# Rapid Cycle-Accurate Simulator for High-Level Synthesis (FLASH)
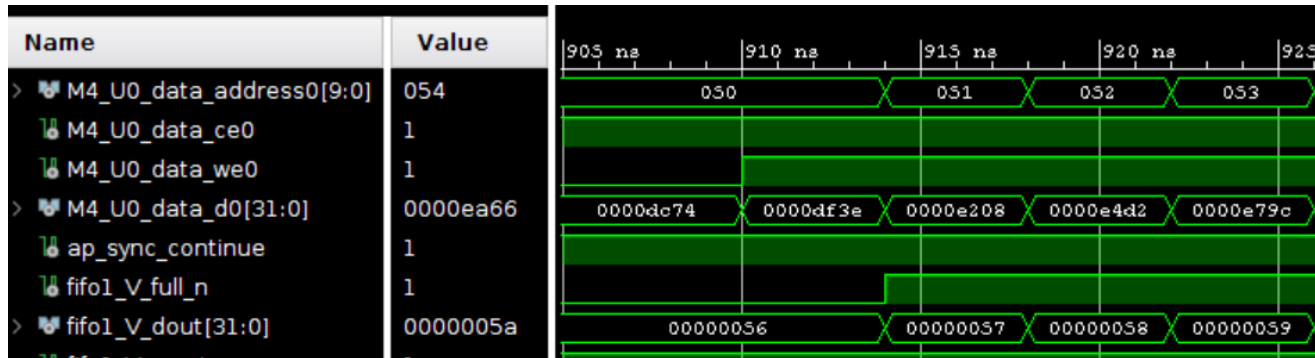
Yuze Chi, <u>Young-kyu Choi</u>, Jason Cong, and Jie Wang

University of California, Los Angeles

# Motivation

- ## RTL co-simulation for HLS



Too **slow**...
(ex matmul: 192s)

Difficult to understand

- ## SW simulation for HLS


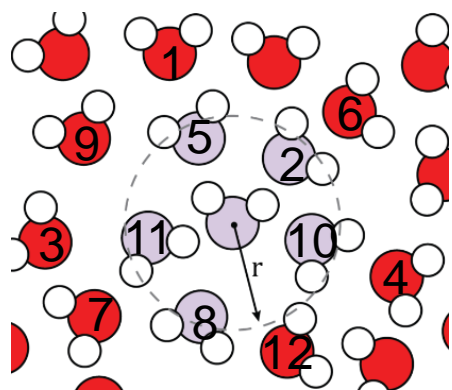
100X to 1000X faster than RTL co-sim
(ex matmul: 0.05s)

Easy to understand

- But can it measure the execution time?
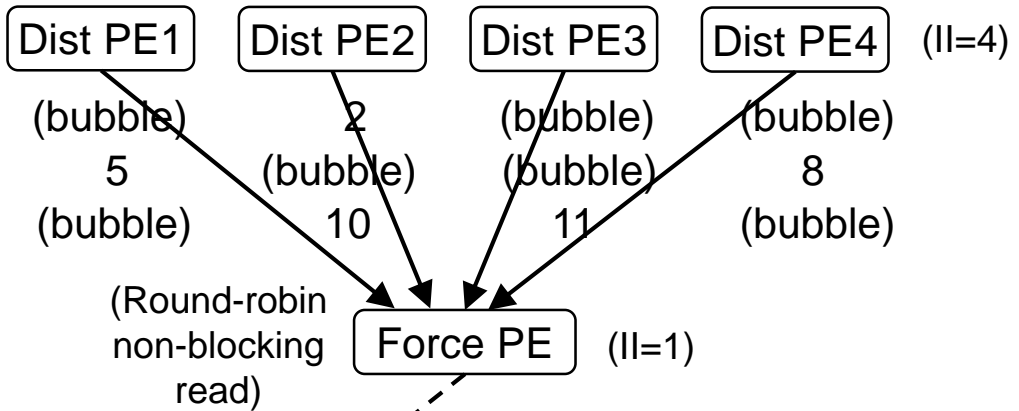- Is it producing the **correct** result?

- # HLS simulation of molecular dynamics



| Dist PE1 | Dist PE2 | Dist PE3 | Dist PE4 | (II=4) |
|---|---|---|---|---|

| | | | | |
|---|---|---|---|---|
| 1st round: | (bubble) | 2 | (bubble) | (bubble) |
| 2nd round: | 5 | (bubble) | (bubble) | 8 |
| 3rd round: | (bubble) | 10 | 11 | (bubble) |

(Round-robin non-blocking read) → Force PE (II=1)

< HLS C code>

```
#pragma HLS dataflow
Dist_PE1();
Dist_PE2();
Dist_PE3();
Dist_PE4();
Force_PE();
```

Simulated in instantiation order
→ Missing bubbles

RTL sim output: 2  5  8  10 11

Does not match!

