# Improvement in a Lazy Context:
# An Operational Theory for Call-By-Need
# (Extended Version)

Andrew Moran          David Sands

Chalmers[1]

## Abstract

The standard implementation technique for lazy functional languages is call-by-need, which ensures that an argument to a function in any given call is evaluated at most once. A significant problem with call-by-need is that it is difficult — even for compiler writers — to predict the effects of program transformations. The traditional theories for lazy functional languages are based on call-by-name models, and offer no help in determining which transformations do indeed optimize a program.

In this article we present an operational theory for call-by-need, based upon an improvement ordering on programs: $M$ is improved by $N$ if in all program-contexts $\mathbb{C}$, when $\mathbb{C}[M]$ terminates then $\mathbb{C}[N]$ terminates at least as cheaply.

We show that this improvement relation satisfies a "context lemma", and supports a rich inequational theory, subsuming the call-by-need lambda calculi of Ariola *et al.* [AFM$^+$95]. The reduction-based call-by-need calculi are inadequate as a theory of lazy-program transformation since they only permit transformations which speed up programs by at most a constant factor (a claim we substantiate); we go beyond the various reduction-based calculi for call-by-need by providing powerful proof rules for recursion, including syntactic continuity — the basis of fixed-point-induction style reasoning, and an improvement theorem, suitable for arguing the correctness and safety of recursion-based program transformations.

## 1 Introduction

Call-by-need optimises call-by-name by ensuring that when evaluating a given function application, arguments

[1]Department of Computing Science, Chalmers University of Technology and the University of Göteborg, S-412 96 Göteborg, Sweden; [andrew,dave]@cs.chalmers.se.

are evaluated at most once. All serious compilers for lazy functional languages implement call-by-need evaluation. Lazy functional languages are believed to be well-suited to high-level program transformations, and some state-of-the-art compilers take advantage of this by applying a myriad of transformations and analyses during compilation [PJS98]. However, it is notoriously difficult, even for those with extremely solid intuitions about call-by-need, to predict the effects of a program transformation on the running time. Since traditional theories for lazy languages are based upon call-by-name models, they give no assurance that a given transformation doesn't lead to an asymptotic slow-down.

**Call-by-need Calculi** The call-by-need lambda calculi [AFM$^+$95, AF97, MOW98] offer a solution to some of these problems. By permitting fewer equations than call-by-name, these calculi enable term-level reasoning without ignoring the key implementation issues underpinning call-by-need. However, they do have some serious limitations. All of the equations in the calculi are, by definition, symmetric. This means that certain useful local transformations cannot be present. In fact, the call-by-need calculi are limited to transformations which change running-times by at most a constant-factor (see section 4.6), independent of the context in which the programs are used. Even within the confines of constant-factor transformations there are significant shortcomings, since none of the calculi have proof rules for recursion; we believe that, as a consequence, almost no interesting equivalences between recursive programs — such as the fusion of recursive functions (*e.g.* via deforestation) — can be justified in the calculi.

**Our Approach** We aim to go beyond these limitations by refining the notion of *observational approximation* between terms, and by establishing algebraic laws (containing the laws of the call-by-need calculi as theorems) and recursion principles for that approximation relation. A key result of [AFM$^+$95] is that the standard observational equivalence and approximation relations, in which one only observes termination, cannot distinguish call-by-need evaluation from call-by-name. To obtain an operational theory which retains the computational distinctions between name and need, we also observe the *cost* of evaluation, in terms of a high-level model of computation steps. Our observational approximation relation, *improvement*, is defined with respect to a fixed opera-

tional semantics by saying that: $M$ is improved by $N$ if in all program-contexts $\mathbb{C}$, when $\mathbb{C}[M]$ terminates then $\mathbb{C}[N]$ terminates at least as fast.

**Summary of Results**  We develop an operational theory for a call-by-need lambda calculus with recursive lets, constructors, and case expressions. The theory is based upon an abstract machine semantics for call-by-need, and is cost-sensitive, and therefore reflects the computational distinctions between call-by-name and call-by-need. We show that the improvement relation has a rich inequational theory, validating the reduction rules of the call-by-need calculi. Most importantly, it supports powerful induction principles for recursive programs. Some specific original results are:

- A *context lemma* for call-by-need, meaning we can establish improvement by considering just computation in a restricted class of contexts, the evaluation contexts;

- A rich inequational theory, the *tick algebra*, which subsumes the call-by-need calculi;

- A *syntactic continuity* property which characterises improvement of a recursive function in terms of its finite unwindings, and forms the basis of fixed-point induction style proofs, and

- Two powerful proof techniques, the *improvement theorem* and *improvement induction*, which are particularly well-suited to inferring the correctness *and* safety of recursion-based program transformations which proceed by local improvements.

**Overview**  We begin with a discussion of related work in section 2. Section 3 then presents the operational semantics (Sestoft's "mark 1" abstract machine for laziness). This is used as the basis for a contextual definition of improvement and cost equivalence, and the context lemma is stated. The inequational theory, known as the tick algebra, is presented in section 4, and the relative power of the algebra and the call-by-need calculi is discussed. Syntactic continuity is presented in section 5 and used to show that an unwinding fixed-point combinator is improved (up to a constant factor) by a knot-tying fixed-point combinator. We also present a syntactic variant of fixed-point fusion for call-by-need, which can be established via syntactic continuity. The improvement theorem is introduced in section 6, along with improvement induction and examples of their use. Section 7 concludes, and we discuss of future avenues of research. For the interested reader, appendix A summarises the technical development and presents the proofs of the main theorems.

## 2  Related Work

Improvement theory and the improvement theorem were originally developed in the call-by-name setting [San91, San96], and generalised to a variety of call-by-name and call-by-value languages in [San97]. Whether this programme could be carried out in a call-by-need setting has long been an open question. An inspiration which gave us confidence in the possibility of a tractable improvement theory for call-by-need is the *call-by-need lambda calculus* presented by Ariola and Felleisen, and Maraist, Odersky and Wadler [AFM+95, AF97, MOW98]. For us, the significance of the call-by-need calculi is that they are based on reduction (and hence equations) between terms in the source language (see figure 7), rather than, say, term-graphs, abstract-machine configurations, or terms plus explicit substitutions. The reduction rules are confluent, and enjoy a deterministic notion of standard reduction. Related concepts appear in other approaches, in particular in the study of so-called optimal reductions *e.g.*, [Fie90, Mar91, Yos93].

One limitation of the original work by Ariola *et al.* is in the treatment of recursive *cycles*; naïve extension of the calculi to deal with recursive lets leads to a loss of confluence [Jef93, AK97]. The original call-by-need calculus considers recursive lets only briefly. To recover confluence, one can simply disallow reductions under cycles, as in *e.g.*, [BLR96, Nie96]. Ariola and Blom give a full study of cyclic recursion in [AB97, AB98], and show that an approximation to confluence can be obtained by equating terms with the same infinite normal-form. Their $\lambda_{\text{oshare}}$ calculus can be seen as the natural successor to the call-by-need calculi.

In general, reduction calculi appear to be a good vehicle for exploring the language design space with regard to call-by-need-like features. Rose's work *e.g.* [Ros96, BLR96] exemplifies this approach in an elegant combination of explicit substitution and combinatory reduction systems. Our view is complementary to the rewriting approaches: once a particular operational semantics (reduction strategy) has been fixed, one can go beyond the confines of the calculi by developing an operational theory.

Apart from the rewriting-based approaches, there have been a few attempts to give a high-level semantics to call-by-need *e.g.* [Jos89, Jef94, SPI96, Lau93, Ses97]. Launchbury's natural semantics, and Sestoft's abstract machine(s) have been adopted by a number of researchers as the formal definition of call-by-need *e.g.* [TWM95, HM95, SPJ97, Gus98]. Since it appears to be a non-controversial choice, we adopt Sestoft's machine — essentially a Krivine-machine [Cur91] with updating of the heap — as the operational model underpinning our theory. As others have observed (*e.g.* [Pit97a]), working with an abstract machine rather than an inductive semantics also has benefits in proofs about computations (examples of this may be seen in the appendix).

## 3  The Operational Semantics

Our language is an untyped lambda calculus with recursive lets, structured data, and case expressions. We work with a restricted syntax in which arguments to functions (including constructors) are always variables:

$$L, M, N ::= x \mid \lambda x.M \mid M\,x \mid c\,\vec{x}$$
$$\mid \text{ let } \{\vec{x} = \vec{M}\} \text{ in } N$$
$$\mid \text{ case } M \text{ of } \{c_i\,\vec{x}_i \rightarrow N_i\}.$$

The syntactic restriction is now rather standard, following its use in core language of the Glasgow Haskell compiler, *e.g.*, [PJPS96, PJS98], and in [Lau93, Ses97].

All constructors have a fixed arity, and are assumed to be saturated. By $c\,\vec{x}$ we mean $c\,x_1\cdots x_n$. The only values are lambda expressions and fully-applied constructors. Throughout, $x, y, z$, and $w$ will range over variables, $c$ over constructor names, and $V$ and $W$ over values. We will write let $\{\vec{x} = \vec{M}\}$ in $N$ as a shorthand for

$$\text{let } \{x_1 = M_1, \ldots, x_n = M_n\} \text{ in } N$$

where the $\vec{x}$ are distinct, the order of bindings is not syntactically significant, and the $\vec{x}$ are considered bound in $N$ *and* the $\vec{M}$ (so our lets are recursive). Similarly we write case $M$ of $\{c_i\,\vec{x}_i \rightarrow N_i\}$ for

$$\text{case } M \text{ of } \{c_1\,\vec{x}_1 \rightarrow N_1\,|\cdots|c_m\,\vec{x}_m \rightarrow N_m\}.$$

where each $\vec{x}_i$ is a vector of distinct variables, and the $c_i$ are distinct constructors. In addition, we will sometimes write *alts* as an abbreviation for case alternatives $\{c_i\,\vec{x}_i \rightarrow N_i\}$.

For examples, working with a restricted syntax can be cumbersome, so it is sometimes useful to lift the restriction. Where we do this it should be taken that

$$M\,N \equiv \text{let } \{x = N\} \text{ in } M\,x, \quad x \text{ fresh}$$

whenever $N$ is not a variable. Similarly for constructor expressions.

The only kind of substitution that we consider is *variable for variable*, with $\sigma$ ranging over such substitutions. The simultaneous substitution of one vector of variables for another will be written $M[\vec{y}/\vec{x}]$, where the $\vec{x}$ are assumed to be distinct (but the $\vec{y}$ need not be).

## 3.1 The Abstract Machine

The semantics presented in this section is essentially Sestoft's "mark 1" abstract machine for laziness [Ses97]. In that paper, he proves his abstract machine semantics sound and complete with respect to Launchbury's natural semantics, and we will not repeat those proofs here.

Transitions are over configurations consisting of a heap, containing bindings, the expression currently being evaluated, and a stack. The heap is a partial function from variables to terms, and denoted in an identical manner to a collection of let-bindings. The stack may contain variables (the arguments to applications), case alternatives, or *update markers* denoted by $\#x$ for some variable $x$. Update markers ensure that a binding to $x$ will be recreated in the heap with the result of the current evaluation; this is how sharing is maintained in the semantics.

We write $\langle\,\Gamma,\ M,\ S\,\rangle$ for the abstract machine configuration with heap $\Gamma$, expression $M$, and stack $S$. We denote the empty heap by $\emptyset$, and the addition of a group of bindings $\vec{x} = \vec{M}$ to a heap $\Gamma$ by juxtaposition: $\Gamma\{\vec{x} = \vec{M}\}$. The stack written $b : S$ will denote the a stack $S$ with $b$ pushed on the top. The empty stack is denoted by $\epsilon$, and the concatenation of two stacks $S$ and $T$ by $ST$ (where $S$ is on top of $T$).

We will refer to the set of variables bound by $\Gamma$ as dom $\Gamma$, and to the set of variables marked for update in a stack $S$ as dom $S$. Update markers should be thought of as binding occurrences of variables. A configuration

is *well-formed* if dom $\Gamma$ and dom $S$ are disjoint. We write dom$(\Gamma, S)$ for their union. For a configuration $\langle\,\Gamma,\ M,\ S\,\rangle$ to be closed, any free variables in $\Gamma$, $M$, and $S$ must be contained in dom$(\Gamma, S)$. For sets of variables $P$ and $Q$ we will write $P \nmid Q$ to mean that $P$ and $Q$ are disjoint, *i.e.*, $P \cap Q = \emptyset$. The free variables of a term $M$ will be denoted $\mathsf{FV}(M)$; for a vector of terms $\vec{M}$, we will write $\mathsf{FV}(\vec{M})$.

The abstract machine semantics is presented in figure 3.1; we implicitly restrict the definition to well-formed configurations. There are seven rules, which can grouped as follows. Rules (*Lookup*) and (*Update*) concern evaluation of variables. To begin evaluation of $x$, we remove the binding $x = M$ from the heap and start evaluating $M$, with $x$, marked for update, pushed onto the stack. Rule (*Update*) applies when this evaluation is finished, and we may update the heap with the new binding for $x$.

Rules (*Unwind*) and (*Subst*) concern function application: rule (*Unwind*) pushes an argument onto the stack while the function is being evaluated; once a lambda expression has been obtained, rule (*Subst*) retrieves the argument from the stack and substitutes it into the body of that lambda expression.

Rules (*Case*) and (*Branch*) govern the evaluation of case expressions. Rule (*Case*) initiates evaluation of the case expression, with the case alternatives pushed onto the stack. Rule (*Branch*) uses the result of this evaluation to choose one of the branches of the case, performing substitution of the constructor's arguments for the branch's pattern variables.

Lastly, rule (*Letrec*) adds a set of bindings to the heap. The side condition ensures that no inadvertent name capture occurs, and can always be satisfied by a local $\alpha$-conversion.

**Definition 3.1 (Convergence)** *For closed configurations* $\langle\,\Gamma,\ M,\ S\,\rangle$,

$$\langle\,\Gamma,\ M,\ S\,\rangle\Downarrow^n \stackrel{\text{def}}{=} \exists\Delta, V.\langle\,\Gamma,\ M,\ S\,\rangle \rightarrow^n \langle\,\Delta,\ V,\ \epsilon\,\rangle,$$

$$\langle\,\Gamma,\ M,\ S\,\rangle\Downarrow \stackrel{\text{def}}{=} \exists n.\langle\,\Gamma,\ M,\ S\,\rangle\Downarrow^n,$$

$$\langle\,\Gamma,\ M,\ S\,\rangle\Downarrow^{\leqslant n} \stackrel{\text{def}}{=} \exists m.\langle\,\Gamma,\ M,\ S\,\rangle\Downarrow^m \wedge\ m \leqslant n.$$

Closed configurations which do not converge are of three types: they either reduce indefinitely, get stuck because of a type error, or get stuck because of a *black-hole* (a self-dependent expression as in let $x = x$ in $x$). All non-converging configurations will be semantically identified.

We will also write $M\Downarrow$, $M\Downarrow^n$ and $M\Downarrow^{\leqslant n}$, identifying closed $M$ with the initial configuration $\langle\,\emptyset,\ M,\ \epsilon\,\rangle$.

