

Modular Operational Semantic Specification of Transport Triggered Architectures

Jon Mountjoy

Department of Computer Science, University of Amsterdam

Kruislaan 403, 1098 SJ Amsterdam, The Netherlands.

Telephone: + 31 20 525 7463, E-mail: jon@wins.uva.nl

Pieter Hartel

Department of Electronics and Computer Science, University of Southampton

Highfield, Southampton SO17 1BJ, United Kingdom

Henk Corporaal

Department of Electrical Engineering, Delft University of Technology

Mekelweg 4, 2628 CD Delft, The Netherlands

Abstract

The formal specification of hardware at the instruction level is a daunting task. The complexity, size and intricacies of most instruction sets makes this task even more difficult. However, the benefits of such a specification can be quite rewarding: a precise, unambiguous description is provided for each instruction, a basis for proving the correctness of code transformations is made available, and the specification can be animated, providing a simulator. This paper proposes a high level structural operational semantic (S.O.S.) specification for the class of transport triggered architectures. These architectures are simple, powerful, flexible and modular and can exploit very fine grained parallelism. The S.O.S. is novel in that it follows the structure of the architecture, and by doing so inherits the modularity of the architecture.

Keywords

Operational semantics; architecture specification; model development

1 INTRODUCTION

The precise definition of programming languages is important; ambiguities in programming language definitions were rife before the introduction of formal techniques of semantic specification. Likewise, the specification of compilers is becoming more rigorous, relying on the correct incremental transformation of programs written in a high level language to some low level machine language.

There is also much work on verifying hardware and formal hardware description languages. However, the formal specification and verification of low level machine languages, lying somewhere between architecture and language, is lagging behind.

It is perhaps due to the size and intricacies of most processors that the formalisation is so difficult. For instance, the specification of a large portion of the MC68020 instruction set described in Boyer & Yu (1996) took approximately 80 pages of text. This difficulty is unfortunate, as a manageable specification has many uses:

1. The specification can serve as a rigorous *hardware description* of the architecture. The formal meaning given to the machine code in this way eliminates ambiguities. This can then serve as documentation for VLSI designers and compiler writers.
2. Having a formal semantics helps us to prove code *transformations* or optimizations correct, or to clarify the conditions under which such transformations can be applied.
3. Animating (implementing) the semantics will yield a prototype of the machine architecture and instruction set. This tool can then serve as the basis for *verification and validation* of hardware and compiler. From another perspective, the hardware can be used to validate the semantic model.

In general, an abstract semantic description of a low level machine architecture may be non-trivial due to the complexity of the architecture. The class of *transport triggered architectures* (TTAs) (Corporaal & van der Arend 1993) were introduced as modular, parallel, application-specific processors. This extensible, VLIW-like architecture, has parallel data moves as the only operations. All that is visible at the instruction level are a set of registers whose values can be moved to other registers. The architecture shows much promise (Corporaal & Mulder 1991, Corporaal 1995), and although seemingly simple, is powerful and flexible. Several physical implementations exist; see Corporaal & van der Arend (1993) for an example. Since the machine itself is at a low level (data transports are exposed, pipelines are visible, instructions have to be scheduled), compilers for the architecture are by necessity more complex (Fisher, Ellis, Ruttenberg & Nicolau 1984) making the precise documentation of the machine even more important.

The inherent modularity of transport triggered architectures lends itself to formal modelling not easily applied to other processors. In this paper we provide an operational semantic description for a class of TTAs. This is based on structural operational semantics (Plotkin 1981). This semantic framework is simple, readable, and easily implemented. It can also be used to model a fine grained parallel system, which makes it suited to the task at hand.

A requirement of the semantics is that it must be easily extensible. If we wish to prototype various configurations of architectures (which is the case since we are dealing with application-specific processors), we need the semantics to be modular with respect to the architecture: small changes in the

architecture should not mean extensive changes to the semantics. With structural operational semantics, one usually builds the semantics following the structure of the syntax of a language. This paper builds the semantics of the architecture after the structure of the architecture; in effect, the semantics is a synthesis of the essential ingredients of the architecture. We believe this approach to be novel.

Another ingredient of a good semantic description is that it models at the 'correct' level. A model should not be too fine, describing circuit level when the user of the model has no use of this information. Likewise, it should not be too coarse, eliminating information that could be useful. We have sought to find a balance which will result in a model capable of fulfilling points (1–3) above. By necessity, the model is then at a much higher level than those found in the formal methods used in high-level design and verification of VLSI circuits. There are *hardware verification* languages, where one takes a gate-level description of a circuit and derives a model thereof. Some kind of model checking can then be used to ensure that the model is consistent with some other specification. There are also languages used for formal *hardware specification* such as VHDL. These languages typically address design issues at the lower end of the spectrum (close to the hardware). The language Ruby (Jones & Sheeran 1990) allows one to specify circuits in a high-level formal language, and 'calculate' circuits. We hope that our higher level specification can serve as a base for linking the specifications at the various levels. In addition we eventually would like to prove that a model based on the hardware circuits fulfills the formal specification in this paper. Our methodology thus involves analysing the architecture at the particular level of interest and extracting a language design. This language is then given a semantics. It is the aim of this paper to show the use of this methodology and the resulting semantics. We believe that this methodology can also be applied to other architectures which share the TTA architecture hierarchy, such as VLIW processors.

We begin by describing the concepts behind transport triggered architectures in section 2. A typical TTA is given, and this serves as the basis for the semantic description provided in section 3. Section 4 shows how several variations of the architecture can be modelled, and how the semantics is animated. Section 5 concludes this paper and examines related work and extensions.

