

Towards Small-step Compilation Schemas for SOS

Ferdinand Vesely

Department of Computer Science, Swansea University, Swansea SA2 8PP, UK
csfvesely@swansea.ac.uk

Abstract

We present work in progress on a method of compiling programs based on SOS specifications. The idea is to compile programs using SOS rules by translation into labelled blocks with explicit exit points, which implement a valid computation in the LTS of the program. Under this approach, a correct compiler can be constructed in a systematic way, based on an SOS specification.

1 Introduction and Background

Small-step SOS is a popular framework for specifying semantics of programming and specification languages. A collection of SOS rules together define the transitions of a labelled transition system (LTS). For programming languages, the states of an LTS usually contain a program term along with auxiliary entities (stores, environments) and labels may contain emitted signals or output streams. Table 1 contains a small example specification. Under such a specification, each program is represented by a concrete LTS. We present a compilation method which can be understood as the translation of the LTS to a corresponding control-flow graph (CFG). The nodes of the CFG are sequences of instructions with behaviour that should be equivalent to the states in the LTS.

Atomic Blocks Our method produces a collection of labelled *atomic blocks* (AB) containing instructions for a virtual machine. Each AB corresponds to a state in the LTS of the program, and is essentially a *basic block*: a sequence of instructions with one entry (at the beginning) and one exit point (at the end) [1]. However, we relax the second condition and allow multiple exit points at the end of the block, while requiring that an AB executes atomically as a single unit.

Target Machine Language We are targeting a simple register machine, with an unlimited supply of temporaries (registers). In this regard it is similar to LLVM [3], which we intend to

$\frac{\rho \vdash s_1 \xrightarrow{l} s'_1}{\rho \vdash \mathbf{let}(i, s_1, s_2) \xrightarrow{l} \mathbf{let}(i, s'_1, s_2)} \quad (1)$	$\frac{\mathbf{Value } v_1 \quad \rho[i \mapsto v_1] \vdash s_2 \xrightarrow{l} s'_2}{\rho \vdash \mathbf{let}(i, v_1, s_2) \xrightarrow{l} \mathbf{let}(i, v_1, s'_2)} \quad (2)$
$\frac{\mathbf{Value } v_2}{\rho \vdash \mathbf{let}(i, v_1, v_2) \xrightarrow{\tau} v_2} \quad (3)$	$\frac{\rho(i) = v}{\rho \vdash \mathbf{bound}(i) \xrightarrow{\tau} v} \quad (4)$
$\frac{\rho \vdash s \xrightarrow{l} s'}{\rho \vdash \mathbf{print}(s) \xrightarrow{l} \mathbf{print}(s')} \quad (5)$	$\frac{\mathbf{Value } v}{\rho \vdash \mathbf{print}(v) \xrightarrow{\mathbf{out } v} \mathbf{skip}} \quad (6)$

Table 1: Example language specification. ‘Value s ’ asserts that ‘ $\langle \rho, s \rangle$ ’ is a value (terminal) state for any ρ . As usual, $\rho \vdash s \xrightarrow{l} s'$ is a shorthand for $\langle \rho, s \rangle \xrightarrow{l} \langle \rho, s' \rangle$.

use as the ultimate target. There is no program counter, instead the program is stored as a set of labelled code blocks β . Each block is (just) a sequence of instructions ‘ $\iota_1 \cdot \iota_2 \cdot \dots$ ’. The basic control-flow instructions are **halt** for stopping the machine immediately, and **jump** for unconditionally jumping to a labelled AB.

$$\frac{}{\beta \vdash \mathbf{halt} \cdot t \xrightarrow{\tau} \mathbf{halt}} \quad \frac{\langle l, t_2 \rangle \in \beta}{\beta \vdash \mathbf{jump} \ l \cdot t_1 \xrightarrow{\tau} t_2}$$

Further instructions will be mentioned in our translation example in the next section.

2 A Small-step Compilation Schema

Let’s take a simple construct like **print**. If there is a sequence of n transitions starting from term t , then the computation starting from **print**(t) will look as follows:

$$\begin{array}{ccccccccccc} t & \xrightarrow{L_1} & t_1 & \xrightarrow{L_2} & \dots & \xrightarrow{L_{n-1}} & t_{n-1} & \xrightarrow{L_n} & v & & \\ \mathbf{print}(t) & \xrightarrow{L_1} & \mathbf{print}(t_1) & \xrightarrow{L_2} & \dots & \xrightarrow{L_{n-1}} & \mathbf{print}(t_{n-1}) & \xrightarrow{L_n} & \mathbf{print}(v) & \xrightarrow{\{\mathbf{out}=v, \dots\}} & \mathbf{skip} \end{array}$$

The ABs for **print**(t) should each corresponds to a term (state) in the lower part of the above sequence. We construct a *translator* which will generate code blocks that implement the steps of the construct. A translator for construct f , \mathbf{tr}_f , is a structure of operations **next**, **code**, and **label**. A *translator state* is constructed by applying \mathbf{tr}_f to translator states for arguments of f . We also write $\llbracket f(t_1, \dots, t_n) \rrbracket$ for $\mathbf{tr}_f(\llbracket t_1 \rrbracket, \dots, \llbracket t_n \rrbracket)$, where n is the arity of f . For a translator tr , **next** tr is the next translator state, **code** tr is a code block corresponding to the current state, and **label** tr assigns a name to the state. The name can be used as a label for the atomic block or as the name of a temporary holding the computed value. The value of **next** tr can be **none** if the current state is final. In that case, the instructions in **code** tr must store a value in the named temporary label tr . The translator for **print**, $\mathbf{tr}_{\mathbf{print}}$, can be defined as follows:

$$\mathbf{code}[\llbracket \mathbf{print}(t) \rrbracket] = \begin{cases} \mathbf{code}[\llbracket t \rrbracket] & \text{if } \mathbf{next}[\llbracket t \rrbracket] \neq \mathbf{none} \\ \mathbf{code}[\llbracket t \rrbracket] \cdot \mathbf{out} \ \mathit{temp} & \text{otherwise, } \mathit{temp} = \mathbf{label}[\llbracket t \rrbracket] \end{cases} \quad (7)$$

$$\mathbf{next}[\llbracket \mathbf{print}(t) \rrbracket] = \begin{cases} \mathbf{tr}_{\mathbf{print}}(\mathbf{next}[\llbracket t \rrbracket]) & \text{if } \mathbf{next}[\llbracket t \rrbracket] \neq \mathbf{none} \\ \llbracket \mathbf{skip} \rrbracket & \text{otherwise} \end{cases} \quad (8)$$

A translator for a value v just has to put (a representation of) the value to a temporary, for which we introduce an instruction **ldval**.

$$\mathbf{code}[\llbracket v \rrbracket] = \mathbf{ldval} \ \mathit{temp} \ v \quad \mathbf{next}[\llbracket v \rrbracket] = \mathbf{none} \quad \mathbf{label}[\llbracket v \rrbracket] = \mathit{temp} \quad (9)$$

(where temp is a fresh temporary name)

The role of a top-level translator $\mathbf{tr}_{\mathbf{top}}$ is to take the code block for each step and turn it into an atomic block by appending an explicit exit point (**jump** or **halt**):

$$\mathbf{code}[\llbracket t \rrbracket]_{\mathbf{top}} = \begin{cases} \mathbf{code}[\llbracket t \rrbracket] \cdot \mathbf{jump} \ l & \text{if } \mathbf{next}[\llbracket t \rrbracket] \neq \mathbf{none} \text{ and } l = \mathbf{label}(\mathbf{next}[\llbracket t \rrbracket]) \\ \mathbf{code}[\llbracket t \rrbracket] \cdot \mathbf{halt} & \text{otherwise} \end{cases} \quad (10)$$

$$\mathbf{next}[\llbracket t \rrbracket]_{\mathbf{top}} = \mathbf{next}[\llbracket t \rrbracket] \quad \mathbf{label}[\llbracket t \rrbracket]_{\mathbf{top}} = \mathbf{label}[\llbracket t \rrbracket] \quad (11)$$

The main compilation function just collects all the atomic blocks for a term, and returns them together with the initial block label:

$$\mathcal{C}(s) = \{\langle \text{label}[[t]]_{\text{top}}, \text{fold}_{\text{tr}}[[t]]_{\text{top}} \rangle\} \quad (12)$$

where

$$\text{fold}_{\text{tr}} tr = \begin{cases} \{\langle \text{label } tr, \text{code } tr \rangle\} \cup \text{fold}_{\text{tr}}(\text{next } tr) & \text{if } tr \neq \text{none} \\ \emptyset & \text{otherwise} \end{cases} \quad (13)$$

As a further illustration, we look at a definition of `code` for **let**. The construct uses an updated context in the premise of the rule in Eq. (2). The resulting code block also has to provide a corresponding context for the sub-block. This context has to be explicitly constructed at the beginning of the premise transition and cleaned up at the end. In this definition we assume a value translator for sets of mappings ‘ $\{i \mapsto v\}$ ’ and machine operations for manipulating environments.

$$\text{code}[[\text{let}(i, t_1, t_2)]] = \begin{cases} \text{code}[[t_1]] & \text{if } \text{next}[[t_1]] \neq \text{none} \\ \text{code } tr_{iv} \cdot \text{push_env } tmp \cdot \text{code}[[t_2]] \cdot \text{pop_env} & \text{if } \text{next}[[t_1]] = \text{none}, \\ & \text{next}[[t_2]] \neq \text{none}, \\ & tr_{iv} = [[\{i \mapsto t_2\}]], \\ & tmp = \text{label}(tr_{iv}) \\ \text{code}[[t_2]] & \text{otherwise} \end{cases} \quad (14)$$

3 Conclusion

We have illustrated a schema for small-step compilation on a few simple programming constructs. For lack of space we didn’t illustrate translations for, e.g., conditional, iterative, or non-deterministic constructs. The approach could be used with a suitable notion of bisimulation: to prove its correctness, to develop a compiler calculation method (following [2]), and to explore (semi-) automatic compiler generation based on SOS rules. To this end, we intend to work with Modular SOS [4], a modular variant of SOS, which places all auxiliary entities into labels of transitions, and the corresponding notions of bisimulation [5]. To deal with inherent inefficiencies (e.g., construction and destruction of contexts in atomic blocks), common optimisation methods, such as peephole optimisation, could be applied to the resulting translations.

References

- [1] A. V. Aho, M. S. Lam, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2006.
- [2] P. Bahr and G. Hutton. Calculating Correct Compilers. To appear in *J. Fun. Program.*, 2015.
- [3] C. Lattner and V. Adve. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. CGO ’04, Washington, DC, USA, 2004. IEEE Computer Society.
- [4] P. D. Mosses. Modular structural operational semantics. *J. Log. Algebr. Program.*, 60-61:195–228, 2004.
- [5] P. D. Mosses and F. Vesely. Weak bisimulation as a congruence in MSOS. In N. Martí-Oliet, P. C. Ölveczky, and C. Talcott, editors, *Logic, Rewriting, and Concurrency*, volume 9200 of *LNCS*, pages 519–538. Springer, 2015.