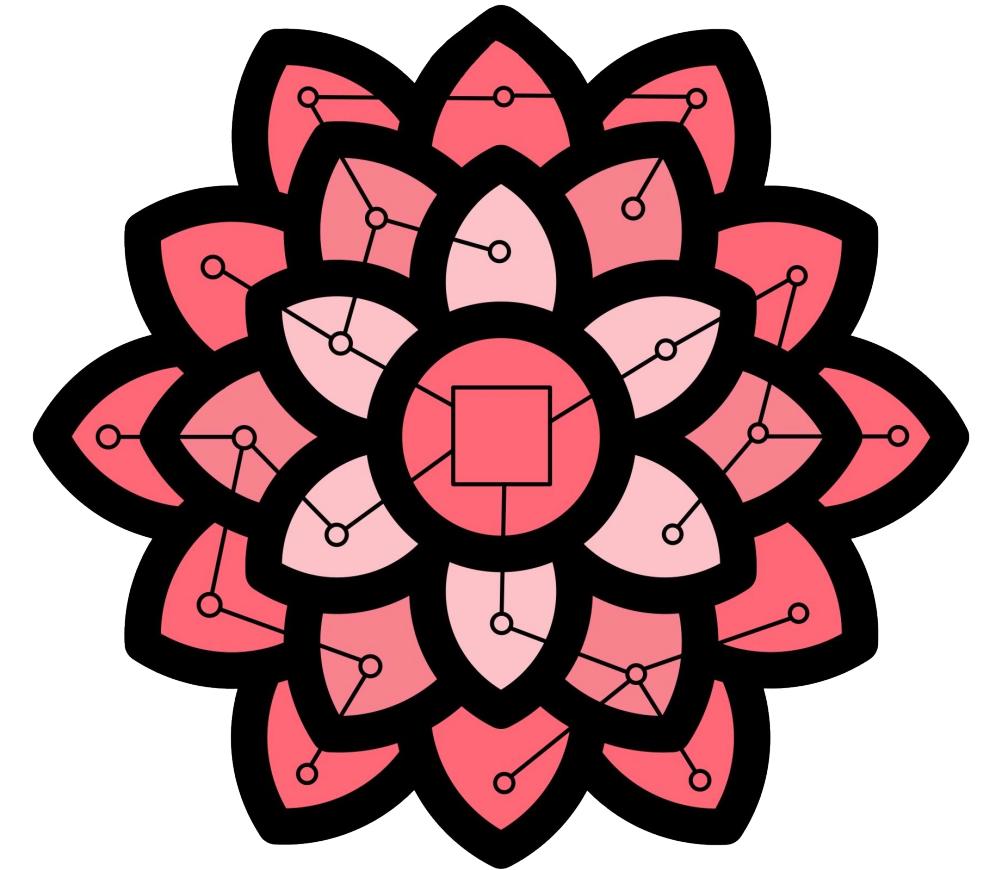


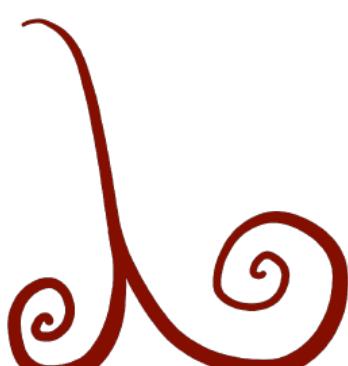
Predictable Accelerator Design with Time-Sensitive Affine Types



Rachit Nigam, Sachille Atapattu, Samuel Thomas, Zhijing Li, Theodore Bauer,
Yuwei Ye, Apurva Koti, Adrian Sampson, Zhiru Zhang

