
A Naïve Time Analysis and its Theory of Cost Equivalence

DAVID SANDS, *DIKU, Department of Computer Science, University of Copenhagen, Universitetsparken 1, DK-2100 Copenhagen, Denmark.*
E-mail: *dave@diku.dk*

Abstract

Techniques for reasoning about extensional properties of functional programs are well understood, but methods for analysing the underlying intensional or operational properties have been much neglected. This paper begins with the development of a simple but useful calculus for time analysis of non-strict functional programs with lazy lists. One limitation of this basic calculus is that the ordinary equational reasoning on functional programs is not valid. In order to buy back some of these equational properties we develop a non-standard operational equivalence relation called *cost equivalence*, by considering the number of computation steps as an ‘observable’ component of the evaluation process. We define this relation by analogy with Park’s definition of bisimulation in CCS. This formulation allows us to show that cost equivalence is a contextual congruence (and thus is substitutive with respect to the basic calculus) and provides useful proof techniques for establishing cost-equivalence laws. It is shown that basic evaluation time can be derived by demonstrating a certain form of cost equivalence, and we give an axiomatization of cost equivalence which is complete with respect to this application. This shows that cost equivalence subsumes the basic calculus. Finally we show how a new operational interpretation of evaluation demands can be used to provide a smooth interface between this time analysis and more compositional approaches, retaining the advantages of both.

Keywords: Time analysis, lazy evaluation, operational semantics.

1 Introduction

An appealing property of functional programming languages is the ease with which the *extensional* properties of a program can be understood—above all the ability to show that operations on programs preserve meaning. Prominent in the study of algorithms in general, and central to formal activities such as program transformation and parallelization, are questions of efficiency, i.e. the running-time and space requirements of programs. These are *intensional* properties of a program—properties of *how* the program computes, rather than *what* it computes. The study of intensional properties is not immediately amenable to the algebraic methods with which extensional properties are so readily explored. Moreover, the declarative emphasis of functional programs, together with some of the features that afford expressive power and modularity, namely higher-order functions and lazy evaluation, serve to make intensional properties *more opaque*. In spite of this, relatively little attention has been given to the development of methods for reasoning about the computational cost of functional programs.

As a motivating example consider the following defining equations for insertion sort (written in a Haskell-like syntax)

2 A Naïve Time Analysis and its Theory of Cost Equivalence

```
isort []      = []
isort (h:t)   = insert h (isort t)

insert x []   = [x]
insert x (h:t) = x:(h:t)  if x ≤ h
               = h:(insert x t) otherwise.
```

As expected, `isort` requires $\mathcal{O}(n^2)$ time to sort a list of length n . However, under lazy evaluation, `isort` enjoys a rather nice modularity property with respect to time: if we specify a program which computes the minimum of a list of numbers, by taking the head of the sorted list,¹

$$\text{minimum} = \text{head} \circ \text{isort}$$

then the time to compute `minimum` is only $\mathcal{O}(n)$. This rather pleasing property of insertion-sort is a well-used example in the context of reasoning about running time of lazy evaluation.

By contrast, the following time property of lazy ‘quicksort’ is seldom reported. A typical definition of a functional quicksort over lists might be:

```
qsort []      = []
qsort (h:t)   = qsort (below h t) ++ (h:qsort (above h t))
```

where `below` and `above` return lists of elements from `t` which are no bigger, and strictly smaller than `h`, respectively, and `++` is infix list-append. Functional accounts of quicksort are also quadratic time algorithms, but conventional wisdom would label quicksort as a better algorithm than insertion-sort because of its better average-case behaviour. A rather less pleasing property of lazy evaluation is that by replacing ‘better’ sorting algorithm `qsort` for `isort` in the definition of `minimum`, we obtain an asymptotically *worse* algorithm, namely one which is $\Omega(n^2)$ in the length of the input.

1.1 Overview

In the first part of the paper we consider the problem of reasoning about evaluation time in terms of a very simple measure of evaluation cost. A simple set of *time rules* are derived very directly from a call-by-name operational model, and concern equations on $\langle e \rangle^H$ (the ‘time’ to evaluate expression e to (weak) head normal form) and $\langle e \rangle^N$ (the time to evaluate e to normal-form). The approach is naïve in the sense that it is non-compositional (in general, the cost of computing an expression is not defined as a combination of the costs of computing its subexpressions), and does not model graph reduction. However, despite (or perhaps because of) its simplicity, the method appears to be useful as a means of formalizing sufficiently many operational details to reason (rigorously, but not necessarily formally) about the complexity of lazy algorithms.

One of the principal limitations of the approach is the fact that the usual meanings of ‘equality’ for programs do not provide equational reasoning in the context of the time rules. This problem motivates development of a non-standard theory of operational equivalence in which the number of computation steps are viewed as an ‘observable’ component of the evaluation process. We

¹The example is originally due to D. Turner; it appears as an exercise in informal reasoning about lazy evaluation in [6][Ch. 6], and in the majority(!) of papers on time analysis of non-strict evaluation.

define this relation by analogy with Park’s definition of bisimulation between processes. This formulation provides a uniform method for establishing cost-equivalence laws, and together with the key result that cost equivalence is a contextual congruence, provides a useful substitutive equivalence with which the time rules can be extended, since if e_1 is cost equivalent to e_2 , then for any syntactic context C , $\langle C[e_1] \rangle^\alpha = \langle C[e_2] \rangle^\alpha$.

In addition we show that the theory of cost equivalence subsumes the time rules, by providing an axiomatization of cost equivalence which is sound and complete (in a certain sense) with respect to simple evaluation time properties of expressions.

Finally, we return to a significant flaw in the time model, namely of its use of call-by-name rather than call-by-need. We sketch a method to alleviate this problem which provides a smooth integration of the simple time analysis here, and the more compositional call-by-need approaches, with some of the advantages of both.

The development of the theory of cost equivalence is somewhat technical, but the paper is written so that the reader interested primarily in the problem of time analysis of programs in a lazy language should be able to skip the bulk of the technical development, but still take advantage of its results, namely cost equivalence. In the remainder of this introduction we summarize the rest of the paper.

Sections 2 to 5 develop a simple time analysis for a first order language with lazy lists. Section 2 gives some background describing approaches to the efficiency analysis of lazy functional programs. In Section 3 we define our language and its operational semantics. Section 4 defines the notion of time cost over this operational model, and introduces the time rules which form the basis of the calculus. Section 5 provides some examples of the use of the time rules in reasoning about the complexity of simple programs.

Section 6 motivates and develops the theory of cost equivalence. Cost equivalence is based upon a cost simulation preordering which is shown to be preserved by substitution into arbitrary program contexts. It is also shown that it is the largest such relation. Section 7 gives some variants of the *co-induction* proof principal which are useful for establishing cost equivalences, and presents an axiomatization of cost equivalence which is complete with respect to the basic time properties of expressions. In Section 8 we extend the language with higher-order functions. Time rules are easily added to the new language, and the theory of cost equivalence is extended in the obvious way by considering an ‘applicative’ cost simulation, which is also shown to have the necessary substitutivity property.

Section 9 presents an example time analysis, illustrating the combined use of time rules and cost equivalences.

Section 10 outlines a flexible approach to increasing the compositionality and accuracy of the time analysis with respect to call-by-need evaluation, via the definition of a family of *evaluators* indexed by representations of strictness properties.

To conclude, we consider related work in the area of intensional semantics.

A preliminary version of this paper appeared as [39], and summarized [36][Ch. 4]. In addition to the inclusion of proofs, further examples and additional technical results, Sections 7, 8, 9 and 10 are new, and contain a number of important extensions to the earlier work.

2 Time analysis: background

A number of researchers have developed prototype (time) complexity analysis tools in which the algorithm under analysis is expressed as a first-order call-by-value functional program [44, 22,

33, 17]. It could be argued that the subject of study in these cases is not functional programming *per se*; the choice of a functional language is motivated by the fact that, for a first-order language with a call-by-value semantics, it is straightforward to construct, mechanically, functions with the same domain as a given function, which describe (recursively) the number of computation steps required by that function. Although this does not by any means trivialize the problem of finding solutions to these equations in terms of some size-measure of the arguments, it gives a simple but formal reading of the program as a description of computational cost. This is because the cost of evaluating some function application $\text{fun}(E)$ can be understood in terms of the cost of evaluating E , plus the cost of evaluating the application of fun to the value of E .

In the case of a higher-order strict language cost is not only dependent on the simple cost of evaluating the argument E , but also on the possible cost of subsequently applying E , applying the result of an application, and so forth. Techniques for handling this problem were introduced in [34], where syntactic structures called *cost closures* were introduced to enable intensional properties to be carried by functions. Additional techniques for reasoning about higher-order functions which complement this approach are described in [36].

A problem in reasoning about the efficiency of programs under lazy evaluation (i.e. call-by-name, or more usually, call-by-need, extended to data structures) is that the cost of computing the subexpression E is dependent entirely on the way in which the expression is used in the function fun . More generally, the cost of evaluating some (sub)expression is dependent on the amount of its value needed by its context.

The compositional approach

One approach to reasoning about the time cost of lazy evaluation is to parameterize the description of the cost of computing an expression by a description of the amount of the result that is *needed* by the context in which it appears. This approach is due to Bjerner [7], where a compositional theory for time analysis of the (primitive recursive) programs of Martin-Löf type-theory is developed. A characterization of ‘need’ (more accurately ‘not-need’) provided by a new form of strictness analysis [43] enabled Wadler to give a simpler account of Bjerner’s approach [42] in the context of a (general) first-order functional language. The strictness-analysis perspective also gives a natural notion of approximation in the description of context information, and gives rise, via *abstract interpretation* to a completely mechanizable analysis for reasoning about (approximate) contexts. In [37, 36] the context information available from such an analysis is used to characterize *sufficient-time* and *necessary-time* equations which together provide bounds on the exact time cost of lazy evaluation, and the method is extended to higher-order functions using a modification of the cost-closure technique.

A problem with these compositional approaches to time analysis remains: the information required about context is itself an uncomputable property in general. The options are to settle either for approximate information via abstract interpretation (or a related approach), or to work with a complete calculus for contexts and hope to find more exact solutions. The former approach, while simplifying the task of reasoning about context (assuming that an implementation is available), can lead to unacceptable approximations in time cost. The latter approach (see [8]) can be impractically cumbersome for many relatively simple problems, and is unlikely to extend usefully to higher-order languages.

The naïve approach

In the following three sections, we explore a complementary approach which begins with a more direct operational viewpoint. We define a small first-order lazy functional language with lists, and define time cost in terms of an operational model. (The treatment of a higher-order language in the naïve approach is just as straightforward, but is postponed in order to simplify the exposition of the theory of cost equivalence.) The simplicity of the chosen semantics (a substitution-based call-by-name model) leads to a correspondingly straightforward definition of time cost, which is refined to give an unsophisticated calculus, in the form of *time rules* with which we can analyse time cost. We illustrate the utility of the naïve approach before going on to consider extensions and improvements.

3 A simple operational model

We initially consider a first-order language with lists. For simplicity we present an untyped semantics, but the syntax will be suggestive of a typed version. List construction is sugared with an infix *cons* ‘:’, and lists are examined and decomposed via a case-expression. Programs are closed expressions in the context of function definitions

$$f_i(x_1, \dots, x_{n_i}) = e_i.$$

We also assume some strict primitive functions over the atomic constants of the language (booleans, integers etc.). Expressions are described by the grammar in Fig. 1.

$e ::=$	$f(e_1, \dots, e_n)$	(function call)
	$p(e_1, \dots, e_n)$	(primitive function call)
	if e_1 then e_2 else e_3	(conditional)
	$\left(\begin{array}{l} \text{case } e_1 \text{ of} \\ \quad \text{nil} \Rightarrow e_2 \\ \quad x : xs \Rightarrow e_3 \end{array} \right)$	(list-case expression)
	$e_1 : e_2$	(cons)
	x	(identifier)
	c	(constant)

FIG. 1. Expression syntax

3.1 Semantic rules

It is possible to reason about time-complexity of a closed expression by reasoning directly about the ‘steps’ in the evaluation of an expression. The problem with this approach is that it requires us to have the machinery of an operational semantics at our fingertips in order to reason in a formal manner. The degree of operational reasoning necessary can be minimized by an appropriately abstract choice of semantics. In particular, simplicity motivates the choice of a call-by-name calling mechanism—shortcomings and improvements to this model are discussed in Section 10. The semantics is defined via *two* types of evaluation rule: one describing evaluation to head-normal-form, and one for evaluation to normal-form. Including rules for evaluation to

normal-form is somewhat non-standard for the semantics of a lazy language. To talk about the complete evaluation of programs it is usual to define a *print-loop* to describe accurately the top-level behaviour of a program. Since our motivation is the analysis of time cost, the rules for evaluation to normal-form give a convenient approximation to the printing mechanism (since in the case of non-terminating programs we would need to place them in some ‘terminating context’ to describe their time behaviour anyway).

Fun	$\frac{e_i\{e_1/x_1 \cdots e_{n_i}/x_{n_i}\} \rightarrow_\alpha u}{f_i(e_1, \dots, e_{n_i}) \rightarrow_\alpha u}$
Prim	$\frac{e_1 \rightarrow_H c_1 \cdots e_{n_k} \rightarrow_H c_k \quad (v = \mathbf{apply}_p(p_i, c_1, \dots, c_{n_i}))}{p_i(e_1, \dots, e_{n_k}) \rightarrow_\alpha v}$
Cond	$\frac{e_1 \rightarrow_H \mathbf{true} \quad e_2 \rightarrow_\alpha u}{\text{if } e_1 \text{ then } e_2 \text{ else } e_3 \rightarrow_\alpha u} \quad \frac{e_1 \rightarrow_H \mathbf{false} \quad e_3 \rightarrow_\alpha u}{\text{if } e_1 \text{ then } e_2 \text{ else } e_3 \rightarrow_\alpha u}$
Cons	$\frac{e_1 \rightarrow_N v_1, e_2 \rightarrow_N v_2}{e_1 : e_2 \rightarrow_N v_1 : v_2} \quad \frac{}{e_1 : e_2 \rightarrow_H e_1 : e_2}$
Const	$\frac{}{c \rightarrow_\alpha c}$
Case	$\frac{e_1 \rightarrow_H \mathbf{nil} \quad e_2 \rightarrow_\alpha u}{\left(\begin{array}{l} \text{case } e_1 \text{ of} \\ \quad \mathbf{nil} \Rightarrow e_2 \\ \quad x : xs \Rightarrow e_3 \end{array} \right) \rightarrow_\alpha u} \quad \frac{e_1 \rightarrow_H e_h : e_t \quad e_3\{e_h/x, e_t/xs\} \rightarrow_\alpha u}{\left(\begin{array}{l} \text{case } e_1 \text{ of} \\ \quad \mathbf{nil} \Rightarrow e_2 \\ \quad x : xs \Rightarrow e_3 \end{array} \right) \rightarrow_\alpha u}$

FIG. 2. Dynamic semantics

We define the operational semantics via rules which allow us to make judgements of the form: $e \rightarrow_N v$ and $e \rightarrow_H h$. These can be read as ‘expression e evaluates to normal form v ’ and ‘expression e evaluates to head-normal form² h ’ respectively. There is no rule for evaluating a variable—evaluation is only defined over closed expressions. These rules are presented in Fig. 2, using meta-variable α to range over labels H and N . Normal-forms, ranged over by v, v_1, v_2 etc. (sometimes referred to simply as *values*) are the fully evaluated expressions i.e. either fully evaluated lists, or atomic constants (c):

$$v ::= c \mid v_1 : v_2.$$

Head-normal forms, ranged over by h, h_1, h_2 etc., are simply the constants and arbitrary expressions:

$$h ::= c \mid e_1 : e_2.$$

A brief explanation of the semantic rules is given below:

²The use of the term *head-normal-form* is not to be confused with the corresponding notion in the (pure) lambda calculus; we use this term as a first-order manifestation of the notion of *weak head-normal-form* from the terminology of lazy functional languages [31].

- To describe function application we perform direct substitution of parameters. We use the notation $e\{e'/x\}$ to mean expression e with free occurrences of x replaced by the expression e' (see comments below).
- We assume the primitive functions are strict functions on constants, and are given meaning by some partial function **apply** _{p} . Since the constants are included in the head-normal forms it is sufficient to evaluate the arguments to primitive functions with \rightarrow_H . In a lower-level semantics in which errors (eg. ‘divide-by-zero’) are distinguished from non-termination this choice would be significant, but here it makes no difference.
- To evaluate a case-expression either to normal or head-normal-form, we must evaluate the list-expression e_1 to determine which branch to take. However, we do not evaluate the expression any further than the first cons-node.

Notation

We summarize some of the notation used in the remainder of the paper.

Variables and substitution A list of zero or more variables x_1, \dots, x_n will often be denoted \vec{x} , and similarly for a list of expressions.

We use the notation $e\{e_1, \dots, e_n/x_1, \dots, x_n\}$ to mean expression e with free occurrences of x_1, \dots, x_n simultaneously replaced by the expressions e_1, \dots, e_n . In a case expression

$$\begin{aligned} &\text{case } e_1 \text{ of} \\ &\quad \text{nil} \Rightarrow e_2 \\ &\quad x : xs \Rightarrow e_3 \end{aligned}$$

the variables x and xs are considered bound in e_3 . A formal definition of substitution is omitted, but is standard (see e.g. [4]). We will also assume the following substitution property (the standard ‘substitution lemma’): if variables \vec{x} and \vec{y} are distinct, then

$$p\{\vec{q}/\vec{x}\}\{\vec{r}/\vec{y}\} = p\{\vec{r}/\vec{y}\}\{\vec{q}\{\vec{r}/\vec{y}\}/\vec{x}\}$$

where $\vec{q}\{\vec{r}/\vec{y}\} = q_1\{\vec{r}/\vec{y}\}, \dots, q_n\{\vec{r}/\vec{y}\}$

The idea of a *context*, ranged over by C, C_1 , etc. will be used (informally) to denote an expression with a ‘hole’, $[]$, in the place of a subexpression; $C[e]$ is the expression produced by replacing the hole with expression e . Generally we will assume that the expression is closed, and so this notation can be considered shorthand for substitution.³

Relations If R is a relation, then we will usually write $a R b$ to mean $(a, b) \in R$. A relation R is a *preorder* if it is transitive and reflexive, and an *equivalence relation* if it is also symmetric. Syntactic equivalence up to renaming of bound variables will be denoted \equiv . The maximum relation on closed expression will be denoted by \top .

