

# Functional Programming in Sublinear Space

Ugo Dal Lago<sup>1</sup> and Ulrich Schöpp<sup>2,\*</sup>

<sup>1</sup> University of Bologna, Italy

<sup>2</sup> University of Munich, Germany

**Abstract.** We consider the problem of functional programming with data in external memory, in particular as it appears in sublinear space computation. Writing programs with sublinear space usage often requires one to use special implementation techniques for otherwise easy tasks, e.g. one cannot compose functions directly for lack of space for the intermediate result, but must instead compute and recompute small parts of the intermediate result on demand. In this paper, we study how the implementation of such techniques can be supported by functional programming languages.

Our approach is based on modeling computation by interaction using the Int construction of Joyal, Street & Verity. We derive functional programming constructs from the structure obtained by applying the Int construction to a term model of a given functional language. The thus derived functional language is formulated by means of a type system inspired by Baillot & Terui's Dual Light Affine Logic. We assess its expressiveness by showing that it captures LOGSPACE.

## 1 Introduction

A central goal in programming language theory is to design programming languages that allow a programmer to express efficient algorithms in a convenient way. The programmer should be able to focus on algorithmic issues as much as possible and the programming language should give him the means to delegate inessential implementation details to the computer.

In this paper we study programming language constructs that are useful for expressing sublinear space algorithms. Sublinear space algorithms use less memory space than would be needed to store their input. They are useful for computing with large inputs that do not fit into memory.

When writing programs with sublinear space usage, one must often use special techniques for tasks that would normally be simple. A typical example is the composition of two algorithms. In order to remain in sublinear space, one cannot run them one after the other, as there may not be enough space to store the intermediate result. Instead, one can implement composition by storing at any time only a small part of the intermediate value and by (re)computing small parts only as they are needed.

Since it is easy to make mistakes in the implementation of such on-demand recomputation of intermediate values, we believe that programming language support should be very useful for such tasks. Instead of implementing composition with on-demand

---

\* This work was carried out while Ulrich Schöpp was supported by a fellowship of the Institute of Advanced Studies at the University of Bologna.

recomputation by hand, the programmer should be able to write function composition in the usual way and have a compiler generate the complicated program.

The possibilities of doing this have been explored in work on implicit characterisations of LOGSPACE. A number of characterisations of this complexity class have been explored in terms of function algebras [14] and linear logics [15] (there is more work for LOGSPACE-predicates, e.g. [10,5], but we focus on functions here). However, these characterisations are still far away from being real programming languages.

Here we work towards the goal of making the abstractions explored in this line of work usable in programming. In contrast to previous work [14,15], we aim to enrich the language with constructs for working with on-demand recomputation conveniently, rather than hiding it completely. This is in line with the fact that sublinear space algorithms are usually not used in isolation and will generally appear within larger programs. Language support for writing sublinear space algorithms should not become a hindrance in other parts of the program that do not operate on large data.

In this paper we present the functional programming language INTML for programming with sublinear space. This language is derived from an instance of the Int construction [7]. Our thesis is that the Int construction naturally captures space bounded computation and thus exposes mathematical structure that is useful for writing space bounded programs. Even though the type system of INTML is quite simple when compared to earlier higher-order type systems for LOGSPACE [15], INTML allows LOGSPACE algorithms to be written in a natural way, as is discussed in Sec. 3.2.

## 2 Space-Bounded Computation

Our approach to representing sublinear space computation in a functional programming language is best explained by analysing the definition of space complexity classes.

In the definition of space complexity classes, in particular sublinear space complexity classes, one uses *Offline Turing Machines* (OTMs) instead of standard Turing Machines. Offline Turing Machines are multi-tape Turing Machines that differ from the standard ones in that the input tape is read-only, the output tape is write-only and the output head may be moved in one direction only; finally the input and output tapes do not count towards the space usage of an Offline Turing Machine.

The definition of Offline Turing Machine captures a special class of Turing Machines that do not store their input or output in memory, but instead have (random) access to some externally stored input, and that give their output as a stream of characters. Since neither input nor output must be stored in memory, it is justified to count only the work tape towards the space usage of an Offline Turing Machine.

More formally, while a normal Turing Machine computes a function  $\Sigma^* \rightarrow \Sigma^*$  on words over an alphabet  $\Sigma$ , an Offline Turing Machine may be seen as a function of type

$$(\text{State} \times \Sigma) + \mathbb{N} \longrightarrow (\text{State} \times \mathbb{N}) + \Sigma ,$$

where  $A + B$  denotes the (tagged) disjoint union  $\{\text{inl}(x) \mid x \in A\} \cup \{\text{inr}(y) \mid y \in B\}$ . An input  $n \in \mathbb{N}$  stands for the request to compute the  $n$ -th character on the OTM's output tape. An output in  $\Sigma$  is a response to this request. Whenever the OTM wants to read a character from its input tape, it outputs a pair  $\langle s, n \rangle \in \text{State} \times \mathbb{N}$ , where  $n$  is the number

of the input character it wants to read and  $s$  is its machine state, comprising finite control state, work tape contents, etc. Receiving this request, the environment looks up the  $n$ -th input character  $i_n$  and restarts the machine with input  $\langle s, i_n \rangle \in \text{State} \times \Sigma$ . It supplies the machine state  $s$  that was part of the input request, so that the machine can resume its computation from the point where it requested an input.

In this way, we can consider each Offline Turing Machine as a normal Turing Machine with the special input/output interface given by the type above. This special machine needs space only to store the work tape(s) of the OTM and the positions of the OTM's input and output heads. This view justifies the exclusion of the input and output tapes in the definition of the space usage of OTMs, as long as the space usage is at least logarithmic. We will not consider OTMs with sublogarithmic space usage here and indeed 'Classes of languages accepted within sublogarithmic space depend heavily on the machine models and the mode of space complexity' [16].

While at first sight, the step from Turing Machines to Offline Turing Machines appears to be just a technicality, it is a step from unidirectional to bidirectional computation. For example, while standard Turing Machines are composed just by running them one after the other, composition of Offline Turing Machines involves bidirectional data flow. This composition is implemented as a dialogue between the machines: one starts by requesting an output character from the second machine and every time this machine queries a character of its input, the first machine is started to compute this character.

Bidirectional computation is thus an integral feature of space-bounded computation, which must be accounted for in a programming language for space-bounded functions. We argue that a good way of accounting for this bidirectionality is to study space-bounded computation in terms of the Int construction of Joyal, Street & Verity [7,6]. The Int construction is a general algebraic method of constructing a bidirectional universe from a unidirectional one. It appears in categorical formulations of the Geometry of Interaction [1,2] and has many applications, e.g. to Attribute Grammars [8].

## 2.1 Structuring Space Bounded Computation

In this section we observe that the step from Turing Machines to Offline Turing Machines can be understood in terms of the Int construction. This simple observation gives us guidance for structuring space-bounded computation, since it allows us to draw on existing work on the structure obtained by the Int construction.

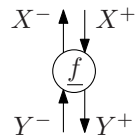
The idea is to start from a computational model, e.g. the partial computable functions as formalised by Turing Machines, and then apply the Int construction to this model. The result is a model that still contains the original one, but in addition also captures bidirectional (or interactive) computation in the style of Offline Turing Machines.

In general, the Int construction starts with a traced monoidal category  $\mathbf{B}$  and yields a category  $\text{Int}(\mathbf{B})$  that represents bidirectional computation in  $\mathbf{B}$ . For the argument in this paper, it suffices to describe just one particular instance of  $\text{Int}(\mathbf{B})$ , that where  $\mathbf{B}$  is the category  $\mathbf{Pfn}$  of sets and partial functions. In this example we drop computability for the sake of simplicity and assume that our 'computational' model consists just of partial functions between arbitrary sets.

If we apply the Int construction to  $\mathbf{Pfn}$  with respect to tagged disjoint union as monoidal structure, then we obtain the following category  $\text{Int}(\mathbf{Pfn})$ . Its objects are

pairs  $(X^-, X^+)$  of sets  $X^-$  and  $X^+$ . A morphism  $f$  from  $X = (X^-, X^+)$  to  $Y = (Y^-, Y^+)$  is a partial function  $\underline{f}: X^+ + Y^- \rightarrow Y^+ + X^-$ .

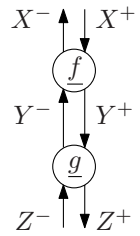
Morphisms capture bidirectional computation. Let us think of the partial function  $\underline{f}$  as a message-passing node with two input wires and two output wires as drawn on the right. When an input value arrives on an input wire then the function  $\underline{f}$  is applied to it and the resulting value is passed along the corresponding output wire.



We will combine the two edges for  $X^-$  and  $X^+$  into a single edge in which messages may be passed both ways (and likewise for  $Y$ ). Thus we obtain the node on the right, whose edges are bidirectional in the sense that an edge with label  $X$  allows any message from  $X^+$  to be passed in the forward direction and any message from  $X^-$  to be passed in the backwards direction.



Composition in  $\text{Int}(\mathbf{Pfn})$  allows one to build message passing networks out of such nodes. The composition  $g \circ f: X \rightarrow Z$  of  $f: X \rightarrow Y$  and  $g: Y \rightarrow Z$  is obtained simply by connecting the two nodes. The underlying partial function  $\underline{g \circ f}: X^+ + Z^- \rightarrow Z^+ + X^-$  is most easily described in terms of message passing. An input in  $X^+$  is given to  $\underline{f}$  and one in  $Z^-$  to  $\underline{g}$ . If either  $\underline{f}$  or  $\underline{g}$  give an output in  $X^-$  or  $Z^+$  then this is the output of  $\underline{g \circ f}$ . If, however,  $\underline{f}$  (resp.  $\underline{g}$ ) outputs a message on  $Y^+$  (resp.  $Y^-$ ), this message is given as an input to  $\underline{g}$  (resp.  $\underline{f}$ ). This may lead to a looping computation and  $\underline{g \circ f}$  may be partial even if  $\underline{g}$  and  $\underline{f}$  are both total.