capra.cs.cornell.edu/dahlia

Computer Architecture and
Programming Abstractions Group





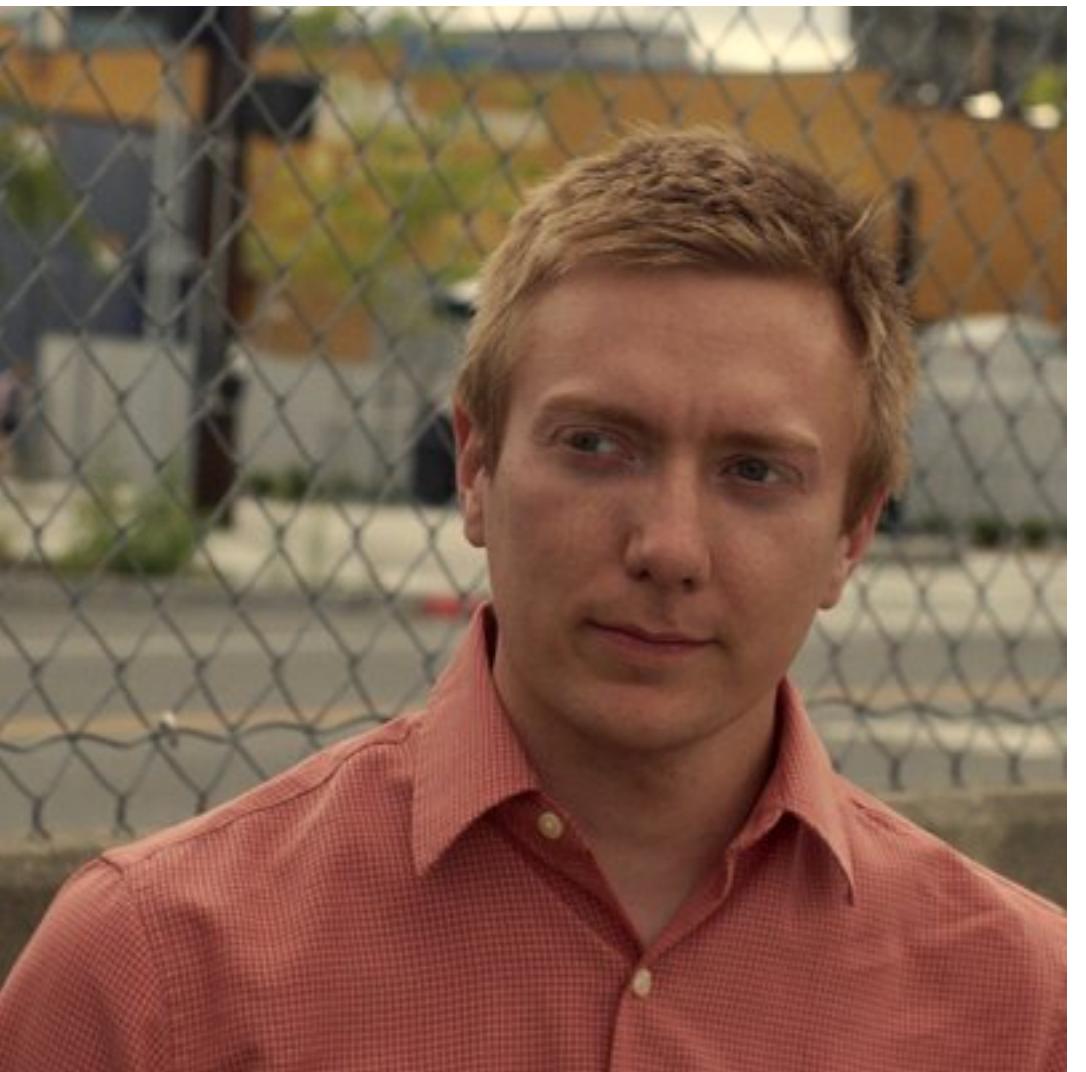
Rachit Nigam

- Second year PhD
- Computer Architect by day
- Programming languages by night



Rachit Nigam

- Second year PhD
- Computer Architect by day
- Programming languages by night



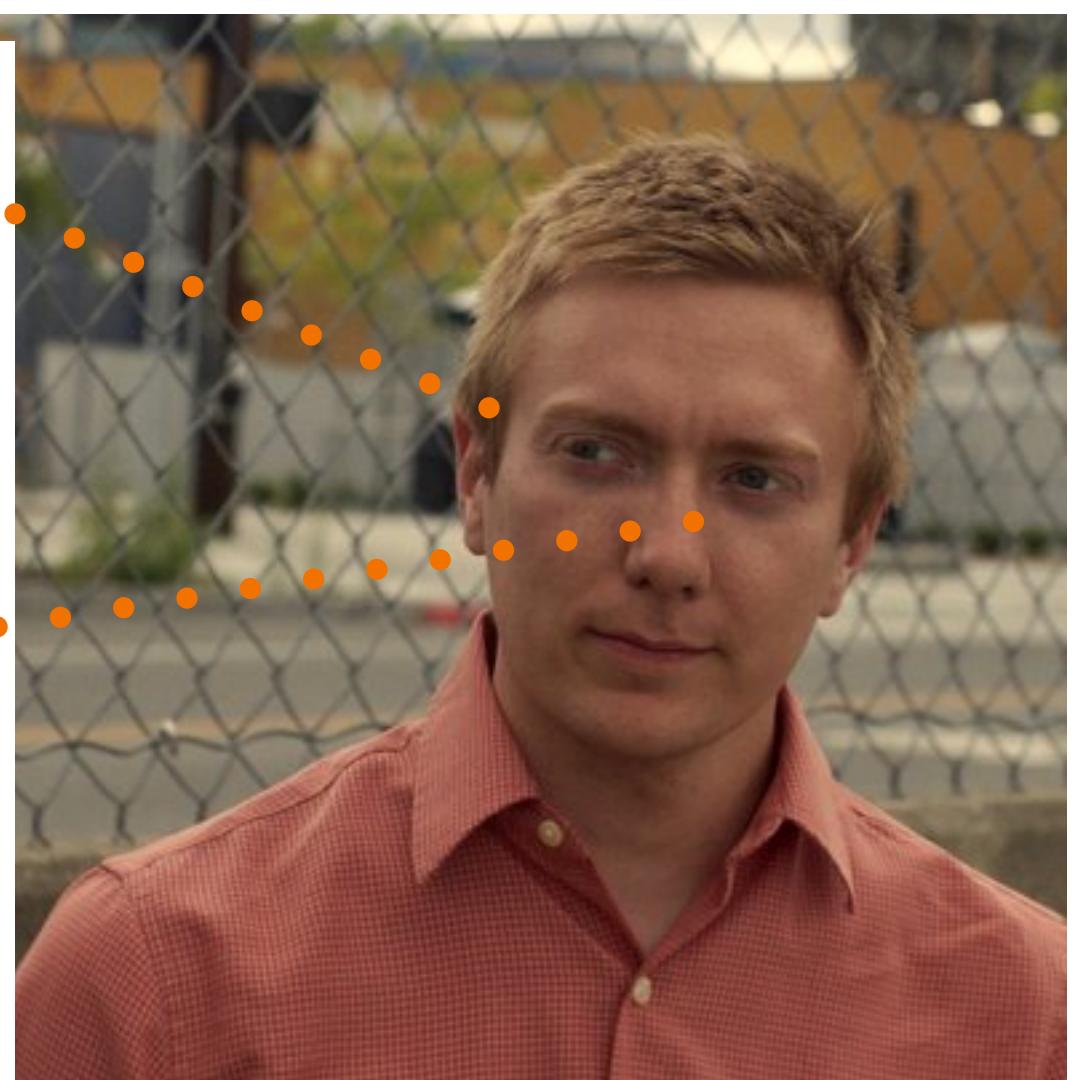
Adrian Sampson

- Nth year PhD
- PL & Compilers & Architecture



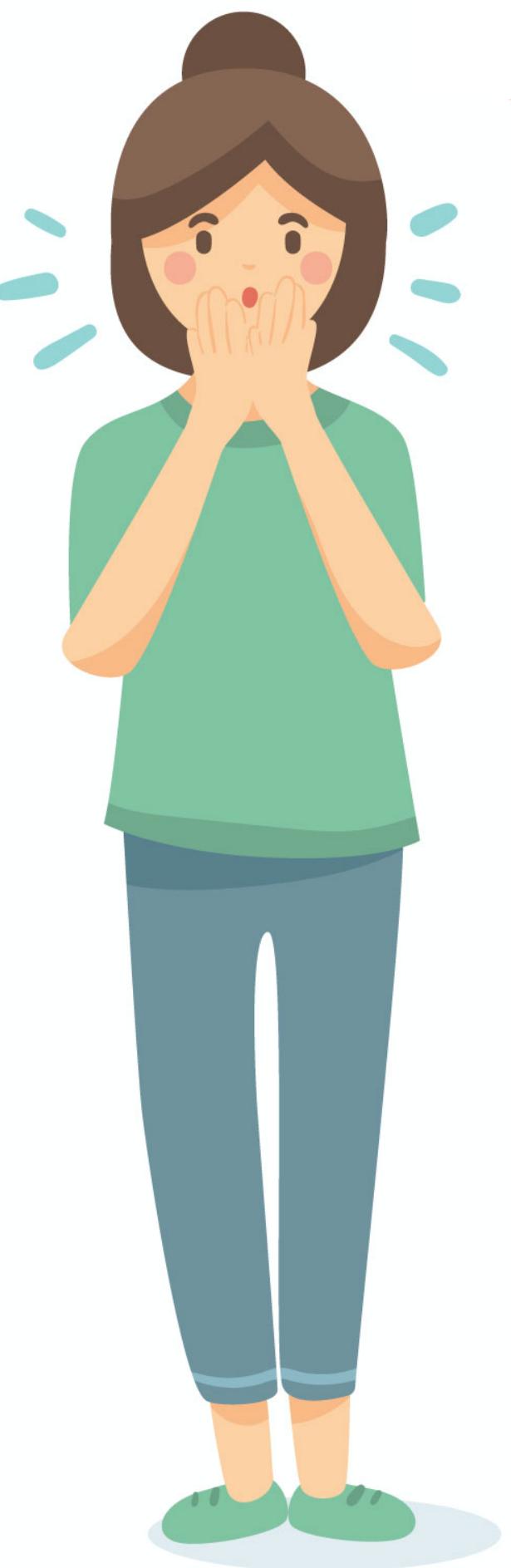
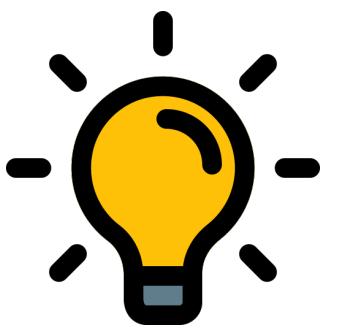
Rachit Nigam

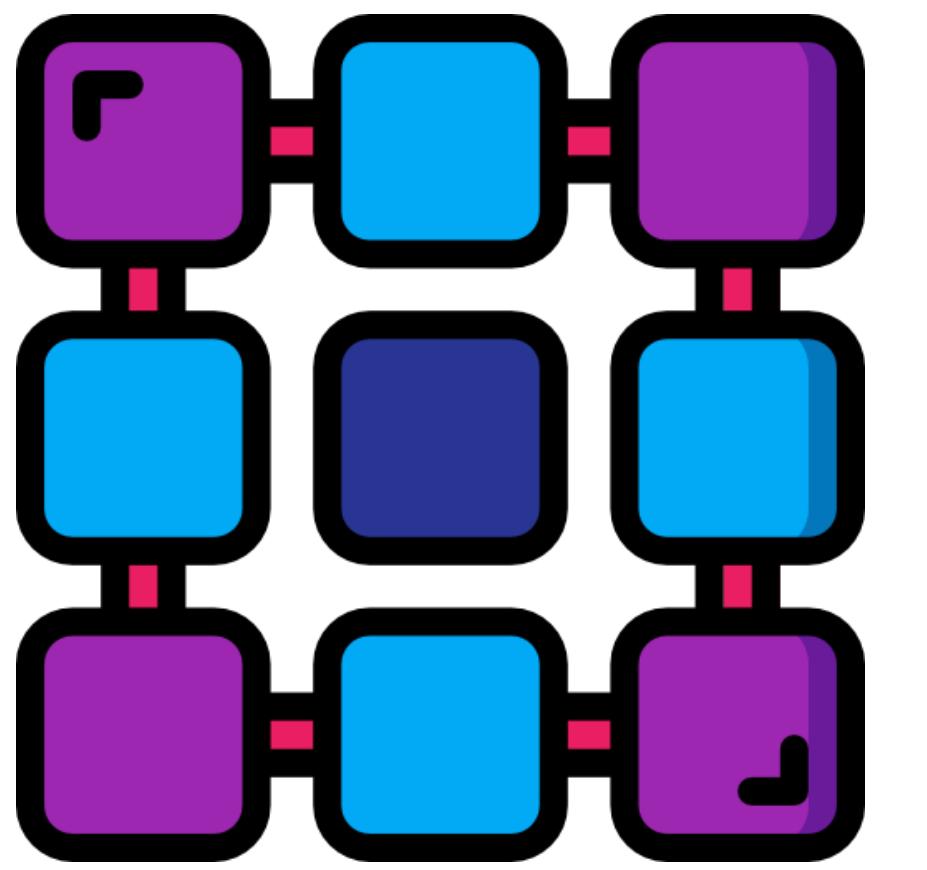
- Second year PhD
- Computer Architect by day
- Programming languages by night



Adrian Sampson

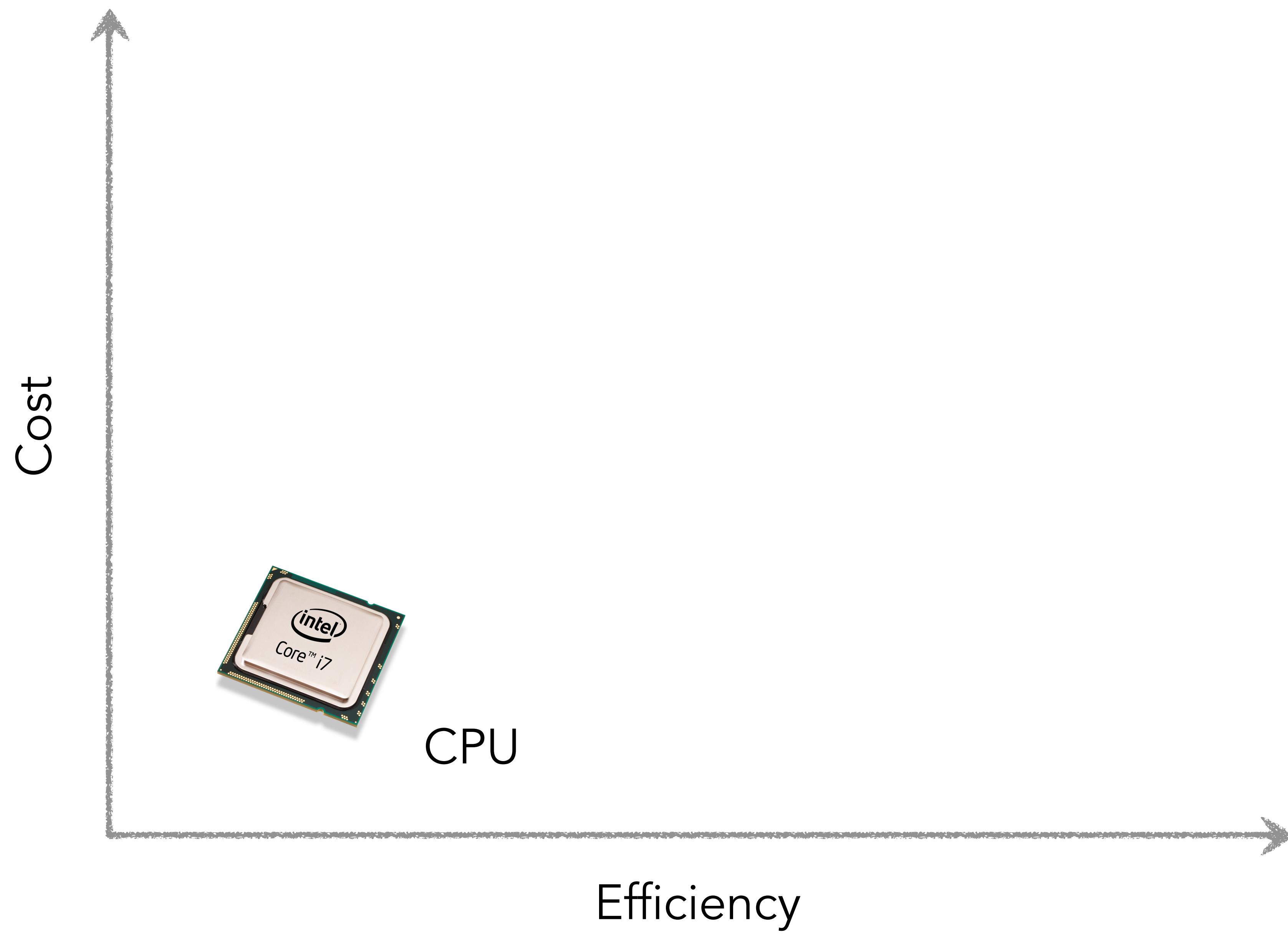
- Nth year PhD
- PL & Compilers & Architecture