4 Deriving time-rules

We wish to reason about the time cost of evaluating an expression. For simplicity we express this property in terms of the number of non-primitive function calls occurring in the evaluation of the expression.

³ Formal definitions usually allow free variables in e to be captured by C . We will revert to the substitution notation when we need to be more formal, and consider the special case of variable capture explicitly.

8 A Naïve Time Analysis and its Theory of Cost Equivalence

For the operational semantics given, the evaluation process is understood in terms of (the construction of) a proof of some judgement according to the semantic rules. The above property of an evaluation corresponds to the number of instances of the rule **Fun** in the proof of $e \rightarrow_\alpha u$ for some closed expression e , whenever such a proof exists for some u . In order to extract rules for reasoning about this property, we rely on some basic properties of the semantics:

- The rules describe deterministic computation: if $e \rightarrow_\alpha u$ and $e \rightarrow_\alpha u'$ then $u \equiv u'$.
- Proofs are unique: if Δ and Δ' are proofs of $e \rightarrow_\alpha u$ then Δ and Δ' are identical.

In the following let $S, S_1 \dots$ range over judgements of the form $e \rightarrow_\alpha u$, and let $\Delta, \Delta_1 \dots$ range over proofs of judgements.

DEFINITION 4.1

Let $T(\Delta)$ be the number of instances of rule **Fun** in a given proof Δ .

Since all proofs are finite,⁴ assuming the inferences are labeled, we can define T inductively in the structure of the proof, according to the last rule applied:

$$T\left(\frac{S}{\Delta_1, \dots, \Delta_k} r\right) = \begin{cases} 1 + T(\Delta_1) + \dots + T(\Delta_k) & \text{if } r = \mathbf{Fun} \\ T(\Delta_1) + \dots + T(\Delta_k) & \text{otherwise.} \end{cases}$$

To define equations for reasoning about time we can abstract away from the structure of the proof, and express this property in terms of the structure of *expressions*, since the last rule used in the proof of some judgement S is determined largely by the expression syntax. Using this principle we define equations for $\langle e \rangle^N$, the time to compute the normal-form, and $\langle e \rangle^H$, the time to compute head-normal-form of expression e . The rules for $\langle \rangle^N$ and $\langle \rangle^H$ are given in Fig. 3.

When we write $\langle e \rangle^\alpha = M$ we mean that this is provable from the time rules, together with standard arithmetic identities. Since we do not include an axiomatization of integers and their addition, this statement is not completely formal, but will be sufficient for non-automated reasoning.

The rules are adequate in the following sense:

PROPOSITION 4.2

For all expressions e , if Δ is a proof of $e \rightarrow_\alpha u$, for some u , then

$$T(\Delta) = n \iff \langle e \rangle^\alpha = n$$

The proof of this proposition is a straightforward induction in the structure of Δ . Notice that the premiss of the proposition makes a termination assumption about the evaluation of e . In this sense the time rules are partially correct. We could then refine this correctness statement further by treating run-time errors separately from non-termination.

5 Direct time analysis

The time rules in Fig. 3, although simple, are sufficient to reason directly about the cost of evaluating closed expressions in this language. We illustrate the utility of this approach with

⁴Proofs correspond to terminating computations; our calculus will therefore allow us only to conclude time-properties under a termination assumption. For a further discussion of this point see [36].

$$\begin{aligned}
 \langle f_i(e_1, \dots, e_{n_i}) \rangle^\alpha &= 1 + \langle e_i \{e_1/x_1, \dots, e_{n_i}/x_{n_i}\} \rangle^\alpha \\
 \langle p(e_1, \dots, e_k) \rangle^\alpha &= \langle e_1 \rangle^H + \dots + \langle e_k \rangle^H \\
 \langle \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \rangle^\alpha &= \langle e_1 \rangle^H + \begin{cases} \langle e_2 \rangle^\alpha & \text{if } e_1 \rightarrow_H \text{ true} \\ \langle e_3 \rangle^\alpha & \text{if } e_1 \rightarrow_H \text{ false} \end{cases} \\
 \left\langle \begin{array}{l} \text{case } e_1 \text{ of} \\ \text{nil} \Rightarrow e_2 \\ x : xs \Rightarrow e_3 \end{array} \right\rangle^\alpha &= \langle e_1 \rangle^H + \begin{cases} \langle e_2 \rangle^\alpha & \text{if } e_1 \rightarrow_H \text{ nil} \\ \langle e_3 \{e_h/x, e_t/xs\} \rangle^\alpha & \text{if } e_1 \rightarrow_H e_h : e_t \end{cases} \\
 \langle e_1 : e_2 \rangle^N &= \langle e_1 \rangle^N + \langle e_2 \rangle^N \\
 \langle e_1 : e_2 \rangle^H &= \langle e \rangle^\alpha = 0
 \end{aligned}$$

FIG. 3. Time rules

some small examples. The examples are chosen to emphasize features of non-strict evaluation, rather than to present interesting asymptotic analyses, for which the reader may consult a standard text on the analysis of algorithms (e.g. [20, 2]) or associated mathematical techniques (e.g. [15]).

To reason about complexity we consider expressions containing some non-specified *input value* (i.e. normal-form), which we will denote by a (meta-)variable (written in an italic font). We sometimes also allow meta-variables to range over arbitrary expressions, although usually this is more awkward since the calculus is not compositional.

5.1 Example

Consider the functions over lists given in Fig. 4.

$$\begin{aligned}
 \text{append}(xs,ys) &= \text{case } xs \text{ of} \\
 &\quad \text{nil} \Rightarrow ys \\
 &\quad h:t \Rightarrow h:\text{append}(t,ys) \\
 \text{reverse}(xs) &= \text{case } xs \text{ of} \\
 &\quad \text{nil} \Rightarrow \text{nil} \\
 &\quad h:t \Rightarrow \text{append}(\text{reverse}(t), h:\text{nil}) \\
 \text{head}(xs) &= \text{case } xs \text{ of} \\
 &\quad \text{nil} \Rightarrow \text{undefined} \\
 &\quad h:t \Rightarrow h
 \end{aligned}$$

FIG. 4. Some list-manipulating functions

Now we wish to consider the cost of evaluating the expression

$$\text{head}(\text{reverse}(v))$$

which computes the last element of some non-empty list-value $v \equiv v_h : v_t$. Applying the definitions in Fig. 3

$$\langle \text{head}(\text{reverse}(v)) \rangle^N = 1 + \langle \text{reverse}(v) \rangle^H + \langle e_h \rangle^N$$

where $\text{reverse}(v) \rightarrow_H e_h : e_t$ for some e_h, e_t . It is not hard to show that e_h is a value (using the fact that v is, and by induction on its length), and hence that $\langle e_h \rangle^N = 0$. Now since $v \equiv v_h : v_t$, and hence $v \rightarrow_H v_h : v_t$ we have that

$$\begin{aligned} & \langle \text{reverse}(v) \rangle^H \\ &= 1 + \langle v \rangle^H + \langle \text{append}(\text{reverse}(v_t), v_h : \text{nil}) \rangle^H \\ &= 1 + \langle \text{append}(\text{reverse}(v_t), v_h : \text{nil}) \rangle^H \\ &= 1 + 1 + \langle \text{reverse}(v_t) \rangle^H \\ & \quad + \begin{cases} \langle \text{nil} \rangle^H, & \text{if } \text{reverse}(v_t) \rightarrow_H \text{nil} \\ \langle e'_h : \text{append}(e'_t, v_h : \text{nil}) \rangle^H, & \\ \text{if } \text{reverse}(v_t) \rightarrow_H e'_h : e'_t \end{cases} \\ &= 2 + \langle \text{reverse}(v_t) \rangle^H. \end{aligned}$$

We now have the recurrence equations parameterized by the input value:

$$\begin{aligned} \langle \text{reverse}(\text{nil}) \rangle^H &= 1 \\ \langle \text{reverse}(v : vs) \rangle^H &= 2 + \langle \text{reverse}(vs) \rangle^H \end{aligned}$$

whose solution is $\langle \text{reverse}(v) \rangle^H = 1 + 2n$, where n is the length of the list v . Thus we have a total cost of

$$\langle \text{head}(\text{reverse}(v)) \rangle^N = 2(1 + n)$$

where n is the length of the list v , i.e. linear time complexity (compare with quadratic complexity for the call-by-value reading, and for $\langle \text{reverse}(v) \rangle^N$).

5.2 Example

Here we present the example from the introduction (in the syntax we have defined) which shows the rather pleasing property that one can compute the smallest element of a list in linear time by taking the first element of the (insertion-) sort of the list. The equations in Fig. 5 define an insertion-sort function (`isort`).

The time to compute the head-normal-form of insertion-sort given some list-value (normal-form) v is easily calculated from the time rules. First consider the insertion function. Any exhaustive application of the time rules together with a few minor simplifications allow us to conclude that for integer valued expressions e_1 , and integer-list valued expressions e_2 ,

$$\begin{aligned} \langle \text{insert}(e_1, e_2) \rangle^H &= 1 + \langle e_2 \rangle^H \\ &+ \begin{cases} 0 & \text{if } e_2 \rightarrow_H \text{nil} \\ \langle e_1 \rangle^H + \langle h \rangle^H & \text{if } e_2 \rightarrow_H h : t \end{cases} \end{aligned}$$

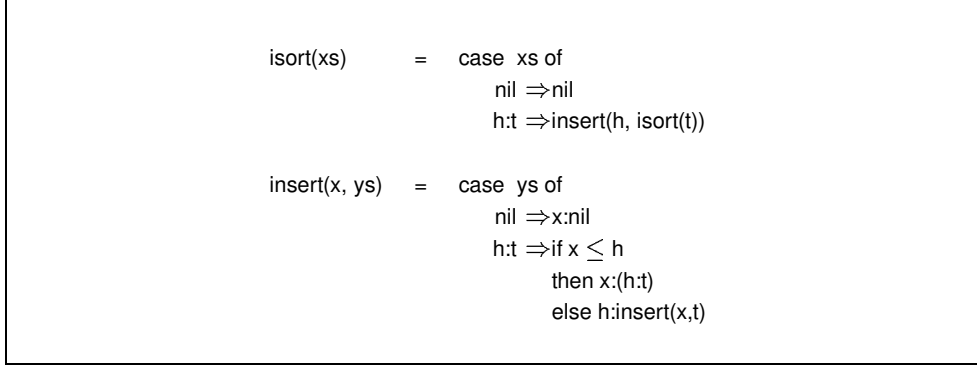


FIG. 5. Insertion sort

Now consider computing the head-normal-form of insertion-sort applied to some list of integers $v_n : \dots : v_1 : \text{nil}$ where $n \geq 1$. To aid notation, let V_0 denote the list nil and, for each $i < n$, let V_{i+1} denote the list $v_{i+1} : V_i$.

$$\begin{aligned}
 \langle \text{isort}(V_0) \rangle^H &= 1 \\
 \langle \text{isort}(V_{i+i}) \rangle^H &= 1 + \langle \text{insert}(v_{i+1}, \text{isort}(V_i)) \rangle^H \\
 &= 1 + 1 + \langle \text{isort}(V_i) \rangle^H \\
 &+ \begin{cases} 0 & \text{if } \text{isort}(V_i) \rightarrow_H \text{ nil} \\ \langle v_{i+1} \rangle^H + \langle h \rangle^H & \text{if } \text{isort}(V_i) \rightarrow_H h : t \end{cases} .
 \end{aligned}$$

Clearly $\langle v_{i+1} \rangle^H = 0$. A simple induction in i establishes that if $\text{isort}(V_i) \rightarrow_H h : t$ then $\langle h \rangle^H = 0$ also. This leaves us with the simple recurrence

$$\begin{aligned}
 \langle \text{isort}(V_0) \rangle^H &= 1 \\
 \langle \text{isort}(V_{i+i}) \rangle^H &= 2 + \langle \text{isort}(V_i) \rangle^H
 \end{aligned}$$

giving $\langle \text{isort}(V_n) \rangle^H = 2n + 1$.

5.3 Example

Consider the following (somewhat non-standard⁵) definition of Fibonacci:

$$\begin{aligned}
 \text{fib}(n) &= \text{f}(n,0) \\
 \text{f}(n, r) &= \text{if } n=0 \text{ then } 1 \\
 &\quad \text{else } r + \text{f}(n-1, \text{f}(n-2,0))
 \end{aligned}$$

⁵Note that under a call-by-value semantics fib is divergent for any $n > 0$.

Consider the time to compute an instance of fib:

$$\begin{aligned}
\langle \text{fib}(k) \rangle^N &= 1 + \langle f(k, 0) \rangle^N \\
\langle f(k, 0) \rangle^N &= 1 + \left\langle \begin{array}{l} \text{if } k=0 \text{ then } 1 \\ \text{else } 0 + f(k-1, f(k-2, 0)) \end{array} \right\rangle^N \\
&= 1 + \begin{cases} \langle 1 \rangle^N & \text{if } k=0 \rightarrow_N \text{ true} \\ \langle 0 + f(k-1, f(k-2, 0)) \rangle^N & \text{if } k=0 \rightarrow_N \text{ false} \end{cases} \\
&= 1 + \begin{cases} 0 & \text{if } k=0 \rightarrow_N \text{ true} \\ \langle f(k-1, f(k-2, 0)) \rangle^N & \text{if } k=0 \rightarrow_N \text{ false} \end{cases} .
\end{aligned}$$

Instantiating k we have

$$\langle f(0, 0) \rangle^N = 1 \tag{5.1}$$

$$\begin{aligned}
\langle f(1, 0) \rangle^N &= 1 + \langle f(0, f(1-2, 0)) \rangle^N \\
&= 2 \tag{5.2}
\end{aligned}$$

$$\begin{aligned}
\langle f(k+2, 0) \rangle^N &= 1 + \langle f(k+1, f(k, 0)) \rangle^N \\
&= 1 + 1 + \langle f(k, 0) \rangle^N \tag{5.3}
\end{aligned}$$

$$+ \langle f(k, f(k-1, 0)) \rangle^N . \tag{5.4}$$

Now $\langle f(k+1, 0) \rangle^N = 1 + \langle f(k, f(k-1, 0)) \rangle^N$, so

$$\langle f(k+2, 0) \rangle^N = 1 + \begin{array}{l} \langle f(k, 0) \rangle^N \\ + \langle f(k+1, 0) \rangle^N . \end{array} \tag{5.5}$$

Equations (5.1), (5.2) and (5.5) give a linear recurrence-relation that can be solved (exactly) using standard techniques e.g., [15]; the asymptote for $\langle f(k, 0) \rangle^N$ has the form:

$$a(1 + \sqrt{5})/2)^k + b(1 - \sqrt{5})/2)^k$$

for some constants a and b .

6 A theory of cost-equivalence

6.1 Motivation

The previous example subtly illustrates some potential problems in reasoning about cost using the equations for $\langle \cdot \rangle^H$ and $\langle \cdot \rangle^N$. The use of the time rules in the previous examples follows a simple pattern of case analysis (instantiation) and simplification, leading to the construction of a recurrence by a simple syntactic matching. In the simplification process, it is tempting to make simplifications which are not directly justifiable from the time rules.

The potential problems stem from the fact that if we know that two expressions are extensionally equivalent, $e_1 = e_2$, it is clearly *not* the case that $\langle e_1 \rangle^N = \langle e_2 \rangle^N$ in general since we expect, with any reasonable definition of extensional equivalence, that an expression and its normal-form

(supposing one exists) will be equivalent. More generally given any context C , we expect that $C[e_1] = C[e_2]$, but **not** that $\langle C[e_1] \rangle^\alpha = \langle C[e_2] \rangle^\alpha$ in general, so ordinary equational reasoning is not valid within the time rules. Similarly if $\langle e_1 \rangle^H = \langle e_2 \rangle^H$ then we cannot expect in general that $\langle C[e_1] \rangle^H = \langle C[e_2] \rangle^H$.

However, even in the last example above we *have* used simple equalities such as (line 5.4) $(k + 1) - 1 = k$ in precisely this way (albeit benignly) to simplify cost expressions in order to construct a recurrence.⁶ In this instance the simplification is obviously correct, but with the current calculus we cannot justify it.

To take a less contrived example, where the limitations of the method are more significant, consider the quicksort program given in the introduction. The following equations (Fig. 6) define a simple functional version of quicksort (qs) using auxiliary functions below and above, and append (as defined earlier) written here as an infix function $++$. Primitive functions for integer comparison have also been written infix to aid readability. The definition for above has been omitted, but is like that of below with the comparison ‘>’ in place of ‘≤’.

qs(xs)	=	case xs of nil \Rightarrow nil h:t \Rightarrow qs(below(h,t)) $++$ (h : qs(above(h,t)))
below(x,ys)	=	case ys of nil \Rightarrow nil h:t \Rightarrow if $h \leq x$ then h : below(x,t) else below(x,t)

FIG. 6. Functional quicksort

The aim will be fairly modest: to show that quicksort exhibits its worst-case $O(n^2)$ behaviour even when we only require the first element of the list to be computed, in contrast to the earlier insertion sort example which always takes linear time to compute the first element of the result. First consider the general case:

$$\langle \text{qs}(e) \rangle^H = 1 + \langle e \rangle^H + \begin{cases} 0 & \text{if } e \rightarrow_H \text{ nil} \\ \left\langle \begin{array}{l} \text{qs}(\text{below}(y, z)) \\ ++(y:\text{qs}(\text{above}(y, z))) \end{array} \right\rangle^H & \text{if } e \rightarrow_H y:z \end{cases}$$

From the time rules and the definition of append, this simplifies to

$$\langle \text{qs}(e) \rangle^H = 1 + \langle e \rangle^H + \begin{cases} 0 & \text{if } e \rightarrow_H \text{ nil} \\ 1 + \langle \text{qs}(\text{below}(y,z)) \rangle^H & \text{if } e \rightarrow_H y:z \end{cases} \quad (6.1)$$

Proceeding to the particular problem, it is not too surprising that we will use non-increasing lists v to show that $\langle \text{qs}(v) \rangle^H = \Omega(n^2)$. Towards this goal, fix an arbitrary family of integer

⁶Spelling this simplification out, we have an application of a primitive function (subtraction) to the (meta) constant ‘ $k + 1$ ’ (i.e., the constant one larger than the meta-constant k), which we simplify to ‘ k ’.

