# Catapult HLS Modeling using SystemC and Matchlib

## Stuart Swan

HLS IP/Platform Architect

11 September 2019

**Mentor®**
A Siemens Business

# Why use HLS?

- Automatic scheduling, resource sharing, pipelining

- Easier retargeting to different silicon technologies

- More abstract design models

- Enable more DV and debugging to occur at higher abstraction level

- Make architectural exploration feasible

**Mentor®**
A Siemens Business

# Why use C++ for HLS?

- Templates, classes, inheritance, …
- Easy integration of existing C/C++ models
- Simulation speed
- …

**Mentor**®
A Siemens Business

# Why use SystemC + Matchlib for HLS?

- SystemC provides HDL semantics on top of C++ language
  — time, module structure, hierarchy, channels, HW semantics, resets, signals

- SC enables time based behaviors to be cleanly modeled and verified prior to synthesis (as compared to pure C++ HLS)

- SC integrates nicely with HDL simulators

- Matchlib enables models to be "throughput accurate" pre-HLS

**Mentor®**
A Siemens Business

# How does Catapult SC synthesis differ from C++ synthesis ?

- Primary difference is in input language
- Once code is read in by Catapult, 90% of flow and usage is same.

**Mentor®**
A Siemens Business

# SC Combinational Process

```
 3 #pragma once
 4
 5 #include "stdlib.h"
 6 #include "mc_trace.h"
 7
 8 #pragma hls_design top
 9 class and_gate : public sc_module {
10 public:
11   sc_in<bool>  INIT_S1(in1);
12   sc_in<bool>  INIT_S1(in2);
13   sc_out<bool> INIT_S1(out1);
14
15   SC_CTOR(and_gate)
16   {
17     SC_METHOD(run);
18     sensitive << in1 << in2;
19   }
20
21   void run() {
22     out1 = in1.read() & in2.read();
23   }
24 };
```

```
10 //
11 // ---------------------------------------------------
12 //  Design Unit:    and_gate
13 // ---------------------------------------------------
14
15
16 module and_gate (
17   in1, in2, out1
18 );
19   input in1;
20   input in2;
21   output out1;
22
23
24
25   // Interconnect Declarations for Component Instantiations
26   assign out1 = in1 & in2;
27 endmodule
28
29
```

HLS Input

HLS Output

# SC Sequential Process

```
10 #pragma hls_design top
11 class flop : public sc_module {
12 public:
13   sc_in<bool>      INIT_S1(clk);
14   sc_in<bool>      INIT_S1(rst_bar);
15   sc_in<uint32_t>  INIT_S1(in1);
16   sc_out<uint32_t> INIT_S1(out1);
17
18   SC_CTOR(flop)
19   {
20     SC_THREAD(process);
21     sensitive << clk.pos();
22     async_reset_signal_is(rst_bar, false);
23   }
24
25   void process() {
26     // this is the reset state:
27     out1 = 0;
28     wait();
29
30     // this is the non-reset state:
31     while (1) {
32       out1 = in1.read();
33       wait();
34     }
35   }
36 };
```

```verilog
16 module flop_process (
17   clk, rst_bar, in1, out1
18 );
19   input clk;
20   input rst_bar;
21   input [31:0] in1;
22   output [31:0] out1;
23   reg [31:0] out1;
24
25
26
27   // Interconnect Declarations for Component Instantiations
28   always @(posedge clk or negedge rst_bar) begin
29     if ( ~ rst_bar ) begin
30       out1 <= 32'b00000000000000000000000000000000;
31     end
32     else begin
33       out1 <= in1;
34     end
35   end
36 endmodule
37
```