2 TRANSPORT TRIGGERED ARCHITECTURES

TTAs can be compared to VLIW architectures. In both cases the instructions are horizontally encoded; i.e. each instruction has a number of fields. Whereas fields for VLIWs specify RISC like operations, for TTAs they specify the required data transports. These transports may trigger operations as side effects. Programming transports adds an extra level of control to the code generator, and enables new optimizations; in particular, it allows us to eliminate many superfluous data transports to and from the register files and to reduce the on-chip connectivity. TTAs are fully introduced in Corporaal & van der Arend (1993) and Hoogebrugge (1996).

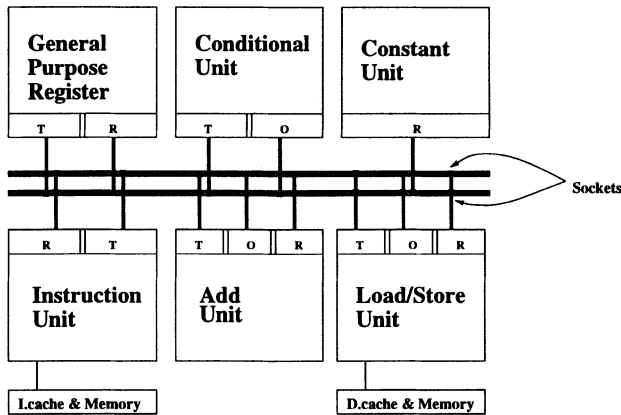


Figure 1 An Abstract Operational view of a TTA with Trigger(T), Operand(O) and Result(R) registers indicated.

A compiler views a TTA as a collection of function units (FUs), register files (RFs), *move buses*, and *sockets*; see figure 1. FUs perform operations, RFs provide general purpose registers for temporary fast accessible storage, the network of move buses performs data transports between the FUs and RFs, and sockets interface FUs and RFs to move buses. Normally, each socket is connected to a different FU input/output or RF port.

The *instruction unit* plays the special rôle of dispatching the operations of the current instruction to the appropriate buses, and of storing the program counter.

To illustrate TTA programming, consider the execution of the following three instructions on a conventional machine:

```
add r1, r2, r3 /* r1 = r2 + r3 */
sub r4, r2, r6 /* r4 = r2 - r6 */
st r4, r1 /* store r4 at address r1 */
```

These operations can be translated into the following two TTA instructions, given a suitable architecture (unit latency) with four buses:

```
r2→ add_o || r3→ add_t || r2→ sub_o || r6→sub_t
add_r→ st_o || sub_r→ st_t || nop || nop
```

Each instruction is composed of four operations and each operation is destined for a single bus. These operations will then be executed in parallel. An increase in the number of buses or FUs increases the amount of potential parallelism. In the first instruction the four operands of the add and subtract operations are moved from general purpose registers to the inputs of the functional units which perform addition and subtraction. The add and subtract

unit then, have two input sockets each. Output sockets, which will hold the result of the operation, are also available. In the second instruction the results of the add and subtract operations are moved from the FUs that performed them to the FU that performs the store operation.

Conditional execution is provided by guarding every move with the result of a conditional unit. Control flow is provided by exposing the program counter as a register of the instruction unit. Writing to it with an address, forces a jump to this address (since the instruction unit dispatches the instruction found at its program counter location).

2.1 A Detailed Look

Figure 1 gives an operational view of a possible transport triggered architecture. It consists of a number of functional and register units and an interconnection network. In this example, we have only provided units for an addition, comparison, general purpose register and load/store. We also assume a fully connected network of two buses. A more general architecture may have additional functional units for multiplication, division and logical operations. Instead of providing a single register, it would more typically provide a register file of several registers. Note that these would be independent and modular extensions: all that has to change is the addition of the functional unit and connections to the buses. More buses can be added, and the interconnection network need not be fully connected. In addition, functional units may be duplicated. For instance, we may have two independent load/store units attached to the buses, and address each by a unique identifying name.

The architecture presented here is sufficiently general to cover all of the important aspects in the modelling. Section 4 addresses extensions. We will adopt the convention of labeling the input and output sockets of the functional units by the unit name and a letter, either 'O', 'T' or 'R'. We will also refer to these as registers. An operand register (O), is used to supply a datum to a functional unit. The trigger register (T), is also used to supply a datum to the functional unit. In addition, it may set a flag enabling the functional unit to start its computation. When the functional unit produces its result, it is usually placed in a result register (R). We will always refer to registers provided by a register unit for temporary storage as general purpose registers to avoid ambiguity.

It is important to note that the only aspects of the functional unit made available to the programmer are the registers and data transports between them. The architecture is thus *transport triggered*, as opposed to the more conventional *operation triggered* architectures where operations are issued and data transports occur as possible side effects of these operations.

Small constants are made available by embedding them in the operation word. A constant unit (seen operationally) then reads the instruction off the bus and makes the embedded constant available on its result register. A regis-

ter unit provides a register. Writing a value to the trigger will make it available for reading from its result register in the next cycle.

Since it is not always possible to fill an instruction completely with useful operations, a NOP operation is provided.

Control Flow and Conditional Execution

Conditional execution is provided by *guarding* every data transport. The guard condition is taken as the output of the conditional functional unit. The output of a conditional unit is a boolean which is not made available as a register. Instead, it is used internally to control squash lines on the various buses. If a bus is squashed, then the data transport taking place on that bus is cancelled. Flow control is made possible by making the program counter visible as a register, (the instruction unit trigger and result registers). This program counter automatically increments every cycle (unless written to). Writing to the instruction unit trigger register (INS^T) forces a jump to the location written.

