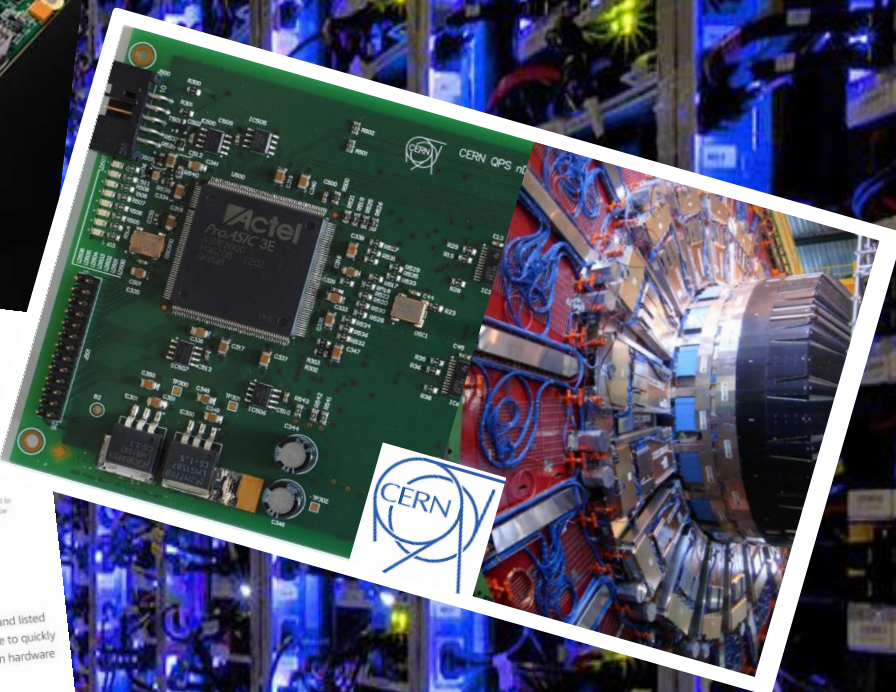# From C/C++ to Dynamically Scheduled Circuits

Lana Josipović

September 2022

**ETH**zürich
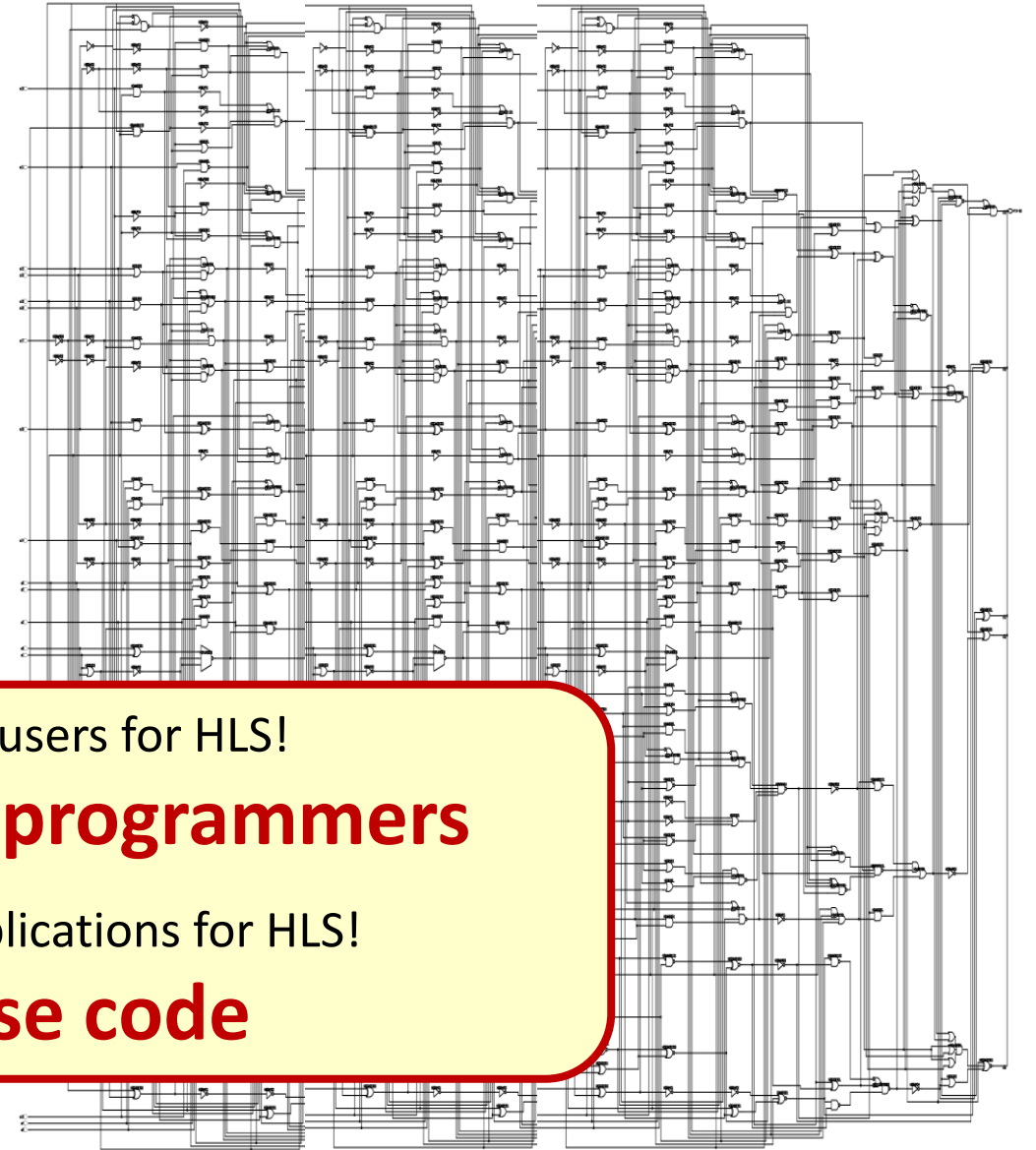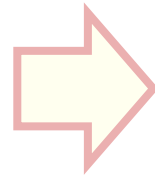
How to perform hardware design?

High parallelism and energy efficiency

# High-Level Synthesis: From Programs to Circuits

```c
#define PI 3.14159265358979932384626434

complex* DFT_naive(complex* x, int N) {
  complex* X = (complex*) malloc(sizeof(struct complex_t) * N);
  int k, n;
  for(k = 0; k < N; k++) {
    X[k].re = 0.0;
    X[k].im = 0.0;
    for(n = 0; n < N; n++) {
      X[k] = add(X[k], multiply(x[n],
                            conv_from_polar(1,
                                -2*PI*n*k/N)));
    }
  }

  return X;
}
```
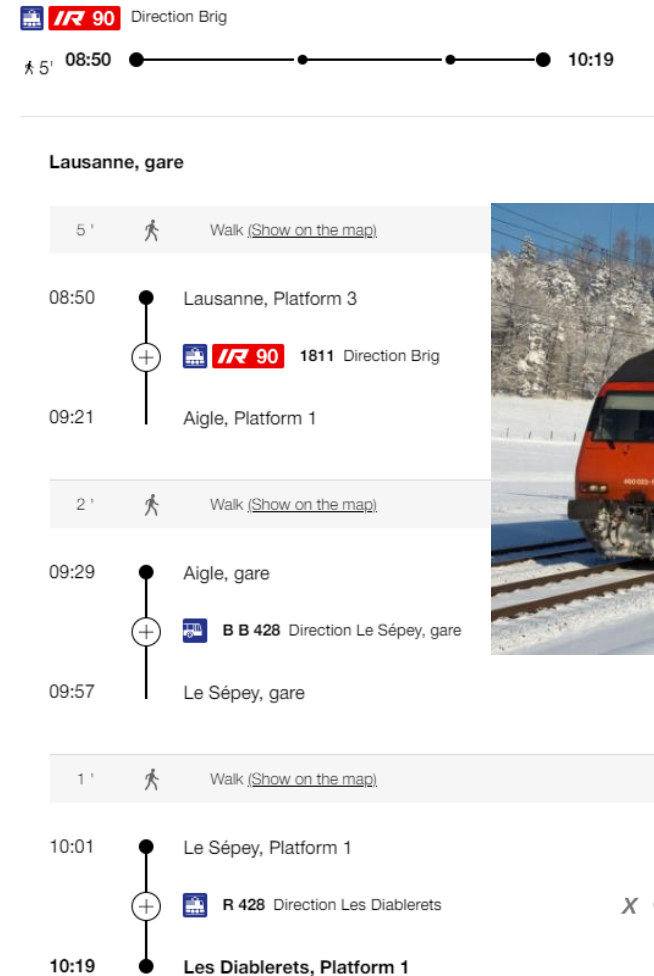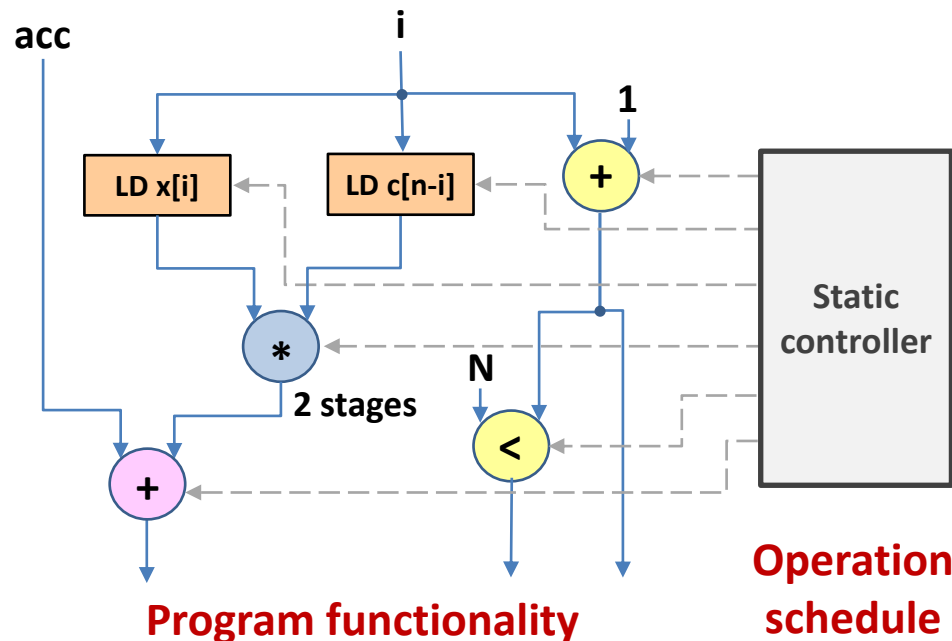
A completely new type of users for HLS!

**Software application programmers**

A completely new type of applications for HLS!

**General-purpose code**

# Standard HLS

- **Create a datapath** suitable to implement the required computation
- Create a **fixed schedule at compile time** to activate the datapath components

```
for (i=0; i<n; i++) {
        acc += x[i] * c[n-i];
}
```



**Program functionality**

**Operation schedule**

# Standard HLS

- **Create a datapath** suitable to implement the required computation
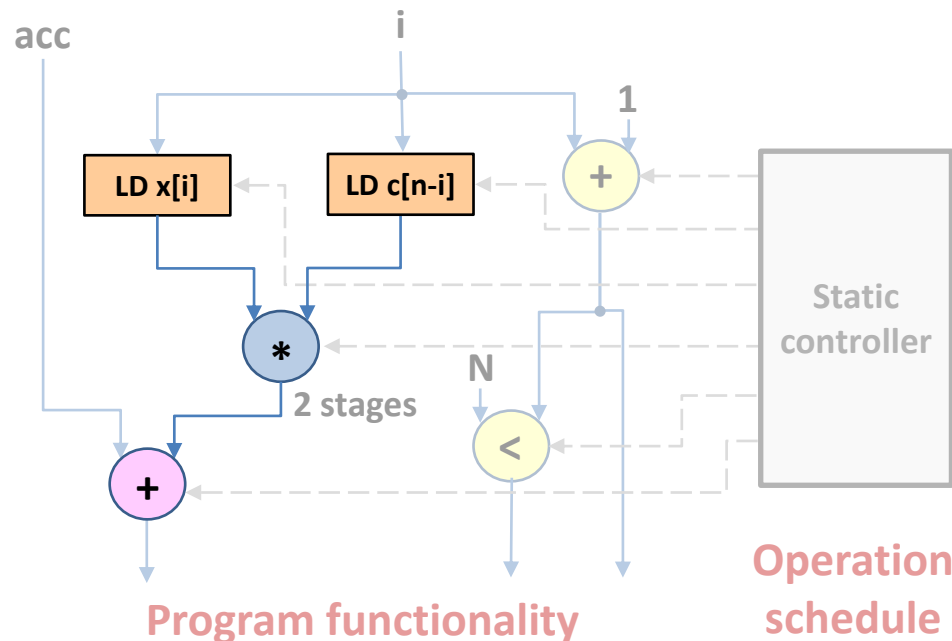- Create a **fixed schedule at compile time** to activate the datapath components

```
for (i=0; i<n; i++) {
        acc += x[i] * c[n-i];
}
```



**Program functionality**  **Operation schedule**
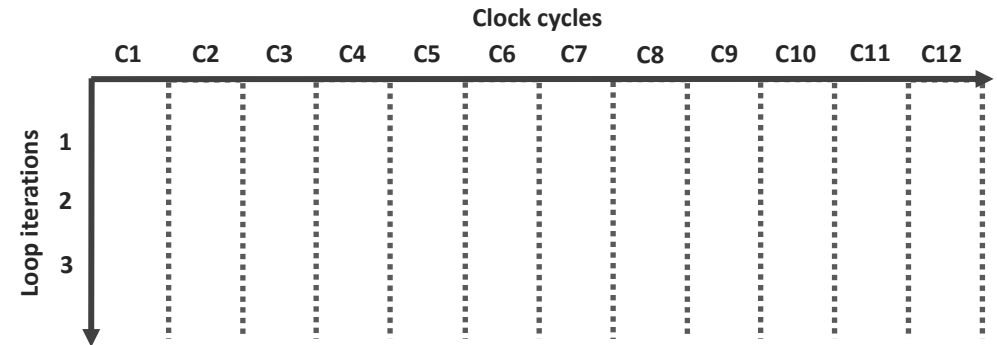
**Naïve schedule:**

# Standard HLS

- **Create a datapath** suitable to implement the required computation
- Create a **fixed schedule at compile time** to activate the datapath components

```
for (i=0; i<n; i++) {
        acc += x[i] * c[n-i];
}
```
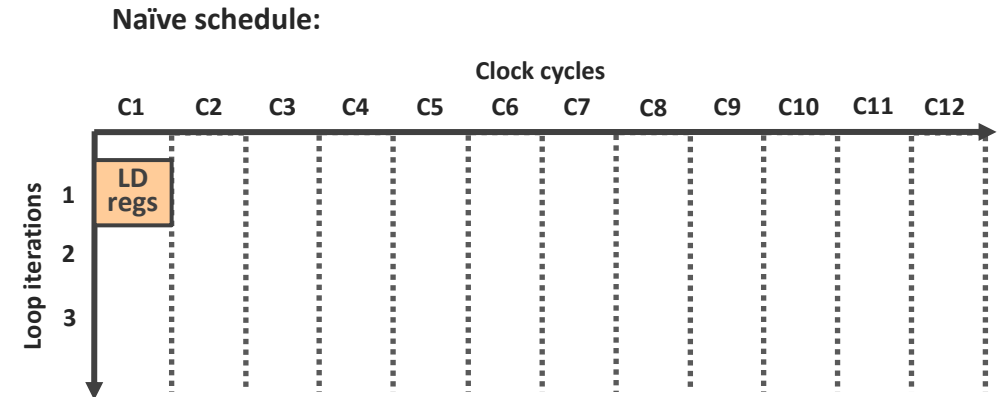