HLS Input

HLS Output

Mentor
A Siemens Business
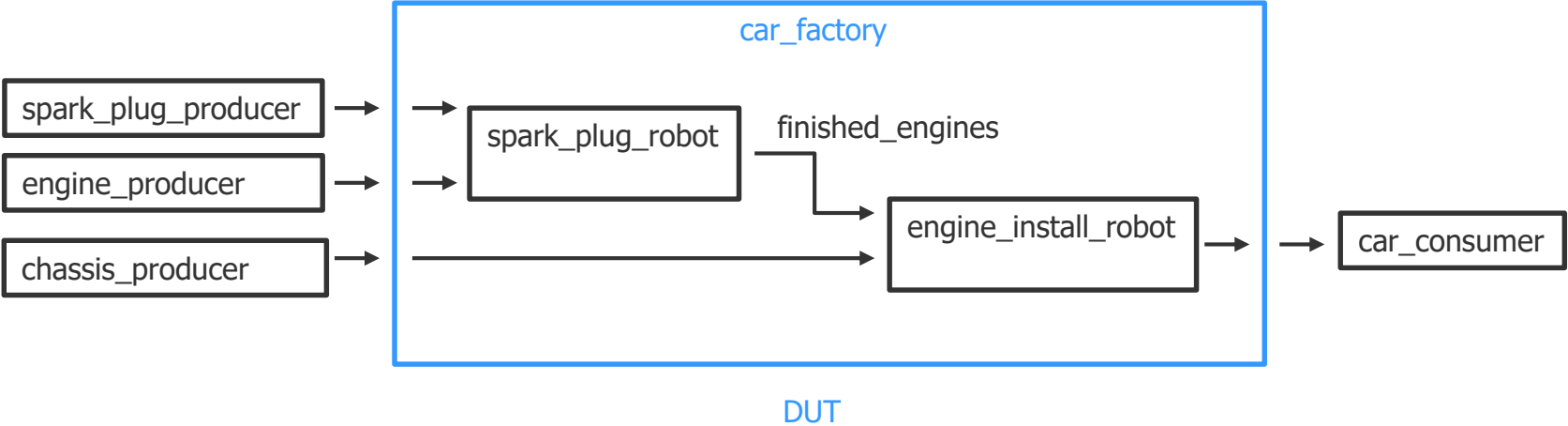
# Important to understand how HLS sees your SC model:

- Set of processes communicating **only** thru signals (sc_signal<>)
- For synthesis purposes, hierarchy is irrelevant (it is preserved in the RTL however)
- HLS synthesis occurs on each process one at a time, in complete isolation
  — No analysis/optimizations across processes
  — Huge designs can be synthesized thru HLS, as long as each process is not too large
- Each combinational process (SC_METHOD) becomes combinational logic in RTL
- Each sequential process becomes exactly one FSM + Datapath in RTL

# How HLS sees your SC model (continued)

- Primary optimizations done by HLS are in construction of the FSM and Datapath (e.g. resource sharing).

- Unlike RTL synthesis, HLS is (usually) allowed to add clock cycles (latency) into design to improve QOR
  - e.g. often loops are pipelined to maintain thruput while latency is added by HLS

- Properly constructed models and proper usage of HLS should show no functional differences (aside from latency differences) pre and post HLS.

- For SC HLS, reset behaviors and IO protocols are present in pre-HLS model and synthesized into RTL together with rest of model

# Simple Example of HW Architectural Model using SC + Matchlib

# spark_plug_producer

```cpp
15 class spark_plug_producer : public sc_module {
16 public:
17   sc_in<bool>                    INIT_S1(clk);
18   Connections::Out<spark_plug_t> INIT_S1(spark_plugs);
19
20   SC_CTOR(spark_plug_producer) {
21     SC_THREAD(main);
22     sensitive << clk.pos();
23   }
24
25   void main() {
26     int count=0;
27
28     while (1) {
29       spark_plug_t spark_plug;
30       spark_plug.spark_plug = count++;
31       spark_plugs.Push(spark_plug);
32       wait(3);
33       if (rand() & 1)
34         wait(3);
35     }
36   }
37 };
```

produces new spark_plug every 3-6 seconds

# engine_producer

```
39 class engine_producer : public sc_module {
40 public:
41   sc_in<bool>                 INIT_S1(clk);
42   Connections::Out<engine_t> INIT_S1(engines);
43
44   SC_CTOR(engine_producer) {
45     SC_THREAD(main);
46     sensitive << clk.pos();
47   }
48
49   void main() {
50     int count=0;
51
52     while (1) {
53       engine_t engine;
54       engine.engine = count++;
55       engines.Push(engine);
56       wait(20);
57     }
58   }
59 };
```

produces new engine every 20 seconds

**Mentor®**
A Siemens Business

# chassis_producer

```
61 class chassis_producer : public sc_module {
62 public:
63   sc_in<bool>                    INIT_S1(clk);
64   Connections::Out<chassis_t>    INIT_S1(chassis_out);
65
66   SC_CTOR(chassis_producer) {
67     SC_THREAD(main);
68     sensitive << clk.pos();
69   }
70
71   void main() {
72     int count=0;
73
74     while (1) {
75       chassis_t chassis;
76       chassis.chassis = count++;
77       chassis_out.Push(chassis);
78       wait(25);
79     }
80   }
81 };
```

produces new chassis every 25 seconds

# car_consumer

```cpp
83 class car_consumer : public sc_module {
84 public:
85   sc_in<bool>              INIT_S1(clk);
86   Connections::In<car_t>  INIT_S1(cars);
87
88   SC_CTOR(car_consumer) {
89     SC_THREAD(main);
90     sensitive << clk.pos();
91   }
92
93   void main() {
94     int count = 0;
95     while (1) {
96       cars.Pop();
97       ++count;
98       LOG("got car # " << count);
99       if (count == 10)
100       {
101         LOG("total cars produced: " << count);
102         LOG("time per car: " << sc_time_stamp() / count);
103         sc_stop();
104       }
105     }
106   }
107 };
```