14 A Naïve Time Analysis and its Theory of Cost Equivalence

values $\{v_i\}_{i>0}$ such that $v_i \leq v_j$ whenever $i \leq j$. Now define the family of non-increasing lists $\{A_i\}_{i \geq 0}$ by induction on i :

$$\begin{aligned} A_0 &= \text{nil} \\ A_{k+1} &= v_{k+1} : A_k. \end{aligned}$$

The goal is now to show that $\langle \text{qs}(A_n) \rangle^H$ is quadratic in n . Instantiating the general equation (6.1) in the ‘interesting’ case when $e = A_{k+2}$ for some $k \geq 0$ we obtain:

$$\begin{aligned} &\langle \text{qs}(A_{k+2}) \rangle^H \\ &= 1 + \langle A_{k+2} \rangle^H + 1 + \langle \text{qs}(\text{below}(y,z)) \rangle^H \\ &\quad \text{where } A_{k+2} \rightarrow_H y:z \\ &= 2 + \langle \text{qs}(\text{below}(v_{k+2}, A_{k+1})) \rangle^H. \end{aligned}$$

At this point we can only further manipulate the expression $\langle \text{qs}(\text{below}(v_{k+2}, A_{k+1})) \rangle^H$:

$$\begin{aligned} &\langle \text{qs}(\text{below}(v_{k+2}, A_{k+1})) \rangle^H \\ &= 1 + \langle \text{below}(v_{k+2}, A_{k+1}) \rangle^H + 1 + \langle \text{qs}(\text{below}(y,z)) \rangle^H \\ &\quad \text{where } \text{below}(v_{k+2}, A_{k+1}) \rightarrow_H y:z \\ &= 3 + \langle \text{qs}(\text{below}(v_{k+1}, \text{below}(v_{k+2}, A_k))) \rangle^H. \end{aligned}$$

So now we have the equation

$$\langle \text{qs}(A_{k+2}) \rangle^H = 5 + \langle \text{qs}(\text{below}(v_{k+1}, \text{below}(v_{k+2}, A_k))) \rangle^H.$$

This should be sufficient to convince the reader that quadratic-time behaviour is a possibility, since in the successive recursive calls to qs (with respect to this time equation) we can see that the arguments become increasingly complex, and it becomes increasingly costly to compute their respective head-normal-forms.

With the current calculus we can do little more than give this intuition. What we are not able to do is to simplify or generalize the calls to below to obtain a simple recurrence equation. We will return to this example.

The remainder of this section is devoted to developing a stronger notion of equivalence of expressions which respects cost, and allows a richer form of equational reasoning on expressions within the calculus. We will conclude the above example in Section 7.3.

What is needed is an appropriate characterization of (the weakest) equivalence relation $=_{\langle \rangle}$ which satisfies

$$e =_{\langle \rangle} e' \implies \langle C[e] \rangle^\alpha = \langle C[e'] \rangle^\alpha.$$

To develop this general ‘contextual congruence’ relation, we use a notion of *simulation* similar to the various simulations developed in process algebras such as Milner’s calculus of communicating systems [26]. In the theory of concurrency a central idea is that processes that cannot be distinguished by observation should be identified. This ‘observational’ viewpoint is adopted in the ‘lazy’ λ -calculus [1], where an equivalence called *applicative bisimulation* is introduced. In the lazy λ -calculus, the observable properties are just the convergence of untyped lambda-terms. For our purposes we need to treat cost as an observable component of the evaluation process, and so we develop a suitable notion of *cost (bi)simulation*.

6.2 Cost simulation

The partial functions \rightarrow_H and \rightarrow_N together with $\langle \rangle^N$ and $\langle \rangle^H$ are not sufficient to characterize completely the cost behaviour of expressions in all contexts, since we need to characterize possibly infinite ‘observations’ on expressions which arise in our language because of the non-strict list-constructor (c.f. untyped weak head-normal-forms in [1]).

Roughly speaking, the notion of equivalence we want satisfies:

- e and e' are equivalent iff $\langle e \rangle^H = \langle e' \rangle^H$ and their head-normal-forms are either*
1. *identical constants, or*
 2. *cons-expressions, whose corresponding components are equivalent.*

Unfortunately, although this is a property that we would like our equivalence to obey, it does not constitute a definition (to see why, note that we not only wish to relate expressions having normal-forms, but also those which are ‘infinite’), so following [26] we use a technique due to Park [30] for identifying processes—the notion of a *bisimulation* and its related proof technique. We will develop the equivalence relation we require in terms of preorders called *cost simulations*—we will then say that two expressions are cost equivalent if they simulate each other.

To simplify our presentation we add some notation:

DEFINITION 6.1

If \mathcal{R} is a binary relation on closed expressions, then \mathcal{R}^\cdot is the binary relation on head-normal-forms such that

$$(h \mathcal{R}^\cdot h') \iff \begin{array}{l} \text{either } h = h' = c \\ \text{or } h = e_1 : e_2, h' = e'_1 : e'_2 \text{ and} \\ e_1 \mathcal{R} e'_1, \text{ and } e_2 \mathcal{R} e'_2. \end{array}$$

DEFINITION 6.2

The cost-labelled transition, $\xrightarrow{\alpha}_t$, $t \in \mathbf{N}$ is defined

$$e \xrightarrow{\alpha}_t u \stackrel{\text{def.}}{=} e \rightarrow_{\alpha} u \text{ and } \langle e \rangle^{\alpha} = t.$$

Now we define a basic notion of cost simulation, by analogy with Park’s (bi)simulation:

DEFINITION 6.3 (Cost Simulation)

A binary relation \mathcal{R} on closed expressions is a **cost simulation** if, whenever $e \mathcal{R} e'$

$$e \xrightarrow{t}_H h \implies (e' \xrightarrow{t}_H h' \text{ and } h \mathcal{R}^\cdot h').$$

DEFINITION 6.4

For each relation Q on closed expressions, define $\mathcal{F}(Q)$ to be the relation on closed expressions that relates e and e' exactly when

$$e \xrightarrow{t}_H h \implies (e' \xrightarrow{t}_H h' \text{ and } h Q^\cdot h').$$

Now we can easily see that

- \mathcal{F} is monotonic, i.e., $\mathcal{R} \subseteq \mathcal{S} \implies \mathcal{F}(\mathcal{R}) \subseteq \mathcal{F}(\mathcal{S})$
- \mathcal{S} is a cost simulation iff $\mathcal{S} \subseteq \mathcal{F}(\mathcal{S})$ (since $\mathcal{S} \subseteq \mathcal{F}(\mathcal{S})$ means that $\forall e_1, e_2, e_1 \mathcal{S} e_2 \implies e_1 \mathcal{F}(\mathcal{S}) e_2$, and by expanding the definition of \mathcal{F} we recover Definition 6.3).

DEFINITION 6.5

Let \preceq denote the maximum cost simulation

$$\bigcup \{ \mathcal{S} : \mathcal{S} \subseteq \mathcal{F}(\mathcal{S}) \}.$$

PROPOSITION 6.6

\preceq is the maximal fixed point of \mathcal{F} .

PROOF. Follows from the Knaster–Tarski fixed point theorem, by the fact that \mathcal{F} is a monotone function on a complete lattice. \blacksquare

With these results we have the following useful proof technique: to show that $e \preceq e'$ it is necessary and sufficient to exhibit *any* cost simulation containing (i.e. relating) the pair (e, e') . This technique will be illustrated later in the proof that \preceq is a precongruence.

6.3 Expressing \rightarrow_N in terms of \rightarrow_H

The above definition of cost-simulation is described in terms of evaluation to head-normal-form only. For this to be sufficient to describe properties of evaluation to normal-form we need some properties relating \rightarrow_N and \rightarrow_H . The following property allows us to factor evaluation to normal-form through evaluation to head normal form, while preserving cost behaviour:

PROPOSITION 6.7

For all closed e, e'

- $e \rightarrow_N c \iff e \rightarrow_H c$.
- $e \xrightarrow[t]{t} v \iff e \xrightarrow[t_1]{t_1} h$ and $h \xrightarrow[t_2]{t_2} v$ and $t_1 + t_2 = t$ for some h .

And finally we have

LEMMA 6.8

If $e \preceq e'$ then if $e \xrightarrow[t]{t} u$ then $e' \xrightarrow[t]{t} u$.

The proofs are outlined in Appendix A.

Remark The (first) implication in the lemma cannot be reversed. For example, if l is an identity function, then the expressions $l(\text{nil}):\text{nil}$ and $\text{nil}:l(\text{nil})$ take the same time to reach identical normal forms but are not cost simulation comparable.

6.4 Precongruence

Now we are ready to prove the key property that we demand of cost simulation: cost-simulation is a precongruence, i.e. it is a substitutive preordering (the fact that \preceq is a preorder is easily established).

Some notation: for convenience we abbreviate some indexed family of expressions $\{e_j : j \in J\}$ by \tilde{e} . Similarly we will abbreviate the substitution $\{e_j/x_j : j \in J\}$ by $\{\tilde{e}/\tilde{x}\}$, and when, for all $j \in J$, $(e_j Q e'_j)$ for some relation Q , we write $(\tilde{e} Q \tilde{e}')$.

DEFINITION 6.9

\mathfrak{R} is defined to be the relation

$$\mathfrak{R} = \{(e\{\tilde{e}/\tilde{x}\}, e'\{\tilde{e}'/\tilde{x}\}) \mid \tilde{x} \subseteq FV(e), \tilde{e} \preceq \tilde{e}'\}.$$

LEMMA 6.10

\mathfrak{R} is a cost simulation.

PROOF. Assume that $\tilde{e} \preceq \tilde{e}'$, for some closed expressions \tilde{e}, \tilde{e}' . Abbreviate substitutions $\{\tilde{e}/\tilde{x}\}$ and $\{\tilde{e}'/\tilde{x}\}$ by σ and σ' , respectively, and assume that e is any expression containing at most variables \tilde{x} . Assume that $e\sigma \xrightarrow{t}_H h$. The lemma now requires us to prove that $e\sigma' \xrightarrow{t}_H h'$ for some h' such that $h \mathfrak{R} h'$. We establish this by induction on the structure of the proof of $e\sigma \xrightarrow{t}_H h$, and by cases according to the structure of expression e . We give a couple of illustrative cases:

$e \equiv x$ Observe that \preceq is contained in \mathfrak{R} , and the result follows.

$e \equiv f(e_1 \dots e_n)$

Assume that f is defined by $f(y_1 \dots y_n) = e_f$. Since $f(e_1 \dots e_n)\sigma \equiv f(e_1\sigma \dots e_n\sigma)$, the last rule in the above inference must be an instance of **Fun**, and so we must have $e_f\{e_1\sigma/y_1 \dots e_n\sigma/y_n\} \xrightarrow{t-1}_H h$. We can take variables in \tilde{x} to be distinct from $y_1 \dots y_n$, and so

$$e_f\{e_1\sigma/y_1 \dots e_n\sigma/y_n\} \equiv (e_f\{e_1/y_1 \dots e_n/y_n\})\sigma.$$

Now since $(e_f\{e_1/y_1 \dots e_n/y_n\})\sigma \xrightarrow{t-1}_H h$ by a smaller proof, the inductive hypothesis gives

$$(e_f\{e_1/y_1 \dots e_n/y_n\})\sigma' \xrightarrow{t-1}_H h' \text{ where } h \mathfrak{R} h'.$$

So by rule **Fun**, together with Definition 4.1 we can conclude that

$$f(e_1\sigma', \dots, e_n\sigma') \xrightarrow{t}_H h'$$

with $h \mathfrak{R} h'$ as required. ■

THEOREM 6.11 (Precongruence)

If $\tilde{e} \preceq \tilde{e}'$ for some commonly indexed families of closed expressions \tilde{e}, \tilde{e}' , then for all expressions e containing at most variables \tilde{x}

$$e\{\tilde{e}/\tilde{x}\} \preceq e\{\tilde{e}'/\tilde{x}\}.$$

PROOF. The relation \mathfrak{R} , given above, is such that $(e\{\tilde{e}/\tilde{x}\}, e\{\tilde{e}'/\tilde{x}\}) \in \mathfrak{R}$ whenever $\tilde{e} \preceq \tilde{e}'$ and e contains at most variables \tilde{x} . Lemma 6.10 establishes that \mathfrak{R} is a cost simulation, i.e. that $\mathfrak{R} \subseteq \preceq$, and so we must also have $(e\{\tilde{e}/\tilde{x}\}, e\{\tilde{e}'/\tilde{x}\}) \in \preceq$. ■

Although in this case we can see that \mathfrak{R} is identically \preceq , we express the proof in this way since it illustrates a general method for establishing cost simulations. Now we can define our notion of cost equivalence to be the equivalence relation:

DEFINITION 6.12 (cost equivalence)

$(=_{\text{()}}) \equiv (\preceq \cap \preceq^{-1})$, i.e.

$$(e_1 =_{\text{()}} e_2) \iff (e_1 \preceq e_2) \ \& \ (e_2 \preceq e_1)$$

So two expressions are *cost equivalent* if they cost simulate each other. Now we have as the main corollary of precongruence

COROLLARY 6.13

For all contexts C , and closed expressions e and e' , if $e =_{\emptyset} e'$, then

$$\langle C[e] \rangle^\alpha = \langle C[e'] \rangle^\alpha$$

whenever $C[e] \Downarrow_\alpha$

PROOF. $\alpha = H$, immediate from Theorem 6.11 and the definition of \preceq , viewing a context as an expression containing a single free variable; $\alpha = N$, immediate from Theorem 6.11 and Lemma 6.8 ■

Open expressions In the obvious way we can extend cost simulation to open expressions by saying that $e \preceq e'$ if for all closing substitutions σ , $e\sigma \preceq e'\sigma$. As a consequence we can show that, on open expressions,

$$e \preceq e' \ \& \ e_1 \preceq e_2 \Rightarrow e\{e_1/x\} \preceq e'\{e_2/x\}.$$

We can extend the congruence property to open expressions (where free variables may be captured by the context) by showing that, for any expressions e_1 and e_2 containing at most free variables x and xs , and for any closed expressions e and e' such that $e_1 \preceq e_2$

$$\left(\begin{array}{c} \text{case } e \text{ of} \\ \text{nil} \Rightarrow e' \\ x : xs \Rightarrow e_1 \end{array} \right) \preceq \left(\begin{array}{c} \text{case } e \text{ of} \\ \text{nil} \Rightarrow e' \\ x : xs \Rightarrow e_2 \end{array} \right).$$

Open endedness A statement of cost equivalence involving some function symbol f naturally assumes a particular defining equation, so strictly speaking, cost equivalence should be parameterized by a set of function definitions. The semantic rule for function application is really a rule schema, but in the proof that cost simulation is a precongruence, it is not necessary to assume a particular set of definitions. As a result, adding a new function definition (i.e., a defining equation for a new function name) does not invalidate earlier cost equivalences; furthermore, the maximality results of the next section imply that such an extension of the language must be conservative with respect to cost equivalence.

6.5 Cost simulation as the largest contextual cost congruence

We have shown that cost simulation is a precongruence, which was sufficient for it to be substitutive with respect to the time rules. A remaining question is whether it is the largest possible precongruence with respect to the time rules. i.e. are there expression pairs e, e' such that in all contexts C , $\langle C[e] \rangle^\alpha = \langle C[e'] \rangle^\alpha$ but for which $e \not\preceq e'$?

In this section we outline the result that, under a mild condition on the constructs of the language (outlined below), \preceq is indeed the largest such relation, i.e.

$$(\forall C. C[e] \Downarrow_\alpha \Rightarrow \langle C[e] \rangle^\alpha = \langle C[e'] \rangle^\alpha) \Rightarrow e \preceq e'.$$

Roughly speaking, we say that two expressions e, e' are *cost distinguishable* whenever there exists a context C such that $\langle C[e] \rangle^\alpha \neq \langle C[e'] \rangle^\alpha$. A necessary (and, as we will show, sufficient) condition for the above implication to hold is that every pair of distinct constants in the language are cost distinguishable. We refer to this as the *CD condition*. The CD condition is not a particularly strong one since it is satisfied if, for example, we assume a primitive function providing an equality test over the constants.

We prove the above result by exploring the relationship between cost simulation and various contextual congruences. We summarize these results below, and refer the reader to Appendix B for more details.

- We define a cost congruence preorder \leq_c between closed expressions such that $e_1 \leq_c e_2$ if and only if for all contexts C , $C[e_1] \mathcal{F}(\top) C[e_2]$. i.e. if the results of evaluation to head-normal-form are produced in the same number of steps, and the results have the same ‘outermost’ form (see Definition 6.4).
- We show that $\preceq = \leq_c$.
- We define a pure cost congruence preorder \leq_{pc} which does not take into account the actual head-normal-forms produced: $e_1 \leq_{pc} e_2$ if and only if for all contexts C , whenever $C[e_1] \xrightarrow{t}_H u_1$ then there exists u_2 such that $C[e_2] \xrightarrow{t}_H u_2$.
- Assuming the CD condition, we show that $\leq_c = \leq_{pc}$. As a corollary, we have that $(\forall C. C[e_1] \downarrow_H \Rightarrow \langle C[e_1] \rangle^H = \langle C[e_2] \rangle^H) \Rightarrow e \preceq e'$. The extension of this to include evaluation to normal-form is straightforward.

7 Proof principles and an axiomatization of cost equivalence

The definition of cost simulation comes with a useful proof technique for establishing instances. In the first part of this section we outline some simple variations of this technique.

