

Modular Hardware Design with Timeline Types

ANONYMOUS AUTHOR(S)

Modular design is a key challenge for enabling large-scale reuse of hardware modules. Unlike software, however, hardware designs correspond to physical circuits and inherit constraints from them. Timing constraints—which cycle a signal arrives, when an input is read—and structural constraints—how often a multiplier accepts new inputs—are fundamental to hardware interfaces. Existing hardware design languages do not provide a way to encode these constraints; a user must read a module’s implementation to understand how to use it. We present Filament, a language for modular hardware design that supports the specification and enforcement of timing and structural constraints. Filament uses *timeline types*, which describe the intervals of clock-cycle time when a given signal is available or required. Filament enables *safe composition* of hardware modules, ensures that the resulting designs are correctly pipelined, and predictably lowers them to efficient hardware.

1 INTRODUCTION

Like software languages, interfaces in hardware description languages (HDLs) simply consist of arguments and simple datatypes. Unlike software, however, hardware inherits constraints from the underlying physical circuits—the inputs are used and outputs are available during specific cycles, and new inputs may only be provided when the circuit can process them. The rudimentary interfaces of existing HDLs fail to capture these constraints making modular design difficult. To approach the reusability of software library ecosystems, HDLs need a systematic way to encode the requirements for hardware modules.

Modern languages for hardware design fall into three categories. Embedded HDLs [2, 13, 24, 29, 35] (eHDLs) use software host languages for metaprogramming. Accelerator design languages (ADLs) [17, 20, 21, 26, 33, 44] are higher-level languages that expose new abstractions and compile to HDLs. Finally, traditional HDLs, such as SystemVerilog and VHDL, are the de facto standard for hardware design and interacting with hardware toolchains. ADLs and eHDLs must be compiled to HDLs to interact with hardware design toolchains and integrate proprietary black box implementations of optimized primitives. Furthermore, ADLs are often limited to a single domain, such as image processing, and can therefore benefit from defining a foreign function interface (FFI) to interact with eHDLs and other ADLs.

Composition in HDLs is challenging because interfaces only expose the names and value types of input–output ports. However, efficient integration requires knowledge of timing behavior: number of cycles required to produce an output, number of cycles required to consume an input, and whether the module can be pipelined. In current HDLs, this timing information is latent. It appears only in natural-language comments or—depressingly often—nowhere at all, requiring programmers to learn how to use each module by reading its source code.

The key to an ecosystem of reusable hardware is a low-level mechanism to *safely* and *efficiently* compose hardware modules. *Safe composition* requires module interfaces to specify timing details such as latency and pipelinability, while *efficiency* requires that the interfaces do not add substantial overheads. The effect is a clear way to integrate hardware, regardless of whether it was written in an eHDL, generated by an ADL, or implemented as a proprietary black box module.

Our solution is *timeline types*, which compactly encode latency and throughput properties of static hardware pipelines. Static pipelines have data-independent timing behavior and encompass a large class of efficient hardware designs [17, 20, 21, 25]. Our type system, inspired by separation logic [39], proves that pipelined execution of a module is safe, i.e., there are no resource conflicts. Our contributions are as follows:

- We provide a characterization of pipelining constraints for static pipelines and model them using *timeline types* in an HDL called Filament.

- We formalize these pipelining constraints using *log-based semantics* of hardware and prove our type system is sound with respect to the model.
- We demonstrate that timeline types can express the timing behavior of modules generated from several state-of-the-art hardware generators [17, 25, 41] and how we can integrate them in Filament designs.
- We show that Filament designs use fewer resource and run at faster frequencies than those generated by hardware generators.

2 EXAMPLE

We'll discuss the challenges associated with compositional hardware design by implementing a pipelined arithmetic logic unit (ALU).

2.1 Traditional Hardware Description Languages

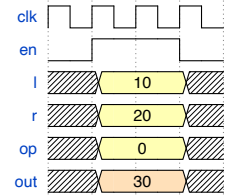
We start off by implementing an ALU that can perform addition or multiplication in a traditional HDL. The ALU's circuit consists of an adder, a multiplier, and a multiplexer which selects the output of either the adder or the multiplier. The interfaces for the module specify the inputs and outputs along with their bitwidths. The ALU itself simply instantiates the modules and forwards the input signals: the adder and multiplier perform computations in parallel and the multiplexer selects the output based on the op signal.

```

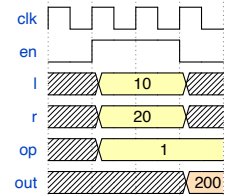
module Add(a: 32, b: 32) -> (o: 32)
module Mul(a: 32, b: 32) -> (o: 32)
module Mux(
  sel: 1, a: 32, b: 32) -> (o: 32)
module ALU(
  op: 1, l: 32, r: 32) -> (o: 32) {
  Mul M(l, r); Add A(l, r);
  Mux Mx(op, A.out, M.out);
  o = Mx.out; }

```

We will use waveform diagrams to understand the execution behavior of this module. A waveform diagram explains the flow of signals in the circuit over time and usually with respect to the global clock signal. For example, the following waveform diagram shows that when provided with the inputs 10 and 20 and the op code 0, the ALU produces the output 30 in the same cycle.



However, if we provide set op to 1, the timing behavior changes—the product is produced two cycles after the input is provided. Additionally, if the op is not asserted for an additional cycle, the output is wrong. The problem is that an adder is *combinational*—it produces its output in the same cycle as the inputs—while a multiplier is *sequential*—it takes several cycles to produce its output. op is required for an extra cycle because the multiplier output is produced later than the adder and the multiplexer needs to select the correct output in a later cycle using the op input.



The interfaces for ALU, the adder, and the multiplier do not capture these details. One option to sidestep this problem is to “wrap” every module in a *latency-insensitive* interface, such as ready-valid handshaking. But these interfaces incur overhead that can be prohibitive for fine-grained composition [32]. This paper aims to specify efficient, *latency-sensitive* interfaces based on clock cycles and to statically rule out misuses of these interfaces.

2.2 Filament

Filament is an HDL that allows users to directly *specify* and *check* the timing behavior of their modules. Each component can be parameterized by multiple events which are used to specify the timing behavior of the module. Our ALU implementation has unpredictable timing behavior because adders and multipliers have different timing behavior. Filament allows us to encode their timing behavior explicitly using *events* which parameterize modules:

```

99  extern component Add<T>(
100    @interface[T] go: 1, @[T, T+1] left: 32, @[T, T+1] right: 32) -> (@[T, T+1] out: 32);
101  extern component Mult<T>(
102    @interface[T] go: 1, @[T, T+1] left: 32, @[T, T+1] right: 32) -> (@[T+2, T+3] out: 32);

```

Both components use the event T to specify their timing behavior. The adder is *combinational*—it produces outputs in the same cycle as the inputs. This fact is encoded by the *availability intervals* of the inputs and outputs: the inputs are provided in the half-open interval $[T, T + 1)$, which corresponds to the first cycle of execution of the component, and the output is produced during the same interval. In contrast, a multiplier is *sequential*—it takes two cycles to produce its output. This is encoded by stating that the output is available in the interval $[T + 2, T + 3)$, two cycles after the inputs are provided in the interval $[T, T + 1)$. In order to signal that the event T has occurred, a user of these modules must set the *interface port* `go` to 1, provide the inputs according to their required intervals and read the output when they are available. Multiplexers (not shown) are also combinational and take all their inputs in the same cycle.

Like our HDL implementation, our Filament implementation of the ALU explicitly instantiates all the hardware resources it needs to use. The key difference is how Filament expresses the use of the hardware instances through *invocations*. An invocation schedules the execution of a hardware instance using a particular set of events and provides all inputs. For example, the invocation `a0` of the adder `A` is scheduled using the event `G`. By naming uses, Filament can check the timing behavior of the module. There is no assignment for the `go` port of the adder—it is automatically inserted by the compiler using the scheduling event `G`. Invocations are a logical construct that are compiled away by Filament (Section 5). Similarly, the multiplier and multiplexer are also scheduled using the event `G`. Instead of using outputs from the instance, the multiplexer uses outputs the invocations which reflect output from a specific use.

```

component ALU<G>(
  @interface[G] en: 1
  @[G, G+1] op: 1,
  @[G, G+1] l: 32,
  @[G, G+1] r: 32,
) -> (@[G+2, G+3] o: 32) {
  A := new Add; M := new Mult;
  Mx := new Mux;
  a0 := A<G>(l, r);
  m0 := M<G>(l, r);
  mux := Mux<G>(
    op, m0.out, a0.out);
  o = mux.out; }

```

2.3 Checking Timing Behavior

However, when we attempt to compile this program, Filament gives us the following error:

```

131  mux := Mux<G>(op, m0.out, a0.out);
132  Available for [G+2, G+3) but required during [G, G+1)

```



The error states that the use of our multiplexer expects all of its inputs during $[G, G + 1)$ while the multiplier's output, `m0.out`, is available in $[G + 2, G + 3)$. Filament requires that all inputs be available for at least as long as the corresponding argument's requirement. This was the problem in our original HDL design (Section 2.1)—the output of the adder is available in a different cycle from the multiplier which results in unexpected timing behavior. Filament's type system statically catches this error.

The solution is to use *registers* to store values and make them available in future cycles. A register's signature captures its timing behavior—the output is available one cycle after the input:

```

142  component Reg<G>(@interface[G] en: 1, @[G, G+1] in: 32) -> (@[G+1, G+2] out: 32)

```

The corrected implementation uses two registers to make the sum available in the same cycle as the multiplier. The first register's output is available in $[G + 1, G + 2)$ while

```

component ALU<G>(@[G, G+3] op: 32, ...) {
  a0 := A<G>(l, r); R0 := new Reg; R1 := new Reg;
  r0 := R0<G>(a0.out); r1 := R1<G+1>(r0.out);
  mux := Mux<G+2>(op, r1.out, m0.out); ... }

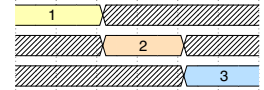
```

99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147

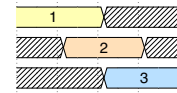
148 the second register's output in $[G + 2, G + 3)$, the same cycle when the `m0.out` is available. We
 149 schedule the execution of the multiplexer in cycle $G + 2$ when both the outputs are available. This
 150 design is still problematic because `op` is only available in $[G, G + 1)$ while the multiplexer reads it
 151 in $[G + 2, G + 3)$. We fix this by making `op` signal available in $[G, G + 3)$. This results in a correct
 152 ALU implementation. However, it is not clear when the ALU is ready to accept new inputs: should
 153 we wait till outputs are produced or can the module process multiple inputs in parallel?
 154

155 2.4 Pipelining

156 *Pipelining* is a common optimization that enables hardware to process
 157 multiple inputs in parallel. For example, **diagram (a)** shows how a se-
 158 quential module processes its inputs—one at a time. **Diagram (b)**, shows
 159 pipelined execution that can overlap the processing of multiple inputs.
 160 Pipelining is challenging because it requires reasoning about the inter-
 161 action between multiple, concurrent executions of the same physical
 162 resources—correctly pipelining requires using values from the correct
 163 pipeline stage, and ensuring there are no *control hazards*, i.e., there are
 164 no conflicting uses of internal components.



(a) Sequential processing

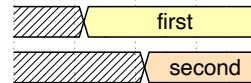


(b) Pipelined processing

```
component Add<T: 1>(…)
component Mult<T: 3>(…)
component ALU<G: 1>(…)
```

165 Filament presents a concise solution: each event has an associated
 166 *delay* that specifies how many cycles to wait before accepting new
 167 inputs. We can update the signature of the adder and multiplier to
 168 reflect this. Since the adder is combinational, it can accept new inputs every cycle. However, the
 169 multiplier accepts new inputs every 3 cycles. For user-level module, Filament ensures that the
 170 *delay* of the event is correct through type-checking. We'll redesign our ALU to be pipelined and
 171 accept new inputs every cycle by specifying that the delay of G is 1. Since we know our design is
 172 not pipelined, Filament will generate errors explaining why the design cannot be pipelined.
 173

```
174 component ALU<G: 1>(
175     Event may retrigger every cycle
176     @[G, G+3] op: 1, Signal lasts for 3 cycles
```



177 Our first problem is that the signature requires input signal `op` to be available for three cycles
 178 whereas the pipeline may trigger every cycle. The waveform diagram demonstrates the problem—
 179 the input for `op` from the first iteration will overlap with the input for the second iteration. How-
 180 ever, `op` is a *physical port* in a circuit and can only hold one value at a time; this is a fundamental
 181 physical constraint of hardware design. Filament requires that the delay of an event is longer than
 182 the length of any availability interval that uses it; we must make `op`'s availability interval 1-cycle
 183 long. We choose $[G + 2, G + 3)$ since the the multiplexer uses `op` during this interval.