consumes cars as quickly as possible

# Simple car_factory

```
139 #if defined(SIMPLE)
140 class car_factory : public sc_module {
141 public:
142   sc_in<bool>                    INIT_S1(clk);
143   Connections::In<spark_plug_t> INIT_S1(spark_plugs);
144   Connections::In<engine_t>     INIT_S1(engines);
145   Connections::In<chassis_t>    INIT_S1(chassis);
146   Connections::Out<car_t>       INIT_S1(cars);
147
148   Connections::Combinational<engine_t> INIT_S1(finished_engines);
149
150   spark_plug_robot      INIT_S1(spark_plug_robot1);
151   engine_install_robot  INIT_S1(engine_install_robot1);
152
153   SC_CTOR(car_factory)
154   {
155     spark_plug_robot1.clk(clk);
156     spark_plug_robot1.spark_plugs(spark_plugs);
157     spark_plug_robot1.engines_in(engines);
158     spark_plug_robot1.engines_out(finished_engines);
159
160     engine_install_robot1.clk(clk);
161     engine_install_robot1.chassis(chassis);
162     engine_install_robot1.engines(finished_engines);
163     engine_install_robot1.cars(cars);
164   }
165 };
166
```

# spark_plug_robot

```cpp
71 class spark_plug_robot : public sc_module {
72 public:
73    sc_in<bool>                      INIT_S1(clk);
74    Connections::In<spark_plug_t>    INIT_S1(spark_plugs);
75    Connections::In<engine_t>        INIT_S1(engines_in);
76    Connections::Out<engine_t>       INIT_S1(engines_out);
77    SC_SIG(bool, busy);
78    SC_SIG(bool, maintenance);
79
80    SC_CTOR(spark_plug_robot)
81    {
82      SC_THREAD(main);
83      sensitive << clk.pos();
84    }
85
86    void main() {
87      int count = 0;
88      while (1) {
89        engine_t engine_in = engines_in.Pop();
90        for (int i=0; i < engine_t::plugs; i++)
91          engine_in.spark_plugs[i] = spark_plugs.Pop();
92        busy = 1;
93        wait(60);
94        busy = 0;
95        engines_out.Push(engine_in);
96        if ((count++ & 1) && (rand() & 3))
97        {
98          maintenance = 1;
99          wait(60);
100         maintenance = 0;
101       }
102     }
103   }
104 };
```

Consumes 4 spark_plugs and 1 unfinished engine
Produces finished_engine after 60 seconds
After every other engine, 75% of time
needs 60 seconds of maintenance (ie idle time)

# engine_install_robot

```
106 class engine_install_robot : public sc_module {
107 public:
108    sc_in<bool>                    INIT_S1(clk);
109    Connections::In<chassis_t>   INIT_S1(chassis);
110    Connections::In<engine_t>    INIT_S1(engines);
111    Connections::Out<car_t>      INIT_S1(cars);
112    SC_SIG(bool, busy);
113
114    SC_CTOR(engine_install_robot)
115    {
116      SC_THREAD(main);
117      sensitive << clk.pos();
118    }
119
120    void main() {
121      while (1) {
122        car_t car;
123        car.chassis = chassis.Pop();
124        car.engine = engines.Pop();
125        busy = 1;
126        wait(30);
127        busy = 0;
128        cars.Push(car);
129      }
130    }
131 };
```

