

# Automatic generation of RISC-V cores through speculative loop pipelining

Colloque GDR-SoC

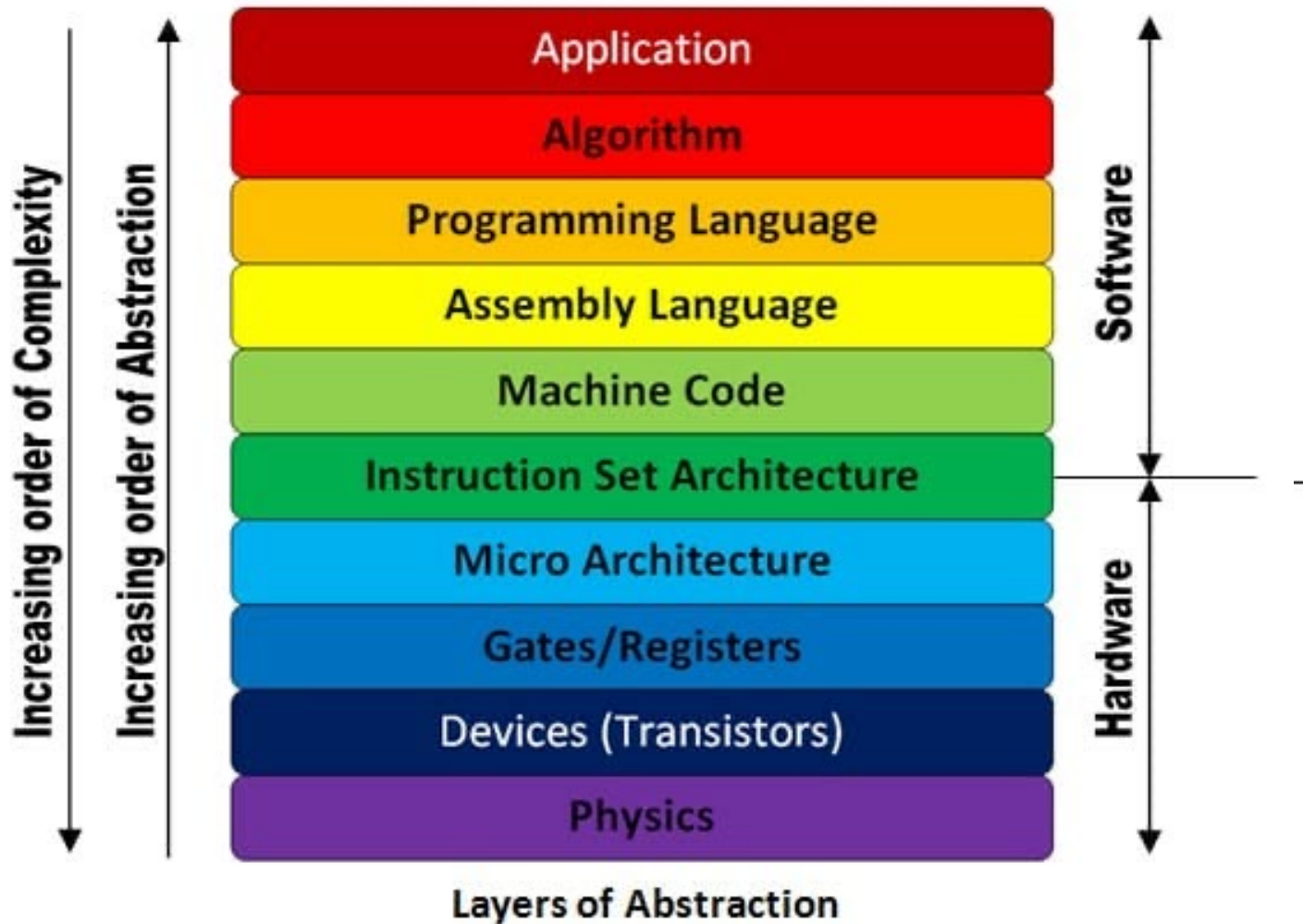
Steven Derrien

ARCAD Team, UBO/LabSTICC

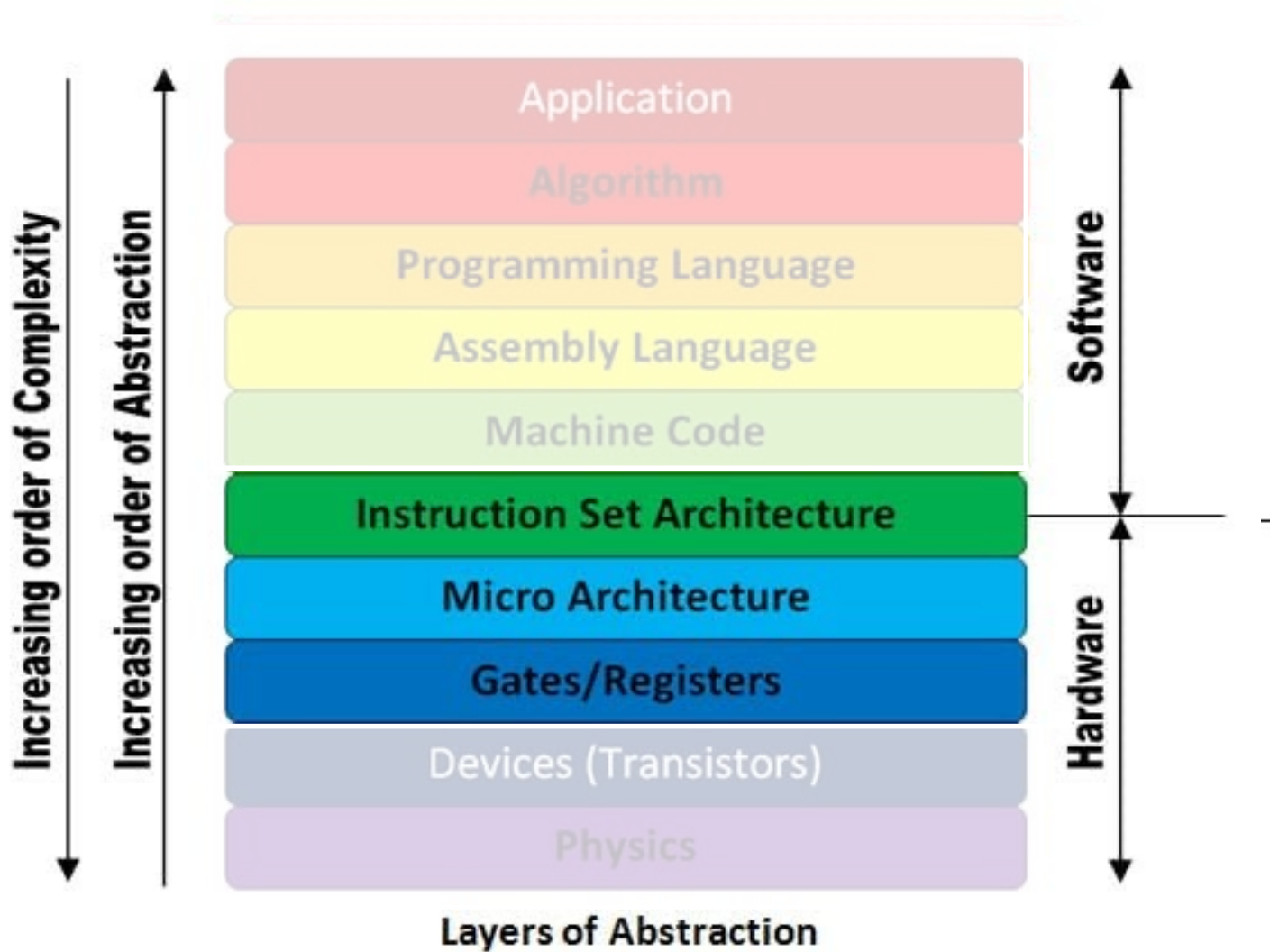
(with contributions of former colleagues from IRISA/INRIA : Simon Rokicki,  
Jean-Michel Gorius<sup>1</sup>, Dylan Leothaud, Thibaut Marty, Tomofumi Yuki)



# Abstraction layers in computer systems



# Abstraction layers in computer systems



# Instruction Set Architectures

Paper  
specification



Executable  
specification

imm[11:0]		rs1	010	rd	0000011	LW
imm[11:0]		rs1	100	rd	0000011	LBU
imm[11:0]		rs1	101	rd	0000011	LHU
imm[11:5]	rs2	rs1	000	imm[4:0]	0100011	SB
imm[11:5]	rs2	rs1	001	imm[4:0]	0100011	SH
imm[11:5]	rs2	rs1	010	imm[4:0]	0100011	SW
imm[11:0]		rs1	000	rd	0010011	ADDI
imm[11:0]		rs1	010	rd	0010011	SLTI
imm[11:0]		rs1	011	rd	0010011	SLTIU
imm[11:0]		rs1	100	rd	0010011	XORI
imm[11:0]		rs1	110	rd	0010011	ORI
imm[11:0]		rs1	111	rd	0010011	ANDI
0000000	shamt	rs1	001	rd	0010011	SLLI
0000000	shamt	rs1	101	rd	0010011	SRLI
0100000	shamt	rs1	101	rd	0010011	SRAI
0000000	rs2	rs1	000	rd	0110011	ADD
0100000	rs2	rs1	000	rd	0110011	SUB
0000000	rs2	rs1	001	rd	0110011	SLL

```

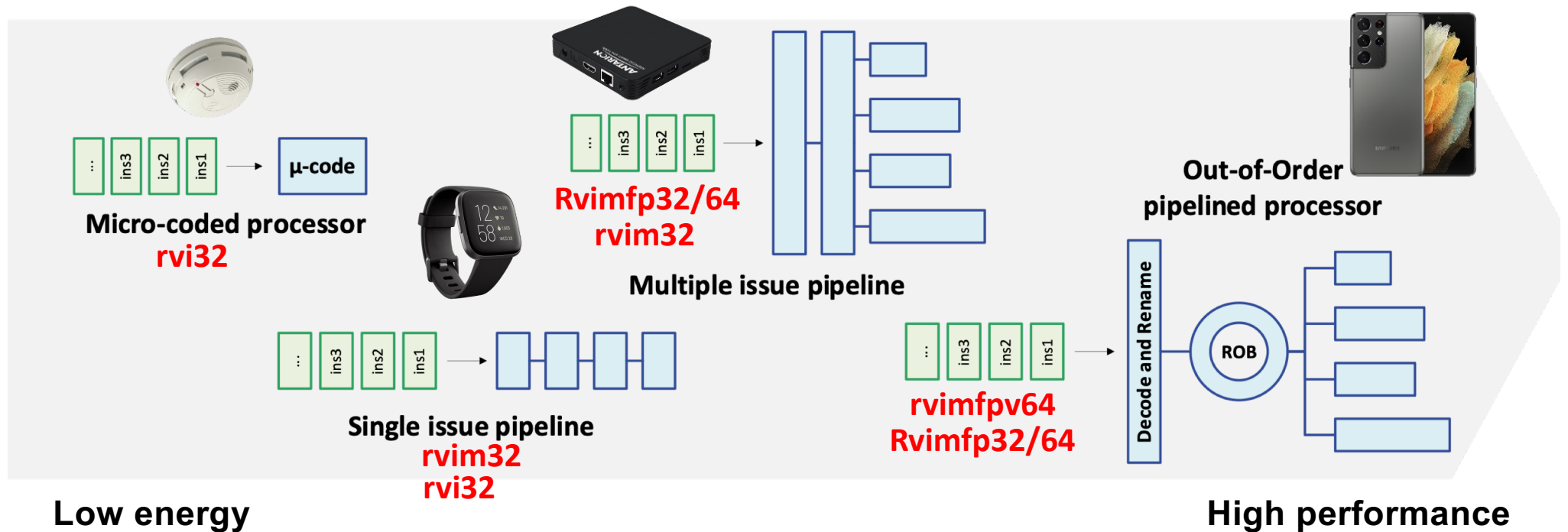
while(1) {
  IR = mem[PC], nextPC = PC + 4;
  a = rs1(IR), b = rs2(IR), d = rd(IR);
  switch (op(IR)) {
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);
              break;
    case MUL: X[d] = X[a]*X[b];
              break;
    case MULI: X[d] = X[a]*ext32(imm(IR));
              break;
    case LDW: X[d] = Mem[X[b]+imm(IR)];
              break;
    case STW: Mem[X[b]+imm(IR)] = X[a];
              break;
    case BEQ:
      nextPC += X[a]==X[b] ? imm(IR) : 0;
      break;
    // ...
  }
  PC = nextPC;
}
    
```

Example instruction	Instruction name	Meaning
add x1,x2,x3	Add	Regs[x1] ← Regs[x2]+Regs[x3]
addi x1,x2,3	Add immediate unsigned	Regs[x1] ← Regs[x2]+3

Instruction Set Simulator (ISS)

# Micro-architectures

- $\mu$ architecture = *how* the machine executes instructions

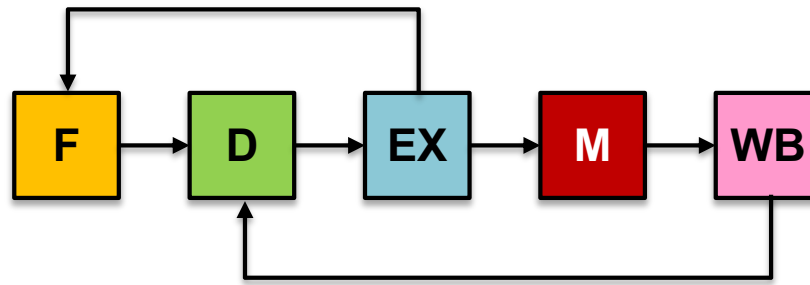


## Takeaway

Need for highly customized micro-architecture (there is no *one size fits all*)

# What is *processor* pipelining ?

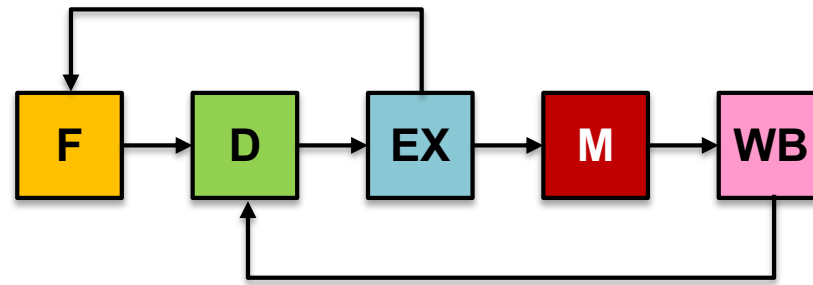
- Overlapping the execution of consecutive instructions
  - An instruction execution is decomposed into stages