184 Next, Filament complains that while our ALU pipeline
 185 may accept new inputs every cycle, the multiplier `M` can
 186 accept new inputs every 3 cycles. This is a fundamental
 187 limitation of the multiplier circuit we're using; to fix it, we
 188 must use a different multiplier. Filament catches yet an-
 189 other pipelining bug that arises from composition: every
 190 subcomponent used in a pipeline must be able to process inputs at least as often as the pipeline
 191 itself. Fixing this will result in a correct, fully pipelined ALU.

```
component Mult<T: 3>(
  Event may retrigger every 3 cycles
component ALU<G: 1>(
  Event may retrigger every cycle
  m0 := M<G>(l, r);
  Cannot safely pipeline
```

```
192 component ALU<G: 1>(@interface[G] en: 1, @[G+2, G+3] op: 1, ...) {
193   A := new Add; Mx := new Mux; R0 := new Reg; R1 := new Reg; FM := new FastMult; // delay = 1
194   a0 := A<G>(l, r); r0 := R0<G>(a0.out); r1 := R1<G+1>(r0.out); m0 := FM<G>(l, r);
195   mux := Mux<G>(op, r1.out, m0.out); o = mux.out; }
```

2.5 Area-Throughput Trade-offs with Filament

We show how Filament enables us to safely explore *area-throughput* trade-offs by implementing three different versions of a divider using a restoring division algorithm. The `Init` component generates the initial quotient and accumulator values from the input dividend. The `Next` component computes the next quotient (QN) and accumulator (AN) using the divisor, current quotient, and accumulator. For an 8-bit value, we must apply `Next` 8 times.

Combinational divider. A combinational divider computes its output in same cycle when the inputs are provided. Such a design is generally inefficient because it schedules a lot of complex logic in the same clock cycle and forces the design to operate at a low frequency. However, combinational designs are a good starting point to ensure that our algorithm is correct. All `Next` instances are scheduled using the event `G` which mean they'll execute in the same cycle.

```
component Comb<G: 1>(…) -> (
  @[G, G+1] q: 8) {
  I := new Init;
  i := I<G>(left);
  N0 := new Next;
  s0 := N0<G>(i.A, div, i.Q);
  ...
  s7 := N7<G>(s6.AN, div, s6.QN)
  q = s7.QN
```

Pipelined divider. To make our design run at a higher frequency, we can pipeline it by scheduling each `Next` instance to execute in a different cycle. We'll add registers to store the values generated by each instance and use it in the next cycle. The delay of the module remains 1, allowing it to process a new value every cycle, the latency is now 7 cycles unlike the combinational module.

```
component Pipe<G: 1>(…) -> (
  @[G+7, G+8] q: 8) {
  I := new Init; i := I<G>(left);
  N0 := new Next; s0 := N0<G>(…);
  RA0 := new Reg; RQ0 := new Reg
  ra0 := RA0<G, G+1>(s0.AN);
  rq0 := RQ0<G, G+1>(s0.QN); ...
  s7 := N7<G+7>(ra6.out, div, rq6.out)
```

Iterative divider. Both the combinational and pipelined inputs can process a new input every cycle but require a large amount of hardware since they instantiate 8 instances of the `Next` component. We can instead use the same `Next` component iteratively and reduce the size of our design. We start with our original combinational design and replace each `Ni` with `N` to express resource sharing.

Filament tells us that this design is buggy. We're attempting to send two different inputs into the `Next` instance in the same cycle. However, `Next` is a physical circuit and can only process one input every cycle. Therefore, we must schedule the uses of the instance in different cycles. We can start with our pipelined design, which schedules different instances of `Next` in different cycles and replace each instance with `N`.

```
component Next<T:1>(…)
  Delay requires uses to be 1 cycle apart
  s0 := N<G>(i.A, div, i.Q); First use
  s1 := N<G>(s0.AN, div, s0.QN); Second use
```

Filament complains with a new error message. Since we're sharing the instance `Next` over 8 cycles, we cannot start a new pipelined execution every cycle. Again, this is because `Next` is a physical circuit that can only process one input a cycle. To fix this, we can change the delay to 8 cycles which guarantees to Filament that the instance will only be run every 8 cycles. This ensures that all iterations using the instance `N` complete before new inputs are provided. Implicitly, Filament showed us that the reusing the instance is a trade-off: while we use fewer resources, our throughput is also reduced since our iterative implementation can only process a new input every 8 cycles compared to every cycle for the pipelined implementation.

```
component Iter<G:1>(…)
  Event may trigger every cycle
  causing shared uses to conflict
  s0 := N<G>(i.A, div, i.Q);
  First use
  s7 := N<G+7>(s6.AN, div, s6.QN);
  Last use
```

197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245

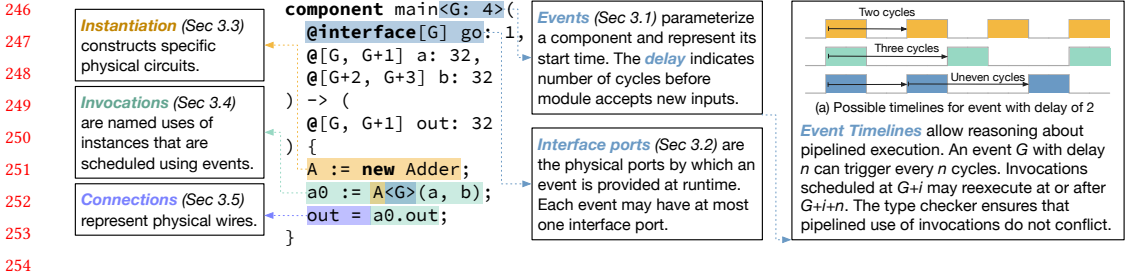


Fig. 2. Overview of the Filament language. Programs are a sequence of *component* definitions which correspond to individual modules. The signature of the component is parameterized using *events* and the body of component consists of three types of statements: Instantiations, connections, and invocations.

2.6 Summary

Filament is an HDL for safe design and composition of static pipelines. Specifically, Filament programs can *specify* and *check* timing properties of hardware modules and ensure that:

- (1) Values on ports and wires are only read when they are *semantically* valid.
- (2) Hardware instances are not used in a conflicting manner.

These properties ensure that the resulting pipelines are safe, i.e., there are no resource conflicts, and efficient, i.e., they can overlap computation as specified by their interface without any overhead. Filament's utility extends to components defined outside the language as well. By giving external modules a type signature, users can safely compose modules. Section 3 overview the constructs in Filament, Section 4 explains how Filament's type system checks pipeline safety, and Section 5 shows how Filament's high-level constructs are compiled to efficient hardware.

3 THE FILAMENT LANGUAGE

Figure 2 gives an overview of the Filament language. Filament's level of abstraction is comparable to *structural* HDLs where computation must be explicitly mapped onto hardware. Filament only has four constructs: components, instantiation, connections, and invocations. The first three have direct analogues in traditional HDLs while invocations are a novel construct.

3.1 Events and Timelines

Events are the core abstraction of time in Filament. Instead of using a clock signal, designs use events to schedule computation. The Filament compiler generates efficient, pipelined finite state machines to reify events (Section 5.2).

Defining events. There are only two ways to define events: (1) component signatures bind *event variables* like G , and (2) users can write *event expressions* such as $G+n$ where n is a constant. Events have a direct relationship to clock: if G occurs at clock cycle i , then $G+n$ occurs at clock cycle $i+n$.¹ This relationship with clock is crucial since it allows Filament to represent timing properties of components defined in clock-based HDLs. Adding event variables ($G_0 + G_1$) is disallowed since events correspond to particular clock cycles, and it is meaningless to add them together.

Timeline interpretation of events. In order to capture potential resource conflicts from pipelined execution, Filament interprets events as a set of possible timelines. A timeline for an event G with a delay n is any infinite sequence of 1-cycle pulses clock cycles such that each pulse is at least n cycles

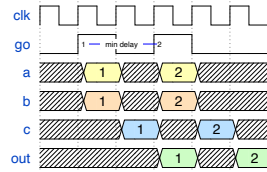
¹All event variables operate in the same clock domain, but this limitation can be removed in the future.

295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343

```

component AddMult<G: 2>(
  @interface[G] go: 1,
  @[G, G+1] a: 32, @[G, G+1] b: 32, @[G+1, G+2] c: 32
) -> (@[G+2, G+3] out: 32) { commands }
    
```

(a) Component's signature in Filament



(b) Waveform showing pipelined use

Fig. 3. Signature and waveform diagram. The component allows pipelined execution or reuse after two cycles allowing overlapped execution. Shaded regions represent unknown values.

apart. Figure 2a shows a set of valid timelines for an event with delay 2. By imbuing events with a timeline interpretation, Filament can reason about *repeated execution* and consider how pipelined executions may affect each other. By reasoning about such properties, we can define and enforce safety properties for pipelined execution of hardware. Furthermore, the timeline interpretation has a direct relationship to hardware: the delay of an event represents how many cycles a user must wait before providing a new set of inputs. This is usually referred to as the *initiation interval* of a pipeline by hardware designers (Section 4.3).

3.2 Component

Filament programs are organized in terms of *components* which describe timing behavior in their signatures and their circuit using a set of commands. Figure 3 shows the signature of a component in Filament (Figure 3a) and a waveform diagram visualizing two sets of inputs being processed in parallel (Figure 3b). The component is parameterized using the event G with a delay of 2 which means that pipelined use can begin two cycles after the previous use.

Interface ports. Hardware components typically have *control ports* which signal when values on *data ports* are valid and that the computation should be performed. Values on control ports are always considered semantically valid while values on data ports are only valid when the corresponding control port is high. Filament distinguishes control ports by defining them as *interface ports*. Interface ports are 1-bit ports that are associated with a particular event. When an interface port is set to 1, it signals to the component that the corresponding event has occurred. For example, setting `go` to 1 on an `AddMult` instance (Figure 3a) makes the module start processing the inputs. The availability intervals of all ports that use an event are relative to when the corresponding event's interface port is set to 1. If an event does not have an interface port, then the module can assume that the event triggers every n cycles where n is the event's delay.

Availability intervals. The input and output ports of the component describe their availability in terms of the events bound by a component. For `AddMult` (Figure 3a), all ports use the event G . Availability intervals are *half-open*: for example, the input port `a` is available during $[G, G + 1)$ which means it is available during the first cycle when the component is invoked. Inside the body of a component, an input's availability interval represents a *guarantee* while an output's availability requires a *requirement* that the body must fulfill. When using a component, this is reversed: inputs have requirements that must be fulfilled by the user while outputs have guarantees.

3.3 Instance

All computations in a hardware design must be explicitly mapped onto physical circuits. Filament's new keyword allows instantiation of subcomponents.

```

component Add<T: 1>(@[T, T+1] left: 32, @[T, T+1] right) -> (@[T, T+1] o: 32);
    
```

```
component AddTwo<G: 1>(…) { A0 := new Add; A1 := new Add; }
```

The above program instantiates two instances of the `Add` component named `A0` and `A1` that can be used independently. Note that the instantiations do not provide a binding for the event `T` of for `Add`; *invocations* are responsible for providing those and scheduling the execution of an instance.