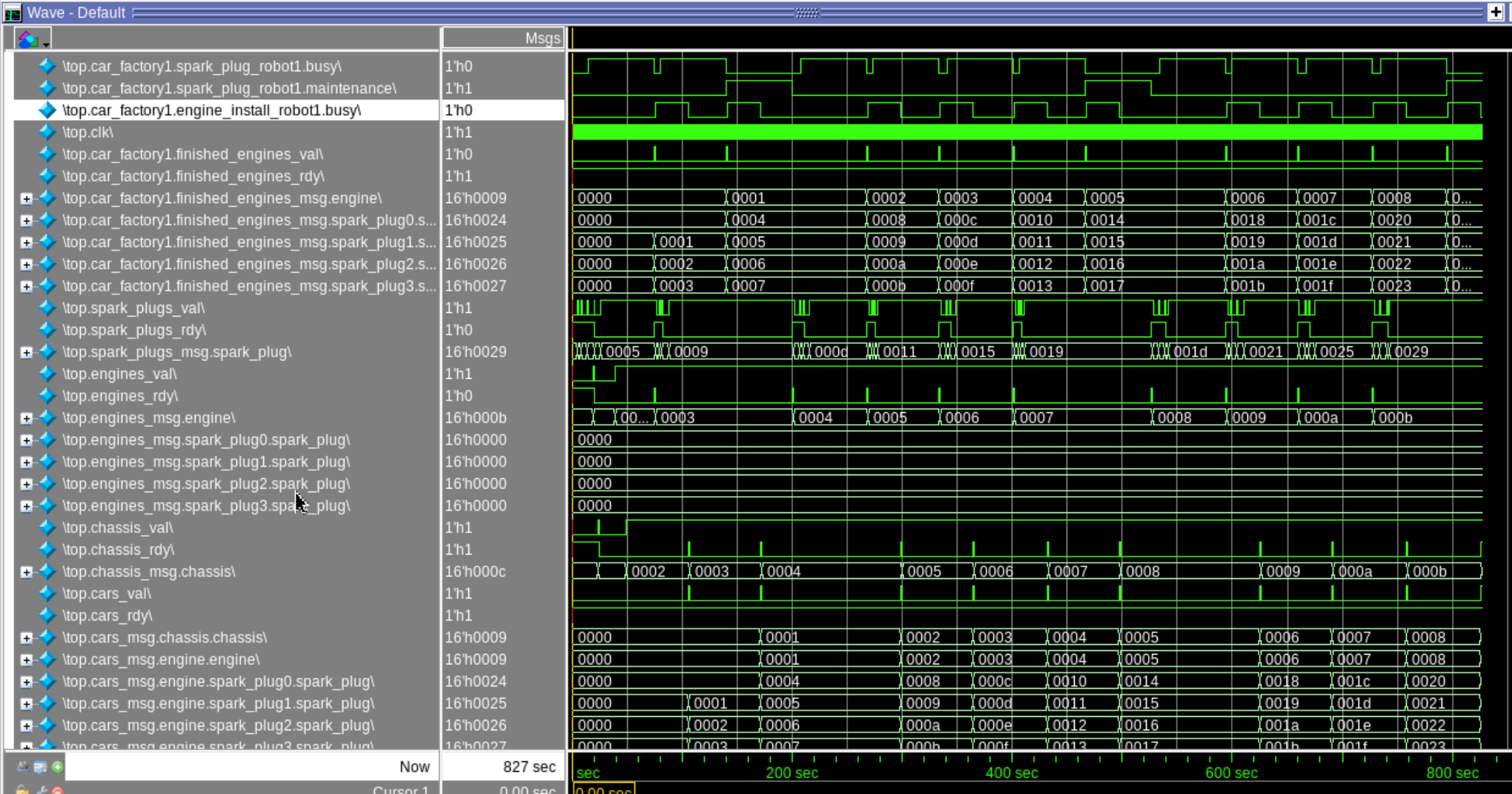Consumes 1 chassis and 1 finished_engine
Produces car after 30 seconds

# Running simple car_factory

```
107 s top.car_consumer1 got car # 1
172 s top.car_consumer1 got car # 2
300 s top.car_consumer1 got car # 3
365 s top.car_consumer1 got car # 4
433 s top.car_consumer1 got car # 5
498 s top.car_consumer1 got car # 6
626 s top.car_consumer1 got car # 7
691 s top.car_consumer1 got car # 8
759 s top.car_consumer1 got car # 9
827 s top.car_consumer1 got car # 10
827 s top.car_consumer1 total cars produced: 10
827 s top.car_consumer1 time per car: 82700 ms

Info: /OSCI/SystemC: Simulation stopped by user.
```