## 3.2 Program Contexts

The starting point for an operational theory is usually an approximation and an equivalence defined in terms of *program contexts*. Program contexts are usually introduced as "programs with holes", the intention being that an expression is to be "plugged into" all of the holes in the context. The central idea is that to compare the behaviour of two terms one should compare their behaviour in all program contexts.

$$\langle \Gamma\{x = M\},\ x,\ S \rangle \to \langle \Gamma,\ M,\ \#x : S \rangle \qquad\qquad (\textit{Lookup})$$

$$\langle \Gamma,\ V,\ \#x : S \rangle \to \langle \Gamma\{x = V\},\ V,\ S \rangle \qquad\qquad (\textit{Update})$$

$$\langle \Gamma,\ M\ x,\ S \rangle \to \langle \Gamma,\ M,\ x : S \rangle \qquad\qquad (\textit{Unwind})$$

$$\langle \Gamma,\ \lambda x.M,\ y : S \rangle \to \langle \Gamma,\ M[\textstyle\frac{y}{x}],\ S \rangle \qquad\qquad (\textit{Subst})$$

$$\langle \Gamma,\ \mathsf{case}\ M\ \mathsf{of}\ \textit{alts},\ S \rangle \to \langle \Gamma,\ M,\ \textit{alts} : S \rangle \qquad\qquad (\textit{Case})$$

$$\langle \Gamma,\ c_j\ \vec{y},\ \{c_i\ \vec{x}_i \to N_i\} : S \rangle \to \langle \Gamma,\ N_j[\textstyle\frac{\vec{y}}{\vec{x}_j}],\ S \rangle \qquad\qquad (\textit{Branch})$$

$$\langle \Gamma,\ \mathsf{let}\ \{\vec{x} = \vec{M}\}\ \mathsf{in}\ N,\ S \rangle \to \langle \Gamma\{\vec{x} = \vec{M}\},\ N,\ S \rangle \quad \vec{x} \not\in \mathrm{dom}(\Gamma, S) \qquad (\textit{Letrec})$$

Figure 1: The abstract machine semantics for call-by-need.

We will use contexts of the following form:

$$\mathbb{C}, \mathbb{D} ::= [\cdot] \mid x \mid \lambda x.\mathbb{C} \mid \mathbb{C}\ x \mid c\ \vec{x}$$
$$\mid\ \mathsf{let}\ \{\vec{x} = \vec{\mathbb{C}}\}\ \mathsf{in}\ \mathbb{D} \mid \mathsf{case}\ \mathbb{C}\ \mathsf{of}\ \{c_i\ \vec{x}_i \to \mathbb{D}_i\}.$$

Our contexts may contain zero or more occurrences of the hole, and as usual the operation of filling a hole with a term can cause variables in the term to become captured.

The relationship between terms and configurations is characterised by a translation function from configurations to terms, defined inductively on the stack:

$$\mathsf{trans}\langle \emptyset,\ M,\ \epsilon \rangle = M$$
$$\mathsf{trans}\langle \{\vec{x} = \vec{M}\},\ N,\ \epsilon \rangle = \mathsf{let}\ \{\vec{x} = \vec{M}\}\ \mathsf{in}\ N$$
$$\mathsf{trans}\langle \Gamma,\ M,\ x : S \rangle = \mathsf{trans}\langle \Gamma,\ M\ x,\ S \rangle$$
$$\mathsf{trans}\langle \Gamma,\ M,\ \#x : S \rangle = \mathsf{trans}\langle \Gamma\{x = M\},\ x,\ S \rangle$$
$$\mathsf{trans}\langle \Gamma,\ M,\ \textit{alts} : S \rangle = \mathsf{trans}\langle \Gamma,\ \mathsf{case}\ M\ \mathsf{of}\ \textit{alts},\ S \rangle$$

The following lemma clarifies the relationship:

**Lemma 3.1 (Translation)** *For all* $\Gamma$, $\mathbb{C}$, $S$, *there exists* $k \geqslant 0$ *such that for any* $M$, $\langle \emptyset,\ \mathsf{trans}\langle \Gamma,\ \mathbb{C}[M],\ S \rangle,\ \epsilon \rangle \to^k \langle \Gamma,\ \mathbb{C}[M],\ S \rangle.$

PROOF. Simple induction on the size of $S$. $\qquad\square$

### 3.3 Improvement

We define observational approximation and equivalence via contexts in the standard way [AO93].

**Definition 3.2 (Observational Approximation)** *We say that* $M$ *observationally approximates* $N$, *written* $M \mathrel{\underset{\sim}{\sqsubseteq}} N$, *if for all* $\mathbb{C}$ *such that* $\mathbb{C}[M]$ *and* $\mathbb{C}[N]$ *are closed,*

$$\mathbb{C}[M]{\Downarrow} \implies \mathbb{C}[N]{\Downarrow}.$$

We say that $M$ and $N$ are *observationally equivalent*, written $M \cong N$, when $M \mathrel{\underset{\sim}{\sqsubseteq}} N$ and $N \mathrel{\underset{\sim}{\sqsubseteq}} M$.

We know that $\cong$ coincides with its call-by-name counterpart, so this tells us nothing new. We need to incorporate more intensional information if we are to build an operational theory that retains the distinction between name and need. Since call-by-need may be thought of as an optimisation of call-by-name, a natural intensional property to compare is how many reduction steps are required for termination.

**Definition 3.3 (Improvement)** *We say that* $M$ *is improved by* $N$, *written* $M \mathrel{\underset{\sim}{\gtrsim}} N$, *if for all* $\mathbb{C}$ *such that* $\mathbb{C}[M]$ *and* $\mathbb{C}[N]$ *are closed,*

$$\mathbb{C}[M]{\Downarrow}^n \implies \mathbb{C}[N]{\Downarrow}^{\leqslant n}.$$

We say that $M$ and $N$ are *cost equivalent*, written $M \mathrel{\underset{\sim}{\Leftrightarrow}} N$, when $M \mathrel{\underset{\sim}{\gtrsim}} N$ and $N \mathrel{\underset{\sim}{\gtrsim}} M$.

This definition suffers from the same problem as any contextual definition: to prove that two terms are related requires one to examine their behaviour in *all* contexts. For this reason, it is common to seek to prove a *context lemma* [Mil77] for an operational semantics: one tries to show that to prove $M$ observationally approximates $N$, one only need compare their behaviour with respect to a much smaller set of contexts.

We have established the following context lemma for call-by-need:

**Lemma 3.2 (Context Lemma)** *For all terms* $M$ *and* $N$, *if for all* $\Gamma$ *and* $S$,

$$\langle \Gamma,\ M,\ S \rangle{\Downarrow}^n \implies \langle \Gamma,\ N,\ S \rangle{\Downarrow}^{\leqslant n}$$

*then* $M \mathrel{\underset{\sim}{\gtrsim}} N$.

It says that we need only consider configuration contexts of the form $\langle \Gamma,\ [\cdot],\ S \rangle$ where the hole $[\cdot]$ appears only once. This corresponds exactly to a subset of term contexts called *evaluation contexts*, in which the hole is the subject of evaluation. We shall make this correspondence precise in the section 4.2.

Note that the context lemma applies to open terms $M$ and $N$. It is more common to restrict one's attention to closed terms, and then show that the preorder in question is closed under (general) substitution.

### 3.4 Strong Improvement

The improvement relation, like the notion of operational approximation which it refines, also increases the termination of programs, so if $M \mathrel{\underset{\sim}{\gtrsim}} N$ then $N$ may also terminate "more often" than $M$. In the context of compiler optimisations it is natural to ask for a stronger notion of improvement which does not permit any change in termination behaviour.
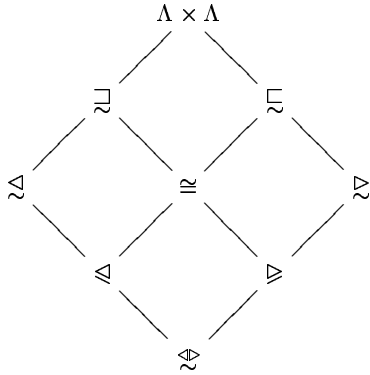
**Definition 3.4 (Strong Improvement)** *We say that $M$ is strongly improved by $N$, written $M \unrhd N$, if*

$$M \overset{\sim}{\unrhd} N \wedge N \overset{\sim}{\sqsubseteq} M$$

$M$ is strongly improved by $N$ if it is improved by $N$, and $N$ has identical termination behaviour (note that we need only have $N \overset{\sim}{\sqsubseteq} M$ in the definition since $M \overset{\sim}{\unrhd} N \implies M \overset{\sim}{\sqsubseteq} N$).

For simplicity of presentation we emphasise improvement rather than strong improvement. However, almost all the laws and proof rules presented in subsequent sections also hold for strong improvement. Notable exceptions being the "strictness laws" concerning $\Omega$, the divergent term. The syntactic continuity proof principle is sound for strong improvement, but degenerates to a trivial rule.

The following Hasse-diagram illustrates the relationships between the various approximations and equivalences introduced in this section:



The diagram is a $\cap$-semi-lattice of relations on terms. In other words, the greatest lower bound of any two relations in the diagram is equal to their set-intersection.

## 4   The Tick Algebra

Consider the following improvement:

$$(\lambda x.M)\, y \overset{\sim}{\unrhd} M[y/x] \qquad\qquad (*)$$

Clearly, for any $\Gamma$ and $S$:

$$\langle\, \Gamma,\ (\lambda x.M)\, y,\ S \,\rangle \to \langle\, \Gamma,\ \lambda x.M,\ y : S \,\rangle$$
$$\to \langle\, \Gamma,\ M[y/x],\ S \,\rangle,$$

so $(*)$ follows from the context lemma. But we can say more: $(\lambda x.M)\, y$ *always* takes exactly two more steps to converge than $M[y/x]$. If we had some syntactic way of slowing the right-hand side down, $(*)$ could be written as a cost equivalence, which would be preferable, since it is a more informative statement. This motivates the introduction of the "tick", written $\checkmark$, which we will use to add a dummy step to a computation. Now we can write $(*)$ as

$$(\lambda x.M)\, y \overset{\sim}{\Lleftarrow\!\!\Rrightarrow} {}^{\checkmark\checkmark} M[y/x] \qquad\qquad (\beta)$$

We can define $\checkmark$ within the language[2] as an empty let binding, thus:

$$\checkmark M \overset{\text{def}}{=} \text{let } \{\} \text{ in } M.$$

Clearly, $\checkmark$ adds one unit to the cost of evaluating $M$ without otherwise changing its behaviour. Note that:

$$M\Downarrow \iff \checkmark M\Downarrow$$
$$M\Downarrow^{n} \iff \checkmark M\Downarrow^{n+1}$$

We will write ${}^{k\checkmark} M$ to mean that $M$ has been slowed down by $k$ ticks. The following inference rule and axiom, known collectively as "tick elimination" are crucial when establishing improvement or cost equivalence.

$$\frac{\checkmark M \overset{\sim}{\unrhd} \checkmark N}{M \overset{\sim}{\unrhd} N} \qquad\qquad \checkmark M \overset{\sim}{\unrhd} M \qquad\qquad (\checkmark\text{-}elim)$$

Their validity follows from the definition of $\overset{\sim}{\unrhd}$.

We can easily prove a number of improvements and cost equivalences modulo tick, and we present a selection of the more useful ones in the following sections. Throughout, we will follow the standard convention that all bound variables in the statement of a law are distinct, and that they are disjoint from the free variables. Together with $(\checkmark\text{-}elim)$, the laws presented in figures 2, 3, 4, 5, and figure 6 are known collectively as the *tick algebra*.

### 4.1   Beta Laws

The first set of laws, presented in figure 2, are important in that they allow us to mimic evaluation within the algebra. We have already seen $(\beta)$; $(case\text{-}\beta)$ is the analogous law for case expressions. In $(value\text{-}\beta)$, one may replace occurrences of a variable, which is bound to some value $V$, with ${}^{2\checkmark} V$. The ticks reflect the fact that by replacing $x$ with its value, we are short-circuiting one lookup and one update step.

The proofs of validity of $(value\text{-}\beta)$, $(var\text{-}\beta)$, $(var\text{-}abs)$, and $(var\text{-}subst)$ rely upon general techniques that are outlined in the appendix.

There are also two derived beta laws, corresponding to unrestricted versions of $(\beta)$ and $(case\text{-}\beta)$. We can derive the following cost equivalence (modulo tick):

$$(\lambda x.M)\, N \overset{\sim}{\Lleftarrow\!\!\Rrightarrow} {}^{2\checkmark} \text{let } \{x = N\} \text{ in } M \qquad (\beta')$$

where $N$ is not a variable. There is a similar derived law for general case expressions.

### 4.2   Laws for Evaluation Contexts

An *evaluation context* is a context in which the hole is the target of evaluation; in other words, evaluation cannot

---

[2] We could introduce a new syntactic construct instead, with the semantics of a dummy step. Since we are working with contextual approximations and equivalences, this would not necessarily be a conservative extension: we would have to prove that it could be represented in the original language. By defining $\checkmark$ in the language, we neatly sidestep such considerations.

# Laws of the Tick Algebra

Throughout, we follow the standard convention that all bound variables in the statement of a law are distinct, and that they are disjoint from the free variables.

$$(\lambda x.M)\, y \;\cong^{2^\checkmark}\; M[y/x] \tag{$\beta$}$$

$$\text{case } c_j\ \vec{y} \text{ of } \{c_i\ \vec{x}_i \to M_i\} \;\cong^{2^\checkmark}\; M_j[\vec{y}/\vec{x}_j] \tag{$case\text{-}\beta$}$$

$$\text{let } \{x = V, \vec{y} = \vec{\mathbb{D}}[x]\} \text{ in } \mathbb{C}[x] \;\cong^{2^\checkmark}\; \text{let } \{x = V, \vec{y} = \vec{\mathbb{D}}[^{2^\checkmark}V]\} \text{ in } \mathbb{C}[^{2^\checkmark}V] \tag{$value\text{-}\beta$}$$

$$\text{let } \{x = z, \vec{y} = \vec{\mathbb{D}}[x]\} \text{ in } \mathbb{C}[x] \;\gtrsim\; \text{let } \{x = z, \vec{y} = \vec{\mathbb{D}}[z]\} \text{ in } \mathbb{C}[z] \tag{$var\text{-}\beta$}$$

$$\text{let } \{x = z, \vec{y} = \vec{\mathbb{D}}[^{2^\checkmark}z]\} \text{ in } \mathbb{C}[^{2^\checkmark}z] \;\gtrsim\; \text{let } \{x = z, \vec{y} = \vec{\mathbb{D}}[x]\} \text{ in } \mathbb{C}[x] \tag{$var\text{-}abs$}$$

$$\text{let } \{x = z, \vec{y} = \vec{M}\} \text{ in } N \;\gtrsim\; \text{let } \{x = z, \vec{y} = \vec{M}[z/x]\} \text{ in } N[z/x] \tag{$var\text{-}subst$}$$

Figure 2: Beta laws for call-by-need.

$$\mathbb{E}[\text{case } M \text{ of } \{pat_i \to N_i\}] \;\cong\; \text{case } M \text{ of } \{pat_i \to \mathbb{E}[N_i]\} \tag{$case\text{-}\mathbb{E}$}$$

$$\mathbb{E}[\text{let } \{\vec{x} = \vec{M}\} \text{ in } N] \;\cong\; \text{let } \{\vec{x} = \vec{M}\} \text{ in } \mathbb{E}[N] \tag{$let\text{-}\mathbb{E}$}$$

Figure 3: Laws for evaluation contexts.

$$\text{let } \{\vec{x} = \vec{L}, \vec{y} = \vec{M}\} \text{ in } N \;\cong\; \text{let } \{\vec{y} = \vec{M}\} \text{ in } N, \quad \text{if } \vec{x} \not\in \mathsf{FV}(\vec{M}, N) \tag{$gc$}$$

$$\text{let } \{\vec{x} = \vec{L}\} \text{ in let } \{\vec{y} = \vec{M}\} \text{ in } N \;\cong^{\checkmark}\; \text{let } \{\vec{x} = \vec{L}, \vec{y} = \vec{M}\} \text{ in } N \tag{$let\text{-}flatten$}$$

$$\text{let } \{x = \text{let } \{\vec{y} = \vec{L}, \vec{z} = \vec{M}\} \text{ in } N\} \text{ in } N' \;\cong\; \text{let } \{x = \text{let } \{\vec{z} = \vec{M}\} \text{ in } N, \vec{y} = \vec{L}\} \text{ in } N' \tag{$let\text{-}let$}$$

$$^{\checkmark}(\lambda x.\text{let } \{\vec{y} = \vec{V}, \vec{z} = \vec{M}\} \text{ in } N) \;\cong\; \text{let } \{\vec{y} = \vec{V}\} \text{ in } \lambda x.\text{let } \{\vec{z} = \vec{M}\} \text{ in } N \tag{$let\text{-}float\text{-}val$}$$

$$\text{let } \{x = M, \vec{y} = \vec{N}\} \text{ in } x \;\cong\; \text{let } \{x = M, \vec{y} = \vec{N}\} \text{ in } ^{2^\checkmark}M, \quad \text{if } x \notin \mathsf{FV}(M, \vec{N}) \tag{$inline$}$$

$$\text{let } \{\vec{x} = \vec{V}\sigma_1, \vec{y} = \vec{V}\sigma_2, \vec{z} = \vec{M}\} \text{ in } N \;\cong\; \text{let } \{\vec{x} = \vec{V}\sigma_2\sigma_3, \vec{z} = \vec{M}\sigma_3\} \text{ in } N\sigma_3, \quad \sigma_1 = [\vec{y}/\vec{w}], \sigma_2 = [\vec{x}/\vec{w}], \sigma_3 = [\vec{x}/\vec{y}], \tag{$value\text{-}copy$}$$

