

Formally Deriving an STG Machine*

Alberto de la Encina
Departamento de Sistemas Informáticos y
Programación
Universidad Complutense de Madrid, Spain
Avda. Complutense s/n, 28040 Madrid
albertoe@sip.ucm.es

Ricardo Peña
Departamento de Sistemas Informáticos y
Programación
Universidad Complutense de Madrid, Spain
Avda. Complutense s/n, 28040 Madrid
ricardo@sip.ucm.es

ABSTRACT

Starting from P. Sestoft semantics for lazy evaluation, we define a new semantics in which normal forms consist of variables pointing to lambdas or constructions. This is in accordance with the more recent changes in the *Spineless Tagless G-machine* (STG) machine, where constructions only appear in closures (lambdas only appeared in closures already in previous versions). We prove the equivalence between the new semantics and Sestoft's. Then, a sequence of STG machines are derived, formally proving the correctness of each derivation. The last machine consists of a few imperative instructions and its distance to a conventional language is minimal.

The paper also discusses the differences between the final machine and the actual STG machine implemented in the Glasgow Haskell Compiler.

Categories and Subject Descriptors

D.3.1 [Programming Languages]: Formal Definitions and Theory—*semantics, syntax*; D.3.2 [Programming Languages]: Language Classifications—*applicative (functional) languages*; D.3.4 [Programming Languages]: Processors—*code generation, compilers*; F.3.2 [Logics and meanings of programs]: Semantics of Programming Languages—*operational semantics*

General Terms

Theory, Languages, Verification

Keywords

Functional programming, abstract machines, operational semantics, compiler verification

*Work partially supported by the Spanish project TIC 2000-0738.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PPDP'03, August 27–29, 2003, Uppsala, Sweden.

Copyright 2003 ACM 1-58113-705-2/03/0008 ...\$5.00.

1. INTRODUCTION

The *Spineless Tagless G-machine* (STG) [6] is at the heart of the *Glasgow Haskell Compiler* (GHC) [7] which is perhaps the Haskell compiler generating the most efficient code. For a description of Haskell language see [8]. Part of the secret for that is the set of analysis and transformations carried out at the intermediate representation level. Another part of the explanation is the efficient design and implementation of the STG machine.

A high level description of the STG can be found in [6]. If the reader is interested in a more detailed view, then the only available information is the Haskell code of GHC (about 80.000 lines, 12.000 of which are devoted to the implementation of the STG machine) and the C code of its different runtime systems (more than 40.000 lines)[1].

In this paper we provide a step-by-step derivation of the STG machine, starting from a description higher-level than that of [6] and arriving at a description lower-level than that.

Our starting point is a commonly accepted operational semantics for lazy evaluation provided by Peter Sestoft in [10] as an improvement of John Launchbury's well-known definition in [4]. Then, we present the following refinements:

1. A new operational semantics, which we call semantics *S3*—acknowledging that semantics 1 and 2 were defined by Mountjoy in a previous attempt [5]—, where normal forms may appear only in bindings.
2. A first machine, called STG-1, derived from *S3* in which explicit replacement of pointers for variables is done in expressions.
3. A second machine STG-2 introducing environments in closures, **case** alternatives, and in the control expression.
4. A third machine, called ISTG (I stands for *imperative*) with a very small set of elementary instructions, each one very easy to be implemented in a conventional language such as C.
5. A translation from the language of STG-2 to the language of ISTG in which the data structures of STG-2 are represented (or implemented) by the ISTG data structures.

$e \rightarrow x$	-- variable
$\lambda x.e$	-- lambda abstraction
$e x$	-- application
letrec $\overline{x_i = e_i}$ in e	-- recursive let
$C \overline{x_i}$	-- constructor application
case e of $\overline{C_i \overline{x_{ij}} \rightarrow e_i}$	-- case expression

Figure 1: Launchbury’s normalized λ -calculus

At each refinement, a formal proof of the soundness and completeness of the lower level with respect to the upper one is carried out¹. In the end, the final implementation is shown correct with respect to Sestoft’s operational semantics.

The main contribution of the work is showing that an efficient machine such as STG can be presented, understood, and formally reasoned about at different levels of abstraction. Also, there are some differences between the machine we arrive at and the actual STG machine implemented in the Glasgow Haskell Compiler. We argue that some design decisions in the actual STG machine are not properly justified.

The plan of the paper is as follows: after this introduction, in Section 2, a new language called FUN is introduced and the semantics $S3$ for this language is defined. Two theorems relating Launchbury’s original language and semantics to the new ones are presented. Section 3 defines the two machines STG-1 and STG-2. Some propositions show the consistency between both machines and the correctness and completeness of STG-1 with respect to $S3$, even though the latter creates more closures in the heap and produces different (but equivalent) normal forms. Section 4 defines machine ISTG and Section 5 defines the translation from STG-2 expressions to ISTG instructions. Two invariants are proved which show the correctness of the translation. Section 6 discusses the differences between our translation and the actual implementation done by GHC. Finally, Section 7 concludes.

2. A NEW SEMANTICS FOR LAZY EVALUATION

We begin by reviewing the language and semantics given by Sestoft as an improvement to Launchbury’s semantics. Both share the language given in Figure 1 where $\overline{A_i}$ denotes a vector A_1, \dots, A_n of subscripted entities. It is a normalized λ -calculus, extended with recursive **let**, constructor applications and **case** expressions. Sestoft’s normalization process forces constructor applications to be saturated and all applications to only have variables as arguments. Weak head normal forms are either lambda abstractions or constructions. Throughout this section, w will denote (weak head) normal forms.

Sestoft’s semantic rules are given in Figure 2. There, a judgement $\Gamma : e \Downarrow_A \Delta : w$ denotes that expression e , with its free variables bound in heap Γ , reduces to normal form w and produces the final heap Δ . When fresh pointers are created, freshness is understood w.r.t. $(\text{dom } \Gamma) \cup A$, where A contains the addresses of the closures under evaluation

¹The details of the proofs can be found in a technical report at one of the author’s page <http://dalila.sip.ucm.es/~albertoe>.