**Program functionality**

**Operation schedule**

**Naïve schedule:**

# Standard HLS

- **Create a datapath** suitable to implement the required computation
- Create a **fixed schedule at compile time** to activate the datapath components
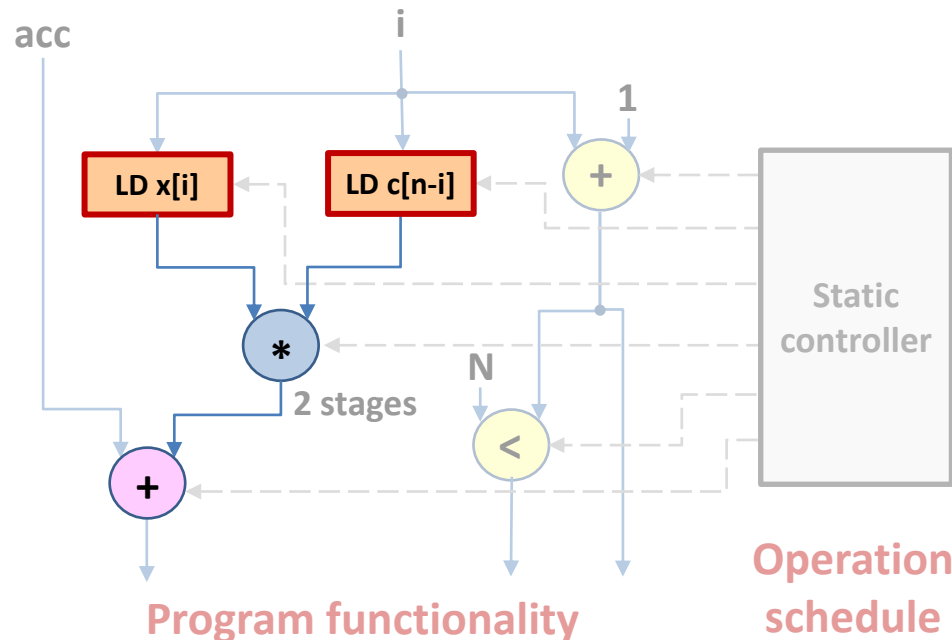
```
for (i=0; i<n; i++) {
        acc += x[i] * c[n-i];
}
```



**Program functionality**

**Operation schedule**

# Standard HLS

- **Create a datapath** suitable to implement the required computation
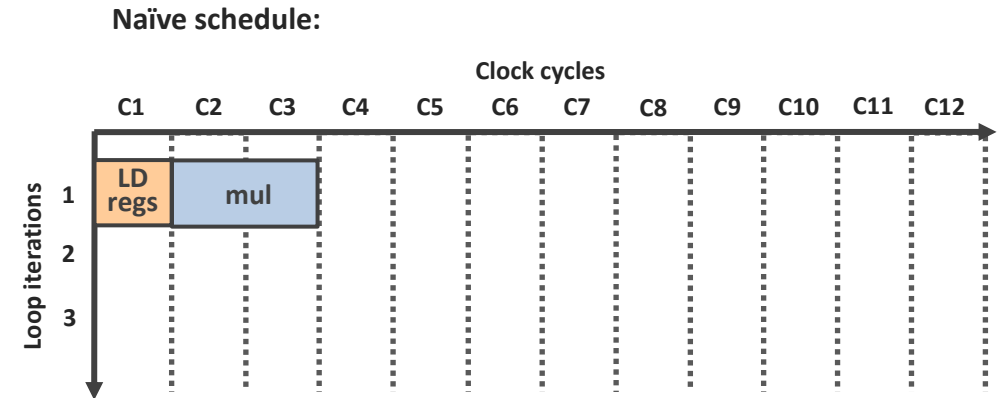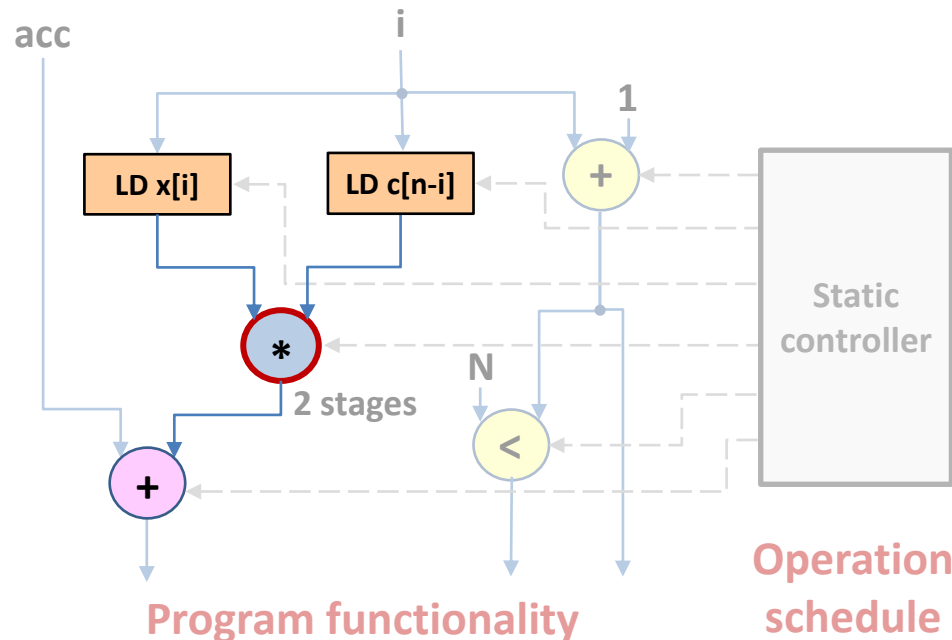- Create a **fixed schedule at compile time** to activate the datapath components

```
for (i=0; i<n; i++) {
        acc += x[i] * c[n-i];
}
```



**Program functionality**

**Operation schedule**

# Standard HLS

- **Create a datapath** suitable to implement the required computation
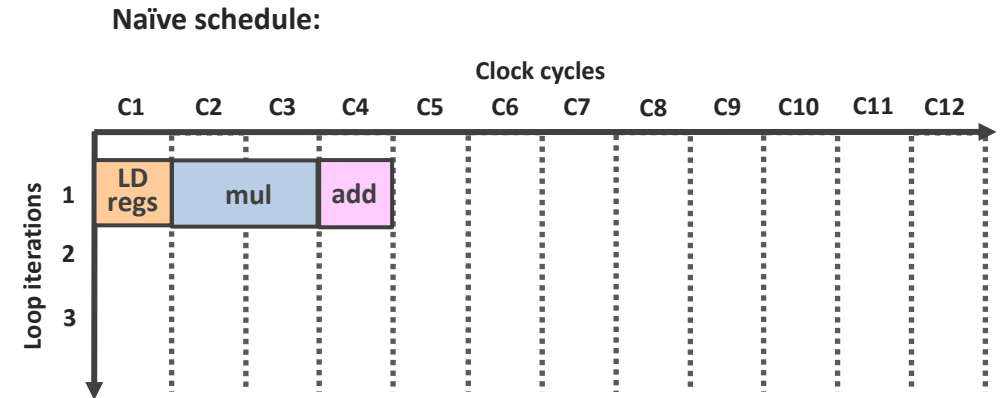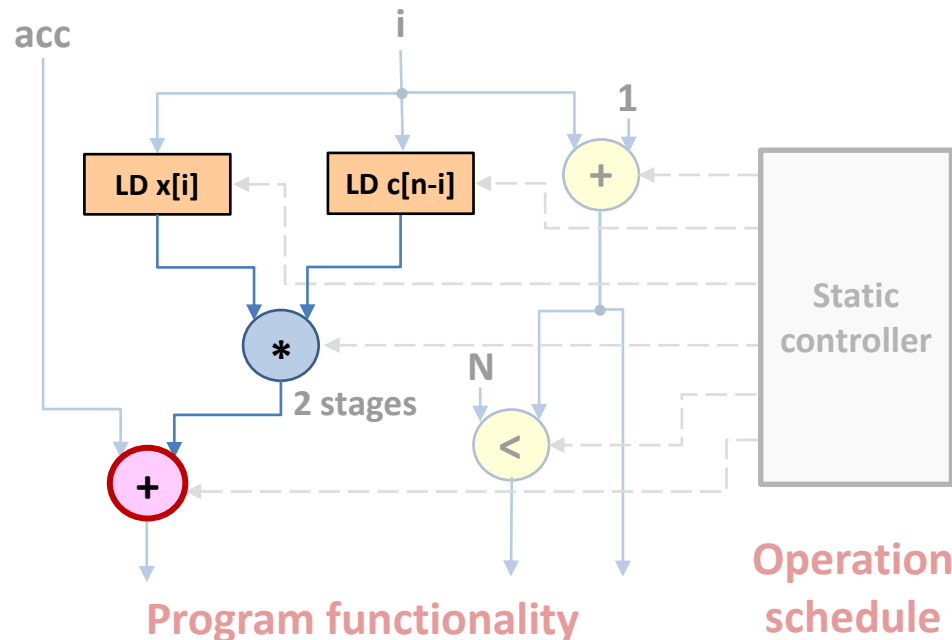- Create a **fixed schedule at compile time** to activate the datapath components

```
for (i=0; i<n; i++) {
        acc += x[i] * c[n-i];
}
```



acc  i

LD x[i]   LD c[n-i]   1   +

*

2 stages   N

+   <

Static controller

**Program functionality**   **Operation schedule**



Naïve schedule:

Clock cycles

C1  C2  C3  C4  C5  C6  C7  C8  C9  C10  C11  C12

Loop iterations

1  | LD regs | mul | add |
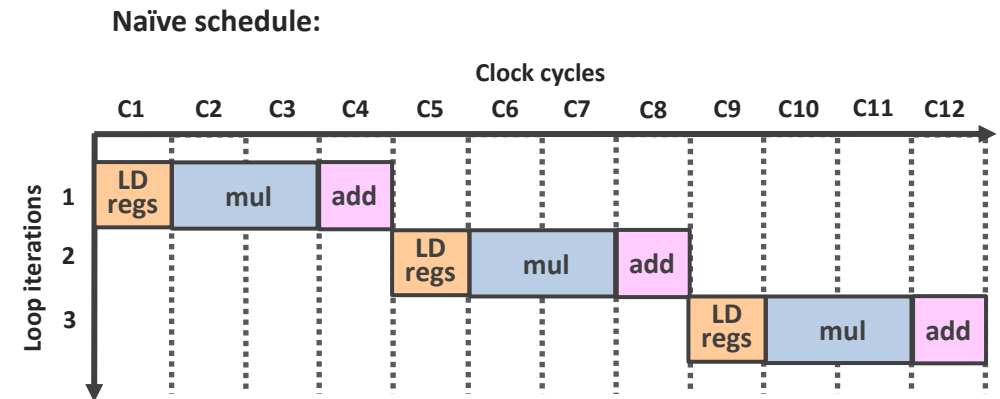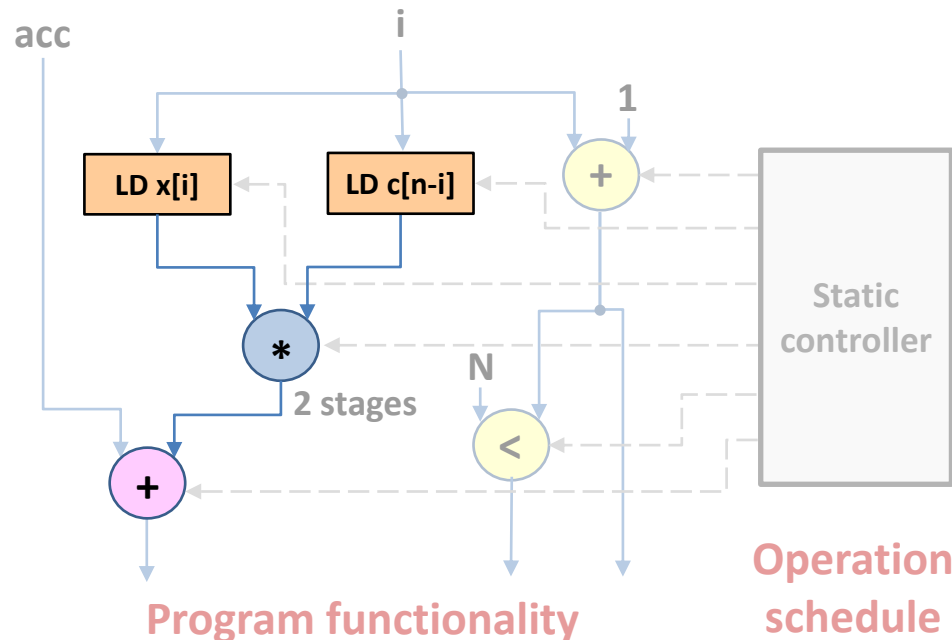2  | LD regs | mul | add |
3  | LD regs | mul | add |

# Standard HLS

- **Create a datapath** suitable to implement the required computation
- Create a **fixed schedule at compile time** to activate the datapath components

```
for (i=0; i<n; i++) {
        acc += x[i] * c[n-i];
}
```
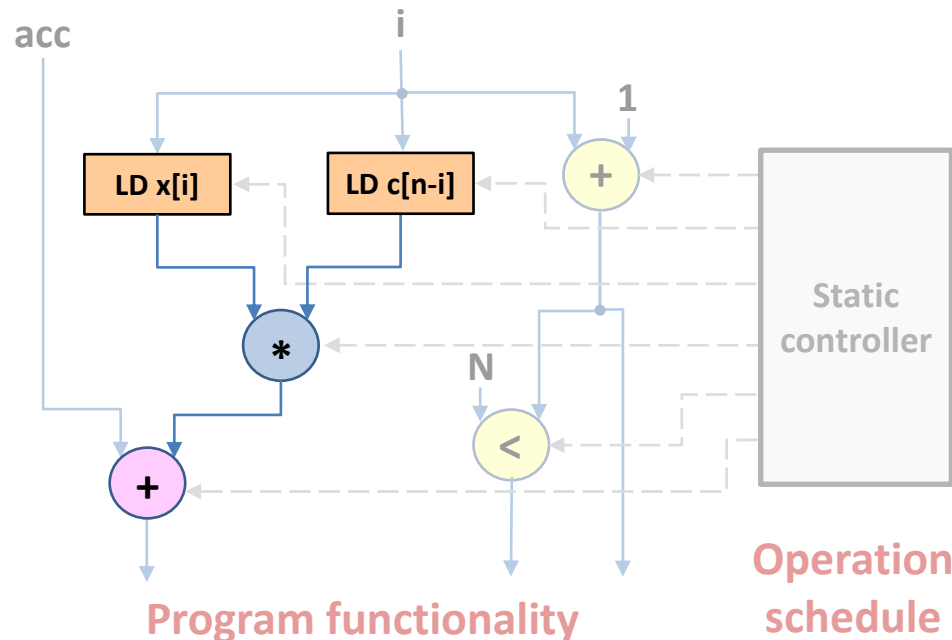


acc    i

LD x[i]    LD c[n-i]    1    +

*

2 stages    N    Static controller

+    <

**Program functionality**    **Operation schedule**

Naïve schedule:

Clock cycles



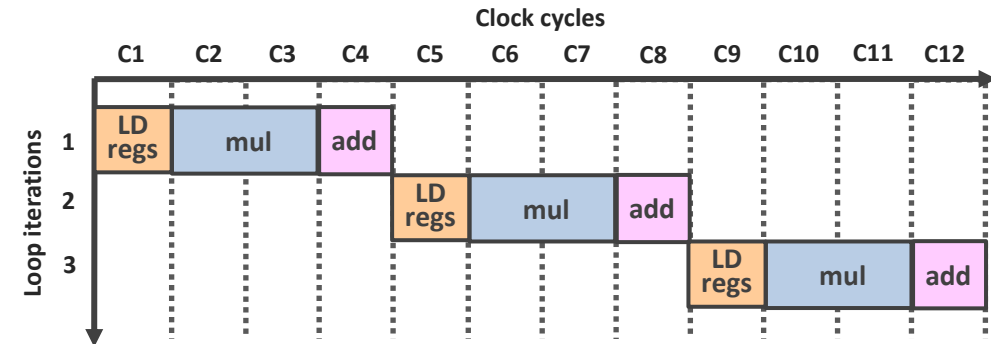**Low throughput: slow execution**

# Standard HLS

- **Create a datapath** suitable to implement the required computation
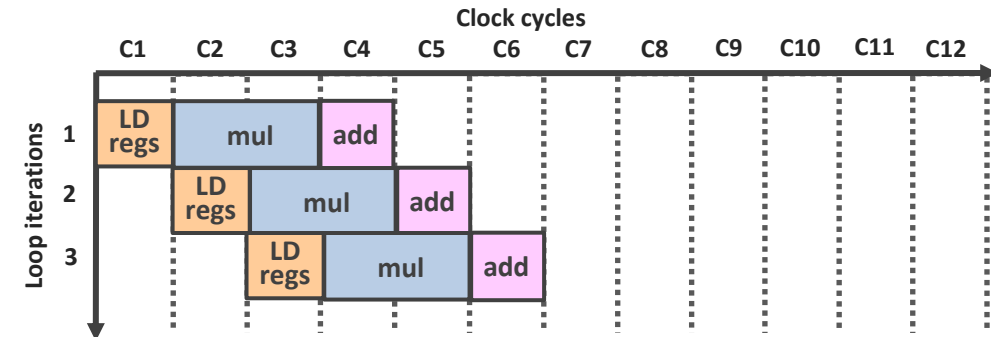- Create a **fixed schedule at compile time** to activate the datapath components

```
for (i=0; i<n; i++) {
        acc += x[i] * c[n-i];
}
```