In the second part of this section we show that cost equivalence subsumes the time rules (i.e. can be viewed as the basis of a time calculus independently) by giving a complete axiomatization of cost equivalence with respect to the cost-labelled transition \xrightarrow{t}_H .

7.1 Co-induction principles

We motivated the theory of cost equivalence with a need for substitutive laws (i.e. cost-equivalence schemas) with which to augment the time rules. Some example laws are given below:

PROPOSITION 7.1

$$(i) \quad p(c_1, \dots, c_n) = () \quad c \text{ if } \mathbf{apply}(p, c_1, \dots, c_n) = c;$$

$$(ii) \quad (e_1 + e_2) + e_3 = () \quad e_1 + (e_2 + e_3);$$

$$(iii) \quad \text{case} \left(\begin{array}{l} \text{case } e_0 \text{ of} \\ \quad \text{nil} \Rightarrow e_1 \\ \quad y : ys \Rightarrow e_2 \end{array} \right) \text{ of} \\ \quad \text{nil} \Rightarrow e_3 \\ \quad x : xs \Rightarrow e_4 \\ = () \\ \text{case } e_0 \text{ of} \\ \quad \text{nil} \Rightarrow (\text{case } e_1 \text{ of} \\ \quad \quad \text{nil} \Rightarrow e_3 \\ \quad \quad x : xs \Rightarrow e_4) \\ \quad y : ys \Rightarrow (\text{case } e_2 \text{ of} \\ \quad \quad \text{nil} \Rightarrow e_3 \\ \quad \quad x : xs \Rightarrow e_4)$$

The proof of Theorem 6.11 illustrates a general technique for establishing cost-equivalence laws such as the above, where we construct a suitable relation (i.e. containing all instances of the law), and show that it is a cost simulation. Recall from the previous section the functional $\mathcal{F}(_)$: for each relation Q on closed expressions,

$$e \mathcal{F}(Q) e' \iff e \xrightarrow{t}_H h \implies (e' \xrightarrow{t}_H h' \text{ and } h Q h')$$

The definition of \preceq as the maximal fixed point of $\mathcal{F}(_)$ comes with the following useful proof technique, which following [27] we call *co-induction*:

To show that $e \preceq e'$ it is necessary and sufficient to exhibit *any* relation R containing (i.e. relating) the pair (e, e') and such that R is a cost simulation (i.e. $R \subseteq \mathcal{F}(R)$).

Some minor variations of this technique also turn out to be useful.

PROPOSITION 7.2

To prove R is a cost simulation, it is sufficient to prove either of the following conditions (*cost simulation modulo S* , and *cost simulation up to cost equivalence* respectively):

1. $R \subseteq \mathcal{F}(R \cup S)$ for some cost simulation S .
2. $R \subseteq \mathcal{F}(=_{\emptyset}; R; =_{\emptyset})$.

PROOF. 1. $R \subseteq \mathcal{F}(R \cup S)$ implies

$$\begin{aligned} R \cup \preceq &\subseteq \mathcal{F}(R \cup S) \cup \preceq \\ &= \mathcal{F}(R \cup S) \cup \mathcal{F}(\preceq) \quad (\preceq \text{ is a f. p.}) \\ &\subseteq \mathcal{F}(R \cup S \cup \preceq) \quad (\text{monotonicity}) \\ &= \mathcal{F}(R \cup \preceq) \quad (S \subseteq \preceq) \end{aligned}$$

which implies that $(R \cup \preceq) \subseteq \preceq$ and hence $R \subseteq \preceq$.

2. For arbitrary relations A and B it is not hard to show that

$$\mathcal{F}(A); \mathcal{F}(B) \subseteq \mathcal{F}(A; B).$$

Using the fact that $=_{\emptyset}$ is transitive, and a fixed point of $\mathcal{F}(_)$ we can show that

$$(=_{\emptyset}; \mathcal{F}(=_{\emptyset}; R; =_{\emptyset}); =_{\emptyset}) \subseteq \mathcal{F}(=_{\emptyset}; R; =_{\emptyset}).$$

Now $R \subseteq \mathcal{F}(=_{\emptyset}; R; =_{\emptyset})$ implies

$$\begin{aligned} =_{\emptyset}; R; =_{\emptyset} &\subseteq =_{\emptyset}; \mathcal{F}(=_{\emptyset}; R; =_{\emptyset}); =_{\emptyset} \\ &\subseteq \mathcal{F}(=_{\emptyset}; R; =_{\emptyset}). \end{aligned}$$

Hence $(=_{\emptyset}; R; =_{\emptyset}) \subseteq \preceq$, and since $=_{\emptyset}$ is greater than the identity relation, $R \subseteq \preceq$. ■

Method (i) is typically used with S taken to be the relation of syntactic equivalence. Method (ii) can be viewed as a ‘semantic’ co-induction principle. For example, part (iii) of Proposition 7.1 is proved by showing that the relation containing all instances is a cost simulation modulo syntactic equivalence. This goes through by a simple case analysis on the possible outcomes of the conditionals, and is left as an exercise.

7.2 An axiomatization of cost equivalence

The language is too expressive to expect a complete set of cost equivalence laws. However we can give a set which is complete with respect to the time rules in a sense that we will make precise below.

The key to this axiomatization is the use of an identity function to represent a single ‘tick’ of computation time.

DEFINITION 7.3

Let l be an identity function given by a program definition $l(x) = x$. For any integer $n \geq 0$, write $l^n(exp)$ for the expression given by n applications of the function l to exp :

$$l(\underbrace{\dots l}_{n}(exp)\dots).$$

We will write $l^1(exp)$ as simply $l(exp)$.

In Fig. 7 we state a set K of cost-equivalence laws. We write

Fun.l	$f_i(e_1, \dots, e_{n_i}) = ()$	$l(e_i\{e_1/x_1 \dots e_{n_i}/x_{n_i}\})$
Prim.l	$p_i(e_1, \dots, l(e_j), \dots, e_{n_k}) = ()$	$l(p_i(e_1, \dots, e_j, \dots, e_{n_k}))$
Prim	$p_i(c_1, \dots, c_{n_k}) = ()$	v if $v = \mathbf{apply}_p(p_i, c_1, \dots, c_{n_i})$
Cond.l	if $l(e_1)$ then e_2 else $e_3 = ()$	$l(\text{if } e_1 \text{ then } e_2 \text{ else } e_3)$
Cond.true	if true then e_2 else $e_3 = ()$	e_2
Cond.false	if false then e_2 else $e_3 = ()$	e_3
Case.l	$\left(\begin{array}{l} \text{case } l(e_1) \text{ of} \\ \quad \text{nil} \Rightarrow e_2 \\ \quad x : xs \Rightarrow e_3 \end{array} \right) = ()$	$l \left(\begin{array}{l} \text{case } e_1 \text{ of} \\ \quad \text{nil} \Rightarrow e_2 \\ \quad x : xs \Rightarrow e_3 \end{array} \right)$
Case.nil	$\left(\begin{array}{l} \text{case nil of} \\ \quad \text{nil} \Rightarrow e_2 \\ \quad x : xs \Rightarrow e_3 \end{array} \right) = ()$	e_2
Case.cons	$\left(\begin{array}{l} \text{case } e_h : e_t \text{ of} \\ \quad \text{nil} \Rightarrow e_2 \\ \quad x : xs \Rightarrow e_3 \end{array} \right) = ()$	$e_3\{e_h/x, e_t/xs\}$

FIG. 7. Cost-equivalence laws K

$$\vdash_K e_1 = () e_2$$

22 A Naïve Time Analysis and its Theory of Cost Equivalence

if $e_1 =_{\langle \rangle} e_2$ is provable from the cost-equivalence laws K together with the facts that $=_{\langle \rangle}$ is a congruence relation (i.e. reflexivity, transitivity and substitutivity rules) and the following tick-elimination rule:

$$\mathbf{l\text{-}elim} \quad \frac{l(e_1) =_{\langle \rangle} l(e_2)}{e_1 =_{\langle \rangle} e_2} .$$

Now we have the following soundness and completeness results for the system \vdash_K , with respect to the cost-labelled transition relation \xrightarrow{t}_H :

THEOREM 7.4 (completeness)

For all closed expressions e , if $e \xrightarrow{m}_H h$ then $\vdash_K e =_{\langle \rangle} l^m(h)$.

The proof is given in Appendix A

THEOREM 7.5 (soundness)

For all closed expressions e , and head-normal-forms h_1 , if

$$\vdash_K e =_{\langle \rangle} l^m(h_1)$$

then $e \xrightarrow{m}_H h_2$ for some h_2 such that $\vdash_K h_1 =_{\langle \rangle} h_2$.

PROOF. By definition of l , and the fact that h_1 is a head-normal-form, $l^m(h_1) \xrightarrow{m}_H h_1$. By definition of cost equivalence it follows that $e \xrightarrow{m}_H h_2$ for some h_2 such that $h_1 =_{\langle \rangle} h_2$. It remains to show that this cost equivalence is provable:

$$\begin{aligned} e \xrightarrow{m}_H h_2 &\Rightarrow \vdash_K l^m(h_2) =_{\langle \rangle} e && \text{(Theorem 7.4)} \\ &\Rightarrow \vdash_K l^m(h_1) =_{\langle \rangle} l^m(h_2) && \\ &\Rightarrow \vdash_K h_1 =_{\langle \rangle} h_2 && \text{(l-elim)}. \end{aligned}$$

■

7.3 Example (continued)

Now we conclude the example from the beginning of Section 6.1, illustrating the use of cost equivalence, together with its proof techniques and axiomatization.

Recall the definition of quicksort (qs) from figure 6. From the time rules and the earlier analysis of the append function we obtained the time equation

$$\langle \text{qs}(e) \rangle^H = 1 + \langle e \rangle^H + \begin{cases} 0 & \text{if } e \rightarrow_H \text{nil} \\ 1 + \langle \text{qs}(\text{below}(y,z)) \rangle^H & \text{if } e \rightarrow_H y:z . \end{cases}$$

Then we considered the special case of non-increasing lists $\{A_i\}_{i \geq 0}$, and we showed that

$$\langle \text{qs}(A_{k+2}) \rangle^H = 5 + \langle \text{qs}(\text{below}(v_{k+1}, \text{below}(v_{k+2}, A_k))) \rangle^H .$$

The key to showing that $\langle \text{qs}(A_n) \rangle^H$ is quadratic in n is the identification of a cost equivalence which allows us to simplify a general instance of below. Define the family of lists $\{A_i^a\}_{i \geq 0, a \geq 0}$

inductively as follows:⁷

$$\begin{aligned} A_0^a &= \mathsf{l}^a(\text{nil}) \\ A_{k+1}^a &= \mathsf{l}^a(v_{k+1} : A_k^a). \end{aligned}$$

So the list-valued expression A_j^a is just like A_j except that each cons-cell (or nil) that is needed takes a evaluation steps to produce.

PROPOSITION 7.6

For all $a \geq 0$, and i, j such that $0 \leq j \leq i$,

$$\text{below}(v_i, A_j^a) =_{\langle \rangle} A_j^{a+1}.$$

PROOF. We sketch two proofs of the proposition. The first illustrates the basic cost-simulation proof techniques. Cost-simulation proofs are quite low level since they reason directly from the operational semantics. The second proof is more ‘calculational’ in style, and serves to illustrate the practical application of the axiomatization of cost equivalence.

1. We construct a family of relations containing all instances of the proposed cost equivalence, and show that each member (and hence their union) is a cost simulation. For each $i > 0$ $a \geq 0$, let X_i^a be the symmetric closure of the following relation:

$$\{(\text{filter}(\lambda x. x \leq v_i, A_j^a), A_j^{a+1}) \mid j \leq i\}$$

Now it is sufficient to show that each X_i^a is contained in a cost simulation, since this implies that their union is also contained in a cost simulation. To do this we show that X_i^a is a cost simulation modulo identity, i.e., that $X_i \subseteq \mathcal{F}(X_i^a \cup \equiv)$. Each pair of related elements in X_i^a has the form $(\text{below}(v_i, A_j^a), A_j^{a+1})$ (or vice versa). We proceed by cases according to the value of j . Suppose $j = 0$. Then from the definitions we have that $\text{below}(v_i, A_0^a) \xrightarrow{a+1}_H \text{nil}$ and $A_0^{a+1} \xrightarrow{a+1}_H \text{nil}$, and we are done. Suppose $j = k + 1$ for some $k \geq 0$. Then by calculation from the definitions

$$\begin{array}{ccc} \text{below}(v_i, A_{k+1}^a) & \xrightarrow{2+2a}_H & v_{k+1} : \text{below}(v_i, A_k^a) \\ A_{k+1}^{a+1} & \xrightarrow{2(a+1)}_H & v_{k+1} : A_k^{a+1}. \end{array}$$

Now the heads are related by the identity, and the tails by X_i^a , so the results are related by $(X_i^a \cup \equiv)$ and we are done.

2. By induction on j , using the cost equivalence laws K (Fig. 7).

Base: $j = 0$

$$\begin{aligned} \text{below}(v_i, A_0^a) &=_{\langle \rangle} \mathsf{l}(\text{case } A_0^a \text{ of } \dots) \\ &=_{\langle \rangle} \mathsf{l}(\mathsf{l}^a(\text{case } A_0 \text{ of } \dots)) \\ &=_{\langle \rangle} \mathsf{l}^{a+1}(\text{nil}) \\ &= A_0^{a+1}. \end{aligned}$$

⁷Recalling from earlier in the section that l is just the identity function $\mathsf{l}(x) = x$, and $\mathsf{l}^n(e)$ denotes n applications of the identity function to e , with the convention that $\mathsf{l}^0(e)$ is just e .

Induction: $j = k + 1$

$$\begin{aligned}
\text{below}(v_i, A_{k+1}^a) &=_{\langle \rangle} \text{l(case } A_{k+1}^a \text{ of } \dots) \\
&=_{\langle \rangle} |^{a+1}(\text{case } v_{k+1} : A_k^a \text{ of } \dots) \\
&=_{\langle \rangle} |^{a+1}(\text{if } v_i \leq v_{k+1} \text{ then } v_{k+1}.\text{below}(v_i, A_k^a) \text{ else } \dots) \\
&=_{\langle \rangle} |^{a+1}(v_{k+1}.\text{below}(v_i, A_k^a)) \\
&=_{\langle \rangle} |^{a+1}(v_{k+1}.A_k^{a+1}) && \text{(Hypothesis)} \\
&= A_{k+1}^{a+1}.
\end{aligned}$$

■

Remark While a proof using the cost-equivalence laws and simple induction may be preferable, the proof method of constructing a cost simulation is strictly more powerful—for example, it allows the proposition be generalized to include possibly infinite non-increasing lists (although this generalization is not relevant in the context of the sorting example).

Now we can return to quicksort. We consider the more general case of $\langle \text{qs}(A_j^a) \rangle^H$. Considering the cases when $j = 0$ and $j = k + 1$, and instantiating the general time equation gives

$$\begin{aligned}
\langle \text{qs}(A_0^a) \rangle^H &= 1 + \langle A_0^a \rangle^H \\
&= 1 + a \\
\langle \text{qs}(A_{k+1}^a) \rangle^H &= 1 + \langle A_{k+1}^a \rangle^H + 1 + \langle \text{qs}(\text{below}(v_{k+1}, A_k^a)) \rangle^H \\
&= 2 + a + \langle \text{qs}(\text{below}(v_{k+1}, A_k^a)) \rangle^H.
\end{aligned}$$

We can simplify the right-hand side by the proposition, to give

$$\langle \text{qs}(A_{k+1}^a) \rangle^H = 2 + a + \langle \text{qs}(A_k^{a+1}) \rangle^H.$$

Again the recurrence is easily solved; a simple induction is sufficient to check that

$$\langle \text{qs}(A_n^a) \rangle^H = \frac{n(n+5)}{2} + a(n+1) + 1.$$

Since the A_n are just A_n^0 , we finally have

$$\langle \text{qs}(A_n) \rangle^H = \frac{n(n+5)}{2} + 1.$$

8 Higher-order functions

One advantage of a simple operational approach to reasoning about programs is the relative ease with which we can handle higher-order functions. In this section we show how the time rules can be easily extended to cope with the incorporation of the terms and evaluation rules of the lazy lambda calculus [1]. The only potentially difficult part is the extension of the theory of cost equivalence. We sketch how the precongruence proof for cost simulation can be extended, with few modifications, to handle lambda terms and their application.

8.1 The lazy lambda calculus

We consider an extension to the language with the terms and evaluation rules of the lazy lambda calculus [1, 29]. The lazy lambda calculus, Λ , shares the syntax of the pure untyped lambda

calculus, but has an operational semantics which is consistent with implementations of higher-order functional languages, namely, there is no evaluation ‘under a lambda’.

The usual definitions of free and bound variables in lambda terms apply, and we do not repeat the definitions here. The evaluation rules for application and lambda terms are given below:

$$\begin{array}{c} \mathbf{lambda} \lambda x. e \rightarrow_{\alpha} \lambda x. e \\ \\ \mathbf{apply} \frac{e_1 \rightarrow_H \lambda x. e \quad e\{e_2/x\} \rightarrow_{\alpha} u}{e_1 e_2 \rightarrow_{\alpha} u} . \end{array}$$

Remark Notice that we do not evaluate under a lambda even in the case of evaluation to ‘normal-form’. This is consistent with the printing mechanisms provided for higher-order functional languages that allow functions to be the top-level results of programs. However, in the sequel we focus purely on the \rightarrow_H relation.

In the analysis of cost we choose additionally to count the number of times we invoke the **apply** rule in the evaluation of a term. The extension of the time rules is completely obvious:

$$\begin{array}{l} \langle \lambda x. e \rangle^{\alpha} = 0 \\ \langle e_1 e_2 \rangle^{\alpha} = 1 + \langle e_1 \rangle^H + \langle e\{e_2/x\} \rangle^{\alpha}, \text{ if } e_1 \rightarrow_H \lambda x. e . \end{array}$$

8.2 Applicative cost simulation

We will sketch the following:

- the extension of the definition of cost simulation to *applicative* cost simulation;
- the proof that applicative cost simulation is a precongruence.