$\Gamma : \lambda x.e \Downarrow_A \Gamma : \lambda x.e$	<i>Lam</i>
$\Gamma : C \overline{p_i} \Downarrow_A \Gamma : C \overline{p_i}$	<i>Cons</i>
$\frac{\Gamma : e \Downarrow_A \Delta : \lambda x.e' \quad \Delta : e'[p/x] \Downarrow_A \Theta : w}{\Gamma : e p \Downarrow_A \Theta : w}$	<i>App</i>
$\frac{\Gamma : e \Downarrow_{A \cup \{p\}} \Delta : w}{\Gamma \cup [p \mapsto e] : p \Downarrow_A \Delta \cup [p \mapsto w] : w}$	<i>Var</i>
$\frac{\Gamma \cup [\overline{p_i} \mapsto \overline{e_i}] : \hat{e} \Downarrow_A \Delta : w}{\Gamma : \text{letrec } \overline{x_i = e_i} \text{ in } e \Downarrow_A \Delta : w}$ where $\overline{p_i}$ fresh	<i>Letrec</i>
$\frac{\Gamma : e \Downarrow_A \Delta : C_k \overline{p_j} \quad \Delta : e_k[\overline{p_j}/\overline{x_{kj}}] \Downarrow_A \Theta : w}{\Gamma : \text{case } e \text{ of } \overline{C_i \overline{x_{ij}} \rightarrow e_i} \Downarrow_A \Theta : w}$	<i>Case</i>

Figure 2: Sestoft’s natural semantics

(see rule *Var*). The notation \hat{e} in rule *Letrec* means the replacement of the variables x_i by the fresh pointers p_i . This is the only rule where new closures are created and added to the heap. We use the term *pointers* to refer to dynamically created free variables, bounded to expressions in the heap, and the term *variables* to refer to (lambda-bound, let-bound or case-bound) program variables. We consistently use p, p_i, \dots to denote free variables and x, y, \dots to denote program variables.

J. Mountjoy’s [5] had the idea of changing Launchbury-Sestoft’s language and semantics in order to get closer to the STG language, and then to derive the STG machine from the new semantics.

He developed two different semantics: In the first one, which we call semantics $S1$, the main change was that normal forms were either constructions (as they were in Sestoft’s semantics) or variables pointing to closures containing λ -abstractions, instead of just λ -abstractions. The reason for this was to forbid a λ -abstraction in the control expression as it happens in the STG machine. Another change was to force applications to have the form $x x_1$, i.e. consisting of a variable in the functional part. This is also what the STG language requires. These changes forced Mountjoy to modify the source language and to define a normalization from Launchbury’s language to the new one. Mountjoy proved that the normalization did not change the normal forms arrived at by both semantics.

The second semantics, which we call semantics $S2$, forced applications to be done at once to n arguments instead of doing it one by one. Correspondingly, λ -abstractions were allowed to have several arguments. This is exactly what the STG machine requires. Semantics $S2$ was informally derived and contained some mistakes. In particular, (cf. [5, pag. 171]) rule *App_M* makes a λ -abstraction to appear in the control expression, in contradiction with the desire of having λ -abstractions only in the heap.

Completing and correcting Mountjoy’s work we have defined a new semantics $S3$ in which the main changes in the source language w.r.t. Mountjoy’s are the following:

1. We force constructor applications to appear only in bindings, i.e. in heap closures. Correspondingly, normal forms are variables pointing to either λ -abstractions or constructions. We will use the term *lambda forms* to refer to λ -abstractions or constructions alike. The motivation for this decision is to generate more efficient code as it will be seen in Section 5.1.

$\Gamma[p \mapsto \lambda \bar{x}_i^n . e] : p \downarrow \Gamma : p$	<i>Lam_{S3}</i>
$\Gamma[p \mapsto C \bar{p}_i] : p \downarrow \Gamma : p$	<i>Cons_{S3}</i>
$\frac{\Gamma : e \downarrow \Delta[p \mapsto \lambda \bar{x}_i^n . \lambda \bar{y}_i^m . e'] : p}{\Gamma : e \bar{p}_i^n \downarrow \Delta \cup [q \mapsto \lambda \bar{y}_i^m . e'[\bar{p}_i/\bar{x}_i^n]] : q} \quad m, n > 0, \text{ q fresh}$	<i>App_{S3}</i>
$\frac{\Gamma : e \downarrow \Delta[p \mapsto \lambda \bar{x}_i^m . e'] : p \quad \Delta : e'[\bar{p}_i/\bar{x}_i^m] p_{m+1} \dots p_n \downarrow \Theta[q \mapsto w] : q}{\Gamma : e \bar{p}_i^n \downarrow \Theta : q} \quad n \geq m$	<i>App'_{S3}</i>
$\frac{\Gamma : e \downarrow \Delta[q \mapsto w] : q}{\Gamma \cup [p \mapsto e] : p \downarrow \Delta \cup [p \mapsto w] : q}$	<i>Var_{S3}</i>
$\frac{\Gamma \cup [\bar{p}_i \mapsto \hat{l}f_i] : \hat{e} \downarrow \Delta[p \mapsto w] : p \quad \bar{p}_i \text{ fresh}}{\Gamma : \mathbf{letrec} \ \bar{x}_i = \bar{l}f_i \ \mathbf{in} \ e \downarrow \Delta : p}$	<i>Letrec_{S3}</i>
$\frac{\Gamma : e \downarrow \Delta[p \mapsto C_k \bar{p}_j] : p \quad \Delta : e_k[\bar{p}_j/\bar{y}_{kj}] \downarrow \Theta[q \mapsto w] : q}{\Gamma : \mathbf{case} \ e \ \mathbf{of} \ C_i \ \bar{y}_{ij} \rightarrow e_i \downarrow \Theta : q}$	<i>Case_{S3}</i>

Figure 4: Semantics S3

PROOF. By induction on the number of reductions of Launchbury expressions. \square

Now we prove the equivalence between the two semantics. We consider only FUN expressions because it has been proved that the normalization does not change the meaning of an expression.