Figure 4: Laws for dealing with lets.

$$\Omega \;\gtrsim\; M \tag{$\Omega$}$$

$$M \;\gtrsim\; \Omega, \quad \text{iff } M \cong \Omega \tag{$imp\text{-}\Omega$}$$

$$M \cong \Omega, \quad \text{iff } M \gtrsim {}^{\checkmark}M \tag{$diverge$}$$

$$\text{let } \{x = \Omega, \vec{y} = \vec{\mathbb{D}}[x]\} \text{ in } \mathbb{C}[x] \;\cong\; \text{let } \{x = \Omega, \vec{y} = \vec{\mathbb{D}}[\Omega]\} \text{ in } \mathbb{C}[\Omega] \tag{$\Omega\text{-}\beta$}$$

$$^{\checkmark}(\lambda x.\text{let } \{y = \Omega, \vec{z} = \vec{M}\} \text{ in } N) \;\cong\; \text{let } \{y = \Omega\} \text{ in } \lambda x.\text{let } \{\vec{z} = \vec{M}\} \text{ in } N \tag{$let\text{-}float\text{-}\Omega$}$$

$$\mathbb{C}[^{\checkmark}M] \;\gtrsim\; {}^{\checkmark}\mathbb{C}[M], \quad \text{if } \mathbb{C} \text{ is strict} \tag{$\checkmark\text{-}float$}$$

Figure 5: Laws for $\Omega$ and strictness.

$$\text{let } \{x = M, \vec{y} = \vec{\mathbb{D}}[^{2^\checkmark}M]\} \text{ in } \mathbb{C}[^{2^\checkmark}M] \;\gtrsim\; \text{let } \{x = M, \vec{y} = \vec{\mathbb{D}}[x]\} \text{ in } \mathbb{C}[x] \tag{$\beta\text{-}expand$}$$

Figure 6: Beta expansion conjecture.

proceed until the hole is filled. Evaluation contexts have the following form:

$$\mathbb{E} ::= \mathbb{A} \mid \mathsf{let} \; \{\vec{x} = \vec{M}\} \; \mathsf{in} \; \mathbb{A}$$
$$\mid \; \mathsf{let} \; \{\vec{y} = \vec{M}, x_0 = \mathbb{A}_0[x_1], x_1 = \mathbb{A}_1[x_2], \ldots, x_n = \mathbb{A}_n\}$$
$$\mathsf{in} \; \mathbb{A}[x_0]$$
$$\mathbb{A} ::= [\cdot] \mid \mathbb{A} \, x \mid \mathsf{case} \; \mathbb{A} \; \mathsf{of} \; \{c_i \; \vec{x}_i \to M_i\}.$$

$\mathbb{E}$ ranges over evaluation contexts, and $\mathbb{A}$ over what we call *applicative contexts*. Our evaluation contexts are strictly contained in those mentioned in the letrec extension of Ariola and Felleisen [AF97]: there they allow $\mathbb{E}$ to appear anywhere we have an $\mathbb{A}$. Our "flattened" definition corresponds exactly to configuration contexts (with a single hole) of the form $\langle \Gamma, [\cdot], S \rangle$, as stated by the following lemma, where $\Lambda_\mathbb{E}$ is the set of all evaluation contexts.

**Lemma 4.1** $\Lambda_\mathbb{E} = \{\mathsf{trans}\langle \Gamma, [\cdot], S \rangle \mid all \; \Gamma, S\}$.

The two laws in figure 3 are very useful indeed: they allow us to move cases and lets in and out of evaluation contexts. A common motif in proofs using the tick algebra is the use of (*case*-$\mathbb{E}$) and (*let*-$\mathbb{E}$) to expose the sub-term of interest. Their validity follows easily from a simple lemma (given in the appendix).

We can derive a law which allows us to move ticks in and out of evaluation contexts:

$$\mathbb{E}[\checkmark M] \underset{\sim}{\gtrless} \checkmark \mathbb{E}[M] \qquad\qquad (\checkmark\text{-}\mathbb{E})$$

Since $\checkmark$ is just a let expression, this is an instance of (*let*-$\mathbb{E}$). Another derived law is the following:

$$\mathsf{let} \; \{x = M\} \; \mathsf{in} \; \mathbb{E}[x] \underset{\sim}{\gtrless} {}^{3\checkmark}\mathbb{E}[M], \quad x \notin \mathsf{FV}(M, \mathbb{E})$$
$$(inline\text{-}\mathbb{E})$$

allowing us to inline $x$ when used but once in an evaluation context. It follows by (*let*-$\mathbb{E}$), (*inline*), and (*gc*).

## 4.3 Concerning Lets

Some of the laws that allow us to manipulate lets are presented in figure 4. Law (*gc*) corresponds to garbage collection: it allows us to add or remove superfluous bindings. A derived garbage collection cost equivalence is the following:

$$\mathsf{let} \; \{x = M\} \; \mathsf{in} \; N \underset{\sim}{\gtrless} \checkmark N, \quad x \notin \mathsf{FV}(N)$$

It follows from (*gc*) and the definition of $\checkmark$. Laws (*let-flatten*) and (*let-let*) allow bindings to move across each other, and law (*let-float-val*) concerns the movement of value bindings across $\lambda$s (along with (*let-float*) below, it forms the essence of the full-laziness transformation, as noted in [PJPS96]). (*inline*) is simple inlining. The last law, (*value-copy*) says that if we have two copies of a strongly-connected component of the heap (composed solely of values), then we may remove one of them, provided we perform some renaming.

Note that in, for example, the (*let-let*) axiom, the variable convention ensures that the $\vec{z}$ do not occur free in the $\vec{L}$; in (*let-float-val*), the convention guarantees that $x$ is not free in the $\vec{V}$.

All of the let laws except (*value-copy*) follow via similar arguments to that for ($\beta$) above. (*value-copy*) requires the use of the same general techniques needed to justify the more complex $\beta$ laws (proof given in the appendix).

## 4.4 Divergence and Strictness

Let $\Omega$ denote any closed term which does not converge. For example, the "black-hole" term, $\mathsf{let} \; x = x \; \mathsf{in} \; x$, would suffice as a definition for $\Omega$. The first three laws in figure 5 concern $\Omega$ and its relationship with $\underset{\sim}{\gtrless}$. ($\Omega$-$\beta$) and (*let-float*-$\Omega$) are similar to (*value*-$\beta$) and (*let-float-val*) except that $\Omega$ is used in place of a value. All of these laws follow in a straightforward manner from the context lemma and the fact that call-by-name termination behaviour is preserved in call-by-need.

We say that a context $\mathbb{C}$ is *strict* if and only if $\mathbb{C}[\Omega] \cong \Omega$. Given this definition, we can float ticks out of any strict context, as stated by ($\checkmark$-*float*). The proof follows by the same techniques used to prove (*value*-$\beta$).

It turns out that this tick floating property can be used as a characterisation of strictness: for all $\mathbb{C}$, if $\mathbb{C}[\checkmark x] \gtrsim \checkmark \mathbb{C}[x]$, $x$ fresh, then $\mathbb{C}$ is strict. This follows since, by congruence,

$$\mathsf{let} \; x = \Omega \; \mathsf{in} \; \mathbb{C}[\checkmark x] \gtrsim \mathsf{let} \; x = \Omega \; \mathsf{in} \; \checkmark \mathbb{C}[x]$$

which implies, by ($\Omega$-$\beta$), and (*gc*), that $\mathbb{C}[\checkmark \Omega] \gtrsim \checkmark \mathbb{C}[\Omega]$. But since $\checkmark\Omega \underset{\sim}{\not\gtrless} \Omega$ (by ($\Omega$) and (*imp*-$\Omega$)), $\mathbb{C}[\Omega] \gtrsim \checkmark\mathbb{C}[\Omega]$. Therefore, by (*diverge*), $\mathbb{C}[\Omega] \cong \Omega$.

## 4.5 Beta Expansion: A Conjecture

In analogy to (*value*-$\beta$), we have ($\beta$-*expand*) where values are replaced by general terms:

$$\mathsf{let} \; \{x = M, \vec{y} = \vec{\mathbb{D}}[{}^{2\checkmark} M]\} \; \mathsf{in} \; \mathbb{C}[{}^{2\checkmark} M]$$
$$\gtrsim \mathsf{let} \; \{x = M, \vec{y} = \vec{\mathbb{D}}[x]\} \; \mathsf{in} \; \mathbb{C}[x] \qquad (\beta\text{-}expand)$$

The intuition here is that the rule undoes a call-by-name computation step (a beta-reduction). This is an improvement providing we can pay for the potential gain that the computation step might have made — which is at most two ticks at each occurrence of the variable which is unfolded.

Unfortunately we lack a satisfactory proof for ($\beta$-*expand*). The context lemma seems inadequate to establish this property. This seems to be linked to the fact that the axiom embodies the essential difference between call-by-name and call-by-need evaluation, and thus it may be possible to adapt techniques based on redex-marking [MOW98].

The conjecture can also be used to "tie the knot" when deriving cyclic programs. This possible since we allow $x$ to occur free in $M$. See the last step of the proof of proposition 5.4 for an example of the use of ($\beta$-*expand*) in this context.

Using the conjecture, we can also establish the following:

$$\checkmark(\lambda x.\mathsf{let} \; \{\vec{y} = \vec{L}, \vec{z} = \vec{M}\} \; \mathsf{in} \; N)$$
$$\gtrsim \mathsf{let} \; \{\vec{y} = \vec{L}\} \; \mathsf{in} \; \lambda x.\mathsf{let} \; \{\vec{z} = \vec{M}\} \; \mathsf{in} \; N \qquad (let\text{-}float)$$

which concerns moving non-value bindings out of $\lambda$s (where the variable convention ensures that $x$ does not occur free in the $\vec{L}$). As noted above, this is an essential part of the full-laziness transformation. Another consequence of the conjecture is standard common sub-expression elimination:

$$^{\checkmark}\mathbb{C}[^{2\checkmark}M] \mathrel{\substack{\rhd \\ \approx}} \text{let } \{x = M\} \text{ in } \mathbb{C}[x] \qquad (cse)$$

Again, the convention ensures that any free variables of $M$ are not captured by context $\mathbb{C}$.

## 4.6 Constant Factors and the Calculi

We reproduce the axioms of the call-by-need calculus of [AFM$^+$95], in figure 7[3].

The laws collected in figures 2, 3, and 4 subsume the call-by-need lambda calculi (in both cases minus the symmetry law): each calculus rewrite rule of the form $L \rightarrow R$ turns out to be either an outright improvement, *i.e.* $L \mathrel{\rhd\!\!\sim} R$, or, in the case of (let-$A$), an improvement "modulo tick":

$$^{\checkmark}\text{let } y = (\text{let } x = L \text{ in } M) \text{ in } N \mathrel{\rhd\!\!\sim} \text{let } x = L \\ \text{in let } y = M \text{ in } N.$$

This follows by (*let-flatten*), (*let-let*), and ($\checkmark$-*elim*). The extra tick is needed in this case is because without it the right-hand side might (in the case when $y$ is not used) perform one extra heap-allocation step.

What is more interesting is that in each case we can reverse the improvement modulo tick. In other words, there exists an $R'$, obtained from $R$ by inserting ticks, such that $R' \mathrel{\rhd\!\!\sim} L$. This fact will enable us to prove that any two terms related by these calculi compute within a constant factor of each other in any program context. Thus the best (worst) speedup (resp. slowdown) program obtainable in these calculi is linear.

First it is natural to generalise the idea of improvement modulo ticks.

**Definition 4.1 (Imp. within a Constant Factor)** *We say that $M$ is improved by $N$ within a constant factor, written $M \mathrel{\substack{\rhd \\ \approx}} N$, if there exists a $k$ such that for all $\mathbb{C}$ such that $\mathbb{C}[M]$ and $\mathbb{C}[N]$ are closed,*

$$\mathbb{C}[M]\Downarrow^n \implies \mathbb{C}[N]\Downarrow^{\leqslant k(n+1)}.$$

So $M \mathrel{\substack{\rhd \\ \approx}} N$ means that $N$ is never more than a constant factor slower than $M$ (but it might still be faster by a non-constant factor). Note that the constant factor is independent of the context of use.

It can be seen that $\mathrel{\substack{\rhd \\ \approx}}$ is a precongruence relation (to show transitivity requires a small calculation) and clearly contains the improvement relation.

Now we consider a special case of $\mathrel{\substack{\rhd \\ \approx}}$, namely programs which only differ by ticks. Let $M \mathrel{\overset{\checkmark}{\rightarrow}} N$ mean that $N$ can be obtained from $M$ by removing some ticks (from

anywhere within the term), and $M \mathrel{\overset{\checkmark}{\sim}} N$ mean that there exists an $L$ such that $M \mathrel{\overset{\checkmark}{\rightarrow}} L$ and $N \mathrel{\overset{\checkmark}{\rightarrow}} L$. Clearly $\mathrel{\overset{\checkmark}{\rightarrow}}$ is a precongruence and $\mathrel{\overset{\checkmark}{\sim}}$ is a congruence.

**Lemma 4.2** $M \mathrel{\overset{\checkmark}{\sim}} N \implies M \mathrel{\substack{\rhd \\ \approx}} N$.

PROOF. (Sketch) Clearly $\mathrel{\overset{\checkmark}{\rightarrow}} \subseteq \mathrel{\substack{\rhd \\ \approx}}$, so it suffices to show that $M \mathrel{\overset{\checkmark}{\rightarrow}} N \implies N \mathrel{\substack{\rhd \\ \approx}} M$. First show that the nesting of ticks in a configuration never increases as computation proceeds (easy to see since the rules never substitute terms for variables). Then let $k$ be the maximum nesting of ticks in $M$, and show by induction on the length of the computation that $\mathbb{C}[N]\Downarrow^n$ implies $\mathbb{C}[N]\Downarrow^{k(n+1)}$ (strengthening this statement to configurations). $\square$

With this lemma we can establish the following:

**Theorem 4.3** *For all terms $N$ and $M$ (of our restricted syntax) if $M =_{\text{NEED}} N$ then $M \mathrel{\substack{\rhd \\ \approx}} N$.*

PROOF. (Sketch) By induction on the proof of $M =_{\text{NEED}} N$. The base case requires us to show that the (oriented) equations are contained in $\mathrel{\substack{\rhd \\ \approx}}$. This follows easily since they are all either improvements or improvements modulo tick. In the inductive cases, the congruence and transitivity rules follow from the inductive hypothesis since $\mathrel{\substack{\rhd \\ \approx}}$ is a precongruence. The only difficult case is symmetry. It will be sufficient to prove that reversed equations are contained in $\mathrel{\substack{\rhd \\ \approx}}$. For each equation $L =_{\text{NEED}} R$ we have from the laws an $R'$ such that $R' \mathrel{\overset{\checkmark}{\rightarrow}} R$ and $R' \mathrel{\rhd\!\!\sim} L$. By lemma 4.2 we know that $R \mathrel{\substack{\rhd \\ \approx}} R'$, so $R \mathrel{\substack{\rhd \\ \approx}} L$ follows from the fact that $\mathrel{\rhd\!\!\sim} \subseteq \mathrel{\substack{\rhd \\ \approx}}$ and transitivity of $\mathrel{\substack{\rhd \\ \approx}}$. $\square$

**Corollary 4.4** *The call-by-need calculus of [AFM$^+$95] cannot improve (or worsen) a program by more than a constant factor.*

We are confident that this result can be extended to Ariola and Blom's sharing calculus $\lambda_{\text{oSHARE}}$ [AB97] since almost all the rules are represented more or less directly in the collection of improvement laws. It is interesting to note that we assembled our collection of laws "by need", considering what was required to tackle a number of examples, and it was encouraging to find that we had already covered almost all of Ariola and Blom's rules. As it stands however, our (*value-copy*) cost equivalence is not as expressive as Ariola and Blom's value-copy rule.[4] We believe that Ariola and Blom's value-copy rule *is* a cost equivalence, but their formulation of the rule is rather indirect, so it is not obvious to us how to prove this.