	1	2	3	4	5	6	7	8	9	10	11	12	13
100: add r6,r5,r2	F	DC	EX	M	WB								
104: add r3,r1,r6		F	DC			EX	M	WB					
108: lw r5,(r1)			F			DC	EX	M	WB				
10B: bne r2,r2,100						F	DC	EX	M	WB			
110: lw r1,16(r5)							F	DC	<del>EX</del>	<del>M</del>	<del>WB</del>		
114: lw r2,4(r5)								F	<del>DC</del>	<del>EX</del>	<del>M</del>	<del>WB</del>	
100: add r6,r5,r2									F	DC	EX	M	WB

# What is *processor* pipelining ?

- Overlapping the execution of consecutive instructions
  - An instruction execution is decomposed into stages



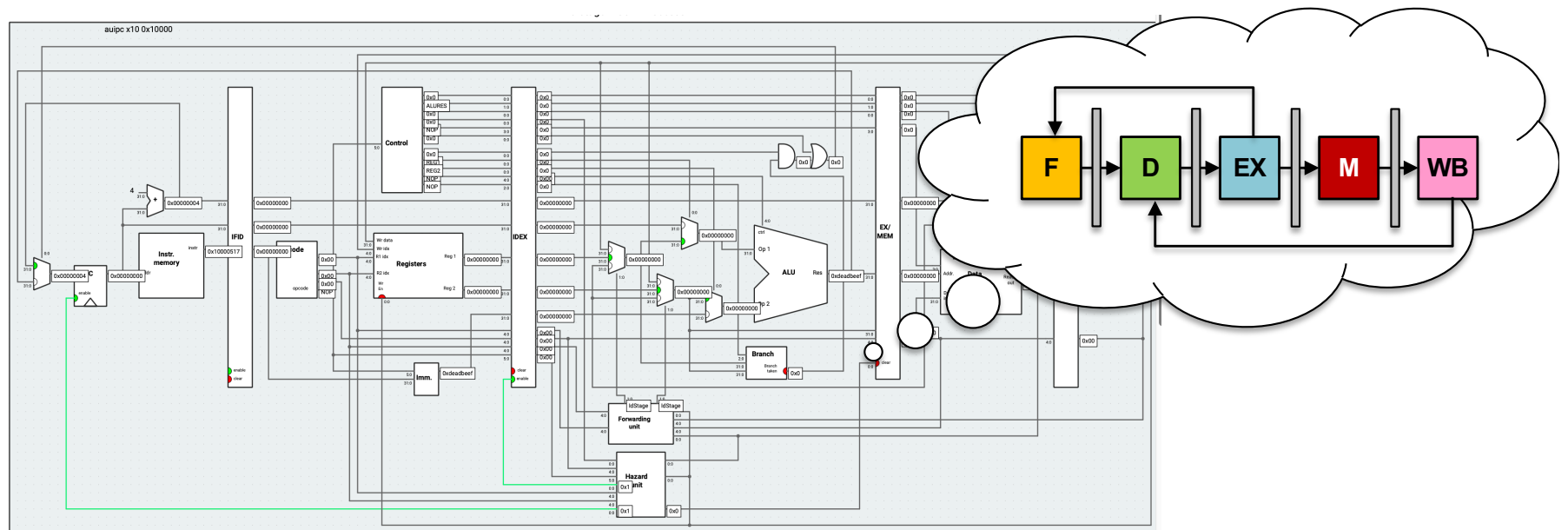
	1	2	3	4	5	6	7	8	9	10	11	12	13
100: add r6, r5, r2	F	DC	EX	M	WB								
104: add r3, r1, r6		F	DC			EX	M	WB					
108: lw r5, (r1)			F			DC	EX	M	WB				
10B: bne r2, r2, 100						F	DC	EX	M	WB			
110: lw r1, 16(r5)							F	DC					
114: lw r2, 4(r5)								F					

**Takeaway**

Hazard management is what make processor pipelining difficult

# Implementing the processor as a circuit

- Logic gates + registers implementing the chosen march
  - Described using Hardware Description Languages



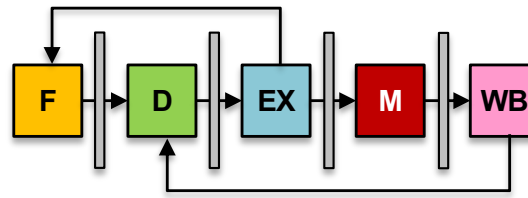
- Must be optimized for energy, area and clock speed
  - and must still be faithful to the reference marchitecture !

# Existing CPU Design Flows

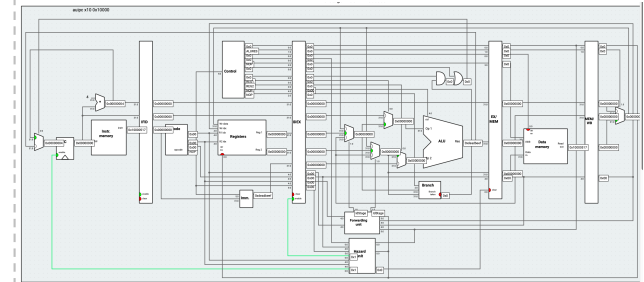
## Instruction set

Example instruction	Instruction name	Meaning
add x1,x2,x3	Add	Regs[x1] ← Regs[x2]+Regs[x3]
addi x1,x2,3	Add immediate unsigned	Regs[x1] ← Regs[x2]+3
lui x1,42	Load upper immediate	Regs[x1] ← 0 <sup>32</sup> ##42##0 <sup>12</sup>
sll x1,x2,5	Shift left logical	Regs[x1] ← Regs[x2] << 5
slt x1,x2,x3	Set less than	if (Regs[x2] < Regs[x3]) Regs[x1] ← 1 else Regs[x1] ← 0

## Micro-architecture



## Circuit



## DSL specification

```
while (1) {
  opcode,ra,rb,rc,imm = mem[pc];
  pc++;
  switch (opcode) {
    case LOAD: reg[rc] = mem[reg[ra] + imm]; break;
    case ADD: reg[rc] = reg[ra] + reg[rb]; break;
    case JUMP: pc = reg[ra]; break;
    case BNZ: pc = rb ? reg[ra] : pc; break;
    ...
  }
}
```

## DSL specification

The Codal logo is shown with a red 'X' over a blue arrow pointing to it. Below the logo is a snippet of DSL specification code:

```
// Instruction-set grammar
opn my_core (arith_inst | ctrl_inst);
opn arith_inst (a:alu_inst,
               d: div_inst, l:load_store_inst);
...
stage EX2:
  RA[z] = C @alu;
syntax: op " RA" y " ", RB" x " , RA" z ;
image: "0":op:x:y:z;
```

## Verilog, Chisel

```
1 module control(clk,
2   opcode, isalump,
3   bus_addr,
4   do_fetch,
5   bus_addr, bus_
6   state);
7   /* Control
8   module control(clk,
9   do 2
10  opcode, isalump,
11  bus_addr, bus_
12  do_fetch, do_r
13  state);
14   /* Control module. */
15   `include "parameters.vh"
16   input wire clk;
17   /* Control mod
18
19   `include "parameters.vh"
20
21   input wire clk;
```

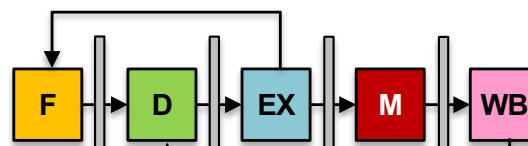
DSL = Domain Specific Languages

# Existing CPU Design Flows

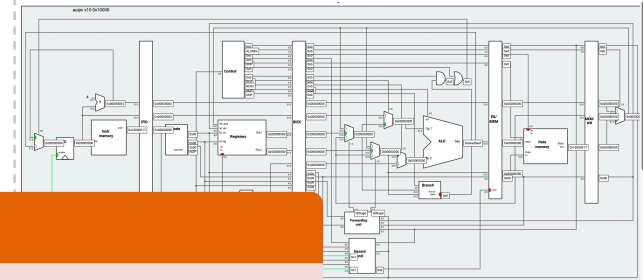
## Instruction set

Example instruction	Instruction name	Meaning
add x1,x2,x3	Add	Regs[x1] ← Regs[x2]+Regs[x3]
addi x1,x2,3	Add immediate unsigned	Regs[x1] ← Regs[x2]+3
lui x1,42	Load upper immediate	Regs[x1] ← 0 <sup>32</sup> ##42##0 <sup>12</sup>
sll x1,x2,5	Shift left logical	Regs[x1] ← Regs[x2] << 5
slt x1,x2,x3	Set less than	if (Regs[x1] < Regs[x2]) Regs[x1] ← 1

## Micro-architecture



## Circuit



**Takeaway**  
In practice, micro architectural DSLs are RTL in disguise

## DSL specification

```
while (1) {
  opcode,ra,rb,rc,imm = mem[pc];
  pc++;
  switch (opcode) {
    case LOAD: reg[rc] = mem[reg[ra] + imm]; break;
    case ADD: reg[rc] = reg[ra] + reg[rb]; break;
    case JUMP: pc = reg[ra]; break;
    case BNZ: pc = rb ? reg[ra] : pc; break;
    ...
  }
}
```



## DSL specification... Verilog, Chisel

```
// Instruction-set grammar
opn my_core (arith) {
  module control(clk,
    bus_addr,
    bus_read, bus_write, bus_state);
  do_fetch, do_regload, do_aluop, do_state);
  // Control module.
  include "parameters.vh"
  input wire clk;
  // Control module.
  include "parameters.vh"
  input wire clk;
}
```

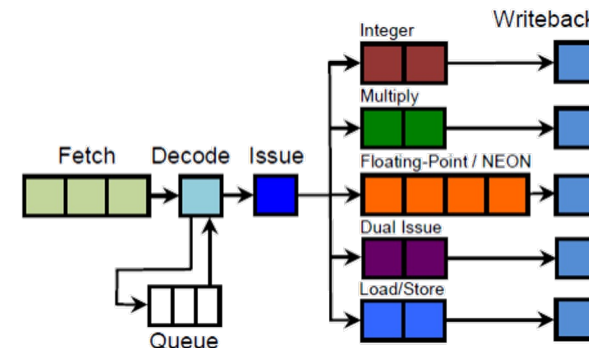
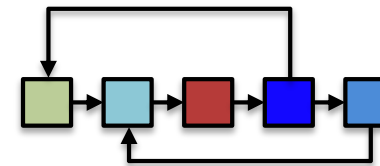
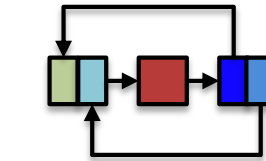
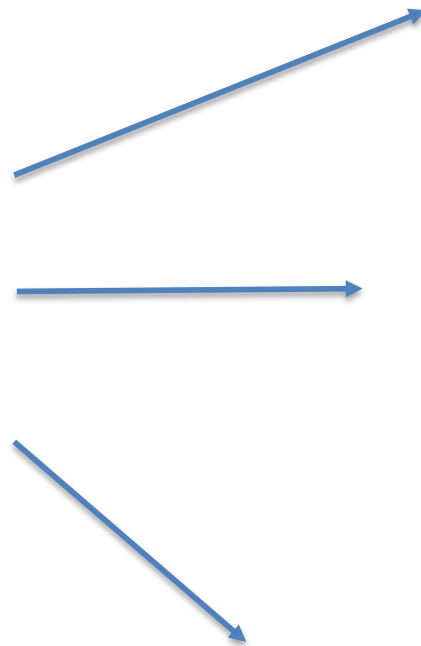
DSL = Domain Specific Languages

# Punchline : let's try to *compile* CPUs

- Generate CPU circuits straight out from the ISA simulator!

```
while(1) {  
  IR = mem[PC], nextPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ: nextPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nextPC;  
}
```

Instruction Set Simulator (ISS)

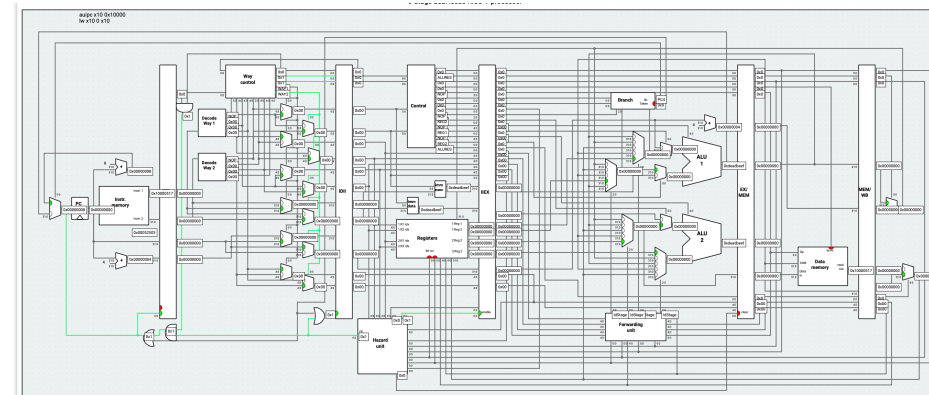
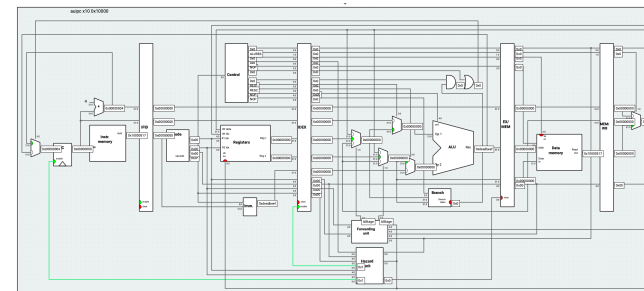
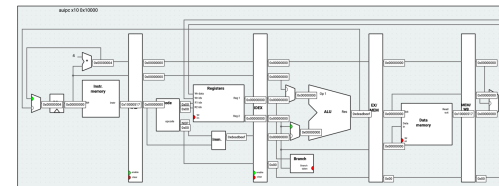
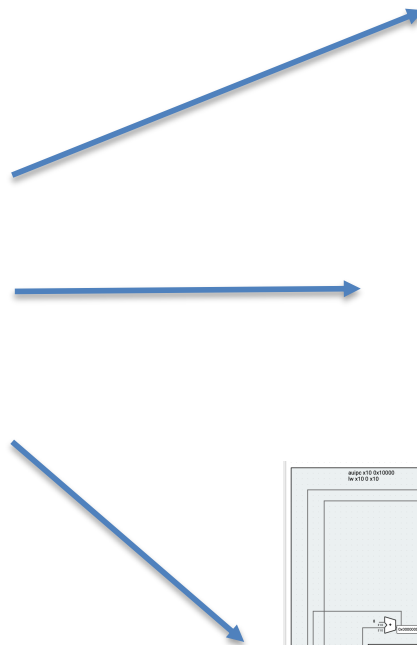


# Punchline : let's try to *compile* CPUs

- Generate CPU circuits straight out from the ISA simulator!

```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:  
    nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```

Instruction Set Simulator (ISS)



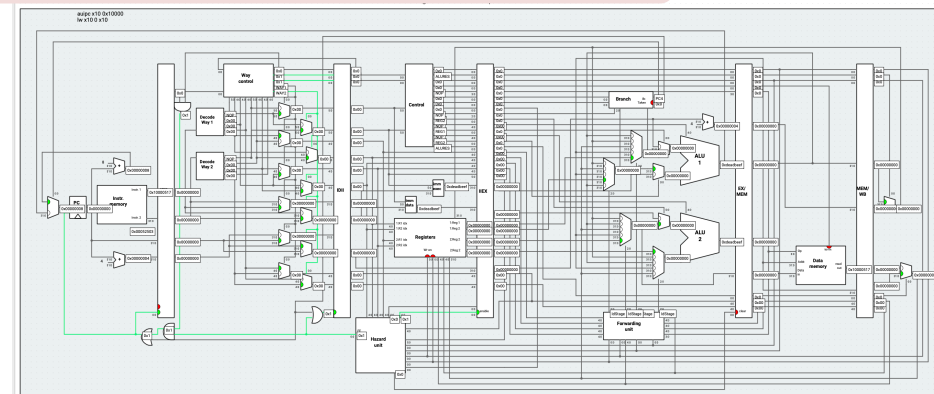
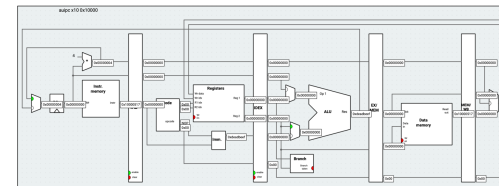
# Punchline : let's try to *compile* CPUs

- Generate CPU circuits straight out from the ISA simulator!