Goal is to produce each car in smallest amount of time

# Running simple car_factory

# Sequential car_factory

```cpp
#elif defined(SEQUENTIAL)
class car_factory : public sc_module {
public:
  sc_in<bool>                 INIT_S1(clk);
  Connections::In<spark_plug_t> INIT_S1(spark_plugs);
  Connections::In<engine_t>    INIT_S1(engines);
  Connections::In<chassis_t>   INIT_S1(chassis);
  Connections::Out<car_t>      INIT_S1(cars);

  SC_CTOR(car_factory)
  {
    SC_THREAD(main);
    sensitive << clk.pos();
  }

  SC_SIG(bool, spark_plug_robot_busy);
  SC_SIG(bool, spark_plug_robot_maintenance);
  SC_SIG(bool, engine_install_robot_busy);

  void main() {
    spark_plug_t plugs[engine_t::plugs];
    int plug_count = 0;

    engine_t unfinished_engine;
    int unfinished_engine_count = 0;
    engine_t finished_engine;
    int finished_engine_count = 0;
    chassis_t chassis_inst;
    int chassis_count = 0;
    int spark_plug_robot_count = 0;

    while (1) {
      if (plug_count < engine_t::plugs)
        if (spark_plugs.PopNB(plugs[plug_count]))
          ++plug_count;

      if (unfinished_engine_count == 0)
        if (engines.PopNB(unfinished_engine))
          ++unfinished_engine_count;

      if (chassis_count == 0)
        if (chassis.PopNB(chassis_inst))
          ++chassis_count;
```

```cpp
      if ((unfinished_engine_count == 1) && (plug_count == engine_t::plugs)
          && (finished_engine_count == 0))
      {
        finished_engine = unfinished_engine;
        for (int i=0; i < engine_t::plugs; i++)
          finished_engine.spark_plugs[i] = plugs[i];
        spark_plug_robot_busy = 1;
        wait(60);
        spark_plug_robot_busy = 0;
        if ((spark_plug_robot_count++ & 1) && (rand() & 3))
        {
          spark_plug_robot_maintenance = 1;
          wait(60);
          spark_plug_robot_maintenance = 0;
        }
        finished_engine_count = 1;
        plug_count = 0;
        unfinished_engine_count = 0;
      }

      if ((finished_engine_count == 1) && (chassis_count == 1))
      {
        car_t car;
        car.chassis = chassis_inst;
        car.engine = finished_engine;
        engine_install_robot_busy = 1;
        wait(30);
        engine_install_robot_busy = 0;
        cars.Push(car);
        finished_engine_count = 0;
        chassis_count = 0;
      }

      wait();
    }
  }
};
```
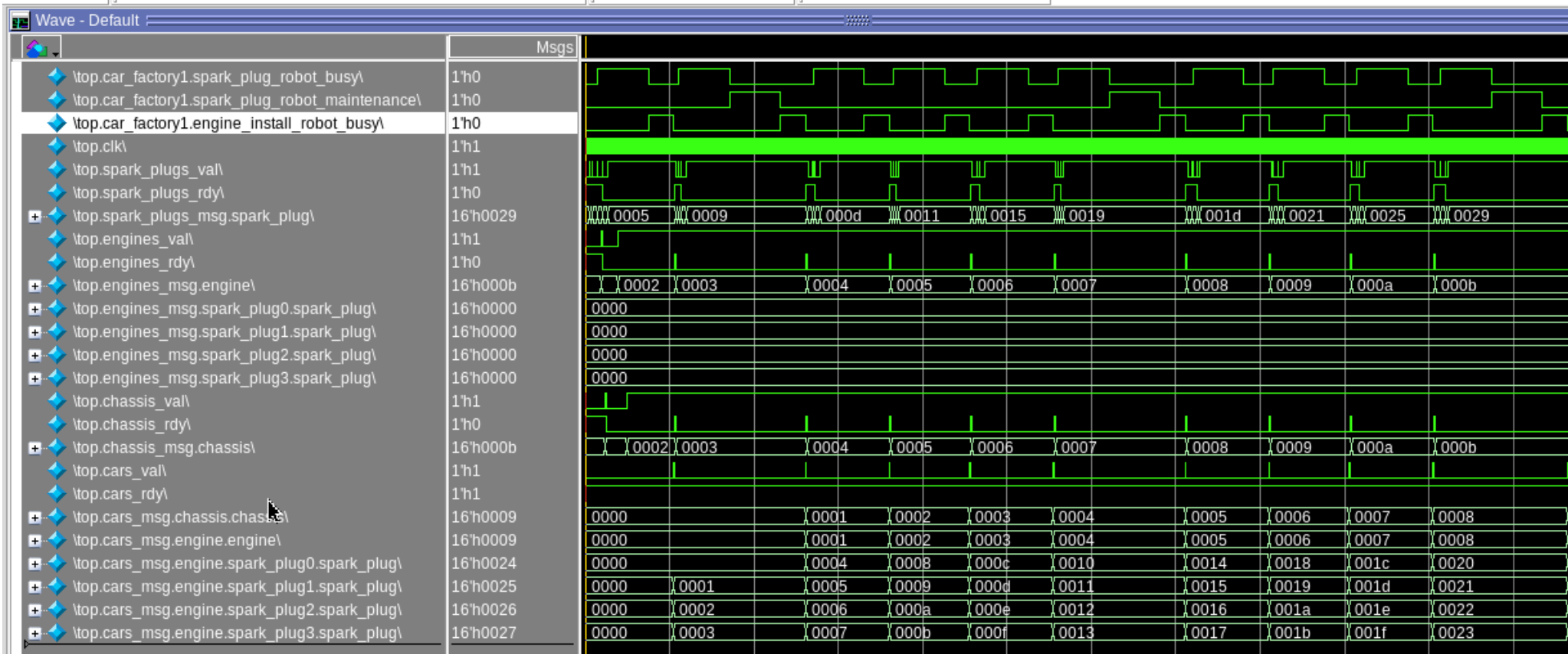
# Running sequential car_factory

```
106 s top.car_consumer1 got car # 1
262 s top.car_consumer1 got car # 2
361 s top.car_consumer1 got car # 3
457 s top.car_consumer1 got car # 4
556 s top.car_consumer1 got car # 5
712 s top.car_consumer1 got car # 6
811 s top.car_consumer1 got car # 7
907 s top.car_consumer1 got car # 8
1006 s top.car_consumer1 got car # 9
1165 s top.car_consumer1 got car # 10
1165 s top.car_consumer1 total cars produced: 10
1165 s top.car_consumer1 time per car: 116500 ms

Info: /OSCI/SystemC: Simulation stopped by user.
```

Car production time got worse!