## 5 Syntactic Continuity

We wish to say something meaningful about recursive functions with this theory, and a natural starting point is to attempt to mimic the fixed-point induction Scott-style

---

[3]In the original paper $V$ ranges over variables as well as values. In addition, Ariola and Felleisen [AF97] restrict $\mathbb{C}$ in (let-$V$) to be evaluation contexts.

[4]Thanks to Stefan Blom for providing an example, and to Zena Ariola for pointing out an error in the use of an earlier formulation of our value-copy rule.

$$(\lambda x.M)\,N =_{\text{NEED}} \text{let } x = N \text{ in } M \tag{let-$I$}$$

$$\text{let } x = V \text{ in } \mathbb{C}[x] =_{\text{NEED}} \text{let } x = V \text{ in } \mathbb{C}[V] \tag{let-$V$}$$

$$(\text{let } x = L \text{ in } M)\,N =_{\text{NEED}} \text{let } x = L \text{ in } M\,N \tag{let-$C$}$$

$$\text{let } y = (\text{let } x = L \text{ in } M) \text{ in } N =_{\text{NEED}} \text{let } x = L \text{ in let } y = M \text{ in } N \tag{let-$A$}$$

Figure 7: Axioms of the call-by-need calculus [AFM$^{+}$95].

denotational semantics. Examples of this kind of operational analogue to Scott induction for other languages may be found in *e.g.*, [Pit97b, Smi91, MST96, San97, Las98]; we present the first such result for a call-by-need semantics.

We will use the following mechanism to describe the syntactic unwindings of a recursive function. In the definition, the $f_i$ are distinct, new variables.

**Definition 5.1** $f \overset{0}{=} V \overset{\text{def}}{=} f_0 = \Omega$,
$$f \overset{n+1}{=} V \overset{\text{def}}{=} f \overset{n}{=} V, f_{n+1} = V[f_n/f].$$

Then, for an $f$ defined by let $\{f = V\}$ in $f$, we define the $n^{\text{th}}$ unwinding as let $\{f \overset{n}{=} V\}$ in $f_n$. If we expand the definition of $f \overset{n}{=} V$, we see that this is really

$$\text{let } \{f_0 = \Omega, f_1 = V[f_0/f], \ldots, f_n = V[f_{n-1}/f]\}$$
$$\text{in } f_n.$$

Note that we have restricted our attention to $f$ whose defining body is a value; this unwinding trick would not work for general cycles (since loss of sharing would render the exercise pointless). To extend the method to cycles would require some extension to the language, but this would lead to the problem of showing that the extension is conservative with respect to the improvement relation.

The point is that the functions let $\{f \overset{n}{=} V\}$ in $f_n$ completely characterise the behaviour of let $\{f = V\}$ in $f$. This is the essence of Scott induction. The main property that justifies this is a syntactic notion of continuity, which says that let $\{f = V\}$ in $f$ is the least upper bound of chain $\{$let $\{f \overset{n}{=} V\}$ in $f_n\}_{n \geqslant 0}$ and that any $M$ which uses $f$ preserves this property.

We first show that $\{$let $\{f \overset{n}{=} V\}$ in $M[f_n/f]\}_{n \geqslant 0}$ does indeed form a chain with respect to $\gtrsim$, and that let $\{f = V\}$ in $M$ is an upper bound of that chain.

**Lemma 5.1** $\forall n.$ let $\{f \overset{n}{=} V\}$ in $M[f_n/f]$
$$\gtrsim \text{let } \{f \overset{n+1}{=} V\} \text{ in } M[f_{n+1}/f]$$
$$\gtrsim \text{let } \{f = V\} \text{ in } M.$$

PROOF. We prove only the second improvement, that for all $n$, let $\{f \overset{n}{=} V\}$ in $M[f_n/f] \gtrsim$ let $\{f = V\}$ in $M$. The first follows by a similar argument. We proceed by induction on $n$. The base case follows easily by $(gc)$ and the $\Omega$ laws, and the inductive case follows by $(\checkmark\text{-}elim)$, since

$$\checkmark \text{let } \{f \overset{n}{=} V\} \text{ in } M[f_n/f] \gtrsim \checkmark \text{let } \{f = V\} \text{ in } M$$

by the derivation in figure 8. $\qquad\square$

To establish syntactic continuity, we will need the following lemma (proof given in the appendix). It says that if let $\{f = V\}$ in $M$ converges then there must exist some unwinding that does so with the same cost.

**Lemma 5.2 (Unwinding)** *For all* $\Gamma, S$, *and* $n$,

$$\langle \Gamma, \text{ let } \{f = V\} \text{ in } M, \ S \rangle\Downarrow^n \implies$$
$$\exists m.\langle \Gamma, \text{ let } \{f \overset{m}{=} V\} \text{ in } M[f_m/f], \ S \rangle\Downarrow^n.$$

**Theorem 5.3 (Syntactic Continuity)** *The following is a sound proof rule:*

$$\frac{\forall n.\text{let } \{f \overset{n}{=} V\} \text{ in } M[f_n/f] \ \gtrsim \ N}{\text{let } \{f = V\} \text{ in } M \ \gtrsim \ N}$$

PROOF. Assume $\langle \Gamma, \text{ let } \{f = V\} \text{ in } M, \ S \rangle\Downarrow^n$. Then by the Unwinding lemma, there exists some $m$ such that $\langle \Gamma, \text{ let } \{f \overset{m}{=} V\} \text{ in } M[f_m/f], \ S \rangle\Downarrow^n$. By the premise, we have that $\langle \Gamma, \ N, \ S \rangle\Downarrow^{\leqslant n}$, and the result follows by the context lemma. $\qquad\square$

Syntactic continuity is also valid for mutually recursive functions. This proof rule is sound for strong improvement, but note that the base case of the premise requires that $N$ be contextually equivalent to $\Omega$. This tends to limit the applicability of the strong improvement version of syntactic continuity.

As an example of the use of syntactic continuity, we show that an unwinding fixed-point combinator is improved within a constant factor by a "knot-tying" fixed-point combinator.

**Proposition 5.4** *If* $(\beta\text{-}expand)$ *is valid, then*

$$\text{let } rec = (\lambda f.\text{let } x = rec\ f \text{ in } f\ x) \text{ in } rec$$
$$\underset{\approx}{\gtrsim} \text{let } fix = (\lambda f.\text{let } x = f\ x \text{ in } x) \text{ in } fix.$$

PROOF. Let $V = \lambda f.\text{let } x = \checkmark rec\ f \text{ in } {}^{2\checkmark} f\ x$, and abbreviate $V[rec_n/rec]$ by $V_n$. We will show that for all $n$, let $rec \overset{n}{=} V$ in $rec_n \gtrsim {}^{3\checkmark} \lambda f.\text{let } x = f\ x \text{ in } x$. Then the result will then follow by syntactic continuity, since

$${}^{3\checkmark} \lambda f.\text{let } x = f\ x \text{ in } x$$
$$\underset{\approx}{\gtrsim} {}^{2\checkmark} \text{let } fix = (\lambda f.\text{let } x = f\ x \text{ in } x) \quad (gc)$$
$$\text{in } \lambda f.\text{let } x = f\ x \text{ in } x$$
$$\underset{\approx}{\gtrsim} \text{let } fix = (\lambda f.\text{let } x = f\ x \text{ in } x) \text{ in } fix \quad (value\text{-}\beta)$$

We proceed via induction on $n$. The base case follows trivially by $(imp\text{-}\Omega)$ and $(\Omega)$ since let $rec_0 = \Omega$ in $rec_0 \cong \Omega$, and the inductive case follows by the derivation in figure 9. We have $\underset{\approx}{\gtrsim}$ and not $\gtrsim$ because we use a slightly slower version of $rec$. $\qquad\square$

$$\checkmark \text{let } \{f \stackrel{n}{=} V, f_{n+1} = V[f_n/f]\} \text{ in } M[f_{n+1}/f]$$

$$\gtrsim \text{let } \{f \stackrel{n}{=} V\} \text{ in let } \{f_{n+1} = V[f_n/f]\} \text{ in } M[f_{n+1}/f] \qquad (\textit{let-let})$$

$$\equiv \text{let } \{f \stackrel{n}{=} V\} \text{ in let } \{g = V[f_n/f]\} \text{ in } M[g/f] \qquad (\text{rename})$$

$$\gtrsim \text{let } \{f = V\} \text{ in let } \{g = V[f/f]\} \text{ in } M[g/f] \qquad (\text{I.H.})$$

$$\gtrsim \checkmark \text{let } \{f = V, g = V\} \text{ in } M[g/f] \qquad (\textit{let-let})$$

$$\gtrsim \checkmark \text{let } \{f = V\} \text{ in } M \qquad (\textit{value-copy}), (\textit{gc})$$

Figure 8: The inductive case for lemma 5.1.

$$\text{let } rec \stackrel{n}{=} V, rec_{n+1} = V_n \text{ in } rec_{n+1}$$

$$\gtrsim \text{let } rec \stackrel{n}{=} V, rec_{n+1} = V_n \text{ in } {}^{2\checkmark}V_n \qquad (\textit{value-}\beta)$$

$$\gtrsim \text{let } rec \stackrel{n}{=} V \text{ in } {}^{2\checkmark}\lambda f.\text{let } x = {}^{\checkmark}rec_n \, f \text{ in } {}^{2\checkmark}f \, x \qquad (\textit{gc}), (\text{defn. of } V_n)$$

$$\gtrsim {}^{3\checkmark}\lambda f.\text{let } rec \stackrel{n}{=} V, x = {}^{\checkmark}rec_n \, f \text{ in } {}^{2\checkmark}f \, x \qquad (\textit{let-float-val}), (\textit{let-float-}\Omega)$$

$$\gtrsim {}^{3\checkmark}\lambda f.\text{let } x = (\text{let } rec \stackrel{n}{=} V \text{ in } rec_n) \, f \text{ in } {}^{2\checkmark}f \, x \qquad (\textit{let-let}), (\textit{let-}\mathbb{E})$$

$$\gtrsim {}^{3\checkmark}\lambda f.\text{let } x = ({}^{3\checkmark}\lambda g.\text{let } y = g \, y \text{ in } y) \, f \text{ in } {}^{2\checkmark}f \, x \qquad (\text{I.H.}), (\text{rename})$$

$$\gtrsim {}^{3\checkmark}\lambda f.\text{let } x = {}^{5\checkmark}\text{let } y = f \, y \text{ in } y \text{ in } {}^{2\checkmark}f \, x \qquad (\beta)$$

$$\gtrsim {}^{3\checkmark}\lambda f.\text{let } x = {}^{6\checkmark}y, y = f \, y \text{ in } {}^{2\checkmark}f \, x \qquad (\textit{let-let})$$

$$\gtrsim {}^{3\checkmark}\lambda f.\text{let } x = y, y = f \, y \text{ in } {}^{2\checkmark}f \, y \qquad (\checkmark\text{-}elim), (\textit{var-subst})$$

$$\gtrsim {}^{3\checkmark}\lambda f.\text{let } x = f \, x \text{ in } {}^{2\checkmark}f \, x \qquad (\textit{gc}), (\text{rename})$$

$$\gtrsim {}^{3\checkmark}\lambda f.\text{let } x = f \, x \text{ in } x \qquad (\beta\text{-}expand)$$

Figure 9: The inductive case for proposition 5.4.

The converse of the proposition is false, since the knot-tying fixed-point combinator can give asymptotically better programs.

We can also use syntactic continuity to establish the following proof rule, which is a syntactic, call-by-need version of what is called *fixed-point fusion* in [MFP91]. In the statement, $\mathbb{V}$ and $\mathbb{W}$ range over value contexts.

**Theorem 5.5 (Improvement Fusion)** *If $\mathbb{C}$ is strict, and $\mathbb{C}[\mathbb{V}[x]] \gtrsim \mathbb{W}[\mathbb{C}[x]]$ where $x \notin \mathsf{FV}(\mathbb{V}, \mathbb{W}, \mathbb{C}) \cup \mathsf{CV}(\mathbb{V}, \mathbb{W}, \mathbb{C})$, then for all $\mathbb{D}$ such that $x \notin \mathsf{FV}(\mathbb{D}) \cup \mathsf{CV}(\mathbb{D})$,*

$$\text{let } \{x = \mathbb{V}[x]\} \text{ in } \mathbb{D}[\mathbb{C}[x]] \gtrsim \text{let } \{x = \mathbb{W}[x]\} \text{ in } \mathbb{D}[x].$$

PROOF. Assume $\mathbb{C}$ is strict, and that $\mathbb{C}[\mathbb{V}[x]] \gtrsim \mathbb{W}[\mathbb{C}[x]]$. By syntactic continuity, it suffices to show, for all $n$ and all $\mathbb{D}$ such that $x \notin \mathsf{FV}(\mathbb{D}) \cup \mathsf{CV}(\mathbb{D})$,

$$\text{let } \{x \stackrel{n}{=} \mathbb{V}[x]\} \text{ in } \mathbb{D}[\mathbb{C}[x_n]] \gtrsim \text{let } \{x = \mathbb{W}[x]\} \text{ in } \mathbb{D}[x].$$

The base case follows by this calculation:

$$\text{let } \{x_0 = \Omega\} \text{ in } \mathbb{D}[\mathbb{C}[x_0]]$$

$$\gtrsim \text{let } \{x_0 = \Omega\} \text{ in } \mathbb{D}[\mathbb{C}[\Omega]] \qquad (\Omega\text{-}\beta)$$

$$\gtrsim \text{let } \{x_0 = \Omega\} \text{ in } \mathbb{D}[\Omega] \qquad (\mathbb{C} \text{ strict})$$

$$\gtrsim \checkmark \mathbb{D}[\Omega] \qquad (\textit{gc})$$

$$\gtrsim \text{let } \{x = \mathbb{W}[x]\} \text{ in } \mathbb{D}[\Omega] \qquad (\textit{gc})$$

$$\gtrsim \text{let } \{x = \mathbb{W}[x]\} \text{ in } \mathbb{D}[x] \qquad (\Omega \gtrsim x), (\text{cong.})$$

and for the inductive case:

$$\text{let } \{x \stackrel{n+1}{=} \mathbb{V}[x]\} \text{ in } \mathbb{D}[\mathbb{C}[x_{n+1}]]$$

$$\gtrsim \text{let } \{x \stackrel{n}{=} \mathbb{V}[x]\} \text{ in } \mathbb{D}[\mathbb{C}[{}^{2\checkmark}\mathbb{V}[x_n]]] \qquad (\textit{value-}\beta), (\textit{gc})$$

$$\gtrsim \text{let } \{x \stackrel{n}{=} \mathbb{V}[x]\} \text{ in } \mathbb{D}[{}^{2\checkmark}\mathbb{C}[\mathbb{V}[x_n]]] \qquad (\mathbb{C} \text{ strict})$$

$$\gtrsim \text{let } \{x \stackrel{n}{=} \mathbb{V}[x]\} \text{ in } \mathbb{D}[{}^{2\checkmark}\mathbb{W}[\mathbb{C}[x_n]]] \qquad (\text{assumption})$$

$$\gtrsim \text{let } \{x = \mathbb{W}[x]\} \text{ in } \mathbb{D}[{}^{2\checkmark}\mathbb{W}[x]] \qquad (\text{I.H.})$$

$$\gtrsim \text{let } \{x = \mathbb{W}[x]\} \text{ in } \mathbb{D}[x] \qquad (\textit{value-}\beta)$$

$\square$

Fixed-point fusion can be used to establish a number of general fusion laws. It is also central to Tullsen and Hudak's [TH98] approach to program transformation in Haskell.