Offline Turing Machines appear as morphisms of type  $(\text{State} \times \mathbb{N}, \text{State} \times \Sigma) \rightarrow (\mathbb{N}, \Sigma)$  in  $\text{Int}(\mathbf{Pfn})$ . This follows immediately from the definition of morphisms and the discussion in Sec. 2. What we gain from viewing OTMs as morphisms in  $\text{Int}(\mathbf{Pfn})$  is that we can use the well-known structure of this category for constructing and manipulating them. For example,  $\text{Int}(\mathbf{Pfn})$  is compact closed and therefore allows us to use linear lambda calculus and higher-order functions for the manipulation of OTMs.

We list the structure in  $\text{Int}(\mathbf{Pfn})$  that we use in the definition of INTML.

**Partial Functions.** First we note that inside  $\text{Int}(\mathbf{Pfn})$  we still find  $\mathbf{Pfn}$ . For any set  $A$  we have an object  $\mathcal{I}A = (\emptyset, A)$ . A morphism from  $\mathcal{I}A$  to  $\mathcal{I}B$  is a partial function of type  $A + \emptyset \rightarrow \emptyset + B$ , so that the morphisms of that type are in one-to-one correspondence with the partial functions from  $A$  to  $B$ .

**Thunks.** Also useful is the object  $[A] = (\{*\}, A)$ , where a singleton replaces the empty set. We will use  $*$  as question that signals an explicit request for a value of type  $A$ . Thus, one may think of  $[A]$  as a type of thunks that are evaluated on demand.

**Higher-order Functions.**  $\text{Int}(\mathbf{Pfn})$  has a monoidal structure  $\otimes$  that on objects is defined by  $X \otimes Y = (X^- + Y^-, X^+ + Y^+)$ . In addition, there is a dualising operation  $(-)^*$  that exchanges question and answer sets, i.e.  $(X^-, X^+)^* = (X^+, X^-)$ . Together,  $\otimes$  and  $(-)^*$  make  $\text{Int}(\mathbf{Pfn})$  a compact closed category. As a consequence, we obtain a linear function space by letting  $X \multimap Y = X^* \otimes Y$ .

**Indexed Tensor.** Of further use is an indexed tensor product  $\bigotimes_A X$ . The object  $\bigotimes_A X$  is isomorphic to  $X \otimes \dots \otimes X$  ( $|A|$  times). It is defined by  $(\bigotimes_A X)^- = A \times X^-$  and  $(\bigotimes_A X)^+ = A \times X^+$ . The first component in the messages indicates which component of the tensor product we are communicating with.

To see how this structure is useful for working with OTMs, consider a morphism of type

$$\bigotimes_{\text{State}} (\mathcal{I}\mathbb{N} \multimap \mathcal{I}\Sigma) \longrightarrow (\mathcal{I}\mathbb{N} \multimap \mathcal{I}\Sigma) \quad (1)$$

in  $\text{Int}(\mathbf{Pfn})$ . By definition, it is a partial function from  $(\text{State} \times (\emptyset + \Sigma)) + (\mathbb{N} + \emptyset)$  to  $(\text{State} \times (\mathbb{N} + \emptyset)) + (\emptyset + \Sigma)$ , which, if we remove the superfluous empty sets, is the same as the type of an Offline Turing Machine, as described in Sec. 2.

As a morphism of type (1), an Offline Turing Machine is modelled simply as a map from inputs to outputs, which are both modelled as (linear) functions  $\mathcal{I}\mathbb{N} \multimap \mathcal{I}\Sigma$ . This encoding of words as functions reflects the fact that Offline Turing Machines do not have access to the whole input at once, but rather can read only a single character at a time. Reading the  $n$ -th character just corresponds to applying the input function of type  $\mathcal{I}\mathbb{N} \multimap \mathcal{I}\Sigma$  to the natural number  $n$ . In  $\text{Int}(\mathbf{Pfn})$  we may therefore use the input as if it were a single function from natural numbers to characters, even though in reality the input can only be queried character by character. We argue that this is more convenient than the access to the input by means of explicit questions as in the direct implementation of Offline Turing Machines. For example we can use  $\lambda$ -calculus to manipulate the input functions.

### 3 A Functional Language for Logarithmic Space

In the rest of this paper we develop a functional programming language INTML that is based on understanding space-bounded computation in terms of the Int construction. For the definition of INTML we start with a standard functional programming language. In order to program sublinear space algorithms in this language, we would like to implement message passing networks in it and we would like to manipulate these networks using the structure that we have found in the Int construction. We obtain INTML by extending the original programming language with primitives for constructing and manipulating such message passing networks.

INTML provides a syntax for the structure that one obtains from applying the Int construction to a *term model* of the standard functional programming that we start with. In terms of the outline above, one should replace all the partial functions in  $\mathbf{Pfn}$  with terms in the functional programming language. That means that the message passing networks captured by the Int construction are now not just partial functions but partial functions implemented in the functional programming language.

INTML is thus a typed functional programming language with two classes of terms and types: one for the functional language that we start with and one for the structure that we obtain by applying the Int construction to a term model of this language. We call the former the *working class* and the latter the *upper class*. The upper class part can be seen as a definitional extension of the working class part. Upper class terms do not compute themselves, but rather represent message passing networks that are implemented by working class terms.

In this paper, we choose for the working class part of INTML a simple first order functional language with finite types. We have chosen such a simple language because the main novelty of INTML is the upper class calculus. We stress that the working class calculus may be replaced by a more expressive language, like PCF for example.

The upper class provides constructs for space-bounded programming. It is a linear type system inspired by Dual Light Affine Logic (DLAL) [4]. INTML treats the indexed tensor  $\otimes_A$  much like exponential modality  $!$  is treated in DLAL. Just as the exponential modality may only appear in negative positions in DLAL, i.e. one can have  $!X \multimap Y$  but not  $X \multimap !Y$ , INTML only accounts for types of the form  $\otimes_A X \multimap Y$ . In the syntax we write  $A \cdot X \multimap Y$  for them. The restriction to negative occurrences of  $\otimes_A$  simplifies the type system without being too limiting in applications. In fact, we do not know of an example where  $\otimes_A$  in a positive position would be useful.

### 3.1 Type System

Working class types are first-order types with type variables. They may appear in the upper class types, which represent structure obtained from the Int construction.

$$\begin{array}{ll} \text{Working class} & A, B ::= \alpha \mid 1 \mid A \times B \mid A + B \\ \text{Upper class} & X, Y ::= [A] \mid X \otimes Y \mid A \cdot X \multimap Y \end{array}$$

Instead of  $1 \cdot X \multimap Y$  we write just  $X \multimap Y$ . For working class types we define coherent type isomorphism to be the least congruence generated by  $1 \times A \cong A \times 1 \cong A$  and  $A \times (B + C) \cong (A \times B) + (A \times C)$  and  $(B + C) \times A \cong (B \times A) + (C \times A)$ .

The terms of INTML are formed by the grammars below. We write  $c, d$  for working class variables and use  $f, g, h$  to range over working class terms. Upper class variables are ranged over by  $x, y$  and upper class terms by  $s, t$ . The terms  $\text{loop}(c.f)(g)$  and  $\text{hack}(c.f)$  bind the variable  $c$  in  $f$ .

$$\begin{array}{ll} \text{Working class} & f, g, h ::= c \mid \text{min}_A \mid \text{succ}_A(f) \mid \text{eq}_A(f, g) \\ & \mid \text{inl}(f) \mid \text{inr}(f) \mid \text{case } f \text{ of } \text{inl}(c) \Rightarrow g \mid \text{inr}(d) \Rightarrow h \\ & \mid * \mid \langle f, g \rangle \mid \text{fst}(f) \mid \text{snd}(f) \mid \text{loop}(c.f)(g) \mid \text{unbox}(t) \\ \text{Upper class} & s, t ::= x \mid \langle t, t \rangle \mid \text{let } s \text{ be } \langle x, y \rangle \text{ in } t \mid \lambda x. t \mid s t \\ & \mid [f] \mid \text{let } s \text{ be } [c] \text{ in } t \mid \text{case } f \text{ of } \text{inl}(c) \Rightarrow s \mid \text{inr}(d) \Rightarrow t \\ & \mid \text{copy } t \text{ as } x, y \text{ in } t \mid \text{hack}(c.f) \end{array}$$

The working class terms include the standard terms for  $1$ ,  $\times$  and  $+$ . In addition there are constants  $\text{min}_A$ ,  $\text{succ}_A$  and  $\text{eq}_A$  for any type  $A$ . These constants provide a total ordering and decidable equality on any type. For example, the values of the type  $(1+1) \times (1+1)$  can be ordered as  $\langle \text{inl}(*), \text{inl}(*) \rangle$ ,  $\langle \text{inr}(*), \text{inl}(*) \rangle$ ,  $\langle \text{inl}(*), \text{inr}(*) \rangle$ ,  $\langle \text{inr}(*), \text{inr}(*) \rangle$  and with  $\text{min}$  and  $\text{succ}$  we can access such an ordering generically for any type without having to define it by hand. Finally, there is a term  $\text{loop}(c.f)$  for iteration. It is a simple syntax for a trace operator with respect to  $+$ . The intended operational semantics of  $\text{loop}$  is  $\text{loop}(c.f)(\text{inl}(v)) \longrightarrow \text{loop}(c.f)(f[\text{inl}(v)/c])$  and  $\text{loop}(c.f)(\text{inr}(v)) \longrightarrow v$ , where in both cases  $v$  has already been reduced to a value (see Section 4.1).