# Running sequential car_factory

# Will HLS fix sequential car_factory?

- Can HLS split the single sequential process into 2 different processes to improve the utilization / QOR?
  - No. HLS always generates a single FSM+Datapath per SC process in the input model.

- OK, then can HLS generate single FSM that is the product of the 2 simpler state machines?
  - No. Even if it did, you would have "state explosion", leading to bad QOR due to a huge FSM.

**Mentor®**
A Siemens Business

# The "state explosion" problem

### 5.1.1 State Explosion

A flat FSM suffers from *state explosion*, which occurs when multiple independent activities interfere in a single model. Assume that an FSM has to capture two independent activities, each of which can be in one of three states. The resulting FSM,

133

called a *product state-machine*, needs nine states to represent the overall model. The product state-machine needs to keep track of the current state from two independent state machines at the same time. Due to conditional state transitions, one state machine can remain in a single state while the other state machine proceeds to the next state. This results in multiple intermediate states such as A1, A2, and A3. Figure 5.1 illustrates the effect of state explosion in a product state-machine. Two state machines, FSM1 and FSM2, need to be merged into a single product state-machine FSM1xFSM2. In order to represent all individual states, 9 states are needed in total. The resulting number of state transitions (and state transition conditions) is even higher. Indeed, if we have $n$ independent state transition conditions in the individual state machines, the resulting product state-machine can have up to $2^n$ state transition conditions.

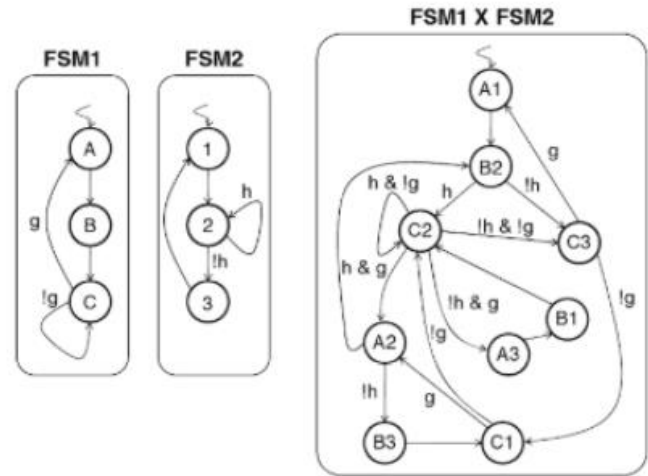134                                              5 Microprogrammed Architectures



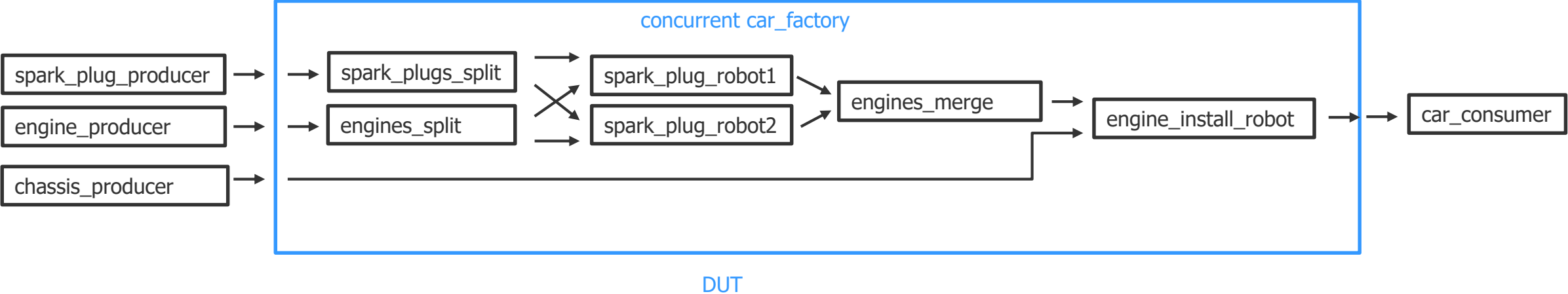**Fig. 5.1** State explosion in FSM when creating a product state-machine