## 6 The Improvement Theorem

In this section we introduce a second key technique for reasoning about recursion, the improvement theorem. In [San96] a call-by-name improvement theorem was introduced as a means to prove the extensional correctness of recursion-based program transformations. In this section we show how these results carry over to the call-by-need setting.

## 6.1 The Problem of Transformations

As a motivation for the improvement theorem, consider the correctness problem for recursion-based program transformations such as *unfold-fold*; the correctness of such transformations does not follow from the simple fact that the basic transformation steps are equivalences. To take a simple example to illustrate the problem, consider the following "transformation by equivalence-preserving steps". Start with the recursive function *repeat* which produces the "infinite" list of its argument:

$$repeat\ x = x : (repeat\ x)$$

The following property can be easily deduced: $repeat\ x \cong \mathsf{tail}(repeat\ x)$. Now suppose that we use this "local equivalence" to transform the body of the function to obtain a new version of the function:

$$repeat\ x = x : (\mathsf{tail}\,(repeat\ x))$$

This definition is not equivalent to the original, since it can never produce more than first element in the list. How did equivalence-preserving local steps produce a non-equivalent function? Analysing such transformations more carefully we see that while it is true that

$$M \cong N \implies \mathsf{let}\ \{x = M\}\ \mathsf{in}\ L \cong \mathsf{let}\ \{x = N\}\ \mathsf{in}\ L \tag{6.1}$$

it is no longer the case when the transformation from $M$ to $N$ depends on the recursive definition of $x$ itself:

$$\mathsf{let}\ \{x = M\}\ \mathsf{in}\ M \cong \mathsf{let}\ \{x = M\}\ \mathsf{in}\ N$$
$$\not\!\!\implies \mathsf{let}\ \{x = M\}\ \mathsf{in}\ L \cong \mathsf{let}\ \{x = N\}\ \mathsf{in}\ L.$$

But in order to reason about "interesting" program transformations (*e.g.* unfold-fold, recursion-based deforestation, partial evaluation with memoization), inference (6.1) is simply not sufficient.

The improvement theorem comes to the rescue:

$$\frac{\mathsf{let}\ \{x = M\}\ \mathsf{in}\ M \gtrsim \mathsf{let}\ \{x = M\}\ \mathsf{in}\ N}{\mathsf{let}\ \{x = M\}\ \mathsf{in}\ L \gtrsim \mathsf{let}\ \{x = N\}\ \mathsf{in}\ L} \tag{6.2}$$

This is sufficient to establish the correctness of recursion-based transformations by requiring — rather naturally — that the local transformation steps are also improvements. This was proved for an improvement theory based on call-by-name, so the fact that the theorem gives "improved" programs as well as correctness is not considered to be particularly significant.

A question left open was whether the improvement theorem holds for a call-by-need improvement theory. We can now supply the answer:

**Theorem 6.1 (Improvement Theorem)** *The following proof rule is sound:*

$$\frac{\mathsf{let}\ \{f = V\}\ \mathsf{in}\ V \gtrsim \mathsf{let}\ \{f = V\}\ \mathsf{in}\ W}{\mathsf{let}\ \{f = V\}\ \mathsf{in}\ N \gtrsim \mathsf{let}\ \{f = W\}\ \mathsf{in}\ N}$$

*The inference is also sound when $\gtrsim$ is replaced throughout with $\stackrel{\scriptscriptstyle\triangleleft}{\scriptscriptstyle\triangleright}$ (the cost equivalence theorem).*

The improvement theorem and the cost equivalence theorem can also be stated for a set of mutually recursive definitions. The proof of the theorem is appendix.

**Notation** In establishing a premise of the improvement theorem, in the context of some recursive declarations $\vec{g} = \vec{V}$, a derivation of the form

$$\mathsf{let}\ \{\vec{g} = \vec{V}\}\ \mathsf{in}\ M_1 \gtrsim \mathsf{let}\ \{\vec{g} = \vec{V}\}\ \mathsf{in}\ M_2$$
$$\gtrsim \mathsf{let}\ \{\vec{g} = \vec{V}\}\ \mathsf{in}\ M_3 \ldots$$

will be written in the following abbreviated form:

$$\vec{g} \vdash\ M_1 \gtrsim M_2$$
$$\gtrsim M_3 \ldots$$

when the declarations $\vec{g}$ are clear from the context.

The following example illustrates the use of the proof rule, which shows that a representation of the standard lambda-calculus fixed-point combinator $Y = \lambda f.f\,((\lambda x.f\,(x\,x))\,\lambda x.f\,(x\,x))$ (suitably converted to the restricted syntax) is cost equivalent to the non-cyclic version *rec* from proposition 5.4.

**Proposition 6.2**

$$\mathsf{let}\ Y = \lambda f.\mathsf{let}\ d = \lambda y.\mathsf{let}\ z = y\,y\ \mathsf{in}\ f\,z$$
$$x = d\,d$$
$$\mathsf{in}\ f\,x$$
$$\mathsf{in}\ Y$$
$$\stackrel{\scriptscriptstyle\triangleleft}{\scriptscriptstyle\triangleright}\ \mathsf{let}\ rec = (\lambda f.\mathsf{let}\ x = rec\,f\ \mathsf{in}\ f\,x)\ \mathsf{in}\ rec.$$

PROOF. By a simple 6-step derivation, we have that

$$Y \vdash \lambda f.\mathsf{let}\ d = \lambda y.\mathsf{let}\ z = y\,y\ \mathsf{in}\ f\,z$$
$$x = d\,d$$
$$\mathsf{in}\ f\,x$$
$$\stackrel{\scriptscriptstyle\triangleleft}{\scriptscriptstyle\triangleright}\ \lambda f.\mathsf{let}\ x = Y\,f\ \mathsf{in}\ f\,x$$

Then the result follows by the Cost Equivalence Theorem. □

**Improvement Theorem vs. Syntactic Continuity** Suppose one wants to establish an improvement of the form

$$\mathsf{let}\ \{f = V\}\ \mathsf{in}\ N \gtrsim \mathsf{let}\ \{f = W\}\ \mathsf{in}\ N.$$

If the left-hand side is non-recursive (in $f$) then syntactic continuity is of no help, since the unwindings ($> 0$) of the left-hand side will all be identical; conversely, if the right-hand side is non recursive (in $f$) then the improvement theorem is not immediately useful, since proving the premise amounts to directly proving the conclusion of the rule. There are, however, many examples which can be proved by both methods. In these cases the improvement theorem is often preferable since it is more calculational in style.

## 6.2 Improvement Induction

Finally, we mention one last proof rule which is closely allied to the improvement theorem (in the sense that a closely-related rule can be derived from the improvement theorem); this corresponds to what we called *improvement induction* in [San97], where it was established for any call-by-name or call-by-value language with SOS rules fitting a certain syntactic rule-format.

**Theorem 6.3 (Improvement Induction)** *For any $M$, $N$, $\mathbb{C}$, and substitution $\sigma$, the following proof rule is sound:*

$$\frac{M \mathrel{\underset{\sim}{\rhd}} {}^{\checkmark}\mathbb{C}[M\sigma] \quad N \mathrel{\underset{\sim}{\lessgtr}} {}^{\checkmark}\mathbb{C}[N\sigma]}{M \mathrel{\underset{\sim}{\rhd}} N}$$

The proof is quite straightforward, and is given in the appendix.

Let us take a standard example to illustrate the proof technique: the associativity of append (the list-concatenation function). In order to show that it is an improvement, we need to insert a tick into the recursive branch; this is a consequence of the fact that our cost-measure is rather fine-grained. To ease the notation, we will make use of the syntactic identity for general application from section 3, and we will use an infix form of append.

Given the recursive declaration

$$(+\!\!+) = \lambda xs.\lambda ys.\text{case } xs \text{ of}$$
$$\text{nil} \quad \to ys$$
$$h : t \to h : {}^{\checkmark}(t +\!\!+ ys),$$

then $(+\!\!+) \vdash (x +\!\!+ y) +\!\!+ z \mathrel{\underset{\sim}{\rhd}} x +\!\!+ (y +\!\!+ z)$. To show this, by improvement induction it is sufficient to find a context $\mathbb{C}$ and substitution $\sigma$ such that

$$(+\!\!+) \vdash (x +\!\!+ y) +\!\!+ z \mathrel{\underset{\sim}{\rhd}} {}^{\checkmark}\mathbb{C}[((x +\!\!+ y) +\!\!+ z)\sigma],$$
$$(+\!\!+) \vdash x +\!\!+ (y +\!\!+ z) \mathrel{\underset{\sim}{\lessgtr}} {}^{\checkmark}\mathbb{C}[(x +\!\!+ (y +\!\!+ z))\sigma].$$

The solution is to take $\mathbb{C}$ to be

$${}^{5\checkmark}\text{case } x \text{ of}$$
$$\text{nil} \quad \to {}^{\checkmark}(y +\!\!+ z)$$
$$h : t \to \text{let } r = [\cdot] \text{ in } h : r$$

and $\sigma = [{}^{t}/{}_{x}]$. This context is easily derived by performing a cost equivalence calculation on the right-hand side until a recurring instance of $x +\!\!+ (y +\!\!+ z)$ is discovered. We omit the derivations.

By using suitably slowed versions of append on the right-hand side, we can also show that this is only a linear speedup in all contexts, *i.e.* that

$$(+\!\!+) \vdash x +\!\!+ (y +\!\!+ z) \mathrel{\underset{\approx}{\rhd}} (x +\!\!+ y) +\!\!+ z.$$

Using the above properties we have been able to prove, with the help of the improvement theorem and an adaptation of the argument given in [San96], that a mechanisable append-elimination transformation described in [Wad88] can never slow-down programs by more than a constant factor. What is interesting about this example is that it shows that the improvement theorem can obtain asymptotic speedups using only linear ones, since in particular the transformation is able to turn the naïve quadratic-time definition of reverse into a linear-time version.

## 7  Conclusions and Future Work

We have presented a rich operational theory for a call-by-need based on an improvement ordering on programs. The theory subsumes the (oriented) call-by-need lambda

calculi of Ariola *et al.* [AFM$^+$95]. The most important extensions are proof techniques for reasoning about recursion. Syntactic continuity allows us to prove properties of recursive programs via a kind of fixed-point induction, without sacrificing information about intensional behaviour, like sharing. The improvement theorem and improvement induction are rules for recursion which support more calculational proofs. Both are particularly useful in proving the safety of program transformations.

Our unit of cost is the abstract machine reduction step. This choice simplifies the technical development. A drawback is that it is very fine-grained, so one must carefully track costs of optimisation steps. This bookkeeping can, in larger examples, become rather tedious. However, it should be possible to follow our programme with a much more abstract cost measure. For example, one possibility would be to count only the number of lookups[5]. This should significantly simplify the tick algebra, without compromising the ability to characterise relative efficiency within a constant factor.

An obvious further application of the theory is to formalise arguments about the running time of programs, following Sands' use of call-by-name cost equivalence for this purpose [San95, San98b].

Another direction for future work would be to consider the time-safety of a larger-scale program transformation, such as deforestation [Wad90]. In such a transformation we must inevitably consider conditions under which we can unfold function calls. It is straightforward to define simple syntactic conditions on contexts which guarantee that

$$\text{let } \{\vec{x} = \vec{M}\} \text{ in } \mathbb{C}[\vec{x}] \mathrel{\underset{\sim}{\rhd}} \text{let } \{\vec{x} = \vec{M}\} \text{ in } \mathbb{C}[\vec{M}],$$

but in the case where holes occur under $\lambda$-abstractions a more global form of information is required: one needs to know that the lambda expression in question will not be applied more than once. The type system of [TWM95] provides just such global information, so it would be interesting to prove that their system (and generalisations to full recursive lets [Gus98]) does indeed satisfy the desired improvement property above. We saw in section 4.4 that the strictness property of a context can be characterised exactly by

$$\mathbb{C}[{}^{\checkmark}x] \mathrel{\underset{\sim}{\rhd}} {}^{\checkmark}\mathbb{C}[x],$$

where $x$ is fresh. Could it be the case that the "used at most once" property might be semantically characterised by ${}^{\checkmark}\mathbb{C}[x] \mathrel{\underset{\sim}{\rhd}} \mathbb{C}[{}^{\checkmark}x]$?

---

[5]Of course, one would need to change the definition of $\checkmark$ accordingly.

## References

[AB97]     Z. M. Ariola and S. Blom, *Cyclic lambda calculi*, Proc. TACS'97, LNCS, vol. 1281, Springer-Verlag, February 1997, pp. 77–106.

[AB98]     Z. M. Ariola and S. Blom, *Lambda calculi plus letrec*, Tech. report, Dept. of Computer Science, University of Oregon, 1998, Extended version of [AB97]; submitted for publication.

[AF97]     Z. M. Ariola and M. Felleisen, *The call-by-need lambda calculus*, Journal of Functional Programming **7** (1997), no. 3, 265–301.

[AFM$^+$95] Z. Ariola, M. Felleisen, J. Maraist, M. Odersky, and P. Wadler, *A call-by-need lambda calculus*, Proc. POPL'95, the $22^{nd}$ ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, ACM Press, January 1995, pp. 233–246.

[AK97]     Z. M. Ariola and J. W. Klop, *Lambda calculus with explicit recursion*, Information and Computation **139** (1997), no. 2, 154–233.

[AMST97]   G. A. Agha, I. A. Mason, S. F. Smith, and C. L. Talcott, *A foundation for actor computation*, Journal of Functional Programming **7** (1997), 1–72.

[AO93]     S. Abramsky and C.-H. L. Ong, *Full abstraction in the lazy lambda calculus*, Information and Computation **105** (1993), 159–267.

[BLR96]    Z.-E.-A. Benaissa, P. Lescanne, and K. H. Rose, *Modeling sharing and recursion for weak reduction strategies using explicit substitution*, Proc. PLILP'96, the $8^{th}$ International Symposium on Programming Languages, Implementations, Logics, and Programs, LNCS, vol. 1140, Springer-Verlag, 1996, pp. 393–407.

[Cur91]    P.-L. Curien, *An abstract framework for environment machines*, Theoretical Computer Science **82** (1991), no. 2, 389–402.

[Fie90]    J. Field, *On laziness and optimality in lambda interpreters: Tools for specification and analysis*, Proc. POPL'90, the $17^{th}$ ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, ACM Press, January 1990, pp. 1–15.

[GP98]     A. D. Gordon and A. M. Pitts (eds.), *Higher Order Operational Techniques in Semantics*, Publications of the Newton Institute, Cambridge University Press, 1998.

[Gus98]    J. Gustavsson, *A type based sharing analysis for update avoidance and optimisation*, Proc. ICFP'98, the $3^{rd}$ ACM SIGPLAN International Conference on Functional Programming, September 1998, pp. 39–50.

[HM95]     J. Hughes and A. K. Moran, *Making choices lazily*, Proc. FPCA'95, ACM Conference on Functional Programming Languages and Computer Architecture, ACM Press, June 1995, pp. 108–119.

[Jef93]    A. Jeffrey, *A fully abstract semantics for concurrent graph reduction*, Tech. Report 93:12, School of Cognitive and Computing Sciences, University of Sussex, 1993.

[Jef94]    A. Jeffrey, *A fully abstract semantics for concurrent graph reduction*, Proc. LICS'94, the $9^{th}$ IEEE Symposium on Logic in Computer Science, IEEE Computer Society Press, July 1994, pp. 82–91.

[Jos89]    M. B. Josephs, *The semantics of lazy functional languages*, Theoretical Computer Science **68** (1989), no. 1, 105–111.