The typing rules for working class terms appear in Fig. 1. They derive working class sequents of the form  $\Sigma \vdash f : A$  asserting that  $f$  has type  $A$  in context  $\Sigma$ . The context  $\Sigma$  assigns working class types to a finite number of working class variables. As usual, the comma in contexts corresponds to  $\times$ , even though  $\text{loop}$  corresponds to a trace with respect to  $+$ . The typing rules should be unsurprising, except perhaps that for  $\text{unbox}$ . One may think of a term  $t : [A]$  as a thunk that can be evaluated with  $\text{unbox}$ .

$$\begin{array}{c}
 \frac{}{\Sigma, c:A \vdash c: A} \quad \frac{}{\Sigma \vdash \text{min}_A: A} \quad \frac{\Sigma \vdash f: A}{\Sigma \vdash \text{succ}_A(f): A} \quad \frac{\Sigma \vdash f: A \quad \Sigma \vdash g: A}{\Sigma \vdash \text{eq}_A(f, g): 1 + 1} \\
 \\
 \frac{}{\Sigma \vdash *: 1} \quad \frac{\Sigma \vdash f: A \quad \Sigma \vdash g: B}{\Sigma \vdash \langle f, g \rangle: A \times B} \quad \frac{\Sigma \vdash f: A \times B}{\Sigma \vdash \text{fst}(f): A} \quad \frac{\Sigma \vdash f: A \times B}{\Sigma \vdash \text{snd}(g): B} \\
 \\
 \frac{\Sigma \vdash f: A}{\Sigma \vdash \text{inl}(f): A + B} \quad \frac{\Sigma \vdash f: A + B \quad \Sigma, c:A \vdash g: C \quad \Sigma, d:B \vdash h: C}{\Sigma \vdash \text{case } f \text{ of inl}(c) \Rightarrow g \mid \text{inr}(d) \Rightarrow h: C} \\
 \\
 \frac{\Sigma \vdash f: B}{\Sigma \vdash \text{inr}(f): A + B} \quad \frac{\Sigma, c:A + B \vdash f: C + B \quad \Sigma \vdash g: A + B}{\Sigma \vdash \text{loop}(c.f)(g): C} \quad \frac{\Sigma \mid \vdash t: [A]}{\Sigma \vdash \text{unbox}(t): A}
 \end{array}$$

**Fig. 1.** Working Class Typing Rules

Upper class terms denote message passing networks that will be implemented by working class terms. An upper class typing sequent has the form  $\Sigma \mid \Gamma \vdash t: X$ , where  $\Sigma$  is a working class context. The context  $\Gamma$  is a finite list of declarations of the form  $x_1: A_1 \cdot X_1, \dots, x_k: A_k \cdot X_k$ . As usual, we assume that no variable is declared more than once in  $\Gamma$ . The term  $t$  denotes a network with a single (bidirectional) output edge of type  $X$  and (bidirectional) input edges of types  $\bigotimes_{A_1} X_1, \dots, \bigotimes_{A_k} X_k$ . Informally, one may think of a declaration  $x: A \cdot X$  as a declaration of  $A$ -many copies of a value in  $X$ , i.e. one copy for each value  $v:A$ . Having multiple copies is useful because our message passing networks are stateless. When we send a query to  $X$  and later receive an answer, then we may not know what to do with that answer, since we have forgotten all that was computed earlier. However, we can use  $A \cdot X$  instead of  $X$  in order to remember a value of type  $A$ . If we want to remember a value  $v:A$  then we simply query the  $v$ -th copy of  $X$ .

For any upper class context  $\Gamma$  and any working class type  $A$ , we define an upper class context  $A \cdot \Gamma$  by  $A \cdot \langle \rangle = \langle \rangle$  and  $A \cdot (\Delta, x: B \cdot X) = (A \cdot \Delta), x: (A \times B) \cdot X$ .

The upper class typing rules appear in Fig. 2. While upper class terms denote message passing networks, at a first reading they may be understood without knowing precisely the networks they denote. In particular, the reader may wish to look at the reduction rules in Fig. 5, which are soundly implemented by the translation of terms to message passing networks. In Sec. 4 we describe which networks the terms denote.

The upper class rules represent a choice of the structure that the Int construction adds to the working class calculus. It does not capture this rich structure completely, however. Therefore we add the ‘hacking’ rule below that allows one to implement message passing nodes directly, much like one can use inline assembler in C.

$$(\text{HACK}) \frac{\Sigma, c: X^- \vdash f: X^+}{\Sigma \mid \Gamma \vdash \text{hack}(c.f): X}$$

In this rule,  $X^-$  and  $X^+$  denote the negative and positive parts of  $X$  defined by:

$$\begin{array}{l}
 [A]^- = 1 \qquad [A]^+ = A \\
 (X \otimes Y)^- = X^- + Y^- \qquad (X \otimes Y)^+ = X^+ + Y^+ \\
 (A \cdot X \multimap Y)^- = A \times X^+ + Y^- \qquad (A \cdot X \multimap Y)^+ = A \times X^- + Y^+
 \end{array}$$

$$\begin{array}{c}
\text{(WEAK)} \frac{\Sigma \mid \Gamma \vdash s : Y}{\Sigma \mid \Gamma, x : A \cdot X \vdash s : Y} \quad \text{(EXCH)} \frac{\Sigma \mid \Gamma, x : A \cdot X, y : B \cdot Y, \Delta \vdash s : Z}{\Sigma \mid \Gamma, y : B \cdot Y, x : A \cdot X, \Delta \vdash s : Z} \\
\text{(LWEAK)} \frac{\Sigma \mid \Gamma, x : A \cdot X \vdash s : Y}{\Sigma \mid \Gamma, x : (B \times A) \cdot X \vdash s : Y} \quad \text{(CONGR)} \frac{\Sigma \mid \Gamma, x : A \cdot X \vdash s : Y}{\Sigma \mid \Gamma, x : B \cdot X \vdash s : Y} A \cong B \\
\text{(VAR)} \frac{}{\Sigma \mid \Gamma, x : A \cdot X \vdash x : X} \quad \text{(\otimes I)} \frac{\Sigma \mid \Gamma \vdash s : X \quad \Sigma \mid \Delta \vdash t : Y}{\Sigma \mid \Gamma, \Delta \vdash \langle s, t \rangle : X \otimes Y} \\
\text{(\otimes E)} \frac{\Sigma \mid \Gamma \vdash s : X \otimes Y \quad \Sigma \mid \Delta, x : A \cdot X, y : A \cdot Y \vdash t : Z}{\Sigma \mid \Delta, A \cdot \Gamma \vdash \text{let } s \text{ be } \langle x, y \rangle \text{ in } t : Z} \\
\text{(\multimap I)} \frac{\Sigma \mid \Gamma, x : A \cdot X \vdash s : Y}{\Sigma \mid \Gamma \vdash \lambda x. s : A \cdot X \multimap Y} \quad \text{(\multimap E)} \frac{\Sigma \mid \Gamma \vdash s : A \cdot X \multimap Y \quad \Sigma \mid \Delta \vdash t : X}{\Sigma \mid \Gamma, A \cdot \Delta \vdash s t : Y} \\
\text{(CONTR)} \frac{\Sigma \mid \Gamma \vdash s : X \quad \Sigma \mid \Delta, x : A \cdot X, y : B \cdot X \vdash t : Y}{\Sigma \mid \Delta, (A + B) \cdot \Gamma \vdash \text{copy } s \text{ as } x, y \text{ in } t : Y} \\
\text{(CASE)} \frac{\Sigma \vdash f : A + B \quad \Sigma, c : A \mid \Gamma \vdash s : X \quad \Sigma, d : B \mid \Gamma \vdash t : X}{\Sigma \mid \Gamma \vdash \text{case } f \text{ of } \text{inl}(c) \Rightarrow s \mid \text{inr}(d) \Rightarrow t : X} \\
\text{([\ ] I)} \frac{\Sigma \vdash f : A}{\Sigma \mid \Gamma \vdash [f] : [A]} \quad \text{([\ ] E)} \frac{\Sigma \mid \Gamma \vdash s : [A] \quad \Sigma, c : A \mid \Delta \vdash t : [B]}{\Sigma \mid \Gamma, A \cdot \Delta \vdash \text{let } s \text{ be } [c] \text{ in } t : [B]}
\end{array}$$

Fig. 2. Upper Class Typing Rules

**Examples.** We give an example derivation to illustrate that while the upper class type system is linear, working-class variables can be copied arbitrarily:

$$\begin{array}{c}
\vdots \\
\frac{}{c : \alpha \vdash c : \alpha} \\
\frac{c : \alpha \mid f : 1 \cdot ([\alpha] \multimap [\alpha] \multimap [\beta]) \vdash f [c] : [\alpha] \multimap [\beta]}{c : \alpha \mid \vdash [c] : [\alpha]} \\
\frac{\mid x : 1 \cdot [\alpha] \vdash x : [\alpha] \quad c : \alpha \mid f : 1 \cdot ([\alpha] \multimap [\alpha] \multimap [\beta]) \vdash f [c] [c] : [\beta]}{f : (\alpha \times 1) \cdot ([\alpha] \multimap [\alpha] \multimap [\beta]), x : 1 \cdot [\alpha] \vdash \text{let } x \text{ be } [c] \text{ in } f [c] [c] : [\beta]} \text{([\ ] E)} \\
\frac{f : \alpha \cdot ([\alpha] \multimap [\alpha] \multimap [\beta]) \vdash \lambda x. \text{let } x \text{ be } [c] \text{ in } f [c] [c] : [\alpha] \multimap [\beta]}{f : \alpha \cdot ([\alpha] \multimap [\alpha] \multimap [\beta]) \vdash \lambda f. \lambda x. \text{let } x \text{ be } [c] \text{ in } f [c] [c] : \alpha \cdot ([\alpha] \multimap [\alpha] \multimap [\beta]) \multimap [\alpha] \multimap [\beta]} \text{(\multimap I)}
\end{array}$$

Even though upper class terms can be understood as if they were implemented by the reduction rules in Fig. 5, it is important to understand that they will be (in Sec. 4) compiled down to (large) working class terms that implement certain message passing networks. The network for the upper class term in the conclusion of the above derivation, for example, represents a message passing network with a single (bidirectional) output wire. It behaves as follows: if it receives a request for the value of the result in  $[\beta]$ , it first requests the value of  $x$ . Upon receipt of this value, the network will then ask the function  $f$  for its return value. Since  $f$  has a type of the form  $\alpha \cdot X$ , the network has access to  $\alpha$ -many copies of  $f$ . It chooses the copy indexed by the value of  $x$ , so that, even though the network is stateless, the value of  $x$  will be available once an answer from  $f$  arrives. If  $f$  answers with a value in  $[\beta]$ , then this answer is forwarded as the



final answer of the whole network. If  $f$  answers with a request for one of its arguments, then the network gives the value of  $x$  as reply to  $f$ .