```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI:  
    break;  
    case LDW: X  
    break;  
    case STW: M  
    break;  
    case BEQ:  
    nxtPC += X  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```

Instruction Set Simulator (ISS)

**Key challenge**  
We need to infer a microarchitecture from the ISA specification



# Wait, can't I already do that ?

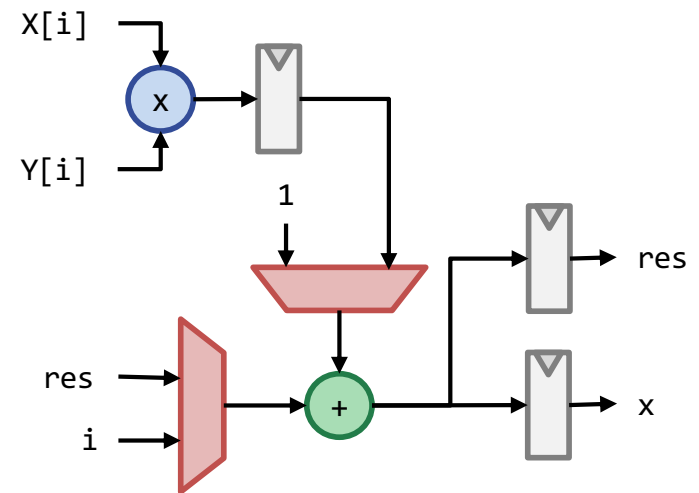
- Yes, High Level Synthesis tools generate circuits from C/C++

Synthesizing circuits from an algorithmic specification

```
int X[N], Y[N];
int tmp, res = 0;
#pragma HLS mult=1, adder=1
for(int i = 0; i < N; ++i) {
    tmp = X[i] * Y[i];
    res += tmp;
}
```

HLS

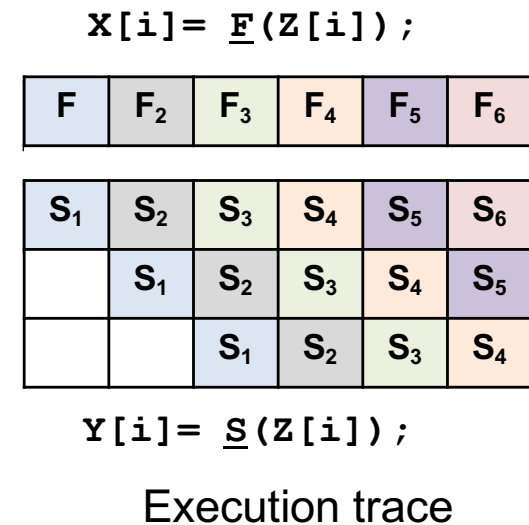
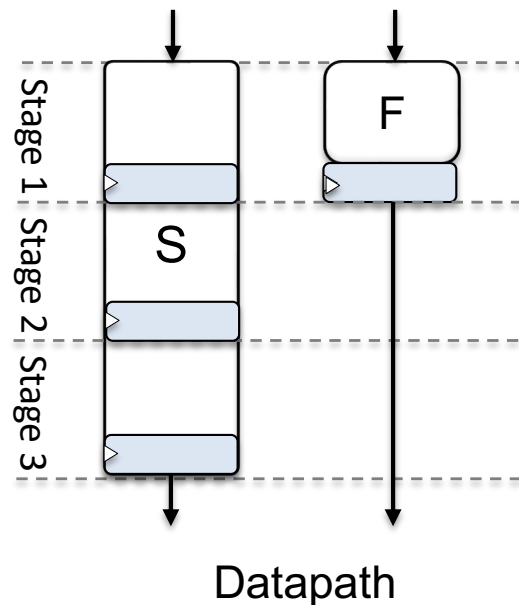
Efficient design synthesis for regular access patterns and computations



# Automatic pipelining in HLS

- Overlaps operations from distinct iterations

```
for (i=4; i<...; i++) {
    X[i] = F(Z[i]);
    Y[i] = S(Z[i]);
}
loop kernel
```



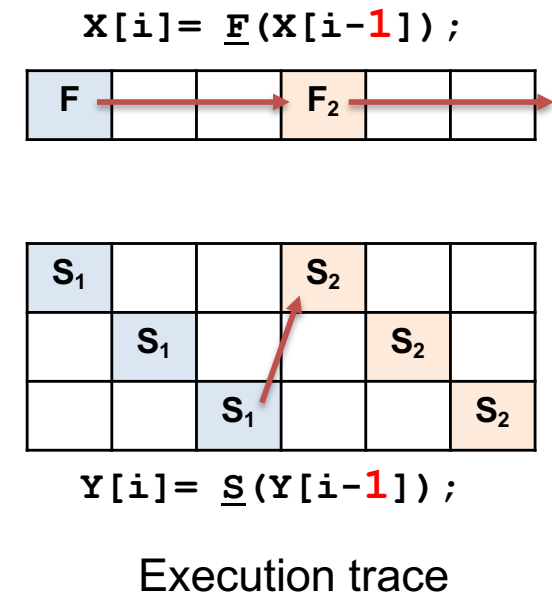
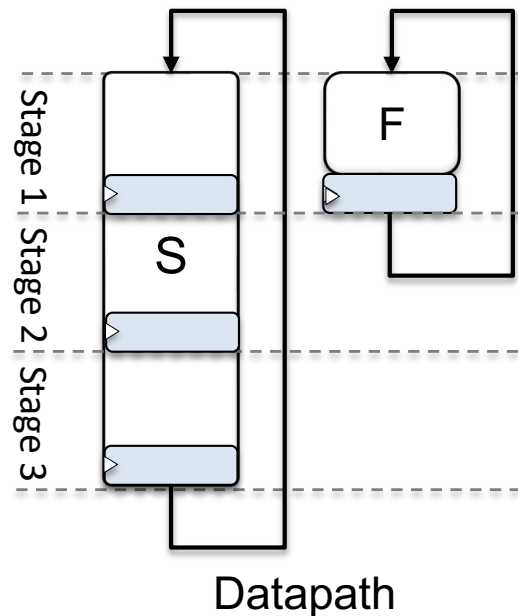
- May result in high clock speed, low CPI

# Pipelining and reuse distance

- Overlapping distinct iterations is not always possible

```
for (i=4; i<...; i++) {
  X[i] = F(X[i-1]);
  Y[i] = S(Y[i-1]);
}
```

loop kernel



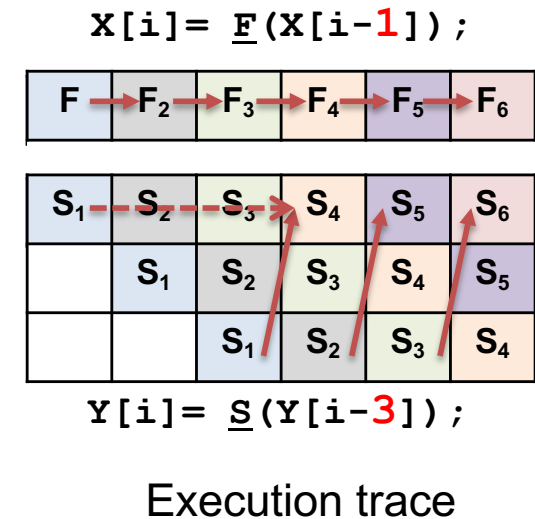
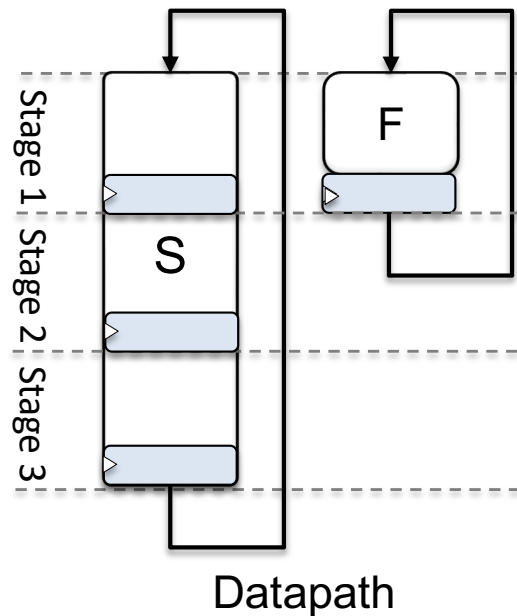
- Constrained by loop carried dependencies

# Pipelining and reuse distance

- What if the reuse distance is  $>1$  ?

```
for (i=4; i<...; i++) {
  X[i] = F(X[i-1]);
  Y[i] = S(Y[i-3]);
}
```

loop kernel



Reuse distance  $\Rightarrow$  pipeline depth  $\Rightarrow$  clock speed  $T_{clk}$

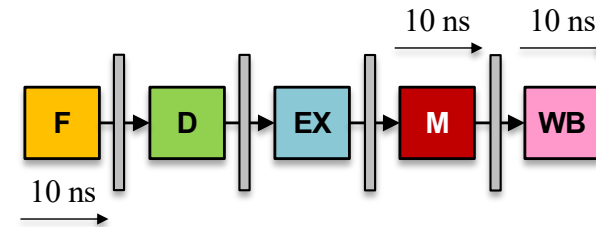
# Synthesizing an ISS using HLS

- Taking advantage of pipelining capabilities of HLS tools
  - Based on “software pipelining”, a well known compiler optimization

```

while(1) {
  IR = mem[PC], nextPC = PC + 4;
  a = rs1(IR), b = rs2(IR), d = rd(IR);
  switch (op(IR)) {
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);
              break;
    case MUL: X[d] = X[a]*X[b];
              break;
    case MULI: X[d] = X[a]*ext32(imm(IR));
              break;
    case LDW: X[d] = Mem[X[b]+imm(IR)];
              break;
    case STW: Mem[X[b]+imm(IR)] = X[a];
              break;
    case BEQ:
      nextPC += X[a]==X[b] ? imm(IR) : 0;
      break;
    // ...
  }
  PC = nextPC;
}
    
```

$T_{clk}=10ns$



	1	2	3	4	5	6	7	8	9	10	11	12	13
1	F	DC	EX	M	WB								
2		F	DC			EX	M	WB					
3			F			DC	EX	M	WB				
4						F	DC	EX	M	WB			
5							F	DC	EX	M	WB		
6								F	DC	EX	M	WB	
7									F	DC	EX	M	WB

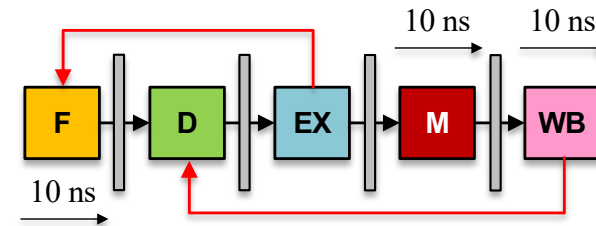
# Synthesizing an ISS using HLS

- Taking advantage of pipelining capabilities of HLS tools
  - Based on “software pipelining”, a well known compiler optimization

```

while(1) {
  IR = mem[PC], nextPC = PC + 4;
  a = rs1(IR), b = rs2(IR), d = rd(IR);
  switch (op(IR)) {
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);
              break;
    case MUL: X[d] = X[a]*X[b];
              break;
    case MULI: X[d] = X[a]*ext32(imm(IR));
              break;
    case LDW: X[d] = Mem[X[b]+imm(IR)];
              break;
    case STW: Mem[X[b]+imm(IR)] = X[a];
              break;
    case BEQ:
      nextPC += X[a]==X[b] ? imm(IR) : 0;
      break;
    // ...
  }
  PC = nextPC;
}
    
```

$T_{clk}=10ns$



	1	2	3	4	5	6	7	8	9	10	11	12	13
1	F	DC	EX	M	WB								
2						F	DC	EX	M	WB			
3											F	DC	EX
4													
5													

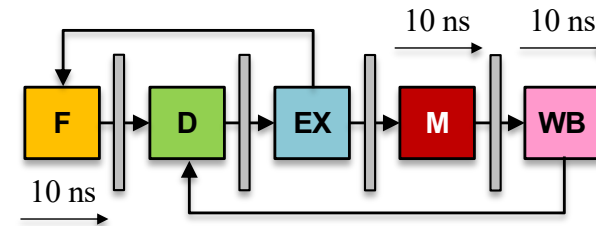
# Synthesizing an ISS using HLS

- Taking advantage of pipelining capabilities of HLS tools
  - Based on “software pipelining”, a well known compiler optimization

```

while(1) {
  IR = mem[PC], nextPC = PC + 4;
  a = rs1(IR), b = rs2(IR), d = rd(IR);
  switch (op(IR)) {
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);
              break;
    case MUL: X[d] = X[a]*X[b];
              break;
    case MULI: X[d] = X[a]*ext32(imm(IR));
              break;
    case LDW: X[d] = Mem[X[b]+imm(IR)];
              break;
    case STW: Mem[X[b]+imm(IR)] = X[a];
              break;
    case BEQ:
      nextPC += X[a]==X[b] ? imm(IR) : 0;
      break;
    // ...
  }
  PC = nextPC;
}
    
```