[Las98]    S. B. Lassen, *Relational reasoning about functions and nondeterminism*, Ph.D. thesis, Department of Computer Science, University of Aarhus, May 1998.

[Lau93]    J. Launchbury, *A natural semantics for lazy evaluation*, Proc. POPL'93, the $20^{th}$ ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, ACM Press, January 1993, pp. 144–154.

[Mar91]    L. Maranget, *Optimal derivations in weak lambda-calculi and in orthogonal term rewriting systems*, Proc. POPL'91, the $18^{th}$ ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, ACM Press, January 1991, pp. 255–269.

[MFP91]    E. Meijer, M. Fokkinga, and R. Paterson, *Functional programming with bananas, lenses, envelopes and barbed wire*, Proc. FPCA'91, ACM Conference on Functional Programming Languages and Computer Architecture (J. Hughes, ed.), LNCS, vol. 523, Springer-Verlag, August 1991, pp. 124–144.

[Mil77]    R. Milner, *Fully abstract models of the typed $\lambda$-calculus*, Theoretical Computer Science **4** (1977), 1–22.

[Mor98]    A. K. Moran, *Call-by-name, call-by-need, and McCarthy's Amb*, Ph.D. thesis, Department of Computing Sciences, Chalmers University of Technology, Sweden, September 1998.

[MOW98]    J. Maraist, M. Odersky, and P. Wadler, *The call by need lambda calculus*, Journal of Functional Programming **8** (1998), no. 3, 275–317.

[MST96]    I. A. Mason, S. F. Smith, and C. L. Talcott, *From operational semantics to domain theory*, Information and Computation **128** (1996), no. 1, 26–47.

[MT91]     I. Mason and C. Talcott, *Equivalence in functional languages with effects*, Journal of Functional Programming **1** (1991), no. 3, 287–327.

[Nie96]    J. Niehren, *Functional computation as concurrent computation*, Proc. POPL'96, the $23^{rd}$ ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, ACM Press, January 1996, pp. 333–343.

[Pit94]     A. M. Pitts, *Some notes on inductive and co-inductive techniques in the semantics of functional programs*, Notes Series BRICS-NS-94-5, BRICS, Department of Computer Science, University of Aarhus, December 1994.

[Pit97a]    A. M. Pitts, *Operational semantics for program equivalence*, March 1997, Invited talk at MFPS XIII, the 13$^{th}$ Conference on Mathematical Foundations of Programming Semantics, slides available at http://www.cl.cam.ac.uk/users/ap/talks/mfps97.ps.gz.

[Pit97b]    A. M. Pitts, *Operationally-based theories of program equivalence*, Semantics and Logics of Computation (P. Dybjer and A. M. Pitts, eds.), Publications of the Newton Institute, Cambridge University Press, 1997, pp. 241–298.

[PJPS96]    S. Peyton Jones, W. Partain, and A. Santos, *Let-floating: moving bindings to give faster programs*, Proc. ICFP'96, the 1$^{st}$ ACM SIGPLAN International Conference on Functional Programming, ACM Press, May 1996, pp. 1–12.

[PJS98]     S. Peyton Jones and A. Santos, *A transformation-based optimiser for Haskell*, Science of Computer Programming **32** (1998), no. 1–3, 3–47.

[Ros96]     K. H. Rose, *Operational reduction models for functional programming languages*, Ph.D. thesis, DIKU, University of Copenhagen, Denmark, February 1996, available as DIKU report 96/1.

[San91]     D. Sands, *Operational theories of improvement in functional languages (extended abstract)*, Proc. 1991 Glasgow Functional Programming Workshop, Workshops in Computing Series, Springer-Verlag, August 1991, pp. 298–311.

[San95]     D. Sands, *A naïve time analysis and its theory of cost equivalence*, Journal of Logic and Computation **5** (1995), no. 4, 495–541.

[San96]     D. Sands, *Total correctness by local improvement in the transformation of functional program*, ACM Transactions on Programming Languages and Systems (TOPLAS) **18** (1996), no. 2, 175–234.

[San97]     D. Sands, *From SOS rules to proof principles: An operational metatheory for functional languages*, Proc. POPL'97, the 24$^{th}$ ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, ACM Press, January 1997.

[San98a]    D. Sands, *Computing with contexts: A simple approach*, Proc. HOOTS II, the 2$^{nd}$ Workshop on Higher Order Operational Techniques in Semantics (A. D. Gordon, A. M. Pitts, and C. L. Talcott, eds.), Electronic Notes in Theoretical Computer Science, vol. 10, Elsevier Science Publishers B.V., 1998, at http://www.elsevier.nl/cas/tree/store/tcs/free/noncas/pc/menu.htm.

[San98b]    D. Sands, *Improvement theory and its applications*, In Gordon and Pitts [GP98], pp. 275–306.

[Ses97]     P. Sestoft, *Deriving a lazy abstract machine*, Journal of Functional Programming **7** (1997), no. 3, 231–264.

[Smi91]     S. F. Smith, *From operational to denotational semantics*, Proc. MFPS VII, the 7$^{th}$ Conference on Mathematical Foundations of Programming Semantics (S. Brookes, M. Main, A. Melton, M. Mislove, and D. Schmidt, eds.), LNCS, vol. 598, Springer-Verlag, March 1991, pp. 54–76.

[SPI96]     J. Seaman and S. Purushothaman Iyer, *An operational semantics of sharing in lazy evaluation*, Science of Computer Programming **27** (1996), no. 3, 289–322.

[SPJ97]     P. Sansom and S. Peyton Jones, *Formally-based profiling for higher-order functional languages*, ACM Transactions on Programming Languages and Systems (TOPLAS) **19** (1997), no. 1, 334–385.

[Tal98]     C. L. Talcott, *Reasoning about functions with effects*, In Gordon and Pitts [GP98], pp. 347–390.

[TH98]      M. Tullsen and P. Hudak, *An intermediate meta-language for program transformation*, YALEU/DCS/RR 1154, Yale University, June 1998.

[TWM95]     D. N. Turner, P. Wadler, and C. Mossin, *Once upon a type*, Proc. FPCA'95, ACM Conference on Functional Programming Languages and Computer Architecture, ACM Press, June 1995, pp. 1–11.

[Wad88]     P. Wadler, *The concatenate vanishes*, Tech. report, University of Glasgow (UK), 1988, appeared as a note on an FP electronic mailing list, December 1987.

[Wad90]     P. Wadler, *Deforestation: Transforming programs to eliminate trees*, Theoretical Computer Science **73** (1990), 231–248.

[Yos93]     N. Yoshida, *Optimal reduction in weak-lambda-calculus with shared environments*, Proc. FPCA'93, ACM Conference on Functional Programming Languages and Computer Architecture, ACM Press, June 1993, pp. 243–254.

## A Proofs of Main Theorems

This appendix gives an outline of the technical development and proofs of the main results. Most proofs follow a direct style reasoning which is reminiscent of proofs about functional languages with effects by Mason and Talcott *et al.* [MT91, AMST97, Tal98]. In order to make this style of proof rigorous we generalise the abstract machine semantics so that it works on *configuration contexts* — configurations with holes. To ensure that transitions on configuration contexts are consistent with hole filling one must work with a more general representation of contexts. One such approach is described in [Tal98]. We use an alternative approach to generalising contexts which is due to Pitts [Pit94].

### A.1 Substituting Contexts

Following Pitts [Pit94], we use second-order syntax to represent (and generalise) the traditional definition of contexts given in section 3.2. We give a fuller description in [San98a]; other examples of their use are to be found in [Las98, Mor98]. The idea is that instead of holes [·] we use *second-order variables*, ranged over by $\xi$, applied to some vector of variables. The syntax of generalised contexts is:

$$\mathbb{C}, \mathbb{D} \quad ::= \quad \xi \cdot \vec{x}$$
$$| \quad x \mid \lambda x.\mathbb{C} \mid \mathbb{C}\,x \mid c\,\vec{x}$$
$$| \quad \text{let } \{\vec{x} = \vec{\mathbb{D}}\} \text{ in } \mathbb{C}$$
$$| \quad \text{case } \mathbb{C} \text{ of } \{c_i\,\vec{x}_i \to \mathbb{D}_i\}.$$

$\mathbb{V}$ and $\mathbb{W}$ will range over value contexts, $\Gamma$ and $\Delta$ over heap contexts, and $\mathbb{S}$ and $\mathbb{T}$ over stack contexts. Each "hole variable" $\xi$ has a fixed arity, and ranges over meta-abstractions of the form $(\vec{x})M$ where the length of $\vec{x}$ is the arity of $\xi$. In the meta-abstraction $(\vec{x})M$, the variables $\vec{x}$ are bound in $M$. Hole-filling is now a general non-capturing substitution: $[(\vec{x})M/\xi]$. The effect of a substitution is as expected (remembering that the $\vec{x}$ are considered bound in $(\vec{x})M$). Coupled with the meta-abstraction is of course meta-application, written $\xi \cdot \vec{x}$. We restrict application of $\xi$ to variables so that hole-filling cannot violate the restriction on syntax. In the definition of substitution we make the following identification:

$$(\vec{x})M \cdot \vec{y} \equiv M[\vec{y}/\vec{x}].$$

This definition of context generalises the usual definition since we can represent a traditional context $\mathbb{C}$ by $\mathbb{C}[\xi \cdot \vec{x}]$ where $\vec{x}$ is a vector of the capture-variables of $\mathbb{C}$; filling $\mathbb{C}$ with a term $M$ is then represented by $(\mathbb{C}[\xi \cdot \vec{x}])[(\vec{x})M/\xi]$.

**Example** The traditional context let $x = [\cdot]$ in $\lambda y.[\cdot]$ can be represented by let $x = \xi \cdot (x, y)$ in $\lambda y.\xi \cdot (x, y)$. Filling the hole with the term $x\,y$ is represented by:

$$(\text{let } x = \xi \cdot (x, y) \text{ in } \lambda y.\xi \cdot (x, y))[(x, y)\,x\,y/\xi]$$
$$\equiv \text{let } z = (x, y)\,x\,y \cdot (z, y) \text{ in } \lambda w.(x, y)\,x\,y \cdot (x, w)$$
$$\equiv \text{let } z = z\,y \text{ in } \lambda w.x\,w$$

which is $\alpha$-equivalent to what we would have obtained by the usual hole-filling with capture. Note that the generalised representation permits contexts to be identified up to $\alpha$-conversion.

Henceforth we work only with generalised contexts. We will write $\mathbb{C}[(\vec{x})M]$ to mean $\mathbb{C}[(\vec{x})M/\xi]$ when $\mathbb{C}$ contains just a single hole variable $\xi$. We assume that the arities of hole variables are always respected.

We implicitly generalise our definitions of improvement to work with generalised contexts. This is not quite identical to the earlier definition since with generalised contexts, when placing a term in a hole we obtain a substitution instance of the term. This means in particular that improvement is now closed under substitution (variable-for-variable) by definition — a useful property. This difference is a relatively minor technicality which we will gloss over in this appendix.

### A.2 Open Uniform Computation

The basis of our proofs will be to compute with configurations containing holes and free variables. Thanks to the capture-free representation of contexts, this means that normal reduction can be extended to contexts with ease. See [San98a] for a thorough treatment of generalised contexts and how they support generalisation of inductive definitions over terms.

Firstly, in order to fill the holes in a configuration we need to identify configurations up to renaming of the heap variables (recalling that update-markers on the stack are also binding occurrences of heap variables).

We tacitly extend the operational semantics to open configurations with holes. Note that holes can only occur in the stack within the branches of case alternatives. In what follows $\theta$ will range over substitutions composed of variable for variable substitutions and substitutions of the form $[(\vec{x}_i)M_i/\xi_i]$.

We have the following key property.

**Lemma A.1 (Extension)** *If* $\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle \to^k \langle \Delta, \mathbb{D}, \mathbb{T} \rangle$ *then*

*(i) for all $\Gamma'$ and $\mathbb{S}'$ such that $\langle \Gamma'\Gamma, \mathbb{C}, \mathbb{S}\mathbb{S}' \rangle$ is well-formed, $\langle \Gamma'\Gamma, \mathbb{C}, \mathbb{S}\mathbb{S}' \rangle \to^k \langle \Gamma'\Delta, \mathbb{D}, \mathbb{T}\mathbb{S}' \rangle$.*

*(ii) for all $\theta$, $\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle\theta \to^k \langle \Delta, \mathbb{D}, \mathbb{T} \rangle\theta$*

PROOF. (i) follows by inspection of possible open reductions over configuration contexts. (ii) amounts to the standard substitution lemma; see [San98a] for a general argument. □

**Definition A.1** $\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle\Downarrow^{\text{MAY}} \stackrel{\text{def}}{=} \exists\theta.\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle\theta\Downarrow$, *where $\theta$ is a closing substitution.*

The following *open uniform computation* property is central. It allows us to evaluate open configuration contexts until either the computation is finished, or we find ourselves in an "interesting" case.

**Lemma A.2 (Open Uniform Computation)** *Given some $\Gamma, \mathbb{C}, \mathbb{S}$, if there exist $\Gamma'$ and $\mathbb{S}'$ such that $\langle \Gamma\Gamma', \mathbb{C}, \mathbb{S}\mathbb{S}' \rangle$ is well-formed, and $\langle \Gamma\Gamma', \mathbb{C}, \mathbb{S}\mathbb{S}' \rangle\Downarrow^{\text{MAY}}$ then $\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle$ reduces to a configuration context of one of the following forms:*

*(i)* $\langle \mathbb{A}, \mathbb{V}, \epsilon \rangle$,

*(ii)* $\langle \mathbb{A}, \xi_i \cdot \vec{y}, \mathbb{T} \rangle$, *for some hole* $\xi_i$, *or*

*(iii)* $\langle \mathbb{A}, x, \mathbb{T} \rangle$, $x \in \mathrm{dom}\,\Gamma'$.

PROOF. Assume there exists some $\theta$ such that $\langle \Gamma\Gamma', \mathbb{C}, \mathbb{S}\mathbb{S}' \rangle \theta \Downarrow^n$. We consider the reduction of $\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle$ and proceed by induction on $n$ with cases on the structure of $\mathbb{C}$. We show three illustrative cases only. The others are similar.

$\mathbb{C} \equiv \xi_i \cdot \vec{y}$. This is a type (ii) context, so we are done.

$\mathbb{C} \equiv x$. Since we have termination, $x$ must be bound in either $\Gamma$ or $\Gamma'$. In the former case, $\Gamma \equiv \mathbb{A}\{x = \mathbb{D}\}$. By (*Lookup*), $\langle \mathbb{A}\{x = \mathbb{D}\}, x, \mathbb{S} \rangle$ reduces to $\langle \mathbb{A}, \mathbb{D}, \#x : \mathbb{S} \rangle$, and $\langle \mathbb{A}\Gamma', \mathbb{D}, \#x : \mathbb{S}\mathbb{S}' \rangle \theta \Downarrow^{n-1}$, by extension. By the inductive hypothesis, we know that $\langle \mathbb{A}, \mathbb{D}, \#x : \mathbb{S} \rangle$ reduces to a configuration context of type (i), (ii), or (iii), and therefore $\langle \mathbb{A}\{x = \mathbb{D}\}, x, \mathbb{S} \rangle$ does also, as required. In the latter case, $\langle \Gamma, x, \mathbb{S} \rangle$ is a type (iii) context, and we are done.

$\mathbb{C} \equiv \mathbb{V}$. There are four sub-cases, depending upon the structure of $\mathbb{S}$; we consider only the case when $\mathbb{S} \equiv x : \mathbb{T}$. Since we have termination, $\mathbb{V} \equiv \lambda y.\mathbb{D}$, and by (*Subst*), $\langle \Gamma, \lambda y.\mathbb{D}, x : \mathbb{T} \rangle$ reduces to $\langle \Gamma, \mathbb{D}[x/y], \mathbb{T} \rangle$, and $\langle \Gamma\Gamma', \mathbb{D}[x/y], \mathbb{T}\mathbb{S}' \rangle \Downarrow^{n-1}$. The inductive hypothesis applies, and the result follows as above. $\square$

In what follows we will use $\Sigma$ range over configuration contexts $\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle$.