3.4 Invocation

Resource reuse in hardware designs is *time-multiplexed*, i.e., different uses of the same resources are scheduled to occur at different times. This is done by building a finite state machine (FSM) using a register and using the output of the register to select which inputs to use. The example program computes $(l \times r)^2$ using a single multiplier using the FSM `F` to forward the inputs `l` and `r` into the multiplier in the first cycle and the output of the multiplier in the second cycle. However, the assignment to `M.left` incorrectly forwards the value from `M.out` in the first cycle. Mistakes in the control logic for the FSM do not lead to any visible errors; this error will lead to the data getting silently corrupted and propagating into to other parts of the system.

```
F := new Reg; // FSM
F.in = F.out == 0 ? l : 0;
M := new Mult; A := new Add;
M.right = F.out == 0 ? r : M.out
M.left = F.out == 1 ? l : M.out
```

In contrast, every use of an instance in `Filament` must be explicitly named and scheduled through an invocation. The first invocation of the multiplier `M` is scheduled using the event `G`, uses the inputs `l` and `r`, and is named `m0`. The second invocation, scheduled one cycle later at `G + 1`, can then use `m0.out` to refer to the output of the first execution and pass it into the multiplier as an input. Because the second invocation is scheduled one cycle later, the input ports have a different requirement: the inputs must be available in the interval $[G + 1, G + 2)$ as opposed to $[G, G + 1)$ in the first invocation. This allows `Filament` to check that `m0.out` is semantically valid when it is used as an input to `m1` and that the two uses of the multiplier are scheduled to occur at different times, allowing the compiler to generate correct FSMs to schedule instance reuse. Each invocation only provides inputs for the data ports and elides inputs for the interface ports. During compilation, `Filament`'s compiler automatically infers assignments for the input ports and generates efficient, pipelined FSM to schedule the invocations (Section 5).

```
component Square<T: 1>(
  @[T, T+1] left: 32,
  @[T, T+1] right: 32
) -> (
  @[T+1, T+2] out: 32);
M := new Mult;
m0 := M<G>(l, r)
m1 := M<G+1>(
  m0.out, m0.out)
```

3.5 Connection

`Filament` programs allow ports to be connected and requires that the source is semantically valid for at least as long as the destination.

```
component Add<G>(@[G, G+3] source: 32) -> (@[G, G+1] dest: 32) { dest = source; }
```

Connections are physically implemented as wires connecting two ports in the circuit and are continuously active.

3.6 Interfacing with External Components

`Filament`'s `extern` keyword allows the user to provide type-safe wrappers for black box modules by specifying a type signature without a body. `Filament`'s standard library, which provides signatures for components like multipliers and registers, is defined using `extern` components.

Phantom events. Phantom events allow `Filament` to model the behavior of components like adders which are *continuously* active and do not take an explicit enable signal. In the following signature, the event `G` is a phantom event because there is no corresponding interface port for it in the signature. Section 5.4 describes how user-level components can use phantom events.

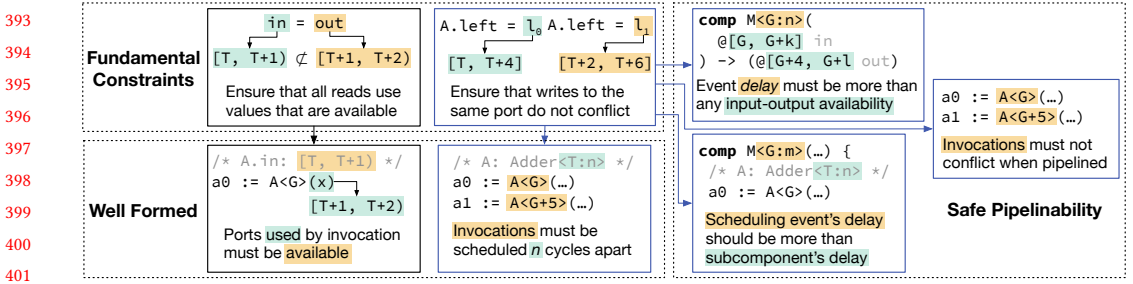


Fig. 4. Overview of the Filament type system. The fundamental constraints of hardware design imply other constraints. Well-formedness ensures that *one execution* of a component is correct. Safe pipelining ensures that *pipelined executions* of the component are correct.

```
extern component Add<G: 1>(@[G, G+1] l: 32, @[G, G+1] r: 32) -> (@[G, G+1] o: 32))
```

Ordering constraints. In order to capture the full expressivity of external components, Filament allows defining ordering constraints between events. For example, combinational components can provide a valid output for more than one cycle if the inputs are provided for multiple cycles. Therefore, a more precise interface of a combinational adder is:

```
component Add<G: L-G, L: 1>(@[G, L] l: 32, @[G, L] r: 32) -> (@[G, L] o: 32) where L > G
```

The events G and L mark the start and end for the input and output availability intervals. In order to ensure that the interval $[G, L]$ is well-formed, the signature requires $L > G$. The component guarantees that the output is provided for as long as the inputs are provided.

Parameteric delays. The new signature of adder additionally specifies a *parameteric delay* of $L - G$ cycles to signal that the adder may not be reused while it is processing a set of inputs. In order to generate *static pipelines* which have input-independent timing behavior, Filament requires all such expressions evaluate to a constant value. Like the *example*, an invocation of `Add` must provide some binding of the form $G = T + i$ and $L = T + k$ such that $k > i$, ensuring that the delay for the corresponding invocation is a compile-time constant $k - i$ and the ordering constraint $L > G$ is satisfied.

```
A := new Add;
// delay = (G+3)-G = 3
a0 := A<G, G+3>(x, y);
```

The signature of registers in Filament allows them to provide the output for as long as needed, similar to an adder. However, because a register is a state element, it only requires its input for one cycle. Furthermore, the delay signals that the register can accept a new write during the last cycle when the output is available.

```
component Register<G: L-(G+1), L: 1>(
  @interface[G] go: 1, @[G, G+1] in: 32) -> (@[G+1, L] out: 32) where L > G+1;
```

4 TYPE SYSTEM

Filament’s type system enforces two fundamental restrictions of hardware design:

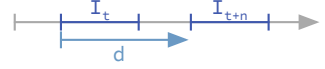
- (1) All reads only use *semantically* valid values. A port or wire will always have a value on it. Filament’s availability intervals mark when the values are semantically valid.
- (2) Writes do not conflict. This is a corollary of the property that uses of a resource must not conflict because use of a resource is represented through a write.

Filament ensures these properties using two checking phases: well-formedness checking, which ensures that a *single execution* of a component is correct, and *safe pipelining*, which ensures that *pipelined executions* of a component are correct.

4.1 Delay Well-Formedness

The delay of an event encapsulates all possible conflicts between parts of the pipeline scheduled using it. Filament requires that the delay of an event is at least as long as each interval that mentions it which ensures that instance reuse does not create conflicts between its input and output ports.

The proof is straightforward: for two invocations at time t and $t + n$ such that $n \geq d$ where d is the delay, let I_t and I_{t+n} be the availability intervals of the input i . Since we know that the start times of the intervals are at least d cycles apart ($I_{t+n} - I_t \geq d$), and that length of the intervals is bounded by d ($|I_t| \leq d$) we can conclude that they do not overlap.



4.2 Well-Formedness

Valid reads. In order to ensure this property, Filament needs to make sure that port values are only read when they are semantically valid. Signals are used in two places:

- (1) *Connections* (Section 3.5) forward a value from one port to another. Filament ensures that the availability of the output port is at least as long as the requirement of the input port.
- (2) *Invocations* (Section 3.4) schedule the use of a component instance using a set of events. Checking the validity of an invocation boils down to two steps: the requirements of the instance's input ports can be computed by binding the event variables in its signature to the invocation's event. Next, each argument essentially represents a connection between the instance's input and the argument and is checked using the criteria for connections.

Conflict-free. If an invocation schedules an instance with delay d using the event G , the instance may not be reused between $[G, G+d)$. This both ensures that there are no conflicts between input and output ports (Section 4.1) and that none of the subcomponents conflict. The latter property holds because safe-pipelining constraints ensure that a valid delay can correctly encapsulate all possible conflicts between subcomponents (Section 4.4). In the example program, the two invocations of M overlap causing Filament to reject this program.

```
component Mult<T: 3>(...);
component main<G:10>() {
  M := new Mult;
  // busy b/w [G, G+3]
  a0 := M<G>(a, b);
  // busy b/w [G+1, G+4]
  a1 := M<G+1>(a0.out, b);
```

4.3 Initiation Intervals

Pipelining is an important optimization since it allows a module to process multiple inputs in parallel. For example, a multiplier with a three cycle latency but an *initiation interval* of one cycle to compute an output but can accept new inputs every cycle. In Filament, the delay of an event corresponds to initiation interval. While hardware designers talk about initiation intervals of a component, Filament generalizes it by allowing a component to have multiple events. In this case, each event specifies the initiation interval of some part of the internal pipeline. Filament ensures that the delay of a module describes a valid initiation interval, defined as follows:

Definition 4.1 (Initiation Interval). Let $P(t)$ be the execution of pipeline P at time t . $P(t_0) \perp P(t_1)$ states that the pipeline executions of P at t_0 and t_1 do not have resource conflicts. Then I is a valid initiation interval of pipeline P if and only if

$$\forall n \geq 0 P(t) \perp P(t + I + n)$$

This definition requires that the pipeline be able to accept new inputs after *any* amount of time after the initiation interval. There might be other delays smaller than the initiation interval which allow the pipeline to accept new inputs in a small window of time before becoming invalid again. This would correspond to the following definition of an initiation interval I :

$$\forall k \neq 0 P(t) \perp P(t + k \times I)$$

Filament uses the first definition because delays are also used to check for well-formedness constraints of a component. If we used the second definition, the well-formedness constraint would require that if an instance is scheduled at time t , it only be scheduled again at other times $k * t$ which we think is less compositional. Regardless, this is not a fundamental limitation since both definitions can be encoded and enforced.

4.4 Safe Pipelining

While well-formedness ensures that one execution of a module is correct, i.e., all reads use valid values and there are no conflicts, safe pipelining must ensure that *pipelined executions* of the component do not create any additional conflicts. Checking that pipelined execution is very similar to checking that invocations of the same instance do not conflict. This is because pipelined execution is exactly the same—an instance being reused after a period of time. Filament must show that for an invocation scheduled using event G , another invocation scheduled at any time after $G + d$ (where d is the delay) does not conflict with the first invocation. The following checks are sufficient to prove this.

Triggering Subcomponents. Filament requires that when an event is used to invoke a subcomponent, the event's delay must be at least as long as the delay of the subcomponent's event.

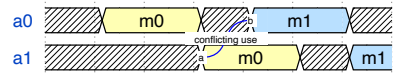
```
component Mult<G: 3>(@interface[G] go: 1, ...)
component main<T: 1>(@interface[T] go: 1, ...) { M := new Mult; m0 := M<T+2>(…) }
```

The event $T + 2$ is used to schedule the invocation of instance M which has a delay of 3. However, $T + 2$ has a delay of 1, same as T . This is problematic because `main` may trigger every cycle while M can only support computations every 3 cycles. Filament therefore rejects this program.

Reusing Instances. Previous checks already ensure that: (1) shared invocations do not conflict during one execution of the pipeline, and (2) pipelined execution of an invocation does not conflict with itself. However, we also need to ensure that pipelined invocations of a shared instance do not conflict with each other.

```
component Mult<G: 3>(…)
component main<T: 3>(…) {
  M := new Mult;
  m0 := M<T+2>(…);
  m1 := M<T+10>(…);
}
```

The example program will pass all our previous checks but is erroneous: executing the pipeline at time T and $T + 10$ will cause the $m1$ from the time T execution to



conflict with $m0$ from the time $T + 10$ execution. Because Filament's definition of initiation interval allows reexecution at *any time* in the future, we must require that all invocations of a shared instance complete before the pipelined execution begins. The following is sufficient to ensure this: the delay must be greater than the number of cycles between the start of the earliest invocation and the end of the last invocation of a shared instance.