Program functionality

Operation schedule



Naïve schedule:

Pipelined schedule:

**High throughput: fast execution**

# The Limitations of Static Scheduling

```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```
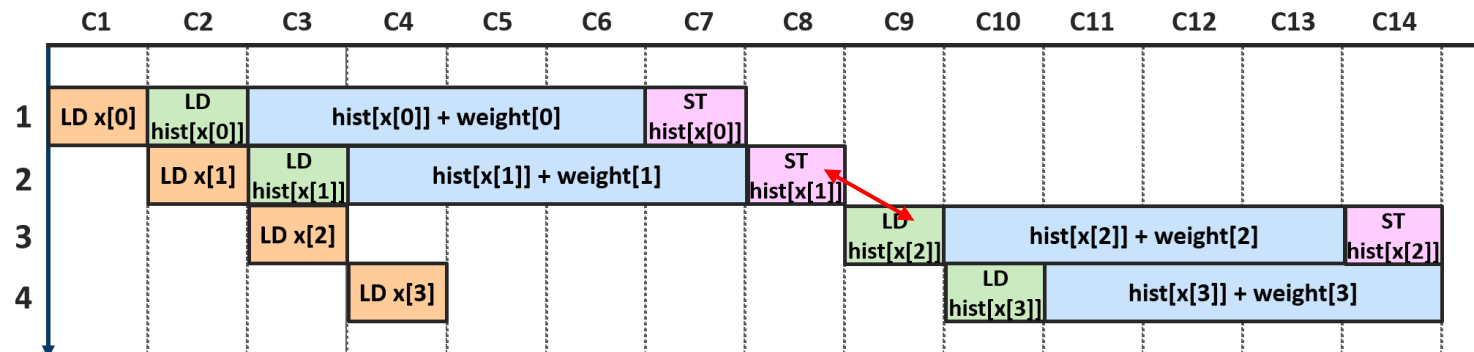
```
1: x[0]=5 → ld hist[5]; st hist[5];
2: x[1]=4 → ld hist[4]; st hist[4];
3: x[2]=4 → ld hist[4]; st hist[4];
```

**RAW dependency**

- ## Static scheduling (standard HLS tool)
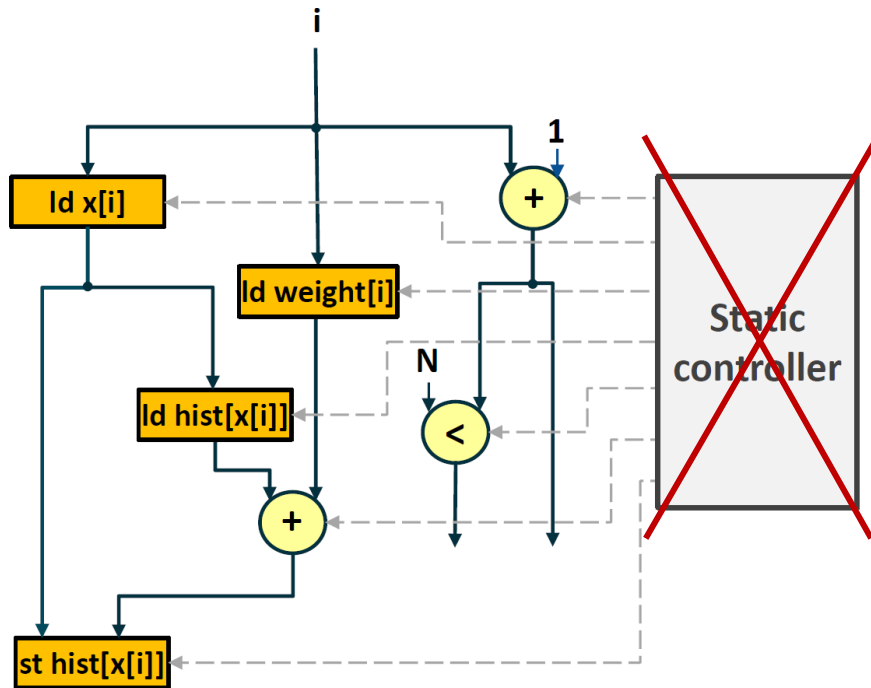  - Inferior when memory accesses cannot be disambiguated at compile time



- ## Dynamic scheduling
  - Maximum parallelism: Only serialize memory accesses on actual dependencies

# A Different Way to Do HLS

**Static scheduling** (standard HLS tool): decide at **compile time** when each operation executes

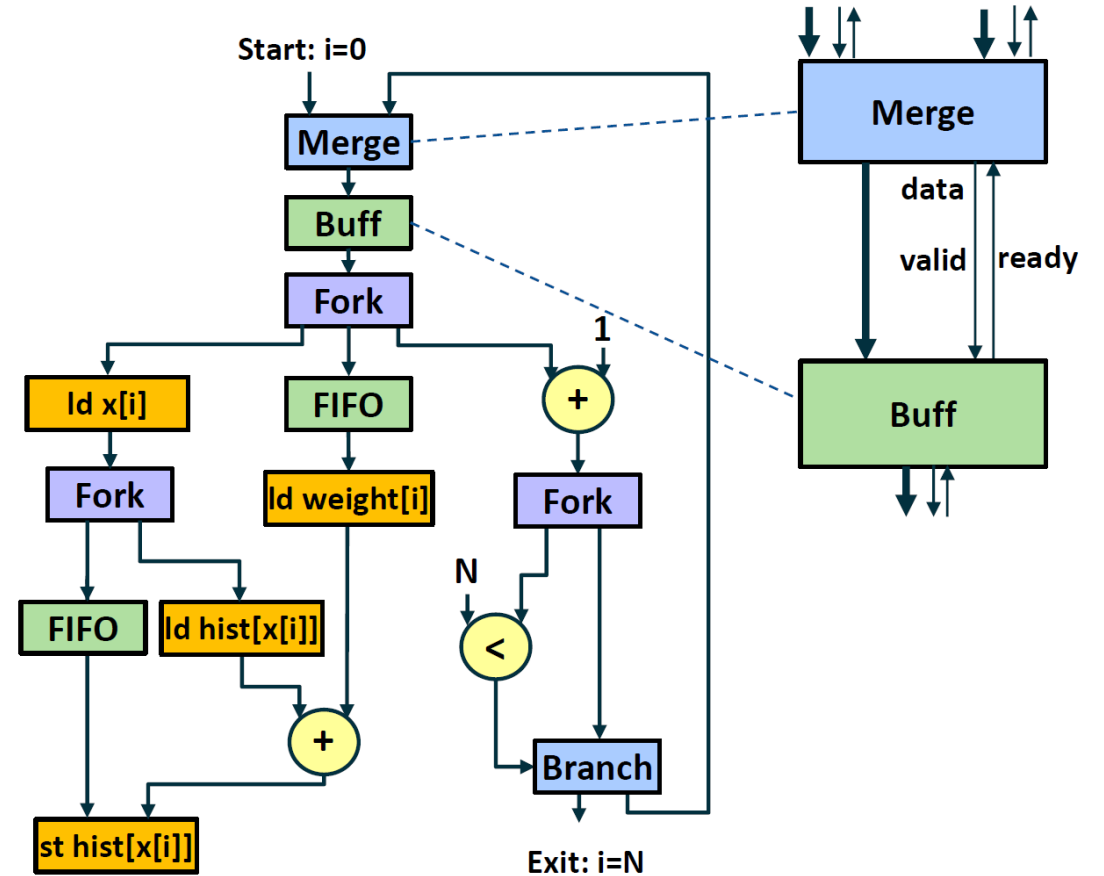**Dynamic scheduling** (our HLS approach): decide at **runtime** when each operation executes

# A Different Way to Do HLS

**Static scheduling** (standard HLS tool): decide at **compile time** when each operation executes

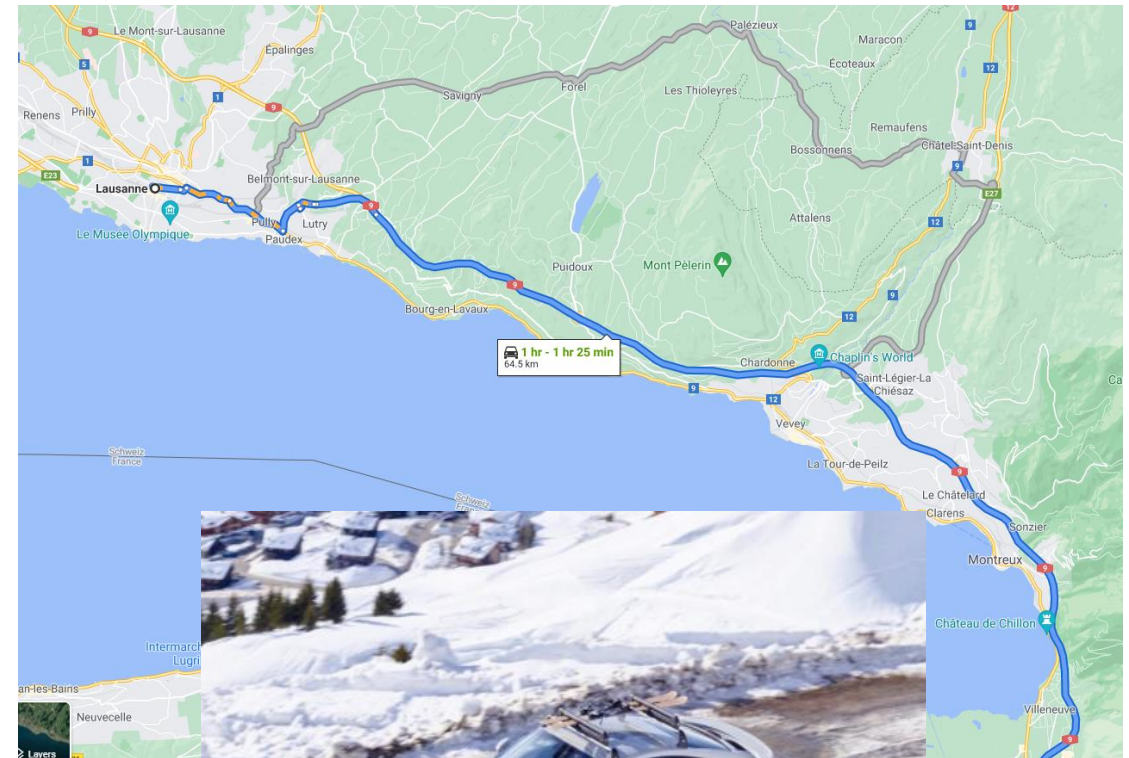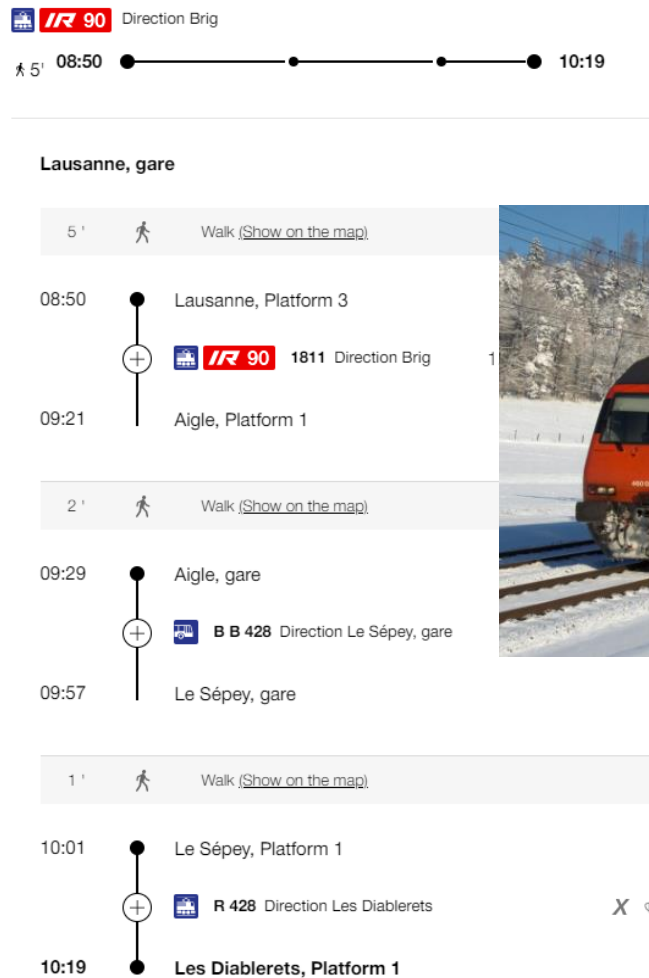**Dynamic scheduling** (our HLS approach): decide at **runtime** when each operation executes

# A Different Way to Do HLS

**Static scheduling** (standard HLS tool): decide at **compile time** when each operation executes

**Dynamic scheduling** (our HLS approach): decide at **runtime** when each operation executes
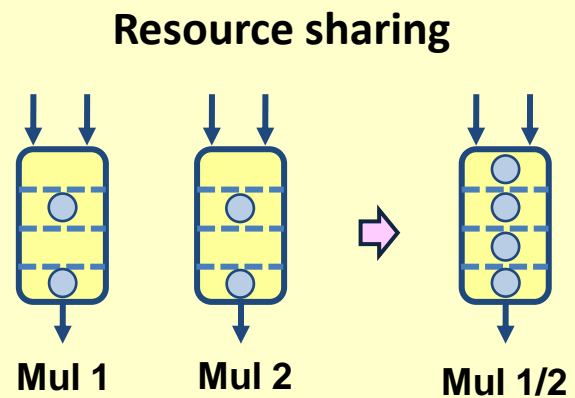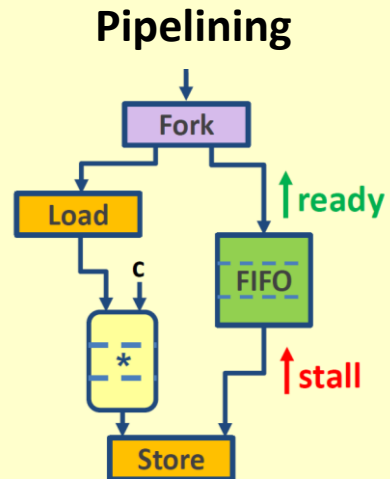
# Dataflow Circuits

- **Asynchronous circuits**: operators triggered when inputs are available
  - Budiu et al. Dataflow: A complement to superscalar. ISPASS'05.

- Dataflow, latency-insensitive, elastic: the **synchronous** version of it
  - Cortadella et al. Synthesis of synchronous elastic architectures. DAC'06.
  - Carloni et al. Theory of latency-insensitive design. TCAD'01.
  - Jacobson et al. Synchronous interlocked pipelines. ASYNC'02.
  - Vijayaraghavan and Arvind. Bounded dataflow networks and latency-insensitive circuits. MEMOCODE'09.

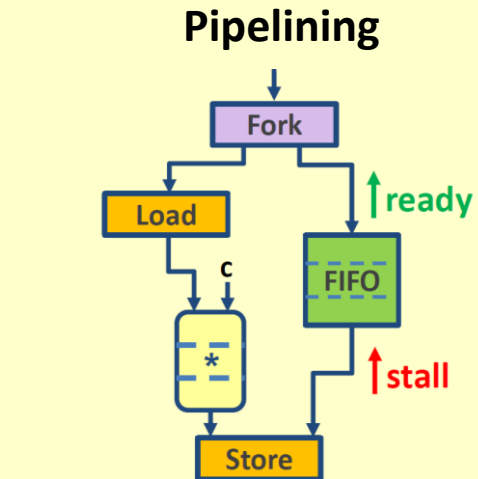**High-level synthesis of
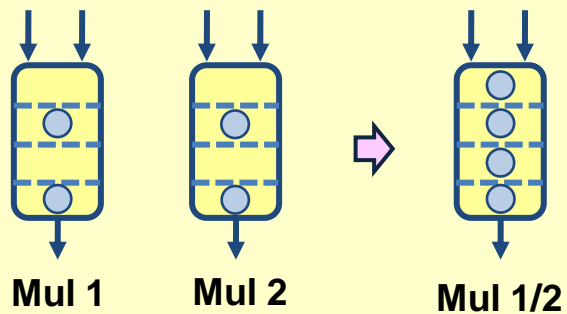dataflow circuits**

# HLS of Dynamically Scheduled Circuits

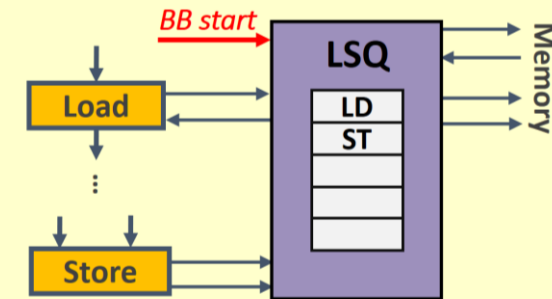# HLS of Dynamically Scheduled Circuits

## Catching up with static HLS

**Pipelining**



**Resource sharing**



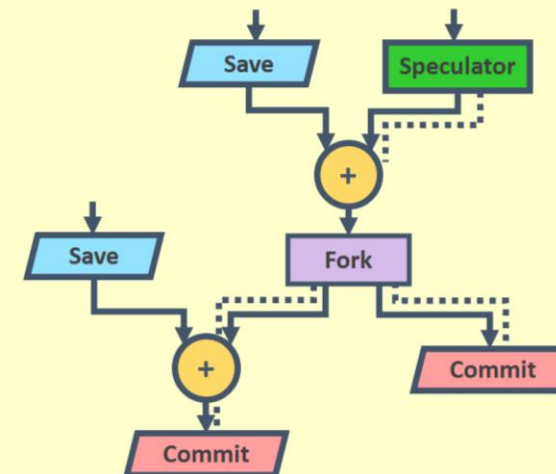Mul 1     Mul 2     Mul 1/2

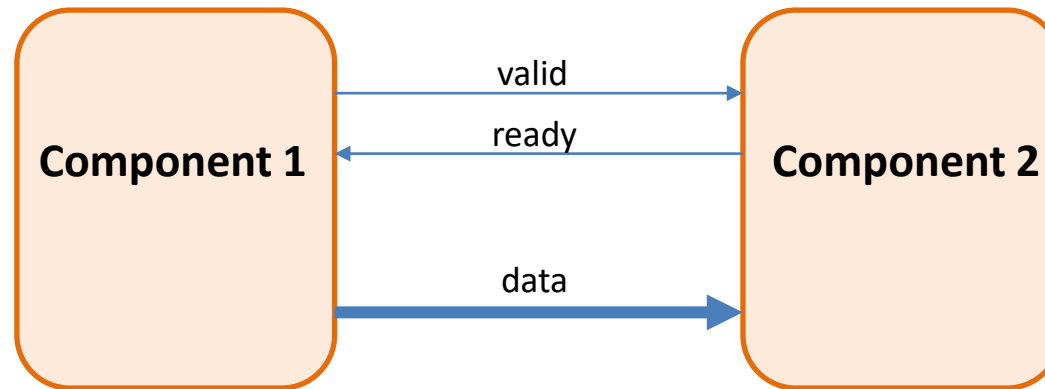## Reaping the benefits of dynamic scheduling

**Out-of-order memory**



**Speculative execution**

# Dataflow Circuits

- We use the **SELF (Synchronous ELastic Flow)** protocol
  - Cortadella et al. Synthesis of synchronous elastic architectures. DAC'06.
- Every component communicates via a pair of handshake signals
- **Make scheduling decisions at runtime**
  - As soon as all conditions for execution are satisfied, an operation starts
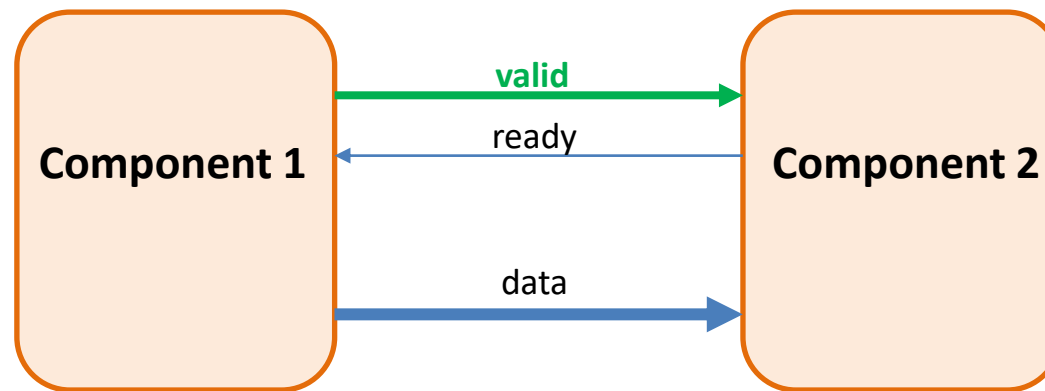
# Dataflow Circuits

- We use the **SELF (Synchronous ELastic Flow)** protocol
  - Cortadella et al. Synthesis of synchronous elastic architectures. DAC'06.
- Every component communicates via a pair of handshake signals
- **Make scheduling decisions at runtime**
  - As soon as all conditions for execution are satisfied, an operation starts
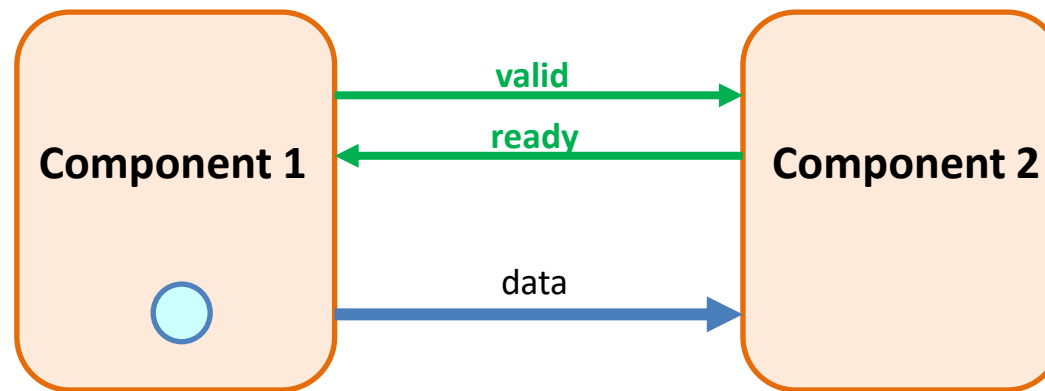
# Dataflow Circuits

- We use the **SELF (Synchronous ELastic Flow)** protocol
  - Cortadella et al. Synthesis of synchronous elastic architectures. DAC'06.
- Every component communicates via a pair of handshake signals
- **Make scheduling decisions at runtime**
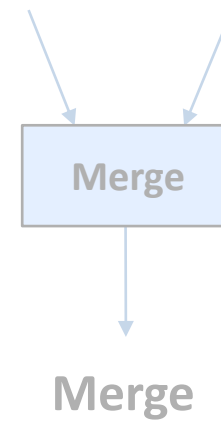  - As soon as all conditions for execution are satisfied, an operation starts
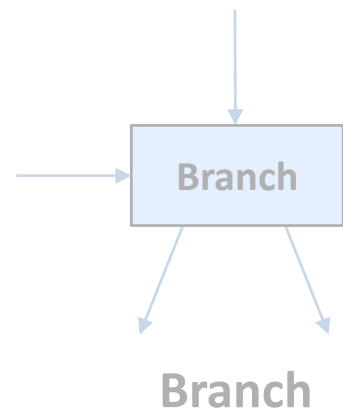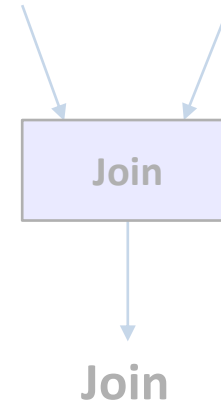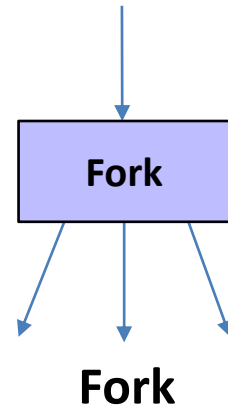
# Dataflow Components



Fork

Join

Branch

Merge

# Dataflow Components



Fork

Fork

Join

Join

Branch

Branch

Merge

Merge

data valid ready

# Dataflow Components

**Fork**

Join

Branch

Merge

# Dataflow Components

# Dataflow Components



Fork

Join

Branch

Merge

# Dataflow Components

Fork

Join

Branch

Merge

# From Program to Dataflow Circuit



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

# From Program to Dataflow Circuit



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

Josipović, Ghosal, and Ienne. Dynamically Scheduled High-Level Synthesis. FPGA 2018 **Best Paper Award Nominee**
Josipović, Brisk, and Ienne. From C to Elastic Circuits. Asilomar 2017

# From Program to Dataflow Circuit



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

# From Program to Dataflow Circuit



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

Josipović, Ghosal, and Ienne. Dynamically Scheduled High-Level Synthesis. FPGA 2018 **Best Paper Award Nominee**
Josipović, Brisk, and Ienne. From C to Elastic Circuits. Asilomar 2017

# From Program to Dataflow Circuit



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

Josipović, Ghosal, and Ienne. Dynamically Scheduled High-Level Synthesis. FPGA 2018 **Best Paper Award Nominee**
Josipović, Brisk, and Ienne. From C to Elastic Circuits. Asilomar 2017

# From Program to Dataflow Circuit



Single token on cycle, in-order tokens in noncyclic paths

Josipović, Ghosal, and Ienne. Dynamically Scheduled High-Level Synthesis. FPGA 2018 **Best Paper Award Nominee**
Josipović, Brisk, and Ienne. From C to Elastic Circuits. Asilomar 2017

# From Program to Dataflow Circuit



Backpressure from slow paths prevents pipelining

# HLS of Dynamically Scheduled Circuits

## Catching up with static HLS

**Pipelining**



## Resource sharing



Mul 1    Mul 2    Mul 1/2

## Reaping the benefits of dynamic scheduling

### Out-of-order memory



### Speculative execution

# Inserting Buffers