The following fragment of code demonstrates these ideas. The first instruction simultaneously moves the value 10 to the conditional units operand, and 20 to the trigger, starting the conditional unit comparison. If the result of this comparison is true, then a jump is made to location 7.

$$\begin{array}{ll} 10 \rightarrow CND^O & || \quad 20 \rightarrow CND^T \\ \vdots & \vdots \\ [Check] 7 \rightarrow INS^T & || \quad NOP \end{array}$$

We will use the word “Check” to denote that the move is conditionally guarded on a value being true, or “Always” to denote that the move is not guarded and will always succeed. For brevity, we sometimes omit the annotation if the guard is unconditional. We will also assume that the conditional unit returns true if the operand register is less than or equal to the trigger register.

Pipelined Functional Units

A functional unit may be pipelined. We will assume a semi virtual-time latching scheme which means that pipelines run synchronously to the instruction stream; each time an instruction is issued, the FU pipelines progress a step. In addition, only moves to the trigger register start new operations. Note that in this scheme a result register can be intentionally overwritten without ever being read. Pipelining and its implications for the semantics is discussed in section 4.2.

2.2 Deriving a Syntax

Our methodology of assigning a structural operational semantics to the architecture is based on extracting details from the architecture at an appropriate

$Operation \ni O$::= NOP	$Guard \ni G$::= Check
	$G : S : T$		Always
$Source \ni S$::= CNST-c	$Target \ni T$::= ADD ^O
	ADD ^R		ADD ^T
	INS ^R		INS ^T
	REG ^R		REG ^T
	LD ^R		CND ^O
			CND ^T
$Constants \ni c$::= 0 1 ... 255		ST ^O
$Instruction \ni I$::= (O O)		ST ^T
$Program \ni P$::= (I, I, ..., I)		LD ^T

Figure 2 Abstract Syntax of the Architectural Language

level of interest. Analysing the architecture, we can immediately see a hierarchy involved in the execution of an instruction: an instruction is composed of operations, operations are guarded moves, and moves are from source to target registers. This tidy separation is unfortunately more difficult to find in conventional architectures.

We will ultimately assign meanings to programs written in the machine language, and so we synthesize a syntax for this language keeping the hierarchy of the architecture in mind. A program then, is represented as a sequence of instructions, an instruction as two operations (recall that the architecture we are focusing on illustrated in figure 1 has two buses), an operation being either a NOP or a guarded move and a move by a source and target register specification. Figure 2 illustrates this.

It is a strength of the architecture that it allows this hierarchical division. Different instantiations of the architecture do not change this fundamental structure, but only the components. Thus adding a bus or FU is mirrored easily in the syntax.

Note that the load/store unit is capable of two actions (loading and storing); the physical architecture uses the opcode in the operation to distinguish between them. We have separated them at the syntax (and thus semantics) level as there is no reason to complicate the system by modelling this explicitly. The point to be noted is that they are the same physical resource, and so cannot both be used at the same time. One can prevent this by using a *static semantics*, as explained in subsection 3.4.

We are now in a position to give a formal semantics to the language.

3 SEMANTIC DESCRIPTION

The TTA architecture is parallel (many moves can be executed at once) and modular (it can be extended by adding functional units or buses). The execution of an instruction consists of the parallel execution of the operations of that instruction. The abstract syntax of the architecture reflects these issues. After the execution of an instruction the functional units update their state (if

there is information in the pipeline, or if a result register has been written to). Initially we will make certain assumptions about the architecture to simplify the description, and address the extensions required to the semantics to cope with the full architecture in a later section.

A structural operational semantics specifies transition relations on various configurations. A configuration is usually a tuple, consisting of the instruction (or part thereof) to be executed and a state; sometimes it is just a value. The state is a representation of the state of the registers and memory locations.

3.1 Sources

We begin by defining the binary transition relation which concerns configurations containing source registers:

$$\xrightarrow{s} \in (\textit{Source} \times \textit{State}) \times \textit{Value}$$

Here *Source* is the set of source registers taken from figure 2. That is, we define a transition relation \xrightarrow{s} between a source register and state pair, and a value. This relation, together with a relation on the target register, is used to build up the meaning of an operation.

We can now begin to define the rules for the various source constructs. The meaning of a constant c is the value c of the constant. The constant unit therefore acts as the identity on the value. This is exactly what happens in the architecture where the constant is read from the instruction word. We write:

$$\frac{S = \text{CNST-}c}{\langle S, s \rangle \xrightarrow{s} c}$$

This can be read as: If the source register is a CNST- c (the premise), then the meaning of the configuration $\langle S, s \rangle$ is the value c (the conclusion). We always write the premise above the rule, and conclusion beneath it. Section 4.1 explores an alternative approach to constants.

If we are dealing with the source register of the add unit, then we define the meaning of the result register as the value of the result register in the state:

$$\frac{S = \text{ADD}^R}{\langle S, s \rangle \xrightarrow{s} s[\text{ADD}^R]}$$

We use the notation $s[r]$ to denote the value of r in state s . The value returned above then is the value of the result register of the add unit. We will later return to defining the actual components in the state. For now, we can consider it as a mapping between registers and values.

All other source register moves are modelled similarly. For instance, we define the rule for the load unit as:

$$\frac{S = \text{LD}^R}{\langle S, s \rangle \xrightarrow{s} s[\text{LD}^R]}$$

3.2 Targets

The transition relation for target operands is given by:

$$\overset{t}{\rightarrow} \in (\textit{Target} \times \textit{Value} \times \textit{State}) \times \textit{Substitution}$$

Target is the set of target registers. We will define the rules for the add unit.