PROPOSITION 4. (*Sestoft* \Rightarrow S3, completeness of S3) For all $e \in \text{FUN}$ we have:

$$\{ \} : e \downarrow \Delta : w \Rightarrow \begin{cases} \{ \} : e \downarrow \Delta' [p \mapsto w'] : p \\ \exists \alpha. \quad \alpha \ w = w' \\ \quad \Delta' \subseteq_{\alpha} \Delta \end{cases}$$

PROOF. By induction on the number of reductions of FUN expressions. \square

PROPOSITION 5. (S3 \Rightarrow Sestoft, soundness of S3) For all $e \in \text{FUN}$ we have:

$$\{ \} : e \downarrow \Delta [p \mapsto w] : p \Rightarrow \begin{cases} \{ \} : e \downarrow \Delta' : w' \\ \exists \alpha. \quad \alpha \ w' = w \\ \quad \Delta' \subseteq_{\alpha} \Delta \end{cases}$$

PROOF. By induction on the number of reductions of FUN expressions. \square

Once adapted the source language to the STG language, we are ready to derive an STG-like machine from semantics S3.

3. A VERY SIMPLE STG MACHINE

Following a similar approach to Sestoft MARK-1 machine [10], we first introduce a very simple STG machine, which we will call STG-1, in which explicit variable substitutions are done. A configuration in this machine is a triple (Γ, e, S) where Γ represents the heap, e is the control expression and

S is the stack. The heap binds pointers to lambda forms which, in turn, may reference other pointers. The stack stores three kinds of objects: arguments p_i of pending applications, case alternatives *alts* of pending pattern matchings, and marks $\#p$ of pending updates.

In Figure 5, the transitions of the machine are shown. They look very close to the lazy semantics S3 presented in Section 2. For instance, the single rule for **letrec** in Figure 5 is a literal transcription of the *Letrec_{S3}* rule of Figure 4. The semantic rules for **case** and applications are split each one into two rules in the machine. The semantic rule for variable is also split into two in order to take care of updating the closure. So, in principle, an execution of the STG-1 machine could be regarded as the linearization of the semantic derivation tree by introducing an auxiliary stack.

But sometimes appearances are misleading. The theorem below shows that in fact STG-1 builds less closures in the heap than the semantics and it may arrive to different (but semantically equivalent) normal forms. In order to prove the soundness and completeness of STG-1, we first enrich the semantics with a stack parameter S in the rules. The new rules for \downarrow_S (only those which modify S) are shown in Figure 6. It is trivial to show that the rules are equivalent to the ones in Figure 4 as the stack is just an observation of the derivations. It may not influence them. The following theorem establishes the correspondence between the (enriched) semantics and the machine.

PROPOSITION 6. Given Γ, e and S , then $\Gamma : e \downarrow_S \Delta [p \mapsto w] : p$ iff $(\Gamma, e, S) \rightarrow^* (\Delta', p', S')$, where

1. $\Delta' \subseteq \Delta$
2. if $\Delta [p \mapsto C \bar{p}_i^n]$ then $S = S', p = p'$ and $\Delta' [p' \mapsto C \bar{p}_i^n]$
3. if $\Delta [p \mapsto \lambda \bar{x}_i^n . e']$ then there exists $m \geq 0$ s.t. $\Delta' [p' \mapsto \lambda \bar{y}_i^m . \lambda \bar{x}_i^n . e'']$ and $S' = \bar{q}_i^m : S$ and $e' = e''[\bar{q}_i/\bar{y}_i^m]$

	Heap	Control	Stack	rule
	Γ	letrec $\{\overline{x_i = lf_i}\}$ in e	S	letrec ⁽¹⁾
\Rightarrow	$\Gamma \cup [\overline{p_i \mapsto lf_i[p_j/x_j]}]$	$e[\overline{p_i/x_i}]$	S	
	Γ	case e of $alts$	S	case1
\Rightarrow	Γ	e	$alts : S$	
	$\Gamma[p \mapsto C_k \overline{p_i}]$	p	$\overline{C_j \overline{y_{ji}} \rightarrow e_j} : S$	case2
\Rightarrow	Γ	$e_k[\overline{p_i/y_{ki}}]$	S	
	Γ	$e \overline{p_i}^n$	S	app1
\Rightarrow	Γ	e	$\overline{p_i}^n : S$	
	$\Gamma[p \mapsto \lambda \overline{x_i}^n . e]$	p	$\overline{p_i}^n : S$	app2
\Rightarrow	Γ	$e[\overline{p_i/x_i}]^n$	S	
	$\Gamma \cup [p \mapsto e']$	p	S	var1
\Rightarrow	Γ	e'	$\#p : S$	
	$\Gamma[p \mapsto \lambda \overline{x_i}^k . \lambda \overline{y_i}^m . e]$	p	$\overline{p_i}^k : \#q : S$	var2
\Rightarrow	$\Gamma \cup [q \mapsto p \overline{p_i}^k]$	p	$\overline{p_i}^k : S$	
	$\Gamma[p \mapsto C_k \overline{p_i}]$	p	$\#q : S$	var3
\Rightarrow	$\Gamma \cup [q \mapsto C_k \overline{p_i}]$	p	S	

(¹) $\overline{p_i}$ are distinct and fresh w.r.t. Γ , **letrec** $\{\overline{x_i = lf_i}\}$ **in** e , and S

Figure 5: The STG-1 Machine

PROOF. By induction on the number of reductions of FUN expressions. \square

The proposition shows that the semantic rule App_{S3} of Figure 4 is not literally transcribed in the machine. The machine does not create intermediate lambdas in the heap unless an update is needed. Rule *app2* in Figure 5 applies a lambda always to *all* its arguments provided that they are in the stack. For this reason, a lambda with more parameters than that of the semantics may be arrived at as normal form of a functional expression. Also for this reason, an update mark may be interspersed with arguments in the stack when a lambda is reached (see rule *var2* in Figure 5). A final implication is that the machine may stop with m pending arguments in the stack and a variable in the control expression pointing to a lambda with $n > m$ parameters. The semantics always ends a derivation with a variable as normal form and an empty stack.

Again, following Sestoft and his MARK-2 machine, once we have proved the soundness and completeness of STG-1, we introduce STG-2 having environments instead of explicit variable substitutions. Also, we add trimmers to this machine so that environments kept in closures and in **case** alternatives only reference the free variables of the expression instead of *all* variables in scope. A configuration of STG-2 is a quadruple (Γ, e, E, S) where E is the environment of e , the alternatives are pairs $(alts, E)$, and a closure is a pair (lf, E) . Now expressions and lambda forms keep their original variables and the associated environment maps them to pointers in the heap. The notation $E \upharpoonright^t$ means the trimming of environment E to the trimmer t . A trimmer is just a collection of variable names. The resulting machine is shown in Figure 7.