The extension of the definition of cost simulation to handle the case where an expression evaluates to a lambda expression follows the definition of applicative (bi)simulation [1].

DEFINITION 8.1

If \mathcal{R} is a binary relation on closed expressions, then \mathcal{R}^{λ} is the binary relation on lambda expressions such that $(\lambda x. e_1 \mathcal{R}^{\lambda} \lambda y. e_2)$ if and only if for all closed expressions e , $(\lambda x. e_1) e \mathcal{R} (\lambda y. e_2) e$

As in Definition 6.5 we define applicative cost simulation as the maximal fixed point of a monotone function: For each binary relation R on closed expressions, define the relation $\mathcal{A}(R)$ by

$$e \mathcal{A}(R) e' \iff \begin{array}{l} \text{if } e \xrightarrow{t}_H h \\ \text{then } e' \xrightarrow{t}_H h' \text{ for some } h' \\ \text{such that } h(\mathcal{R} \cup \mathcal{R}^{\lambda})h' . \end{array}$$

Now we say that a relation S is an *applicative cost simulation* if $S \subseteq \mathcal{A}(S)$.

DEFINITION 8.2 (Applicative cost simulation)

Let \sqsubseteq denote the largest applicative cost simulation, the maximum fixed point of $\mathcal{A}(\cdot)$, given by $\bigcup\{S : S \subseteq \mathcal{A}(S)\}$

It is again straightforward to show that \sqsubseteq is a preorder. We prove that \sqsubseteq is preserved by substitution into arbitrary (closed) contexts by a direct extension of the proof of that for \preceq (Lemma 6.10 and Theorem 6.11). As before we construct a relation which contains \sqsubseteq and all closed substitution instances and show that it is a cost simulation.

THEOREM 8.3 (Precongruence II)

If $\tilde{e} \sqsubseteq \tilde{e}'$ for some commonly indexed families of closed expressions \tilde{e}, \tilde{e}' , then for all expressions e containing at most variables \tilde{x}

$$e\{\tilde{e}/\tilde{x}\} \sqsubseteq e\{\tilde{e}'/\tilde{x}\}.$$

The proof is sketched in Appendix A.

Again the use of the term *congruence* could be challenged since we do not consider open expressions whose free variables are captured by the context. As before we can extend applicative cost simulation to open expressions e and e' , by saying $e \sqsubseteq e'$ if $e\sigma \sqsubseteq e'\sigma$ for all closing substitutions σ . It is then easy to show that, for example, $\lambda x.e \sqsubseteq \lambda x.e'$.

9 A further example

In this section we present a final example. It gives a good illustration of the use (and proof) of cost equivalence in the derivation of a time property. The reader is invited to attempt an analysis without the use of cost equivalence.

9.1 Maxtails

Figure 8 defines some functions including `max`, which computes the first element, according to dictionary order, of a list of words. Words are represented as lists of characters. `max` employs an auxiliary binary comparison on words, `dmax`, which in turn employs a primitive function `precedes` which tests if one character precedes another.

Two abbreviations have been adopted to aid presentation: parentheses have been elided in the application of unary functions, and a function name f (n -ary) written directly denotes the obvious abstraction $\lambda x_1 \dots \lambda x_n.f(x_1 \dots x_n)$.

In what follows we will denote words just by the concatenation of the elements (so, for example, `aab` represents the list `a:a:b:nil`). We can thus illustrate the functionality of `dmax` by saying that $\text{dmax}(a,aa) \rightarrow_N a$ and $\text{dmax}(aa,ab) \rightarrow_N aa$.

The object of the example is the function `maxtails`, which computes the dictionary maximum of the non-empty *tails* of a list. For example, the tails of `aba` are `aba`, `ba` and `a`, of which `a` is the ‘maximum’. The objective is, at first sight a modest one. We wish to show that `maxtails` can take quadratic time (in the length of the list argument) to produce a normal form.

The quadratic time result is not obvious because of the interaction of the lazy evaluation order and the (lazy) lexicographic ordering, which very often gives good performance. For example, `maxtails` is linear on the following classes of lists:

- lists of strictly ‘decreasing’ elements eg. `abcde...`,
- lists of strictly ‘increasing’ elements, eg. `zyx...`, and
- stationary lists, eg `bbb...`

The proofs of these properties are left as exercises.

maxtails xs	=	max (tails xs)
max xs	=	case xs of nil \Rightarrow undefined h:t \Rightarrow foldr(dmax,h,t)
dmax(xs,ys)	=	case ys of nil \Rightarrow nil h:t \Rightarrow case xs of nil \Rightarrow nil h':t' \Rightarrow if h=h' then h: dmax(t',t) else if (h precedes h') then h:t else h':t'
tails ys	=	case ys of nil \Rightarrow nil h:t \Rightarrow (h:t): tails t
foldr(f,b,xs)	=	case xs of nil \Rightarrow b h:t \Rightarrow f h foldr(f,b,t)

FIG. 8. Maxtails

9.1.1 Overview

The remainder of this section builds a proof of the above quadratic time property. We break the proof down into a number of distinct steps, each of which illustrates some techniques for reasoning using cost equivalence.

The first step is to find a simpler representation of the problem via a cost equivalence: we derive a recursive function and recast the problem in terms of properties of this new function. The second step is to find a family of lists that will yield the quadratic time result (we just give some informal motivation at this point). Now, as in the quicksort example, we find a crucial simplifying cost equivalence relating to this family of lists. Given these steps the final time analysis is straightforward.

9.2 An equivalent problem

From the cost-equivalence laws, including the case law from Proposition 7.1, we have that

$$\begin{aligned}
 \max (\text{tails } x) &=_{() } l(\text{case } (\text{tails } x) \text{ of} \\
 &\quad \text{nil} \Rightarrow \text{undefined} \\
 &\quad \text{h:t} \Rightarrow \text{foldr}(\text{dmax}, \text{h}, \text{t})) \\
 &=_{() } l^2(\text{case } x \text{ of} \\
 &\quad \text{nil} \Rightarrow \text{undefined} \\
 &\quad \text{h:t} \Rightarrow \text{foldr}(\text{dmax}, (\text{h:t}), \text{tails } t))
 \end{aligned}$$

Now we wish to proceed by analysing the expression $\text{foldr}(\text{dmax},(h:t),\text{tails } t)$. We derive a recursive function for the slightly more general expression $\text{foldr}(\text{dmax},y,\text{tails } xs)$. The function fot we derive will satisfy $\text{fot}(y,xs) =_{\text{()}} \text{l}(\text{foldr}(\text{dmax},y,\text{tails } xs))$. Initially we can satisfy this by defining

$$\text{fot}(y,xs) = \text{foldr}(\text{dmax},y,\text{tails } xs).$$

We consider this to be an initial specification of fot , and proceed by transforming the right-hand side in the manner of unfold fold transformation [12], maintaining cost equivalence:

$$\begin{aligned}
& \text{foldr}(\text{dmax},y,\text{tails } xs) \\
=_{\text{()}} & \text{l}(\text{case } (\text{tails } xs) \text{ of} && \text{(unfold foldr)} \\
& \quad \text{nil} \Rightarrow y \\
& \quad \text{h:t} \Rightarrow \text{l}^2(\text{dmax}(\text{h},\text{foldr}(\text{dmax},y,t))) \\
=_{\text{()}} & \text{l}^2(\text{case } xs \text{ of} && \text{(unfold tails, case law)} \\
& \quad \text{nil} \Rightarrow y \\
& \quad \text{h:t} \Rightarrow \text{l}^2(\text{dmax}((\text{h:t}),\text{foldr}(\text{dmax},y,\text{tails } t)))) \\
=_{\text{()}} & \text{l}^2(\text{case } xs \text{ of} && \text{(unfold dmax)} \\
& \quad \text{nil} \Rightarrow y \\
& \quad \text{h:t} \Rightarrow \text{l}^3(\text{case } \text{foldr}(\text{dmax},y,\text{tails } t) \text{ of } \dots)) \\
=_{\text{()}} & \text{l}^2(\text{case } xs \text{ of} && \text{(case law)} \\
& \quad \text{nil} \Rightarrow y \\
& \quad \text{h:t} \Rightarrow \text{l}^2(\text{case } \text{l}(\text{foldr}(\text{dmax},y,\text{tails } t)) \text{ of } \dots)) \\
=_{\text{()}} & \text{l}^2(\text{case } xs \text{ of} && \text{(fold dmax)} \\
& \quad \text{nil} \Rightarrow y \\
& \quad \text{h:t} \Rightarrow \text{l}(\text{dmax}((\text{h:t}),\text{l}(\text{foldr}(\text{dmax},y,\text{tails } t)))) \\
=_{\text{()}} & \text{l}^2(\text{case } xs \text{ of} && \text{(fot spec.)} \\
& \quad \text{nil} \Rightarrow y \\
& \quad \text{h:t} \Rightarrow \text{l}(\text{dmax}((\text{h:t}),\text{fot}(y,t))))
\end{aligned}$$

So we obtain a recursive definition

$$\begin{aligned}
\text{fot}(y,xs) = \text{l}^2(\text{case } xs \text{ of} \\
& \quad \text{nil} \Rightarrow y \\
& \quad \text{h:t} \Rightarrow \text{l}(\text{dmax}((\text{h:t}),\text{fot}(y,t))))
\end{aligned}$$

PROPOSITION 9.1

$$\text{fot}(y,xs) =_{\text{()}} \text{l}(\text{foldr}(\text{dmax},y,\text{tails } xs))$$

PROOF. The above derivation constitutes a proof, although the fact that this *is* a proof needs some further justification, and depends critically on the fact that the steps are cost equivalences—see [40], but we can also prove it directly by the usual method of showing that it is a simulation. The details are left as an exercise. ■

We are interested in computing the normal-form of maxtails e . From the above cost equivalences, we have that:

$$\begin{aligned}
 \langle \text{maxtails } e \rangle^N &= 1 + \left\langle l^2 \left(\begin{array}{l} \text{case } e \text{ of} \\ \text{nil} \Rightarrow \text{undefined} \\ \text{h:t} \Rightarrow \text{foldr}(\text{dmax}, (\text{h:t}), \text{tails } t) \end{array} \right) \right\rangle^N \\
 &= 1 + \left\langle l \left(\begin{array}{l} \text{case } e \text{ of} \\ \text{nil} \Rightarrow l \text{ undefined} \\ \text{h:t} \Rightarrow l \text{ foldr}(\text{dmax}, (\text{h:t}), \text{tails } t) \end{array} \right) \right\rangle^N \\
 &= 2 + \left\langle \begin{array}{l} \text{case } e \text{ of} \\ \text{nil} \Rightarrow l \text{ undefined} \\ \text{h:t} \Rightarrow \text{fot}((\text{h:t}), t) \end{array} \right\rangle^N .
 \end{aligned}$$

9.3 A quadratic case

Informally speaking, dmax evaluates enough of its arguments to determine which is the answer. So the amount of evaluation is bounded by the length of the answer. This suggests that to obtain worst-case inputs for maxtails, the size of the result should be $\mathcal{O}(n)$, where n is the length of the input. Furthermore, to force dmax to ‘recurse’ often, the various tails of the input should be, as far as possible, element-wise equal. Inputs of the form $a \dots ab$ satisfy these requirements—the result is the input itself, and any pair of tails, eg. $aaaab$ and aab , are element-wise equal up to but not including the last element of the shorter.

DEFINITION 9.2

For $k \geq 0$, let $a^k b$ denote the list consisting of k as followed by a single b .

We will show that this family of lists gives rise to quadratic time performance.

The following family of functions will be instrumental in expressing a key cost equivalence

DEFINITION 9.3

The functions $\{\text{T}_k\}_{k \geq 0}$ are given by the following scheme:

$$\begin{aligned}
 \text{T}_0 \text{ xs} &= \text{xs} \\
 \text{T}_{k+1} \text{ xs} &= \text{case } \text{xs} \text{ of} \\
 &\quad \text{nil} \Rightarrow \text{nil} \\
 &\quad \text{h:t} \Rightarrow \text{h}:(\text{T}_k \text{ t})
 \end{aligned}$$

For list-valued arguments, T_k ‘traverses’ its argument up to a depth k . The following properties of the T_k follow easily:

PROPOSITION 9.4

1. $\text{T}_0 e = () \text{ l}(e)$.
2. $\text{T}_{k+1} (e_1 : e_2) = () \text{ l}(e_1 : \text{T}_k e_2)$.
3. For all expressions e such that $e \rightarrow_N v_1 : \dots : v_j : \text{nil}$, $j \geq 0$,

$$\langle \text{T}_k e \rangle^N = 1 + \text{max}(j, k) + \langle e \rangle^N .$$

Before we state the key cost equivalence we need one technical construction:

DEFINITION 9.5

A closed expression is *element cheap* if there exists a set E containing the expression, such that the following property is satisfied. For all $e \in E$, if $e \rightarrow_H e_1 : e_2$ then $e =_{\langle \rangle} l^k(v : e')$ for some $k \geq 0$, some normal-form v and some $e' \in E$.

The intuition behind the definition of element-cheap expressions is that if it evaluates to a cons-expression, then the head will already be in normal-form, and that this property holds for the tail of the list, the tail of the tail and so on. The following gives an example, and will be needed later:

PROPOSITION 9.6

If e_1 and e_2 are element cheap, then so is $\text{fot}(e_1, e_2)$.

PROOF. From the definition we need to construct a set E containing $\text{fot}(e_1, e_2)$ and satisfying the condition of the proposition. The set we take is

$$E = \{\text{fot}(e, e') \mid e, e' \text{ are element cheap}\} \cup \{e'' \mid e'' \text{ is element cheap}\}.$$

Now assume that $\text{fot}(e_1, e_2)$ evaluates to a cons (otherwise we are done). The remainder of the proof is a straightforward case analysis on the evaluation of $\text{fot}(e_1, e_2)$, and we omit the details. ■

PROPOSITION 9.7

For all n, k such that $0 \leq k < n$, if e is element cheap and $e \rightarrow_N a^n b$, then

$$\text{dmax}(a^k b, e) =_{\langle \rangle} T_k e =_{\langle \rangle} \text{dmax}(e, a^k b).$$

PROOF. We prove the first cost equivalence; the second is similar (but not symmetrical). Since $n > 0$ and e is element cheap, then $e =_{\langle \rangle} l^j(a : e')$, for some j . Unfolding dmax , and applying the tick-laws to the outer case-expression:

$$\begin{aligned} \text{dmax}(a^k b, e) &=_{\langle \rangle} l^{j+1} \text{ case } a^k b \text{ of} \\ &\quad \text{nil} \Rightarrow \text{nil} \\ &\quad \text{h':t'} \Rightarrow \text{if } a = \text{h'} \text{ then } a: \text{dmax}(t', e') \\ &\quad \quad \text{else if (a precedes h')} \text{ then } a:e' \text{ else h':t'} \end{aligned}$$

Now we show, by induction on k , that this is cost equivalent to $T_k e$.

Base($k = 0$): Then $a^k b$ is just $b:\text{nil}$, so the above simplifies to

$$\begin{aligned} l^{j+1} \text{ if } a = b \text{ then } a: \text{dmax}(\text{nil}, e') &=_{\langle \rangle} l^{j+1} (a:e') \\ \text{else if (a precedes b) then } a:e' \text{ else } b:\text{nil} &=_{\langle \rangle} l(e) \\ &=_{\langle \rangle} T_0 e. \end{aligned}$$

Induction($k = m + 1$): In this case $a^k b$ is $(a:a^m b)$, so the above simplifies to

$$\begin{aligned} l^{j+1} \text{ if } a = a \text{ then } a: \text{dmax}(a^m b, e') &=_{\langle \rangle} l^{j+1} (a:\text{dmax}(a^m b, e')) \\ \text{else if (a precedes a) then } a:e' \text{ else } a:a^m b &=_{\langle \rangle} l^{j+1} (a:\text{dmax}(a^m b, e')). \end{aligned}$$

Now e' has normal form $a^{n-1}b$, and is element cheap (since e is), so we can apply the inductive hypothesis:

$$\begin{aligned} l^{j+1} (a:\text{dmax}(a^m b, e')) &=_{\langle \rangle} l^{j+1} (a:T_m e') \\ &=_{\langle \rangle} l^j T_{m+1} (a:e') \\ &=_{\langle \rangle} T_{m+1} l^j (a:e') \\ &=_{\langle \rangle} T_{m+1} e. \end{aligned}$$

■

9.4 The final analysis

We now draw the components together to show that $\text{maxtails } a^n b$ is $\mathcal{O}(n^2)$. First, assume $n > 0$. From the simplifying cost equivalence in Section 9.2 we have that

$$\begin{aligned} \langle \text{maxtails } a^n b \rangle^N &= 2 + \left\langle \begin{array}{l} \text{case } a^n b \text{ of} \\ \text{nil} \Rightarrow \text{! undefined} \\ \text{h:t} \Rightarrow \text{fot}((\text{h:t}), \text{t}) \end{array} \right\rangle^N \\ &= 2 + \langle \text{fot}(a^n b, a^{n-1} b) \rangle^N. \end{aligned}$$

Now we take advantage of our knowledge of the extensional properties of maxtails (without proof)—in particular, that the normal-form of $\text{fot}(a^n b, a^m b)$ for any $m < n$ is $a^n b$. In the general case where $n > m > 0$ we have that

$$\begin{aligned} \langle \text{fot}(a^n b, a^m b) \rangle^N &= 3 + \langle \text{dmax}(a^m b, \text{fot}(a^n b, a^{m-1} b)) \rangle^N \\ &= 3 + \langle \text{T}^m \text{fot}(a^n b, a^{m-1} b) \rangle^N && (\text{Prop (9.6), Prop (9.7)}) \\ &= 3 + m + \langle \text{fot}(a^n b, a^{m-1} b) \rangle^N && (\text{Prop (9.4)}). \end{aligned}$$

The case when $n > m = 0$ is given by direct calculation from the time rules:

$$\langle \text{fot}(a^n b, b) \rangle^N = 8.$$

So we can solve this recurrence when $n > m \geq 0$ to give

$$\langle \text{fot}(a^n b, a^m b) \rangle^N = \left(\sum_1^m i \right) + 3m + 8.$$

So returning to the main problem,

$$\begin{aligned} \langle \text{maxtails } a^n b \rangle^N &= 2 + \langle \text{fot}(a^n b, a^{n-1} b) \rangle^N \\ &= 2 + \left(\sum_1^{n-1} i \right) + 3(n-1) + 8 \\ &= \frac{1}{2}n(n+1) + 2n + 7. \end{aligned}$$

Final remark Some of the intermediate steps are more general than necessary to obtain this result. In particular the technicalities Definition 9.5–Proposition 9.7 regarding element-cheap expressions and the properties of dmax could be eliminated by a more direct proof in the final analysis above, but they help make the result robust with respect to, for example, the order in which the tails are ‘folded’ together. For example it is an easy exercise to show that replacing foldr by a ‘fold from the left’ does not essentially change this quadratic-time case.