Dynamic Reuse. Since Filament components can be parameterized by multiple events, it is possible to invoke an instance using two different events. In the example program, the type-checker would have to prove that the intervals $[G, G + 3)$ and $[L, L + 3)$ do not overlap to enforce conflict freedom. The constraint $L \geq G + 3$ is sufficient to

```
component Dyn<G: ??, L: ??>() {
  M := new Mult;
  a0 := M<G>(a, b);
  a1 := M<L>(a0.out, b); }
}
```

491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

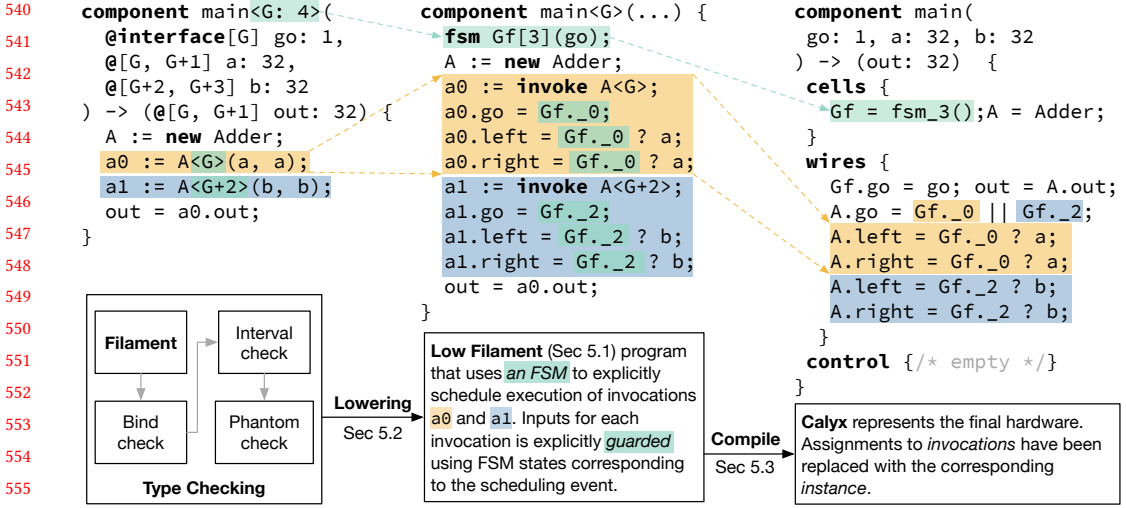


Fig. 5. Compilation Flow. Filament programs are type checked (Section 4) and lowered to *Low Filament* (Section 5.1) programs. Lowering (Section 5.2) instantiates explicit FSMs to schedule invocation. Finally, low filament programs are compiled to Calyx [34] which optimizes the design and generates hardware circuits.

prove this. However, there is no way to statically pipeline this module: the delay of G is *dynamic*, it depends on the exactly which cycle L is provided which cannot be known a priori. There is no compile-time constant value that can express the delays for both events. This is because delays describe the timeline for a single event whereas dynamic modules require relating multiple events. Filament’s solution is to disallow ordering constraints between events in user-level components which disallows the example program. External components (Section 3.6) can still use ordering constraints, but such constraints can only be satisfied using the natural order defined on $G + n$ events. This means in a well-typed program:

- (1) All delays evaluate to compile-time constants.
- (2) Invocations of a shared component all use the same event.

These constraints allow the compiler to generate efficient, statically timed pipelined from well-typed programs. Extending Filament with safe dynamic pipelines is an avenue for future work.

5 COMPILATION

Figure 5 shows an overview of the compilation flow. The primary goal of Filament’s compilation pipeline is to transform the abstract schedules of invocations into explicit, pipelined control logic. The compiler first lowers programs into *low-level Filament* which is an untyped extension of the Filament language that explicitly uses pipelined finite state machines (FSMs) to coordinate the execution of a module. Next, the compiler translates the program into the Calyx intermediate language [34] which performs generic optimizations and generates circuits.

5.1 Low-level Filament

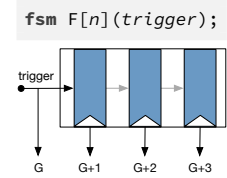
Low-level Filament is an untyped version of Filament that introduces new constructs to explicitly represent the pipelined execution of a module.

Explicit Invocations. Low-level Filament requires all ports corresponding to an invocation be explicitly assigned. This includes interface ports, which high-level Filament manages implicitly.

589 **Guarded Assignment.** Filament uses *guarded assignments* to express multiplex-
 590 ing of signals and correspond directly to guarded assignments in Calyx [34]. The
 591 assignment only forwards the value from `out` when the guard is active. Other-
 592 wise, the value forwarded to `in` is undefined. Calyx’s well-formedness condition requires that only
 593 one of the guards is active at a time for any given source port.

```
in = g1 ? out;  
in = g2 ? out;
```

594
 595 **Finite state machines.** Low-level Filament also provides the `fsm` construct
 596 to explicitly instantiate a pipeline FSM. This defines the FSM `F` with n states
 597 and a single input port `trigger` which triggers its execution. This generates
 598 a shift-register of size n with ports: `F._0`, ..., `F._{n-1}`. If `trigger` is set to 1
 599 at event G , the port `F._i` will become active at event $G + i$.



600 5.2 Generating Explicit Schedules

601
 602 The compilation from Filament to low-level Filament ensures that all high-
 603 level invocations have been compiled into explicit invocations. Figure 5 shows the compilation
 604 process for a program that uses an adder (`A`) through two invocations (`a0` and `a1`).

605 **FSM Generation.** The compiler instantiates an FSM for each event parameterizing the module. The
 606 example program uses event G to schedule the invocations. The compiler walks over all expressions
 607 $G + i$ in the program to compute the number of stages for the pipelined FSM. While the original
 608 program does not explicitly mention the event $G + 3$, it is implied by the output port `a1.out` which
 609 is active in the interval $[G + 2, G + 3]$. The compiler instantiates the FSM `Gf` with 3 states triggered
 610 by the `go` signal. Note that the delay of the FSM *does not* affect the generation of the FSM.

611
 612 **Triggering Interface ports.** The compiler then lowers the invocations by generating explicit as-
 613 signments to the adder’s interface port `go`. The first invocation, scheduled at G , uses the port `Gf._0`
 614 to trigger the invocation while the second invocation, scheduled at $G + 2$, uses the port `Gf._2`.

615
 616 **Guard Synthesis.** In order to ensure that assignments from the two invocations to the data ports
 617 `left` and `right` do not conflict, the compiler synthesizes guards for the assignments. If the input
 618 port of an invocation require inputs during the interval $[G + s, G + e]$, the compiler generates
 619 the guard `Gf._s || ... || Gf._e` for the guard. Since the program is well-typed, the guard
 620 expressions for each invocation are guaranteed to not conflict (Section 4).

621 5.3 Lowering to Calyx

622
 623 Low-level Filament is intentionally designed to be close to Calyx, so compilation is straightforward.
 624 For each FSM size n , we generate a Calyx component and instantiate it for the corresponding Fila-
 625 ment component. The FSM is simply a sequence of registers connected together. Since assignments
 626 to all ports are explicit in low-level Filament, we can simply compile the invocations by replacing
 627 them with the corresponding instance name. In the example program, assignments to both `a0.`
 628 `left` and `a1.left` are compiled to assignments to `A.left`. Since Filament guarantees that the
 629 generated guards are disjoint, we can be sure that Calyx will generate correct FSMs.

631 5.4 Optimizing Continuous Pipelines

632
 633 Continuous pipelines do not make use of a signal to indicate when their inputs are valid and instead,
 634 they continuously process inputs. We can express such pipelines in Filament using *phantom events*
 635 (Section 3.6). Phantom events do not have a corresponding interface port and therefore cannot be
 636 used trigger invocations. Filament ensures that a phantom event is used correctly through its
 637 phantom check analysis which ensures:

<pre> 638 $x \in vars \quad t \in events \quad p, q \in ports$ 639 640 $M ::=$ 641 def $C \langle t : n \rangle (p_1 : \pi_1, \dots, p_j : \pi_j) \{c\}$ 642 $c ::= c_1 \cdot c_2 \mid p_d = p_s \mid x := \mathbf{new} \ C$ 643 $\mid x := \mathbf{invoke} \ x \langle T \rangle (p_1, \dots, p_j)$ 644 $T ::= t \mid T + n \quad \pi ::= [T_1, T_2]$ 645 $\tau ::= \forall \langle t : n \rangle (p_1 : \pi_1, \dots, p_j : \pi_j)$ 646 647 648 (a) Abstract syntax 649 </pre>	<pre> $\llbracket c \rrbracket : \mathcal{L} \rightarrow \mathcal{L} \quad \mathcal{L} : \mathcal{T} \rightarrow \mathcal{R} \times \mathcal{W}$ $\llbracket p_d = p_s \rrbracket (L) = \text{map}(\lambda(R, W). \text{if } p_s \in W$ $\text{then } (R\{p_s/p_d\}, W) \text{ else } (R, W), L)$ $\llbracket c_1 \cdot c_2 \rrbracket (L) = \llbracket c_1 \rrbracket (L) \cup \llbracket c_2 \rrbracket (L)$ </pre> <p style="text-align: center;">(b) Log-transformer semantics</p> $\frac{\Delta, \Lambda_1, \Gamma \vdash c_1 \dashv \Lambda'_1, \Gamma_1 \quad \Delta, \Lambda_2, \Gamma \vdash c_2 \dashv \Lambda'_2, \Gamma_2}{\Delta, \Lambda_1 * \Lambda_2, \Gamma \vdash c_1 \cdot c_2 \dashv \Lambda'_1 * \Lambda'_2, \Gamma_1 \cup \Gamma_2}$ <p style="text-align: center;">(c) COMPOSITION judgement</p>
--	--

Fig. 6. Formal semantics of Filament where command is defined as a *log-transformer*. Typing judgements track the active timeline of an instance and ensure they are used in a disjoint manner.

Definition 5.1 (Phantom Check). A phantom event G is used correctly if:

- (1) It is not used to share any instances.
- (2) It is only used to invoke subcomponents that use phantom events.

First, resource sharing is disallowed because any pipeline that shares an instance must use some signal to trigger an internal FSM and track which use of the instance is currently active. Second, a phantom event is only available at the type-level and cannot be reified since there is no interface port. Therefore, only components that use phantom events can be invoked with a phantom event.

Filament defines two state primitives: a *register* and a *delay* component.

```

665 component Register<G: L-(G+1), L: 1>(
666   @interface[G] en: 1, @[G, G+1] in: 32
667 ) -> (@[G+1, L] out: 32) where L > G+1;

```

```

665 component Delay<G: 1>(
666   @[G, G+1] in: 32
667 ) -> (@[G+1, G+2] out: 32);

```

As the type signatures denote, the difference is that a register can hold onto a value for an arbitrary amount of time while a delay can only hold onto a value for a single cycle. Because the delay component is continuously accepting inputs, it can only hold onto a value for a single cycle. In contrast, the register can use the en signal to hold onto a value for an arbitrary amount of time.

Compilation. The compiler does not instantiate FSMs or synthesize guards for invocations triggered using phantom events. Since Phantom Check ensures that all subcomponent themselves do not have an interface port, the compiler does not have to generate assignments for them. Filament generated code for continuous pipelines matches expert-written code.

6 FORMALIZATION

Figure 6a presents a simplified syntax for Filament: all components can be parameterized using exactly one constraint there cannot specify any ordering constraints between events. Since Filament disallows any form of event interaction in user-level components, multi-event user-level components are not fundamentally more expressive. Multi-event external components are more expressive but not supported in our formalism. A Filament program (\mathcal{P}) is a sequence of components which define a signature and a body in terms of commands: composition, connection, instantiation, and invocation.

6.1 Semantics

Figure 6b presents Filament’s semantics which is defined as functions over logs (\mathcal{L}). A log maps events (\mathcal{T}) to a set of ports that are read from (\mathcal{R}) and a multiset of ports that are written to (\mathcal{W}). Intuitively, a log captures all the reads and writes performed during every cycle of a component’s execution. We track the multiset of writes to capture conflicts—if there are multiple writes to the same port in the same cycle, then the program has a resource conflict.

Concrete logs are generated by the semantics of component definitions (Appendix A.2) while commands simply transform them. For example, a port connection forwards the value from the source port p_s to the destination port p_d . We model this by substituting all occurrences of p_d to p_s in the read set \mathcal{R} when p_s is defined in the write set \mathcal{W} and mapping it over all defined events in the log. Composition reflects the parallel nature of hardware—it simply unions the two logs together. Write conflicts can appear due to composition. The semantics of a program is the log generated by executing a distinguished main component with the empty log.² Using the log-based semantics, we can formalize the well-formedness (Section 4.2) and safe pipelining (Section 4.4) constraints of the type system.

Definition 6.1 (Well-Formedness). A component M is well-formed if and only if its log is well-formed. A log L is well-formed if and only if, for each event:

- There are no conflicting writes: $W_s = W$ where W_s is the deduplicated set of writes.
- Reads are a subset of writes: $R \subseteq W_s$

Definition 6.2 (Safe Pipelining). If a component M has an event T with delay d , and $\llbracket M \rrbracket_G$ represents its log where T is replaced with the event G , then M is safely pipelined if and only if all logs L_n are well-formed: $L_n = \forall n \geq d \ L_n = \llbracket M \rrbracket_T \cup \llbracket M \rrbracket_{T+n}$

6.2 Type System

Filament implements a type system inspired by separation logic [39] to enforce the well-formedness and safe pipelining constraints. Our presentation focuses on the specific typing judgement that ensures that there are no conflicting uses of an instance. Appendix A.3 provides the full type system. At a high level, our typing judgement for composition (Figure 6c) mirrors the parallel composition rule used in concurrent separation logic [7]—the two commands are checked under two disjoint resource contexts. Our insight is adapting the definition of separating split to timelines of instances and ensuring that instance reuse does not conflict. The typing judgements have the form: $\Delta; \Lambda; \Gamma \vdash c \dashv \Lambda'; \Gamma'$. Γ is the standard type environment, Δ tracks each event’s delay, and Λ is the *resource context*.