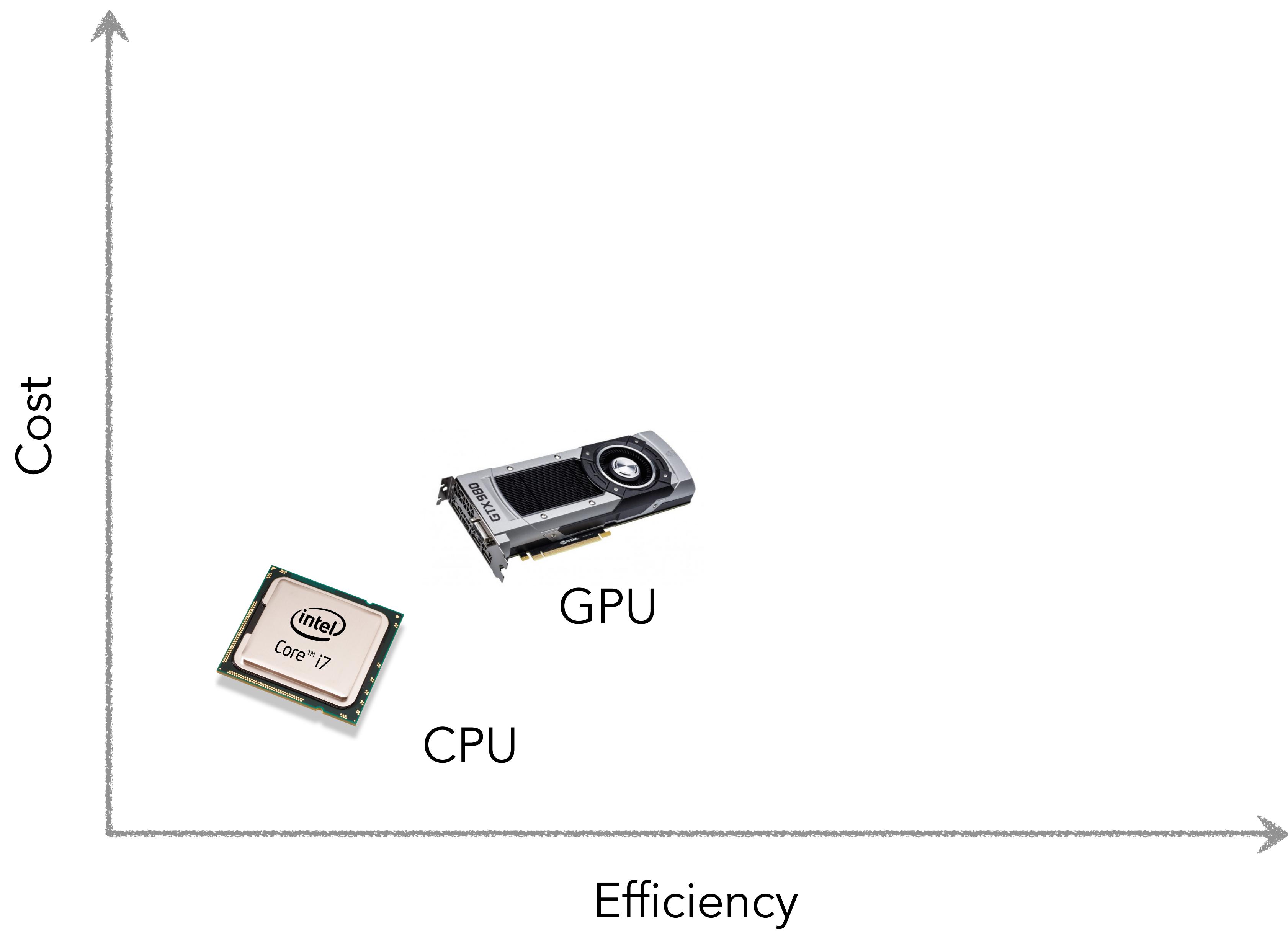


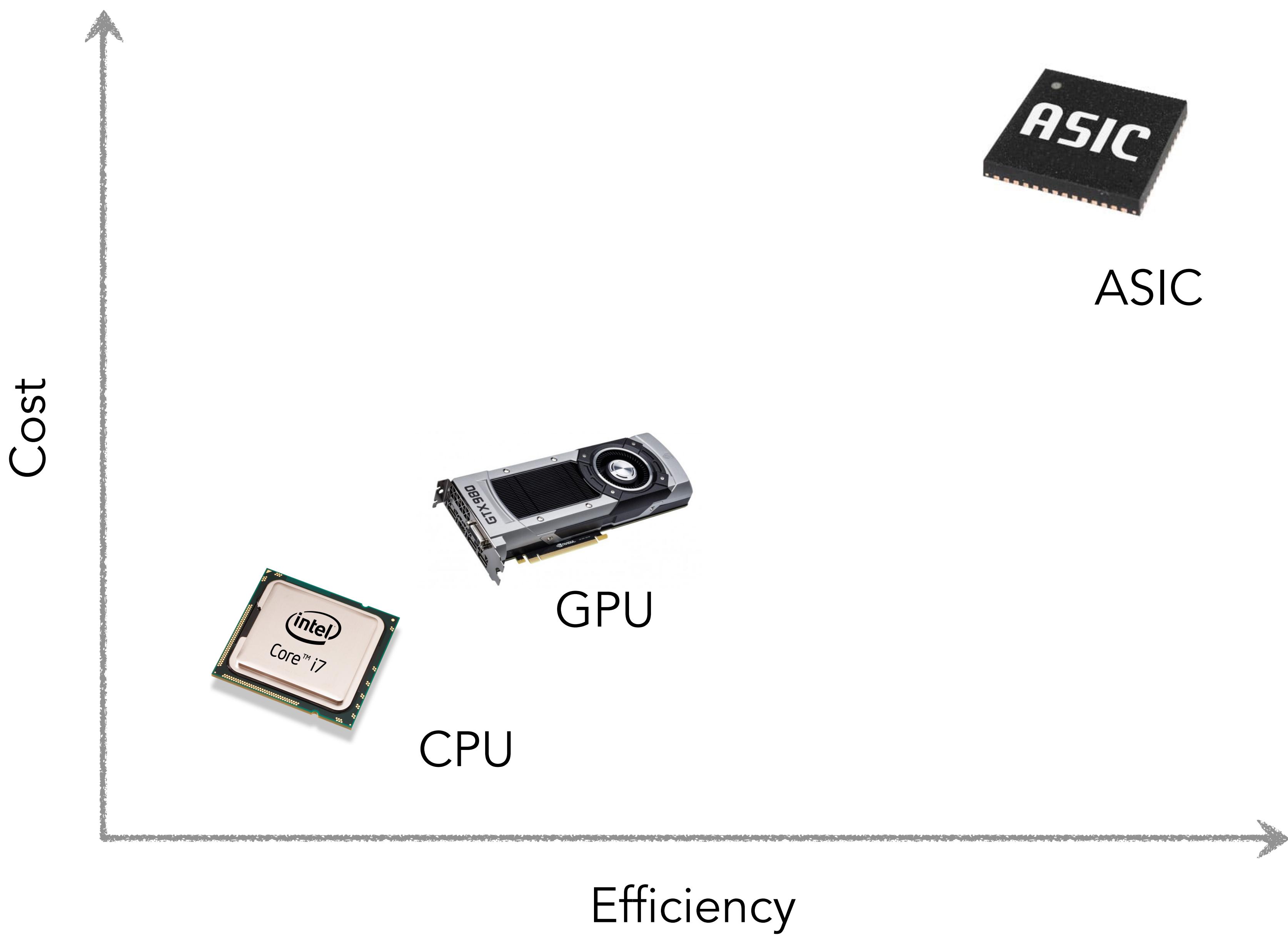


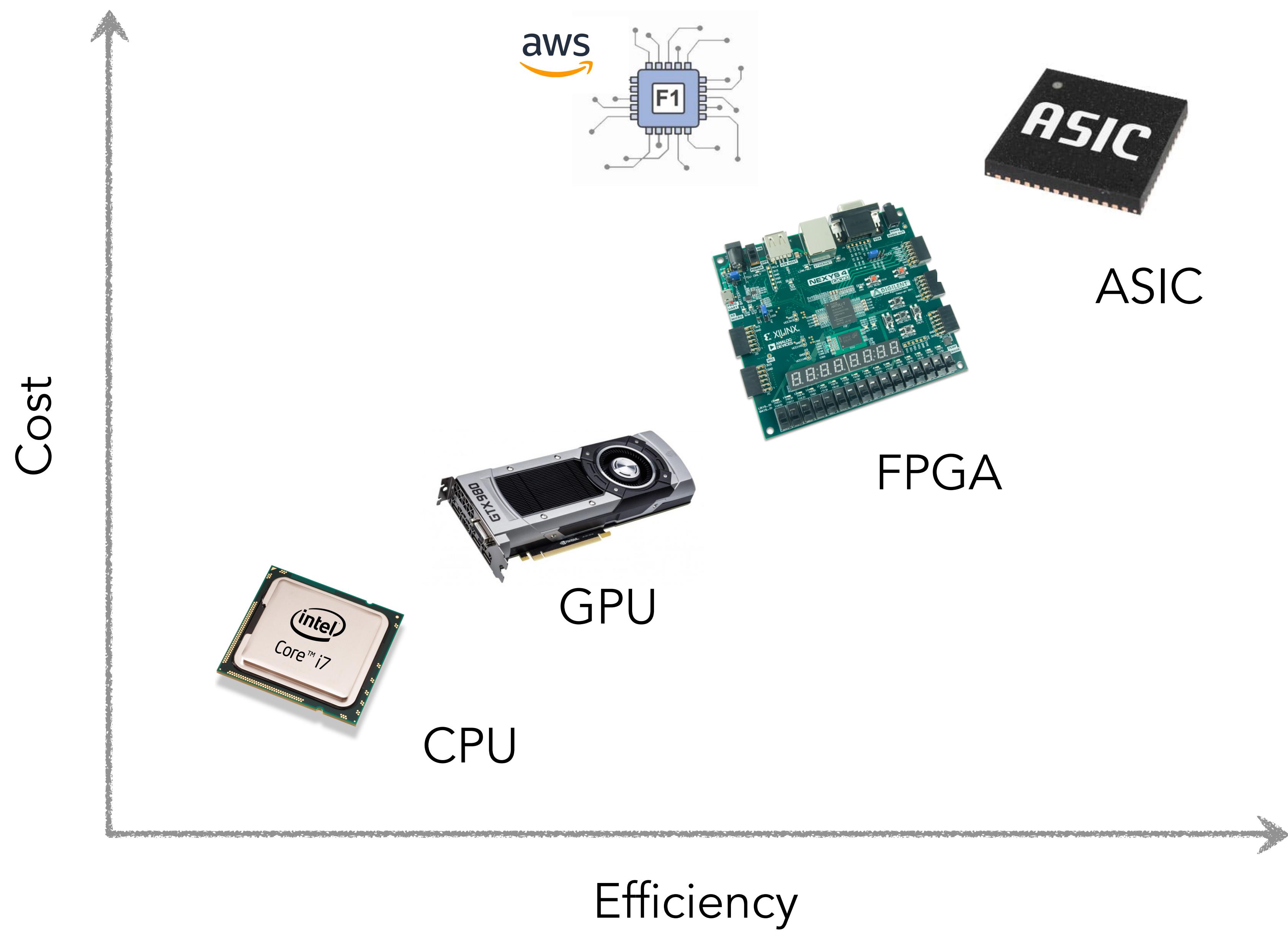


Efficiency

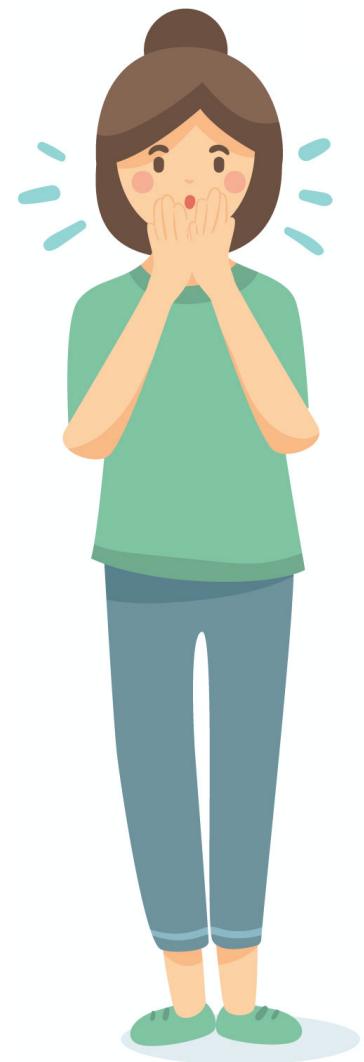








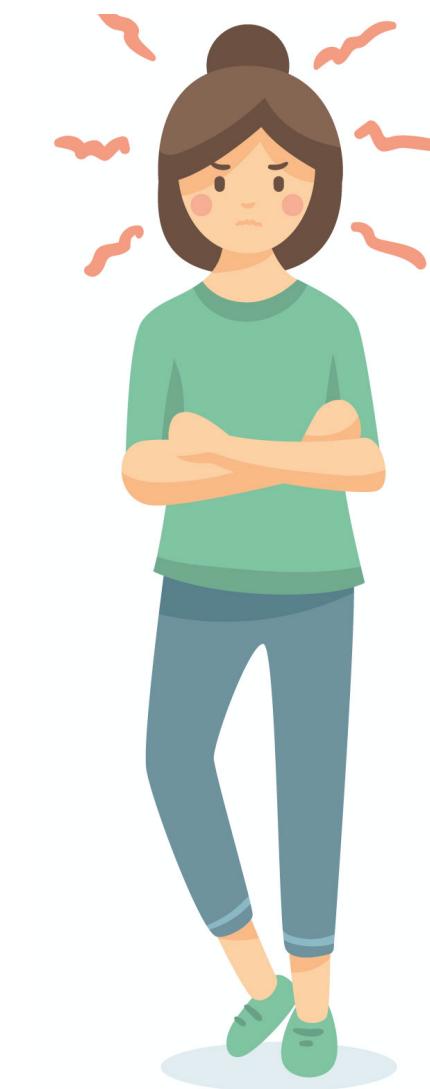
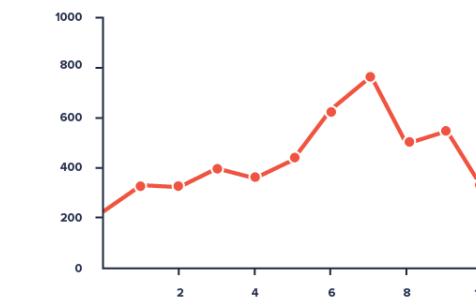
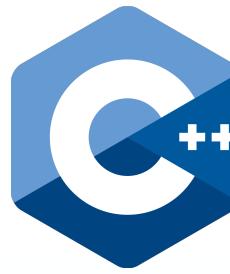
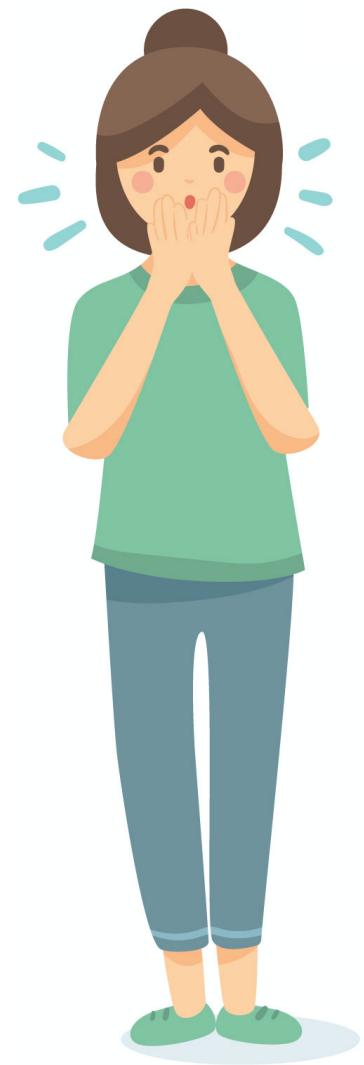
Ada's Journey