$T_{clk}=10ns$



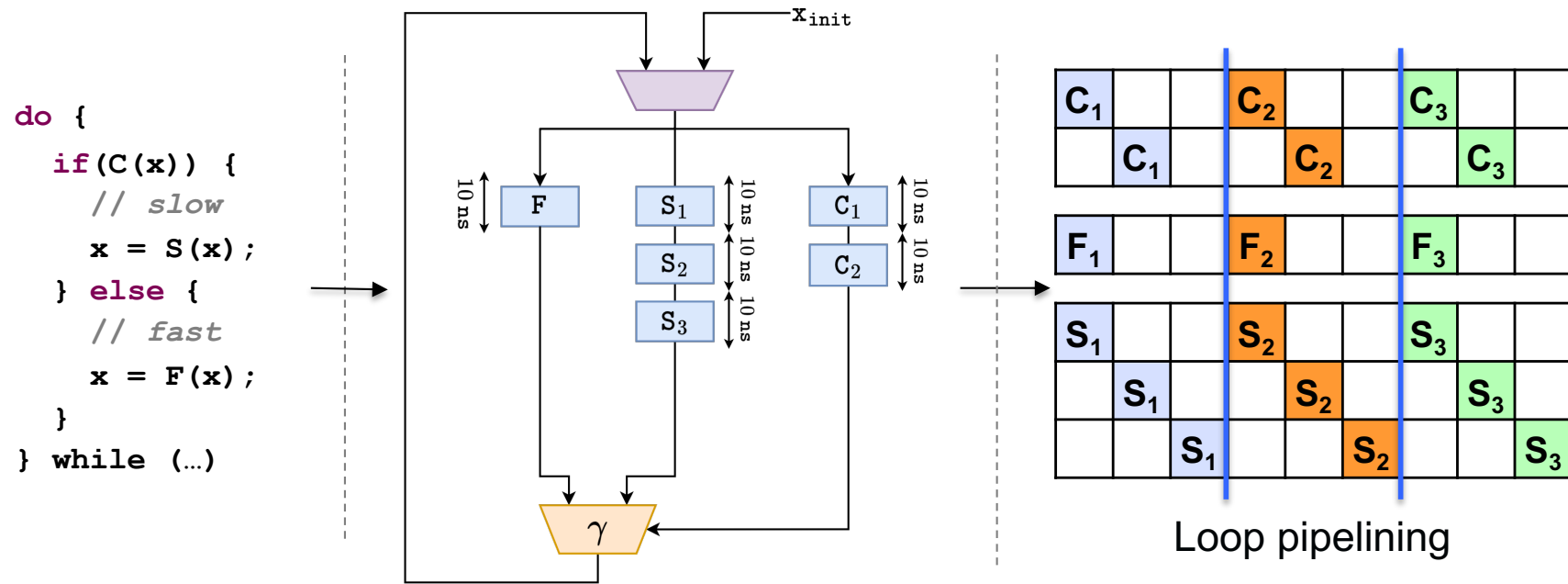
	1	2	3	4	5	6	7	8	9	10	11	12	13
Instruction 1	F	DC	EX	M	WB								
Instruction 2						F	DC	EX	M	WB			
Instruction 3											F	DC	EX
Instruction 4													
Instruction 5													

## Takeaway

Existing static HLS scheduling techniques are unfit for CPU design

# Speculative Loop Pipelining (SLP)

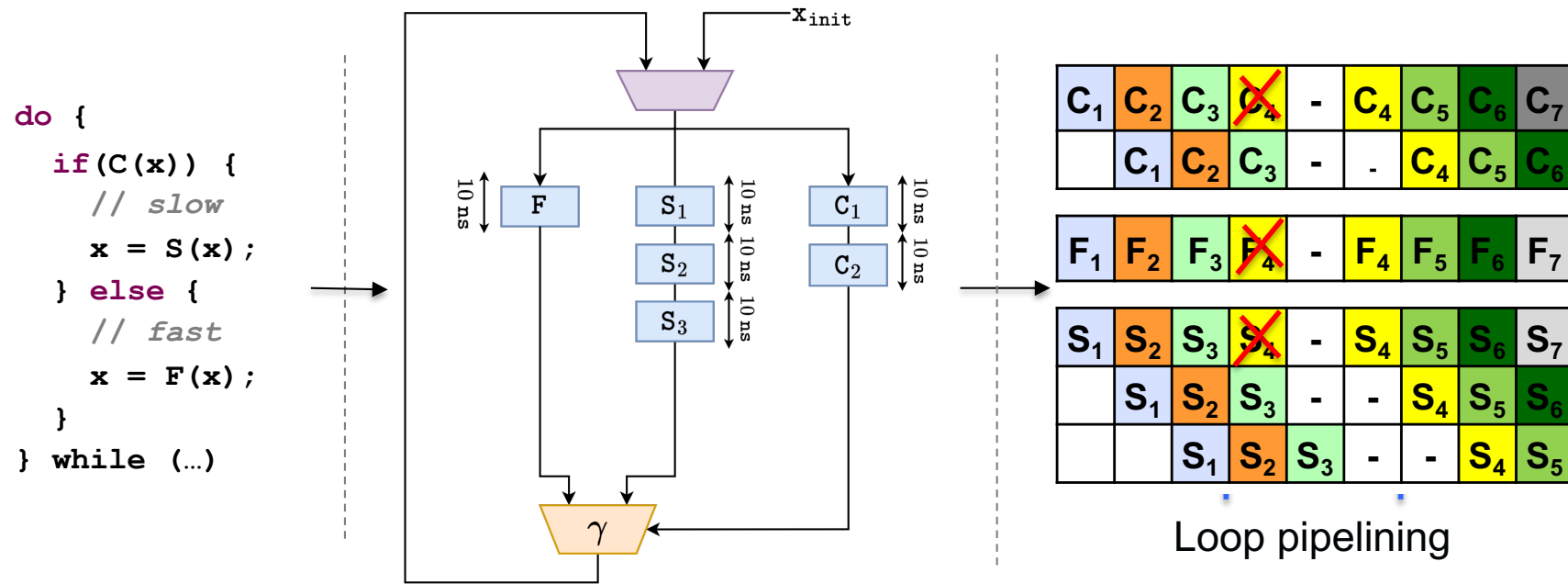
- Extends classic loop pipelining with speculative execution
  - Speculatively overlap the execution of distinct iterations



- Requires a complete rehaul of HLS scheduling techniques

# Speculative Loop Pipelining (SLP)

- Extends classic loop pipelining with speculative execution
  - Speculatively overlap the execution of distinct iterations



- Requires a complete overhaul of HLS scheduling techniques

Steven Derrien, Thibaut Marty, Simon Rokicki, and Tomofumi Yuki. "Toward Speculative Loop Pipelining for High-Level Synthesis". In: IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems 39.11 (Nov. 2020), pp. 4229–4239.

# SpecHLS : a tool for SLP in Vitis HLS

Input code

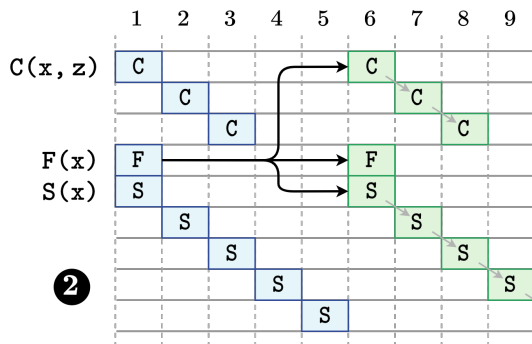
```
do {
  if(C(x, z))
    x = S(x);
  else
    x = F(x);
} while(!x);
```

1

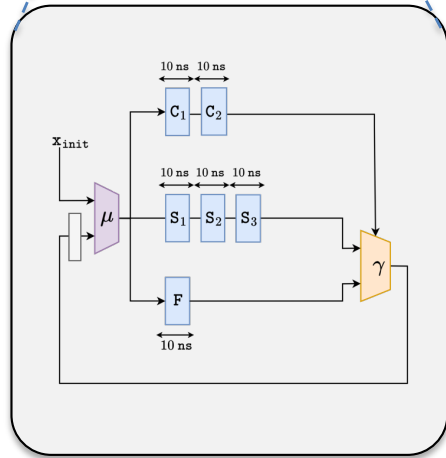


```
#pragma hls distance mis_x=5
#pragma hls distance ctrl=3
do {
  #pragma hls pipeline II=1
  ctrl[t] = C(s_x[t-1], z[t-1]);
  mis_x[t] = S(s_x[t-1]);
  cs = nextstate(cs, ctrl[t-3]);
  s_x[t] = cs.selSlow ?
    mis_x[t-5] :
    F(s_x[t-1]);
  if(cs.rollback) {
    s_x[t] = s_x[t-5];
    z[t] = z[t-5];
  }
  if(cs.commit)
    x = s_x[t-4];
  t += 1;
} while(!x && cs.commit);
```

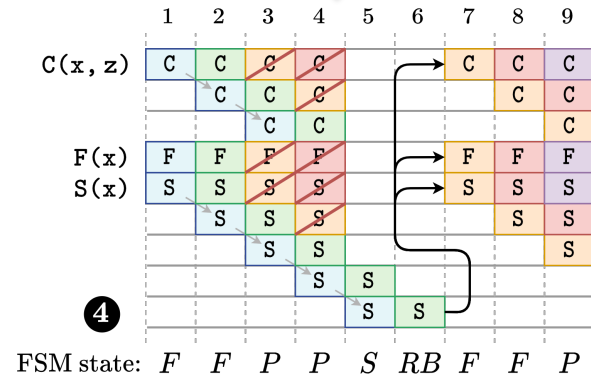
3



2



Program IR

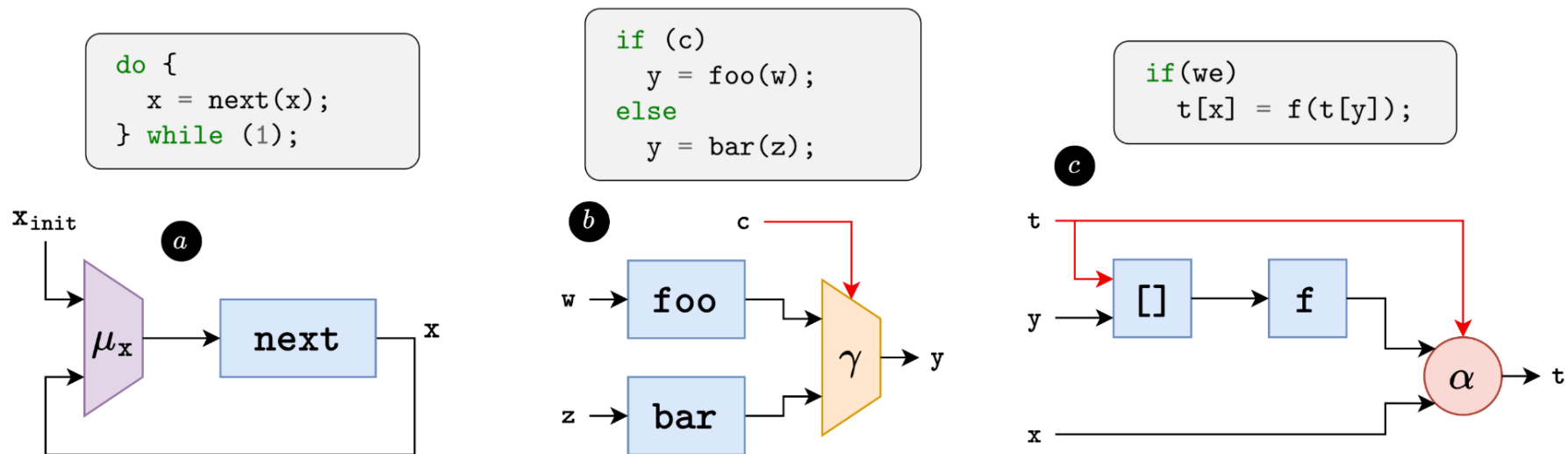


4



# Program representation in Spec-HLS

- The way the compiler sees the program to optimize**
  - Ours is a variant of Gated Static Single Assignment (GSSA)



<sup>1</sup>Peng Tu and David Padua, “Gated SSA-based demand-driven symbolic analysis for parallelizing compilers,” in *Proceedings of the 9th international conference on Supercomputing (ICS '95)*, 1995

**Takeaway**

In GSSA,  $\gamma$  nodes capture speculation opportunities

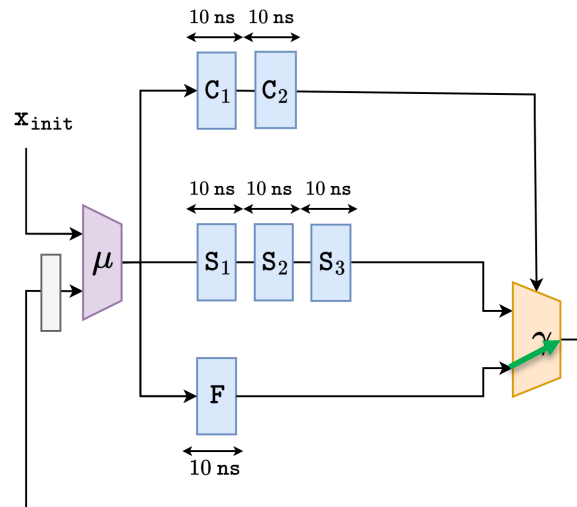
# Speculative Loop Pipelining (SLP)

- Building the program IR + identify speculation points

Input code

```
do {  
  if(C(x, z))  
    x = S(x);  
  else  
    x = F(x);  
} while(!x);
```

1



Speculation decision based on  $\gamma$  profiling

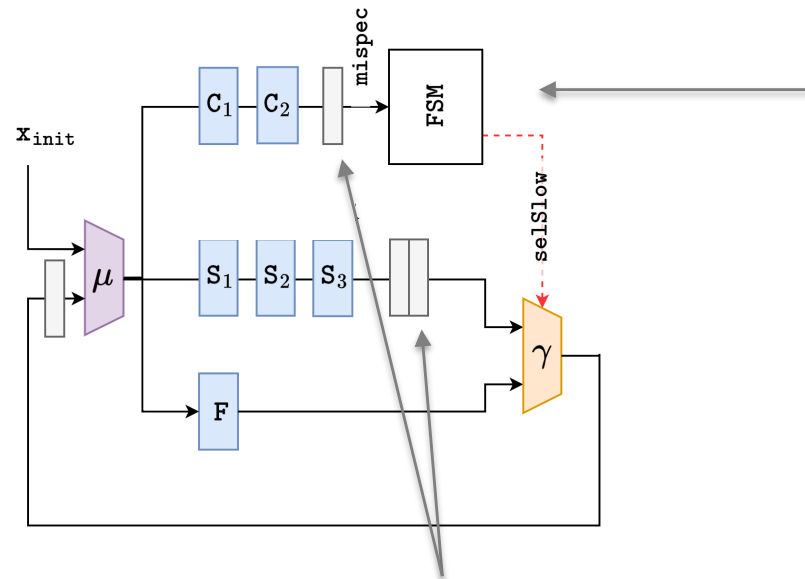
Steven Derrien, Thibaut Marty, Simon Rokicki, Tomofumi Yuki. Toward Speculative Loop Pipelining for High-Level Synthesis. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*. 2020