```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```



Buffers as registers to break combinational paths

Josipović, Sheikhha, Guerrieri, Ienne, and Cortadella. Buffer Placement and Sizing for High-Performance Dataflow Circuits. FPGA 2020 **Best paper award**
Rizzi, Guerrieri, Ienne, and Josipović. A Comprehensive Timing Model for Accurate Frequency Tuning in Dataflow Circuits. FPL 2022

# Inserting Buffers

```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```



Buffers as FIFOs to regulate throughput

Josipović, Sheikhha, Guerrieri, Ienne, and Cortadella. Buffer Placement and Sizing for High-Performance Dataflow Circuits. FPGA 2020 **Best paper award**
Rizzi, Guerrieri, Ienne, and Josipović. A Comprehensive Timing Model for Accurate Frequency Tuning in Dataflow Circuits. FPL 2022

# Inserting Buffers

```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

Josipović, Sheikhha, Guerrieri, Ienne, and Cortadella. Buffer Placement and Sizing for High-Performance Dataflow Circuits. FPGA 2020 **Best paper award**
Rizzi, Guerrieri, Ienne, and Josipović. A Comprehensive Timing Model for Accurate Frequency Tuning in Dataflow Circuits. FPL 2022

# Inserting Buffers

Josipović, Sheikhha, Guerrieri, Ienne, and Cortadella. Buffer Placement and Sizing for High-Performance Dataflow Circuits. FPGA 2020 **Best paper award**
Rizzi, Guerrieri, Ienne, and Josipović. A Comprehensive Timing Model for Accurate Frequency Tuning in Dataflow Circuits. FPL 2022

# Inserting Buffers

**NOW
(with buffers)**



Mixed integer linear programming (MILP) model based on **Petri net theory**
- Analyze token flow through the circuit
- Determine **buffer placement and sizing**
- **Maximize throughput** for a target clock period

Josipović, Sheikhha, Guerrieri, Ienne, and Cortadella. Buffer Placement and Sizing for High-Performance Dataflow Circuits. FPGA 2020 **Best paper award**
Rizzi, Guerrieri, Ienne, and Josipović. A Comprehensive Timing Model for Accurate Frequency Tuning in Dataflow Circuits. FPL 2022

# HLS of Dynamically Scheduled Circuits



**Catching up with static HLS**

Pipelining

Resource sharing

Mul 1    Mul 2    Mul 1/2

**Reaping the benefits of dynamic scheduling**

Out-of-order memory

Speculative execution

# Saving Resources through Sharing

- Static HLS: share units between operations which execute in **different clock cycles**

- Dynamic HLS: share units based on their **average utilization** with tokens

```
for (i = 0; i < N; i++) {
    a[i] = a[i]*x;
    b[i] = b[i]*y;
}
```



M1          M2

Josipović, Marmet, Guerrieri, and Ienne. Resource Sharing in Dataflow Circuits. FCCM 2022. **Best Paper Award Nominee**

# Saving Resources through Sharing

- Static HLS: share units between operations which execute in **different clock cycles**
- Dynamic HLS: share units based on their **average utilization** with tokens

```
for (i = 0; i < N; i++) {
    a[i] = a[i]*x;
    b[i] = b[i]*y;
}
```

**M1**      **M2**

**Units fully utilized
(high throughput)**

**Sharing not possible without
damaging throughput**

**Use MILP (performance optimization)
information to decide what to share**

Josipović, Marmet, Guerrieri, and Ienne. Resource Sharing in Dataflow Circuits. FCCM 2022. **Best Paper Award Nominee**

# Saving Resources through Sharing

- Static HLS: share units between operations which execute in **different clock cycles**
- Dynamic HLS: share units based on their **average utilization** with tokens

```
for (i = 0; i < N; i++) {
    a[i] = a[i]*x;
    b[i] = b[i]*y;
}
```



**M1**      **M2**           **M1/2**

**Sharing possible without damaging throughput**

**Units underutilized (low throughput)**

**Use MILP (performance optimization) information to decide what to share**

Josipović, Marmet, Guerrieri, and Ienne. Resource Sharing in Dataflow Circuits. FCCM 2022. **Best Paper Award Nominee**

# Saving Resources through Sharing

- Static HLS: share units between operations which execute in **different clock cycles**
- Dynamic HLS: share units based on their **average utilization** with tokens

```
for (i = 0; i < N; i++) {
    a[i] = a[i]*x;
    b[i] = b[i]*y;
}
```



**M1**      **M2**        **M1/2**

**Units underutilized (low throughput)**

*Inputs of M1, M2*

Token order: M1, M2

M1/2    *    FIFO

Branch

**Sharing mechanism for deadlock-free execution**

Josipović, Marmet, Guerrieri, and Ienne. Resource Sharing in Dataflow Circuits. FCCM 2022. **Best Paper Award Nominee**

# Inserting Buffers

```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```



**Backpressure from slow paths prevents pipelining**

# Inserting Buffers

```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```



**Buffers for high throughput**

# Inserting Buffers

```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

1: x[0]=5 → ld hist[5]; st hist[5];
2: x[1]=4 → ld hist[4]; st hist[4];
3: x[2]=4 → ld hist[4]; st hist[4];

**RAW dependency**



**RAW dependency
not honored!**

**What about memory?**

# HLS of Dynamically Scheduled Circuits

## Catching up with static HLS

**Pipelining**

**Resource sharing**

Mul 1    Mul 2    Mul 1/2

## Reaping the benefits of dynamic scheduling

**Out-of-order memory**

*BB start*

**Speculative execution**

# We Need a Load-Store Queue (LSQ)!

- Traditional processor LSQs allocate memory instructions **in program order**