10 Call-by-need and compositionality

The calculus is betrayed by its simple operational origins because it describes a *call-by-name* evaluation mechanism, when most actual implementations of lazy evaluation use *call-by-need*. For example, consider the definition

$$\text{average}(xs) = \text{divide}(\text{sum}(xs), \text{length}(xs))$$

where divide is a primitive function, and sum and length are the obvious functions on lists. In reasoning about the evaluation of an instance, $\text{average}(e)$, our method will overestimate the

evaluation time, because of the duplication of expression e on substitution into the body of `average`. Assuming e has a normal-form which is a non-empty list of integers, to compute the necessary calls to `sum` and `length` each `cons` in the result of e must be computed, but this work will be performed independently by `sum` and `length` using their own copies of e , whereas under call-by-need this evaluation (and hence its cost) will be shared.

One solution is to work with an operational model that takes into account the sharing of expressions and evaluations [5]. Unfortunately this may overly complicate the calculus, and is likely to be impractical—although there are some promising (less general) approaches to modelling sharing and storage, [3, 25], which may prove usable.

Another solution is to move towards the compositional approaches mentioned in the introduction. A suitable interface between the compositional approach in [37] (which differs from that of [42] in its use of genuine *strictness* rather than *absence* information) from and the operational approach of this paper is via Burn’s notion of an *evaluator* [10]. An evaluator is an operational concept which provides a link from information provided by (list-based extensions of) strictness analysis, to the operational semantics. In particular, strictness analysis [11] will tell us that when an application of `average` is being evaluated, it is safe to evaluate the argument to normal-form (since this evaluation will occur anyway). In terms of our calculus, by taking into account (in advance) the amount of evaluation that must be performed on the argument, we obtain a more compositional analysis, and a better approximation to call-by-need, using:

$$\langle \text{average}(e) \rangle^H = \langle e \rangle^N + \langle \text{average}(v) \rangle^H, \text{ where } e \rightarrow_N v.$$

In this section we describe a new formalization of evaluators appropriate for providing a smooth interface between compositional and non-compositional (call-by-need and call-by-name) approaches to time analysis. The development is for the first-order language, although it can be used within higher-order programs.

10.1 Demands

Burn’s formalization of evaluators is couched in terms of *reduction strategies* in a typed lambda calculus with constants. An evaluator is defined, relative to a particular (Scott) closed set of denotations, as any reduction strategy which fails to terminate exactly when the denotation of the term is in the set. For example, the leftmost outermost reduction strategy fails to terminate if and only if the denotation of the term is in the Scott-closed set $\{\perp\}$.

However, the definition of an evaluator is not constructive. Given a Scott-closed set, Burn does not provide an operational definition of an evaluator for that set.

In the approach we have taken to the operational semantics of the language, issues of reduction strategy are internalized by the evaluation relations. In order to use the evaluator concept in reasoning about evaluation cost we need a constructive definition of evaluators. We will define a family of evaluators, and show how the relations \rightarrow_N and \rightarrow_H can be viewed as instances. The starting point of this definition is a language of *demands*.

We interpret a demand as a specification of a degree of evaluation, and define a demand-parameterized evaluation relation which realises these demands. Bjerner and Hölmstrom [8] use a particularly simple language of demands in the context of compositional time analysis. Their interpretation of a demand only makes sense relative to a particular expression: a demand on an expression is a representation of an approximation to its denotation. Their language of demands is too ‘precise’ for our purposes. The language of demands which we use will be closer to that of, for example, Lindström’s lattice of demands [23] or Dybjer’s formal opens

[14], and our interpretation of demand, as in [14], will be closely related to strictness. The operational interpretation of demand as an evaluator is, in turn, reminiscent of Bjerner's definition of computing an expression to one of its *proper evaluation degrees* [7].

DEFINITION 10.1

The language of demands, $d \in D$ is given by

$$D = \epsilon \mid \kappa \mid D_1 :: D_2 \mid \kappa + D_1 :: D_2 \mid \mu x.(D) \mid x.$$

The demand ϵ is the zero-demand which is satisfied by any expression. The demand κ is satisfied by any expression which evaluates to a constant, including nil (we could easily extend the demands to include demands for each individual constant). The cons-demand $D_1 :: D_2$ is satisfied by any expression which evaluates to a cons, and whose head and tail satisfy D_1 and D_2 respectively. We have a restricted form of disjunctive demand for either a constant or a cons, and finally we add a recursive demand useful for specifying demands over lists.

This informal reading of a demand can easily be formalized, and the set of expressions which satisfy a demand can be shown to be *open* with respect to the usual operational preorder (c.f.[14]). We do not pursue this formalization here, but just focus on the operational interpretation of demands as evaluators.

DEFINITION 10.2

For each closed demand d we define an associated evaluator \Longrightarrow_d . The family of evaluators is defined by inductively as the least relations between closed expressions that satisfy the following rules:

$$\begin{array}{c} \frac{}{e \Longrightarrow_{\epsilon} e} \qquad \frac{e \rightarrow_H c}{e \Longrightarrow_{\kappa} c} \\ \\ \frac{e \rightarrow_H e_1 : e_2 \quad e_1 \Longrightarrow_{d_1} e'_1 \quad e_2 \Longrightarrow_{d_2} e'_2}{e \Longrightarrow_{d_1 :: d_2} e'_1 : e'_2} \\ \\ \frac{e \Longrightarrow_{\kappa} e'}{e \Longrightarrow_{\kappa + d_1 :: d_2} e'} \qquad \frac{e \Longrightarrow_{d_1 :: d_2} e'}{e \Longrightarrow_{\kappa + d_1 :: d_2} e'} \\ \\ \frac{e \Longrightarrow_{d\{\mu x.(d)/x\}} e'}{e \Longrightarrow_{\mu x.(d)} e'}. \end{array}$$

Note that the evaluators are defined using the basic reduction engine, \rightarrow_H . It is not difficult to show that evaluators are deterministic (i.e., each \Longrightarrow_d is a partial function) and that proofs of evaluation judgements are unique. As for \rightarrow_{α} , we write $e \xrightarrow{t} e'$ when $e \Longrightarrow_d e'$ and its proof uses $t \geq 0$ instances of rule **Fun**.

Now we can show that the previous evaluation relations can be viewed as instances of evaluators:

PROPOSITION 10.3

1. $e \xrightarrow{t}_{\kappa + \epsilon :: \epsilon} e' \iff e \xrightarrow{t}_H e'$
2. $e \xrightarrow{t}_{\mu x.(\kappa + x :: x)} e' \iff e \xrightarrow{t}_N e'$

PROOF. 1. The first part is simple. (\Rightarrow) follows from the definition of evaluators, (\Leftarrow) by considering the two cases for e' .

2. (\Rightarrow) Induction on the size of the proof of $e \xrightarrow{\mu x.(\kappa+x::x)}^t e'$. Abbreviate demand $\mu x.(\kappa+x::x)$ by A . Clearly either $e' = c$, and hence $e \xrightarrow{t}_H c$, in which case $e \xrightarrow{t}_N c$ follows from Proposition 6.7, or $e' = e_1 : e_2$, and the proof has the form

$$\frac{\frac{e \xrightarrow{t'}_H e'_1 : e'_2 \quad e'_1 \xrightarrow{t1}_A e_1 \quad e'_2 \xrightarrow{t2}_A e_2}{e \xrightarrow{t}_A e'_1 : e'_2}}{e \xrightarrow{t}_A e_1 : e_2}}{e \xrightarrow{t}_{\kappa+A::A} e_1 : e_2}}{e \xrightarrow{t}_A e_1 : e_2}$$

where $t = t' + t_1 + t_2$. By the induction hypothesis, $e'_1 \xrightarrow{t1}_N e_1$ and $e'_2 \xrightarrow{t2}_N e_2$. By the rules for evaluation to normal-form, $e'_1 : e'_2 \xrightarrow{t1+t2}_N e_1 : e_2$ and so by Proposition 6.7, since $t = t' + t_1 + t_2$ we can conclude $e \xrightarrow{t}_N e_1 : e_2$, as required.

(\Leftarrow) follows by a routine induction on the structure of normal form e' , appealing again to Proposition 6.7. The details are omitted. \blacksquare

Now we can give a definition of time rules $\langle _ \rangle^d$ which extend the definitions for $\langle _ \rangle^H$ (and subsume the definitions for $\langle _ \rangle^N$). The rules are obvious, but we include them for completeness in Fig. 9. It is also possible to show, (although we do not provide details here) that cost equivalence

$$\begin{aligned} \langle e \rangle^\epsilon &= 0 \\ \langle e \rangle^\kappa &= \langle e \rangle^H \text{ if } e \rightarrow_H c \\ \langle e \rangle^{d_1::d_2} &= \langle e \rangle^H + \langle e_1 \rangle^{d_1} + \langle e_2 \rangle^{d_2} \text{ if } e \rightarrow_H e_1 : e_2 \\ \langle e \rangle^{\kappa+d_1::d_2} &= \langle e \rangle^H + \begin{cases} 0 & \text{if } e \rightarrow_H c \\ \langle e_1 \rangle^{d_1} + \langle e_2 \rangle^{d_2} & \text{if } e \rightarrow_H e_1 : e_2 \end{cases} \end{aligned}$$

FIG. 9. Demand time rules

is congruence with respect to the demanded time rules,

$$e \preceq e' \Rightarrow \forall C. \langle C[e] \rangle^d = \langle C[e'] \rangle^d, \text{ whenever } C[e] \downarrow_a.$$

10.2 Demand use

The weakness in the model with respect to call-by-need computation is that the substitution operation duplicates expressions, and consequently does not ‘share’ any evaluations of that expression to head-normal-form (and subsequent evaluations of sub-expressions).

Let $\langle e \rangle_{\text{need}}^d$ informally denote the call-by-need cost of evaluating expression e up to demand d .

Suppose we know that the following condition (Δ) holds:

Whenever an application of some function f is to be evaluated up to demand d , then its argument must satisfy some demand d' .

In other words, for all arguments e , if $f(e) \Longrightarrow_d e'$ for some e' , then there exists a u such that $e \Longrightarrow_{d'} u$. When this is the case, we can *refine* our call by name model as follows.

$$\langle f(e) \rangle^d \geq \langle e \rangle^{d'} + \langle f(u) \rangle^d \geq \langle f(u) \rangle_{\text{need}}^d.$$

To attempt a rigorous justification of the second inequality, we would naturally need to formalize what we mean by call-by-need evaluation. This nontrivial task is beyond the scope of this paper, and we leave it as an open problem.

Here is an informal explanation of the first inequality. First, it is not too difficult to show that when condition (Δ) holds for f that there exists an expression u such that $e \Longrightarrow_{d'} u$, and that $f(u) \Longrightarrow_d e'$.⁸ Now, since e must satisfy d' whenever $f(e)$ satisfies d , it must be the case that each sub-proof of a judgement of the form $a \rightarrow_H b$ in the proof of $e \Longrightarrow_{d'} u$ must also occur *at least once* in the proof of $f(e) \Longrightarrow_d e'$. This is because no backtracking is needed in proof construction, and so the condition implies that \Longrightarrow_d must perform as much computation on e as $\Longrightarrow_{d'}$. Now consider the proof of $f(u) \Longrightarrow_d e'$. This proof will have essentially the same structure as the proof of the evaluation of $f(e)$ up to d , except that each of the aforementioned sub-proofs will be replaced by (zero-cost) sub-proofs of the form $b \rightarrow_H b$. The *inequality* arises because the costs of some multiply-occurring sub-proofs for an evaluation $a \rightarrow_H b$ in the proof of $e \Longrightarrow_{d'} u$, arising through expression duplication, may be counted only once in $\langle e \rangle^{d'}$.

We have not found a satisfactory proof based on this sketch (or otherwise). It would be nice to have an ‘algebraic’ proof which does not mention proof-trees explicitly. Cost equivalence and related tools such as *improvement* preorderings [38] may be useful here.

10.3 Demand propagation

In order to use this method for refining the call-by-name model we need to establish for some context $C[\]$, the demand an expression placed in its hole satisfy must satisfy, given some demand which must be satisfied by the composite expression. Of course, such demands on the hole are not unique. For example, any expression placed in the hole must satisfy the trivial demand ϵ . However this trivial demand does not allow us to refine the call-by-name cost model. In general we want to determine as ‘large’ a demand on the hole as possible.

In this paper we shall not pursue the problem of demand propagation, but mention some of the connections with the much-studied strictness analysis problem.

As mentioned earlier, demand propagation can be formalized by modelling a demand as the set of expressions which satisfy it. These sets can be shown to be open (right-closed) under an operational preordering on expressions, and so the problem can be viewed as an inverse-image analysis problem [14], albeit expressed in terms of an operational model rather than a denotational one. The corresponding problem in a higher-order language is more easily tackled by a forwards analysis in which information about the *complement* of a demand (all the elements which do not satisfy it) are propagated forwards from sub-expressions to their context. This corresponds to the higher-order strictness analysis described in [11]. So in principle the approach described here can be extended to a higher-order language, but it is not obvious how a higher-order demand (other than simple strictness) can be given an operational interpretation, so there remains an inherently non-compositional aspect (c.f.. the higher-order approach described in [36], Ch. 5).

⁸This is a *constructive* version of Burn’s evaluation transformer theorem—further investigations will be presented elsewhere.

10.4 Comparison with earlier compositional approaches

Here we place the approach outlined in this section in comparison with the compositional methods for time analysis overviewed in Section 2. The use of demands to refine call-by-name evaluation time to give a better approximation to call-by-need is consistent with the use of strictness information in *necessary-time analysis* [35]. The key difference is that strictness information is used by necessary-time analysis to give a lower bound to call-by-need computation cost, whereas it is used here to give an upper-bound (which in turn is bounded above by call-by-name cost). In both cases the quality of demand information determines the tightness of the bound. We would argue that the method here is more useful.

The time analysis described by Wadler [42] (which is the first order instance of *sufficient-time analysis* in [35]) also gives an upper-bound to call-by-need time, but in a rather different manner: it uses information about ‘absence’ (constancy) which describes what parts of an expression will *not* be evaluated. Unfortunately in this case although the quality of this upper bound is determined by the quality of this ‘absence’ information,⁹ this upper-bound may not be well defined even when the program is.

The approach described by Bjerner and Hölmstrom [8] is a fully compositional method for reasoning about first-order functional programs with lazy lists. This approach is harder to compare because the ‘demand’ information is *exact*; the call-by-need time analysis is therefore exact as well (although, as here, the treatment of call-by-need is not formalized). This precision, as the authors note, is also a serious drawback in reasoning about programs, since as a first step one must decide what ‘demand’ to make on the output. Since the language of demands is so precise, in order to reason about computation of a program to normal form one must begin with the demand that *exactly* describes that normal form—in other words (a representation of) the normal form itself. It is not immediately obvious how to introduce ‘approximate’ demands in this approach without running into the definedness problems of Wadler’s method.

To summarize, the approach sketched in this section has the following advantages:

- it provides a safe (well-defined) time bound lying between call-by-need and call-by-name costs;
- it allows flexibility in the use of demand information so it can be targeted to where either
 1. more compositionality is needed to decompose the time analysis of a large program, or
 2. where (it is suspected that) the call-by-name model is too crude.

Although our method and exposition were inspired by formulating a constructive (operational) version of Burn’s evaluators, in retrospect the approach is much closer to Bjerner’s original work on time analysis of programs of type theory, because of its operational basis. The connection to evaluators does not occur in Bjerner’s work because of the absence of non-terminating computations in that language.

11 Related work

Here we review work related to the theory of cost simulation.

⁹Quality should not be confused with *safety*—we assume that the information is always *safe* in the sense that whenever it predicts a property of a program, the property does indeed hold.

11.1 Other non-standard theories of equivalence

Talcott [41] introduced a broad semantic theory of side-effect-free Lisp-like languages, notable for its treatment of both extensional and intensional aspects of computation. In particular a class of preorderings called *comparison relations* were introduced for a side-effect-free Lisp derivative. The method of their definition is similar to the definition of operational approximation as a Park-style simulation, although the cost aspects are not built into this definition directly as they are here. Nonetheless, the class of comparison relations could be said to contain relations analogous to the cost-equivalence relation considered here (this is just the observation that cost simulation can be viewed as a refinement of a pure simulation not involving time properties), and indeed it is suggested (as a topic for further work) that certain comparison relations could be developed to provide soundness and improvement proofs for program transformation laws. However, only the ‘maximal’ comparison relations (essentially, the more usual operational approximation and equivalence relations) and their application are considered in detail. Following on from other aspects of this work, Mason [24] sketches the definition a family of equivalence relations involving a variety of operation execution counts. However, as this is for a pure language with neither higher-order functions nor lazy data structures, the relations have a relatively uninteresting structure.