# Speculative Loop Pipelining (SLP)

- Adding delays + FSM to enable speculative pipelining

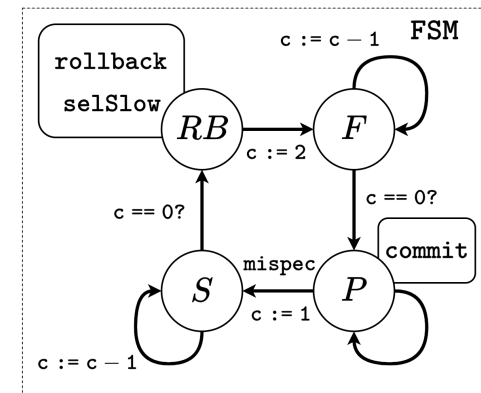
Input code

```
do {
  if(C(x, z))
    x = S(x);
  else
    x = F(x);
} while(!x);
```



Virtual pipeline stage registers

FSM



(RB) Rollback (F) Fill  
(S) Stall (P) Proceed

Steven Derrien, Thibaut Marty, Simon Rokicki, Tomofumi Yuki. Toward Speculative Loop Pipelining for High-Level Synthesis. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*. 2020

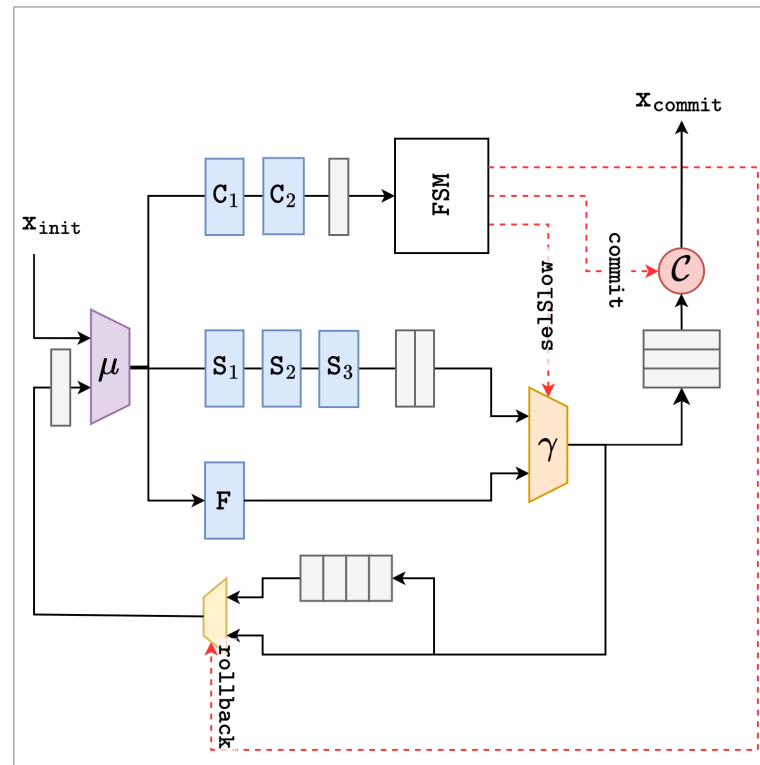
# Speculative Loop Pipelining (SLP)

- Insert commit + rollback logic to fully handle hazards

Input code

```
do {
  if(C(x, z))
    x = S(x);
  else
    x = F(x);
} while(!x);
```

1



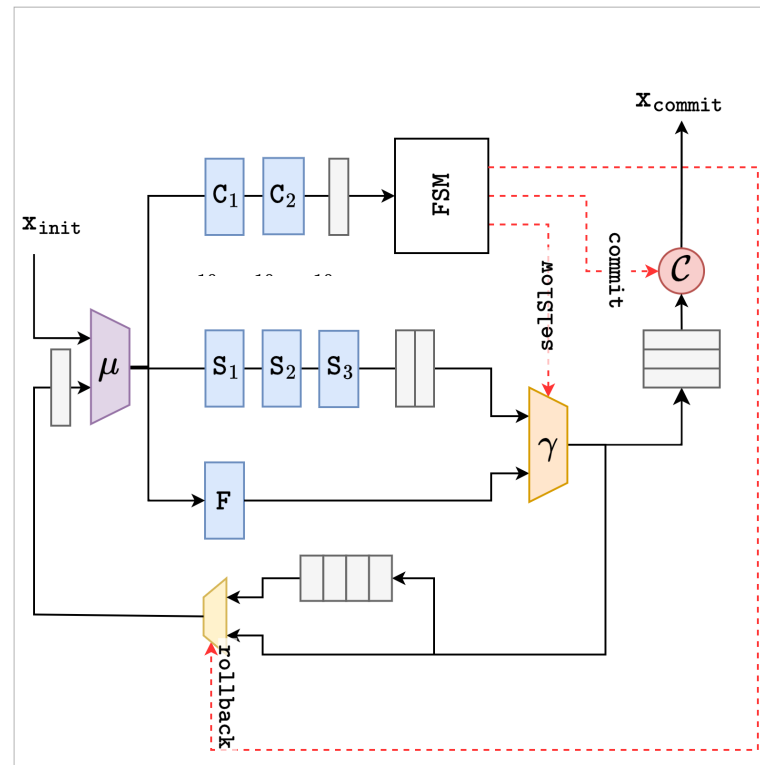
Steven Derrien, Thibaut Marty, Simon Rokicki, Tomofumi Yuki. Toward Speculative Loop Pipelining for High-Level Synthesis. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*. 2020

# Speculative Loop Pipelining (SLP)

- Regenerate C code and let HLS balance the pipeline stage

Input code

```
do {
  if(C(x, z))
    x = S(x);
  else
    x = F(x);
} while(!x);
```



Steven Derrien, Thibaut Marty, Simon Rokicki, Tomofumi Yuki. Toward Speculative Loop Pipelining for High-Level Synthesis. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*. 2020

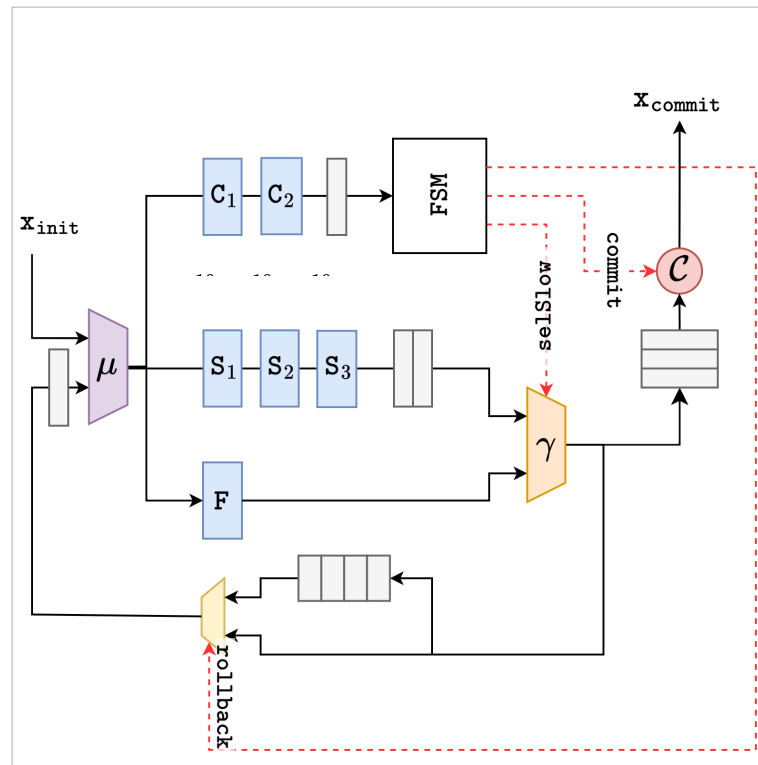
# Speculative Loop Pipelining (SLP)

- Regenerate (cycle accurate) C code

Input code

```
do {
  if(C(x, z))
    x = S(x);
  else
    x = F(x);
} while(!x);
```

1



```
#pragma hls distance mis_x=5
#pragma hls distance ctrl=3
do {
  #pragma hls pipeline II=1
  ctrl[t] = C(s_x[t-1], z[t-1]);
  mis_x[t] = S(s_x[t-1]);
  cs = nextstate(cs, ctrl[t-3]);
  s_x[t] = cs.selSlow ?
           mis_x[t-5] :
           F(s_x[t-1]);
  if(cs.rollback) {
    s_x[t] = s_x[t-5];
    z[t] = z[t-5];
  }
  if(cs.commit)
    x = s_x[t-4];
  t += 1;
} while(!x && cs.commit);
```

3

Steven Derrien, Thibaut Marty, Simon Rokicki, Tomofumi Yuki. Toward Speculative Loop Pipelining for High-Level Synthesis. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*. 2020

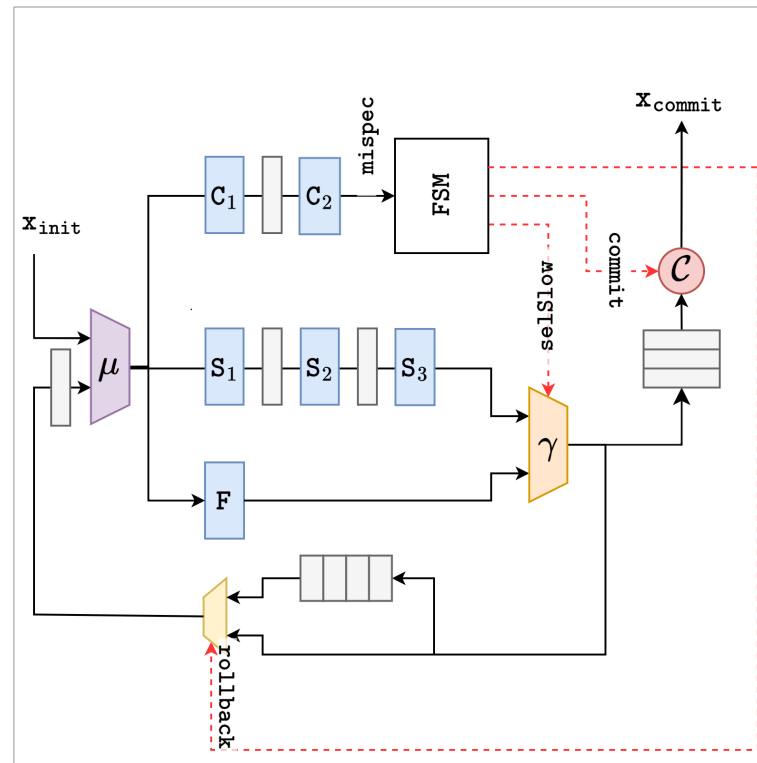
# Speculative Loop Pipelining (SLP)

- Let HLS balance pipeline stages

Input code

```
do {
  if(C(x, z))
    x = S(x);
  else
    x = F(x);
} while(!x);
```

1



```
#pragma hls distance mis_x=5
#pragma hls distance ctrl=3
do {
  #pragma hls pipeline II=1
  ctrl[t] = C(s_x[t-1], z[t-1]);
  mis_x[t] = S(s_x[t-1]);
  cs = nextstate(cs, ctrl[t-3]);
  s_x[t] = cs.selSlow ?
           mis_x[t-5] :
           F(s_x[t-1]);
  if(cs.rollback) {
    s_x[t] = s_x[t-5];
    z[t] = z[t-5];
  }
  if(cs.commit)
    x = s_x[t-4];
  t += 1;
} while(!x && cs.commit);
```

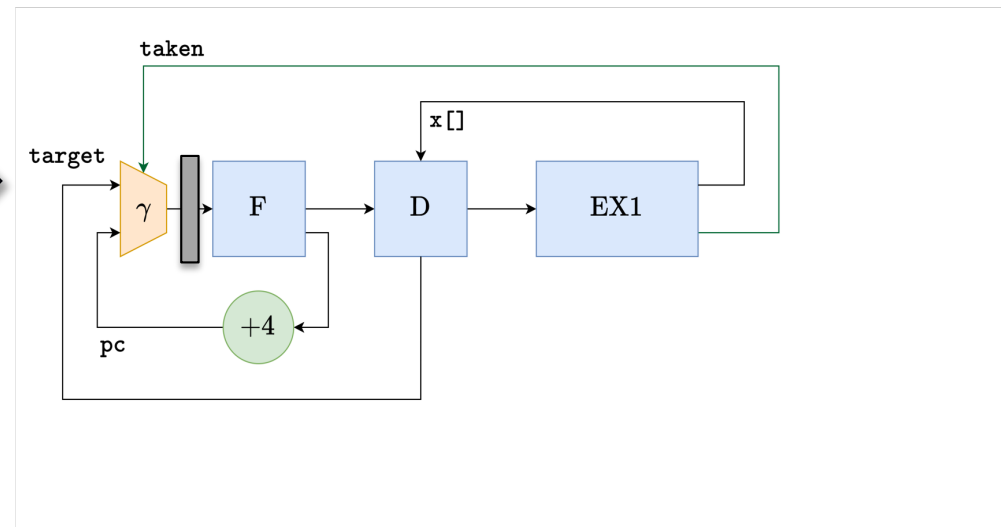
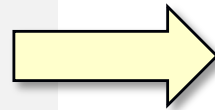
3

Steven Derrien, Thibaut Marty, Simon Rokicki, Tomofumi Yuki. Toward Speculative Loop Pipelining for High-Level Synthesis. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*. 2020

# Synthesizing RISC-V cores from ISSs

- Instruction Set Simulator = behavioral description

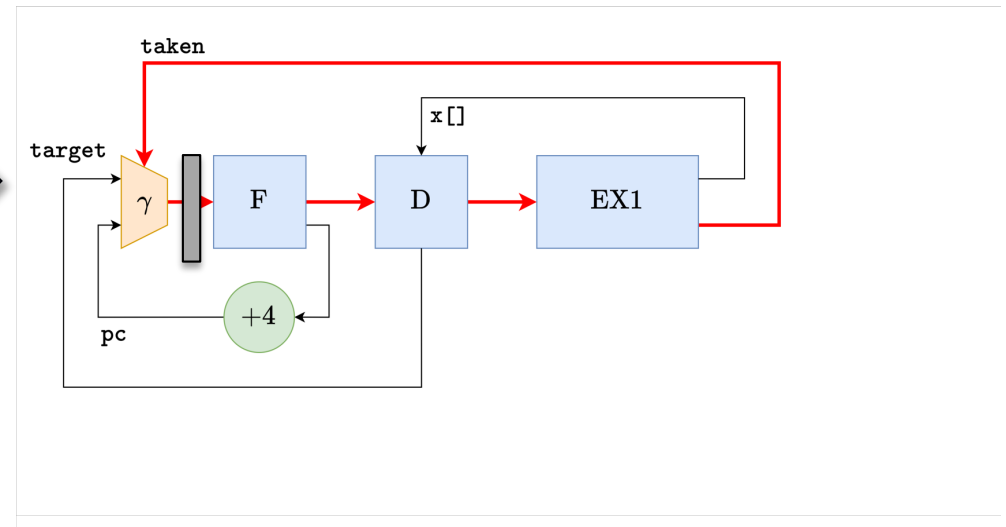
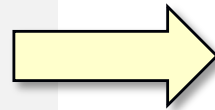
```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:   
      nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```



# Synthesizing RISC-V cores from ISSs

- Instruction Set Simulator = behavioral description

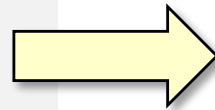
```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:   
      nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```