PROPOSITION 7. *Given a closed expression e_0 . $(\{\}, e_0, []) \xrightarrow{STG-1} (\Delta, q, \overline{p_i}^n)$ where either:*

- $\Delta[q \mapsto C \overline{q_i}^m] \wedge n = 0$

- or $\Delta[q \mapsto \lambda \overline{x_i}^m . e] \wedge m > n \geq 0$

if and only if $(\{\}, e_0, \{\}, []) \xrightarrow{STG-2} (\Gamma, x, E[x \mapsto q], \overline{p_i}^n)$ and either:

- $\Gamma[q \mapsto (C \overline{x_i}^m, \{\overline{x_i \mapsto q_i}^m\})] \wedge n = 0$
- or $\Gamma[q \mapsto (\lambda \overline{x_i}^m . e', E')] \wedge m > n \geq 0 \wedge e = E' e'$.

PROOF. By induction on the number of reductions. \square

4. AN IMPERATIVE STG MACHINE

In this Section we ‘invent’ an imperative STG machine, called ISTG, by defining a set of machine instructions and an operational semantics for them in terms of the state transition that each instruction produces. In fact, this machine tries to provide an intermediate level of reasoning between the STG-2 machine and the final C implementation. In the actual GHC implementation, ‘below’ the operational description of [6] we find only a translation to C. By looking at the compiler and at the runtime system listings, one can grasp some details, but many others are lost. We think that the gap to be bridged is too high. Moreover, it is not possible to reason about the correctness of the implementation when so many details are introduced at once. The ISTG architecture has been inspired by the actual implementation of the STG machine done by GHC, and the ISTG instructions have been derived from the STG-2 machine by analyzing the elementary steps involved in every machine transition.

An ISTG machine configuration consists of a 5-tuple $(is, S, node, \Gamma, cs)$, where is is a machine instruction sequence ended with the instruction ENTER or RETURNCON, S is the stack, $node$ is a heap pointer pointing to the closure under execution (the one to which is belongs to), Γ is the heap and cs is a code store where the instruction sequences resulting from compiling the program expressions are kept.

We will use the following notation: a for pointers to closures in Γ , as and ws for lists of such pointers, and p for

$$\begin{array}{c}
\frac{\Gamma : e \downarrow_{\overline{p_i}^n : S} \Delta[p \mapsto \lambda \overline{x_i}^n . \lambda \overline{y_i}^m . e'] : p}{\Gamma : e \overline{p_i}^n \downarrow_S \Delta \cup [q \mapsto \lambda \overline{y_i}^m . e'[p_i/x_i]] : q} \quad m, n > 0, \text{ q fresh} \\
\text{App}_{S3} \\
\\
\frac{\Gamma : e \downarrow_{\overline{p_i}^n : S} \Delta[p \mapsto \lambda \overline{x_i}^m . e'] : p \quad \Delta : e'[\overline{p_i/x_i}^m] \quad p_{m+1} \dots p_n \downarrow_S \Theta[q \mapsto w] : q}{\Gamma : e \overline{p_i}^n \downarrow_S \Theta : q} \quad n \geq m \\
\text{App}'_{S3} \\
\\
\frac{\Gamma : e \downarrow_{\#p : S} \Delta[q \mapsto w] : q}{\Gamma \cup [p \mapsto e] : p \downarrow_S \Delta \cup [p \mapsto w] : q} \\
\text{Var}_{S3} \\
\\
\frac{\Gamma : e \downarrow_{alts : S} \Delta[p \mapsto C_k \overline{p_j}] : p \quad \Delta : e_k[\overline{p_j/y_{kj}}] \downarrow_S \Theta[q \mapsto w] : q}{\Gamma : \text{case } e \text{ of } alts \downarrow_S \Theta : q} \\
\text{Case}_{S3}
\end{array}$$

Figure 6: The enriched semantics

pointers to code fragments in cs . By $cs[p \mapsto is]$ we denote that the code store cs maps pointer p to the instruction sequence is and, by $cs[p \mapsto \overline{is_i}^n]$, that cs maps p to a vectored set of instruction sequences is_1, \dots, is_n , each one corresponding to an alternative of a **case** expression with n constructors C_1, \dots, C_n . Also, $S!i$ will denote the i -th element of the stack S counting from the top and starting at 0. Likewise, $node_\Gamma!i$ will denote the i -th free variable of the closure pointed to by $node$ in Γ , this time starting at 1.

Stack S may store pointers a to closures in Γ , pointers p to code sequences and code alternatives in cs , and update marks $\#a$ indicating that closure pointed to by a must be updated. A closure is a pair (p, ws) where p is a pointer to an instruction sequence is in cs , and ws is the closure environment, having a heap pointer for every free variable in the expression whose translation is is .

These representation decisions are very near to the GHC implementation. In its runtime system all these elements (stack, heap, node register and code) are present [9]. Our closures are also a small simplification of theirs.

In Figure 8, the ISTG machine instructions and its operational semantics are shown. The machine instructions BUILDENV, PUSHALTS and UPDTMARK roughly correspond to the three possible pushing actions of machine STG-2. The SLIDE instruction has no clear correspondence in the STG-2. As we will see in Section 5, it will be used to change the current environment when a new closure is entered. Instructions ALLOC and BUILDCLS will implement heap closure creation in the *letrec* rule of STG-2. Both BUILDENV and BUILDCLS make use of a list of pairs, each pair indicating whether the source variable is located in the stack or in the current closure. Of course, it is not intended this test to be done at runtime. An efficient translation of these ‘machine’ instructions to an imperative language will generate the appropriate copy statement for each pair.

Instructions ENTER and RETURNCON are typical of the actual STG machine as described in [6]. It is interesting to note that it has been possible to describe our previous STG machines without any reference to them. In our view, they belong to ISTG, i.e. to a lower level of abstraction. Finally, instruction ARGCHECK, which implements updates with lambda normal forms, is here at the same level of abstraction as RETURNCON, which implements updates with constructions normal forms. Predefined code is stored in cs for updating with a partial application and for blackholing

a closure under evaluation. The corresponding code pointers are respectively called p_{pap}^{n+1} and p_{bh} in Figure 8. The associated code is the following:

$$\begin{aligned}
p_{bh} &= [] \\
p_{pap}^{n+1} &= [BUILDENV [(NODE, 1), \dots, (NODE, n+1)], \\
&\quad ENTER]
\end{aligned}$$

The code of a blackhole just blocks the machine as there is no instruction to execute. There is predefined code for partial applications with different values for n . The code just copies the closure into the stack and jumps to the first pointer that is assumed to be pointing to a λ -abstraction closure.