SW sim output: 5  2  11 8  10

- – Reason

| | Xilinx Viv HLS C Sim | Intel OpenCL HLS Sim |
|---|---|---|
| FIFO depth | **Unlimited** | Exact |
| Exec model | **Sequential** | Concurrent |
| Feedback | **Not supported** | Supported |
| Sim speed | ~5 Mcycle/s | **~1 Mcycle/s** |
| Sim order | Deterministic | **Non-deterministic** |
| Cycle-acc | **Not cycle-accurate** | **Not cycle-accurate** |

Christophe Rowley, https://en.wikibooks.org/wiki/Molecular_Simulation/Radial_Distribution_Functions

- # Conventional simulation flows & proposed approach

<HLS design steps>

Allocation

Library

HLS C code → Compilation → Binding → Generation → RTL code

Scheduling

stmt,loop, func, …

Fast, but
1. Output may not be accurate
2. No perf estimation → SW simulator

Proposed simulator (FLASH)

scheduling info

RTL simulator

Accurate, but too slow

- # Overall simulation framework of FLASH*

Input:
Vivado HLS C code → Prepro-cessing → Sim File Generation (w/ ROSE) → New sim file → HLS C sim → Analysis
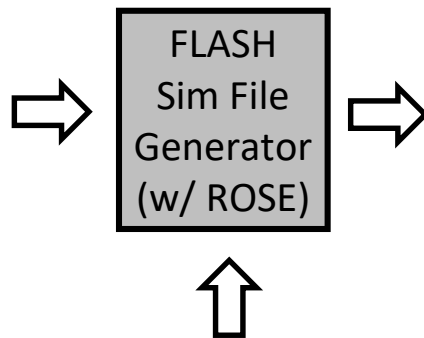
Output:

Scheduling info

HLS Synthesis

*FLASH: Fast, paralleL, Accurate Simulator for HLS

4

- Automated simulation code generation
  - Cycle-accurate simulation
  - Task-level parallelism
  - Pipelined parallelism
  - FIFO simulation & stalls (deadlock)
  - Loop/Func simulation

```
while (i < N){
#pragma HLS pipeline
  if( f1.empty() == false ){
    int temp = f1.read();
    f2.write(temp*711);
    i++;
  }
}
```

<Original HLS C code>

```
* Number of FSM states : 7    * FSM state transitions:
* Pipeline : 1                1 -->
* Pipeline-0 : II = 1, D =              2  / true
* Dataflow Pipeline: 0        2 -->
                                       7  / (!tmp)
                                       3  / (tmp)
<Timing information from       3 -->
  synthesis report>                    4  / true
```

FLASH Sim File Generator (w/ ROSE)

## (Details at poster)

```
static bool p1_en_st3, ...= false;
static int temp_st3, ... temp_st6;
...
if(M2_state == 1){
  ...
  M2_state = 2;
}
else if(M2_state == 2){
  if(p1_en_st6&&f2_wptr==f2_wnum){
    return;
  }
  ...
  if(p1_en_st6 == true){
    p1_en_st6 = false;
    f2_warr[f2_wptr++] = temp_st6;
  }
  ...
  if(p1_en_st3 == true){
    p1_en_st3 = false;
    p1_en_st4 = true;
    temp_st4 = temp_st3;
  }
  ...
  if( i_st2 < N ){
    if( f1_rnum != 0 ){
      p1_en_st3 = true;
      temp_st3=f1_rarr[f1_rptr++];
      i_st2++; ...
} } }
```

Single FSM state simulated per sim func call

Pipeline stall condition

FIFO write

Simulates pipelined parallelism

FIFO empty

FIFO read

<Transformed C code for simulation> 5

- Simulation time comparison

| Benchmark | V C Sim | V RTL Sim | I OCL Sim | FLASH |
|---|---|---|---|---|
| Toy_mpath | 0.602s (1.00X) | 492s (817X) | 4.60s (7.64X) | 0.570s (0.947X) |
| Stencil | 1.46s (1.00X) | 113s (77.4X) | 2.63s (1.80X) | 1.25s (0.856X) |
| MD_sim | 0.0547s (1.00X) | 100s (1,830X) | 0.0921s (1.68X) | 0.0677s (1.24X) |
| Mat_mul | 0.0539s (1.00X) | 192s (3,560X) | 0.201s (3.73X) | 0.0810s (1.50X) |
| AVG | (1.00X) | (1,570X) | (3.71X) | (1.13X) |

Deep (55) pipeline

Frequent FIFO stall
(FIFO depth=1)

The proposed simulator (FLASH):
- runs at a comparable speed with SW simulation (= 1.00X / 1.13X)
- is faster than RTL simulation by 3 orders of magnitude (=1570X/1.13X)
- in some cases, is faster than SW simulation (reason discussed in posters)
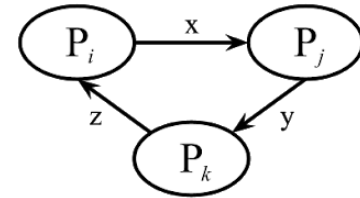- has more overhead with deep pipelines or with frequent FIFO stalls

- Key take-away
  - HLS SW simulation based on the **scheduling information**
    - Can help solve the **correctness** issue and **rapidly** provide **accurate performance estimation**
  - This could substantially **decrease the validation time** of HLS tool customers



Cycle-accurate performance estimation

Correct output data

Detect deadlock situation

  - We hope the presented result could motivate vendors to adopt similar approach in their HLS tools

- Thank you!