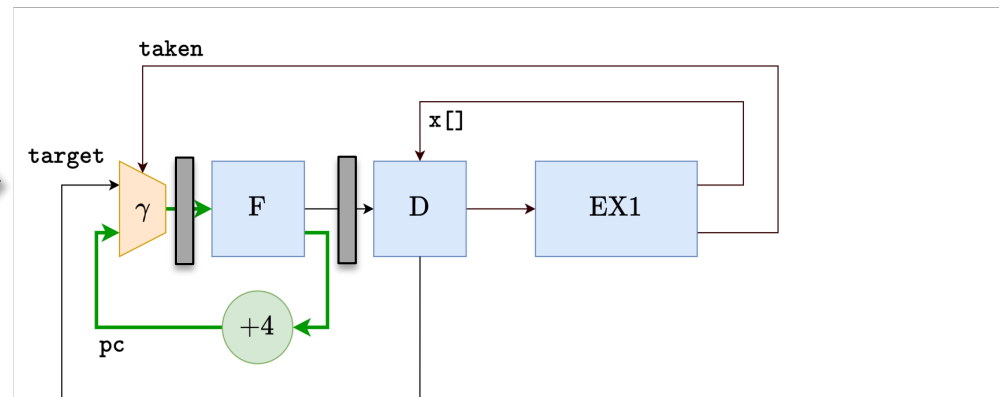
# Synthesizing RISC-V cores from ISSs

- Instruction Set Simulator = behavioral description

```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:  
    nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```



We can speculate that current instruction *is not* a branch



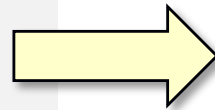
We then have a 1 cycle penalty for every executed branch

This is a 2 stage pipelined CPU

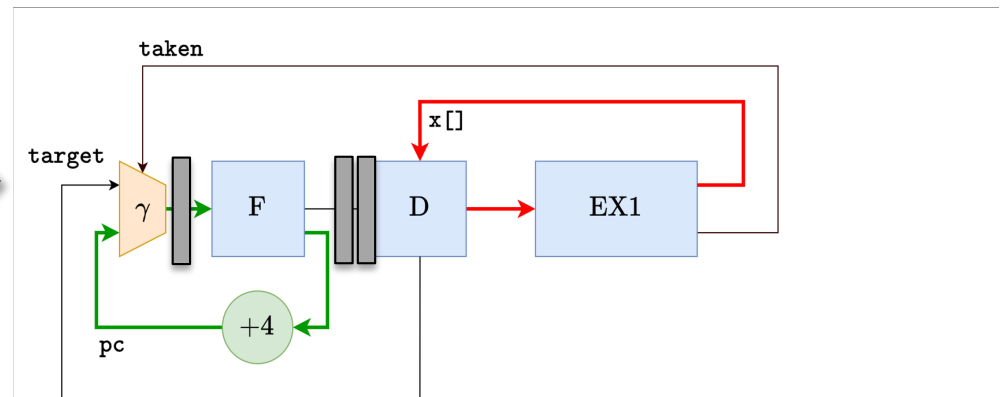
# Synthesizing RISC-V cores from ISSs

- Instruction Set Simulator = behavioral description

```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:  
    nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```



Can we obtain a 3 stage pipeline by increasing the penalty ?



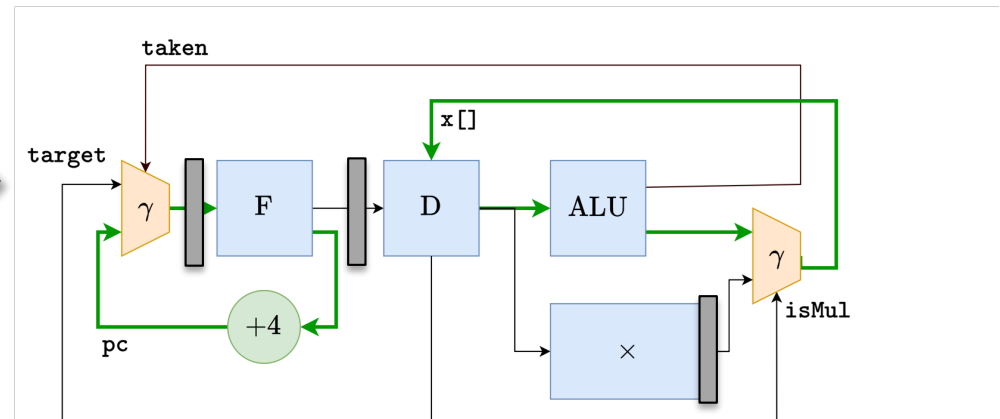
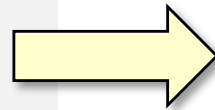
No, because of the dependency on register file  $x[]$



# Synthesizing RISC-V cores from ISSs

- Instruction Set Simulator = behavioral description

```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:  
      nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```

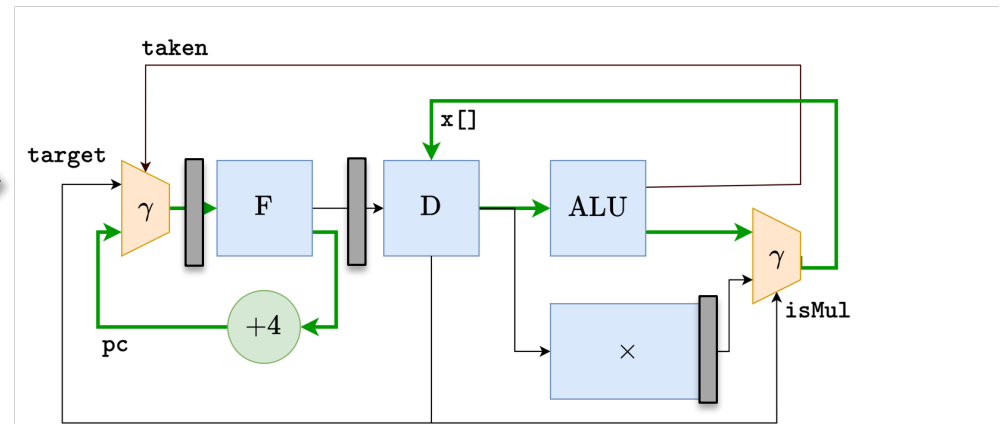
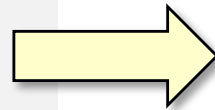


We can speculate that current instruction *is neither* a branch or a multiplication.

# Synthesizing RISC-V cores from ISSs

- Instruction Set Simulator = behavioral description

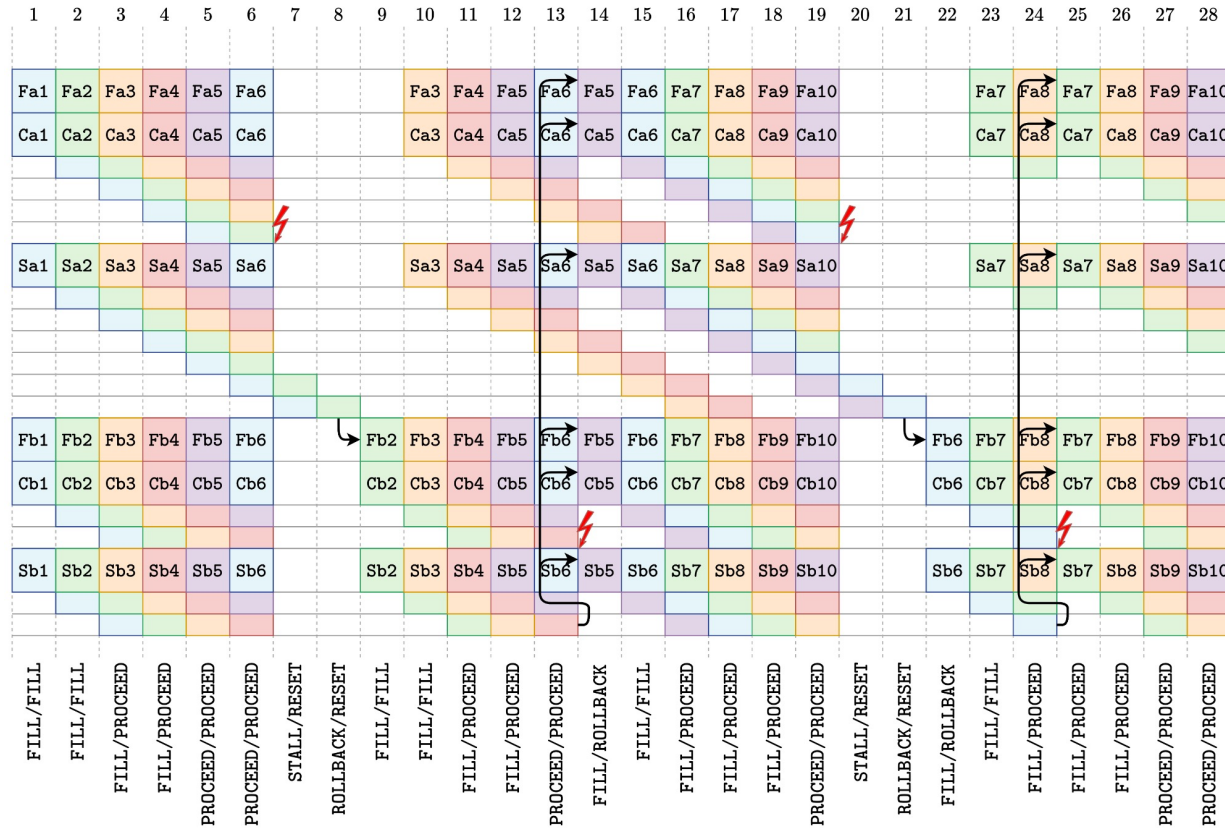
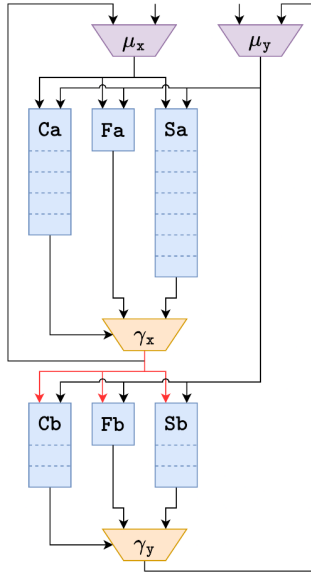
```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:   
      nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC =
```



## Challenge 1

But now, we have to speculate over multiple  $\gamma$  nodes ...

# Challenge 1 : multiple speculations

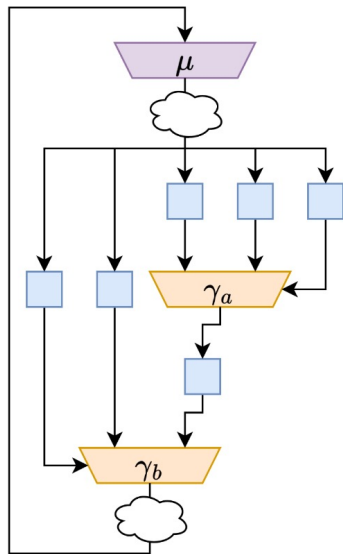


**Interacting speculations**  
 Speculative values of x influence the control decision on y

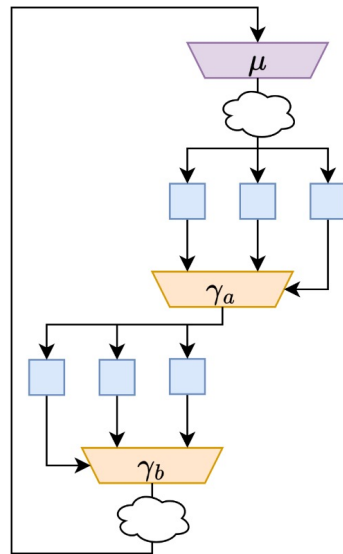
Jean-Michel Gorius, Simon Rokicki, and Steven Derrien. "SpecHLS: Speculative Accelerator Design Using High-Level Synthesis". In: IEEE Micro 42.5 (Sept. 2022), pp. 99–107. doi: [10.1109/MM.2022.3188136](https://doi.org/10.1109/MM.2022.3188136).

# Identifying *speculation patterns*

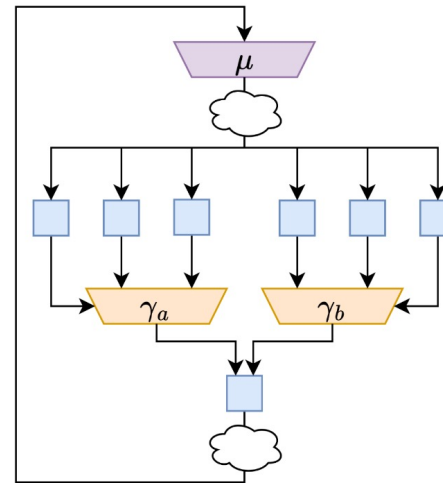
We identify four distinct *speculation patterns*



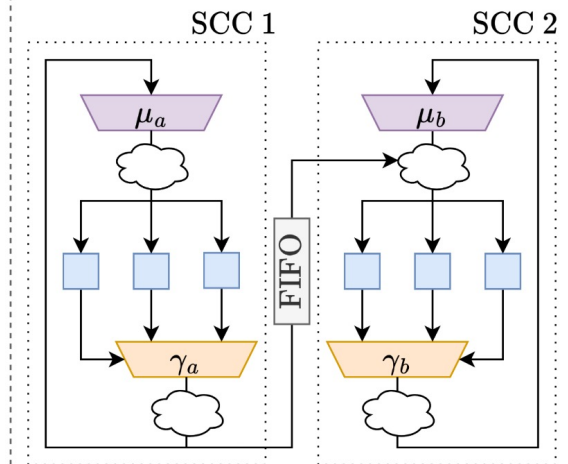
❶ data domination



❷ control domination



❸ independent speculations



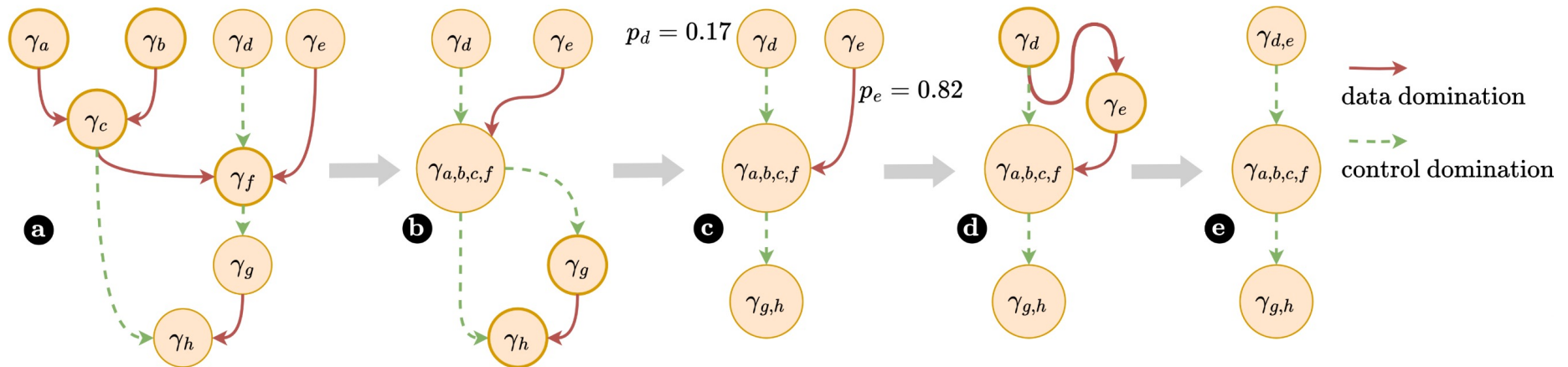
❹ decoupled SCCs

Jean-Michel Gorius, Simon Rokicki, and Steven Derrien. “SpecHLS: Speculative Accelerator Design Using High-Level Synthesis”. In: IEEE Micro 42.5 (Sept. 2022), pp. 99–107. doi: [10.1109/MM.2022.3188136](https://doi.org/10.1109/MM.2022.3188136).

# Iterative Pattern Resolution

**Pattern resolution**

- Iterative algorithm to synthesize the speculation control FSM
- Based on a set of **rewriting rules**, coupled to **profiling**



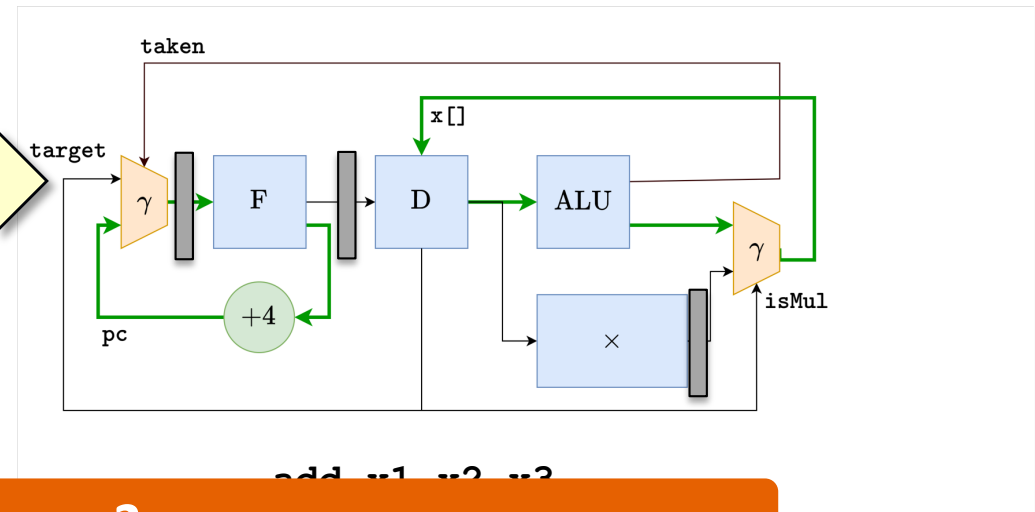
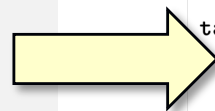
Jean-Michel Gorius, Simon Rokicki, and Steven Derrien. "SpecHLS: Speculative Accelerator Design Using High-Level Synthesis". In: IEEE Micro 42.5 (Sept. 2022), pp. 99–107. doi: [10.1109/MM.2022.3188136](https://doi.org/10.1109/MM.2022.3188136).



# Synthesizing RISC-V cores from ISSs

- Instruction Set Simulator = behavioral description