The translation to C of the 9 instructions of the ISTG should appear straightforward for the reader. For instance, BUILDCLS and BUILDENV can be implemented by a sequence of assignments, copying values to/from the stack and the heap; PUSHALTS, UPDTMARK and ENTER do straightforward stack manipulation; SLIDE is more involved but can be easily translated to a sequence of loops moving information within the stack to collapse a number of stack fragments. The more complex ones are RETURNCON and ARGCHECK. Both contains a loop which updates the heap with normal forms (respectively, constructions and partial applications) as long as they encounter update marks in the stack. Finally, the installation of a new instruction sequence in the control made by ENTER and RETURNCON are implemented by a simple jump.

5. FORMAL TRANSLATION FROM STG-2 TO ISTG

In this Section, we provide first the translation schemes for the FUN expressions and lambda forms and then prove that this translation correctly implements the STG-2 machine on top of the ISTG machine. Before embarking into the details, we give some hints to intuitively understand the translation:

- The ISTG stack will represent not only the STG-2 stack, but also (part of) the current environment E and all the environments associated to pending **case** alternatives. So, care must be taken to distinguish between environments and other objects in the stack.
- The rest of the current environment E is kept in the current closure. The translation knows where each free

	Heap	Control	Environment	Stack	rule
\Rightarrow	Γ	letrec $\{\overline{x_i = lf_i} \mid t_i\}$ in e	E	S	letrec ⁽¹⁾
\Rightarrow	$\Gamma \cup [\overline{p_i \mapsto (lf_i, E' \mid t_i)}]$	e	E'	S	
\Rightarrow	Γ	case e of $alts \mid^t$	E	S	case1
\Rightarrow	Γ	e	E	$(alts, E' \mid^t) : S$	
\Rightarrow	$\Gamma[p \mapsto (C_k \overline{x_i}, \{\overline{x_i \mapsto p_i}\})]$	x	$E\{x \mapsto p\}$	$(alts, E') : S$	case2 ⁽²⁾
\Rightarrow	Γ	e_k	$E' \cup \{\overline{y_{ki} \mapsto p_i}\}$	S	
\Rightarrow	Γ	$e \overline{x_i}^n$	$E\{\overline{x_i \mapsto p_i}^n\}$	S	app1
\Rightarrow	Γ	e	E	$\overline{p_i}^n : S$	
\Rightarrow	$\Gamma[p \mapsto (\lambda \overline{x_i}^n. e, E')]$	x	$E\{x \mapsto p\}$	$\overline{p_i}^n : S$	app2
\Rightarrow	Γ	e	$E' \cup \{\overline{x_i \mapsto p_i}^n\}$	S	
\Rightarrow	$\Gamma \cup [p \mapsto (e, E')]$	x	$E\{x \mapsto p\}$	S	var1
\Rightarrow	Γ	e	E'	$\#p : S$	
\Rightarrow	$\Gamma[p \mapsto (\lambda \overline{x_i}^k. \lambda \overline{y_i}^n. e, E')]$	x	$E\{x \mapsto p\}$	$\overline{p_i}^k : \#q : S$	var2 ⁽³⁾
\Rightarrow	$\Gamma \cup [q \mapsto (x \overline{x_i}^k, E'')]$	x	E	$\overline{p_i}^k : S$	
\Rightarrow	$\Gamma[p \mapsto (C_k \overline{x_i}, E')]$	x	$E\{x \mapsto p\}$	$\#q : S$	var3
\Rightarrow	$\Gamma \cup [q \mapsto (C_k \overline{x_i}, E')]$	x	E	S	

- ⁽¹⁾ $\overline{p_i}$ are distinct and fresh w.r.t. Γ , **letrec** $\{\overline{x_i = lf_i}\}$ **in** e , and S . $E' = E \cup \{\overline{x_i \mapsto p_i}\}$
⁽²⁾ Expression e_k corresponds to alternative $C_k \overline{y_{ki}} \mapsto e_k$ in $alts$
⁽³⁾ $E'' = \{x \mapsto p, \overline{x_i \mapsto p_i}^k\}$

Figure 7: The STG-2 machine

variable is located by maintaining two compile-time environments ρ and η . The first one ρ corresponds to the environment kept in the stack, while the second one η corresponds to the free variables accessed through the *node* pointer.

- The stack can be considered as divided into big blocks separated by code pointers p pointing to case alternatives. Each big block topped with such a pointer corresponds to the environment of the associated alternatives.
- In turn, each big block can be considered as divided into small blocks, each one topped with a set of arguments of pending applications. The compile-time environment ρ is likewise divided into big and small blocks, so reflecting the stack structure.
- When a variable is reached in the current instruction sequence, an ENTER instruction is executed. This will finish the current sequence and start a new one. The upper big block of the stack must be deleted (corresponding to changing the current environment) but arguments of pending applications must be kept. This stack restructuring is accomplished by a SLIDE operation with an appropriate argument.

Definition 2. A stack environment ρ is a list $[(\delta_k, m_k, n_k), \dots, (\delta_1, m_1, n_1)]$ of blocks. It describes the variables in the stack starting from the top. In a block (δ, m, n) , δ is an environment mapping exactly $m - |n|$ program variables to disjoint numbers in the range $1..m - |n|$. The empty environment, denoted ρ_\emptyset is the list $[(\{\}, 0, 0)]$.

A block (δ, m, n) corresponds to a small block in the above explanation. Blocks with $n = -1$, are topped with a code pointer pointing to alternatives. So, they provide a separation between big blocks. The upper big block consists of all

the small blocks up to (and excluding) the first small block with $n = -1$. Blocks with $n > 0$ have $m - n$ free variables and are topped with n arguments of pending applications. The upper block is the only one with $n = 0$ meaning that it is not still closed and that it can be extended.