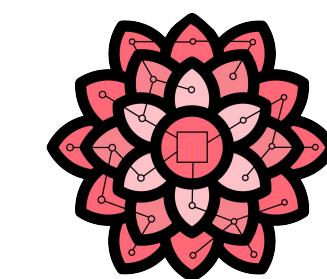
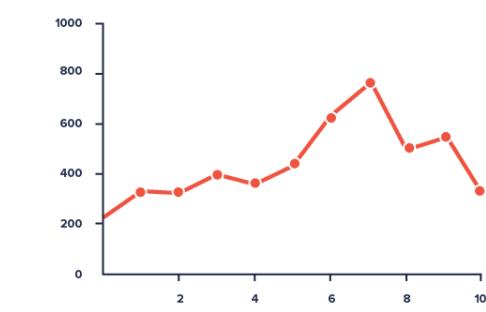
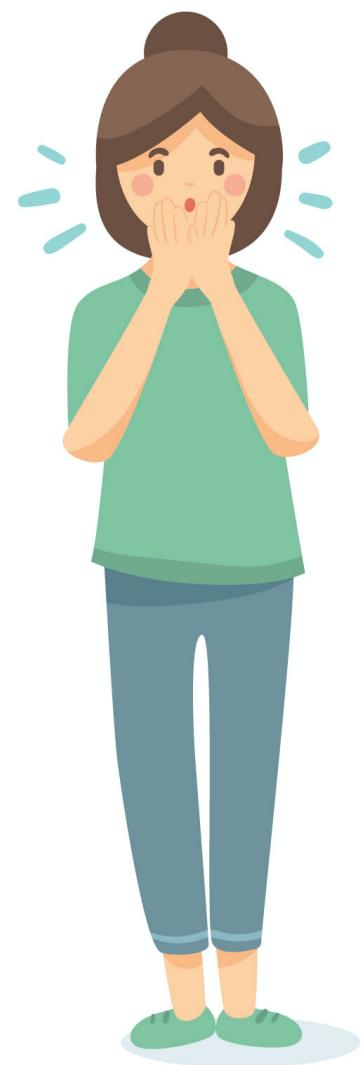
Ada's Journey



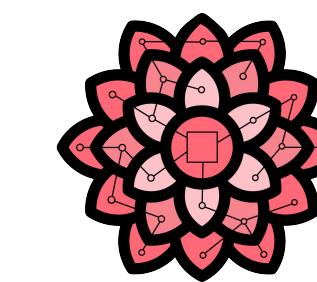
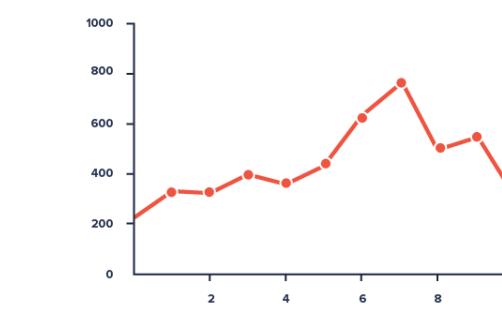
Ada's Journey



Ada's Journey

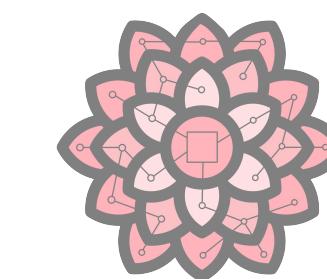
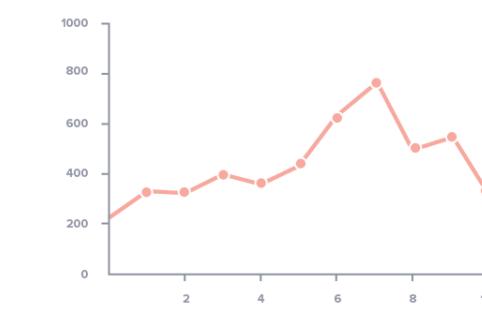
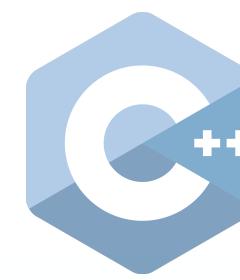