```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:  
    nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```



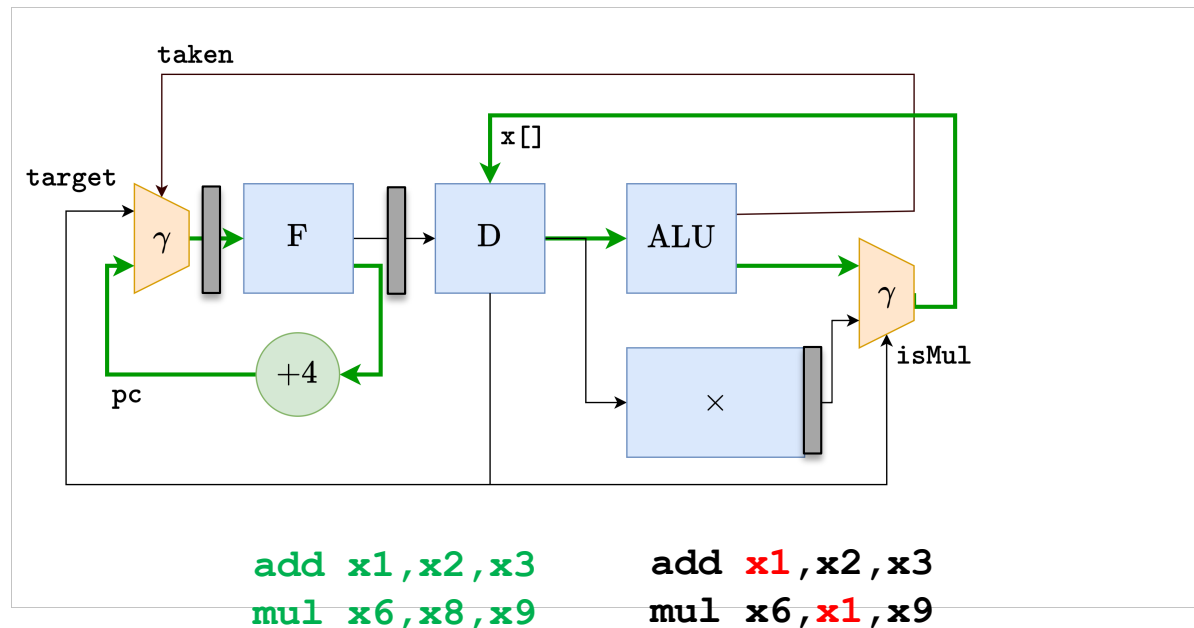
## Challenge 2

How can we speculate on RAW dependencies ?

etter !

# Challenge 2 : memory speculation

- In our program IR Arrays are modeled as values
  - We can check at run-time which version of the array value is needed

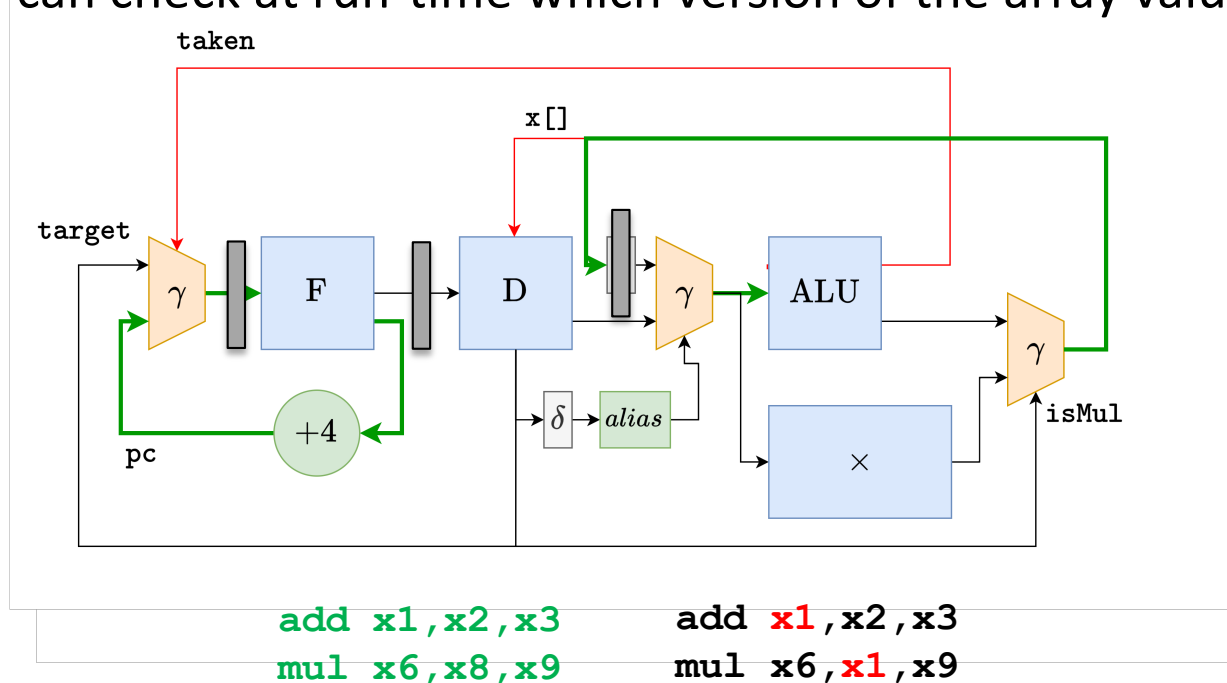


We can now speculate that there is no RAW dependency on `x[]` ...

Jean-Michel Gorius, Simon Rokicki, and Steven Derrien. "A Unified Memory Dependency Framework for Speculative High-Level Synthesis". International Conference on Compiler Construction. CC 2024. New York, NY, USA: Association for Computing Machinery, Feb. 20, 2024, pp. 13–25.

# Challenge 2 : memory speculation

- In our program IR Arrays are modeled as values
  - We can check at run-time which version of the array value is needed



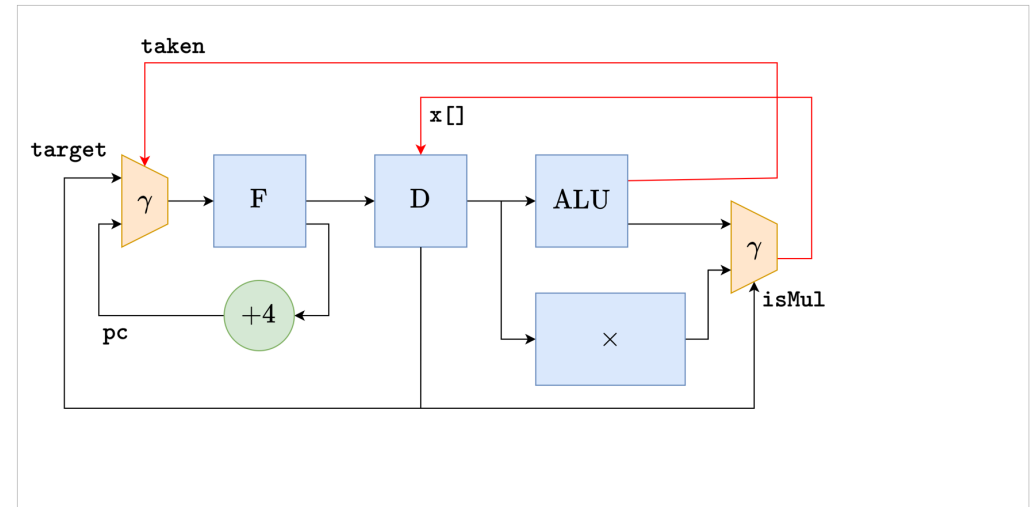
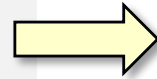
We can now speculate that there is no RAW dependency on x[] ...

Jean-Michel Gorius, Simon Rokicki, and Steven Derrien. "A Unified Memory Dependency Framework for Speculative High-Level Synthesis". International Conference on Compiler Construction. CC 2024. New York, NY, USA: Association for Computing Machinery, Feb. 20, 2024, pp. 13–25.

# Synthesizing RISC-V cores from ISSs

- Instruction Set Simulator = behavioral description

```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:  
    nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```



# Synthesizing RISC-V cores from ISSs

- **Instruction Set Simulator = behavioral description**