Resource contexts and separating split. Λ is the resource context and track the availability of each instance in the form of an interval (π). After instantiation, each instance is available in the interval $[0, \infty)$. The invocation rule (not shown) checks that, for an instance’s event with a delay d , the instance is available in the interval $[G, G+d)$ where G is the scheduling event. The composition rule (Figure 6c) splits the resource context before checking the two commands:

$$\Lambda = \Lambda_1 * \Lambda_2 \text{ iff } \forall (I : \pi) \in \Lambda \Rightarrow \exists \pi_1, \pi_2. (I : \pi_1) \in \Lambda_1 \wedge (I : \pi_2) \in \Lambda_2 \wedge \pi_1 \cap \pi_2 = \emptyset \wedge \pi_1 \cup \pi_2 = \pi$$

A valid split is one where the resulting contexts have disjoint intervals for each instance and the union of the intervals is the original interval. By using this definition of split, Filament ensures that invocations reuse instances in a non-conflicting manner. The remaining typing judgements (not shown) encode the remaining constraints needed to enforce well-formedness and safe pipelining. Using it, we prove (Appendix A.4) the following type soundness theorem:

²The full semantics is provided in Appendix A.2 in the supplementary material.

THEOREM 6.3. *If $\Delta; \Lambda; \Gamma \vdash c \dashv \Lambda'; \Gamma'$ then $\llbracket c \rrbracket$ preserves log well-formedness (Definition 6.1).*

7 EVALUATION

We evaluate Filament’s ability to efficiently express a number of accelerator designs and to express the interfaces generated by state-of-the-art accelerator generators. Our evaluation answers the following questions:

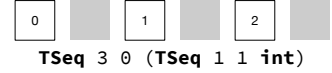
- (1) Can Filament express the interfaces generated by state-of-the-art accelerator generators and integrate with existing tools?
- (2) Can Filament be used to generate efficient accelerators?

Implementation. The Filament is implemented using a pass-based compiler in 5426 lines of Rust, 341 lines of Verilog for the standard library primitives, and the latest version of the Calyx compiler [34] to generate Verilog. All benchmarks compile in under a second.

7.1 Expressivity Evaluation

Our expressivity evaluation focuses on giving type signatures to modules generated from Aetherling [17], a domain-specific language (DSL) for generating streaming image-processing kernels.

Aetherling’s space-time types. Aetherling [17] is a functional, dataflow DSL that generates statically-scheduled, streaming accelerators for image processing tasks. Aetherling’s “space-time” types



enable users to express the shape of the data stream as a sequence of valid and invalid signals. For example, the type `TSeq 1 1` denotes that there will be a stream with one valid element followed by one invalid element. Nesting these types allows users to express more complex shapes: `TSeq 3 0 (TSeq 1 1)` denotes that there will be three valid elements, with no invalid values, each of which has a shape described by `TSeq 1 1`. In our case study, we import 14 designs implementing two kernels: `conv2d` and `sharpen`. Aetherling’s evaluation studies 7 design points for each kernel with different resource-throughput trade-offs. We demonstrate that Filament can express type for all designs and, in the process, finds several bugs in the generated interfaces.

Cycle accurate harness. We implemented a generic, *cycle accurate* harness to test Filament programs. The harness extracts availability intervals for inputs and the initiation interval from the top-level component. It then executes a set of inputs, provided as a JSON array, in a pipelined manner and checks whether the expected outputs are generated during availability interval of the output. The design of this generic harness is reliant on a Filament-like system to document the timing behavior of modules; without Filament, a user would have to manually extract this information from the Verilog code.

Methodology. We compile each Aetherling design to Verilog and use Aetherling’s command line interface to extract the design’s latency information. Each benchmark has five fully-utilized designs, which can accept new inputs every cycle, and two underutilized designs which produce 1/3 and 1/9 pixels per clock cycle and accept new inputs every 3 and 9 cycles. We give each design a type signature and validate its outputs. For designs with mismatched outputs, we change the latency till we get the right answer.

Latency. Table 1 reports the latencies as provided by Aetherling’s command line interface and those that we found to generate correct outputs with Filament’s cycle accurate test harness. Of the 14 designs, Aetherling reports incorrect latencies for 5 designs.

Throughput	Reported	Actual	Throughput	Reported	Actual
16	6	7	16	7	7
8	6	6	8	7	7
4	6	6	4	7	7
2	6	6	2	7	7
1	7	7	1	8	8
1/3	10	12	1/3	11	13
1/9	16	21	1/9	17	20

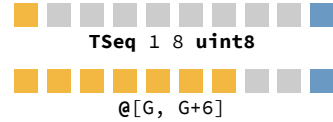
(a) Reported latencies for conv2d

(b) Reported latencies for sharpen

Table 1. Latencies of Aetherling Designs. Highlighted latencies are reported incorrectly by Aetherling.

Underutilized designs. Aetherling explores the utility of *underutilized* designs which produce less than one pixel per clock cycle. Aetherling’s compiler optimizes such designs by sharing compute resources. An Aetherling design that produces 1/9 pixels per clock has the type `TSeq 1 8 uint8` which states that there will be 1 valid datum followed by 8 invalid ones. The type indicates that the design generated by Aetherling should only use its input in the first cycle since the data provided in the next cycles is invalid. However, this interface is incorrect.

```
component Conv2d<G: 9>(
  @[G, G+6] I: 8,
) -> (@[G+21, G+22] O: 8);
```



The Filament type, which reflects the actual interface needed to correctly execute the module, requires the design to hold its input signal for six cycles, i.e., the data element must be valid for six cycles instead of just one; the Aetherling implementation breaks its own interface. The Aetherling test harness does not catch this bug because it always asserts all inputs for 9 cycles. In contrast, Filament’s test harness only asserts the input signal for as long as the corresponding availability interval specifies. Finally, the delay for the phantom event G encodes the fact the design can process a new input every 9 cycles. This illustrates the subtlety of specifying time-sensitive interfaces which accurately describe signal availability and pipelining.

Other designs. We also import designs generated from PipelineC [25], an open-source high-level synthesis compiler that transforms a C-like language into Verilog (Appendix B.2). Providing type signature for these was straightforward since PipelineC always fully pipelines designs and prints out the design’s latency on the command line.

7.2 Accelerator Design with Filament

We study Filament’s efficacy in generating efficient designs and reusing components generated from other language by implementing a two-dimensional convolution in Filament. We build two Filament-based designs and compare them to the Aetherling-generated conv design.

Architecture. Our implementation is directly inspired by the structure of the Aetherling implementation of conv2d that outputs 1 pixel per clock cycle. The design uses a 3×3 filter over a 4×4 matrix. The `Stencil` module (Figure 7a) implements a line buffer to save the last 11 values and outputs 9 values corresponding to the filter start index. The `Conv2D` kernel takes 9 values as inputs and produces an output corresponding to the result of the convolution.

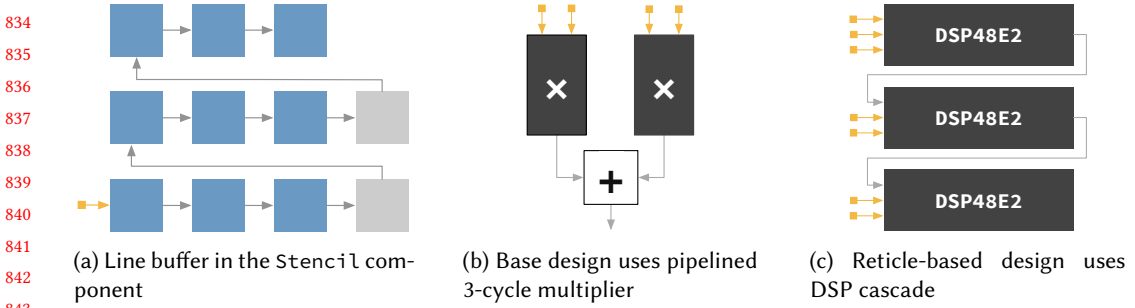


Fig. 7. Components used in the design of Filament-based conv2d convolution. The stencil component provides the last three inputs and is either connected to the naive multiplier or a Reticle-generated DSP cascade.

Stream primitives in Filament. To implement line buffers, we implement a new `Prev` component which outputs the last value stored in it.³ The Verilog implementation of `Prev` is simply a register but Filament gives it a different type signature—it allows access to the output in the same cycle when the input is provided which corresponds to reading the previous value in the register. The component uses a compile-time parameter `SAFE` to indicate whether the first read produces an undefined value. We also define a `ContPrev` component which is similar to a `Prev` component but uses a phantom event and can therefore be used in continuous pipelines (Section 5.4). The stencil component (Figure 7a) as a sequence of `Prev` components.

```
component Prev[SAFE]<G: 1>(
  @interface[G] en: 1,
  @[G, G+1] in: 32,
) -> (@[G, G+1] prev: 32);
```

Design 1: Pipelined multipliers. The base Conv2D kernel uses fully pipelined multipliers with a three cycle latency and combinational adders. The multipliers do not have any associated Verilog implementation—they are implemented using Xilinx’s LogiCORE multiplier generator [23]. However, Filament makes it easy to interface with them by providing a type-safe extern wrapper (Section 3.6).

Design 2: Integrating with Reticle. Our second design uses a dot-product unit generated using Reticle [41], a low-level language for programming FPGAs. Figure 7c shows the architecture Reticle generates to make use of *DSP cascading* which efficiently utilizes resources present on an FPGA. DSP cascading explicitly instantiates low-level FPGA primitives and connects them together to implement the computation: $y = c + \sum_{i=0}^3 a_i \times b_i$. Unlike standard compilation flows which rely on the synthesis tool to infer DSP usage from *behavioral* descriptions, Reticle generates *structural* descriptions that predictably map onto DSPs. We provide a type signature for the Reticle design which indicates that the inputs must be provided in a staggered manner. Note that this is not implementation details leaking through—a DSP cascade that starts a new computation every cycle needs to either register all its inputs or provide them in a staggered manner.

```
component Tdot<G: 1>(
  clk: 1, reset: 1,
  @[G, G+1] a0: 8,
  @[G, G+1] b0: 8,
  @[G+1, G+2] a1: 8,
  @[G+1, G+2] b1: 8,
  @[G+2, G+3] a2: 8,
  @[G+2, G+3] b2: 8,
  @[G+2, G+3] c: 8,
) -> (@[G+5, G+6] y: 8)
```

Evaluation methodology. We validate the correctness of all the designs using our timing-accurate test harness and compare the area and latency of the designs. For each design, we increase the target frequency till we reach worst negative slack of less than $0.1ns$ and synthesize them using Vivado v2020.2. Each design has a throughput of 1 pixel per clock cycle.

³`prev` is a common operator in dataflow and functional reactive languages.

883 **Summary.** Table 2 shows the re-
 884 sults of the comparison: the Fila-
 885 ment design can be synthesized at
 886 a higher frequency and uses fewer
 887 resources than the Aetherling de-
 888 sign. This is because Filament can
 889 safely and directly use low-level
 890 implementation modules which can be directly compiled into a safe and efficient design. In con-
 891 trast, the Aetherling compiler has to generate extra logic when bridging the gap between its high-
 892 level language and low-level circuits. The Reticle-based design uses an order of magnitude fewer
 893 logic resources than the base Filament design or the Aetherling design. This is because unlike
 894 Aetherling, Reticle generates low-level *structural* Verilog which can predictably map onto DSP
 895 resources. This demonstrates the utility of Filament as both an integration and design language—
 896 designs in Filament can use low-level hardware modules safely and compose complex modules
 897 generated from other languages. It also reveals another use case for Filament: instead of directly
 898 generating Verilog, Aetherling-like languages can generate Filament programs and enable perfor-
 899 mance engineers to optimize the designs further and remove abstraction overheads.

900 **Other designs.** We additionally implemented other designs to evaluate Filament (Appendix B.1):
 901 (1) floating-point addition and multiplication and (2) matrix-multiply systolic array [28].