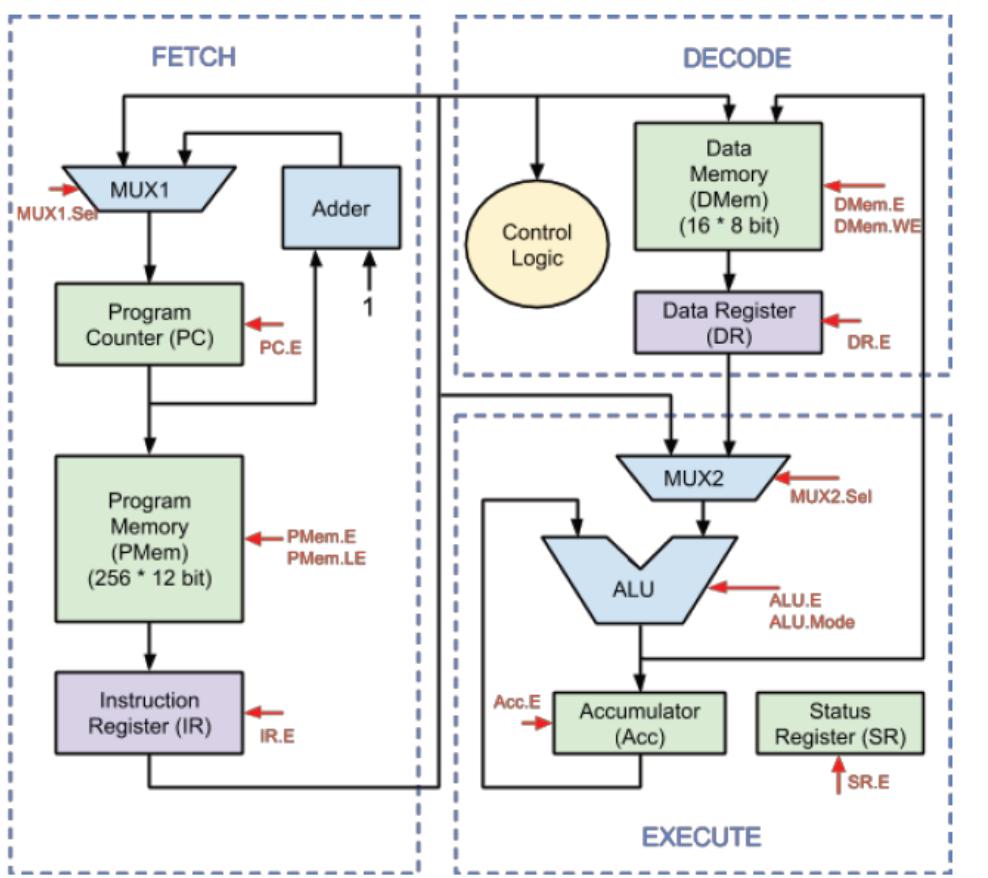
Ada's Journey



Our Research

Ada's Journey

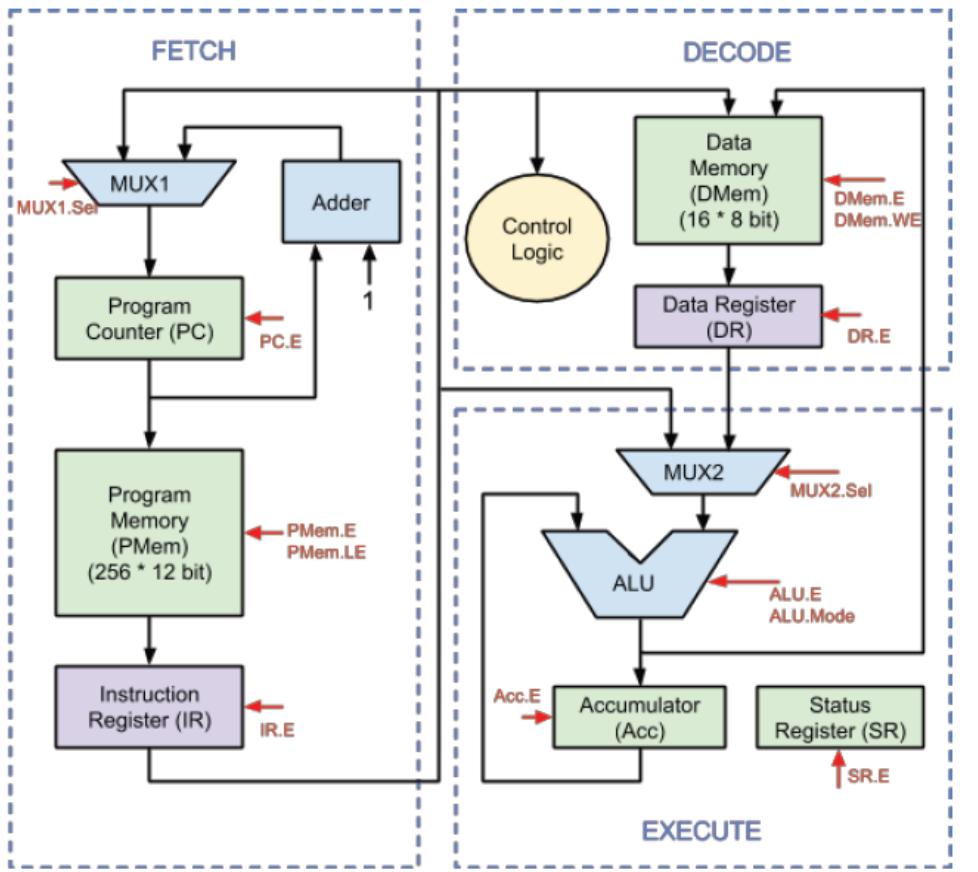




Hardware Design



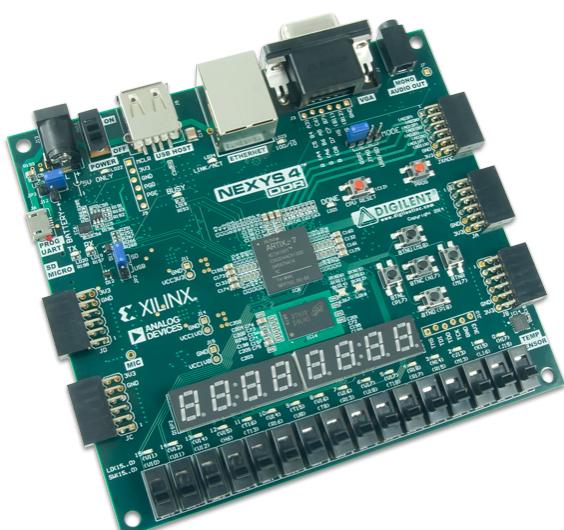
FPGA



Hardware Design



Traditional HDLs



FPGA



Traditional HDLs

```
module ctr (input          up_down,
            clk,
            rstn,
            output reg [2:0]    out);

  always @ (posedge clk)
    if (!rstn)
      out <= 0;
    else begin
      if (up_down)
        out <= out + 1;
      else
        out <= out - 1;
    end
endmodule
```



Traditional HDLs

```
module ctr (input          up_down,
            clk,
            rstn,
            output reg [2:0]    out);

  always @ (posedge clk)
    if (!rstn)
      out <= 0;
    else begin
      if (up_down)
        out <= out + 1;
      else
        out <= out - 1;
    end
endmodule
```

- Time is real



Traditional HDLs

```
module ctr (input          up_down,
            clk,
            rstn,
            output reg [2:0]    out);

  always @ (posedge clk)
    if (!rstn)
      out <= 0;
    else begin
      if (up_down)
        out <= out + 1;
      else
        out <= out - 1;
    end
endmodule
```

- Time is real
- Concurrent semantics



Traditional HDLs

```
neptune_metadata_unpack #(
    .TUSER_W (TUSER_W)
) neptune_metadata_unpack_inst (
    // input
    .metadata      (s_axis_tuser),
    // output
    .std_md        (neptune_std_md),
    .mem_instr_0   (),
    .mem_instr_1   (),
    .mem_instr_2   (),
    .egress_ts     (),
    .fpm           (fpm),
    .reserved      ()
);
```

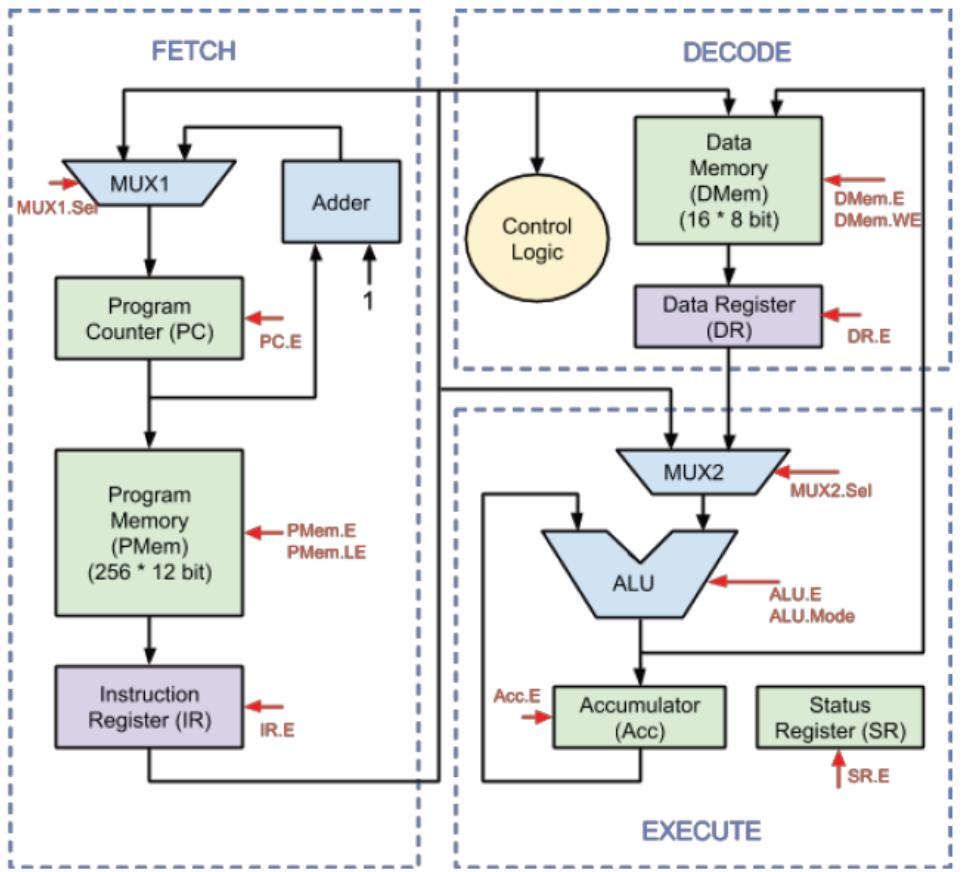
- Time is real
- Concurrent semantics



Traditional HDLs

```
neptune_metadata_unpack #(
    .TUSER_W (TUSER_W)
) neptune_metadata_unpack_inst (
    // input
    .metadata      (s_axis_tuser),
    // output
    .std_md        (neptune_std_md),
    .mem_instr_0   (),
    .mem_instr_1   (),
    .mem_instr_2   (),
    .egress_ts     (),
    .fpm           (fpm),
    .reserved      ()
);
```