That being able to copy working class variables does not make copy superfluous can be seen in the following terms for conversion between  $[\alpha \times \beta]$  and  $[\alpha] \otimes [\beta]$ .

$$\begin{aligned} \lambda y. \text{ copy } y \text{ as } y_1, y_2 \text{ in} \\ \langle \text{let } y_1 \text{ be } [c] \text{ in } [\text{fst}(c)], \text{ let } y_2 \text{ be } [c] \text{ in } [\text{snd}(c)] \rangle & : (\gamma + \delta) \cdot [\alpha \times \beta] \multimap [\alpha] \otimes [\beta] \\ \lambda z. \text{ let } z \text{ be } \langle x, y \rangle \text{ in let } x \text{ be } [c] \text{ in let } y \text{ be } [d] \text{ in } [\langle c, d \rangle] & : \alpha \cdot ([\alpha] \otimes [\beta]) \multimap [\alpha \times \beta] \end{aligned}$$

**Useful Combinators.** The upper class calculus is a simple linear lambda calculus for constructing message passing networks. It appears to be missing many constructs, such as loops, that are required to make it an expressive programming language. Such constructs can be defined as higher-order combinators using hack.

The most important example of such a combinator is a loop iterator that informally satisfies  $\text{loop } f v = w$  if  $f v = \text{inr}(w)$  and  $\text{loop } f v = \text{loop } f w$  if  $f v = \text{inl}(w)$ .

$$\text{loop}: \alpha \cdot (\gamma \cdot [\alpha] \multimap [\alpha + \beta]) \multimap [\alpha] \multimap [\beta]$$

Before we define  $\text{loop}$ , we give a typical example for its use. We define  $\text{fold}_\alpha$ , such that  $\text{fold}_\alpha f y$  computes  $f x_n (\dots (f x_1 (f x_0 y)))$ , where  $x_0 = \text{min}_\alpha$  and  $x_{i+1} = \text{succ}_\alpha(x_i)$  and  $x_n$  is the maximum element of  $\alpha$ , i.e. the element with  $x_n = \text{succ}_\alpha(x_n)$ .

$$\begin{aligned} \text{fold}_\alpha & : (\alpha \times \beta \times \alpha \times \beta) \cdot ([\alpha] \multimap [\beta] \multimap [\beta]) \multimap [\beta] \multimap [\beta] \\ \text{fold}_\alpha = \lambda f. \lambda y. \text{ loop } (\lambda w. \text{ let } w \text{ be } [e] \text{ in let } f \text{ [fst}(e)] \text{ [snd}(e)] \text{ be } [z] \text{ in} \\ & \quad \text{case } eq_\alpha(\text{fst}(e), \text{succ}_\alpha(\text{fst}(e))) \text{ of } \text{inl}(\text{true}) \Rightarrow [\text{inr}(z)] \\ & \quad \quad | \text{inr}(\text{false}) \Rightarrow [\text{inl}(\langle \text{succ}_\alpha(\text{fst}(e)), z \rangle)]) \\ & \quad (\text{let } y \text{ be } [z] \text{ in } [\langle \text{min}_\alpha, z \rangle]) \end{aligned}$$

The definition of  $\text{loop}$  uses hack and therefore makes explicit reference to the translation of upper class terms to message passing networks, which we describe in detail in Sec 4. The  $\text{loop}$ -combinator is defined by  $\text{loop} = \text{hack}(c.l)$ , where  $l$  is a working class term of type  $c:\alpha \times (\gamma \times 1 + (\alpha + \beta)) + (\alpha + 1) \vdash l: \alpha \times (\gamma \times \alpha + 1) + (1 + \beta)$  that implements the following mappings using a nested case expression: (i)  $\text{inr}(\text{inr}(*)) \mapsto \text{inl}(\text{inr}(*))$ ; (ii)  $\text{inr}(\text{inl}(a)) \mapsto \text{inl}(a, \text{inr}(*))$ ; (iii)  $\text{inl}(a, \text{inr}(\text{inr}(b))) \mapsto \text{inr}(\text{inr}(b))$ ; (iv)  $\text{inl}(a, \text{inr}(\text{inl}(a'))) \mapsto \text{inl}(a', \text{inr}(*))$ ; and (v)  $\text{inl}(a, \text{inl}(g, *)) \mapsto \text{inl}(a, \text{inl}(g, a))$ . These assignments can be interpreted as follows: Mapping (i) says that when we get a request for the final value, we start by asking the base case for its value. When an answer from the base cases arrives, we put it in a memory cell (corresponding to  $\alpha \cdot -$  in the type) and ask the step function for its result (ii). Whenever the step function asks for its argument, we supply the value from the memory cell (v). If the step function answers  $\text{inr}(b)$ , then we are done and give  $b$  as output (iii). If the step function answers  $\text{inl}(a')$ , then we overwrite the memory cell content with  $a'$  and restart the step function by asking for its result (iv).

A second useful combinator is a simple version of  $\text{callcc}$ . It can be given the following type for any upper class type  $X$ .

$$\text{callcc}: (\gamma \cdot ([\alpha] \multimap X) \multimap [\alpha]) \multimap [\alpha]$$



safe recursion on notation by computational amnesia [14]. We believe that using `loop` and `callcc` a similar program can be implemented in INTML as a combinator `saferec` taking as arguments a base case  $g: [\alpha] \rightarrow [3]$ , two step functions  $h_0, h_1: ([\alpha] \rightarrow [3]) \rightarrow ([\alpha] \rightarrow [3])$  and a word  $w: [\alpha] \rightarrow [3]$  to recurse on and giving a word as output.

In this way the higher-order features of INTML can be used to abstract away details of a message-passing implementation of functions on words and to write LOGSPACE functions just like in  $BC_\varepsilon^-$  [14]. Moreover, INTML gives access to the implementation details, should the abstraction not be expressive enough. For instance, the proof of LOGSPACE completeness in Thm. 3 below goes by a straightforward encoding of a OTM. It is much simpler than the corresponding encoding in  $BC_\varepsilon^-$  [14] because INTML allows us to manipulate working class values directly.

The higher order approach also works for data types other than strings. For example, graphs can be represented by a type of the form  $([\alpha] \rightarrow [2]) \otimes ([\alpha \times \alpha] \rightarrow [2])$ , where the first component is a predicate that indicates which elements of  $\alpha$  count as graph nodes and the second component is the edge relation.

## 4 Evaluation

In this section we present the evaluation mechanism for INTML. Evaluation of upper class terms is closely related to evaluation in [15], but also to Mackie's Interaction Abstract Machine [11] and to read-back from optimal reduction [12,3]. Indeed, in his 1995 paper [11] Mackie speculates that this form of evaluation could have applications where space usage is important. With INTML we present a calculus that makes space usage analysis possible. INTML differs from the work in loc. cit. in that one can mix working class and upper class terms. Previously only the upper class part was considered.

### 4.1 Reduction of Working-Class Terms

The evaluation of INTML programs is done by reduction of working class terms. Before evaluation of an INTML program, all upper class terms are compiled into working class terms. Since, in particular, any occurrence of `unbox` will be removed, it suffices to define reduction only for `unbox`-free working class terms.

Working class values are defined by:

$$v, w := c \mid * \mid \langle v, w \rangle \mid \text{inl}(v) \mid \text{inr}(v)$$

The reduction of working class terms is explained by a small step reduction relation  $\longrightarrow$  between closed `unbox`-free terms. Closedness here means the absence of both term variables and type variables (which could appear in the type annotations of constants like  $\text{min}_A$ ). The relation  $\longrightarrow$  formalises standard *eager* reduction, see e.g. [17]. We omit standard reduction rules and just explain here how `loop` and the constants  $\text{min}_A$  and  $\text{succ}_A$  are treated. Loops are unfolded by the rules  $\text{loop}(c.f)(\text{inr}(v)) \longrightarrow v$  and  $\text{loop}(c.f)(\text{inl}(v)) \longrightarrow \text{loop}(c.f)(f[\text{inl}(v)/c])$ , in both of which  $v$  must be a value. Constants are unfolded on demand, guided by their type annotation. The minimum elements of all closed types are defined by:

$$\text{min}_1 \longrightarrow * \quad \text{min}_{A+B} \longrightarrow \text{inl}(\text{min}_A) \quad \text{min}_{A \times B} \longrightarrow \langle \text{min}_A, \text{min}_B \rangle$$

In the implementation of  $\text{succ}_A$ , we must be a little careful that reducts do not become too large. We implement  $\text{succ}$  using a new constant  $\text{succmin}$  that informally denotes a function  $\text{succmin}_A : A \rightarrow A + A$  with the following meaning:  $\text{succmin}_A(x) = \text{inl}(y)$  means that  $y$  is the successor of  $x$ ;  $\text{succmin}_A(x) = \text{inr}(y)$  means that  $x$  has no successor and  $y$  is the minimum element of  $A$ . We use the following rules for  $\text{succmin}$ .

$$\begin{aligned}
& \text{succmin}_1(*) \longrightarrow \text{inr}(*) \\
& \text{succmin}_{A+B}(\text{inl}(v)) \longrightarrow \text{case } \text{succmin}_A(v) \text{ of } \text{inl}(x) \Rightarrow \text{inl}(\text{inl}(x)) \\
& \quad \quad \quad | \text{inr}(y) \Rightarrow \text{inl}(\text{inr}(\text{min}_B)) \\
& \text{succmin}_{A+B}(\text{inr}(v)) \longrightarrow \text{case } \text{succmin}_B(v) \text{ of } \text{inl}(x) \Rightarrow \text{inl}(\text{inr}(x)) \\
& \quad \quad \quad | \text{inr}(y) \Rightarrow \text{inr}(\text{inl}(\text{min}_A)) \\
& \text{succmin}_{A \times B}(\langle v, w \rangle) \longrightarrow \text{case } \text{succmin}_A(v) \text{ of } \text{inl}(x) \Rightarrow \text{inl}(\langle x, w \rangle) \\
& \quad \quad \quad | \text{inr}(x) \Rightarrow \text{case } \text{succmin}_B(w) \text{ of } \text{inl}(y) \Rightarrow \text{inl}(\langle x, y \rangle) \\
& \quad \quad \quad | \text{inr}(y) \Rightarrow \text{inr}(\langle x, y \rangle)
\end{aligned}$$

We use these rules instead of the evident rules for  $\text{succ}$  because they are linear in  $v$ ,  $w$ ,  $A$  and  $B$ , which is important for the proof of Prop. 1 below.