903 8 RELATED WORK

904 **Dataflow languages.** Reactive dataflow programming languages [6, 9, 19] provide stream oper-
 905 ators scheduled using logical time steps. Software dataflow languages [38, 40] provide high-level,
 906 declarative operations that can target multiple backends like CPUs and GPUs. Compiling these lan-
 907 guages to hardware requires complex transformations [3, 37]. Filament is lower-level and easier to
 908 compile since it directly reasons about hardware modules and is appropriate as a target language
 909 for hardware generators for these languages.

911 **Accelerator design languages.** Accelerator design languages [14, 17, 20, 21, 26, 43] provide high-
 912 level abstractions to design hardware accelerators. Filament is a low-level HDL that provides a
 913 type system to directly interface with hardware modules and is appropriate both as an integration
 914 language and as a target language for compilers for ADLs.

915 **Embedded HDLs.** Embedded HDLs [1, 2, 4, 13, 24, 29] use software host languages for metapro-
 916 gramming. Most eHDLs simply use the host language’s type system to ensure simple properties
 917 like port width match and signedness. Filament’s type system focuses on expressing structural and
 918 temporal properties of the hardware itself. Rule-based HDLs [5, 35] use *guarded atomic actions* to
 919 provide transactional semantics for hardware specification. To preserve their high-level semantics,
 920 the compiler must generate complex scheduling logic that dynamically aborts conflicting rules. In
 921 contrast, Filament specifies program schedules using invocations which are checked at compile
 922 time and predictably map to efficient, pipelined hardware.

924 **Type systems.** Dahlia [33] is a C-like language that uses a substructural type system to ensure that
 925 high-level programs do not violate hardware constraints. Dahlia’s affine reasoning can be encoded
 926 in Filament’s type system. Kami [10] is a proof-assisted framework for designing hardware. Kami
 927 can be used to prove full functional correctness of hardware modules but requires users to write
 928 proofs while Filament’s type system is focused on timing properties and is automatic. Cordial [16]
 929 is a type system inspired by *session types* to reason about latency-insensitive hardware protocols
 930 while Ghica [18] presents game-semantics-inspired type system to model interfaces; both tools do
 931

Name	LUTs	DSPs	Registers	Freq. (MHz)
Aetherling	104	10	78	769.2
Filament	128	<u>9</u>	<u>11</u>	<u>833.3</u>
Filament Reticle	<u>14</u>	<u>9</u>	20	645.1

Table 2. Resource usage and frequency of conv2d designs. Best values highlighted.

not reason about pipelining. Filament, in contrast, focuses on latency-sensitive interfaces in presence of pipelining. Pi-Ware [36] uses dependent types to ensure low-level circuit properties such as ensuring all ports are connected, and wire sorts [11] check whether module composition can create *combinational loops*, both of which are orthogonal to Filament’s guarantees.

Model checking. Model checking [12] is a popular technique to verify hardware designs. SystemVerilog Assertions [42] provide a linear temporal logic (LTL) based specification language to verify properties of hardware designs [8, 30]. Such systems provide whole program guarantees and can prove more general timing properties than Filament. Filament focuses on providing compositional guarantees and interface specifications. Filament signatures can be potentially compiled to LTL specification and checked by aforementioned tools.

9 FUTURE WORK

Filament represents a first step toward a new class of type systems that can reason about structural and temporal properties of hardware designs.

Dynamic pipelines. While static pipelines encompass a large and interesting set of hardware designs [17, 20, 21, 37], dynamic pipelines are important for expressing complex designs. Filament’s characterization of pipelining constraints (Section 4) provides a formal foundation to reason about dynamic pipelines. A possible solution is allowing modules to *generate* events in addition to consuming them; this would correspond to an existential quantification over the event type as opposed to the universal quantification in the current system. The challenge is showing that this extension produces well-formed dynamic pipelines.

Delay polymorphic and higher-order Filament. HDLs allow programmers to implement modules with parameteric latencies. For example, a shift register implementation can use a parameter to specify the number of stages which affects when the final output is available. Similarly, an iterative divider can have a latency dependent on the width of its inputs. Such an extension is also necessary to support higher-order modules since the latency of a module depends on the latency of the module it is parameterized by.

Type-preserving compilation. High-level synthesis, the process of compiling imperative languages like C to hardware, is notoriously buggy. While previous efforts have focused on building fully verified compilers [22, 27], they fail to capture common pipelining optimizations such as modulo scheduling [15]. We believe that Filament can be used as a typed intermediate language for building a type-preserving compiler [31] that rules out pipeline-related HLS compilation bugs.

10 CONCLUSION

Unlocking the true potential of reusable hardware requires detailed understanding of the structure and timing of the implementation. Filament exposes this knowledge through interfaces and enables any user to fearlessly build high-performance hardware and share it with the world.

REFERENCES

- [1] C. Baaij, M. Kooijman, J. Kuper, A. Boeijink, and M. Gerards. 2010. C_laSH: Structural Descriptions of Synchronous Hardware Using Haskell. In *Euromicro Conference on Digital System Design: Architectures, Methods and Tools*.
- [2] Jonathan Bachrach, Huy Vo, Brian Richards, Yunsup Lee, Andrew Waterman, Rimas Avizienis, John Wawrzynek, and Krste Asanović. 2012. Chisel: constructing hardware in a Scala embedded language. In *Design Automation Conference (DAC)*.
- [3] Gérard Berry. 1992. A hardware implementation of pure Esterel. *Sadhana* (1992).
- [4] Per Bjesse, Koen Claessen, Mary Sheeran, and Satnam Singh. 1998. Lava: hardware design in Haskell. *ACM SIGPLAN Notices* (1998).

- 981 [5] Thomas Bourgeat, Clément Pit-Claudel, and Adam Chlipala. 2020. The essence of Bluespec: a core language for rule-
982 based hardware design. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*.
- 983 [6] Frédéric Boussinot and Robert De Simone. 1991. The ESTEREL language. *Proc. IEEE* (1991).
- 984 [7] Stephen Brookes. 2004. A semantics for concurrent separation logic. In *International Conference on Concurrency*
985 *Theory*. Springer.
- 986 [8] Cadence Inc. 2022. *Jasper Gold FPV App*. Retrieved October 15, 2022 from [https://www.cadence.com/en_US/home/
987 tools/system-design-and-verification/formal-and-static-verification/jasper-gold-verification-platform/formal-
988 property-verification-app.html](https://www.cadence.com/en_US/home/tools/system-design-and-verification/formal-and-static-verification/jasper-gold-verification-platform/formal-property-verification-app.html)
- 989 [9] Paul Caspi, Grégoire Hamon, and Marc Pouzet. 2008. Synchronous functional programming: The Lucid Synchronic
990 experiment. *Real-Time Systems: Description and Verification Techniques: Theory and Tools. Hermes* (2008).
- 991 [10] Joonwon Choi, Muralidaran Vijayaraghavan, Benjamin Sherman, Adam Chlipala, and Arvind. 2017. Kami: A Platform
992 for High-Level Parametric Hardware Specification and Its Modular Verification. (2017).
- 993 [11] Michael Christensen, Timothy Sherwood, Jonathan Balkind, and Ben Hardekopf. 2021. Wire sorts: A language ab-
994 straction for safe hardware composition. In *ACM SIGPLAN Conference on Programming Language Design and Imple-
995 mentation (PLDI)*.
- 996 [12] Edmund M Clarke. 1997. Model checking. In *International Conference on Foundations of Software Technology and*
997 *Theoretical Computer Science*. Springer, 54–56.
- 998 [13] J. Clow, G. Tzimpragos, D. Dangwal, S. Guo, J. McMahan, and T. Sherwood. 2017. A Pythonic approach for rapid
999 hardware prototyping and instrumentation. In *International Conference on Field-Programmable Logic and Applications*
1000 *(FPL)*.
- 1001 [14] Jason Cong and Jie Wang. 2018. PolySA: Polyhedral-based systolic array auto-compilation. In *IEEE/ACM International*
1002 *Conference on Computer-Aided Design (ICCAD)*. IEEE.
- 1003 [15] J. Cong and Zhiru Zhang. 2006. An efficient and versatile scheduling algorithm based on SDC formulation. In *Design*
1004 *Automation Conference (DAC)*.
- 1005 [16] Jan de Muijnck-Hughes and Wim Vanderbauwhede. 2019. A Typing Discipline for Hardware Interfaces. In *European*
1006 *Conference on Object-Oriented Programming (ECOOP)*.
- 1007 [17] David Durst, Matthew Feldman, Dillon Huff, David Akeley, Ross Daly, Gilbert Louis Bernstein, Marco Patrignani,
1008 Kayvon Fatahalian, and Pat Hanrahan. 2020. Type-Directed Scheduling of Streaming Accelerators. In *ACM SIGPLAN*
1009 *Conference on Programming Language Design and Implementation (PLDI)*.
- 1010 [18] Dan R Ghica. 2009. Function interface models for hardware compilation: Types, signatures, protocols. *arXiv preprint*
1011 *arXiv:0907.0749* (2009).
- 1012 [19] Nicholas Halbwegs, Paul Caspi, Pascal Raymond, and Daniel Pilaud. 1991. The synchronous data flow programming
1013 language LUSTRE. *Proc. IEEE* (1991).
- 1014 [20] James Hegarty, John Brunhaver, Zachary DeVito, Jonathan Ragan-Kelley, Noy Cohen, Steven Bell, Artem Vasilyev,
1015 Mark Horowitz, and Pat Hanrahan. 2014. Darkroom: Compiling high-level image processing code into hardware
1016 pipelines. *ACM Transactions on Graphics*.
- 1017 [21] James Hegarty, Ross Daly, Zachary DeVito, Jonathan Ragan-Kelley, Mark Horowitz, and Pat Hanrahan. 2016. Rigel:
1018 Flexible multi-rate image processing hardware. *ACM Transactions on Graphics*.
- 1019 [22] Yann Herklotz, James D. Pollard, Nadesh Ramanathan, and John Wickerson. 2021. Formal Verification of High-Level
1020 Synthesis. In *ACM SIGPLAN Conference on Object Oriented Programming, Systems, Languages and Applications (OOP-
1021 SLA)*.
- 1022 [23] AMD Inc. 2022. *Xilinx LogiCORE IP Multiplier v11.2*. Retrieved October 27, 2022 from [https://docs.xilinx.com/v/u/en-
1023 US/mult_gen_ds255](https://docs.xilinx.com/v/u/en-US/mult_gen_ds255)
- 1024 [24] Jane Street. 2022. *HardCaml: Register Transfer Level Hardware Design in OCaml*. Retrieved October 15, 2022 from
1025 <https://github.com/janestreet/hardcaml>
- 1026 [25] Julian Kemmerer. 2022. *PipelineC*. Retrieved October 15, 2022 from <https://github.com/JulianKemmerer/PipelineC>
- 1027 [26] David Koeplinger, Matthew Feldman, Raghu Prabhakar, Yaqi Zhang, Stefan Hadjis, Ruben Fiszal, Tian Zhao, Luigi
1028 Nardi, Ardavan Pedram, Christos Kozyrakis, and Kunle Olukotun. 2018. Spatial: A language and compiler for appli-
1029 cation accelerators. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*.
- [27] Sudipta Kundu, Sorin Lerner, and Rajesh Gupta. 2008. Validating high-level synthesis. In *International Conference on*
Computer-Aided Verification (CAV).
- [28] Hsiang-Tsung Kung. 1982. Why systolic architectures? *IEEE Computer* (1982).
- [29] Derek Lockhart, Gary Zibrat, and Christopher Batten. 2014. PyMTL: A Unified Framework for Vertically Integrated
Computer Architecture Research. In *IEEE/ACM International Symposium on Microarchitecture (MICRO)*.
- [30] Cristian Mattarei, Makai Mann, Clark Barrett, Ross G. Daly, Dillon Huff, and Pat Hanrahan. 2018. CoSA: Integrated
Verification for Agile Hardware Design. In *Formal Methods in Computer-Aided Design (FMCAD)*.