We first need to define a *substitution*. Intuitively, a substitution holds a list of registers and values to which these registers are to be mapped. Later, we will see how these substitutions are used to update the state of the machine. For now, they can be thought of as reminders that the state is to be updated. We will write ϕ for the empty substitution. A more formal motivation is given in subsection 3.4.

The meaning of a target register and value is given by a substitution.

$$\frac{T = \text{ADD}^O}{\langle T, v, s \rangle \overset{t}{\rightarrow} [(\text{ADD}^O, v)]}$$

Thus in the above rule, the substitution indicates that the operand register should be mapped to the value v . The use of the substitution will become clear once we address the semantics for instructions.

Recall that we are dealing with a semi virtual-time architecture in which pipelines only progress synchronously with virtual time, and if the unit has been ‘triggered’. Moving to an operand register results in a single substitution. Moving to a trigger register results in the additional setting of a flag, which is meant to indicate that the unit has been triggered. A similar mechanism occurs in the hardware where flags are used at the pipeline stages to indicate whether valid data has been latched.

$$\frac{T = \text{ADD}^T}{\langle T, v, s \rangle \overset{t}{\rightarrow} [(\text{ADD}^T, v), (\text{ADD}^F, \textit{True})]}$$

Note that the register ADD^F is something introduced to model the behaviour of the architecture. It is also present in the implementation of the architecture where it is not visible to the user.

The other target transitions are modelled in exactly the same manner.

3.3 Operations

We can now model an operation. This relation is ultimately defined in terms of the above two transitions: the source register transition and the target register transition. The operation transition relation, $\overset{o}{\rightarrow}$, is defined as:

$$\overset{o}{\rightarrow} \in (\textit{Operation} \times \textit{State}) \times \textit{Substitution}$$

The semantic description is formed on the structure of an operation as defined in figure 2. We begin with the NOP. Intuitively, the execution of the NOP should leave the state unchanged. We therefore model the meaning of

this operation, in a state s , as the empty substitution (when this is later applied to the state, the state will remain unchanged).

$$\frac{O = \text{NOP}}{\langle O, s \rangle \xrightarrow{\circ} \phi}$$

The other forms of operations involve guards. If the guard is “Always”, then we will always execute the move. This involves a transition for the source and the target register, and we use the relations defined in the previous two sections. Ultimately, a substitution is yielded.

$$\frac{O = \text{Always} : S : T \quad \langle S, s \rangle \xrightarrow{s} v \quad \langle T, v, s \rangle \xrightarrow{t} nv}{\langle O, s \rangle \xrightarrow{\circ} nv} \quad (1)$$

Thus a transition ($\xrightarrow{\circ}$) is made to nv in the case that the source register produces a value v , and the target transition using this value produces a substitution nv . The rule can be read as: if $\langle S, s \rangle \xrightarrow{s} v$ and $\langle T, v, s \rangle \xrightarrow{t} nv$ then $\langle \text{Always} : S : T, s \rangle \xrightarrow{\circ} nv$.

If the guard is “Check”, then the execution of the move only takes place if the value in the conditional unit is true; otherwise nothing happens. The false case is defined by:

$$\frac{O = \text{Check} : S : T}{\langle O, s \rangle \xrightarrow{\circ} \phi}, \text{ if } \text{Guard}(s) = \text{False}$$

It is easy to see the following:

$$\forall s \in \text{State}. \text{Guard}(s) = \text{False} \Rightarrow \langle \text{Check} : S : T, s \rangle \xrightarrow{\circ} \phi \text{ and } \langle \text{NOP}, s \rangle \xrightarrow{\circ} \phi$$

That is, a transition of an operation with a failed guard and a NOP yield the same (empty) substitution.

The *Guard* function used in the above rules takes the supplied state and yields the result of the conditional units comparison. It is of type $\text{Guard} :: \text{State} \rightarrow \text{Bool}$ and defined as follows:

$$\text{Guard}(s) = \begin{cases} \text{True} & , \text{ if } (s[\text{CND}^R] = 1) \\ \text{False} & , \text{ otherwise} \end{cases}$$

We write $s[\text{CND}^R]$ for the value of the CND^R register in the state s .

A guard which succeeds has the same semantics as a non-conditional move with the same source and target registers, as can be seen by comparing rules 1 and 2.

$$\frac{O = \text{Check} : S : T \quad \langle S, s \rangle \xrightarrow{s} v \quad \langle T, v, s \rangle \xrightarrow{t} nv}{\langle O, s \rangle \xrightarrow{\circ} nv}, \text{ if } \text{Guard}(s) = \text{True} \quad (2)$$

Note that in the definition of the *Guard* function we use the result register of the conditional unit, a register which cannot be explicitly read by the user. Instead of modelling complex squash lines and timings, the hardware has been abstracted and replaced by a ‘pseudo’ register.

3.4 Substitutions and Instructions

The previous sections used substitutions to represent eventual changes to the state. We now motivate the use of substitutions instead of the more conventional techniques discussed below. This then allows us to specify the semantics of instructions.

A key feature of the semantics is that it describes a parallel architecture. It is natural to expect that the semantics also be parallel, but at the same time deterministic. This is a notorious problem, as we have to merge two states which change in parallel.