Definition 3. A closure environment η with n variables is a mapping from these variables to disjoint numbers in the range $1..n$.

Definition 4. The offset of a variable x in ρ from the top of the stack, denoted ρx , is given by

$$\rho x \stackrel{\text{def}}{=} \left(\sum_{i=l}^k m_i \right) - \delta_l x, \text{ being } x \in \text{dom } \delta_l$$

If the initial closed expression to be translated has different names for bound variables, then the compile time environments ρ and η will never have duplicate names. It will be proved below that every free variable of an expression being compiled will necessarily be either in ρ or in η , and never in both. This allows us to introduce the notation $(\rho, \eta) x$ to mean

$$(\rho, \eta) x \stackrel{\text{def}}{=} \begin{cases} (STACK, \rho x) & \text{if } x \in \text{dom } \rho \\ (NODE, \eta x) & \text{if } x \in \text{dom } \eta \end{cases}$$

The stack environment may suffer a number of operations: closing the current small block with a set of arguments, enlarging the current small block with new bindings, and closing the current big block with a pointer to **case** alternatives. These are formally defined as follows.

Definition 5. The following operations with stack environments are defined:

1. $((\delta, m, 0) : \rho) + n \stackrel{\text{def}}{=} (\{\}, 0, 0) : (\delta, m + n, n) : \rho$

	Instructions	Stack	Node	Heap	Code	
control						
\Rightarrow	$[ENTER]$ is	$a : S$ S	$node$ a	$\Gamma[a \mapsto (p, ws)]$ Γ	$cs[p \mapsto is]$ cs	
\Rightarrow	$[RETURNCON C_k^m]$ is_k	$p : S$ S	$node$ $node$	Γ Γ	$cs[p \mapsto \overline{is_i^n}]$ cs	
\Rightarrow	$[RETURNCON C_k^m]$ $[RETURNCON C_k^m]$	$\#a : S$ S	$node$ $node$	$\Gamma[a \mapsto (p_{bh}, as),$ $node \mapsto (p, ws)]$ $\Gamma[a \mapsto (p, ws)]$	cs cs	
\Rightarrow	$ARGCHECK m : is$ is	$\overline{a_i^m} : S$ $\overline{a_i^m} : S$	$node$ $node$	Γ Γ	cs cs	
\Rightarrow	$ARGCHECK m : is$ $ARGCHECK m : is$	$\overline{a_i^n} : \#a : S$ $\overline{a_i^n} : S$	$node$ $node$	$\Gamma[a \mapsto (p_{bh}, ws)]$ $\Gamma[a \mapsto (p_{pap}^{n+1}, node : \overline{a_i^n})]$	cs cs	$n < m$
heap						
\Rightarrow	$ALLOC m : is$ is	S $a_m : S$	$node$ $node$	Γ Γ'	cs cs	(¹)
\Rightarrow	$BUILDCLS i p \overline{z_i^n} : is$ is	S S	$node$ $node$	Γ $\Gamma[S!i \mapsto (p, \overline{a_i^n})]$	cs cs	(²)
stack						
\Rightarrow	$BUILDENV \overline{z_i^n} : is$ is	S $\overline{a_i^n} : S$	$node$ $node$	Γ Γ	cs cs	(²)
\Rightarrow	$PUSHALTS p : is$ is	S $p : S$	$node$ $node$	Γ Γ	cs cs	
\Rightarrow	$UPDTMARK : is$ is	S $\#node : S$	$node$ $node$	$\Gamma[node \mapsto (p, ws)]$ $\Gamma[node \mapsto (p_{bh}, ws)]$	cs cs	
\Rightarrow	$SLIDE (\overline{n_k, m_k})^l : is$ is	$\overline{a_{kj}^{n_k}} : \overline{b_{kj}^{m_k}}^l : S$ $\overline{a_{kj}^{n_k}}^l : S$	$node$ $node$	Γ Γ	cs cs	

(¹) a_m is a pointer to a new closure with space for m free variables, and Γ' is the resulting heap after the allocation

(²) $a_i = \begin{cases} S!i & \text{if } z_i = (STACK, i) \\ node_{\Gamma}!i & \text{if } z_i = (NODE, i) \end{cases}$

Figure 8: The ISTG machine

2. $((\delta, m, 0) : \rho) + (\{\overline{x_i \mapsto j_i^n}\}, n) \stackrel{\text{def}}{=} (\delta \cup \{\overline{x_i \mapsto m + j_i^n}\}, m+n, 0) : \rho$
3. $((\delta, m, 0) : \rho)^+ \stackrel{\text{def}}{=} (\{\}, 0, 0) : (\delta, m+1, -1) : \rho$

5.1 Translation schemes

Functions trE and trA respectively translate a FUN expression and a **case** alternative to a sequence of ISTG machine instructions; function $trAs$ translates a set of alternatives to a pointer to a vectored set of machine instruction sequences in the code store; and function trB translates a lambda form to a pointer to a machine instruction sequence in the code store. The translation schemes are shown in Figure 9.

The notation $\dots \& cs[p \mapsto \dots]$ means that the corresponding translation scheme has a side effect which consists of creating a code sequence in the code store cs and pointing it by the code pointer p .

PROPOSITION 8. (static invariant) Given a closed expression e_0 with different bound variables and an initial call $trE e_0 \rho_0 \{\}$, in all internal calls of the form $trE e \rho \eta$:

1. The stack environment has the form $\rho = (\delta, m, 0) : \rho'$. Moreover, there is no other block (δ', m', n) in ρ' with $n = 0$. Consequently, all environment operations in the above translation are well defined.
2. All free variables of e are defined either in ρ or in η . Moreover, $\text{dom } \rho \cap \text{dom } \eta = \emptyset$.
3. The last instruction generated for e is **ENTER**. Consequently, the main instruction sequence and all sequences corresponding to **case** alternatives and to non-constructor closures, end in an **ENTER**.