© Mentor Graphics Corp.    Company Confidential

# How do we fix the car_factory architecture?

- Primary problem in "simple" car_factory is overutilization of spark_plug_robot

- Obvious solution: add another spark_plug_robot

# concurrent car_factory

# spark_plugs_split and engines_split

```cpp
class spark_plugs_split : public sc_module {
public:
  sc_in<bool>                       INIT_S1(clk);
  Connections::In<spark_plug_t>     INIT_S1(spark_plugs_in);
  Connections::Out<spark_plug_t>    INIT_S1(spark_plugs_out1);
  Connections::Out<spark_plug_t>    INIT_S1(spark_plugs_out2);

  SC_CTOR(spark_plugs_split)
  {
    SC_THREAD(main);
    sensitive << clk.pos();
  }

  void main() {
    while (1) {
      spark_plug_t spark_plug = spark_plugs_in.Pop();
      while (1) {
        if (spark_plugs_out1.PushNB(spark_plug))
          break;
        if (spark_plugs_out2.PushNB(spark_plug))
          break;
        wait();
      }
    }
  }
};
```

```cpp
class engines_split : public sc_module {
public:
  sc_in<bool>                   INIT_S1(clk);
  Connections::In<engine_t>     INIT_S1(engines_in);
  Connections::Out<engine_t>    INIT_S1(engines_out1);
  Connections::Out<engine_t>    INIT_S1(engines_out2);

  SC_CTOR(engines_split)
  {
    SC_THREAD(main);
    sensitive << clk.pos();
  }

  void main() {
    while (1) {
      engine_t engine = engines_in.Pop();
      while (1) {
        if (engines_out1.PushNB(engine))
          break;
        if (engines_out2.PushNB(engine))
          break;
        wait();
      }
    }
  }
};
```

# engines_merge

```cpp
class engines_merge : public sc_module {
public:
  sc_in<bool>                INIT_S1(clk);
  Connections::In<engine_t>      INIT_S1(engines_in1);
  Connections::In<engine_t>      INIT_S1(engines_in2);
  Connections::Out<engine_t>     INIT_S1(engines_out);

  SC_CTOR(engines_merge)
  {
    SC_THREAD(main);
    sensitive << clk.pos();
  }

  void main() {
    while (1) {
      engine_t engine;
      while (1) {
        if (engines_in1.PopNB(engine))
          break;
        if (engines_in2.PopNB(engine))
          break;
        wait();
      }
      engines_out.Push(engine);
    }
  }
};
```
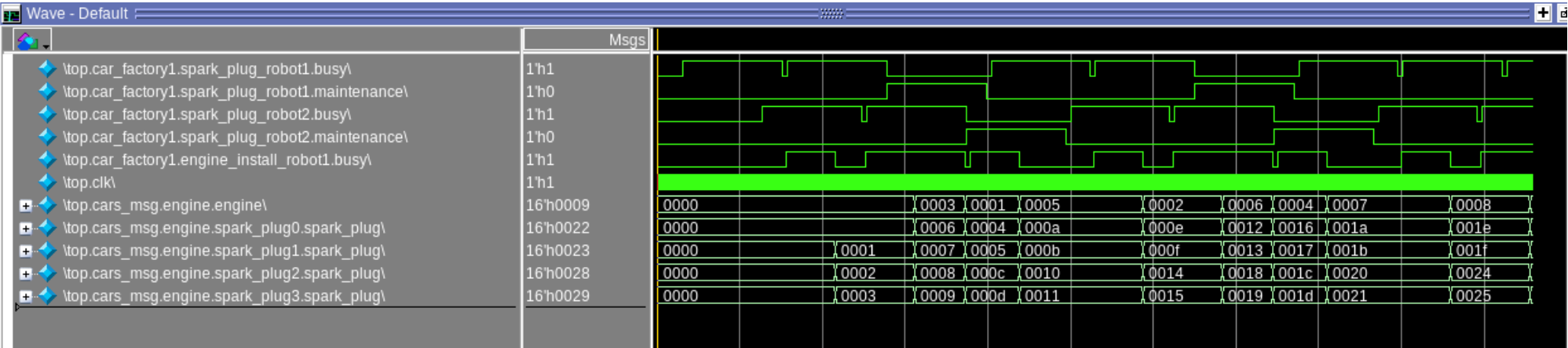
# Running concurrent car_factory

```
109 s top.car_consumer1 got car # 1
157 s top.car_consumer1 got car # 2
187 s top.car_consumer1 got car # 3
220 s top.car_consumer1 got car # 4
295 s top.car_consumer1 got car # 5
343 s top.car_consumer1 got car # 6
373 s top.car_consumer1 got car # 7
406 s top.car_consumer1 got car # 8
481 s top.car_consumer1 got car # 9
529 s top.car_consumer1 got car # 10
529 s top.car_consumer1 total cars produced: 10
529 s top.car_consumer1 time per car: 52900 ms

Info: /OSCI/SystemC: Simulation stopped by user.
```

Big improvement in car production time!