One possible solution is to model the parallel execution of two operations in some state s by providing rules allowing the interleaved execution of either operation. This would be written as:

$$\frac{\langle O_1, s \rangle \rightarrow \langle O'_1, s' \rangle}{\langle O_1 \parallel O_2, s \rangle \rightarrow \langle O'_1 \parallel O_2, s' \rangle} \quad \frac{\langle O_2, s \rangle \rightarrow \langle O'_2, s' \rangle}{\langle O_1 \parallel O_2, s \rangle \rightarrow \langle O_1 \parallel O'_2, s' \rangle}$$

Another solution is to assume the existence of a clever merge operator, written here as $+$, which will only merge those portions of the state that have been changed:

$$\frac{\langle O_1, s \rangle \rightarrow s_1 \quad \langle O_2, s \rangle \rightarrow s_2}{\langle O_1 \parallel O_2, s \rangle \rightarrow s_1 + s_2}$$

We can, however, do better than this. Since we know that each operation will not use the same resources (we have the so called disjointness requirement of Plotkin (1982)), and that the changes to the state that each operation induces are mutually exclusive (see below), we model operations as being substitutions. That is, (name,value) pairs that can later be applied to the state. We thus model the parallel execution as:

$$\frac{\langle O_1, s \rangle \rightarrow nv_1 \quad \langle O_2, s \rangle \rightarrow nv_2}{\langle O_1 \parallel O_2, s \rangle \rightarrow (s \uplus nv_1) \uplus nv_2}$$

The union operator, \uplus , used above takes as arguments a state and substitution, and yields a new updated state. Its type is $\uplus :: State \rightarrow Substitution \rightarrow State$, and its behaviour is such that for a register r :

$$(s \uplus p)[r] = \begin{cases} (p[r]) & \text{if } (r \in p) \\ (s[r]) & \text{otherwise} \end{cases}$$

It is easy to see that if nv_1 and nv_2 are disjoint, then $((s \uplus nv_1) \uplus nv_2)[r] = ((s \uplus nv_2) \uplus nv_1)[r]$.

Armed with this technique, we can now proceed to give a semantics to instructions. The transition relation for instructions is given by:

$$\stackrel{i}{\rightarrow} \in (Instruction \times State) \times State$$

In defining the semantics of an instruction, we are faced with a choice. The TTA is parallel. However, all “well scheduled” operations that execute in an instruction are independent and will use separate resources. Indeed, this is part of the task of the scheduler. We can then either assume that we are given

code which is scheduled correctly, or try and explicitly model incorrect code. In the former scenario, we have:

$$\frac{I = (O_1 \parallel O_2) \quad \langle O_1, s \rangle \xrightarrow{\circ} nv_1 \quad \langle O_2, s \rangle \xrightarrow{\circ} nv_2}{\langle I, s \rangle \xrightarrow{i} (s \uplus nv_1) \uplus nv_2} \quad (3)$$

The meaning of an instruction is the final state achieved by updating the initial state with the substitutions yielded by the individual operations. Note that, because we have assumed nv_1 and nv_2 to be disjoint, the final state update can take place in an arbitrary order. The task of ensuring that operations yield mutually exclusive substitutions can be given to a *static semantics*. This semantics can include rules checking the well-formedness of the instructions. In effect, it checks that the code is conflict free, a property that would be guaranteed if a scheduler were used. Among other things it can ensure that the maximum size of constants is not exceeded, and perform resource checks. If the architecture being modelled never had a fully-connected network, then the check on which moves are valid can also be incorporated in the static semantics. Implementing the semantics with explicit checking would involve replacing rule 3 by rules which check that the substitutions nv_1 and nv_2 are disjoint, and yield a special terminal configuration, \perp , indicating failure. In terms of the hardware, two operations writing to the same register yield an undefined result. For the case of instructions, this would be:

$$\frac{I = (O_1 \parallel O_2) \quad \langle O_1, s \rangle \xrightarrow{\circ} nv_1 \quad \langle O_2, s \rangle \xrightarrow{\circ} nv_2}{\langle I, s \rangle \xrightarrow{i} \perp} \quad , \text{ if } \neg \text{Mutex}(nv_1, nv_2) \quad (3\text{-a})$$

$$\frac{I = (O_1 \parallel O_2) \quad \langle O_1, s \rangle \xrightarrow{\circ} nv_1 \quad \langle O_2, s \rangle \xrightarrow{\circ} nv_2}{\langle I, s \rangle \xrightarrow{i} (s \uplus nv_1) \uplus nv_2} \quad , \text{ if } \text{Mutex}(nv_1, nv_2) \quad (3\text{-b})$$

Adopting this form of semantics allows one to reason about the machine at a slightly lower level. For instance, one can now prove that the result of the following instruction is always undefined:

$$\text{Always} : \text{CNST-23} : \text{ADD}^T \parallel \text{Always} : \text{CNST-5} : \text{ADD}^T$$

As a concrete example of the use of the semantics, we prove that the result of executing the following instruction is only a substitution changing the value of the operand register for the add unit to 10:

$$\text{NOP} \parallel \text{Always} : \text{CNST-10} : \text{ADD}^O$$

We do this by constructing a *derivation tree*, applying the semantic rules:

$$\frac{\frac{\langle \text{CNST-10}, s \rangle \xrightarrow{s} 10 \quad \langle \text{ADD}^O, 10, s \rangle \xrightarrow{i} [(\text{ADD}^O, 10)]}{\langle \text{NOP}, s \rangle \xrightarrow{\circ} \phi} \quad \langle \text{Always} : \text{CNST-10} : \text{ADD}^O, s \rangle \xrightarrow{\circ} [(\text{ADD}^O, 10)]}{\langle \text{NOP} \parallel (\text{Always} : \text{CNST-10} : \text{ADD}^O), s \rangle \xrightarrow{i} (s \uplus \phi) \uplus [(\text{ADD}^O, 10)]}$$

We have thus used the semantics to prove that the above instruction only

introduces one change to the state; the operand of the add unit becomes a 10 (since, by the empty substitution, $(s \uplus \phi) \uplus [(ADD^O, 10)] = s \uplus [(ADD^O, 10)]$).