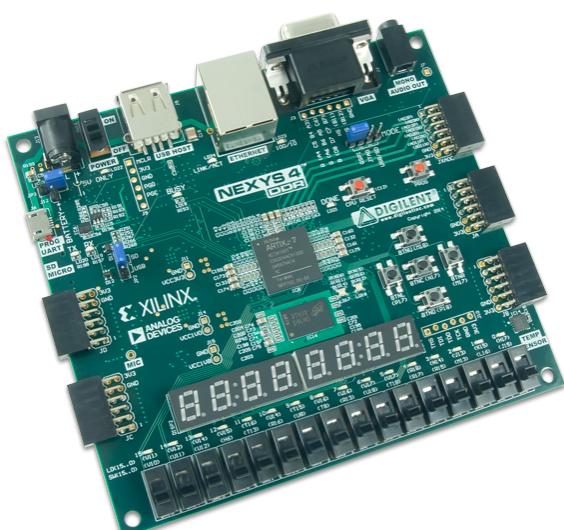
- Time is real
- Concurrent semantics
- **Wires** and hardware modules



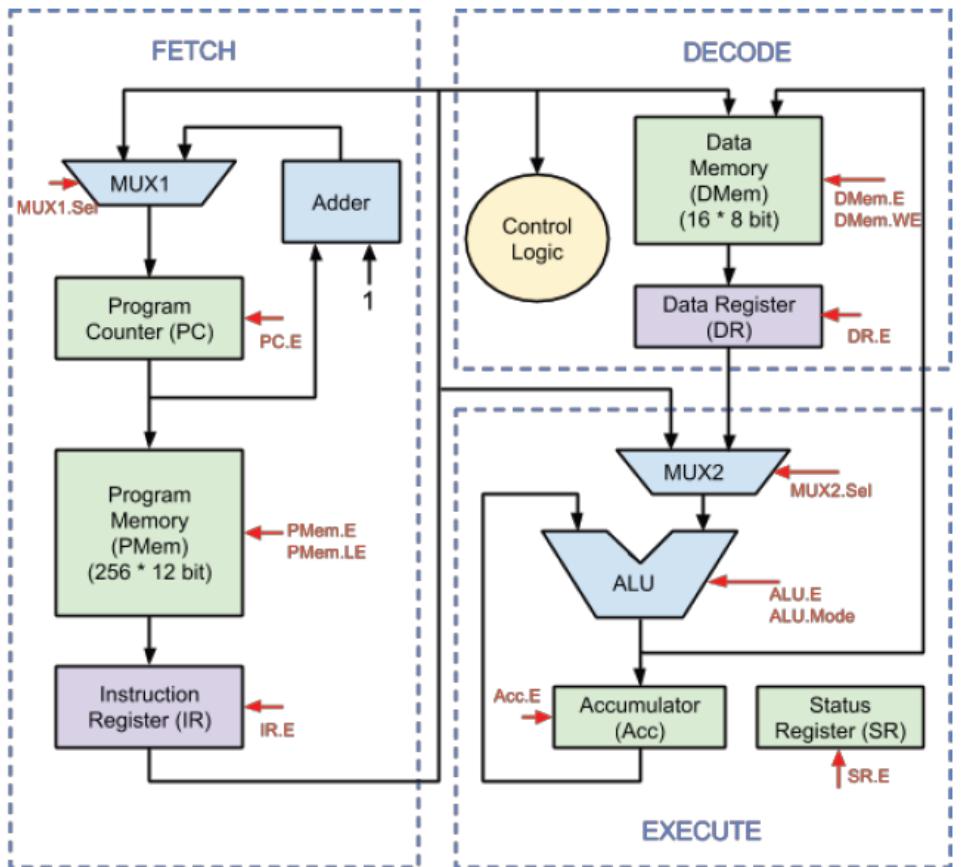
Hardware Design



Traditional HDLs



FPGA



Hardware Design

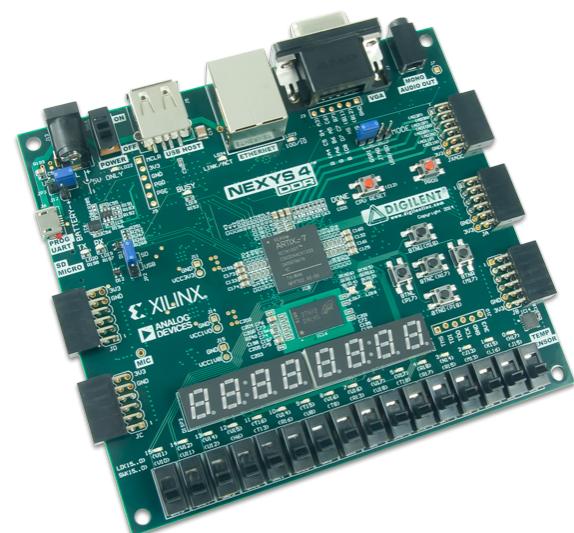


Traditional HDLs

bluespec
CHISEL



Modern HDLs



FPGA

Modern HDLs

```
// Generalized FIR filter parameterized by the convolution coefficients
class FirFilter(bitWidth: Int, coeffs: Seq[UInt]) extends Module {

    val zs = Reg(Vec(coeffs.length, UInt(bitWidth.W)))
    zs(0) := io.in
    for (i <- 1 until coeffs.length) {
        zs(i) := zs(i-1)
    }

    // Do the multiplies
    val products = VecInit.tabulate(coeffs.length)(i => zs(i) * coeffs(i))

}
```

Modern HDLs

```
// Generalized FIR filter parameterized by the convolution coefficients
class FirFilter(bitWidth: Int, coeffs: Seq[UInt]) extends Module {

    val zs = Reg(Vec(coeffs.length, UInt(bitWidth.W)))
    zs(0) := io.in
    for (i <- 1 until coeffs.length) {
        zs(i) := zs(i-1)
    }

    // Do the multiplies
    val products = VecInit.tabulate(coeffs.length)(i => zs(i) * coeffs(i))
}

}
```

High level constructs!

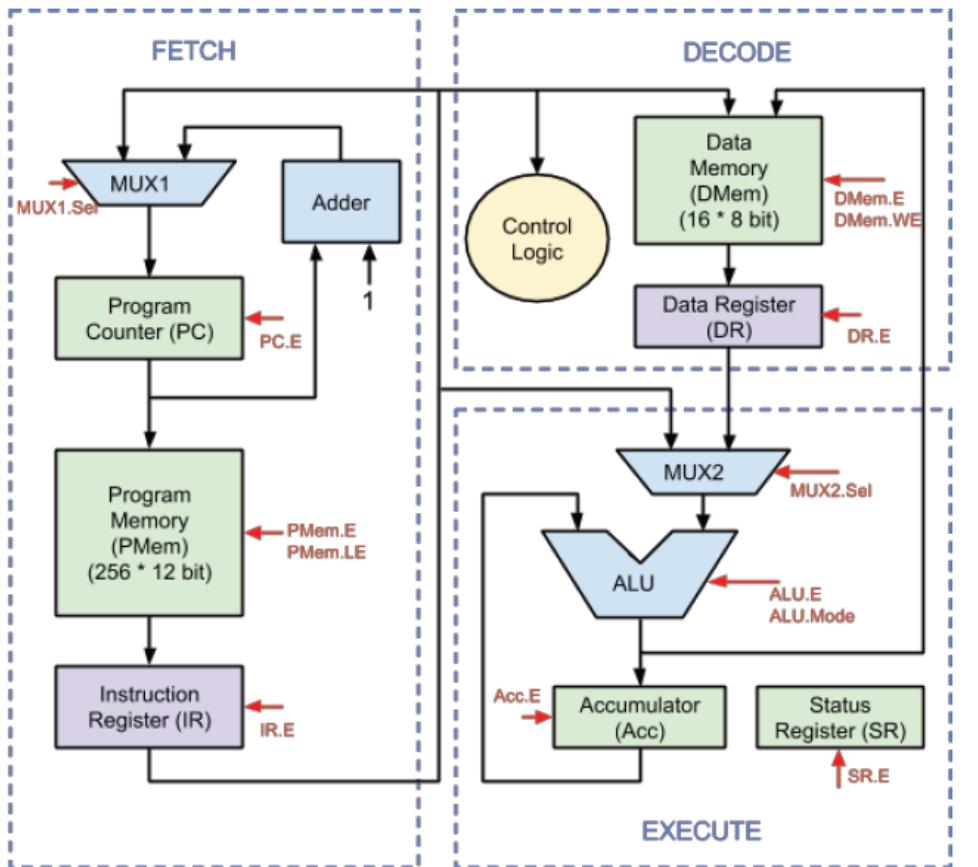
Modern HDLs

```
// Generalized FIR filter parameterized by the convolution coefficients
class FirFilter(bitWidth: Int, coeffs: Seq[UInt]) extends Module {
    val io = IO(new Bundle {
        val in = Input(UInt(bitWidth.W))
        val out = Output(UInt(bitWidth.W))
    })
    // Create the serial-in, parallel-out shift register
    val zs = Reg(Vec(coeffs.length, UInt(bitWidth.W)))
    zs(0) := io.in
    for (i <- 1 until coeffs.length) {
        zs(i) := zs(i-1)
    }
    // Do the multiplies
    val products = VecInit.tabulate(coeffs.length)(i => zs(i) * coeffs(i))

    // Sum up the products
    io.out := products.reduce(_ + _)
}
```

High level constructs!