# Running concurrent car_factory



Both spark_plug_robots busy at same time

Better utilization

Output order of engines and plugs no longer matches input order

# Summing it up…

- In practice, SC modeling for HLS is a series of stepwise refinements that get to a good architectural model.

- You can do a lot of refinement & architectural analysis in pure SC model by analyzing the pre-HLS simulations.

- Use HLS to find latency/area/power issues and iterate with changes to the SC model source code.

**Mentor**®
A Siemens Business

# SystemC HLS Process Structure and Blocking/Non-Blocking IO

- A process (SC_THREAD) is like an "always" block in Verilog that is sensitive to the clock.
  - Usually it is "bigger" than a Verilog always block because it is an implicit state machine, may be pipelined by catapult, etc.
  - Usually every "block" in the HW architecture diagram will imply at least one independent SC_THREAD in the HLS model.

- Process structure has a huge impact on HLS runtime and QOR
  - Big processes usually mean "big" FSMs with lots of states and transitions **== bad QOR**
  - So, knowing how to arrive at a "good" process structure is extremely important
  - We saw in the car_factory example how breaking the design down into simpler concurrent activities led to a more efficient implementation

Mentor®
A Siemens Business

# Guidelines for arriving at a good process structure

- Usually IO is the biggest consideration in how to structure SC_THREADs.

- The ideal structure is usually an SC_THREAD with at least one blocking input (ie one Pop operation) and/or at least one blocking output (ie Push operation).

- At least one of the blocking inputs or outputs should be called unconditionally on every iteration of the main loop.

- If functionality can easily be expressed as independent SC_THREADs then it is generally better to do so – they will run independently, have smaller control FSMs, etc.

- But, resources cannot be shared by Catapult across different SC_THREADs, and any communication between 2 SC_THREADs will take at least 1 clock cycle since they are clocked processes.

# Unconditional blocking inputs or outputs..

- **Why is it important to have at least one unconditional blocking input or output?**
  - Usually the blocking IO is tied to each iteration of the main loop (which may be pipelined with for example II=1).
  - The blocking IO semantics and the clear II semantics clearly tie down the expected thruput requirements between the user and the HLS tool.
  - In contrast, if the IO is non-blocking, then even if the II=1 then Catapult may generate a state machine where on some iterations of a pipelined loop no actual IO occurs (since it is non-blocking)
  - Also, with blocking IO, the FSM will be in a clear "blocked" state when no IO can occur, and this makes idle state modeling easier (e.g. for power saving)

# When do you need to use Non-blocking IO?

- Arbitration requires Peek or PopNB  to all arbitrated inputs.

- Time-based splitting and merging of transaction streams requires PushNB and PopNB (respectively)

- There are cases where a process will need ALL non-blocking IO (all PopNB and PushNB), but it should be pretty rare.
  - In this case the process should be kept as small / simple as possible, ideally communicating with other processes that follow the guidelines above.
  - With all non-blocking IO, you will likely be modeling at very close to RTL level, and most likely HLS will just be translating SC RTL into Verilog RTL.

# Summary: Prefer to use blocking IO over nonblocking IO.

- Your models will be simpler and more likely to have a good process structure.

- 100% blocking IO is called KPN (Kahn Process Networks)
  — KPN is deterministic
  — easier to verify.

- Non-blocking IO is sometimes needed, but introduces timing dependent behavior, and can make verification more difficult in some cases if the timing dependent behavior is externally visible.

# Where to use wait() in HLS models

- Use 1 wait() statement at top of main loop to model the reset state.

- Only use wait() in HLS models within loops that have exclusively non-blocking IO.

- If you are modeling low-level protocols using sc_signal, you may need to use wait() (but you should try to avoid doing this!)

- You should not be using wait() anywhere else.

# Always try to think first about the HW architecture of the model

- Put your HW architects in front of a white board and get them to talk!

- What are primary blocks, transaction streams and required thruputs?

- What activities can naturally be expressed as separate concurrent processes?

- Where does arbitration need to occur?

- How can the architecture be modified to reduce undesirable contention?

- How can data storage and lifetime be minimized?

- What are expensive resources that should be shared ?
  - Either add into same process,
  - Or add arbiter to allow multiple processes to access

# What should be left for Catapult to do?

- Automated microarchitecture generation:
    — Detailed scheduling
    — resource allocation and sharing
    — loop unrolling
    — loop pipelining
    — memory access scheduling and interfacing
    — FSM generation
    — meeting timing of target silicon

# Conclusion

- Modeling a "good" process structure is key to getting good QOR thru HLS

- HLS tools such as Catapult can provide analysis feedback to help you identify "bad" process structures, but they will not automatically find a "good" one for you.

- Properly modeling the HW architecture and a good process structure is the responsibility of the HLS modeling engineers and HW architects

www.mentor.com