## A.3 Translation

We can extend the definition of trans to cover open configurations and configuration contexts, and can therefore extend translation thus:

**Lemma A.3 (Translation)** *For all* $\mathbb{D}, \Gamma, \mathbb{C}, \mathbb{S}$ *such that* $\mathbb{D} \equiv \mathsf{trans}\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle$, *there exists* $k \geqslant 0$ *such that* $\langle \emptyset, \mathbb{D}, \epsilon \rangle \to^k \langle \Gamma, \mathbb{C}, \mathbb{S} \rangle$.

PROOF. Simple induction on $\mathbb{S}$. $\square$

## A.4 Proof: the Context Lemma

The proof of the context lemma relies upon two lemmas, the latter of which is the most complex.

**Lemma A.4** $M \mathrel{\underset{\sim}{\gtrsim}} N$ *if and only if for all* $\Sigma$ *and* $n$, $\Sigma[(\vec{x})M]\Downarrow^n \implies \Sigma[(\vec{x})N]\Downarrow^{\leqslant n}$.

PROOF. (Sketch) ($\Leftarrow$). Trivial; let $\Sigma = \langle \emptyset, \mathbb{C}, \epsilon \rangle$.
($\Rightarrow$). By a simple induction on $n$, using translation. $\square$

Note that the proof of the next lemma does not rely upon any specific cost-measure.

**Lemma A.5** *If for all* $\Gamma, S$, *and* $n$

$$\langle \Gamma, (\vec{x})M \cdot \vec{y}, S \rangle \Downarrow^n \implies \langle \Gamma, (\vec{x})N \cdot \vec{y}, S \rangle \Downarrow^{\leqslant n}$$

*then for all* $\Sigma$ *and* $n$, $\Sigma[(\vec{x})M]\Downarrow^n \implies \Sigma[(\vec{x})N]\Downarrow^{\leqslant n}$, *where* $\vec{x} \supseteq \mathsf{FV}(M, N)$.

PROOF. Assume the premise and suppose $\Sigma[(\vec{x})M]\Downarrow^n$. We proceed via induction on $n$. By open uniform computation, $\Sigma$ reduces in $k \geqslant 0$ steps to one of:

$$\textbf{(1) } \langle \mathbb{A}, \mathbb{V}, \epsilon \rangle, \qquad \textbf{(2) } \langle \mathbb{A}, \xi \cdot \vec{y}, \mathbb{S}' \rangle.$$

(There are only two possibilities since $\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle$ is closed.) In case (1), we are done. In case (2), we have

$$\Sigma[(\vec{x})N] \to^k \langle \mathbb{A}[(\vec{x})N], N[\vec{y}/\vec{x}], \mathbb{S}'[(\vec{x})N] \rangle. \qquad \text{(A.1)}$$

By open uniform computation, $\langle \mathbb{A}, M[\vec{y}/\vec{x}], \mathbb{S}' \rangle$ reduces in $k' \geqslant 0$ steps to one of:

$$\textbf{(2.1) } \langle \mathbb{A}', \mathbb{W}, \epsilon \rangle, \qquad \textbf{(2.2) } \langle \mathbb{A}', \xi \cdot \vec{z}, \mathbb{T} \rangle.$$

(Again, there are only two possibilities since $\langle \mathbb{A}, M[\vec{y}/\vec{x}], \mathbb{S}' \rangle$ is closed.) In case (2.1), we have that $\langle \mathbb{A}[(\vec{x})N], (\vec{x})M \cdot \vec{y}, \mathbb{S}'[(\vec{x})N] \rangle$ reduces in $k'$ steps to $\langle \mathbb{A}'[(\vec{x})N], \mathbb{W}[(\vec{x})N], \epsilon \rangle$, so

$$\begin{aligned}
& \langle \mathbb{A}[(\vec{x})N], M[\vec{y}/\vec{x}], \mathbb{S}'[(\vec{x})N] \rangle \Downarrow^{\leqslant n-k} \\
& \implies \langle \mathbb{A}[(\vec{x})N], N[\vec{y}/\vec{x}], \mathbb{S}'[(\vec{x})N] \rangle \Downarrow^{\leqslant n-k} \quad \text{(ass.)} \\
& \implies \Sigma[(\vec{x})N]\Downarrow^{\leqslant n} \qquad\qquad\qquad\qquad \text{(A.1)}
\end{aligned}$$

as required. In case (2.2), we know that $k' > 0$, since $M[\vec{y}/\vec{x}] \not\equiv \xi \cdot \vec{z}$. We have

$$\begin{aligned}
\langle \mathbb{A}[(\vec{x})M], (\vec{x})M \cdot \vec{y}, \mathbb{S}'[(\vec{x})M] \rangle \to^{k'} \\
\langle \mathbb{A}'[(\vec{x})M], (\vec{x})M \cdot \vec{z}, \mathbb{T}[(\vec{x})M] \rangle, \quad \text{(A.2)}
\end{aligned}$$

and

$$\begin{aligned}
\langle \mathbb{A}[(\vec{x})N], (\vec{x})M \cdot \vec{y}, \mathbb{S}'[(\vec{x})N] \rangle \to^{k'} \\
\langle \mathbb{A}'[(\vec{x})N], (\vec{x})N \cdot \vec{z}, \mathbb{T}[(\vec{x})N] \rangle. \quad \text{(A.3)}
\end{aligned}$$

Therefore

$$\begin{aligned}
& \langle \mathbb{A}'[(\vec{x})M], M[\vec{z}/\vec{x}], \mathbb{T}[(\vec{x})M] \rangle \Downarrow^{n-k-k'} && \text{(A.2)} \\
& \implies \langle \mathbb{A}'[(\vec{x})N], M[\vec{z}/\vec{x}], \mathbb{T}[(\vec{x})N] \rangle \Downarrow^{\leqslant n-k-k'} && \text{(I.H.)} \\
& \implies \langle \mathbb{A}'[(\vec{x})N], N[\vec{z}/\vec{x}], \mathbb{T}[(\vec{x})N] \rangle \Downarrow^{\leqslant n-k-k'} && \text{(ass.)} \\
& \implies \langle \mathbb{A}[(\vec{x})N], M[\vec{y}/\vec{x}], \mathbb{S}'[(\vec{x})N] \rangle \Downarrow^{\leqslant n-k} && \text{(A.3)} \\
& \implies \langle \mathbb{A}[(\vec{x})N], N[\vec{y}/\vec{x}], \mathbb{S}'[(\vec{x})N] \rangle \Downarrow^{\leqslant n-k} && \text{(ass.)} \\
& \implies \Sigma[(\vec{x})N]\Downarrow^{\leqslant n} && \text{(A.1)}
\end{aligned}$$

as required. $\square$

The generalised statement of the context lemma is:
*For all terms* $M$ *and* $N$, *if*

$$\forall \Gamma, S, \sigma, n. \langle \Gamma, M\sigma, S \rangle \Downarrow^n \implies \langle \Gamma, N\sigma, S \rangle \Downarrow^{\leqslant n}$$

*then* $M \mathrel{\underset{\sim}{\gtrsim}} N$.

This follows from lemmas A.4 and A.5, and the fact that $M\sigma \equiv (\vec{x})M \cdot \vec{y}$ for $\sigma = [\vec{y}/\vec{x}]$.

## A.5 Validating the Tick Algebra

We present proofs of the validity of (*value-β*) and (*value-copy*), and sketch a proof of the correspondence between evaluation contexts and configuration contexts of the form $\langle \Gamma, [\cdot], S \rangle$. The proofs of the more complex laws (*e.g.* (*var-β*), (*var-abs*), (*var-subst*), and ($\checkmark$-*float*)) have a similar structure to that for (*value-β*), except they require more use of open uniform computation.

**A.5.1  Proof:** (*value-β*)

Recall (*value-β*):

$$\text{let } \{x = V, \vec{y} = \vec{\mathbb{D}}[x]\} \text{ in } \mathbb{C}[x]$$
$$\underset{\approx}{\Leftrightarrow} \text{let } \{x = V, \vec{y} = \vec{\mathbb{D}}[^{2^{\vee}}V]\} \text{ in } \mathbb{C}[^{2^{\vee}}V].$$

Let $W \equiv {}^{2^{\vee}}V$ throughout. It suffices to show

$$\forall \Gamma, \mathbb{S}. \langle \Gamma[x]\{x = V\}, \ \mathbb{C}[x], \ \mathbb{S}[x] \rangle \Downarrow^{n} \iff$$
$$\langle \Gamma[W]\{x = V\}, \ \mathbb{C}[W], \ \mathbb{S}[W] \rangle \Downarrow^{n}$$

where $x \notin \text{dom}(\Gamma, \mathbb{S})$, and the only hole is $[\cdot]$, a non-capturing hole. We prove the forward direction only; the reverse direction is similar.

Suppose $\langle \Gamma[x]\{x = V\}, \ \mathbb{C}[x], \ \mathbb{S}[x] \rangle \Downarrow^{n}$. We proceed by induction on $n$. By open uniform computation, $\langle \Gamma, \mathbb{C}, \mathbb{S} \rangle$ reduces in $k \geqslant 0$ steps to one of

**(1)** $\langle \Delta, \mathbb{V}, \epsilon \rangle$,  **(2)** $\langle \Delta, [\cdot], \mathbb{T} \rangle$,  **(3)** $\langle \Delta, x, \mathbb{T} \rangle$.

In case (1), we are done. In case (2), by extension, (*Lookup*) and (*Update*), we have

$$\langle \Gamma[x]\{x = V\}, \ \mathbb{C}[x], \ \mathbb{S}[x] \rangle \rightarrow^{k} \langle \Delta[x]\{x = V\}, \ x, \ \mathbb{T}[x] \rangle$$
$$\rightarrow^{2} \langle \Delta[x]\{x = V\}, \ V, \ \mathbb{T}[x] \rangle,$$

and by extension and the definition of $W$,

$$\langle \Gamma[W]\{x = V\}, \ \mathbb{C}[W], \ \mathbb{S}[W] \rangle$$
$$\rightarrow^{k} \langle \Delta[W]\{x = V\}, \ W, \ \mathbb{T}[W] \rangle$$
$$\rightarrow^{2} \langle \Delta[W]\{x = V\}, \ V, \ \mathbb{T}[W] \rangle.$$

Since $\langle \Delta[x]\{x = V\}, \ V, \ \mathbb{T}[x] \rangle \Downarrow^{n-(k+2)}$, by the inductive hypothesis we have $\langle \Delta[W]\{x = V\}, \ V, \ \mathbb{T}[W] \rangle \Downarrow^{n-(k+2)}$, and the result follows.

In case (3), we have $\langle \Delta[x]\{x = V\}, \ V, \ \mathbb{T}[x] \rangle \Downarrow^{n-k-2}$, as above. Furthermore, by extension, (*Lookup*) and (*Update*), we have

$$\langle \Gamma[W]\{x = V\}, \ \mathbb{C}[W], \ \mathbb{S}[W] \rangle$$
$$\rightarrow^{k} \langle \Delta[W]\{x = V\}, \ x, \ \mathbb{T}[W] \rangle$$
$$\rightarrow^{2} \langle \Delta[W]\{x = V\}, \ V, \ \mathbb{T}[W] \rangle.$$

From the inductive hypothesis, we have $\langle \Delta[W]\{x = V\}, \ V, \ \mathbb{T}[W] \rangle \Downarrow^{n-k-2}$, and the result follows.

**A.5.2  Proof:** (*value-copy*)

Recall (*value-copy*):

$$\text{let } \{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta, \vec{z} = \vec{M}\} \text{ in } N$$
$$\underset{\approx}{\Leftrightarrow} \text{let } \{\vec{x} = \vec{V}\theta\sigma, \vec{z} = \vec{M}\sigma\} \text{ in } N\sigma,$$

where $\sigma = [\vec{x}/_{\vec{y}}]$ and $\theta = [\vec{x}/_{\vec{w}}]$.

It suffices to show that for all $\Gamma$, $S$, and $n$,

$$\langle \Gamma\{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta\}, \ N, \ S \rangle \Downarrow^{n}$$
$$\iff \langle \Gamma\{\vec{x} = \vec{V}\theta\sigma\}, \ N\sigma, \ S \rangle \Downarrow^{n}.$$

We show only the forward direction. To show the reverse, we need only establish termination, which follows by the fact that call-by-name and call-by-need agree on termination.

Consider the (holeless) open configuration context $\langle \Gamma, \ N, \ S \rangle$, in which the $\vec{x}$ and $\vec{y}$ may appear free. By open uniform computation, this reduces in $k \geqslant 0$ steps to one of:

**(1)** $\langle \Delta, \ W, \ \epsilon \rangle$,  **(2a)** $\langle \Delta, \ x_i, \ T \rangle$,  **(2b)** $\langle \Delta, \ y_i, \ T \rangle$.

In case (1), we are done. In case (2a), by extension, (*Lookup*) and (*Update*), we have

$$\langle \Gamma\{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta\}, \ N, \ S \rangle$$
$$\rightarrow^{k} \langle \Delta\{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta\}, \ x_i, \ T \rangle$$
$$\rightarrow^{2} \langle \Delta\{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta\}, \ V_i[\vec{y}/_{\vec{w}}], \ T \rangle,$$

and furthermore,

$$\langle \Delta\{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta\}, \ V_i[\vec{y}/_{\vec{w}}], \ T \rangle \Downarrow^{n-(k+2)}. \quad \text{(A.4)}$$

Similarly, by extension, (*Lookup*) and (*Update*), we have also that

$$\langle \Gamma\{\vec{x} = \vec{V}\theta\sigma\}, \ N\sigma, \ S \rangle$$
$$\rightarrow^{k} \langle \Delta\{\vec{x} = \vec{V}\theta\sigma\}, \ x_i\sigma, \ T \rangle$$
$$\rightarrow^{2} \langle \Delta\{\vec{x} = \vec{V}\theta\sigma\}, \ V_i\theta\sigma, \ T \rangle.$$

By elementary properties of substitution,

$$V_i[\vec{y}/_{\vec{w}}][\vec{x}/_{\vec{y}}] \equiv V_i[\vec{x}/_{\vec{w}}][\vec{x}/_{\vec{y}}],$$

so the inductive hypothesis applies (with $N \equiv V_i[\vec{y}/_{\vec{w}}]$), yielding the desired result.

In case (2b), by extension, (*Lookup*) and (*Update*), we have

$$\langle \Gamma\{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta\}, \ N, \ S \rangle$$
$$\rightarrow^{k} \langle \Delta\{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta\}, \ y_i, \ T \rangle$$
$$\rightarrow^{2} \langle \Delta\{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta\}, \ V_i\theta, \ T \rangle,$$

and furthermore,

$$\langle \Delta\{\vec{x} = \vec{V}[\vec{y}/_{\vec{w}}], \vec{y} = \vec{V}\theta\}, \ V_i\theta, \ T \rangle \Downarrow^{n-(k+2)}. \quad \text{(A.5)}$$

Similarly, by extension, (*Lookup*) and (*Update*), we have also that

$$\langle \Gamma\{\vec{x} = \vec{V}\theta\sigma\}, \ N\sigma, \ S \rangle$$
$$\rightarrow^{k} \langle \Delta\{\vec{x} = \vec{V}\theta\sigma\}, \ y_i\sigma, \ T \rangle$$
$$\equiv \langle \Delta\{\vec{x} = \vec{V}\theta\sigma\}, \ x_i, \ T \rangle$$
$$\rightarrow^{2} \langle \Delta\{\vec{x} = \vec{V}\theta\sigma\}, \ V_i\theta\sigma, \ T \rangle.$$

The inductive hypothesis applies (with $N \equiv V_i\theta$), yielding the desired result.