Because of the simplicity of the working class calculus, it is possible to give useful upper bounds on how large a term can become during the course of eager reduction directly by induction on the term structure. Essentially, we can bound the size of values in terms of their types and use this to derive a bound on the potential size of a term under reduction by looking at its variables and their types. We state this result in the following proposition, in which  $|g|$  and  $|C|$  denote the size of the abstract syntax trees of  $g$  and  $C$  respectively.

**Proposition 1.** *If  $c:A \vdash f : B$  is derivable then there are constants  $n$  and  $m$  such that  $(f[C/\alpha])[v/c] \longrightarrow^* g$  implies  $|g| \leq n + m \cdot |C|$  for every closed type  $C$  and every closed value  $v$  of type  $A[C/\alpha]$ .*

We will use for  $C$  types of the form  $2 \times \dots \times 2$  ( $k$  times), where 2 denotes  $1 + 1$ . This type represents the numbers from 0 to  $2^k - 1$  in binary and we have  $|C| \in O(k)$ . A unary encoding is also possible with  $C = 1 + \dots + 1$  ( $2^k$  times), but we have  $|C| \geq 2^k$  and values can indeed become as large as this.

## 4.2 Reducing Upper Class to Working Class

Now we explain how closed upper class terms are compiled down to working class terms, so that they can be reduced with the relation  $\longrightarrow$  from the previous section. The compilation works by interpreting upper class terms as message passing circuits as in Sec. 2 and then implementing these circuits by working class terms.

We start by defining the message passing circuits we use in the compilation. These circuits may be understood as a particular instance of string diagrams for monoidal categories [13]. They are also related to proof nets, see [13] for a discussion. Circuits are directed graphs that represent networks in which messages are passed along edges. A node labelling allows us to use nodes with different message passing behaviour. Edges

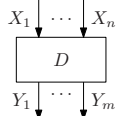
are labelled with types that tell which kind of messages can be passed along them. Edge labels are formed by the grammar below, in which  $A$  ranges over working class types.

$$X, Y ::= [A] \mid [A]^* \mid X \otimes Y \mid \bigotimes_A X$$

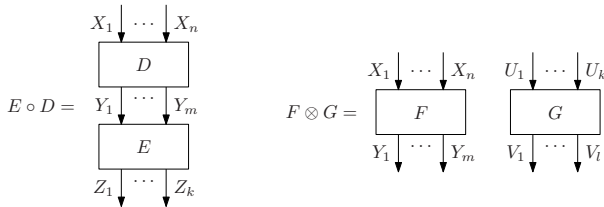
We define an operation  $(-)^*$  on edge labels as follows. It maps  $[A]$  to  $[A]^*$  and  $[A]^*$  to  $[A]$  and is defined on compound expressions by  $(X \otimes Y)^* = X^* \otimes Y^*$  and  $(\bigotimes_A X)^* = \bigotimes_A X^*$ . Note in particular that we have  $X^{**} = X$  for any type label  $X$ .

Circuits are labelled directed graphs that are in addition equipped with a *two-way local ordering*. A two-way local ordering for a graph  $G = (V, E)$  specifies for each node  $v \in V$  a total ordering on both the set  $(\{v\} \times V) \cap E$  of outgoing edges from  $v$  and on the set  $(V \times \{v\}) \cap E$  of incoming edges to  $v$ . The need for a local ordering arises in the treatment of nodes such as  $\otimes E$  below. This node has an incoming edge labelled with  $X \otimes Y$  and two outgoing edges labelled with  $X$  and  $Y$  and we want to distinguish the two outgoing edges even if  $X$  and  $Y$  are the same.

Furthermore, circuits have a number of input and output ports. We capture input and output edges by means of two distinguished nodes: a source and a sink. Edges from the distinguished source are input edges and edges to the sink are output edges. We write  $D: (X_1, \dots, X_n) \rightarrow (Y_1, \dots, Y_m)$  for such a graph  $D$  with input edges of type  $X_1, \dots, X_n$  and output edges of type  $Y_1, \dots, Y_m$  (note that they are ordered because of the local ordering for source and sink). We usually draw  $D$  as shown on the right.



Given two graphs  $D: (\mathbf{X}) \rightarrow (\mathbf{Y})$  and  $E: (\mathbf{Y}) \rightarrow (\mathbf{Z})$ , their sequential composition  $E \circ D: (\mathbf{X}) \rightarrow (\mathbf{Z})$  is defined as depicted below: the  $i$ -th incoming edge to the sink of  $D$  and the  $i$ -th outgoing edge from the source of  $E$  are joined to a single edge. The sink of  $D$  and the source of  $E$  are removed. Furthermore, given  $F: (\mathbf{X}) \rightarrow (\mathbf{Y})$  and  $G: (\mathbf{U}) \rightarrow (\mathbf{V})$ , we write  $F \otimes G$  for the graph of type  $(\mathbf{X}, \mathbf{U}) \rightarrow (\mathbf{Y}, \mathbf{V})$  obtained by putting  $F$  and  $G$  in parallel.

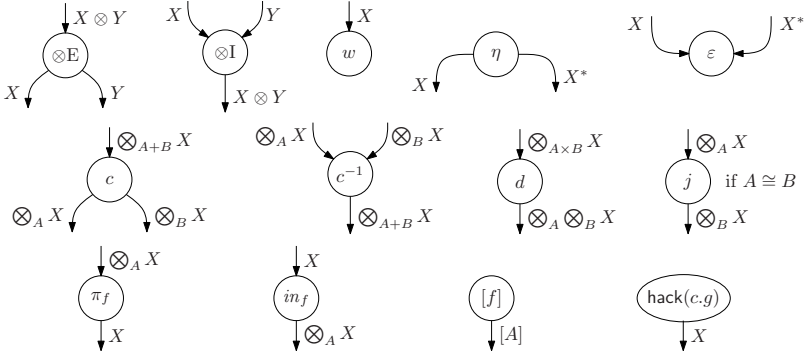


Notice that in a composition the local ordering on the edges labelled with  $\mathbf{Y}$  disappears. For example, if we let  $swap_{X,Y}: (X, Y) \rightarrow (Y, X)$  be the graph on the right, then  $swap_{Y,X} \circ swap_{X,Y}$  is  $id_X \otimes id_Y: (X, Y) \rightarrow (X, Y)$ , where  $id_X: (X) \rightarrow (X)$  is a single edge labelled  $X$  from input to output.



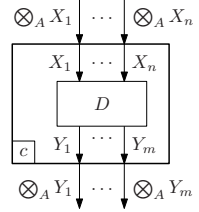
**Definition 1 (Circuit on  $\Sigma$ ).** For any working class context  $\Sigma$ , we define the set of *circuits on  $\Sigma$*  to be the smallest set of two-way locally ordered graphs with input and output edges that satisfies the following conditions.

- For all  $X$  and  $Y$  both  $id_X$  and  $swap_{X,Y}$  are circuits on  $\Sigma$ .
- Each of the following single-node graphs is a circuit on  $\Sigma$ , where  $X$  and  $Y$  may be arbitrary type labels,  $A$  and  $B$  may be arbitrary working class types and  $f$  and  $g$  may be arbitrary terms with  $\Sigma \vdash f : A$  and  $\Sigma, c : X^- \vdash g : X^+$ .



- If  $D : (X) \rightarrow (Y)$  and  $E : (Y) \rightarrow (Z)$  are circuits on  $\Sigma$  then so is  $E \circ D$ .
- If  $D : (X) \rightarrow (Y)$  and  $E : (U) \rightarrow (V)$  are circuits on  $\Sigma$  then so is  $D \otimes E$ .
- If  $D : (X_1, \dots, X_n) \rightarrow (Y_1, \dots, Y_m)$  is a circuit on  $\Sigma$ ,  $c : A$  then the graph

$\otimes_{c:A} D : (\otimes_A X_1, \dots, \otimes_A X_n) \rightarrow (\otimes_A Y_1, \dots, \otimes_A Y_m)$  constructed as follows is a circuit on  $\Sigma$ : each incoming edge of  $D$  is prepended with a node  $\downarrow_{X_i}^c : (\otimes_A X_i) \rightarrow (X_i)$  and to each outgoing edge of  $D$  a node  $\uparrow_{Y_i}^c : (Y_i) \rightarrow (\otimes_A Y_i)$  is appended. For better readability, we do not draw these nodes explicitly and draw a box around  $D$  instead, as depicted on the right.



Given a circuit  $D$  on  $(\Sigma, c:A)$  and a term  $\Sigma \vdash f : A$ , we can form a circuit  $D[f/c]$  on  $\Sigma$  by replacing each node  $\pi_g$  with  $\pi_{g[f/c]}$ , each node  $in_g$  with  $in_{g[f/c]}$ , each node  $[g]$  with  $[g[f/c]]$ , and each node  $hack(d.g)$  with  $hack(d.g[f/c])$ .

To each edge in a circuit  $D$  on  $\Sigma$  we assign a *level*, which is a working class context, by the following requirements. Input and output edges have level  $\Sigma$ . Any two edges incident to the same node have the same level, except if the node is  $\downarrow_X^c : (\otimes_A X) \rightarrow (X)$  or  $\uparrow_X^c : (X) \rightarrow (\otimes_A X)$ . If the incoming edge of  $\downarrow_X^c$  (resp. the outgoing edge of  $\uparrow_X^c$ ) has level  $\Sigma'$  then the outgoing edge (resp. incoming edge) has level  $\Sigma'$ ,  $c:A$ .