PROOF. (1) and (2) are proved by induction on the tree structure of calls to trE ; (3) is proved by structural induction on FUN expressions. \square

In order to prove the correctness of the translation, we only need to consider ISTG machine configurations of the form $(is, S_I, node, \Gamma_I, cs)$ in which is is generated by a call to trE for some expression e and environments ρ, η . We call these stable configurations. We enrich then these configurations with three additional components: the environments

$$\begin{aligned}
trE \ (e \ \overline{x_i^n}) \ \rho \ \eta &= [BUILDENV \ \overline{(\rho, \eta) \ x_i^n}] ++ \\
&\quad trE \ e \ (\rho + n) \ \eta \\
trE \ (\mathbf{case} \ e \ \mathbf{of} \ alts \ |\overline{x_i^n}) \ \rho \ \eta &= [BUILDENV \ zs, PUSHALTS \ p] ++ \ trE \ e \ \rho'^{\#} \ (\eta - xs) \\
\text{where } p &= trAs \ alts \ \rho' \\
\rho' &= \rho + (\{xs_j \mapsto m - j + 1^m\}, m) \\
xs &= [x \mid x \leftarrow \overline{x_i^n} \wedge x \in dom \ \eta] \\
zs &= [(node, \eta \ x) \mid x \leftarrow xs] \\
m &= |xs| \\
trE \ (\mathbf{letrec} \ x_i = lf_i \ |\overline{y_{ij}^{m_i n}} \ \mathbf{in} \ e) \ \rho \ \eta &= [ALLOC \ m_n, \dots, ALLOC \ m_1] ++ \\
&\quad [BUILDCLS \ (i-1) \ p_i \ zs_i^n] ++ \\
&\quad trE \ e \ \rho' \ \eta \\
\text{where } \rho' &= \rho + (\{x_i \mapsto n - i + 1^n\}, n) \\
p_i &= trB \ (lf_i \ |\overline{y_{ij}^{m_i}}), \quad i \in \{1..n\} \\
zs_i &= (\rho', \eta) \ y_{ij}^{m_i}, \quad i \in \{1..n\} \\
trE \ x \ \rho \ \eta &= [BUILDENV \ [(\rho, \eta) \ x], \\
&\quad SLIDE \ ((1, 0) : ms), \\
&\quad ENTER] \\
\text{where } ms &= map \ (\backslash(-, m, n) \rightarrow (n, m - n)) \ (takeWhile \ nn \ \rho) \\
nn \ (-, m, -1) &= False \\
nn \ - &= True \\
trAs \ (\overline{alt_i^n}) \ \rho &= p \ \& \ cs[p \mapsto \overline{trA \ alt_i \ \rho}^n] \\
trA \ (C \ \overline{x_i^n} \rightarrow e) \ \rho &= trE \ e \ \rho \ \{\overline{x_i \mapsto i^n}\} \\
trB \ (C_k^n \ \overline{x_i^n} \ |\overline{x_i^n}) &= p \ \& \ cs[p \mapsto [RETURNCON \ C_k^n]] \\
trB \ (\lambda \overline{x_i^l} . e \ |\overline{y_j^n}) &= p \ \& \ cs[p \mapsto [ARGCHECK \ l] ++ trE \ e \ \rho \ \eta] \\
\text{where } \rho &= [(\{x_i \mapsto l - i + 1^l\}, l, 0)] \\
\eta &= \{y_j \mapsto j^n\} \\
trB \ (e \ |\overline{y_j^n}) &= p \ \& \ cs[p \mapsto [UPDTMARK] ++ trE \ e \ \rho_\emptyset \ \eta] \\
\text{where } \eta &= \{y_j \mapsto j^n\}
\end{aligned}$$

Figure 9: Translation schemes from STG-2 to ISTG

ρ and η used to generate is , and an environment stack S_{env} containing a sequence of stack environments. The environments in S_{env} are in one to one correspondence with **case** pointers stored in S_I . Initially S_{env} is empty. Each time an instruction PUSHALTS is executed (see trE definition for **case**), the environment ρ' the corresponding alternatives are compiled with, is pushed onto stack S_{env} . Each time a RETURNCON pops a **case** pointer, stack S_{env} is also pop-ed. So, enriched ISTG configurations have the form $(is, \rho, \eta, S_I, S_{env}, node, \Gamma_I, cs)$.

Definition 6. A STG-2 environment E is equivalent to an ISTG environment defined by $\rho, \eta, S_I, \Gamma_I$ and $node$, denoted $E \equiv (\rho, S_I, \eta, \Gamma_I, node)$ if $dom \ E \subseteq dom \ \rho \cup dom \ \eta$ and

$$\begin{aligned}
\forall x \in dom \ E \quad E \ x = S_I ! (\rho \ x) &\quad \text{if } x \in dom \ \rho \\
E \ x = node_{\Gamma_I} ! (\eta \ x) &\quad \text{if } x \in dom \ \eta
\end{aligned}$$

Definition 7. A STG-2 stack S is equivalent to a triple (ρ, S_I, S_{env}) of an ISTG enriched configuration, denoted $S \equiv (\rho, S_I, S_{env})$, if

1. Whenever $\rho = (\delta, m, 0) : \rho'$, then $S_I = \overline{a_i^m} : S'_I$ and $S \equiv (\rho', S'_I, S_{env})$
2. Whenever $\rho = (\delta, m, n) : \rho', n > 0$, then $S = \overline{a_i^n} : S'$, $S_I = \overline{a_i^n} : \overline{b_j^{m-n}} : S'_I$ and $S' \equiv (\rho', S'_I, S_{env})$
3. Whenever $\rho = (\delta, m, -1) : \rho'$, then $S = (alts, E) : S'$, $S_I = p_{alts} : S'_I$, $S_{env} = \rho_{alts} : S'_{env}$, $p_{alts} = trAs \ alts \ \rho_{alts}$, $E \equiv (\rho_{alts}, S'_I, -, -, -)$ and $S' \equiv (\rho_{alts}, S'_I, S'_{env})$
4. Whenever $S = \#a : S'$ and $S_I = \#a : S'_I$, then $S' \equiv (\rho, S'_I, S_{env})$
5. Additionally, $[] \equiv (\{\}, [], [])$

Definition 8. A STG-2 heap Γ is equivalent to an ISTG pair (Γ_I, cs) , denoted $\Gamma \equiv (\Gamma_I, cs)$, if for all p we have $\Gamma[p \mapsto (lf \mid \overline{x_i^n}, E)]$ if and only if $\Gamma_I[p \mapsto (q, ws)]$, $cs[q \mapsto is]$, $is = trB (lf \mid \overline{x_i^n})$ and $ws = E \overline{x_i^n}$.