### A.5.3 Proof Sketch: Lemma 4.1

Recall the statement of lemma 4.1: $\Lambda_{\mathbb{E}}$ is equal to $\{\text{trans}\langle \Gamma,\ [\cdot],\ S \rangle \mid \text{all } \Gamma \text{ and } S\}$. So we need to show that:

(i) $\forall \Gamma, S.\ \exists \mathbb{E}.\ \text{trans}\langle \Gamma,\ [\cdot],\ S \rangle \equiv \mathbb{E}$, and

(ii) $\forall \mathbb{E}.\ \exists \Gamma, S.\ \text{trans}\langle \Gamma,\ [\cdot],\ S \rangle \equiv \mathbb{E}$.

First note that $\Lambda_{\mathbb{A}}$ (the set of all applicative contexts) is in 1-1 correspondence to update-marker free stacks, realised by the following isomorphism (writing $[x]$ for the singleton stack):

$$[\cdot]^{\circ} = \epsilon$$
$$(\mathbb{A}\,x)^{\circ} = \mathbb{A}^{\circ}[x]$$
$$(\text{case } \mathbb{A} \text{ of } alts)^{\circ} = \mathbb{A}^{\circ}\,alts$$

$(\cdot)^{\circ}$ takes $\Lambda_{\mathbb{A}}$ into the set of update-marker free stacks. Call its inverse $(\cdot)^{\bullet}$.

Then show that $\langle \Gamma,\ \mathbb{A}[\mathbb{C}],\ \mathbb{S} \rangle \to^{*} \langle \Gamma,\ \mathbb{C},\ \mathbb{A}^{\circ}\mathbb{S} \rangle$ and that $\text{trans}\langle \Gamma,\ \mathbb{A}[\mathbb{C}],\ \mathbb{S} \rangle$ is identical to $\text{trans}\langle \Gamma,\ \mathbb{C},\ \mathbb{A}^{\circ}\mathbb{S} \rangle$ by induction on the structure of $\mathbb{A}$.

To show (i), generalise the statement to show that for all $\Gamma$ and $S$ both $\text{trans}\langle \Gamma,\ \mathbb{A},\ S \rangle$ and $\text{trans}\langle \Gamma\{x_0 = \mathbb{A}_0[x_1], \ldots, x_n = \mathbb{A}_n\},\ \mathbb{A}[x_0],\ S \rangle$ are evaluation contexts by induction on the number of update markers in $S$.

To show (ii), proceed by case analysis on $\mathbb{E}$, and produce a $\Gamma$ and $S$ in each case. The difficult case is when

$$\mathbb{E} \equiv \text{let } \{\vec{y} = \vec{M}, x_0 = \mathbb{A}_0[x_1], \ldots, x_n = \mathbb{A}_n\}$$
$$\text{in } \mathbb{A}[x_0].$$

Here, let $\Gamma$ be $\{\vec{y} = \vec{M}\}$ and let $S$ be

$$\mathbb{A}_n^{\circ}\#x_n \cdots \mathbb{A}_1^{\circ}\#x_1\mathbb{A}_0^{\circ}\#x_0\mathbb{A}^{\circ}.$$

### A.5.4 A lemma for ($case$-$\mathbb{E}$) and ($let$-$\mathbb{E}$)

The following lemma can be used to validate ($case$-$\mathbb{E}$) and ($let$-$\mathbb{E}$). $\mathsf{CV}(\mathbb{E})$ denotes the *capture variables* of $\mathbb{E}$.

**Lemma A.6** *For all $\mathbb{E}$, there exist $\Gamma, S$, such that* $\text{dom}(\Gamma, S) \subseteq \mathsf{CV}(\mathbb{E})$ *and* $\forall \Delta, T.\langle \Delta,\ \mathbb{E},\ T \rangle \to^{*} \langle \Delta\Gamma,\ [\cdot],\ ST \rangle.$

PROOF. By lemma 4.1, there exist $\Gamma$ and $S$ such that $\text{trans}\langle \Gamma,\ [\cdot],\ S \rangle \equiv \mathbb{E}$, so by translation $\langle \emptyset,\ \mathbb{E},\ \epsilon \rangle \to^{*} \langle \Gamma,\ [\cdot],\ S \rangle$, and thus by extension, provided $\text{dom}(\Gamma, S) \subseteq \mathsf{CV}(\mathbb{E})$, $\langle \Delta,\ \mathbb{E},\ T \rangle \to^{*} \langle \Delta\Gamma,\ [\cdot],\ ST \rangle$. □

### A.6 Proof: the Unwinding Lemma

To prove the Unwinding lemma we will need the following lemma, which we state without proof.

**Lemma A.7** *For all $M, \Gamma, S, V$ and $n$,*

$$\langle \Gamma\{x \overset{k}{=} V\},\ M,\ S \rangle\Downarrow^{n} \implies \langle \Gamma\{x \overset{k+1}{=} V\},\ M\sigma,\ S\sigma \rangle\Downarrow^{n}$$

*where $\sigma = [x_{k+1}/x_k]$ and $\{x_i\}_{i=0}^{k+1} \notin \mathsf{FV}(V)$.*

Recall the statement of the Unwinding lemma:

*For all $\Gamma, S$, and $n$,*

$$\langle \Gamma,\ \text{let } \{f = V\} \text{ in } M,\ S \rangle\Downarrow^{n} \implies$$
$$\exists m.\langle \Gamma,\ \text{let } \{f \overset{m}{=} V\} \text{ in } M[f_m/f],\ S \rangle\Downarrow^{n}.$$

It suffices to prove that for all $\Gamma$, $S$, and $n$ such that $\{x_i\}_{i=0}^{n} \notin \mathsf{FV}(\Gamma, S)$,

$$\langle \Gamma\{x = V\},\ M,\ S \rangle\Downarrow^{n} \implies \langle \Gamma\sigma\{x \overset{n}{=} V\},\ M\sigma,\ S\sigma \rangle\Downarrow^{n}$$

where $\sigma = [x_n/x]$ (*i.e.* $m = n$). Suppose $\langle \Gamma\{x = V\},\ M,\ S \rangle\Downarrow^{n}$. We proceed by induction on $n$. By open uniform computation, $\langle \Gamma,\ M,\ S \rangle$ reduces in $k \geqslant 0$ steps to one of

$$\textbf{(1) } \langle \Delta,\ W,\ \epsilon \rangle, \qquad \textbf{(2) } \langle \Delta,\ x,\ T \rangle.$$

(Type (ii) cannot occur, since there is no hole involved.) By extension, the corresponding result holds for $\langle \Gamma,\ M,\ S \rangle\sigma$, and hence for $\langle \Gamma\sigma,\ M\sigma,\ S\sigma \rangle$, since $x_n$ is free in $\langle \Gamma,\ M,\ S \rangle$.

Therefore, in case (1), by extension, $\langle \Gamma\sigma\{x = V\},\ M\sigma,\ S\sigma \rangle$ reduces in $k$ steps to $\langle \Delta\sigma\{x = V\},\ W\sigma,\ \epsilon \rangle$ and we are done, since $k = n$. In case (2), by extension, (*Lookup*), and (*Update*),

$$\langle \Gamma\sigma\{x = V\},\ M\sigma,\ S\sigma \rangle \to^{k} \langle \Delta\sigma\{x = V\},\ x_n,\ T\sigma \rangle$$
$$\to^{2} \langle \Delta\sigma\{x = V\},\ V\sigma,\ T\sigma \rangle.$$

Similarly, $\langle \Gamma\{x = V\},\ M,\ S \rangle$ reduces in $k + 2$ steps to $\langle \Delta\{x = V\},\ V,\ T \rangle$. By the inductive hypothesis, we know that $\langle \Delta\sigma'\{x \overset{k'}{=} V\},\ V\sigma',\ T\sigma' \rangle\Downarrow^{k'}$ where $\sigma' = [x_{k'}/x]$ and $k' = n - k - 2$. By repeated application of lemma A.7, we have that $\langle \Delta\sigma\{x \overset{n}{=} V\},\ V\sigma,\ T\sigma \rangle\Downarrow^{k'}$ and hence $\langle \Gamma\sigma\{x \overset{n}{=} V\},\ M\sigma,\ S\sigma \rangle\Downarrow^{n}$ as required.

### A.7 Proof: the Improvement Theorem

This lemma is used to prove lemma A.9.

**Lemma A.8** *For all $\Gamma, S$, and $n$*

$$\langle \Gamma\{\vec{x} = \vec{V}\},\ M,\ S \rangle\Downarrow^{n} \iff$$
$$\langle \Gamma\sigma\{\vec{x} = \vec{V}, \vec{y} = \vec{V}\sigma\},\ M\sigma,\ S\sigma \rangle\Downarrow^{n}.$$

*where $\sigma = [y/x]$.*

PROOF. (Sketch) ($\Rightarrow$) Simple induction on $n$, with cases of the structure of $M$.

($\Leftarrow$) It is sufficient to show that termination is implied. This is true for the call-by-name theory, and therefore here also. □

To prove the improvement theorem, we will need the following lemma, which is stated without proof. It follows in a straightforward fashion from lemma A.8. (It may seem to follow from the context lemma, but the $\vec{x}$ may appear free in $\Gamma$ in the conclusion, so it does not.)

**Lemma A.9** *If* let $\{\vec{x} = \vec{V}\}$ in $M \gtrsim$ let $\{\vec{x} = \vec{V}\}$ in $N$ *then for all $\Gamma$ and $S$,*

$$\langle \Gamma\{\vec{x} = \vec{V}\},\ M,\ S\rangle\Downarrow^n \implies \langle \Gamma\{\vec{x} = \vec{V}\},\ N,\ S\rangle\Downarrow^{\leqslant n}.$$

We prove the improvement theorem generalised to mutually-recursive definitions:

*The following proof rule is sound:*

$$\frac{\forall j \in I.\ \mathsf{let}\ \{f_i = V_i\}_{i \in I}\ \mathsf{in}\ V_j \gtrsim \mathsf{let}\ \{f_i = V_i\}_{i \in I}\ \mathsf{in}\ W_j}{\mathsf{let}\ \{f_i = V_i\}_{i \in I}\ \mathsf{in}\ N \underset{\sim}{\gtrsim} \mathsf{let}\ \{f_i = W_i\}_{i \in I}\ \mathsf{in}\ N}$$

By the context lemma it suffices to show that for all $\Gamma, S$, and $n$,

$$\langle \Gamma\{\vec{f} = \vec{V}\},\ N,\ S\rangle\Downarrow^n \implies \langle \Gamma\{\vec{f} = \vec{W}\},\ N,\ S\rangle\Downarrow^{\leqslant n}.$$

Assume the premise, and suppose that $\langle \Gamma\{\vec{f} = \vec{V}\},\ N,\ S\rangle\Downarrow^n$. We proceed by induction on $n$. By open uniform computation, $\langle \Gamma,\ N,\ S\rangle$ reduces in $k \geqslant 0$ steps to one of

$$(\mathbf{1})\ \langle \Delta,\ V,\ \epsilon\rangle, \qquad (\mathbf{2})\ \langle \Delta,\ f_i,\ T\rangle.$$

In case (1), we have by extension that $\langle \Gamma\{\vec{f} = \vec{W}\},\ N,\ S\rangle$ reduces in $k$ steps to $\langle \Delta\{\vec{f} = \vec{W}\},\ V,\ \epsilon\rangle$ and $k = n$, so we are done. In case (2),

$$\langle \Gamma\{\vec{f} = \vec{V}\},\ N,\ S\rangle \to^k \langle \Delta\{\vec{f} = \vec{V}\},\ f_i,\ T\rangle$$
$$\to^2 \langle \Delta\{\vec{f} = \vec{V}\},\ V_i,\ T\rangle \qquad (\text{A.6})$$

and

$$\langle \Gamma\{\vec{f} = \vec{W}\},\ N,\ S\rangle \to^k \langle \Delta\{\vec{f} = \vec{W}\},\ f_i,\ T\rangle$$
$$\to^2 \langle \Delta\{\vec{f} = \vec{W}\},\ W_i,\ T\rangle \quad (\text{A.7})$$

so

$$\langle \Delta\{\vec{f} = \vec{V}\},\ V_i,\ T\rangle\Downarrow^{n-(k+2)} \qquad (\text{A.6})$$
$$\implies \langle \Delta\{\vec{f} = \vec{V}\},\ W_i,\ T\rangle\Downarrow^{\leqslant n-(k+2)} \quad (\text{ass., lem. A.9})$$
$$\implies \langle \Delta\{\vec{f} = \vec{W}\},\ W_i,\ T\rangle\Downarrow^{\leqslant n-(k+2)} \quad (\text{I.H.})$$
$$\implies \langle \Gamma\{\vec{f} = \vec{W}\},\ N,\ S\rangle\Downarrow^{\leqslant n}. \qquad (\text{A.7})$$

## A.8   Proof: Improvement Induction

Recall the statement of improvement induction:

*For any $M$, $N$ and substitution $\sigma$, the following proof rule is sound:*

$$\frac{M \gtrsim {}^{\curvearrowleft}\mathbb{C}[M\sigma] \quad N \underset{\sim}{\gtrless} {}^{\curvearrowleft}\mathbb{C}[N\sigma]}{M \underset{\sim}{\gtrsim} N}$$

We generalise $\mathbb{C}[M\sigma]$ to $\mathbb{C}[(\vec{x})M]$. It suffices to show under assumption of the premise, that, for all $\Sigma$ such that $\Sigma[(\vec{x})M]$ and $\Sigma[(\vec{x})N]$ are closed, $\Sigma[(\vec{x})M]\Downarrow^n \implies \Sigma[(\vec{x})N]\Downarrow^{\leqslant n}$.

Suppose $\Sigma[(\vec{x})M]\Downarrow^n$. We proceed by induction on $n$. By open uniform computation, $\Sigma$ reduces in $k \geqslant 0$ to one of

$$(\mathbf{1})\ \langle \Delta,\ \mathbb{V},\ \epsilon\rangle, \qquad (\mathbf{2})\ \langle \Delta,\ \xi \cdot \vec{y},\ \mathbb{T}\rangle.$$

In case (1), we are done. In case (2), first note that, letting $\sigma = [\vec{y}/\vec{x}]$, $\mathbb{C}[(\vec{x})M]\sigma \equiv \mathbb{C}\sigma[(\vec{x})M]$ since $\vec{x} \supseteq \mathsf{FV}(M)$, and similarly for $N$. Then we have that

$$\Sigma[(\vec{x})N] \to^k \langle \Delta[(\vec{x})N],\ N\sigma,\ \mathbb{T}[(\vec{x})N]\rangle \qquad (\text{A.8})$$

and

$$\langle \Delta[(\vec{x})M],\ M\sigma,\ \mathbb{T}[(\vec{x})M]\rangle\Downarrow^{n-k}$$
$$\implies \langle \Delta[(\vec{x})M],\ {}^{\curvearrowleft}\mathbb{C}[(\vec{x})M]\sigma,\ \mathbb{T}[(\vec{x})M]\rangle\Downarrow^{\leqslant n-k} \qquad (\text{ass.})$$
$$\implies \langle \Delta[(\vec{x})M],\ \mathbb{C}[(\vec{x})M]\sigma,\ \mathbb{T}[(\vec{x})M]\rangle\Downarrow^{\leqslant n-(k+1)} \quad (\checkmark)$$
$$\implies \langle \Delta[(\vec{x})N],\ \mathbb{C}\sigma[(\vec{x})N],\ \mathbb{T}[(\vec{x})N]\rangle\Downarrow^{\leqslant n-(k+1)} \quad (\text{I.H.})$$
$$\implies \langle \Delta[(\vec{x})N],\ {}^{\curvearrowleft}\mathbb{C}[(\vec{x})N]\sigma,\ \mathbb{T}[(\vec{x})N]\rangle\Downarrow^{\leqslant n-k} \quad (\checkmark)$$
$$\implies \langle \Delta[(\vec{x})N],\ N\sigma,\ \mathbb{T}[(\vec{x})N]\rangle\Downarrow^{\leqslant n-k} \qquad (\text{ass.})$$
$$\implies \Sigma[(\vec{x})N]\Downarrow^{\leqslant n}. \qquad (\text{A.8})$$