*Message Passing.* Write  $\mathcal{V}_A$  for the set of all closed working class values of type  $A$ . Write  $\mathcal{E}_\Sigma$  for the set of  $\Sigma$ -environments consisting of all functions that map the variables in  $\Sigma$  to closed values of their declared types. The set  $M_{\Sigma,X}$  of messages that can be passed along an edge labelled with  $X$  at level  $\Sigma$  is then defined by  $M_{\Sigma,X} = \mathcal{E}_\Sigma \times (\mathcal{V}_{X^+} \times \{+\} \cup \mathcal{V}_{X^-} \times \{-\})$ . A message  $m \in M_{\Sigma,X}$  is either a *question* or an *answer* depending on whether its third component is ‘-’ or ‘+’. Answers travel in the direction of the edge while questions travel in the opposite direction. Messages are essentially the same as the contexts in context semantics [12].

To define message passing for a circuit  $D : (X) \rightarrow (Y)$  on  $\Sigma$ , let the set  $M_D$  of messages on  $D$  consist of all pairs  $(e, m)$  of an edge  $e$  in  $D$  and a message

$m \in M_{\Sigma(e), X(e)}$ , where  $\Sigma(e)$  is the level of  $e$  and  $X(e)$  is its label. Now we define how each node in  $D$  locally reacts to arriving messages. To this end we define for each node  $v$  a partial function  $\varphi_v$  by the assignments in Fig. 3. For example  $\varphi_\pi$  implements the behaviour of the term in  $([]E)$ , as sketched in Sec. 3.1. An initial question on edge  $o$  becomes a request for the value of type  $[A]$  on edge  $i_1$ . Upon receipt of an answer  $w:A$  from there, this answer is put into the memory cell provided by  $\otimes_A$  and the value of type  $[B]$  is queried on edge  $i_2$ . An answer to this request is then passed as final answer along edge  $o$ . In Fig. 3 we have omitted the cases for  $\varphi_\varepsilon$  and  $\varphi_{c-1}$ , which are just like those for  $\varphi_\eta$  and  $\varphi_c$ . Notice that for different nodes  $v$  and  $w$  in  $D$ , the domains of  $\varphi_v$  and  $\varphi_w$  do not overlap.

The message passing behaviour of the whole circuit is then the partial function  $\varphi_D: M_D \rightarrow M_D$  defined by repeatedly applying the local functions  $\varphi_v$ :

$$\varphi_D = Tr \left( \bigcup_{v \in V(D)} \varphi_v \right) \quad Tr(f)(m) = \begin{cases} Tr(f)(f(m)) & \text{if } f(m) \text{ defined} \\ m & \text{if } f(m) \text{ undefined} \end{cases}$$

It is not hard to see that messages cannot get stuck inside the circuit, i.e.  $\varphi_D(e, m) = (e', m')$  implies that  $e'$  is an input or an output edge. In fact, we will forget about the internal structure of  $D$  and consider just the restriction of  $\varphi_D$  to messages on input or output edges of  $D$ . This restriction corresponds to a partial function  $\Phi_D$  of type  $(X_1^+ + \dots + X_n^+ + Y_1^- + \dots + Y_m^-) \times \mathcal{E}_\Sigma \rightarrow Y_1^+ + \dots + Y_m^+ + X_1^- + \dots + X_n^-$ . We call  $\Phi_D$  *the behaviour of  $D$*  and consider circuits with the same behaviour to be equal. We write  $D \sim E$  if  $D$  and  $E$  are circuits with the same interface and  $\Phi_D = \Phi_E$ .

*Upper Class to Circuits.* We now interpret upper class terms by circuits. To each derivation  $\delta$  ending with sequent  $\Sigma \mid x_1 : A_1 \cdot X_1, \dots, x_n : A_n \cdot X_n \vdash s : Y$  we assign  $\llbracket \delta \rrbracket$ , a circuit on  $\Sigma$  with one output wire of type  $Y$  and  $n$  input wires of type  $(\otimes_{A_1} X_1, \dots, \otimes_{A_n} X_n)$ , where we identify  $A \cdot X \multimap Y$  with  $(\otimes_A X)^* \otimes Y$ . That is, each variable declaration in  $\Gamma$  becomes an input wire of the translated circuit. Abusing notation slightly, we will write  $\Gamma$  also for the list of input types of this circuit.

The definition goes by induction on derivations and is given in the table in Fig. 4. In this table we denote the premises of  $\delta$  by  $\delta_s$  and  $\delta_t$  depending on the term in the premise. We write just  $\otimes_A$  for  $\otimes_{c:A}$  if  $c$  does not appear anywhere. Given an upper class context  $\Gamma$  we write  $w_\Gamma$  for  $w \otimes \dots \otimes w : (\Gamma) \rightarrow ()$ . Similarly, we write  $id_\Gamma : (\Gamma) \rightarrow (\Gamma)$ ,  $d_{A:\Gamma} : (A \cdot \Gamma) \rightarrow (\otimes_A \Gamma)$ ,  $(in_f)_\Gamma : (\Gamma) \rightarrow (\otimes_A \Gamma)$  and  $c_\Gamma : (\Gamma) \rightarrow (\Gamma, \Gamma)$  for the analogous tensorings of  $id$ ,  $d$ ,  $in_f$  and  $c$  (the definition of  $c_\Gamma$  involves evident permutations of the outputs to arrive at the indicated type).

*Implementing Message Passing.* The compilation of upper class terms to message passing circuits almost completes the translation from upper class to working class. It just remains to implement message passing in the working class calculus.

That is, for any circuit  $D$  we construct a closed working class term that implements its behaviour  $\Phi_D$ . In essence, the construction works in the same way as the definition of  $\Phi_D$  above. First note that we can represent the set of messages  $M_{\Sigma, X}$  by the working class type  $(A_1 \times \dots \times A_n) \times (X^- + X^+)$ , where  $\Sigma = x_1 : A_1, \dots, x_n : A_n$ . With this we can represent  $M_D$  in the working class calculus in the form of a big sum type. Then,