Similarly, we can derive:

$$\langle \text{Always} : \text{CNST-10} : \text{ADD}^O \parallel \text{NOP}, s \rangle \xrightarrow{i} s \uplus [(ADD^O, 10)]$$

This shows that these two instructions are equivalent.

We can go further, and prove that for all operations O ,

$$\langle O \parallel \text{NOP} \rangle \xrightarrow{i} s' \text{ and } \langle \text{NOP} \parallel O \rangle \xrightarrow{i} s''$$

and $s' = s''$. We prove this by showing that the resulting states are equal. For the first instruction we have:

$$\frac{\dots}{\langle O, s \rangle \xrightarrow{o} sv \quad \langle \text{NOP}, s \rangle \xrightarrow{o} \phi} \\ \langle O \parallel \text{NOP}, s \rangle \xrightarrow{i} (s \uplus sv) \uplus \phi$$

For the second, in the same state s , we have:

$$\frac{\dots}{\langle \text{NOP}, s \rangle \xrightarrow{o} \phi \quad \langle O, s \rangle \xrightarrow{o} sv} \\ \langle \text{NOP} \parallel O, s \rangle \xrightarrow{i} (s \uplus \phi) \uplus sv$$

By a simple calculation we have $(s \uplus sv) \uplus \phi = s \uplus sv$ and $(s \uplus \phi) \uplus sv = s \uplus sv$, showing that the final states are indeed equal.

3.5 Program

Extending the semantics of an instruction in a natural way, we can provide a semantics for a program. The \xrightarrow{P} relation makes use of an environment, \mathcal{C} , which stores the program. For clarity, we have written the program counter as PC instead of the more explicit $s[\text{INS}^R]$. The transition relation for programs is given by:

$$\begin{array}{c} \xrightarrow{P} \in (\text{State} \times \text{State}) \\ \hline \mathcal{C}[PC] = (O_1 \parallel O_2) \quad s' = s \uplus [(INS^T, PC + 1)] \quad \langle O_1 \parallel O_2, s' \rangle \xrightarrow{i} s'' \quad s'' \xrightarrow{u} s''' \\ \mathcal{C} \vdash s \xrightarrow{P} s''' \\ \text{if } \neg \text{Halt}(PC) \end{array} \quad (4)$$

We can read this as: If in environment \mathcal{C} the program counter of the current state is not a terminating one, then we can look up the instruction at the program counter. A new state, s' , is created which has the program counter updated. This is in accordance to the physical architecture which updates the program counter automatically. Using the new state s' , a transition is made on the instruction, and the resulting state is updated via an update transition (\xrightarrow{u}). Recall that we have been modelling a semi-virtual time system: all updates to the state of the machine due to pipeline changes etc. are performed

before the execution of the next instruction. This is exactly what we have modelled above. In fact, each transition of \xrightarrow{P} can be thought of as a virtual time cycle. Finally, we define some halting condition on the program counter which will indicate when the program has finished execution. An initial state is assumed which maps the program counter to the first location and every other register to zero.

With a suitable update function (one with a latency of 1 for the add unit) we could prove (by constructing a derivation tree) that the final value of the result register in the add unit after executing the following code will be 10:

```

1 : (Always : CNST-5 : ADDO) || (Always : CNST-15 : ADDT)
2 : (Always : CNST-5 : ADDT) || NOP
    
```

3.6 Specifying Updates

Finally, we need to specify the update transition. Since the functional units are independent, an update on the state entails an update on all the components of the architecture (all the FUs and RUs). We will assume that we can split all of the components of the state. That is, all state describing the register unit will be called *REG*, etc. Using this, we can specify the update of the state of the entire machine as the separate modular updates of its components:

$$\frac{s = (INS, \dots, CND) \quad INS \xrightarrow{u^i} INS' \dots CND \xrightarrow{u^c} CND' \quad s' = (INS', \dots, CND')}{s \xrightarrow{u} s'}$$

Since we are using a semi-virtual time pipelining scheme, state updates need only be done at one time in the execution of each instruction, namely before the next instruction issue. We begin by illustrating the specification of the register unit. The triggering of a register unit moves the trigger value to the result value. If the unit has not been triggered, then we do nothing. We use the notation $REG = [(REG^F, False), \dots]$ to indicate the requirement that the state component of the general purpose register maps the flag to *False*.

$$\frac{REG = [(REG^F, False), \dots]}{REG \xrightarrow{u^r} REG} \quad \frac{REG = [(REG^T, v_t), (REG^F, True), \dots]}{REG \xrightarrow{u^r} REG \uplus [(REG^R, v_t), (REG^F, False)]}$$

This can be read as follows. The update of a state *REG*, where the state contains a mapping of REG^F to *False*, is just the state itself. There is no change in the state of the register functional unit if it has not been triggered. If however, the state has the flag register set, then we must yield a new state where this register is reset, and wherein the result register has the value of the trigger register.

In modelling the load/store unit, we need to keep the state of the memory. Instead of dragging an environment around holding this information, we follow Necula & Lee (1996) and consider one register, LD^{MEM} , as special: if

This completes the changes necessary to model an architecture with long constants. As can be seen, only a change in the syntax and two additions to the semantics were necessary, highlighting the benefit of the modular approach.

