text{skip})) \end{aligned}$$

Thus  $\sigma_{t+1} = \sigma'_{t+1}$ , so case 2 holds for  $i = t$ .

Per induction we conclude that  $C\mathcal{R}C'$  for any two computation prefixes  $C \in \mathcal{T}_{pfx}(\pi)$  and  $C' \in \mathcal{T}_{pfx}(\pi')$  of the same length.

Suppose that the program  $\pi$  does not terminate on input  $v$ . Any finite computation prefix of  $\pi$  will have a corresponding computation prefix in  $\pi'$  such that the labels in the computation prefixes are pairwise equivalent. Thus  $\pi'$  will not terminate on input  $v$  either. The reverse direction follows by symmetry. Otherwise  $\pi$  and  $\pi'$  both terminates on input  $v$ . Consider the store at the **write** statement:  $\mathcal{I}_m = \mathcal{I}'_m = (\text{write } y)$ . Suppose  $C\mathcal{R}C'$  for  $t=m$  by case 3. Then by the side condition there exists an  $i > m$  such that  $p_i \models \text{node}(q)$ , i.e.  $p_i = q$ . But by the semantics  $p_i = \text{exit}(\pi) \neq q$  for all  $i > t$  leading to a contradiction. Thus either case 1 or 2 must hold. In either case



$\sigma_t = \sigma'_t$  implying  $\sigma_t(y) = \sigma'_t(y)$ . It follows by Lemma 3 that the transformation is semantics preserving:  $\llbracket \pi \rrbracket = \llbracket \pi' \rrbracket$ .

□

## 9 Summary

In this paper a framework for describing program transformations as conditional rewrite rules has been presented. Specifying program transformations as rewrite rules with temporal logic side conditions is a relatively new approach to proving correctness of program transformations [5, 4].

The main focus in this paper has been on the *correctness* of program transformations rather than orthogonal issues. In order to apply a transformation in the present framework, the overall transformation strategy must decide on a particular program point, before the side condition is verified. In practice this leads to inefficient transformation strategies<sup>2</sup>, if implemented naively. A possible solution[4] is to extend the definition of CTL to admit free variables. The model checker would return a set of substitutions satisfying the side condition, rather than only validating the side condition for a particular label. Each substitution corresponds to an instantiation of the rule schema at particular program points for which the side condition is satisfied. The transformation strategy could then select one of the possible (legal) substitutions and immediately construct the transformed program. The overall development is only slightly more involved than the present paper, and the results from the paper should be easily adaptable. However the simpler framework is sufficient when the main focus is purely on the correctness of program transformations.

The language treated, was a simplistic imperative language without procedure calls or exceptions. In order to allow specification of programs with these language features, the framework would have to be extended to account for the different control flow. This would most likely complicate the proofs since more language constructs are introduced, which would have to be reasoned about in the proofs. However extending the framework to include procedures and exceptions, as well as extensions to object oriented languages, remain an interesting area of future investigations.

## 10 Acknowledgements

The author would like to thank the following people for their kind help during the writing of the present paper. Neil Jones for his guidance and for getting the author started on using temporal logic to prove correctness of program transformations. David Lacey and Eric Van Wyk for inspiration and ideas. Chin Soon Lee for kindly proofreading the paper, resulting in numerous constructive suggestions.

---

<sup>2</sup>The transformation strategy is the order in which a number of program transformations are applied to the given subject program.

## References

- [1] Alfred V. Aho, Ravi Sethi and Jeffrey D. Ullman. Compilers: Principles, Techniques and Tools. Addison-Wesley, 1996.
- [2] Michael R. A. Huth and Mark D. Ryan. Logic in Computer Science: Modelling and reasoning about systems. Cambridge University Press, 1999.
- [3] Joost-Pieter Katoen. Concepts, Algorithms and Tools for Model Checking. Lecture Notes of the Course “Mechanised Validation of Parallel Systems”, Friedrich-Alexander Universität Erlangen-Nürnberg, 1998/1999.
- [4] David Lacey, Neil D. Jones, Eric Van Wyk and Carl C. Frederiksen. Proving Correctness of Compiler Optimizations by Temporal Logic. (Unpublished), 2001.
- [5] David Lacey and Oege de Moor. Imperative Program Transformation by Rewriting. Compiler Construction (CC2001) thread of ETAPS 2001.
- [6] Flemming Nielson, Hanne R. Nielson and Chris Hankin. Principles of Program Analysis. Springer-Verlag, 1999.
- [7] Oege de Moor, David Lacey and Eric Van Wyk. Universal Regular Path Queries. (Submitted to HOSC), 2001.