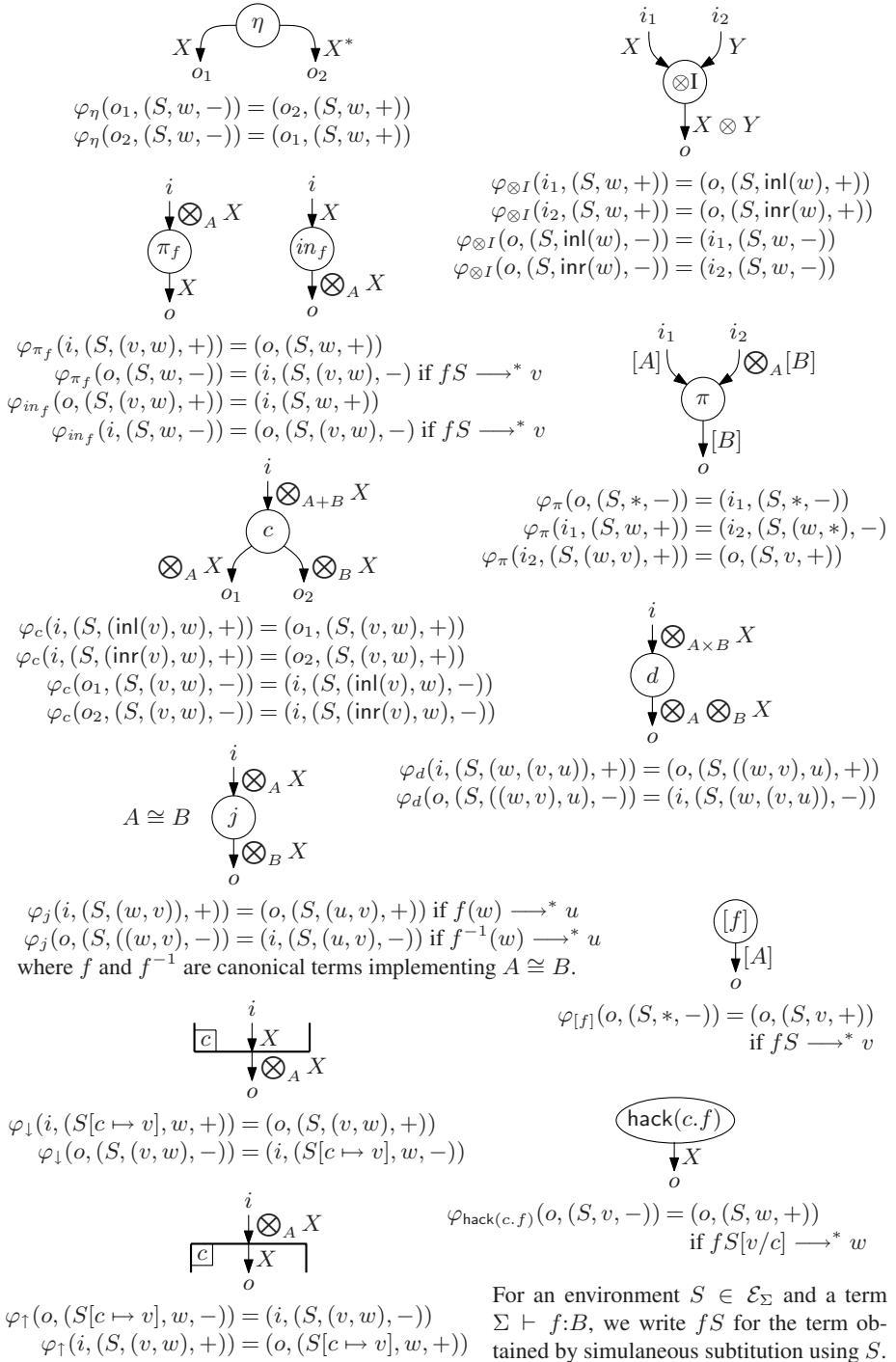


Fig. 3. Local message passing



- (WEAK)  $\llbracket \delta \rrbracket = \llbracket \delta_s \rrbracket \circ (id_\Gamma \otimes w)$   
 (EXCH)  $\llbracket \delta \rrbracket = \llbracket \delta_s \rrbracket \circ (id_\Gamma \otimes swap \otimes id_\Delta)$   
 (CONGR)  $\llbracket \delta \rrbracket = \llbracket \delta_s \rrbracket \circ (id_\Gamma \otimes j)$   
 (VAR)  $\llbracket \delta \rrbracket = w_\Gamma \otimes \pi_{min_A}$   
 ( $\otimes$ I)  $\llbracket \delta \rrbracket = \otimes I \circ (\llbracket \delta_s \rrbracket \otimes \llbracket \delta_t \rrbracket)$   
 ( $\otimes$ E)  $\llbracket \delta \rrbracket = \otimes E \circ (id_\Delta \otimes ((\otimes_{c:A} \llbracket \delta_s \rrbracket) \circ d_{A \cdot \Gamma}))$   
 ( $\rightarrow$ I)  $\llbracket \delta \rrbracket = \otimes I \circ swap \circ (\llbracket \delta_s \rrbracket \otimes id) \circ (id_\Delta \otimes \eta)$   
 ( $\rightarrow$ E)  $\llbracket \delta \rrbracket = (id_Y \otimes \varepsilon) \circ ((swap \circ \otimes E \circ \llbracket \delta_s \rrbracket) \otimes ((\otimes_{c:A} \llbracket \delta_t \rrbracket) \circ d_{A \cdot \Gamma}))$   
 (CONTR)  $\llbracket \delta \rrbracket = \otimes E \circ (id_\Delta \otimes ((\otimes_{c:A+B} \llbracket \delta_s \rrbracket) \circ d))$   
 (CASE)  $\llbracket \delta \rrbracket = \pi_f \circ c^{-1} \circ (\otimes_A \llbracket \delta_s \rrbracket \otimes \otimes_A \llbracket \delta_t \rrbracket) \circ c_\Gamma \circ (inf)_\Gamma$   
 ( $[\ ]$ I)  $\llbracket \delta \rrbracket = [f] \otimes w_\Gamma$   
 ( $[\ ]$ E)  $\llbracket \delta \rrbracket = \pi \circ (\llbracket \delta_s \rrbracket \otimes (\otimes_A \llbracket \delta_t \rrbracket \circ d_{A \cdot \Delta}))$   
 (HACK)  $\llbracket \delta \rrbracket = hack(c, f)$

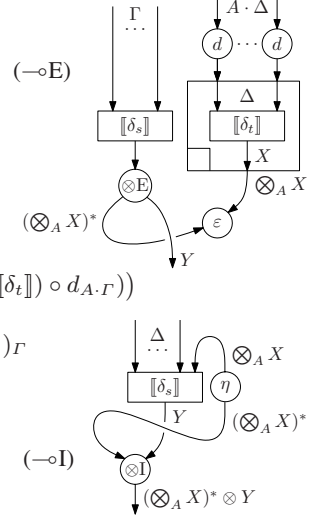


Fig. 4. Compilation to Circuits

$\bigcup_{v \in V(D)} \varphi_v : M_D \rightarrow M_D$  can be implemented easily by a big case distinction and it is easy to obtain  $\varphi_D$  from this using a single loop. From this, finally, we obtain  $\Phi_D$ .

We have thus explained how each upper class term can be translated to a working class term. It remains to say how we deal with terms of the form  $\Sigma \vdash \text{unbox}(s) : A$  where  $\Sigma \vdash s : [A]$ . Note that  $s$  is translated to a circuit  $\llbracket \delta_s \rrbracket$  whose behaviour is a function of type  $\Phi_{\llbracket \delta_s \rrbracket} : 1 \times \mathcal{E}_\Sigma \rightarrow A$ . The working class term implementing this function is just what we need to interpret the term  $\text{unbox}$ .

### 4.3 Soundness

We have defined the evaluation of INTML by translation of upper class to working class, since this allows us to obtain sublinear space bounds (cf. Sec. 5). However, the translation is somewhat complicated and it does not make it obvious just what it is that the upper class terms compute. In this section, we show that one may also consider the upper class as a functional programming language with standard reduction rules, which are then soundly implemented by the translation.

A notion of reduction can be formulated by the rules in Fig. 5. We include a rule for  $\text{loop}$  but note that we cannot give rules for  $\text{hack}$  in general. When we write  $s \longrightarrow t$  we assume  $\vdash s : X$  and  $\vdash t : X$  for some  $X$ . We close these rules under evaluation contexts.

$$E, F ::= [\cdot] \mid \langle E, t \rangle \mid \langle t, E \rangle \mid \text{let } E \text{ be } \langle x, y \rangle \text{ in } t \mid \\ E t \mid t E \mid \text{let } E \text{ be } [c] \text{ in } t \mid \text{copy } E \text{ as } x, y \text{ in } t$$

As usual,  $E[s]$  is the term obtained by substituting  $s$  for the only occurrence of  $[\cdot]$  in  $E$ .

The translation to working class terms is invariant under reduction:

**Theorem 1 (Soundness).** *If  $\delta$  derives  $\vdash s : X$  and  $s \longrightarrow t$ , then there exists a derivation  $\rho$  of  $\vdash t : X$  that satisfies  $\llbracket \delta \rrbracket \sim \llbracket \rho \rrbracket$ .*

$$\begin{aligned}
& (\lambda x.s)t \longrightarrow s[t/x] \\
& \text{let } \langle s, t \rangle \text{ be } \langle x, y \rangle \text{ in } u \longrightarrow u[s/x][t/y] \\
& \text{case } \text{inl}(v) \text{ of } \text{inl}(c) \Rightarrow s \mid \text{inr}(d) \Rightarrow t \longrightarrow s[v/c] \quad \text{if } v \text{ is a value} \\
& \text{case } \text{inr}(v) \text{ of } \text{inl}(c) \Rightarrow s \mid \text{inr}(d) \Rightarrow t \longrightarrow t[v/d] \quad \text{if } v \text{ is a value} \\
& \quad \text{let } [v] \text{ be } [c] \text{ in } t \longrightarrow t[v/c] \quad \text{if } v \text{ is a value} \\
& \quad \text{copy } s \text{ as } x, y \text{ in } t \longrightarrow t[s/x][s/y] \\
& \quad [f] \longrightarrow [g] \quad \text{if } f \longrightarrow g \\
& \text{case } f \text{ of } \text{inl}(c) \Rightarrow s \mid \text{inr}(d) \Rightarrow t \longrightarrow \text{case } g \text{ of } \text{inl}(c) \Rightarrow s \mid \text{inr}(d) \Rightarrow t \quad \text{if } f \longrightarrow g \\
& \quad \text{loop } s \text{ } t \longrightarrow \text{let } t \text{ be } [c] \text{ in let } s [c] \text{ be } [d] \text{ in} \\
& \quad \quad \text{case } d \text{ of } \text{inl}(d) \Rightarrow \text{loop } s [d] \mid \text{inr}(d) \Rightarrow [d]
\end{aligned}$$

Fig. 5. Upper Class Reduction

For the proof we need substitution lemmas, which are proved by induction on  $\rho$ .

**Lemma 1.** *If  $v \in \mathcal{V}_A$  and  $\rho$  derives  $\Sigma, c:A \mid \Gamma \vdash t: X$ , then there is a derivation  $\rho[v/c]$  of  $\Sigma \mid \Gamma \vdash t[v/c]: X$  such that  $\llbracket \rho[v/c] \rrbracket \sim \llbracket \rho \rrbracket[v/c]$  holds.*

**Lemma 2.** *If  $\delta$  derives  $\Sigma \mid \vdash s: X$  and  $\rho$  derives  $\Sigma \mid \Gamma, x:A \cdot X, \Delta \vdash t: Y$ , then there exists a derivation  $\rho[\delta/x]$  of  $\Sigma \mid \Gamma, \Delta \vdash t[s/x]: Y$  that satisfies  $\llbracket \rho[\delta/x] \rrbracket \sim \llbracket \rho \rrbracket \circ (id_\Gamma \otimes \bigotimes_A \llbracket \delta \rrbracket \otimes id_\Delta)$ .*

*Proof (of Theorem 1).* For any reduction  $s \longrightarrow t$  there exist decompositions  $s = E[s']$  and  $t = E[t']$  such that  $s' \longrightarrow t'$  is an instance of one of the reductions in Fig 5.

The proof then goes by induction on the structure of  $E$ . The base case where  $E$  is empty, amounts to showing the assertion for the basic reductions. For lack of space, we just spell out the first case where  $s \longrightarrow t$  has the form  $(\lambda x. q) p \longrightarrow q[p/x]$ .

Since  $\delta$  ends in a sequent with empty context, it cannot end with a structural rule. Hence, the last two rules in  $\delta$  must be  $(\rightarrow E)$  after  $(\rightarrow I)$ . Let  $\sigma$  and  $\tau$  be the derivations of  $x:B \cdot Z \vdash q: Y$  and  $\vdash p: Z$  that derive the premises of these rules. Then we can use  $\otimes E \circ \otimes I \sim id$  and  $(id \otimes \varepsilon) \circ (\eta \otimes id) \sim id$ , which are both easy to show directly, and the substitution lemma to calculate:

$$\begin{aligned}
\llbracket \delta \rrbracket &= (id_Y \otimes \varepsilon) \circ ((\text{swap} \circ \otimes E \circ \otimes I \circ \text{swap} \circ (\llbracket \sigma \rrbracket \otimes id_{\otimes_B Z}) \circ \eta) \otimes \bigotimes_B \llbracket \tau \rrbracket) \\
&\sim (id_Y \otimes \varepsilon) \circ (((\llbracket \sigma \rrbracket \otimes id_{\otimes_B Z}) \circ \eta) \otimes \bigotimes_B \llbracket \tau \rrbracket) \\
&\sim \llbracket \sigma \rrbracket \circ ((id_{\otimes_B Z} \otimes \varepsilon) \circ (\eta \otimes id_{\otimes_B Z})) \circ \bigotimes_B \llbracket \tau \rrbracket \\
&\sim \llbracket \sigma \rrbracket \circ \bigotimes_B \llbracket \tau \rrbracket \sim \llbracket \sigma[\tau/y] \rrbracket
\end{aligned}$$

Since  $\sigma[\tau/y]$  derives  $\vdash q[p/x]: Y$ , this concludes this case.