The new architecture semantics could be compared to the old in various ways. One way is to prove that the architectures are equivalent in some sense. For example, that we could simulate moving a long constant to a register by adding two small constants together in the register. Another comparison would be to look at generated code and compare the efficiency of the two strategies. Unfortunately, the scheduler plays a large rôle in this area, and one cannot provide hard evidence independent of any scheduler. However, an animation of the architectures would provide a test bed for these experiments.

4.2 A Pipelined Add Unit

We examine the consequences of adding a single pipeline stage to the add unit. The only change that has to be made to the semantic description is at the update function for the add unit. We introduce two new pseudo registers, ADD_2^F and ADD_2^R . Intuitively, after these are triggered the addition will take place, and the result will not be placed in the result register but in register ADD_2^R . The associated flag will also be set to indicate that data is sitting in the pipeline stage. The pipeline stages progress with virtual time, and so every update, if a pipeline stage has data, will move the data to the output register. The real architecture keeps a valid bit associated with each pipeline stage, similar to what we do here. Again, the pseudo registers are not visible to the user of the architecture, but similar mechanisms are implemented in the hardware.

The update for the modified functional unit can then be described by four rules. If there is no data in the pipeline, and the unit has not been triggered, then an update is an identity on the state:

$$\frac{ADD = [(ADD^F, False), (ADD_2^F, False), \dots]}{ADD \xrightarrow{u^a} ADD}$$

If, however, the pipeline contains data, then this is fed through to the result register:

$$\frac{ADD = [(ADD^F, False), (ADD_2^R, a2r), (ADD_2^F, True), \dots]}{ADD \xrightarrow{u^a} ADD[(ADD_2^F, False), (ADD^R, a2r)]}$$

The other two rules are similar. These localized changes are all that is needed to model pipelining.

The above technique for modelling pipelines has the disadvantage that we need a number of rules exponential in the depth of the pipeline. This can be avoided by introducing a list of operand/flag values which is cycled on every update. Indeed, this is how the semantics has been animated.

Other pipelining schemes, such as true virtual time, are more easily modelled.

4.3 Locking

If a memory unit uses locking when a cache miss occurs, then the entire machine will stall until the datum is available (a global lock would be in effect). Irrespective of whether a cache miss occurred or not, the machine should still behave in the same way (of course, the machine will take longer to produce the result: a non-functional aspect). Our model only captures virtual time aspects. Modelling the cycle impact of locking could be handled by having a trigger relation returning not only a substitution, but also a cycle count. This cycle count would be greater than zero in the event of some modelled cache miss.

The semantic framework presented here does not model non-functional aspects very well. Indeed, we have purposefully abstracted away from modelling control signals, squash lines and locking. A lower-level semantics of the VLSI circuitry would take this into account. However, animating the semantics provides us with one way of retrieving some non-functional behaviour.

4.4 Determinism and Animation

The above operational semantics can be proved to be *deterministic* by structural induction. By deterministic, we mean that for all choices of I, s, s' and s'' we have:

$$\langle I, s \rangle \xrightarrow{i} s' \ \& \ \langle I, s \rangle \xrightarrow{i} s'' \Rightarrow s' = s''$$

That is, given an instruction and a state we can unambiguously determine the following state. As this is the case, we can capture all of the semantics rules by functions which can easily be implemented. The extended operational semantics given above has been implemented in the functional language Haskell (Peterson et al. 1996), and provides a useful prototype of the architecture. By ‘running’ the semantics we can get an idea of what state changes occur after each virtual cycle and update. We can also gather statistics, and some non-functional aspects such as the number of virtual time cycles.

Here is an excerpt from the code implementing the transition rule on an operation with an ‘Always’ guard (Rule 1), illustrating the proximity of the semantics to the animation:

```
> transition (Oper ALWAYS src trg) s = let v = tran_s src s
>                                       nv = tran_t trg v s
>                                       in nv
```

In addition, techniques are available which can adapt the operational semantics to create reasonably efficient abstract machines (Hannan 1994). The animated semantics can also help provide verification of the compiler, hardware and specification.

(CNST-10 : ADD ^O NOP)	(CNST-10 : ADD ^O CNST-20 : ADD ^T)
(CNST-20 : ADD ^T NOP)	(ADD ^R : ADD ^O CNST-30 : ADD ^T)
(ADD ^R : REG ^T NOP)	(ADD ^R : REG ^T NOP)
(REG ^R : ADD ^O NOP)	
(CNST-30 : ADD ^T NOP)	
(ADD ^R : REG ^T NOP)	
(a)	(b)

Figure 3 (a) Sequential Unoptimized Code (b) Scheduled Optimized Code

4.5 Proving correctness of code transformations

A compiler for a VLIW machine will at some time in the compile cycle produce sequential code which is later optimized and scheduled into parallel code. One optimization which can be applied to TTAs is the removal of redundant register file writes. Figure 3(a) illustrates a sequential piece of code which adds two numbers, writes the result to a register file and then adds 30 to this value writing the final result. We assume a latency of one, and place NOPs to emphasize that the code is unscheduled. We will call this code \mathcal{C}^u .

Figure 3(b) illustrates the code after scheduling and the removing of the redundant register file write. We will call this code \mathcal{C}^o .

These two sequences of code can be shown to be equivalent in the sense that they produce the same final state given the same initial state, modulo the value of the program counter. The proof involves the construction of a number of derivation trees. For each of the sequences of code \mathcal{C}^o and \mathcal{C}^u a number of program transition steps can be derived. To prove the first step in \mathcal{C}^u for instance, we would construct the derivation tree for $\mathcal{C}^u \vdash s_i \xrightarrow{P} s$, which would involve the construction of the derivation tree for the instruction similar to those in subsection 3.4. Since the final states will be the same, we conclude that the two pieces of code are equivalent. This allows us to verify the correctness of the program optimizations and scheduling. However, the current framework only allows us to prove equivalence of given pieces of code. Section 5 discusses enriching this framework.