```
while(1) {  
  IR = mem[PC], nxtPC = PC + 4;  
  a = rs1(IR), b = rs2(IR), d = rd(IR);  
  switch (op(IR)) {  
    case ADD: X[d] = ALU(op(IR), X[a], X[b]);  
    break;  
    case MUL: X[d] = X[a]*X[b];  
    break;  
    case MULI: X[d] = X[a]*ext32(imm(IR));  
    break;  
    case LDW: X[d] = Mem[X[b]+imm(IR)];  
    break;  
    case STW: Mem[X[b]+imm(IR)] = X[a];  
    break;  
    case BEQ:  
    nxtPC += X[a]==X[b] ? imm(IR) : 0;  
    break;  
    // ...  
  }  
  PC = nxtPC;  
}
```



16  $\gamma$  nodes =  
 $\sim 10^9$  solutions to explore

## Challenge 3

How do we figure out which  $\gamma$  nodes to speculate on ?

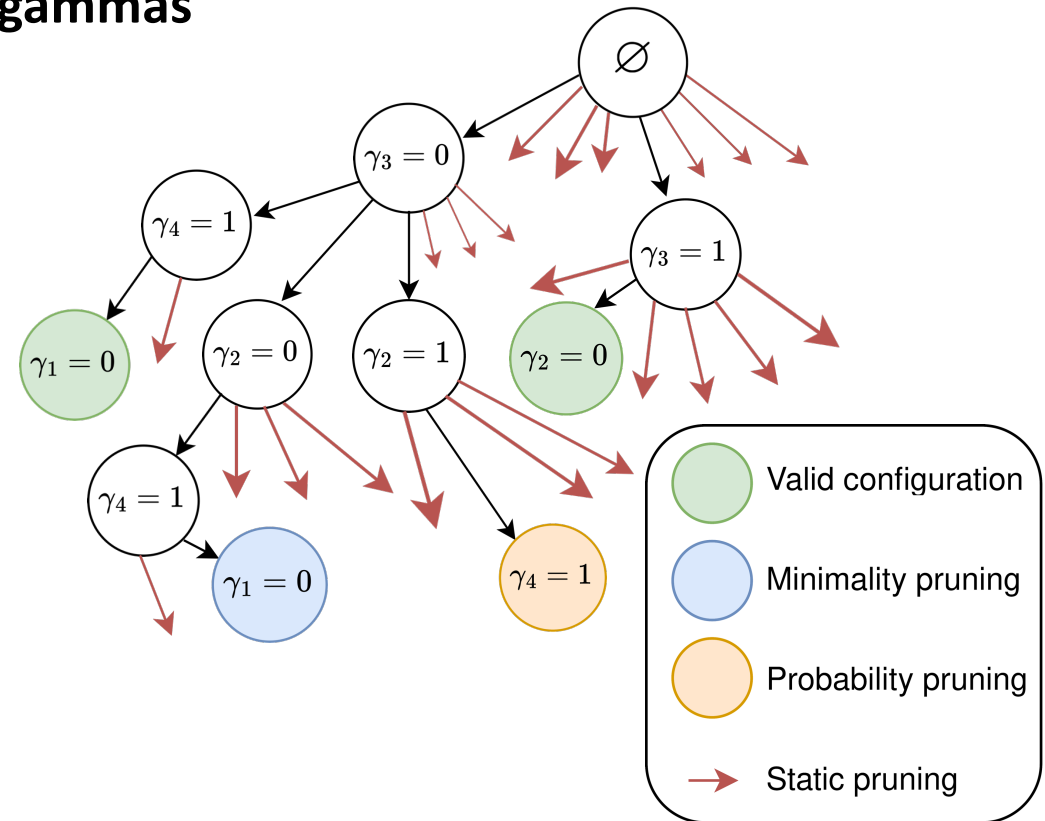
# Challenge 3: which gamma to speculate ?

- **Solution : a branch and bound algorithm**
  - **Use profiling information for gammas**

Branch-and-bound algorithm find every valid configurations

Configuration sorted by increasing speculation profitability (gain + probability)

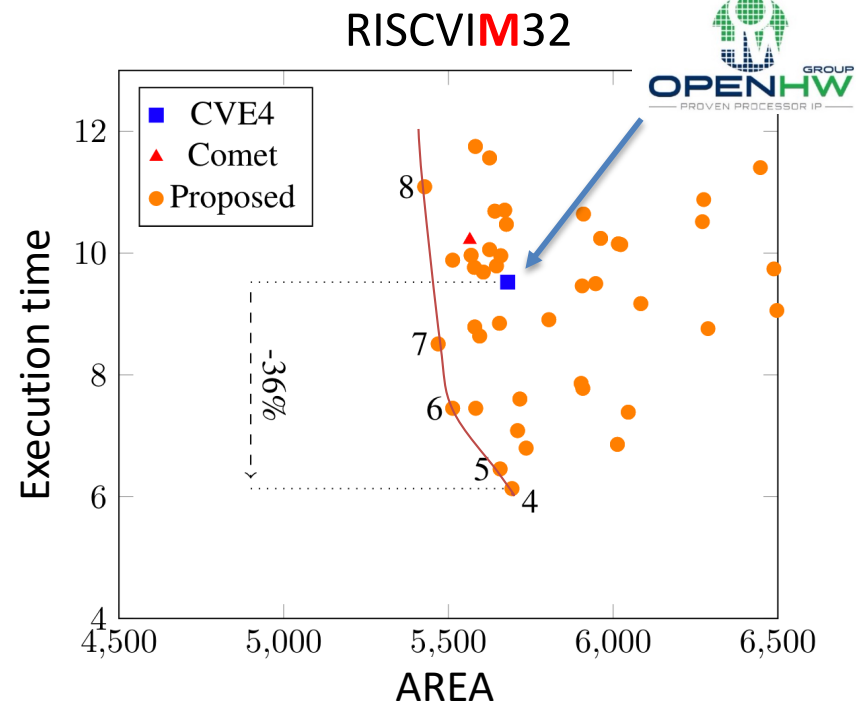
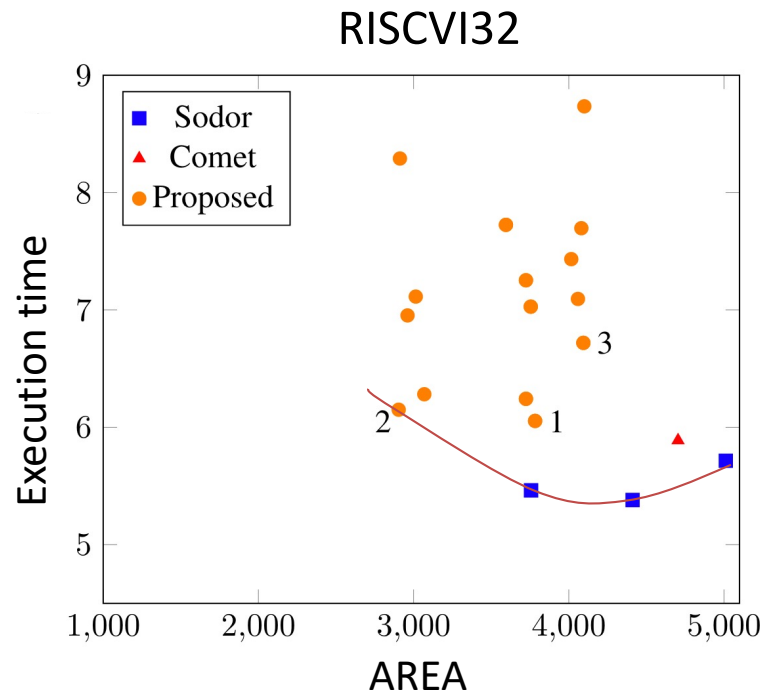
We can efficiently prune the exploration tree



Leothaud, Dylan, et al. "Efficient Design Space Exploration for Dynamic & Speculative High-Level Synthesis." *Field-Programmable Logic and Applications (FPL)*. IEEE, 2024.

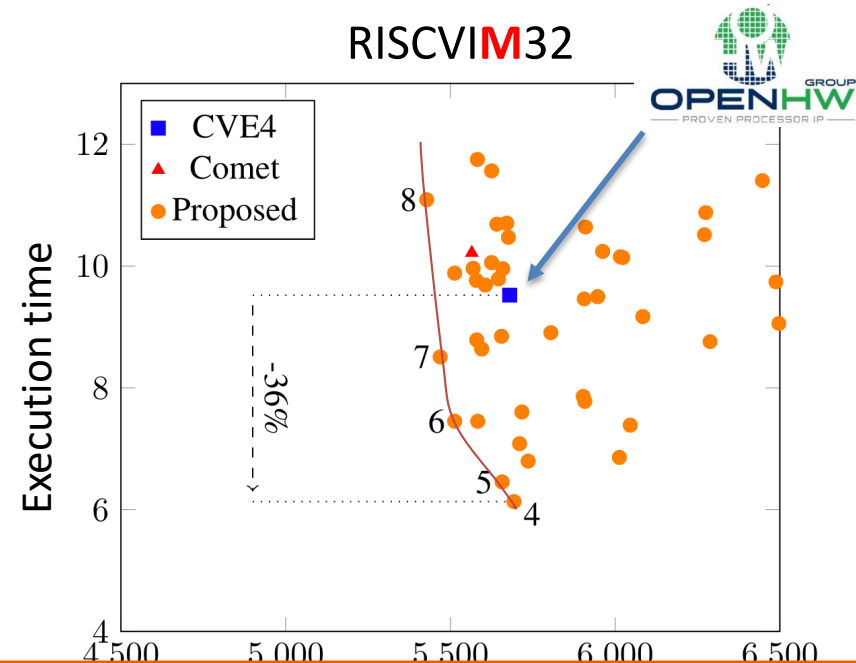
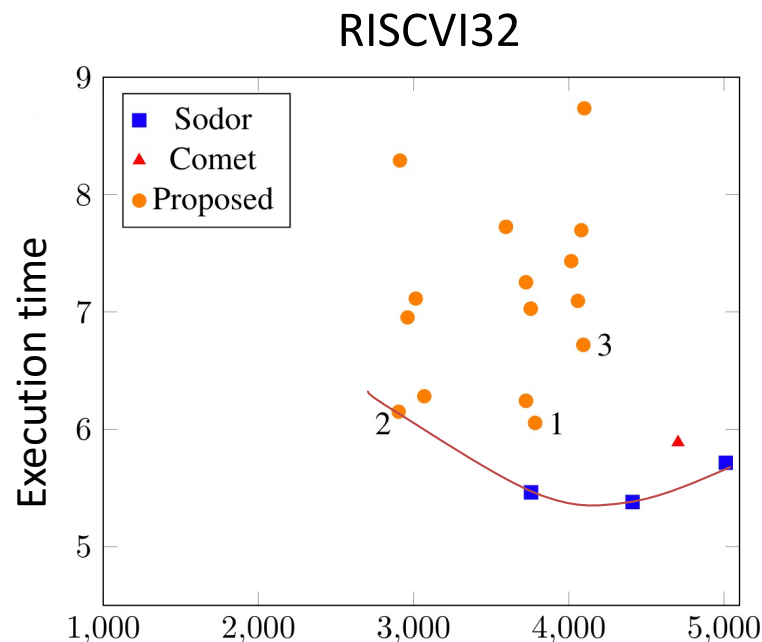
# Early experimental results

- **Evaluation on two simple RISC-V ISAs (RV132/RVIM32)**
  - Searching for Pareto Optimal solutions (performance/area trade-off)
  - Compared against manual designs (CVE4, Sodor)



# Early experimental results

- **Evaluation on two simple RISC-V ISAs (RV132/RVIM32)**
  - Searching for Pareto Optimal solutions (performance/area trade-off)
  - Compared against manual designs (CVE4, Sodor)



## Takeaway

Ability to automatically derive competitive CPU designs

# Future work

- **Pushing toward more complex microarchitectures**
  - Dynamic branch prediction, cache, OoO
- **Formal verification of pipeline hazard logic**
  - Translation validation for SLP
- **Enforcing timing predictability for Real-time systems**
  - Automatic extraction of WCET pipeline models
  - Checking for timing-anomalies
- **Hardening against faults for security**
  - Partial duplication, error detectors/correctors
  - Automatic insertion of CFI/DFI mechanisms

# Current & Future work

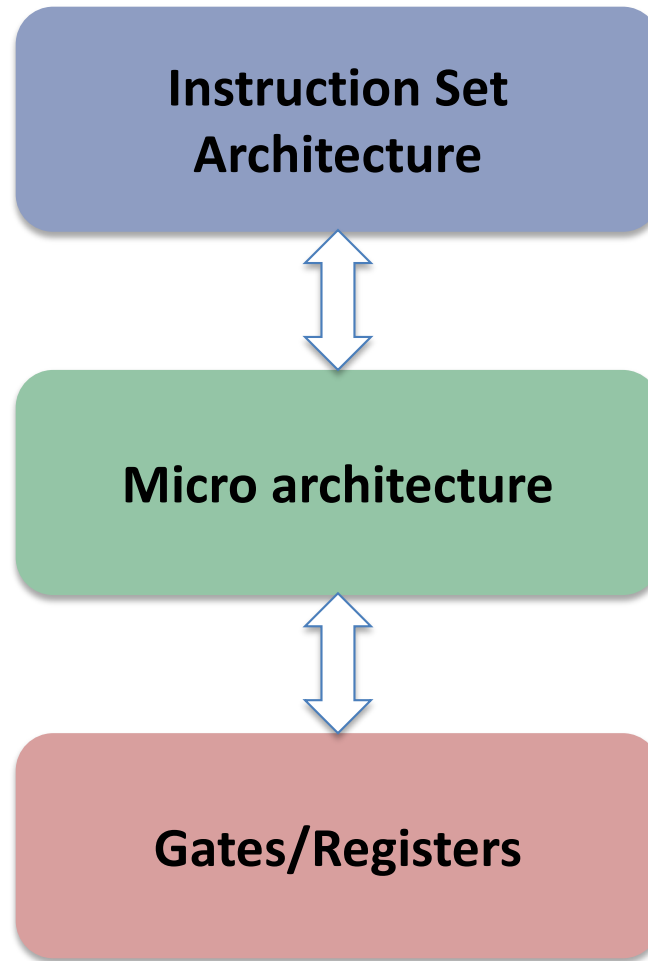
## THE LORD OF THE RISC-V

*One flow to rule them all,  
one flow to check them,  
One flow to cover them all,  
and to the silicon bind them*

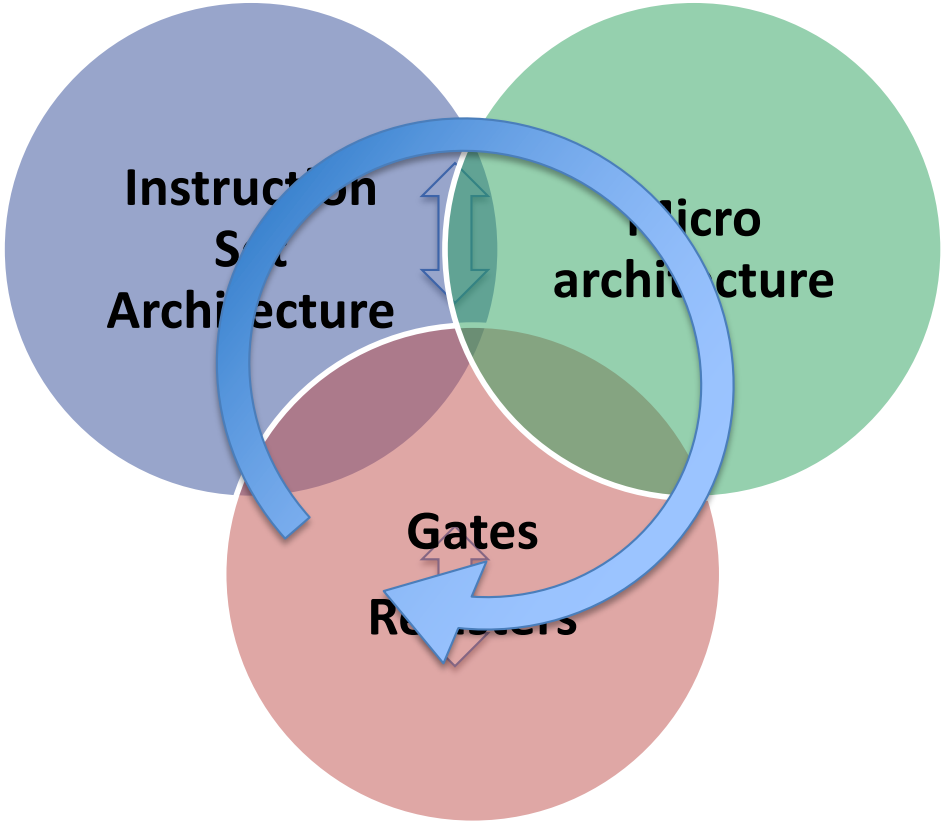


<https://lotr.gitlabpages.inria.fr/website/>

# Take home message



# Take home message



# Modeling a Dynamic Branch Predictor

```
k = hash(pc);
pred_pc = pc + 4;
btb_hit = btb_tag[k] == pc;
// ...
if (btb_hit && btb_state[k] >= WEAK_TAKEN)
    pred_pc = btb_target[k];

switch (instr->opcode) { /* ... */ }

if (pred_pc == next_pc) {
    // Good prediction, fast path
    incr(btb_state[k]);
    pc = pred_pc;
} else {
    // Bad prediction, slow path
    if (btb_hit)
        decr(btb_state[k]);
    else if (/* The instruction was a branch */)
        ↪ {
            btb_state[k] = WEAK_TAKEN;
            btb_tag[k] = pc;
            btb_target[k] = /* Branch target */;
        }
    pc = next_pc;
}
```