Moggi’s categorical semantics of computation [28] is intended to be suitable for capturing broader descriptions of computation than just input–output behaviour. Gurr [16] has studied extensions of denotational semantics to take account of resource use, and has shown how Moggi’s approach can be used to model computation in such a way that program equivalence also captures equivalence of resource requirements. Gurr extends Moggi’s λ_c -calculus (a formal system for reasoning about equivalence) with sequents for reasoning about the resource properties *directly* (although the ability to do this depends on certain ‘representability’ issues, not least of all that the resource itself should correspond to a type in the metalanguage). The resulting calculus is dubbed ‘ λ_{com} ’. We can compare this to the approach taken here, where cost equivalence (a ‘resource’ equivalence) is used in conjunction with a set of *time rules* which are used to reason about the cost property directly. There are further analogies in the details of Gurr’s approach: he defines rules for sequents of the form

$$\Gamma \vdash_{\lambda_{com}} \langle e, v, t \rangle : \tau$$

which says that expression e has value v (of type τ) and consumes resource t . He goes on to show that these rules are redundant since their information can be expressed in the λ_c calculus (with the addition of some specific axioms). In particular the above entailment would imply

$$\Gamma \vdash_{\lambda_c} \text{let } x = t \text{ in } v : \tau$$

where t is the canonical term (of unit type) which consumes time ‘ t ’. Notice the similarity with our axiomatization of cost equivalence, which uses the identity function in much the same manner as the computational let is used in the λ_c calculus.

An important difference is that in our approach these concepts are derived from (and hence correct with respect to) the operational model, so we may argue, at least, that an operational approach provides a more appropriate starting point for a semantic study of efficiency—although the deeper connections between these approaches deserves some further study.¹⁰ Another important difference is that we consider a lazy recursive data structure. One possible method for dealing with this in the context of Moggi’s framework would be to combine it with Pitts’ co-induction

¹⁰ Interestingly, Gurr gives an operational semantics (call-by-value) to terms of the λ_{com} calculus, but leaves the property corresponding to our Theorem 7.5 as a conjecture.

principal for recursively defined domains [32]. Other aspects of Gurr’s work are concerned with giving a semantic framework for some aspects of (asymptotic) complexity, which is outside the scope of this work.

Another, perhaps more abstract category-theoretic approach to intensional semantics is presented by Brookes and Geva [9]. This work places a heavy emphasis on the relationship between the intensional semantics and its underlying extensional one. The key structure is the use of a co-monad, in contrast to Gurr’s use of Moggi’s monadic style, and it’s suitability for formulating the kind of intensional semantics described here is perhaps less obvious.

11.2 Program improvement and generalizations

The theory of cost simulation is significant in its own right since there are many potentially interesting relations (preorders and equivalences) involving time, that can be investigated in a similar manner. For example, consider a *program refinement* relation, \triangleright , which is the largest relation such that whenever $e_1 \triangleright e_2$

$$e_1 \xrightarrow{t}_H h \Rightarrow e_2 \xrightarrow{t'}_H h', \text{ for some } h, h' \text{ such that } t \geq t' \text{ and } (h \triangleright h').$$

Using the same approach to that of cost simulation, it can be shown that \triangleright is a (pre)congruence,¹¹ and consequently that for all contexts C ,

$$e \triangleright e' \implies \langle C[e] \rangle^\alpha \geq \langle C[e'] \rangle^\alpha,$$

which can be restated as ‘ e' is at least as efficient as e in any context’. Using the same proof technique as illustrated in Section 7, we have a systematic means of verifying refinement laws, such as

$$\text{append}(\text{append}(e_1, e_2), e_3) \triangleright \text{append}(e_1, \text{append}(e_2, e_3)).$$

The notion of refinement is a possible semantic criterion (albeit a somewhat exacting one) for the intensional correctness of ‘context-free’ program transformations.

The main foundation for such theories of improvement and equivalence is the definition of an appropriate simulation relation, together with a proof that it satisfies the substitutivity property. So for each variation in the language (such as the addition of new operators) and each new definition of what improvement means, we require these constructions. This is somewhat tedious, so in a separate study[38] the problem of finding a more general formulation of the theory of improvement relations is addressed. For a general class of lazy languages, it is shown how a preorder on computational properties (an ‘improvement’ ordering) induces a simulation-style preordering on expressions (the definitions of cost simulation and program refinement are simple instances). Borrowing some syntactic conventions and semantic techniques from [18], some *improvement extensionality* conditions on the operators of the language are given which guarantee that the improvement ordering is a precongruence. The conditions appear relatively easy to check. Furthermore, the higher-order language given here is studied as a special case. In this context a *computational property* is considered to be a function from the proof of an evaluation judgement (the computation) to a preordered set (the set of properties, ordered by ‘improvement’). The main result of this generalization is that the simulation preorder induced by any computational property is guaranteed to be a congruence whenever the property satisfies a simple monotonicity

¹¹Note that if we consider a strict refinement relation where we replace \geq by $>$ in the definition of \triangleright , then it will *not* be a contextual congruence.

requirement with respect to the rules of the operational semantics. This result can be generalized to Howe’s class of structured computation systems [19].

12 Conclusions

This paper has presented a direct approach to reasoning about time cost under lazy evaluation. The calculus takes the form of *time rules* extracted from a suitably simple operational semantics, together with some equivalence laws which are substitutive with respect to these rules.

The aim of this calculus is to reveal enough of the ‘algorithmic structure’ of operationally opaque lazy functional programs to permit the use of the more traditional techniques developed in the context of the analysis of imperative programs [2], and initial experiments with this calculus suggest that it is of both practical and pedagogical use.

A desire for substitutive equivalences for the time rules led to a theory of cost equivalence, via a non-standard notion of operational approximation called *cost simulation*. Cost equivalence provides useful extensions to the time rules (although, as we showed, from a technical point of view it subsumes them). It is also interesting in its own right, since it suggests an operationally based route to the study of intensional semantics. Initial investigations of this area are reported in [38].

Finally, we proposed an interface of this calculus with a more compositional method which also improves the accuracy of the analysis with respect to call-by-need evaluation, but is able to retain the simplicity of the naïve approach where appropriate. The compositionality is based on an operational interpretation of evaluators to re-order computation on the basis of strictness-like properties. This model can also be used to provide a more constructive formulation of the evaluation transformer theorem [10], which formally connects information from strictness analysis with its associated optimisations. Further work is needed to strengthen the relationship between cost equivalence and the compositional approach. The idea of a *context-dependent bisimulation* between processes, as studied by Larsen (originating with [21]), seems appropriate here since it suggests the introduction of context (= demand)-dependent cost simulation.

Acknowledgements

This work has been partially funded by ESPRIT II BRA 6809 (Semantique), The Danish Natural Sciences Research Council (DART project), The Department of Computer Science at the University of Copenhagen (DIKU), and ESPRIT III BRA 9102 (Coordination).

Thanks for various contributions to Chris Hankin, Sebastian Hunt, Bent Thomsen (while the author was at Imperial College), and to various members of the *Semantique* project at that time for numerous suggestions and constructive criticisms; in particular to John Hughes for suggesting that time analysis might be completely expressed in terms of cost equivalences. More recently, thanks to the TOPPS group at DIKU and to Torben Mogensen and Kristian Nielsen in particular. Special thanks to Richard Bird, who originally suggested to me that a naïve approach could actually be useful for reasoning about call-by-name evaluation, and most recently suggested a number of improvements that have been incorporated into this version, including the maxtails example, offered by way of a ‘challenge problem’.

References

- [1] S. Abramsky. The lazy lambda calculus. In *Research Topics in Functional Programming*, D. Turner, ed., pp. 65–116.

- Addison Wesley, 1990.
- [2] A. V. Aho, J. E. Hopcroft and J. D. Ullman. *The Design and Analysis of Computer Algorithms*. The Addison-Wesley Series in Computer Science and Information Processing. Addison-Wesley, London, 1974.
 - [3] Z. M. Ariola and Arvind. A syntactic approach to program transformation. In *Proceedings of the symposium on Partial Evaluation and Semantics-Based Program Manipulation*, pp. 116–129. ACM press, SIGPLAN notices, 26(9), September 1991.
 - [4] H. P. Barendregt. *The Lambda Calculus*, vol. 103 of *Studies in Logic and the Foundations of Mathematics*. Elsevier Science Publishers B.V., 2nd edition, 1984.
 - [5] H. P. Barendregt, M. C. J. D. van Eekelen, J. R. W. Glauert, J. R. Kennaway, M. J. Plasmeijer and M. R. Sleep. Term graph rewriting. In *PARLE '87 volume II*, number 259 in LNCS, pp. 191–231. Springer Verlag, 1987.
 - [6] R. Bird and P. Wadler. *Introduction to Functional Programming*. Prentice Hall, 1988.
 - [7] B. Bjerner. *Time Complexity of Programs in Type Theory*. PhD thesis, Chalmers University of Technology, 1989.
 - [8] B. Bjerner and S. Holmström. A compositional approach to time analysis of first order lazy functional programs. In *Functional Programming Languages and Computer Architecture, Conference Proceedings*, pp. 157–165. ACM press, 1989.
 - [9] S. Brookes and S. Geva. Computational comonads and intensional semantics. In *Categories in Computer Science*, M. P. Fourman, P. T. Johnstone, and A. M. Pitts, eds, London Mathematical Society Lecture Notes, pp. 1–44. Cambridge University Press, 1992.
 - [10] G. L. Burn. The evaluation transformer model of reduction and its correctness. In *TAPSOFT '91*. Number 494 in LNCS, pp. 458–482. Springer-Verlag, 1991.
 - [11] G. L. Burn, C. L. Hankin and S. Abramsky. The theory and practice of strictness analysis for higher order functions. *Science of Computer Programming*, **7**, 249–278, 1986.
 - [12] R. Burstall and J. Darlington. A transformation system for developing recursive programs. *JACM*, **24**, 44–67, 1977.
 - [13] P. Cousot and R. Cousot. Constructive versions of Tarski's fixed point theorems. *Pacific Journal of Mathematics*, **82**, 43–57, 1979.
 - [14] P. Dybjer. Inverse image analysis generalises strictness analysis. *Information and Computation*, **90**, 194–216, 1991.
 - [15] R. L. Graham, D. E. Knuth, and O. Patashnik. *Concrete Mathematics*. Addison-Wesley, 1989.
 - [16] D. Gurr. *Semantic Frameworks for Complexity*. PhD thesis, Department of Computer Science, Edinburgh, 1991. (Available as reports CST-72-91 and ECS-LFCS-91-130).
 - [17] T. Hickey and J. Cohen. Automating program analysis. *JACM*, **35**, 185–220, 1988.
 - [18] D. J. Howe. Equality in lazy computation systems. In *Fourth annual symposium on Logic In Computer Science*, pp. 198–203. IEEE, 1989.
 - [19] D. J. Howe. On computational open-endedness in Martin-Löf's type theory. In *Sixth annual symposium on Logic In Computer Science*, pp. 162–172, 1991.
 - [20] D. E. Knuth. *Volume 1: Fundamental Algorithms*. The Art of Computer Programming. Addison-Wesley, 1968.
 - [21] K. G. Larsen. *Context-Dependent Bisimulation Between Processes*. PhD thesis, Department of Computing, University of Edinburgh, 1986.
 - [22] D. LeMétayer. An automatic complexity evaluator. *ACM ToPLaS*, 10(2):248–266, April 1988.
 - [23] G. Lindstrom. Static evaluation of functional programs. In *ACM SIGPLAN Symposium on Compiler Construction*, pp. 196–206, 1986.
 - [24] I. Mason. *The Semantics of Destructive Lisp*. Number 5 in CSLI Lecture Notes. CSLI, 1986.
 - [25] I. Mason and C. Talcott. Equivalence in functional languages with effects. *Journal of Functional Programming*, **1**, 287–327, 1991.
 - [26] R. Milner. Calculi for synchrony and asynchrony. *Theoretical Computer Science*, **25**, 267–310, 1983.
 - [27] R. Milner and M. Tofte. Co-induction in relational semantics. *Theoretical Computer Science*, **87**, 209–220, 1991.
 - [28] E. Moggi. Computational lambda-calculus and monads. Technical Report ECS-LFCS-88-63, Department of Computer Science, University of Edinburgh, 1988. (Short version in Fourth LICS Conf. IEEE 1989.)
 - [29] C.-H. Luke Ong. *The Lazy Lambda Calculus: An Investigation into the Foundations of Functional Programming*. PhD thesis, Imperial College, University of London, 1988.
 - [30] D. Park. Concurrency and automata on infinite sequences. In *5th GI Conference on Theoretical Computer Science*. LNCS 104, Springer Verlag, 1980.
 - [31] S. L. Peyton Jones. *The Implementation of Functional Programming Languages*. Prentice-Hall International Series in Computer Science. Prentice-Hall International (UK), London, 1987.

- [32] A. Pitts. A co-induction principal for recursively defined domains. Technical Report 252, University of Cambridge Computer Laboratory, April 1992.
- [33] M. Rosendahl. Automatic complexity analysis. In *Functional Programming Languages and Computer Architecture, Conference Proceedings*, pp. 144–156. ACM press, 1989.
- [34] D. Sands. Complexity analysis for a higher order language. Technical Report DOC 88/14, Imperial College, October 1988.
- [35] D. Sands. Complexity analysis for a lazy higher-order language. In *Proceedings of the Glasgow Workshop on Functional Programming*, Workshops in Computing Series, pp. 56–79. Springer Verlag, August 1989.
- [36] D. Sands. *Calculi for Time Analysis of Functional Programs*. PhD thesis, Department of Computing, Imperial College, University of London, September 1990.
- [37] D. Sands. Complexity analysis for a lazy higher-order language. In *Proceedings of the Third European Symposium on Programming*. Number 432 in LNCS, pp. 361–376. Springer-Verlag, May 1990.
- [38] D. Sands. Operational theories of improvement in functional languages (extended abstract). In *Proceedings of the Fourth Glasgow Workshop on Functional Programming*, pp. 298–311, Skye, August 1991. Springer Workshop Series.
- [39] D. Sands. Time analysis, cost equivalence and program refinement. In *Proceedings of the Eleventh Conference on Foundations of Software Technology and Theoretical Computer Science*. Number 560 in LNCS, pp. 25–39. Springer-Verlag, December 1991.
- [40] D. Sands. Total correctness and improvement in the transformation of functional programs. DIKU, University of Copenhagen, Unpublished (53 pages), May 1994.
- [41] C. L. Talcott. *The Essence of Rum, A Theory of the Intensional and Extensional Aspects of Lisp-type Computation*. PhD thesis, Stanford University, August 1985.
- [42] P. Wadler. Strictness analysis aids time analysis. In *15th ACM Symposium on Principals of Programming Languages*, pp. 119–132, January 1988.
- [43] P. Wadler and R. J. M. Hughes. Projections for strictness analysis. In *1987 Conference on Functional Programming and Computer Architecture*, pp. 385–407, Portland, Oregon, September 1987.
- [44] B. Wegbreit. Mechanical program analysis. *CACM*, **18**, 528–539, 1975.

Appendices

A Proofs

This appendix contains some details of proofs not included in the main text; the corresponding theorems etc. are restated.

PROPOSITION A.1 (PROPOSITION 6.7)

For all closed e, e'

1. $e \rightarrow_N c \iff e \rightarrow_H c$.
2. $e \xrightarrow{t}_N v \iff e \xrightarrow{t_1}_H h$ and $h \xrightarrow{t_2}_N v$ and $t_1 + t_2 = t$ for some h .

PROOF. 1. Follows by routine inductions on the structure of the proofs of $e \rightarrow_N c$ and $e \rightarrow_H c$.

2. (\Rightarrow) Induction on the structure of the proof of $e \rightarrow_N n$. We give an illustrative case:

$\text{Fun} \frac{e_i\{e_1/x_1, \dots, e_{n_i}/x_{n_i}\} \rightarrow_N v}{f_i(e_1, \dots, e_{n_i}) \rightarrow_N v}$
--

Suppose that $f_i(e_1, \dots, e_{n_i}) \xrightarrow{t}_N v$ for some t . Then it follows that

$$e_i\{e_1 \dots e_{n_i}/x_1 \dots x_{n_i}\} \xrightarrow{t-1}_N v,$$

and so the induction hypothesis gives

$$e_i\{e_1 \dots e_{n_i}/x_1 \dots x_{n_i}\} \xrightarrow{t_1}_H h \tag{A.1}$$

$$h \xrightarrow{t_2}_N v \tag{A.2}$$

42 A Naïve Time Analysis and its Theory of Cost Equivalence

for some head-normal-form h , and t_1, t_2 such that $t_1 + t_2 = t - 1$. We can conclude from A.1 by the semantics for application, $f_i(e_1, \dots, e_{n_i}) \xrightarrow{1+t_1}_N v$ which together with A.2 concludes this case.

The other cases follow in a similar manner. Proof in the \Leftarrow direction is again a routine induction on the structure of the inference $e \rightarrow_H h$. \blacksquare

LEMMA A.2 (LEMMA 6.8)

If $e \preceq e'$ then

1. if $e \rightarrow_H h$ then $e' \rightarrow_H h'$ and $\langle e \rangle^H = \langle e' \rangle^H$;
2. if $e \rightarrow_N v$ then $e' \rightarrow_N v$ and $\langle e \rangle^N = \langle e' \rangle^N$.

PROOF. The first part is immediate from the definition of cost simulation. The second part proceeds by induction on the structure of value v .