- 1030 [31] Greg Morrisett, David Walker, Karl Crary, and Neal Glew. 1999. From System F to typed assembly language. *ACM*
1031 *Transactions on Programming Languages and Systems (TOPLAS)* (1999).
- 1032 [32] Kevin E. Murray and Vaughn Betz. 2014. Quantifying the Cost and Benefit of Latency Insensitive Communication on
1033 FPGAs. In *International Symposium on Field-Programmable Gate Arrays (FPGA)*.
- 1034 [33] Rachit Nigam, Sachille Atapattu, Samuel Thomas, Zhijing Li, Theodore Bauer, Yuwei Ye, Apurva Koti, Adrian Samp-
1035 son, and Zhiru Zhang. 2020. Predictable Accelerator Design with Time-Sensitive Affine Types. In *ACM SIGPLAN*
1036 *Conference on Programming Language Design and Implementation (PLDI)*.
- 1037 [34] Rachit Nigam, Samuel Thomas, Zhijing Li, and Adrian Sampson. 2021. A compiler infrastructure for accelerator gen-
1038 erators. In *ACM International Conference on Architectural Support for Programming Languages and Operating Systems*
1039 *(ASPLOS)*.
- 1040 [35] Rishiyur Nikhil. 2004. Bluespec System Verilog: Efficient, correct RTL from high level specifications. In *Conference*
1041 *on Formal Methods and Models for Co-Design (MEMOCODE)*.
- 1042 [36] JP Pizani Flor et al. 2014. *Pi-Ware: An Embedded Hardware Description Language using Dependent Types*. Master's
1043 thesis.
- 1044 [37] Jing Pu, Steven Bell, Xuan Yang, Jeff Setter, Stephen Richardson, Jonathan Ragan-Kelley, and Mark Horowitz. 2017.
1045 Programming heterogeneous systems from an image processing DSL. *ACM Transactions on Architecture and Code*
1046 *Optimization (TACO)*.
- 1047 [38] Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman P. Amarasinghe.
1048 2013. Halide: A language and compiler for optimizing parallelism, locality, and recomputation in image processing
1049 pipelines. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*.
- 1050 [39] John C Reynolds. 2002. Separation logic: A logic for shared mutable data structures. In *IEEE Symposium on Logic in*
1051 *Computer Science*. IEEE.
- 1052 [40] William Thies, Michal Karczmarek, and Saman Amarasinghe. 2002. StreamIt: A language for streaming applications.
1053 In *International Conference on Compiler Construction*. Springer.
- 1054 [41] Luis Vega, Joseph McMahan, Adrian Sampson, Dan Grossman, and Luis Ceze. 2021. Reticle: a virtual machine for
1055 programming modern FPGAs. In *ACM SIGPLAN Conference on Programming Language Design and Implementation*
1056 *(PLDI)*.
- 1057 [42] Srikanth Vijayaraghavan and Meyyappan Ramanathan. 2005. *A practical guide for SystemVerilog assertions*. Springer
1058 Science & Business Media.
- 1059 [43] Xilinx Inc. 2021. *Vivado Design Suite User Guide: High-Level Synthesis. UG902 (v2017.2) June 7, 2017*. Retrieved January
1060 16, 2021 from https://www.xilinx.com/support/documentation/sw_manuals/xilinx2017_2/ug902-vivado-high-level-synthesis.pdf
- 1061 [44] Zhiru Zhang, Yiping Fan, Wei Jiang, Guoling Han, Changqi Yang, and Jason Cong. 2008. AutoPilot: A platform-based
1062 ESL synthesis system. In *High-Level Synthesis*. 99–112.
- 1063
- 1064
- 1065
- 1066
- 1067
- 1068
- 1069
- 1070
- 1071
- 1072
- 1073
- 1074
- 1075
- 1076
- 1077
- 1078

1079	$\llbracket c \rrbracket$	$: \mathcal{L} \rightarrow \mathcal{L} \quad (\mathcal{L} : \mathcal{T} \rightarrow \mathcal{R} \times \mathcal{W})$
1080		
1081	$\llbracket x := \mathbf{new} C \rrbracket$	$= \text{id}$
1082	$\llbracket p_d = p_s \rrbracket$	$= \lambda(R, W). \text{if } p_s \in W \text{ then } (R\{p_s/p_d\}, W) \text{ else } (R, W)$
1083	$\llbracket c_1 \cdot c_2 \rrbracket (L)$	$= \llbracket c_1 \rrbracket (L) \cup \llbracket c_2 \rrbracket (L)$
1084		
1085	$\llbracket x_1 := \mathbf{invoke} x_2 \langle T' \rangle (q_1, \dots, q_m) \rrbracket (L)$	$= \llbracket m \rrbracket \cup \text{connects}(x_1, [q_1, \dots, q_m]) \quad \text{where } \text{bind}(x_2) = m$
1086		
1087	$\llbracket M \rrbracket$	$: \mathcal{P}$
1088		
1089	$\llbracket \mathbf{def} C \langle t : n \rangle (i_0 : \pi_0, \dots, o_0 : \pi_i, \dots) \{c\} \rrbracket$	$= \{\pi_0 \mapsto i_0, \dots\} \times \{\pi_i \mapsto o_0, \dots\}$

Fig. 8. Log-transformer semantics for Filament’s core language. Each command produces a log (\mathcal{L}) which maps events (\mathcal{T}) to multisets of reads (\mathcal{R}) and writes (\mathcal{W}). Component definitions produce partial logs (\mathcal{P}) produce a log by mapping availabilities of inputs to reads and availabilities of outputs to writes.

A FULL FORMALISM

A.1 Syntax

We used a simplified version of Filament for our formalization: all components can be parameterized using exactly one constraint there cannot specify any ordering constraints between events. Neither simplification loses generality because user-level components with multiple events cannot define any form of interaction between them—they are functionally equivalent to multiple components with disjoint events.

A Filament program is a sequence of components M each of which encapsulates the structure and schedule of a pipeline. Commands c include composition, port connections (Section 3.5), component instantiation (Section 3.3), and invocations (Section 3.4). Component are parameterized using exactly one event and invocations allow scheduling using one event.

$x, C \in \text{vars} \quad t \in \text{events}$

$p, q \in \text{ports}$

$M ::=$

$\mathbf{def} C \langle t : n \rangle (p_1 : \pi_1, \dots, p_j : \pi_j) \{c\}$

$c ::= c_1 \cdot c_2 \mid p_d = p_s \mid x := \mathbf{new} C$

$\mid x := \mathbf{invoke} x \langle T \rangle (p_1, \dots, p_j)$

$T ::= t \mid T + n \quad \pi ::= [T_1, T_2]$

A.2 Semantics

Figure 8 provides a denotation of Filament programs as a log (\mathcal{L}) of events to multisets of port reads (\mathcal{R}) and port writes (\mathcal{W}). Since components are allowed to use exactly one event, say T , the log maps events such as $T, T + 1$, etc. to reads and writes to ports defined by the subcomponents. The semantics of commands is given as log-transformers: they take in a log and produce a new log. Instantiation does not affect the logs and composition is the union of the logs produced by the two commands. Connections transform the log by adding the LHS port to the set of writes for event where the RHS port is defined. Invocations produce a partial log using the semantics of the instance and unions them with a log where the inputs are interpreted are treated as connections. The connect metafunction coverts the arguments to an invocation into connections and generates a log by interpreting them using the denotation over commands. The semantics of a program is defined by the log produced by a distinguished top-level component.

1128 **Components and partial logs.** The semantics of a component is derived from its signature: the
 1129 input ports are added to the reads (\mathcal{R}) and the output ports to the writes (\mathcal{W}) for each event
 1130 contained in their corresponding availability interval. This generates a *partial log* which represents
 1131 the requirements of the component signature, in the form of reads, and its guarantees, in the form
 1132 of writes.

1133 For example, a combinational adder and a sequential multiplier with a two-cycle latency produce
 1134 the following logs:

1135
 1136
$$\llbracket \text{def add}\langle G : 1 \rangle(l: [G, G + 1], r: [G, G + 1], \text{out}: [G, G + 1]) \rrbracket = G \rightarrow (\{l, r\}, \{\text{go}, \text{out}\})$$

 1137
$$\llbracket \text{def mul}\langle G : 2 \rangle(l: [G, G + 1], r: [G, G + 1], \text{out}: [G + 2, G + 3]) \rrbracket = G \rightarrow (\{l, r\}, \{\text{go}\})$$

 1138
$$G + 1 \rightarrow (\emptyset, \{\text{go}\})$$

 1139
$$G + 2 \rightarrow (\emptyset, \{\text{out}\})$$

 1140
 1141

1142 Note that the use of the instances is reflected through the writes their interface ports go (not shown
 1143 in the signature). The log indicates that the multiplier accepts new values every 2 cycles by writing
 1144 to the go port in both cycles G and $G + 1$. Because we track multisets of reads and writes, we can
 1145 track conflicting writes to the same port.

1146 Using these semantics, we can define the well-formedness constraint (Section 4.2) on logs:

1147
 1148 *Definition A.1 (Well-Formedness).* A log \mathcal{L} is well-formed if and only if for all events

- 1149 • There are no conflicting writes: $\mathcal{W}_s = \mathcal{W}$ where \mathcal{W}_s is the deduplicated set of writes.
- 1150 • Reads are a subset of writes for every event: $\mathcal{R} \subseteq \mathcal{W}_s$.

1151
 1152 While in real hardware, values are always available on a port or a wire, Filament's semantics only
 1153 track semantically valid values from a read. Usage of hardware resources is denoted by a write,
 1154 and it is physically impossible to write two values to a port; instead, the circuit uses a multiplexer
 1155 to select between the two values. Multiple uses of a resource silently corrupt the data.

1156 The safe pipelining constraints (Section 4.4) can be defined in terms of repeated execution of the
 1157 semantics of a program:

1158
 1159 *Definition A.2 (Safe Pipelining).* If a component M has an event T with delay d , and $\llbracket M \rrbracket_G$ rep-
 1160 represents its log where T is replaced with the event G , then M is safely pipelined if and only if all
 1161 logs L_n are well-formed: $L_n = \forall n \geq d L_n = \llbracket M \rrbracket_T \cup \llbracket M \rrbracket_{T+n}$

1164 A.3 Type System

1165 Our presentation focuses on Filament's substructural type system that is used to track non-conflicting
 1166 use resources as well as signal validity. We elide the description of features that track things such
 1167 as port widths which are standard.

1168 **Typing contexts.** The typing judgements use the following
 1169 typing contexts:

- 1170 • Γ tracks the types for components and instances and $\tau ::= \forall \langle t : n \rangle (p_1 : \pi_1, \dots, p_j : \pi_j)$
 availability of ports. $\Gamma ::= \cdot \mid \Gamma, C : \tau \mid \Gamma, p : \pi$
- 1171 • Δ tracks the delays associated with each event in the $\Delta ::= \cdot \mid \Delta, t : n$
- 1172 • Λ is the *timeline context* and tracks the availability of $\Lambda ::= \cdot \mid \Lambda, I : \pi$
 each instance.

1176

The type context (Γ) and timeline context (Λ) store timelines for ports and instances respectively. Timelines for ports are *reusable* since reading a port does not consume it during that cycle. However, the timeline of an instance is *consumed* when it is used in a cycle. Because of this, timeline contexts also provide a *separating union* inspired by separation logic [39].