... **elaborated** into circuits



Hardware Design

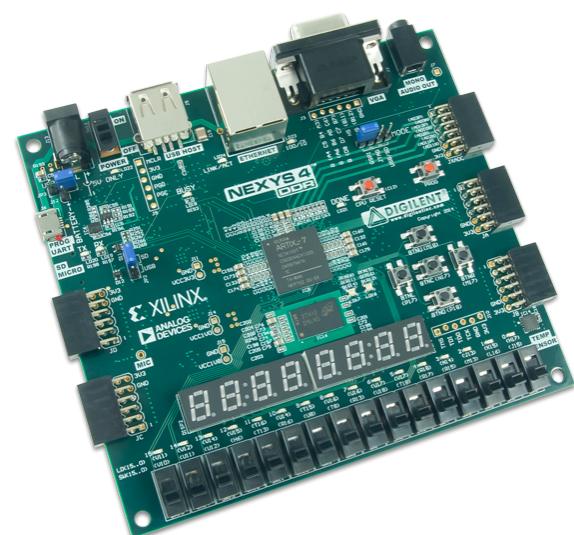


Traditional HDLs

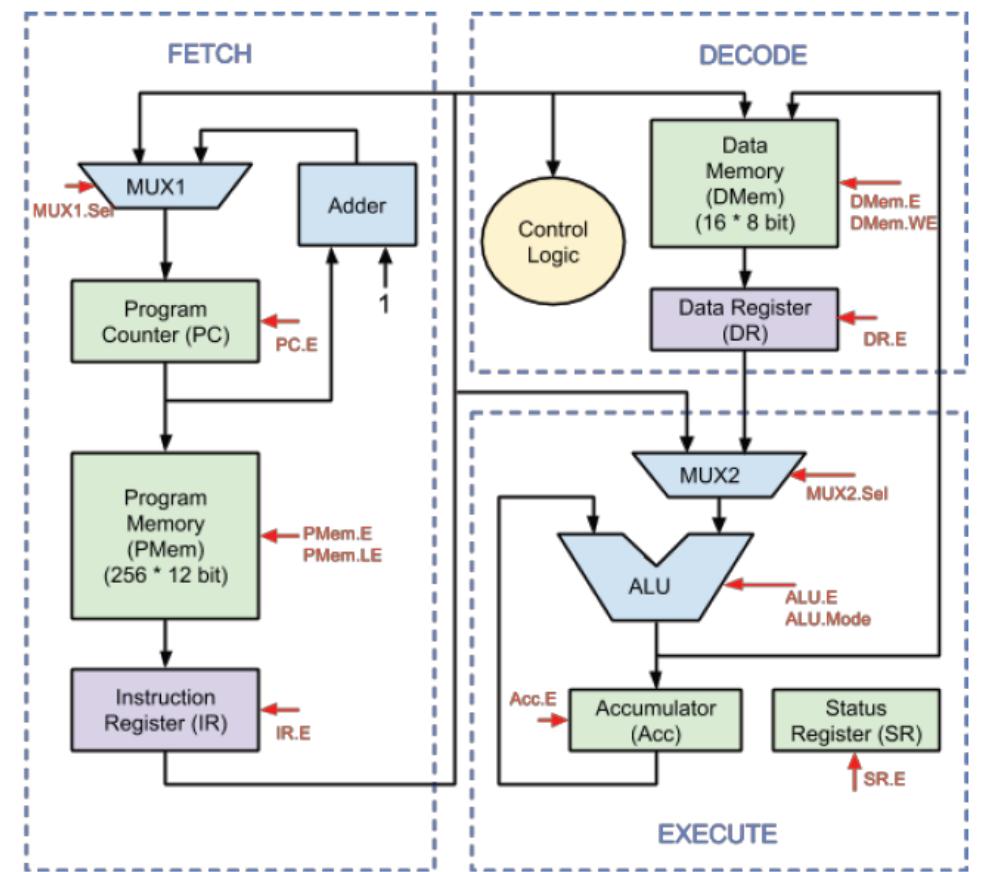
bluespec
CHISEL



Modern HDLs



FPGA



Hardware Design



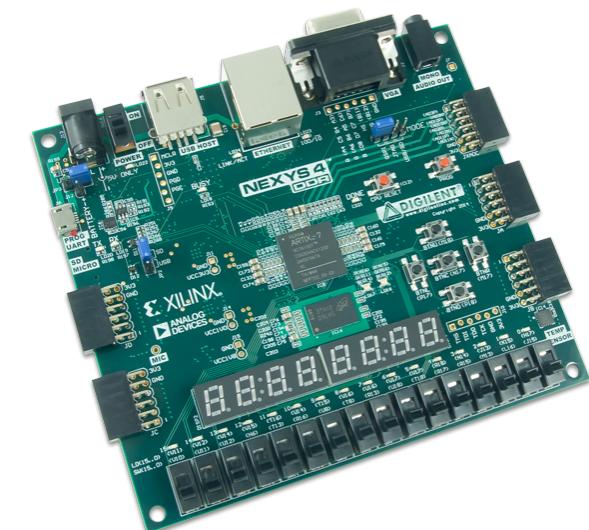
High-Level
Synthesis



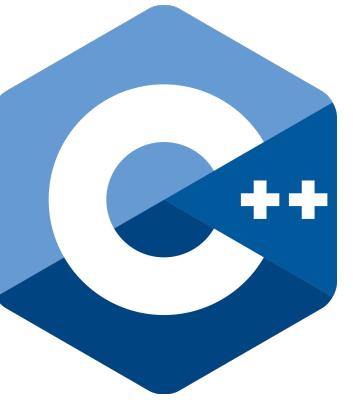
Traditional HDLs



Modern HDLs



FPGA



High-Level Synthesis

```
int m1[512], m2[512], sum;  
for (int i = 0; i < 512; i++) {  
    sum += m1[i] * m2[i]  
}
```



High-Level Synthesis

```
int m1[512], m2[512], sum;  
for (int i = 0; i < 512; i++) {  
    sum += m1[i] * m2[i]  
}
```