- $v \equiv c$ Suppose $e \xrightarrow{t}_N c$. Then by Proposition 6.7 $e \xrightarrow{t}_H c$. By definition of cost simulation, $e' \xrightarrow{t}_H c$, and again from Proposition 6.7 $e' \xrightarrow{t}_N c$.
- $v \equiv v_1 : v_2$ Suppose $e \xrightarrow{t}_N v_1 : v_2$. The induction hypothesis gives: for all t' and for all e_a, e_b such that $e_a \preceq e_b$, $e_a \xrightarrow{t'}_N v_1$ implies $e_b \xrightarrow{t'}_N v_1$ and $e_a \xrightarrow{t'}_N v_2$ implies $e_b \xrightarrow{t'}_N v_2$. From Proposition 6.7 we know that, for some e_1, e_2 , $e \xrightarrow{t'}_H e_1 : e_2$ and $e_1 : e_2 \xrightarrow{t''}_N v_1 : v_2$ with $t = t' + t''$. This implies that

$$e_1 \xrightarrow{t_1}_N v_1 \text{ and } e_2 \xrightarrow{t_2}_N v_2, \text{ where } t'' = t_1 + t_2. \quad (\text{A.3})$$

Since $e \preceq e'$ we know that $e' \rightarrow_H e'_1 : e'_2$ for some e'_1, e'_2 such that $e_1 \preceq e'_1$ and $e_2 \preceq e'_2$. Now we have, by an instance of the induction hypothesis that

$$e'_i \xrightarrow{t_i}_N v_i, i = 1, 2 \quad (\text{A.4})$$

from which we have, by rule **Cons** that $e'_1 : e'_2 \xrightarrow{t''}_N v_1 : v_2$. Finally $e' \xrightarrow{t}_N v_1 : v_2$ follows from Proposition 6.7 (\Leftarrow). \blacksquare

THEOREM A.3 (THEOREM 7.4)

For all closed expressions e , if $e \xrightarrow{m}_H h$ then $\vdash_K e = () \text{ l}^m(h)$.

PROOF. Routine induction on the structure of the proof of $e \xrightarrow{m}_H h$. We argue by cases according to the last rule applied. We sketch a couple of cases only.

Fun

Assume that $e \equiv f(e_1 \dots e_n)$, where f is defined by $f(y_1 \dots y_n) = e_f$. Then the last inference in the proof of $e \xrightarrow{m}_H h$ has the form:

$$\frac{e_f\{e_1/y_1 \dots e_n/y_n\} \rightarrow_H h}{f(e_1, \dots, e_n) \rightarrow_H h}.$$

It follows that $e_f\{e_1/y_1 \dots e_n/y_n\} \xrightarrow{m-1}_H h$, and so by the induction hypothesis,

$$\vdash_K e_f\{e_1/y_1 \dots e_n/y_n\} = () \text{ l}^{m-1}(h).$$

By congruence it follows that $\vdash_K \text{l}(e_f\{e_1\sigma/y_1 \dots e_n\sigma/y_n\}) = () \text{ l}^m(h)$, and finally by **(Fun.)** and transitivity we have $\vdash_K f(e_1 \dots e_n) = () \text{ l}^m(h)$.

Case

We consider only one of the two possible last rules in the evaluation of the case expression: assume the last inference has the form

$$\frac{e_1 \rightarrow_H e_h : e_t \quad e_3\{e_h/x, e_t/xs\} \rightarrow_H h}{\left(\begin{array}{c} \text{case } e_1 \text{ of} \\ \text{nil} \Rightarrow e_2 \\ x : xs \Rightarrow e_3 \end{array} \right) \rightarrow_H h}.$$

It follows that $e_1 \xrightarrow{m_1}_H e_h : e_t$ and $e_3 \{e_h/x, e_t/xs\} \xrightarrow{m_2}_H h$ for some m_1, m_2 satisfying $m_1 + m_2 = m$. The induction hypothesis then gives

$$\begin{aligned} & \vdash_K e_1 =_{\langle \rangle} I^{m_1}(e_h : e_t) \\ \Rightarrow & \vdash_K \left(\begin{array}{c} \text{case } e_1 \text{ of} \\ \text{nil} \Rightarrow e_2 \\ x : xs \Rightarrow e_3 \end{array} \right) =_{\langle \rangle} \left(\begin{array}{c} \text{case } I^{m_1}(e_h : e_t) \text{ of} \\ \text{nil} \Rightarrow e_2 \\ x : xs \Rightarrow e_3 \end{array} \right) \\ \Rightarrow & \vdash_K \left(\begin{array}{c} \text{case } e_1 \text{ of} \\ \text{nil} \Rightarrow e_2 \\ x : xs \Rightarrow e_3 \end{array} \right) =_{\langle \rangle} I^{m_1} \left(\begin{array}{c} \text{case } e_h : e_t \text{ of} \\ \text{nil} \Rightarrow e_2 \\ x : xs \Rightarrow e_3 \end{array} \right) \\ \Rightarrow & \vdash_K \left(\begin{array}{c} \text{case } e_1 \text{ of} \\ \text{nil} \Rightarrow e_2 \\ x : xs \Rightarrow e_3 \end{array} \right) =_{\langle \rangle} I^{m_1}(e_3 \{e_h/x, e_t/xs\}). \end{aligned}$$

The induction hypothesis applied to the second sub-proof gives

$$\vdash_K e_3 \{e_h/x, e_t/xs\} =_{\langle \rangle} I^{m_2}(h)$$

and so by congruence (substituting the right-hand side for the left in the previous equation) we conclude that

$$\vdash_K \left(\begin{array}{c} \text{case } e_1 \text{ of} \\ \text{nil} \Rightarrow e_2 \\ x : xs \Rightarrow e_3 \end{array} \right) =_{\langle \rangle} I^{m_1}(I^{m_2}(h)),$$

and the desired result follows from the fact that $I^{m_1}(I^{m_2}(h)) \equiv I^m(h)$. ■

THEOREM A.4 (THEOREM 8.3)

If $\tilde{e} \sqsubseteq \tilde{e}'$ for some commonly indexed families of closed expressions \tilde{e}, \tilde{e}' , then for all expressions e containing at most variables \tilde{x}

$$e\{\tilde{e}/\tilde{x}\} \sqsubseteq e\{\tilde{e}'/\tilde{x}\}.$$

PROOF. Define the relation

$$\mathcal{I} = \{(e\{\tilde{e}/\tilde{x}\}, e\{\tilde{e}'/\tilde{x}\}) : e \text{ contains at most variables } \tilde{x}, \tilde{e} \sqsubseteq \tilde{e}'\}.$$

It is sufficient to show that $\mathcal{I} \subseteq \mathcal{A}(\mathcal{I})$. Abbreviate substitutions $\{\tilde{e}/\tilde{x}\}$ and $\{\tilde{e}'/\tilde{x}\}$ by σ and σ' , respectively. Assume that $e\sigma \xrightarrow{t}_H h$. It will be sufficient to prove that $e\sigma' \xrightarrow{t}_H h'$ for some h' such that $h(\mathcal{I} \cup \mathcal{I}^\lambda)h'$. We establish this by induction on the structure of the proof of $e\sigma \rightarrow_H h$, and by cases according to the structure of expression e . The cases for all expressions other than lambda abstraction and application are identical to Lemma 6.10. The remaining cases are:

$e \equiv \lambda y.b$ Assume that y is not in the domain of σ (since if it is, we can just rename). By the axiom for lambda expressions, together with its time rule, we have $e\sigma \equiv \lambda y.(b\sigma) \xrightarrow{0}_H \lambda y.(b\sigma)$ and $e\sigma' \equiv \lambda y.(b\sigma') \xrightarrow{0}_H \lambda y.(b\sigma')$. It remains to show that $\lambda y.(b\sigma) \mathcal{I}^\lambda \lambda y.(b\sigma')$. This follows easily from the fact that for all closed a ,

$$(\lambda y.(b\sigma))a \equiv ((\lambda y.b)a)\sigma \mathcal{I} ((\lambda y.b)a)\sigma' \equiv (\lambda y.(b\sigma'))a.$$

$e \equiv e_1 e_2$. By the rule for application it follows that $e_1\sigma \xrightarrow{t_1}_H \lambda y.b$ and that

$$b\{e_2\sigma/y\} \xrightarrow{t_2}_H h$$

for some t_1, t_2 such that $t = 1 + t_1 + t_2$. From the former judgement the inductive hypothesis allows us to conclude that $e_1\sigma' \xrightarrow{t_1}_H \lambda y.b'$ for some b' such that $\lambda y.b \mathcal{I}^\lambda \lambda y.b'$.

So, for each closed expression a , we have substitutions δ, δ' (with a non-empty common domain, and \sqsubseteq -related range) such that $(\lambda y.b)a \equiv d\delta$ and $(\lambda y.b')a \equiv d\delta'$ for some expression d containing at most the variables in the domain of δ . Since a is closed, we can without loss of generality assume that the structure of d satisfies one of the following three cases:

44 A Naïve Time Analysis and its Theory of Cost Equivalence

1. $d \equiv x$, and hence $(\lambda y.b)a \sqsubseteq (\lambda y.b')a$, or
2. $d \equiv xa$ and hence $\lambda y.b \sqsubseteq \lambda y.b'$, or
3. $d \equiv (\lambda y.c)a$, and hence $(\lambda y.b) \mathcal{I} (\lambda y.b')$.

Now consider the judgement $b\{e_2\sigma/y\} \xrightarrow{t_2}_H h$. By the rule for application this is equivalent to the judgement $(\lambda y.b)e_2\sigma \xrightarrow{1+t_2}_H h$. By the case analysis above, taking $a \equiv e_2\sigma$, either

1. $(\lambda y.b)a \sqsubseteq (\lambda y.b')a$, and so by the definition of \sqsubseteq , $(\lambda y.b)e_2\sigma \xrightarrow{1+t_2}_H h'$ with $h(\mathcal{R} \cup \mathcal{R}^\lambda)h'$, and hence $b\{e_2\sigma/y\} \xrightarrow{t_2}_H h'$, or
2. $(\lambda y.b) \sqsubseteq (\lambda y.b')$ and so by definition $(\lambda y.b)a \sqsubseteq (\lambda y.b')a$ and we argue as above, or
3. $(\lambda y.b) \mathcal{I} (\lambda y.b')$ with $(\lambda y.b)a \equiv d\delta$ and $(\lambda y.b')a \equiv d\delta'$. Assume, without loss of generality, that the variables in σ , δ and $\{y\}$ are all distinct. Massaging substitutions (details omitted), it is easy to see that

$$b\{e_2\sigma/y\} \equiv d\delta\{e_2\sigma/y\} \equiv d\{e_2/y\}\delta\sigma \mathcal{I} d\{e_2/y\}\delta'\sigma' \equiv d\delta\{e_2\sigma/y\} \equiv b'\{e_2\sigma/y\}.$$

Since this shows that $b\{e_2\sigma/y\} \mathcal{I} b'\{e_2\sigma/y\}$, by the induction hypothesis in this case we can also conclude that $b'\{e_2\sigma/y\} \xrightarrow{t_2}_H h'$ for some h' such that $h(\mathcal{R} \cup \mathcal{R}^\lambda)h'$.

Using rule for application we conclude that $e\sigma' \xrightarrow{t}_H h'$ for some h' such that $h(\mathcal{R} \cup \mathcal{R}^\lambda)h'$, as required. \blacksquare

B The largest cost precongruence

In this appendix we give the details of the technical development outlined in Section 6.5, which characterizes when cost simulation is the largest cost congruence.

DEFINITION B.1

$$e_1 \leq_c e_2 \iff \forall C.C[e_1] \mathcal{F}(\top) C[e_2].$$

THEOREM B.2

$$\preceq = \leq_c$$

Since \preceq is the maximal fixed point of \mathcal{F} , to prove the ‘difficult’ half of this equality it is sufficient to show that $\leq_c \subseteq \mathcal{F}(\preceq_c)$. However, we have been unable to prove this property by ‘co-induction’ (although the proof for pure operational simulation is straightforward [18]). Instead, we will use an alternative definition of cost simulation as the limit of a descending sequence of relations beginning with \top :

PROPOSITION B.3

$$\preceq = \bigcap_{n \in \omega} \mathcal{F}^n(\top)$$

where \mathcal{F}^0 is the identity function, and $\mathcal{F}^{n+1} = \mathcal{F}^n \circ \mathcal{F}$

PROOF. It is sufficient to show that \mathcal{F} is *anti-continuous* (see e.g. [13]). That is, for every decreasing sequence $R_1 \supseteq R_2 \supseteq \dots \supseteq R_n \supseteq \dots$,

$$\bigcap_n \mathcal{F}(R_n) = \mathcal{F}\left(\bigcap_n R_n\right).$$

The \supseteq half follows directly from monotonicity. Now suppose that $e(\bigcap_n \mathcal{F}(R_n))e'$. The only non-trivial case is when $e \xrightarrow{t}_H e_1 : e_2$: for each n , we have that $e\mathcal{F}(R_n)e'$, so from the definition of \mathcal{F} , we have that $e' \xrightarrow{t}_H e'_1 : e'_2$ with $e_i(R_n)e'_i$, and hence that $e_i(\bigcap_n (R_n))e'_i$, giving $e(\mathcal{F}(\bigcap_n R_n))e'$ as required. \blacksquare

Now we sketch the proof of theorem B.2.

PROOF. $\preceq \subseteq \leq_c$ follows from the fact that \preceq is a precongruence, and that \mathcal{F} is monotone. It remains to show that that $\preceq \supseteq \leq_c$. We show the contrapositive: that for all expressions e_1, e_2 ,

$$e_1 \not\preceq e_2 \Rightarrow e_1 \not\leq_c e_2.$$

For each $n \geq 0$, define $\preceq_n = \mathcal{F}^n(\top)$. Since the \preceq_n form a decreasing chain, by Proposition B.3 it follows that there exists a smallest $k > 0$ (and hence a largest \preceq_k) such that $e_1 \not\preceq_k e_2$. Call such a \preceq_k the *maximum distinguisher* for (e_1, e_2) . Negating the definition of \leq_c , we see that $e_1 \not\leq_c e_2$ if there exists a context C such that $C[e_1] \not\preceq_1 C[e_2]$. Call such a context C a *distinguishing context* for (e_1, e_2) .

We prove by induction on k that for all e_1, e_2 , if \preceq_k the maximum distinguisher for (e_1, e_2) then there exists a distinguishing context C_{e_1, e_2} for (e_1, e_2) .

Base: \preceq_1 Since $e_1 \not\preceq_1 e_2$, it immediately follows that the simple context $[\]$ distinguishes (e_1, e_2) .

Induction: \preceq_{k+1} Since \preceq_{k+1} is maximum for (e_1, e_2) , it follows that $e_1 \xrightarrow{t}_H u_1$ and $e_2 \xrightarrow{t}_H u_2$. Since $k+1 > 1$, it follows, again by maximality, that $u_1 = p_1 : q_1$ and $u_2 = p_2 : q_2$ for some p_1, q_1, p_2, q_2 . Now \preceq_k must be a maximum distinguisher for either (p_1, p_2) or (q_1, q_2) , since if it were not, \preceq_{k+1} would not be maximal for (e_1, e_2) . Suppose that it is maximal for (p_1, p_2) (the other case is similar). By induction, there is a distinguishing context C_{p_1, p_2} , so we can easily construct a context

$$\begin{aligned} &\text{case } [\] \text{ of} \\ &\quad \text{nil} \Rightarrow \text{nil} \\ &\quad x : xs \Rightarrow C_{p_1, p_2}[x] \end{aligned}$$

which distinguishes (e_1, e_2) . ■

DEFINITION B.4

Pure cost congruence, \leq_{pc} , is defined to be the least relation on closed expressions such that $e_1 \leq_{pc} e_2$ if and only if for all contexts C , whenever $C[e_1] \xrightarrow{t}_H u_1$ then there exists u_2 such that $C[e_2] \xrightarrow{t}_H u_2$.

We will say that a context C *cost distinguishes* a pair of expressions e, e' if $C[e] \not\leq_{pc} C[e']$.

DEFINITION B.5

We say that the language satisfies the CD condition if for all constants c , if $e \rightarrow_H c$ and $e \leq_{pc} e'$ then $e' \rightarrow_H c$.

PROPOSITION B.6

$\leq_c = \leq_{pc} \iff$ the CD condition is satisfied.

PROOF. (\Rightarrow) Straightforward from the definition of \leq_c . (\Leftarrow) Clearly $e_1 \leq_c e_2 \Rightarrow e_1 \leq_{pc} e_2$, so it is sufficient to show that $e_1 \not\leq_c e_2 \Rightarrow e_1 \not\leq_{pc} e_2$. Supposing $e_1 \not\leq_c e_2$, then there exists a distinguishing context C for (e_1, e_2) . Now suppose that C is *not* sufficient to *cost* distinguish (e_1, e_2) . In this case we must have $C[e_1] \xrightarrow{t}_H u_1$ and $C[e_2] \xrightarrow{t}_H u_2$ with either

1. $u_1 = c, u_2 \neq c$, or
2. $u_1 = u : u', u_2 = c$ for some constant c .

In the first case $C[e_1] \not\leq_{pc} C[e_2]$ follows from the CD condition, and since \leq_{pc} is a precongruence, we must have $e_1 \not\leq_{pc} e_2$. In the second case a context of the form $\text{case } [\] \text{ of } \dots$ can be used to *cost* distinguish $(C[e_1], C[e_2])$, and hence $e_1 \not\leq_{pc} e_2$ by precongruence. ■

Remark The CD condition is not particularly strong since it is satisfied if a binary equality test over constants is included as one of the primitive functions; for each constant c , consider the context $\text{if } [\] = c \text{ then nil else fail}$, where *fail* is any expression lacking a head normal form. On the other hand, the CD condition fails if we have two distinct constants in the language, call them *error1* and *error2*, over which all primitive functions are undefined (i.e. $\text{APPLY}(p, \dots, \text{error}_i, \dots)$ is undefined). So we could never distinguish between *error1* and *error2* on the basis of cost alone, but because they are distinct they would not be cost equivalent.

Now we have that cost equivalence is the largest congruence with respect to the time rules for head-normal-form (modulo the partial correctness of the rules, which is reflected by convergence conditions).

COROLLARY B.7

If the CD condition is satisfied, then

$$e =_{\emptyset} e' \iff (\forall C. (C[e_1] \downarrow_H \ \& \ C[e_2] \downarrow_H) \Rightarrow \langle C[e_1] \rangle^H = \langle C[e_2] \rangle^H).$$

46 *A Naïve Time Analysis and its Theory of Cost Equivalence*

Lemma 6.7, together with the fact that $\equiv_{\langle \rangle}$ is a congruence allows us to extend the above corollary to include evaluation to normal-form, and the $\langle \rangle^N$ rules. In fact it is also possible to show that we can completely replace ' H ' by ' N ' in the corollary, but this is left as an exercise.

Received 4 May 1993