```
loop: lw  $t2, 0($t4)
      lw  $t3, 100($t4)
      mul $t5, $t2, $t3
      addi $t5, $t5, $t1
      sw  $t5, 100($t4)
      addi $t1, $t1, 4
      bne $t6, $t1, loop
```

Instruction fetch & decode (in order) → Processor datapath (out of order) → Ordering (load-store queue) → Memory

- Dataflow circuits have **no notion of program order**

**How to supply program order to the LSQ?**

Dataflow (out of order)

load x[i]
load x[0]
load y[i]
store x[i]

**???**

Ordering (load-store queue) → Memory

# LSQ Allocation

- An LSQ for dataflow circuits whose only difference is in the **allocation policy**:
  - **Static knowledge** of memory access program order inside each basic block
  - **Dynamic knowledge** of the sequence of basic blocks **from the dataflow circuit**



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

Josipović, Brisk, and Ienne. An Out-of-Order Load-Store Queue for Spatial Computing. CASES 2017 **Best Paper Award Nominee**
Josipović, Bhattacharrya, Guerrieri, and Ienne. Shrink It or Shed It! Minimize the Use of LSQs in Dataflow Designs. FPT 2019

# LSQ Allocation

- An LSQ for dataflow circuits whose only difference is in the **allocation policy**:
  - **Static knowledge** of memory access program order inside each basic block
  - **Dynamic knowledge** of the sequence of basic blocks **from the dataflow circuit**



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

Josipović, Brisk, and Ienne. An Out-of-Order Load-Store Queue for Spatial Computing. CASES 2017 **Best Paper Award Nominee**
Josipović, Bhattacharrya, Guerrieri, and Ienne. Shrink It or Shed It! Minimize the Use of LSQs in Dataflow Designs. FPT 2019

# LSQ Allocation

- An LSQ for dataflow circuits whose only difference is in the **allocation policy**:
  - **Static knowledge** of memory access program order inside each basic block
  - **Dynamic knowledge** of the sequence of basic blocks **from the dataflow circuit**



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

Josipović, Brisk, and Ienne. An Out-of-Order Load-Store Queue for Spatial Computing. CASES 2017 **Best Paper Award Nominee**
Josipović, Bhattacharrya, Guerrieri, and Ienne. Shrink It or Shed It! Minimize the Use of LSQs in Dataflow Designs. FPT 2019
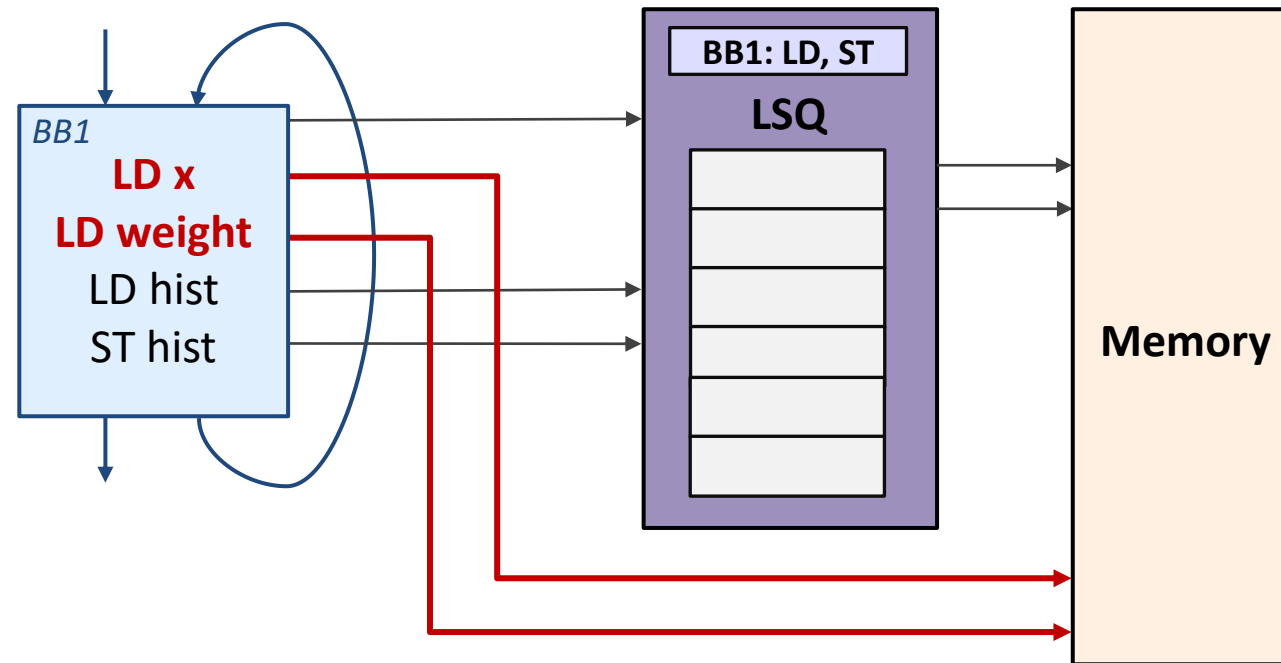
# LSQ Allocation

- An LSQ for dataflow circuits whose only difference is in the **allocation policy**:
  - **Static knowledge** of memory access program order inside each basic block
  - **Dynamic knowledge** of the sequence of basic blocks **from the dataflow circuit**



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

Josipović, Brisk, and Ienne. An Out-of-Order Load-Store Queue for Spatial Computing. CASES 2017 **Best Paper Award Nominee**
Josipović, Bhattacharrya, Guerrieri, and Ienne. Shrink It or Shed It! Minimize the Use of LSQs in Dataflow Designs. FPT 2019
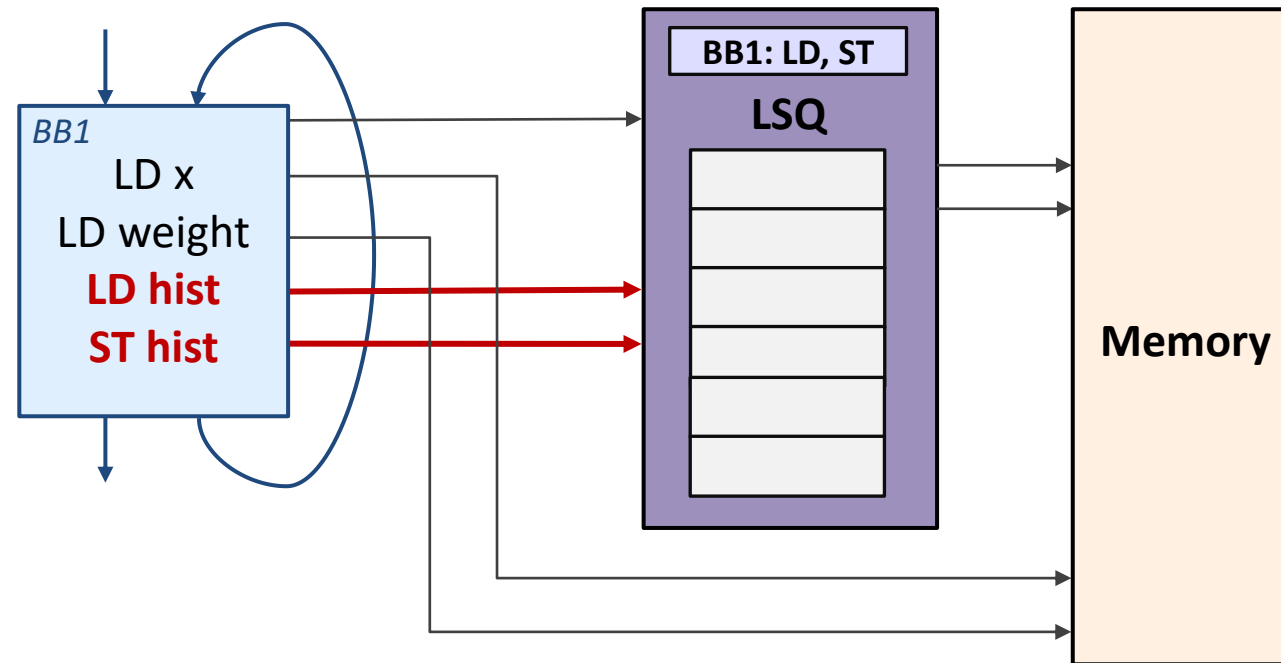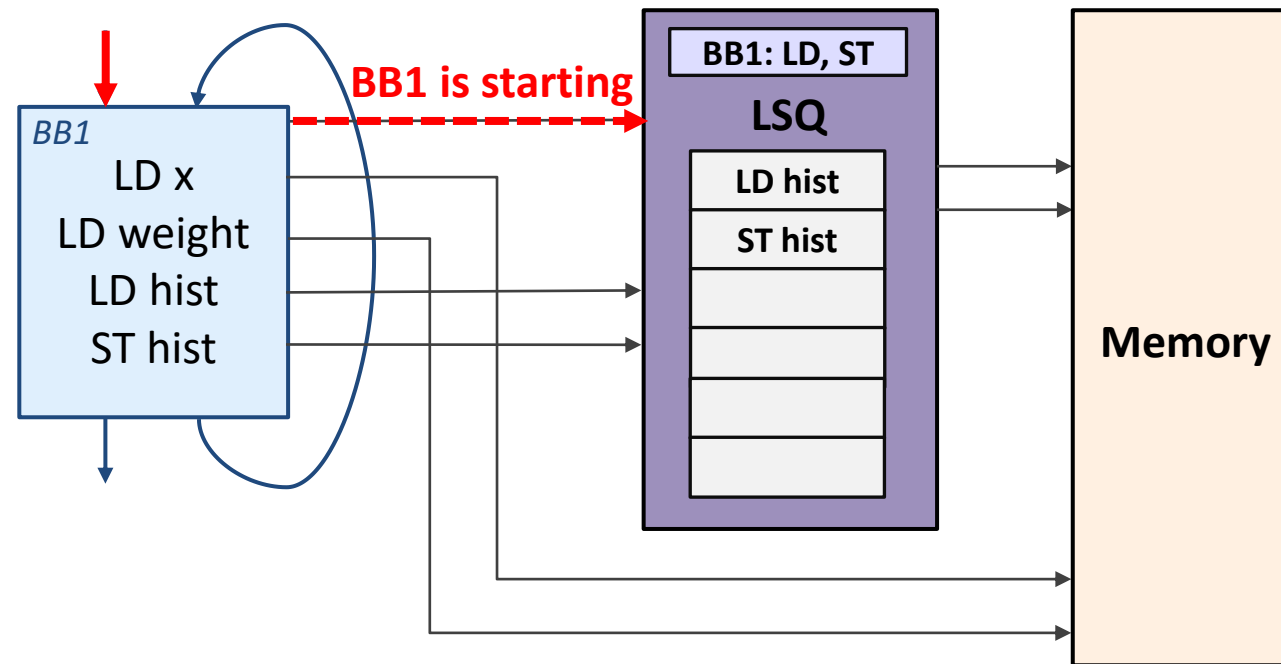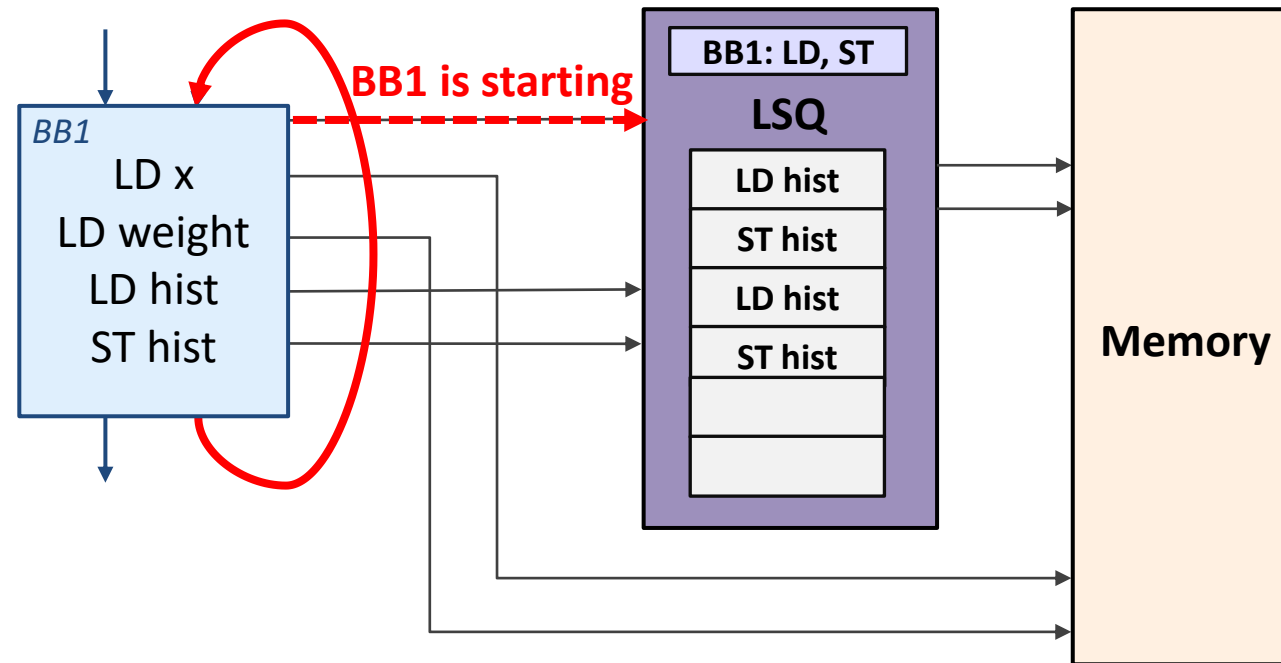
# LSQ Allocation

- An LSQ for dataflow circuits whose only difference is in the **allocation policy**:
  - **Static knowledge** of memory access program order inside each basic block
  - **Dynamic knowledge** of the sequence of basic blocks **from the dataflow circuit**



```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

Josipović, Brisk, and Ienne. An Out-of-Order Load-Store Queue for Spatial Computing. CASES 2017 **Best Paper Award Nominee**
Josipović, Bhattacharrya, Guerrieri, and Ienne. Shrink It or Shed It! Minimize the Use of LSQs in Dataflow Designs. FPT 2019
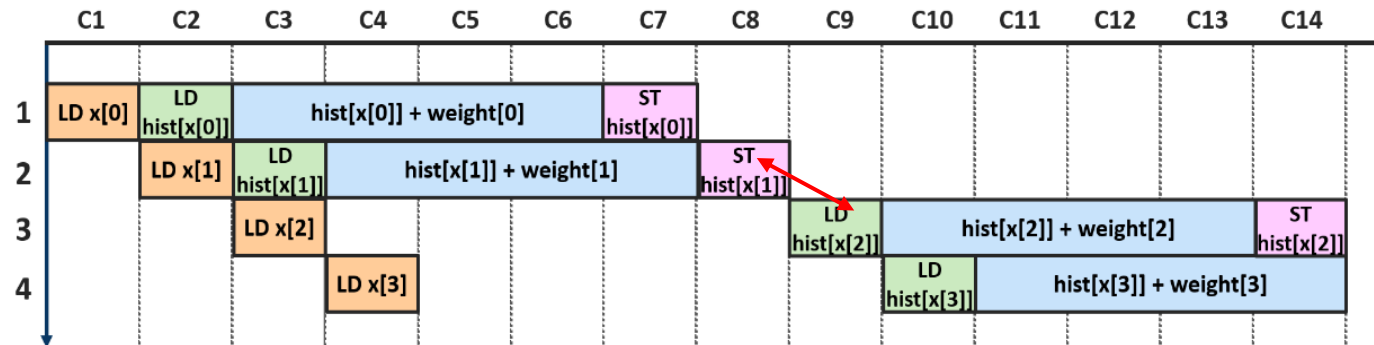
# Dataflow Circuit with the LSQ

```
for (i=0; i<N; i++) {
    hist[x[i]] = hist[x[i]] + weight[i];
}
```

```
1: x[0]=5 → ld hist[5]; st hist[5];
2: x[1]=4 → ld hist[4]; st hist[4];
3: x[2]=4 → ld hist[4]; st hist[4];
```