The induction step follows straightforwardly from the induction hypothesis. If, for example,  $E$  is  $\langle F, u \rangle$ , then  $\delta$  must end in rule  $(\otimes I)$  with two premises  $\vdash F[s']: Y$  and  $\vdash u: Z$ , derived by  $\sigma$  and  $\tau$ . We apply the induction hypothesis to  $\sigma$  to obtain  $\sigma'$ . Then we note that replacing  $\sigma$  with  $\sigma'$  in  $\delta$  gives us a derivation  $\rho$  of the required term. The assertion then follows:  $\llbracket \delta \rrbracket = \otimes I \circ (\llbracket \sigma \rrbracket \otimes \llbracket \tau \rrbracket) \sim \otimes I \circ (\llbracket \sigma' \rrbracket \otimes \llbracket \tau \rrbracket) = \llbracket \rho \rrbracket$ .  $\square$

We can conclude that the compilation in the previous section computes the same results as the standard reduction introduced here.

**Corollary 1.** *If  $\delta$  derives  $\vdash s : [A]$  and  $s \longrightarrow^* [v]$  for some  $v \in \mathcal{V}_A$ , then  $\Phi_{[\delta]}$  is the function of type  $\mathcal{V}_1 \times 1 \rightarrow \mathcal{V}_A$  that maps the unique element of its domain to  $v$ .*

This corollary follows by observing that  $\vdash [v] : [A]$  must be derived by some  $\rho$  with  $([\ ]\mathbf{I})$  as last rule, that  $\Phi_{[\rho]}$  is by definition the function in the corollary and that  $\Phi_{[\rho]} \sim \Phi_{[\delta]}$  follows from the theorem.

## 5 Logarithmic Space

In this section we describe a precise correspondence between the class of functions representable in INTML and the class of functions computable in logarithmic space.

First we must define how to represent functions from binary strings to binary strings in INTML. In principle, such functions can be programmed directly as terms of type  $A \cdot [B] \multimap [B]$ , where  $B$  is a working class type of binary strings. However, this would lead to linear space usage, so we use a more interactive type as discussed in Sec. 2.

With the constants *min* and *succ*, we can view any closed type  $[A]$  as a type of natural numbers that can represent the numbers from 0 to  $|\mathcal{V}_A| - 1$ : the number  $i$  is encoded by the normal form of  $[succ_A^i(\text{min}_A)]$ , for which we write  $\langle i \rangle_A$ .

Then, for every closed type  $B$ , binary strings of length at most  $|\mathcal{V}_B|$  can be represented as terms of type  $\mathbb{B}_A(B) = A \cdot [B] \multimap [3]$ : a binary string  $s$  is encoded by a function that when applied to  $\langle i \rangle_B$  returns:  $b$  if the  $i$ -th symbol in  $s$  is  $b$  and 2 if  $i$  exceeds the length of  $s$ . Write  $\langle s \rangle_{A,B}$  for the encoding of  $s \in \{0, 1\}^{\leq |\mathcal{V}_B|}$  in  $\mathbb{B}_A(B)$ .

Functions that take as inputs strings of arbitrary length can be represented using type variables. Let  $X$  be the upper class type  $A \cdot \mathbb{B}_B(\alpha) \multimap \mathbb{B}_C(D)$ , where  $\alpha$  is a type variable and where  $A, B, C$  and  $D$  are arbitrary types that may also contain  $\alpha$ , but not other type variables. If  $E$  is a closed type then  $X[E/\alpha]$  is a type of functions from strings of length at most  $|\mathcal{V}_E|$  to strings of length at most  $|\mathcal{V}_{D[E/\alpha]}|$ . We say that a term  $\vdash t : X$  represents a function  $\phi : \{0, 1\}^* \rightarrow \{0, 1\}^*$  if and only if:

- For every  $n$  and every  $x \in \{0, 1\}^n$ ,  $|\phi(x)| \leq |\mathcal{V}_{D[n/\alpha]}|$ ;
- For every  $n$ , every  $x \in \{0, 1\}^n$  and every  $E$  with  $|\mathcal{V}_E| \geq n$ , if  $\phi(x) = b_1 \dots b_m$  then for every  $i \leq |\mathcal{V}_{D[n/\alpha]}|$ , the term  $t[E/\alpha]\langle s \rangle_{B[E/\alpha], E} \langle i \rangle_{D[E/\alpha]}$  reduces to  $\langle b_i \rangle_3$  if  $i \leq m$  and to  $\langle 2 \rangle_3$  if  $i > m$ .

We remark that one may also use an alternate definition, for which the following theorems are also valid, but which does not refer to upper class reduction: instead of requiring  $t[E/\alpha]\langle s \rangle_{B[E/\alpha], E} \langle i \rangle_{D[E/\alpha]}$  to reduce to one of  $\langle b_i \rangle_3$  or  $\langle 2 \rangle_3$ , one may ask that the circuit for this term have the same behaviour as either  $\langle b_i \rangle_3$  or  $\langle 2 \rangle_3$ , depending on  $i$ .

**Theorem 2 (Logspace Soundness).** *If  $\phi : \{0, 1\}^* \rightarrow \{0, 1\}^*$  is represented by  $t$ , then  $\phi$  is computable in logarithmic space. Moreover, a LOGSPACE algorithm computing  $\phi$  is given by INTML-evaluation of  $t$ .*

*Proof.* Compiling  $t$  to a working class term yields a term  $x : X^- \vdash f : X^+$ . We now choose  $E_n = 2^{\lceil \log n \rceil}$  to substitute for  $\alpha$ . Clearly,  $|\mathcal{V}_{E_n}| \geq n$ , but  $|E_n| \leq 2(\log n + 1)$ . By Prop. 1, evaluation of  $f[E_n/\alpha]$  (i.e. computation of  $\phi$  on strings of length up to  $n$ ) can be implemented using space linear in  $|2^{\lceil \log n \rceil}|$ , thus logarithmic in  $n$ .

**Theorem 3 (Logspace Completeness).** *Any function  $\phi: \{0, 1\}^* \rightarrow \{0, 1\}^*$  computable in logarithmic space is represented by some upper class term  $t_\phi$ .*

The proof goes by an easy encoding of Offline Turing Machines in INTML. The step function of a given LOGSPACE OTM  $M$  can be mimicked by an upper class term  $| x:[\alpha] \multimap [3] \vdash \text{step}_M: [S(\alpha)] \multimap [S(\alpha) + S(\alpha)]$ , where  $x$  represents the input string and where  $S(\alpha)$  is a type with free variable  $\alpha$  that can encode the state of  $M$ . The term  $t_M$  can be obtained by passing  $\text{step}_M$  to the loop combinator.

## 6 Conclusion

We have found that the Int construction is a good way of structuring space bounded computation that can help us to understand the principles of space bounded functional programming. This view has guided the design of INTML, a simple language capturing LOGSPACE. Initial experience with an experimental implementation of INTML suggests that, with suitable type inference, INTML can be made to be quite usable.

We hope that our systematic approach will be helpful for developing INTML further. For instance, recent results on capturing non-deterministic token machines by the Int construction [9] may perhaps be used to develop a NLOGSPACE-version of INTML.

## References

1. Abramsky, S., Jagadeesan, R.: New Foundations for the Geometry of Interaction. *Inf. Comput.* 111(1), 53–119 (1994)
2. Abramsky, S., Haghverdi, E., Scott, P.J.: Geometry of interaction and linear combinatory algebras. *Math. Struct. in Comput. Sci.* 12(5), 625–665 (2002)
3. Asperti, A., Guerrini, S.: The optimal implementation of functional programming languages. Cambridge University Press, Cambridge (1998)
4. Baillot, P., Terui, K.: Light types for polynomial time computation in lambda calculus. *Inf. Comput.* 207(1), 41–62 (2009)
5. Bonfante, G.: Some programming languages for Logspace and Ptime. In: Johnson, M., Vene, V. (eds.) AMAST 2006. LNCS, vol. 4019, pp. 66–80. Springer, Heidelberg (2006)
6. Hasegawa, M.: On traced monoidal closed categories. *Math. Struct. in Comput. Sci.* 19(2), 217–244 (2009)
7. Joyal, A., Street, R., Verity, D.: Traced monoidal categories. *Math. Proc. Cambridge Philos. Soc.* 119(3), 447–468 (1996)
8. Katsumata, S.: Attribute grammars and categorical semantics. In: Aceto, L., Damgård, I., Goldberg, L.A., Halldórsson, M.M., Ingólfssdóttir, A., Walukiewicz, I. (eds.) ICALP 2008, Part II. LNCS, vol. 5126, pp. 271–282. Springer, Heidelberg (2008)
9. Katsumata, S., Hoshino, N.: Int construction and semibiproducts. RIMS-Technical Report 1676 (2009)
10. Kristiansen, L.: Neat function algebraic characterizations of LOGSPACE and LINSPEC. *Computational Complexity* 14, 72–88 (2005)
11. Mackie, I.: The geometry of interaction machine. In: POPL, pp. 198–208 (1995)

12. Mairson, H.G.: From hilbert spaces to dilbert spaces: Context semantics made simple. In: Agrawal, M., Seth, A.K. (eds.) FSTTCS 2002. LNCS, vol. 2556, pp. 2–17. Springer, Heidelberg (2002)
13. Melliès, P.A.: Functorial boxes in string diagrams. In: Ésik, Z. (ed.) CSL 2006. LNCS, vol. 4207, pp. 1–30. Springer, Heidelberg (2006)
14. Neergaard, P.M.: A functional language for logarithmic space. In: Chin, W.-N. (ed.) APLAS 2004. LNCS, vol. 3302, pp. 311–326. Springer, Heidelberg (2004)
15. Schöpp, U.: Stratified bounded affine logic for logarithmic space. In: LICS, pp. 411–420 (2007)
16. Szepietowski, A.: Turing Machines with Sublogarithmic Space. LNCS, vol. 843. Springer, Heidelberg (1994)
17. Winskel, G.: The Formal Semantics of Programming Languages. MIT Press, Cambridge (1993)