Splitting timelines. A valid separating split of a timeline context $\Lambda = \Lambda_1 * \Lambda_2$ if and only if both Λ_1 and Λ_2 bind all the same instances and for each instance, the timelines are disjoint. Formally:

$$\Lambda = \Lambda_1 * \Lambda_2 \text{ iff } \forall (I : \pi) \in \Lambda \Rightarrow \exists \pi_1, \pi_2. (I : \pi_1) \in \Lambda_1 \wedge (I : \pi_2) \in \Lambda_2 \wedge \pi_1 \cap \pi_2 = \emptyset \wedge \pi_1 \cup \pi_2 = \pi$$

Instantiating components. Instantiating a module binds the signature of the component to the instance and make the resource available throughout the timeline of the program, denoted by $[0, \infty)$.

$$\frac{\Gamma(C) = \tau \quad \Gamma' = \Gamma, I : \tau \quad \Lambda' = \Lambda, I : [0, \infty)}{\Delta, \Lambda, \Gamma \vdash I := \text{new } C \vdash \Lambda', \Gamma'}$$

Port connections. Connecting ports checks that the source port is available for at least as long as the destination port requires:

$$\frac{\Gamma(p_d) \subseteq \Gamma(p_s)}{\Delta, \Lambda, \Gamma \vdash p_d = p_s \vdash \Lambda', \Gamma'}$$

Splitting timelines with composition. The composition rule splits the timeline context using the separating split operator and checks the two commands. Note that that same type context Γ is used for both commands which means previously defined ports are available in both the commands:

$$\frac{\Delta, \Lambda_1, \Gamma \vdash c_1 \vdash \Lambda'_1, \Gamma_1 \quad \Delta, \Lambda_2, \Gamma \vdash c_2 \vdash \Lambda'_2, \Gamma_2}{\Delta, \Lambda_1 * \Lambda_2, \Gamma \vdash c_1 \cdot c_2 \vdash \Lambda'_1 * \Lambda'_2, \Gamma_1 \cup \Gamma_2}$$

Checking invocations. The invocation rule enforces well-formedness and safe-pipelining constraints and is therefore quite verbose. We separate out type checking of invocations into three sets of premises that logically reflect the properties presented in Section 4.

$$\frac{\text{valid reads} \quad \text{no conflicts} \quad \text{safe pipelining}}{\Delta, \Lambda, \Gamma \vdash x := \text{invoke}I\langle T \rangle(q_1, \dots, q_j) \vdash \Lambda, \Gamma''}$$

The first set of premises check that all the reads from all ports mentioned in an invocation are valid, i.e., they are available for at least as long as the instance's signature requires. Finally, invocations bind the availability of all the ports associated with the instance to the type context.

$$\begin{aligned} \Gamma(I) = \forall \langle t : n \rangle (p_1 : \pi_1, \dots, p_j : \pi_j) \quad \Gamma(q_1) = \pi_1, \dots, \Gamma(q_j) = \pi_j \\ \pi'_1 = \pi_1[t/T], \dots, \pi'_j = \pi_j[t/T] \quad \Gamma(p_1) \subseteq \pi'_1, \dots, \Gamma(p_j) \subseteq \pi'_j \\ \Gamma' = x. \{p_1 : \pi'_1, \dots, p_j : \pi'_j\} \quad \Gamma'' = \Gamma \cup \Gamma' \end{aligned}$$

The next set of premises ensure that the instance is available in the current timeline context. This ensures that there are no conflicting uses of the component anywhere else in the design.

$$\begin{aligned} \Gamma(I) = \forall \langle t : n \rangle (p_1 : \pi_1, \dots, p_j : \pi_j) \\ \Lambda(I) = \pi \quad [T, T + n] \subseteq \pi \end{aligned}$$

The composition rule is responsible for selecting a valid split to ensure that the above rule's constraints are satisfied. If there is no such split possible, then the program has conflicting uses of the instance.

A final set checks for the safety of pipelining an invocation (Section 4.4):

$$\begin{aligned} \Gamma(I) = \forall \langle t : n \rangle (p_1 : \pi_1, \dots, p_j : \pi_j) \\ \mathcal{E}(T) = t' \quad \Delta(t') \geq \Delta(t) \end{aligned}$$

1226 A.4 Type Soundness

1227 Our type system guarantees theorem focuses on the well-formedness property (Section 4.2). It
 1228 states that well-typed commands preserve well-formed logs.
 1229

1230 *Definition A.3 (Compatible Partial Log).* A partial log $\mathcal{P} = \mathcal{R}_p \times \mathcal{W}_p$ is *compatible* with a log
 1231 $\mathcal{L} = \mathcal{R}_l \times \mathcal{W}_l$ if the log $\mathcal{L}' = \mathcal{R}_p \times \mathcal{W}_l$ is well-formed.

1232 Partial logs are generated from the semantics of a component definition and are not well-formed.
 1233 We define a *compatibility* relation between partial logs and well-formed logs. Intuitively, a com-
 1234 patible log provides writes to make the reads of the partial log valid. A partial log is *not required*
 1235 to maintain the second well-formedness constraint that there are no conflicting writes.
 1236

1237 **LEMMA A.4 (DISJOINT WRITES).** *If $\Delta; \Lambda; \Gamma \vdash c_1; c_2 \dashv \Lambda'; \Gamma'$, then for logs $(\mathcal{R}_1, \mathcal{W}_1) = \llbracket c_1 \rrbracket (\mathcal{L})$
 1238 and $(\mathcal{R}_2, \mathcal{W}_2) = \llbracket c_2 \rrbracket (\mathcal{L})$, \mathcal{W}_1 and \mathcal{W}_2 are disjoint.*

1239 **PROOF.** The only possibility of conflicts come from writes to an instance's interface port. How-
 1240 ever, because the typing rule for composition requires that the two commands type check in dis-
 1241 joint timeline environments, it means that the writes to the interface ports are disjoint. All other
 1242 writes are specific to invocations and therefore do not conflict with each other. Thus, the lemma
 1243 holds. \square
 1244

1245 **THEOREM A.5 (SOUNDNESS PROPERTY).** *If $\Delta; \Lambda; \Gamma \vdash c \dashv \Lambda'; \Gamma'$ then $\llbracket c \rrbracket$ maps well-formed logs to
 1246 well-formed logs.*

1247 **PROOF.** The proof follows by induction on the typing derivation of c .
 1248

- 1249 • **CASE INSTANTIATE:** by assumption, \mathcal{L} is well-formed.
- 1250 • **CASE CONNECTION:** by assumption \mathcal{L} is well-formed, therefore \mathcal{W} is a set. For the first
 1251 condition, there are two possibilities: either $p_d \in \mathcal{R}$ or $p_d \notin \mathcal{R}$. If the second case holds
 1252 then, $\mathcal{R}\{p_d \mapsto p_s\} = \mathcal{R} \subseteq \mathcal{W}$. If $p_d \in \mathcal{R}$ then, by assumption of the typing rule, the
 1253 availability of p_d is a subset of the availability of p_s and, by well-formedness of the input
 1254 log, $p_s \in \mathcal{W}$, which implies $\mathcal{R}\{p_d \mapsto p_s\} \subseteq \mathcal{W}$.
- 1255 • **CASE COMP:** by the induction hypothesis $\mathcal{L}_1 = \llbracket c_1 \rrbracket (\mathcal{L})$ and $\mathcal{L}_2 = \llbracket c_2 \rrbracket (\mathcal{L})$ are well-
 1256 formed. Additionally, the writes for the two logs are disjoint (Lemma A.4). The first con-
 1257 dition of well-formedness is remains valid since new reads cannot make it invalid. This
 1258 means both conditions of well-formedness hold.
- 1259 • **CASE INVOKE:** The invoked instance M generates a partial log \mathcal{P} . This log is compatible with
 1260 the set of writes produced by the arguments of the invocation because each argument is
 1261 required to be available for at least as long as the signature of the module requires, similar to
 1262 the COMP case. Therefore, the log $\mathcal{P} \cup \mathcal{L}$ satisfies the first condition. The second condition
 1263 holds because the typing judgement requires the that the timeline context contains the
 1264 timeline required by the current invocation to be in the context. This is only possible if a
 1265 prior invocation did not use an overlapping timeline.
 1266

1267 \square

1268 B EVALUATION

1269 B.1 Filament Designs

1270 We implemented several designs in Filament to evaluate its expressivity. Each design has a cor-
 1271 responding baseline that we extracted from either a handwritten implementation or a generated
 1272 design.
 1273
 1274

1275 **Floating-point add.** We extracted a floating point adder implemented in Verilog.⁴ The imple-
 1276 mentation provides both a single-cycle and a pipelined version. We then translated the pipelined
 1277 implementation into a Filament design and used a fuzzing harness to ensure that the outputs of
 1278 the implementation matched to the source. In the process of translation, we found several bugs in
 1279 the original implementation. Each bug was caused by the pipelined design accessing signals from
 1280 a previous pipeline stage. Such bugs were immediately obvious in Filament since the type checker
 1281 disallows using signals from a previous stage to compute the value of a signal in the current stage.
 1282 For example, in the second stage of the pipeline, the adder attempts to use a value from the previous
 1283 stage:

```
1284 // The suffix *_s2 indicates that the signal is used in the second stage.
1285 if (exponent != 0) begin
1286   l1_s2 = {1'b1, Large_s2[22:1]};
1287 end else begin
1288   l1_s2 = Large_s2;
1289 end
```

1290 In Filament, this bug is caught by the type checker:

```
1291 M := new Mux[23];
1292 l1_mant := M<G+1>(e2_is_zero_s2.out, large_mant_s2.out, normed_large.out);
1293 Available in [G, G+1] but required in [G+1, G+2]
```

1294 Using cycle-accurate harness, coupled with a simple fuzzer made it easy to find such bugs by
 1295 differential testing of the combinational, pipelined, and Filament implementations. Our takeaway
 1296 from the study was that *once a Filament design functionally agrees with a pipelined design, any*
 1297 *changes in the Filament code to pipeline the design are unlikely to introduce new bugs due to the type*
 1298 *checking.*

1299 **Systolic Array.** Systolic arrays [28] are a common design pattern for implementing linear algebra
 1300 operations. They mimic the two-dimensional structure of the data by arranging the computation in
 1301 a two-dimensional array of processing elements and exploit the data reuse between computations.
 1302 We took the baseline systolic array reported in the Calyx [34] paper and implemented it in Filament.
 1303 We used the Prev component implemented in Section 7.2 to express the computation of the systolic
 1304 array.

```
1305 component main<G: 1>(  
1306   @interface[G] go: 1, @[G, G+1] l0: 32, @[G, G+1] l1: 32, @[G, G+1] t0: 32, @[G, G+1] t1: 32,  
1307 ) -> (  
1308   @[G, G+1] out00: 32, @[G, G+1] out01: 32, @[G, G+1] out10: 32, @[G, G+1] out11: 32,  
1309 ) {  
1310   // Systolic registers that go from left to right  
1311   r00_01 := new Prev[32, 1]<G>(l0); r00_10 := new Prev[32, 1]<G>(t0);  
1312   r10_11 := new Prev[32, 1]<G>(l1); r01_11 := new Prev[32, 1]<G>(t1);  
1313  
1314   // Connection registers to processing elements  
1315   pe00 := new Process<G>(l0, t0);  
1316   pe01 := new Process<G>(r00_01.prev, t1);  
1317   pe10 := new Process<G>(l1, r00_10.prev);  
1318   pe11 := new Process<G>(r10_11.prev, r01_11.prev);  
1319  
1320   out00 = pe00.out; out01 = pe01.out; out10 = pe10.out; out11 = pe11.out;  
1321 }
```

1320 The processing element Process performs a multiply-accumulate operation:

1322 ⁴Source: <https://github.com/suhasr1991/5-Stage-Pipelined-IEEE-Single-Precision-Floating-Point-Adder-Design>.

```

1324 component Process<G: 1>(
1325   @interface[G] go: 1, @[G, G+1] left: 32, @[G, G+1] right: 32) -> (@[G, G+1] out: 32) {
1326   // If acc does not contain a valid value, use 0
1327   acc := new Prev[32, 0]<G>(add.out);
1328   go_prev := new Prev[1, 1]<G>(go);
1329   mux := new Mux[32]<G>(go_prev.prev, acc.prev, 0);
1330
1331   mul := new MultComb[32]<G>(left, right);
1332   add := new Add[32]<G>(mux.out, mul.out);
1333
1334   out = add.out;
1335 }

```

1335 This process element will run at a lower frequency because it uses a combinational multiply
1336 (MultComb) instead of a pipelined multiply (Mult). However, it is straightforward to change this
1337 by using a pipelined multiply instead since it only changes the latency of the process module. Our
1338 alternate design uses a pipelined multiplier (FastMult) which changes the latency of the design.

```

1339 component Process<G: 1>(
1340   @interface[G] go: 1, @[G, G+1] left: 32, @[G, G+1] right: 32) -> (@[G+3, G+4] out: 32) {
1341   ...
1342   mul := new FastMult<G>(left, right);
1343 }

```

1344 B.2 Type Signatures

1345 We import designs from PipelineC [25] and give them type signatures in Filament.

1346 **Floating-point add.** The floating-point add generated by PipelineC is automatically pipelined by
1347 the compiler based on the frequency target provided by the user. We used the floating-point unit
1348 implemented in Appendix B.1 to differentially test the module and ensure that we provided the
1349 right type signature for the module.
1350

```

1351 component FpAdd<G: 1>(
1352   clk: 1, @[G, G+1] my_pipeline_x: 32, @[G, G+1] my_pipeline_y: 32
1353 ) -> (
1354   @[G+6, G+7] my_pipeline_return_output: 32
1355 );
1356

```

1357 **AES.** We give the type signature for an AES module implemented in PipelineC and automatically
1358 pipelined by the compiler:

```

1359 component AES<G: 1>(
1360   clk: 1, @[G, G+1] state_words: 128, @[G, G+1] keys: 1280) -> (@[G+18, G+19] out_words: 128)
1361

```

1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372