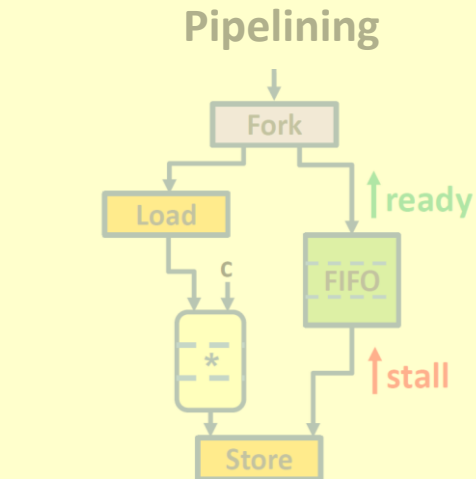**RAW dependency**



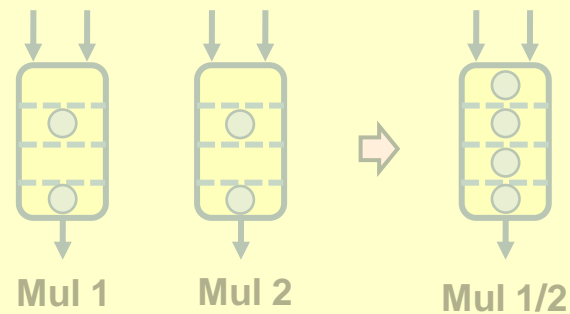**High-throughput pipeline with memory dependencies honored**

# HLS of Dynamically Scheduled Circuits



**Catching up with static HLS**

Pipelining
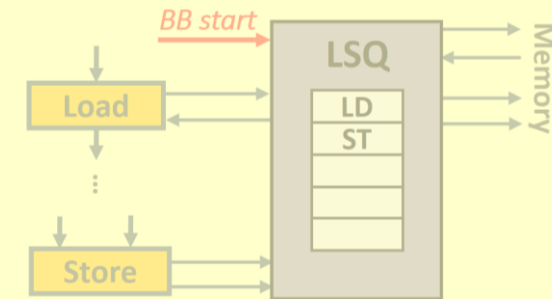
Resource sharing

**Reaping the benefits of dynamic scheduling**

Out-of-order memory

Speculative execution

# Nonspeculative Dataflow Circuit



```
float d=0.0; x=100.0; int i=0;

do {
    d = a[i] + b[i];
    i++;
}
while (d<x);
```

# Nonspeculative Dataflow Circuit



```
float d=0.0; x=100.0; int i=0;

do {
    d = a[i] + b[i];
    i++;
}
while (d<x);
```

# Nonspeculative Dataflow Circuit



```
float d=0.0; x=100.0; int i=0;

do {
    d = a[i] + b[i];
    i++;
}
while (d<x);
```

# Nonspeculative vs. Speculative System



Nonspeculative schedule

Long control flow decision prevents pipelining
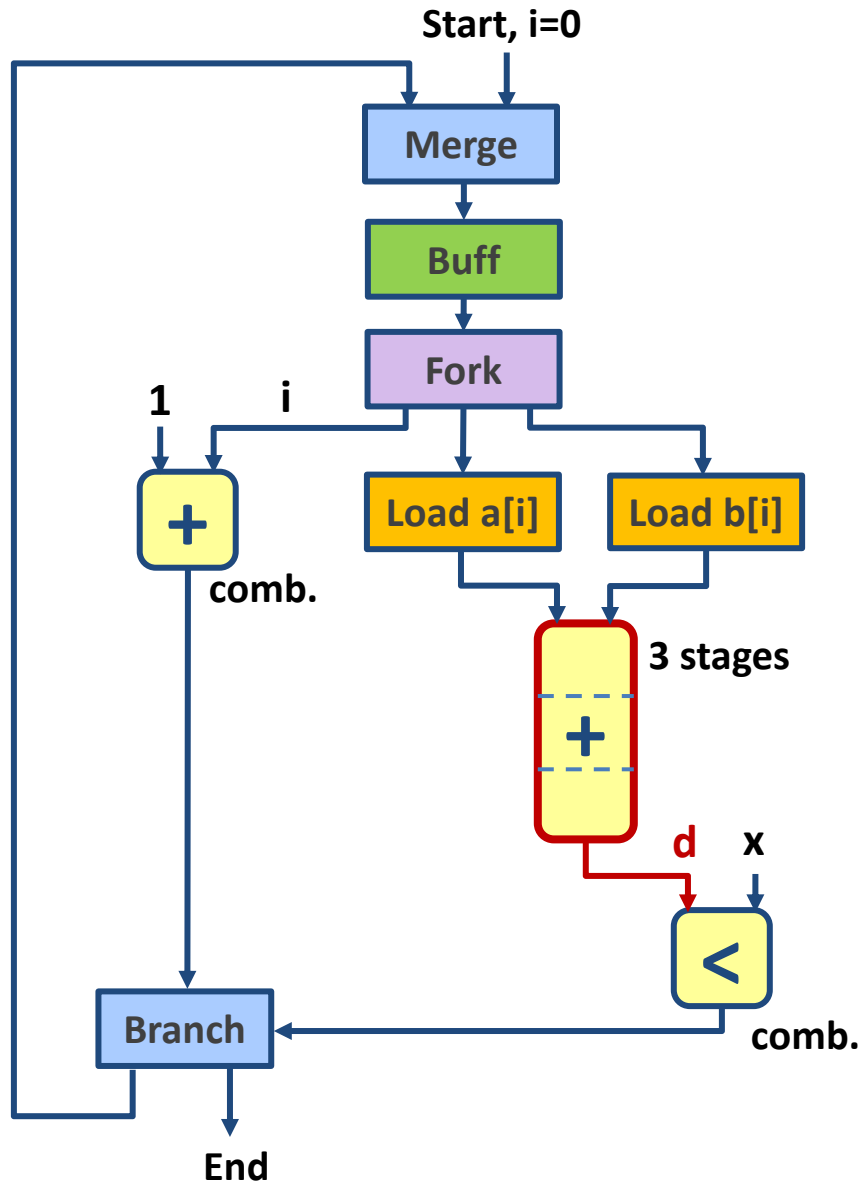
# Nonspeculative vs. Speculative System

# Speculation in Dataflow Circuits

- Contain speculation in a region of the circuit delimited by special components
  - Issue speculative tokens (pieces of data which might or might not be correct)
  - Squash and replay in case of misspeculation



data + handshake
speculative tag

Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019

# Speculation in Dataflow Circuits

- Contain speculation in a region of the circuit delimited by special components
  - Issue speculative tokens (pieces of data which might or might not be correct)
  - Squash and replay in case of misspeculation



Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019

# Speculation in Dataflow Circuits

- Contain speculation in a region of the circuit delimited by special components
  - Issue speculative tokens (pieces of data which might or might not be correct)
  - Squash and replay in case of misspeculation



Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019

# Speculation in Dataflow Circuits

- Contain speculation in a region of the circuit delimited by special components
  - Issue speculative tokens (pieces of data which might or might not be correct)
  - Squash and replay in case of misspeculation



Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019
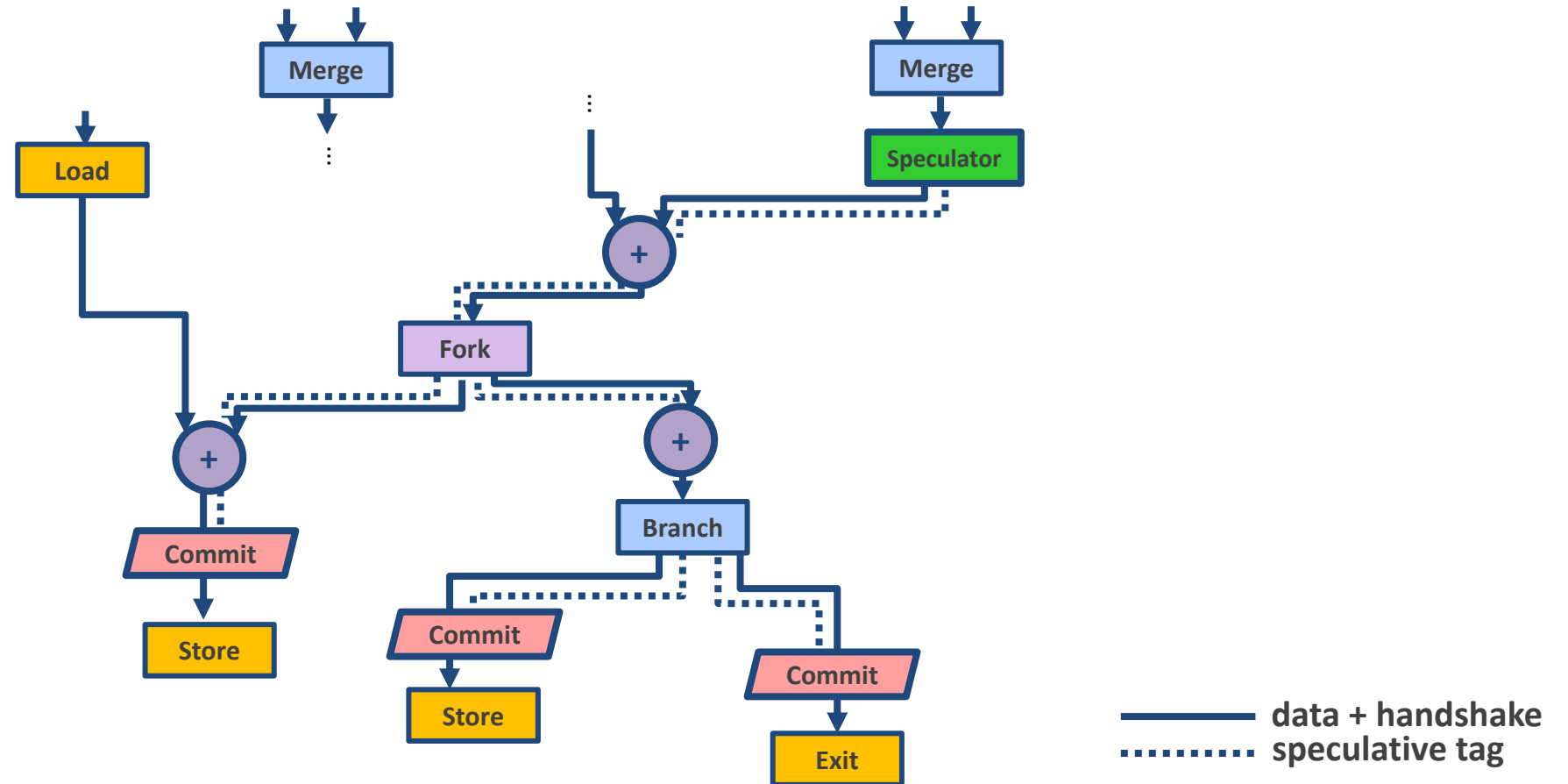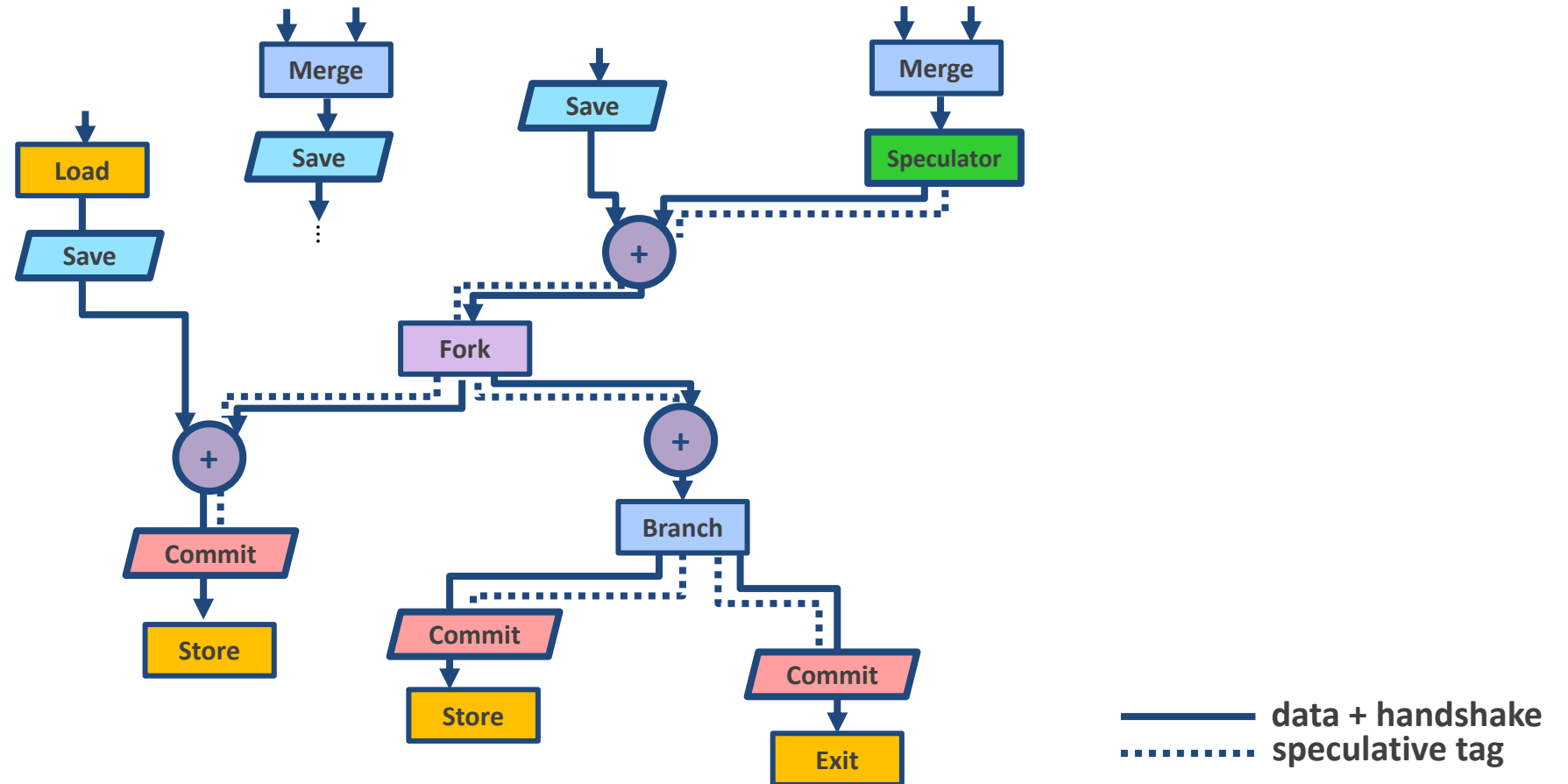
# Speculation in Dataflow Circuits
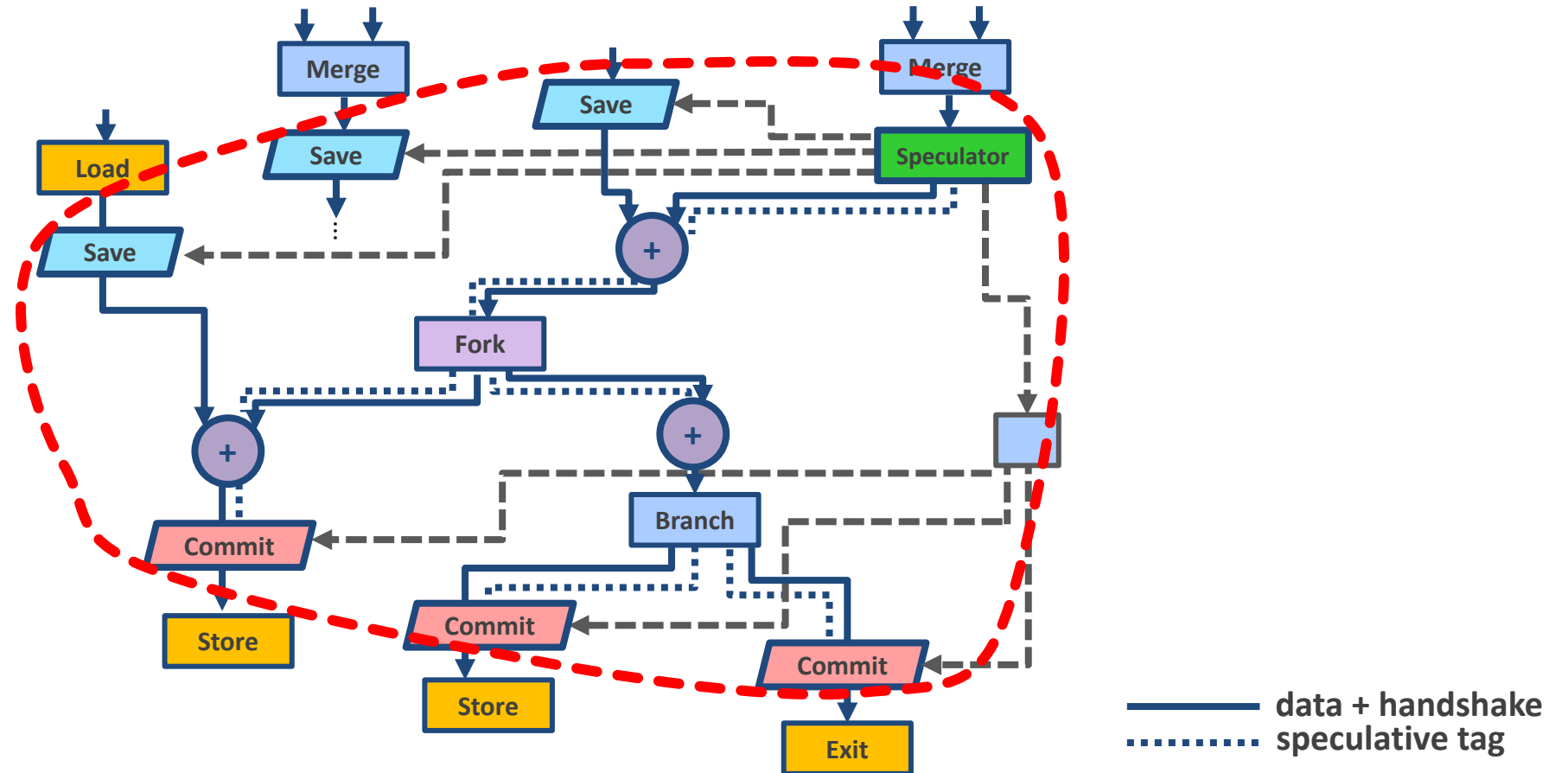
- Contain speculation in a region of the circuit delimited by special components
  - Issue speculative tokens (pieces of data which might or might not be correct)
  - Squash and replay in case of misspeculation



Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019

# Speculative Dataflow Circuit

Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019

# Speculative Dataflow Circuit



Speculator instead of regular Branch

Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019

# Speculative Dataflow Circuit



Input boundary:
Save units

Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019

# Speculative Dataflow Circuit



Output boundary:
Commit units

Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019

# Speculative Dataflow Circuit



Continue computing before condition known

Wait for long-latency condition

BEFORE (without speculation)

# Speculative Dataflow Circuit



**High-throughput speculative pipeline**

Josipović, Guerrieri, and Ienne. Speculative Dataflow Circuits. FPGA 2019

# HLS of Dynamically Scheduled Circuits



**Catching up with static HLS**

Pipelining

Resource sharing

Mul 1    Mul 2    Mul 1/2

**Reaping the benefits of dynamic scheduling**

Out-of-order memory

Speculative execution

## Static HLS vs. dynamic HLS?

# Dynamatic: An Open-Source HLS Compiler

- From C/C++ to synthesizable dataflow circuit description

Josipović, Guerrieri, and Ienne. Dynamatic: From C/C++ to Dynamically Scheduled Circuits. FPGA 2020

# Experimental Results

- Resource utilization and execution time of the dataflow designs, **normalized to the corresponding static designs** produced by Vivado HLS

▲ Dynamic, control dependences

● Dynamic, memory dependences

◆ Dynamic, speculative

✕ Dynamic, no dependences

■ Static (all points)

# Experimental Results

- Resource utilization and execution time of the dataflow designs, **normalized to the corresponding static designs** produced by Vivado HLS



Josipović, Guerrieri, and Ienne. Synthesizing General-Purpose Code into Dynamically Scheduled Circuits. CASM 2021

# Experimental Results

- Resource utilization and execution time of the dataflow designs, **normalized to the corresponding static designs** produced by Vivado HLS

Josipović, Guerrieri, and Ienne. Synthesizing General-Purpose Code into Dynamically Scheduled Circuits. CASM 2021
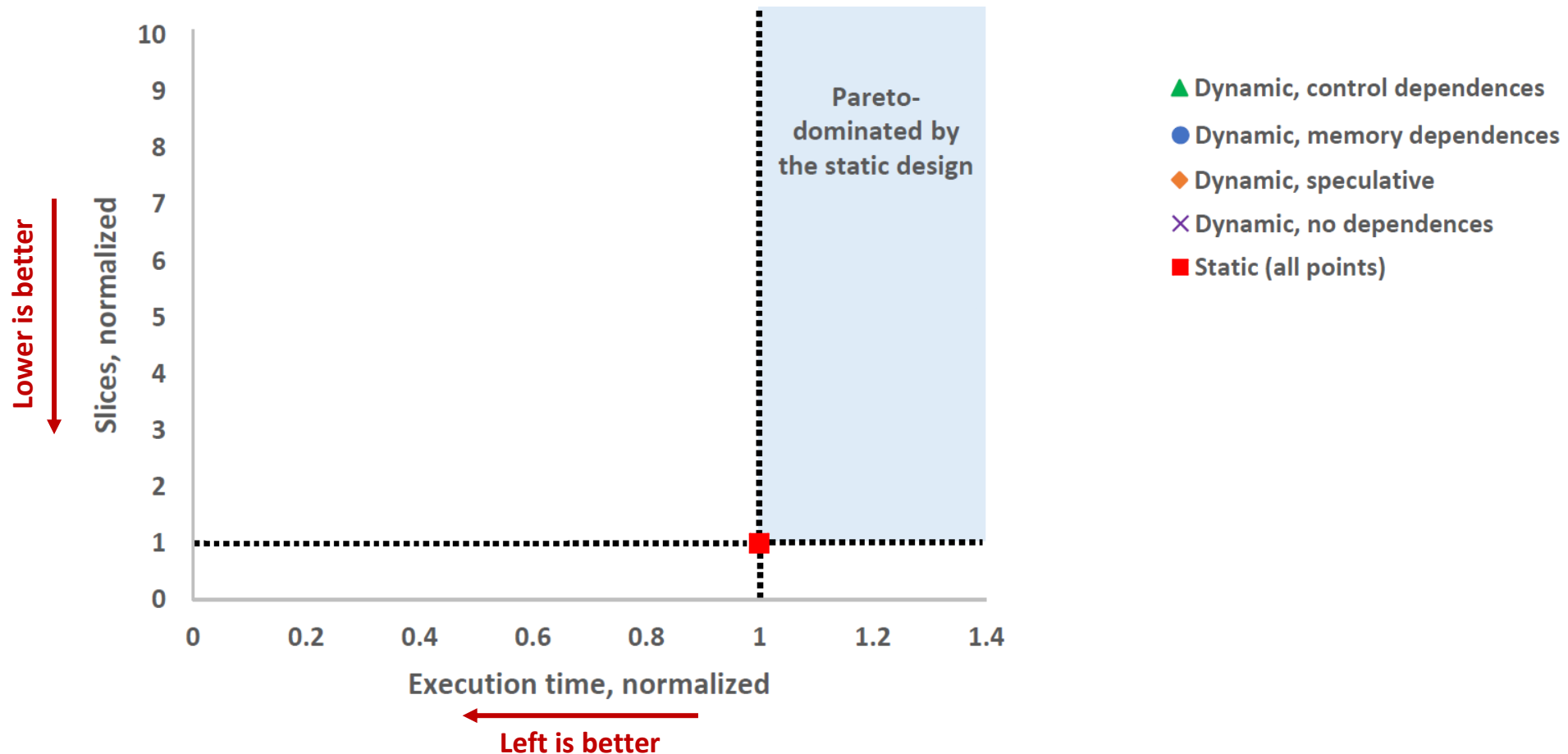
# Experimental Results

- Resource utilization and execution time of the dataflow designs, **normalized to the corresponding static designs** produced by Vivado HLS



Reduced execution time in irregular benchmarks (speedup of up to 14.9X)

Josipović, Guerrieri, and Ienne. Synthesizing General-Purpose Code into Dynamically Scheduled Circuits. CASM 2021
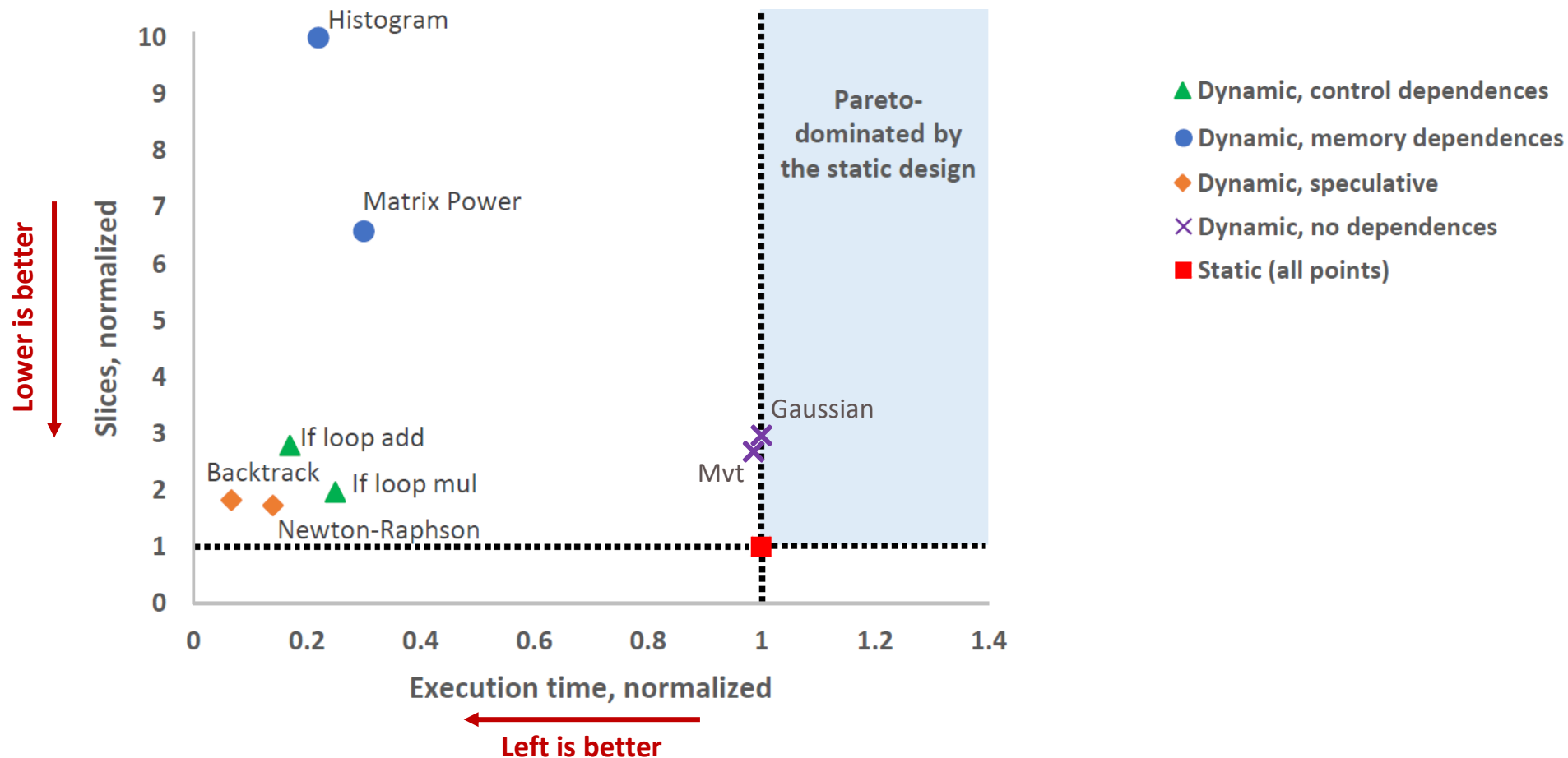
# Experimental Results

- Resource utilization and execution time of the dataflow designs, **normalized to the corresponding static designs** produced by Vivado HLS



**Reduced execution time in irregular benchmarks (speedup of up to 14.9X)**

Josipović, Guerrieri, and Ienne. Synthesizing General-Purpose Code into Dynamically Scheduled Circuits. CASM 2021
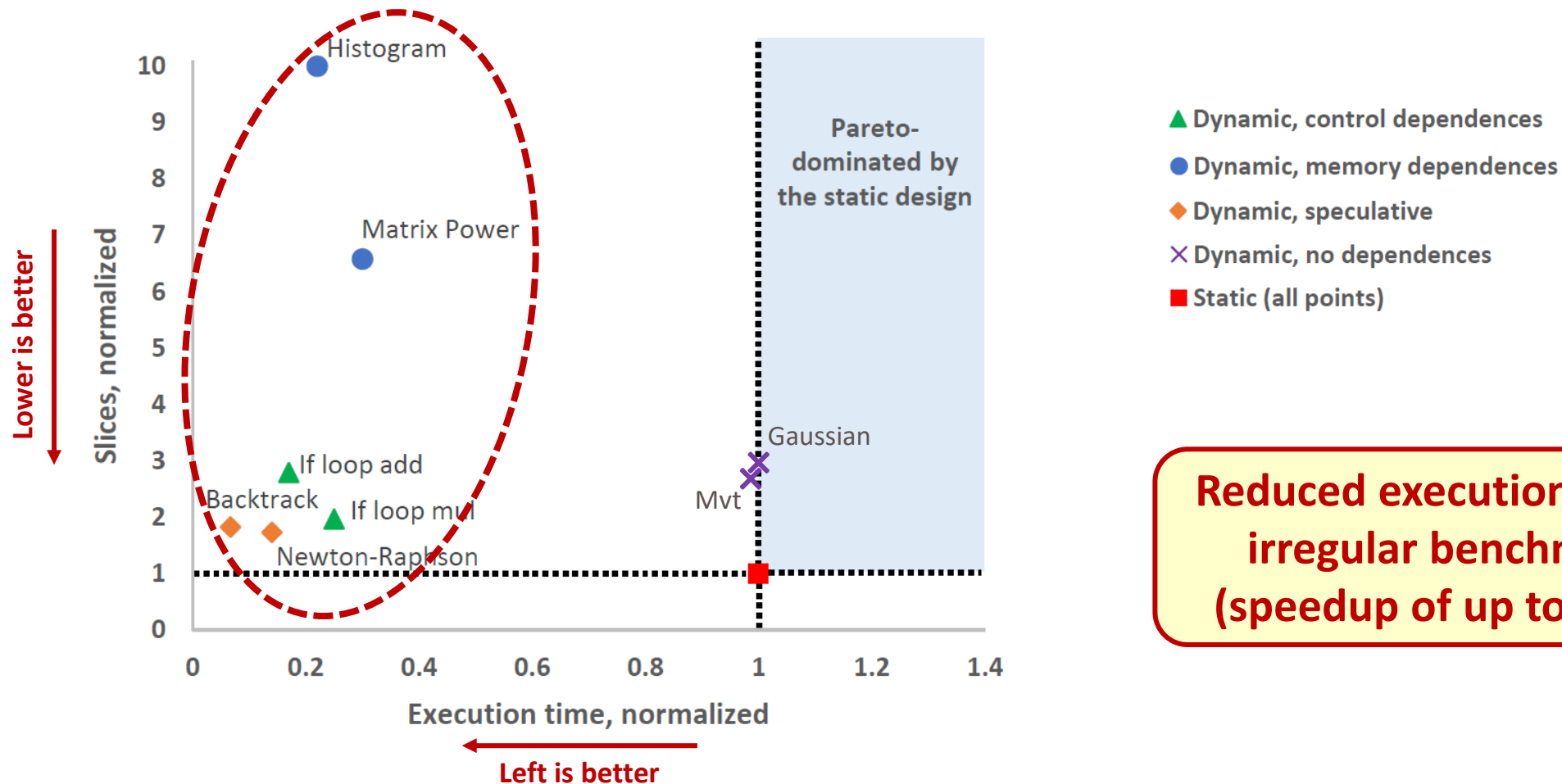
# Experimental Results

- Resource utilization and execution time of the dataflow designs, **normalized to the corresponding static designs** produced by Vivado HLS



**LSQ causes significant resource overheads**

Josipović, Guerrieri, and Ienne. Synthesizing General-Purpose Code into Dynamically Scheduled Circuits. CASM 2021