5 CONCLUSIONS AND RELATED WORK

This paper has introduced a technique for modelling architectures based on a synthesized abstract syntax. The model is based on an operational semantics derived from the structure of the machine, and is in essence a semantics of the instruction set. The hierarchical and modular structure of transport triggered architectures allows for such a modular syntax and semantics. In principle, the technique can be applied to other VLIW-like architectures which share this hierarchical structure. Since the description is manageable and readable, we have provided a precise specification of the architecture which can be used

by compiler writers and as a formal design document. The semantics has been animated, providing a simulator on which to explore the design space of the architecture. It can also be used to verify the behaviour of physical architectures.

There is much related work in the area of hardware description and modelling languages. However, much of this work is based on a lower level of description, for instance VHDL or Boyer-Moore theorem provers applied to low level descriptions (Boyer & Yu 1996, Hunt & Brock 1989). Recently there has been work on hardware description languages with a good formal semantic footing, for instance HML (O'Leary, Linderman, Leeser & Aagaard 1993) and (O'Donnell 1992) where functional languages are used as description languages. We believe that a gap exists in the specification languages lying between languages and hardware, and that the system presented in this paper goes some way towards bridging this.

We have already seen examples of using the given semantics to prove small properties about programs or instructions. Further work lies in proving properties about programs in general. Aiken (Aiken 1988) builds on a transition semantics for parallel execution by constructing a notion of an execution trace. The execution trace is defined as the successive states that a program goes through while executing (similar to our (\xrightarrow{P}) in rule 4). Using this, he develops provably correct code transformations which can be used to prove percolation scheduling correct. Although Aiken's framework has to be extended we believe that our proof framework will be useful in proving similar properties for scheduling techniques for TTAs.

6 ACKNOWLEDGEMENTS

The authors thank Marcel Beemster, Hugh McEvoy and the referees for their comments which greatly improved the paper. The first author is supported by the Netherlands Computer Science Research Foundation with financial support from the Netherlands Organization for Scientific Research (NWO).

REFERENCES

- Aiken, A. (1988), *Compaction-Based Parallelization*, PhD thesis, Department of Computer Science, Cornell University.
- Boyer, R. S. & Yu, Y. (1996), 'Automated proofs of object code for a widely used microprocessor', *Journal of the ACM* 43(1), 166–192.
- Corporaal, H. (1995), *Transport Triggered Architectures: Design and Evaluation*, PhD thesis, Technische Universiteit Delft.
- Corporaal, H. & Mulder, H. J. M. (1991), *Move: A framework for high-performance processor design*, in 'Supercomputing-91', IEEE Computer Society Press, Albuquerque, New Mexico, pp. 692–701.
- Corporaal, H. & van der Arend, P. (1993), 'MOVE32INT, a sea of gates realization of a high performance transport triggered architecture', *Mi-*

- croprocessor and Microprogramming* **38**, 53–60.
- Fisher, J. A., Ellis, J. R., Ruttenberg, J. C. & Nicolau, A. (1984), Parallel processing: A smart compiler and a dumb machine, in 'Proceedings of the ACM SIGPLAN '84 Symposium on Compiler Construction', Vol. 19 of *SIGPLAN Notices*, pp. 37–47.
- Hannan, J. (1994), 'Operational semantics-directed compilers and machine architectures', *Transactions on Programming Languages and Systems* **16**(4), 1215–1247.
- Hoogebrugge, J. (1996), Code Generation for Transport Triggered Architectures, PhD thesis, Technische Universiteit Delft.
- Hunt, W. A. & Brock, B. C. (1989), The verification of a bit-slice ALU, in M. Leeser & G. Brown, eds, 'Workshop on Hardware Specification, Verification and Synthesis: Mathematical Aspects', Vol. 408 of *LNCS*, Springer-Verlag, pp. 282–306.
- Jones, G. & Sheeran, M. (1990), Circuit design in Ruby, in J. Staunstrup, ed., 'Formal methods for VLSI design', North-Holland, pp. 13–70.
- Necula, G. C. & Lee, P. (1996), Safe kernel extensions without run-time checking, in '2nd Operating System Design and Implementation (OSDI)', Seattle, Washington.
- O'Donnell, J. T. (1992), Generating netlists from executable circuit specifications in a pure functional language, in J. Launchbury & P. Sansom, eds, 'Functional Programming, Glasgow 1992', Workshops in Computing, Springer-Verlag.
- O'Leary, J., Linderman, M., Leeser, M. & Aagaard, M. (1993), HML: A hardware description language based on standard ML, in D. Agnew, L. Claesen & R. Camposano, eds, 'Proceedings of the 11th IFIP Conference on Computer Hardware Description Languages and their Applications (CHDL'93)', IFIP, North-Holland, pp. 327–334.
- Peterson, J., Hammond, K., Augustsson, L., Boutel, B., Burton, W., Fasel, J., Gordon, A. D., Hughes, J., Hudak, P., Johnsson, T., Jones, M., Peyton Jones, S., Reid, A. & Wadler, P. (1996), Report on the programming language Haskell (version 1.3), Technical Report YALEU/DCS/RR-1106, Yale University.
- Plotkin, G. D. (1981), A structural approach to operational semantics, Technical Report DAIMI FN-19, Computer Science Department, Aarhus University, Denmark.
- Plotkin, G. D. (1982), An operational semantics for CSP, Technical Report CSR-114-82, University of Edinburgh.