Definition 9. A STG-2 configuration is equivalent to an ISTG enriched stable configuration, denoted $(\Gamma, e, E, S) \equiv (is, \rho, \eta, S_I, S_{env}, node, \Gamma_I, cs)$ if

1. $\Gamma \equiv (\Gamma_I, cs)$
2. $is = trE \ e \ \rho \ \eta$
3. $E \equiv (\rho, S_I, \eta, \Gamma_I, node)$
4. $S \equiv (\rho, S_I, S_{env})$

PROPOSITION 9. (*dynamic invariant*) *Given a closed expression e_0 with different bound variables and initial STG-2 and ISTG configurations, respectively $(\{\}, e_0, \{\}, [])$ and $(trE \ e_0 \ \rho_0 \ \{\}, \rho_0, \{\}, [], [], \perp, \{\}, cs)$, where cs is the code store generated by the whole translation of e_0 , then both machines evolve through equivalent configurations.*

PROOF. By induction on the number of transitions of both machines. Only transitions between ISTG stable configurations are considered. \square

COROLLARY 10. *The translation given in Section 5.1 is correct.*

6. DIFFERENCES WITH THE ACTUAL STG MACHINE

There are some differences between the machine translation presented in Section 5 and the actual code generated by GHC. Some are just omissions, other are non-substantial differences and some other are deeper ones.

In the first group it is the treatment of basic values, very elaborated in GHC (see for example [3]) and completely ignored here. We have preferred to concentrate our study in the functional kernel of the machine but, of course, a formal reasoning about this aspect is a clear continuation of our work.

In the second group it is the optimization of update implementation. In GHC, updates can be done either by indirection or by closure creation, depending on whether there is enough space or not in the old closure to do update in place. This implies to keep closure size information somewhere. GHC keeps it in the so called *info table*, a static part shared by all closures created from the same bind. This table forces an additional indirection to access the closure code. Our model has simplified these aspects. We understand also that stack restructuring, as the one performed by our SLIDE instruction, is not implemented in this way by GHC. Apparently, stubbing of non used stack positions is done instead. An efficiency study could show which implementation is better. The cost of our SLIDE instruction is in $O(n)$, being n the number of arguments to be preserved in the stack when the current environment is discarded.

Perhaps the deeper difference between our derived machine and the actual STG is our insistence in that FUN applications should have the form $e \ \overline{x_i^n}$ instead of $x \ \overline{x_i^n}$ as it is the case in the STG language. This decision is not justified in the GHC papers and perhaps could have a noticeable negative impact in performance. In a lazy language, the functional part of an application should be eagerly evaluated, but GHC does it lazily. This implies constructing a number

of closures that will be immediately entered (and perhaps updated afterwards), with a corresponding additional cost both in space and time. Our translation avoids creating and entering these closures. If the counter-argument were having the possibility of sharing functional expressions, this is always available in FUN since a variable is a particular case of an expression. What we claim is that the normalization process in the Core-to-STG translation should not introduce unneeded sharing.

7. CONCLUSIONS

We have presented a stepwise derivation of a (well known) abstract machine starting from Sestoft's operational semantics, going through several intermediate machines and arriving at an imperative machine very close to a conventional imperative language. This strategy of adding a small amount of detail in each step has allowed us both to provide insight on fundamental decisions underlying the STG design and, perhaps more importantly, to be able to show the correctness of each refinement with respect to the previous one and, consequently, the correctness of the whole derivation. To our knowledge, this is the first time that formal translation schemes and a formal proof of correctness of the STG to C translation has been done.

Our previous work [2] followed a different path: it showed the soundness and completeness of a STG-like machine called STG-1S (laying somewhere between machines STG-2 and ISTG of this paper) with respect to Sestoft's semantics. The technique used was also different: a bisimulation between the STG-1S machine and Sestoft's MARK-2 machine was proved. We got the inspiration for the strategy followed here from Mountjoy [5] and in Section 2 we have explained the differences between his and our work. The previous machines of all these works, including STG-1 and STG-2 of this paper, are very abstract in the sense that they deal directly with functional expressions. The new machine ISTG introduced here is a really low level machine dealing with raw imperative instructions and pointers. Two contributions of this paper have been to bridge this big gap by means of the translations schemes and the proof of correctness of this translation.

Our experience is that formal reasoning about even well known products always reveals new details, give new insight, makes good decisions more solid and provides trust in the behavior of our programs.

8. REFERENCES

- [1] A. at URL: <http://www.haskell.org/ghc/>.
- [2] A. Encina and R. Peña. Proving the Correctness of the STG Machine. In *Implementation of Functional Languages, IFL'01. Selected Papers. LNCS 2312*, pages 88–104. Springer-Verlag, 2002.
- [3] S. P. Jones and J. Launchbury. Unboxed values as first class citizens in a non-strict functional language. *Conference on Functional Programming Languages and Computer Architecture FPCA'91, LNCS 523*, September 1991.
- [4] J. Launchbury. A Natural Semantics for Lazy Evaluation. In *Proc. Conference on Principles of Programming Languages, POPL'93*. ACM, 1993.
- [5] J. Mountjoy. The Spineless Tagless G-machine, Naturally. In *Third International Conference on Functional Programming, ICFP'98, Baltimore*. ACM Press, 1998.

- [6] S. L. Peyton Jones. Implementing Lazy Functional Languages on Stock Hardware: the Spineless Tagless G-machine, Version 2.5. *Journal of Functional Programming*, 2(2):127–202, April 1992.
- [7] S. L. Peyton Jones, C. V. Hall, K. Hammond, W. D. Partain, and P. L. Wadler. The Glasgow Haskell Compiler: A Technical Overview. In *Joint Framework for Inf. Technology, Keele*, pages 249–257, 1993.
- [8] S. L. Peyton Jones and J. Hughes, editors. *Report on the Programming Language Haskell 98*. URL <http://www.haskell.org>, February 1999.
- [9] S. L. Peyton Jones, S. Marlow, and A. Reid. The STG Runtime System (revised). <http://www.haskell.org/ghc/docs>, 1999.
- [10] P. Sestoft. Deriving a Lazy Abstract Machine. *Journal of Functional Programming*, 7(3):231–264, May 1997.