# Experimental Results

- Resource utilization and execution time of the dataflow designs, **normalized to the corresponding static designs** produced by Vivado HLS
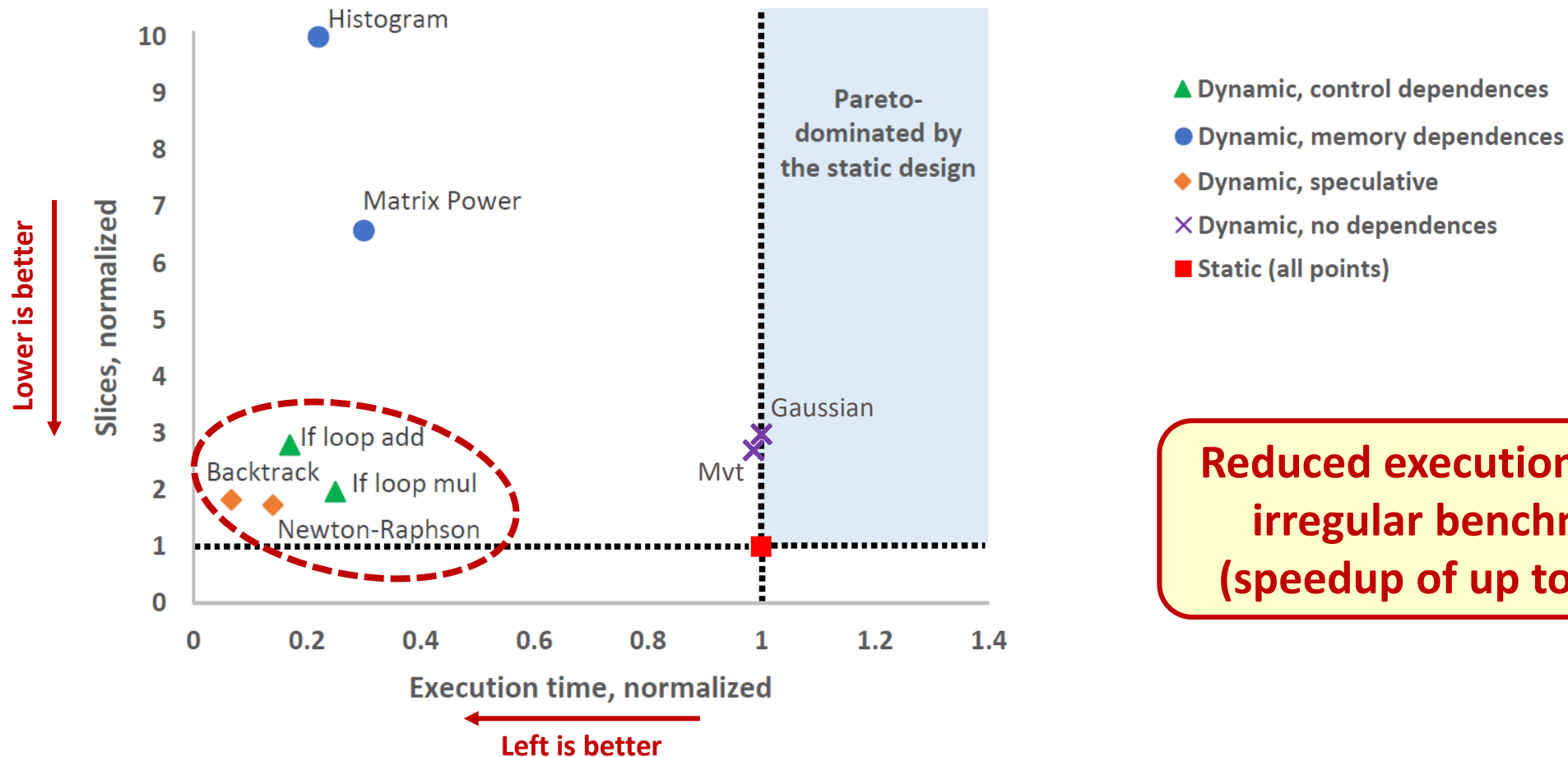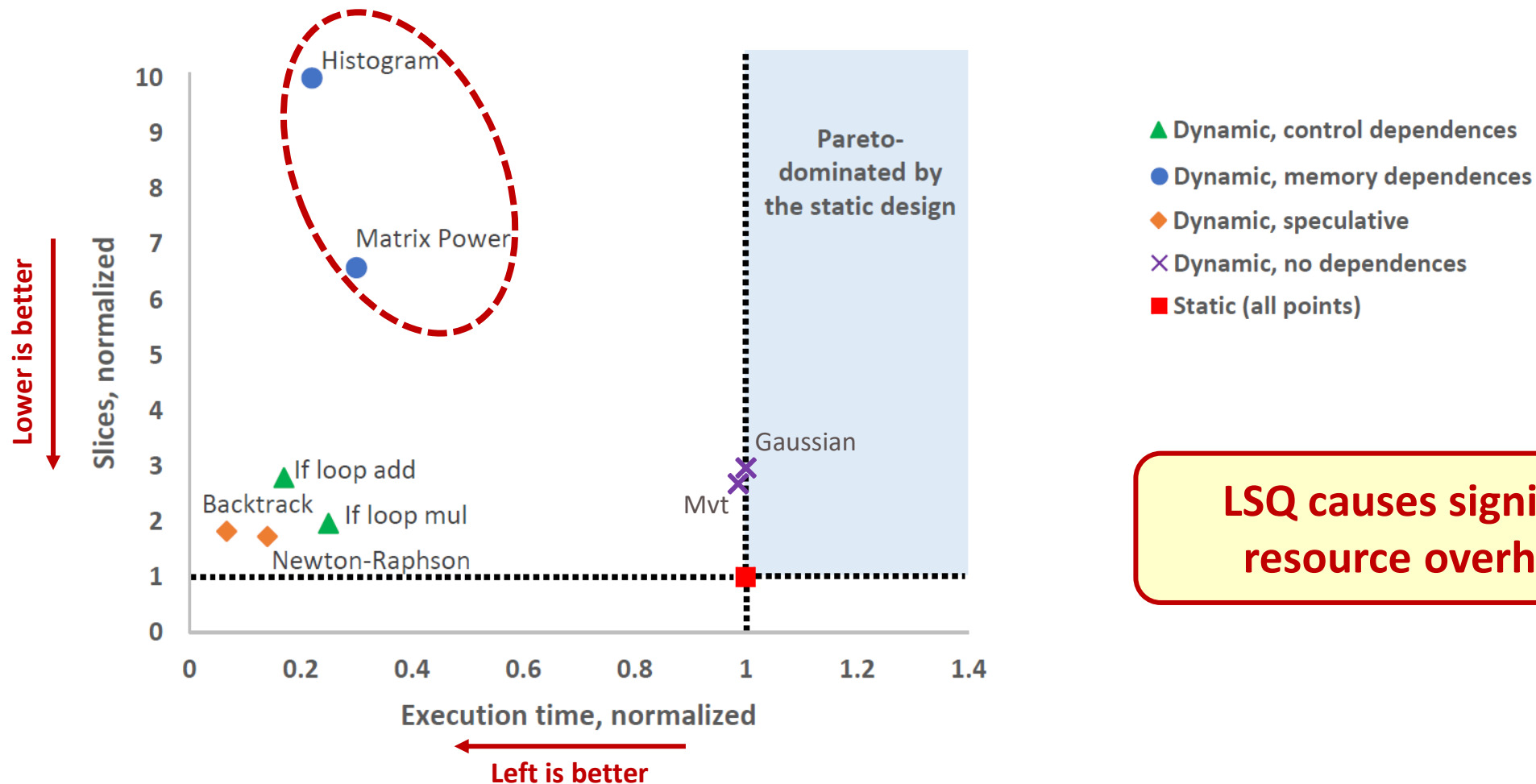
Josipović, Guerrieri, and Ienne. Synthesizing General-Purpose Code into Dynamically Scheduled Circuits. CASM 2021
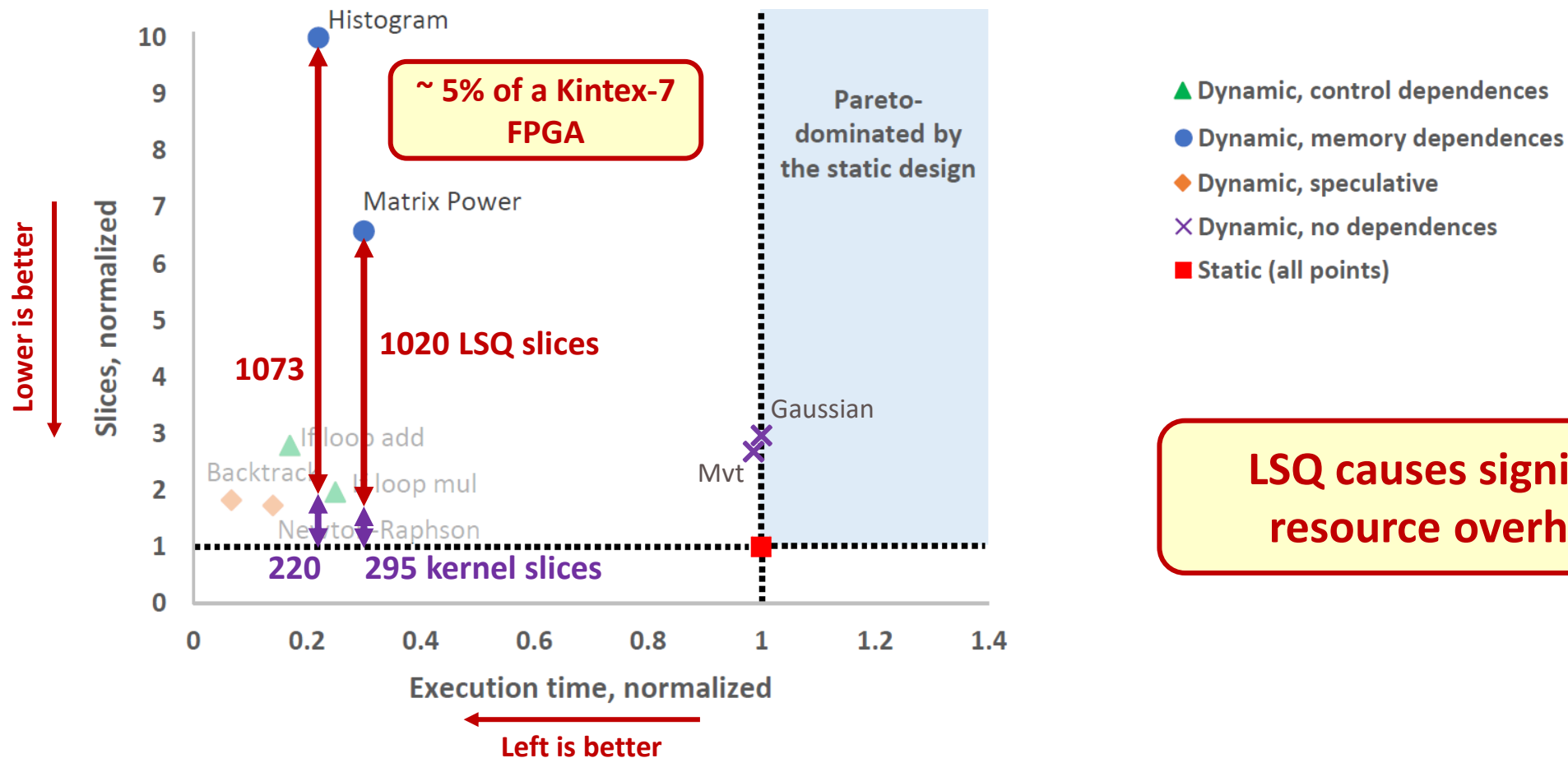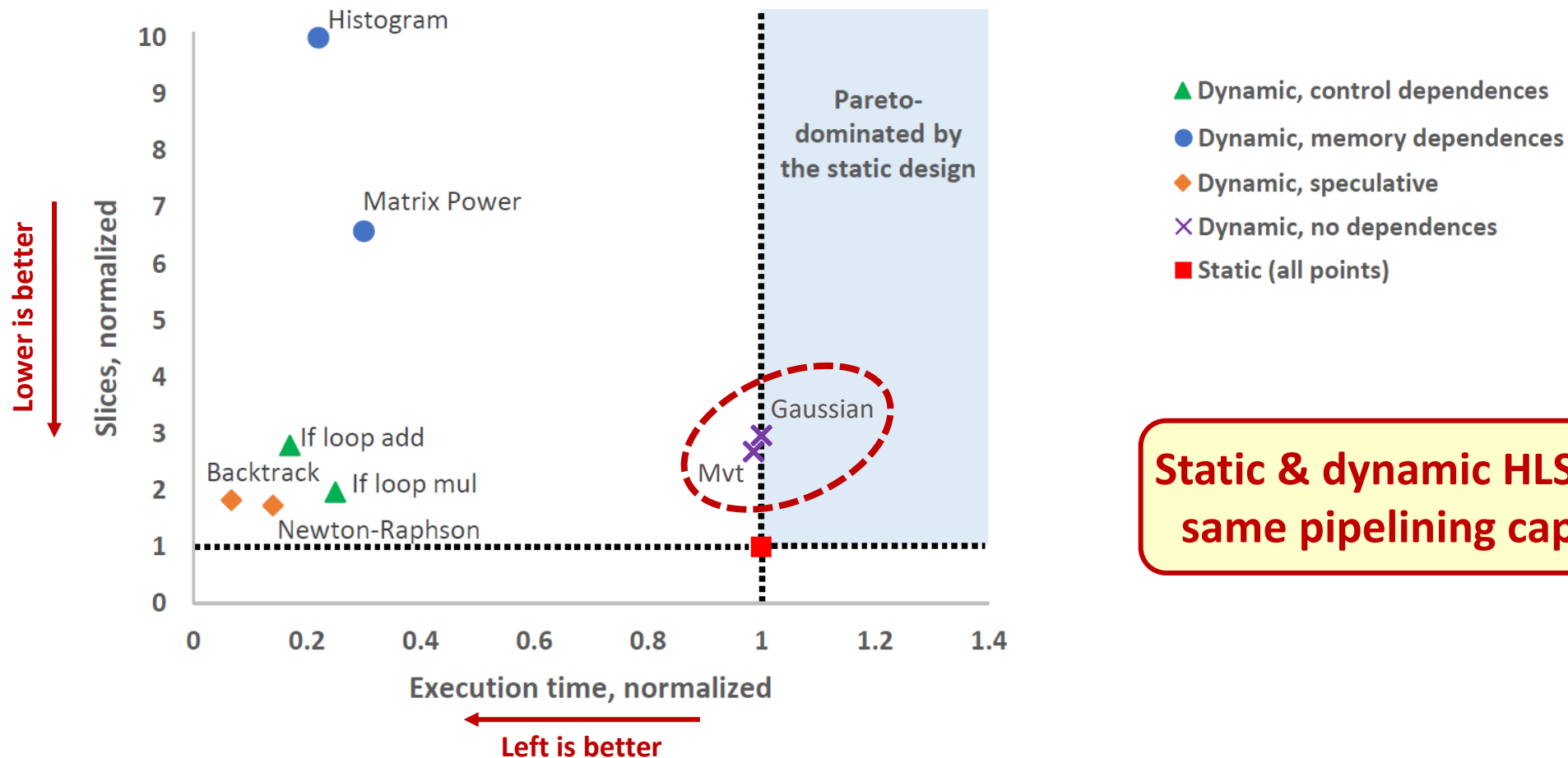
# Experimental Results

- Resource utilization and execution time of the dataflow designs, **normalized to the corresponding static designs** produced by Vivado HLS



Static & dynamic HLS have the same pipelining capabilities

Josipović, Guerrieri, and Ienne. Synthesizing General-Purpose Code into Dynamically Scheduled Circuits. CASM 2021

# Static vs. Dynamic Scheduling



|  | **Statically Scheduled**<br>→ "Compiler does the job" | **Dynamically Scheduled**<br>→ "Hardware does the job" |
|---|---|---|
| Computer<br>Architecture | **VLIW<br>Processors** | **Out-of-Order<br>Superscalar<br>Processors** |
| High-Level<br>Synthesis | **Traditional HLS** | **Dataflow circuits** |

**DSP-oriented applications**

**General-purpose code<br>(new applications and users)**

# Thanks! ☺

**Research group:**



**https://dynamo.ethz.ch/**

**Dynamatic HLS tool:**



**https://dynamatic.epfl.ch/**