Software programs that are
compiled to hardware designs



High-Level Synthesis

Level of abstraction

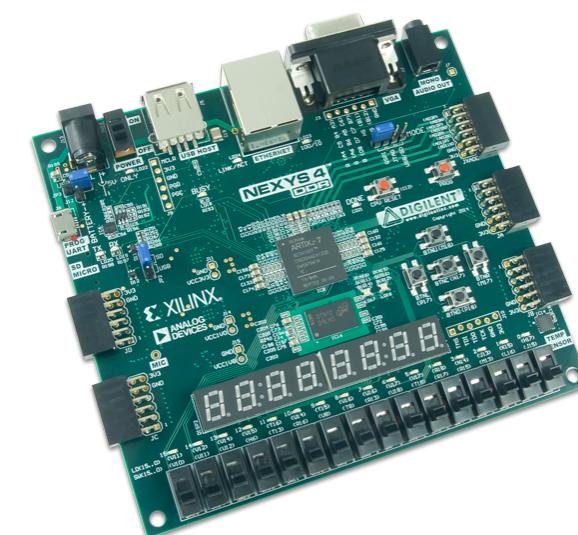
CHISEL
bluespec™

PyMTL

Modern HDLs

SystemVerilog

Traditional HDLs



FPGA



High-Level Synthesis

Level of abstraction

bluespec™
CHISEL

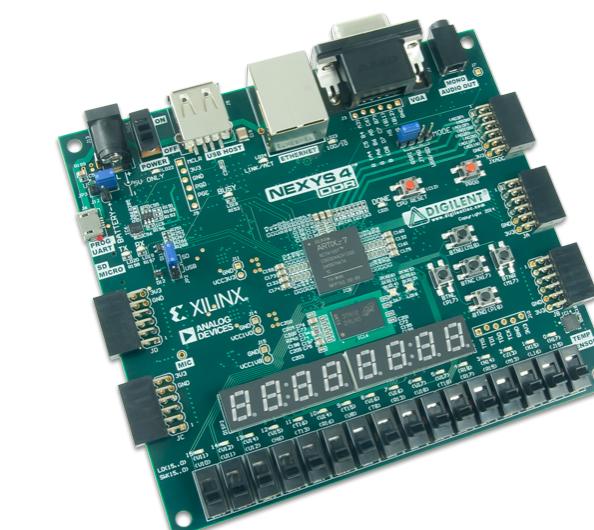
PyMTL

Modern HDLs

Circuit Specification

SystemVerilog

Traditional HDLs



FPGA



Computational Specification

High-Level Synthesis

Level of abstraction

CHISEL
bluespec™

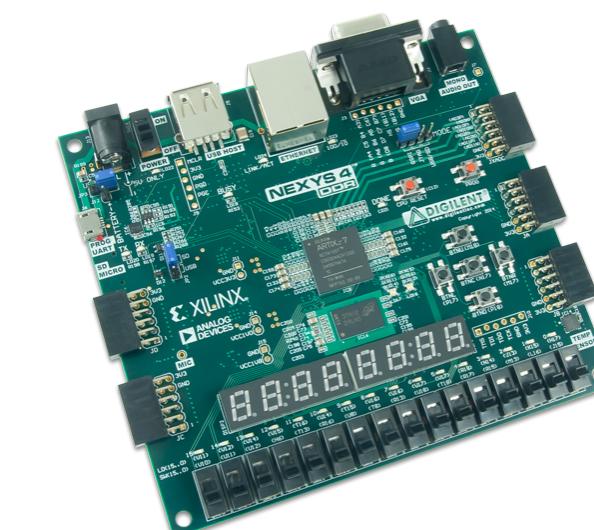
PyMTL

Modern HDLs

Circuit Specification

SystemVerilog

Traditional HDLs



FPGA

loops

method calls

arrays

break

malloc

functions

conditionals



loops

~~method calls~~

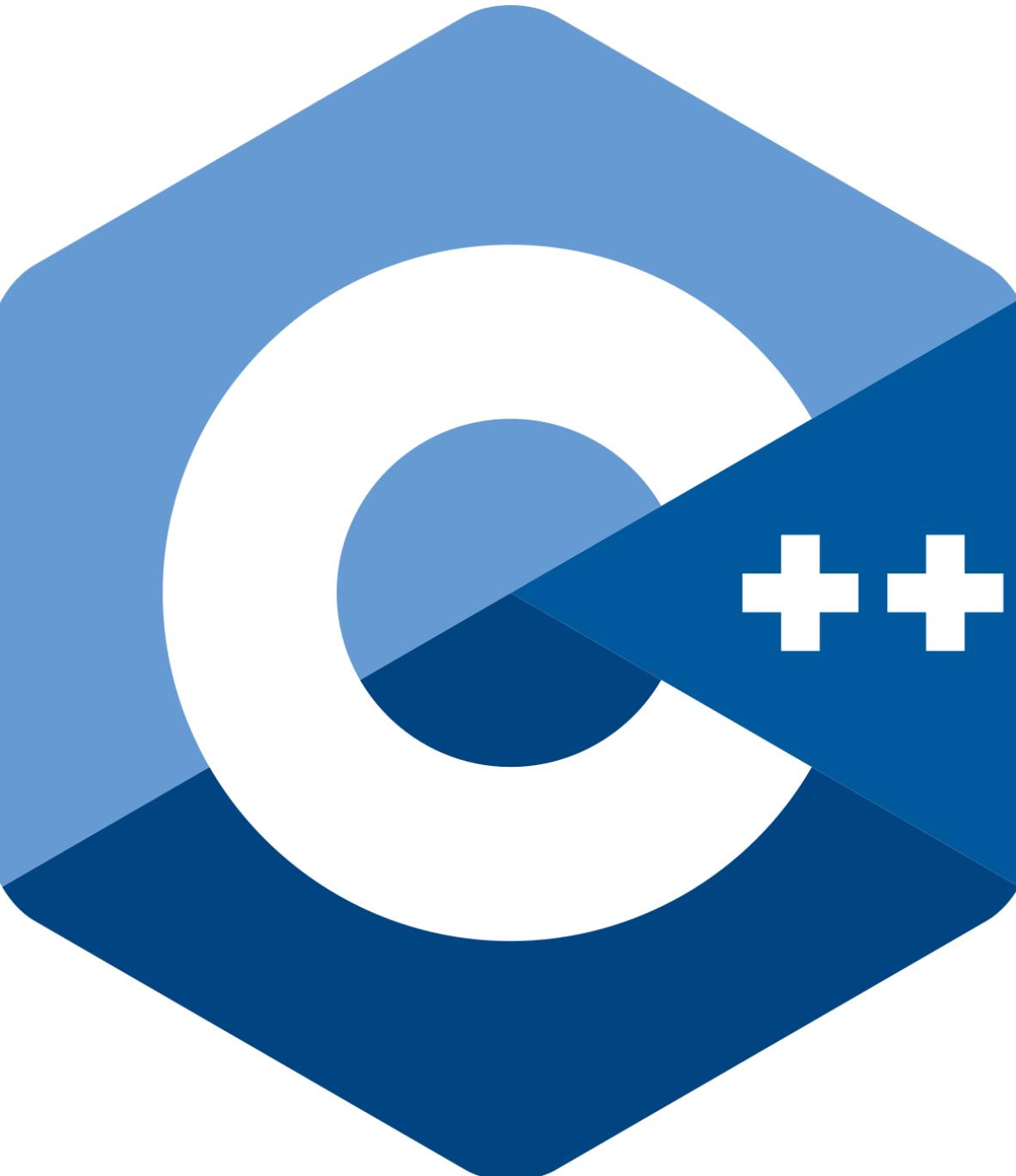
arrays

~~break~~

~~malloc~~

functions^{*}

conditionals



loops

~~method calls~~

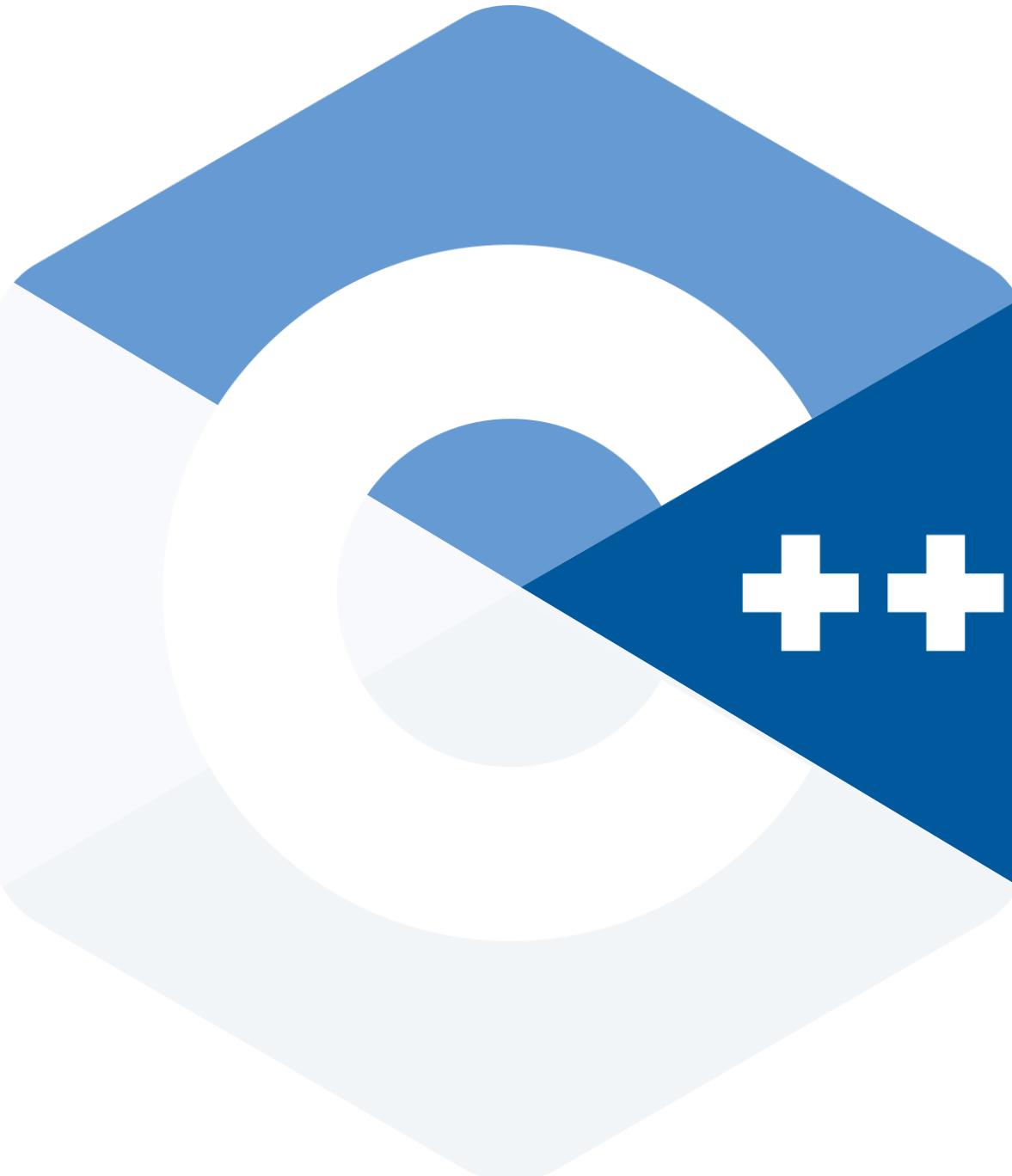
arrays

~~break~~

~~malloc~~

functions^{*}

conditionals



loops

~~method calls~~

arrays

~~break~~

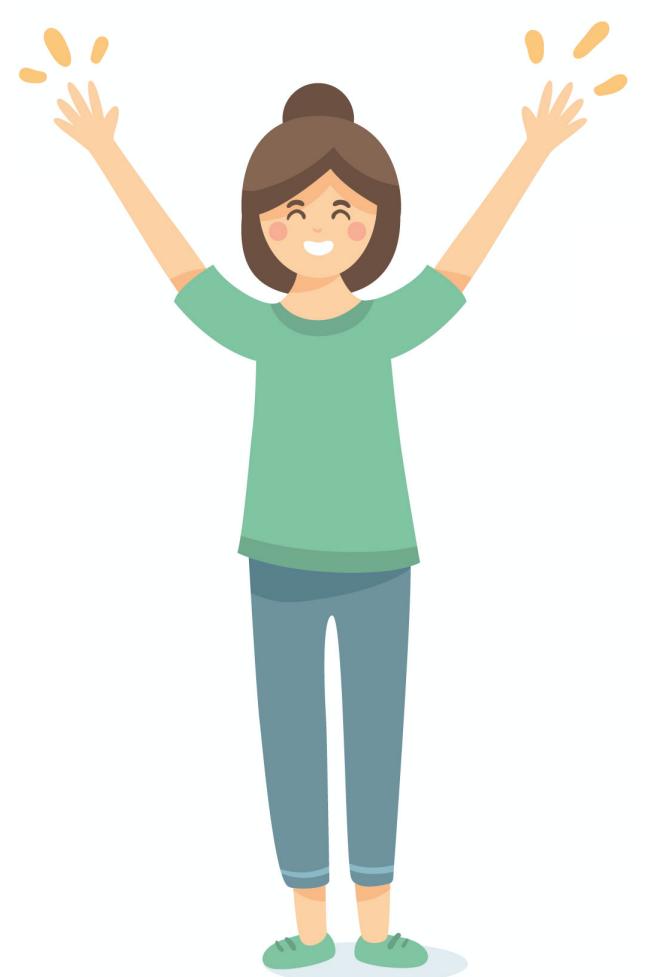
~~malloc~~

functions^{*}

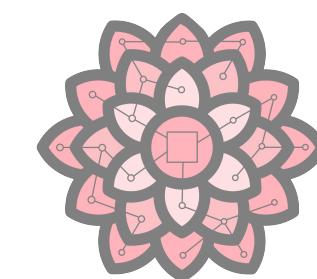
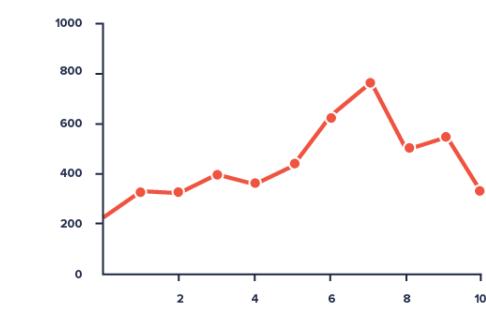
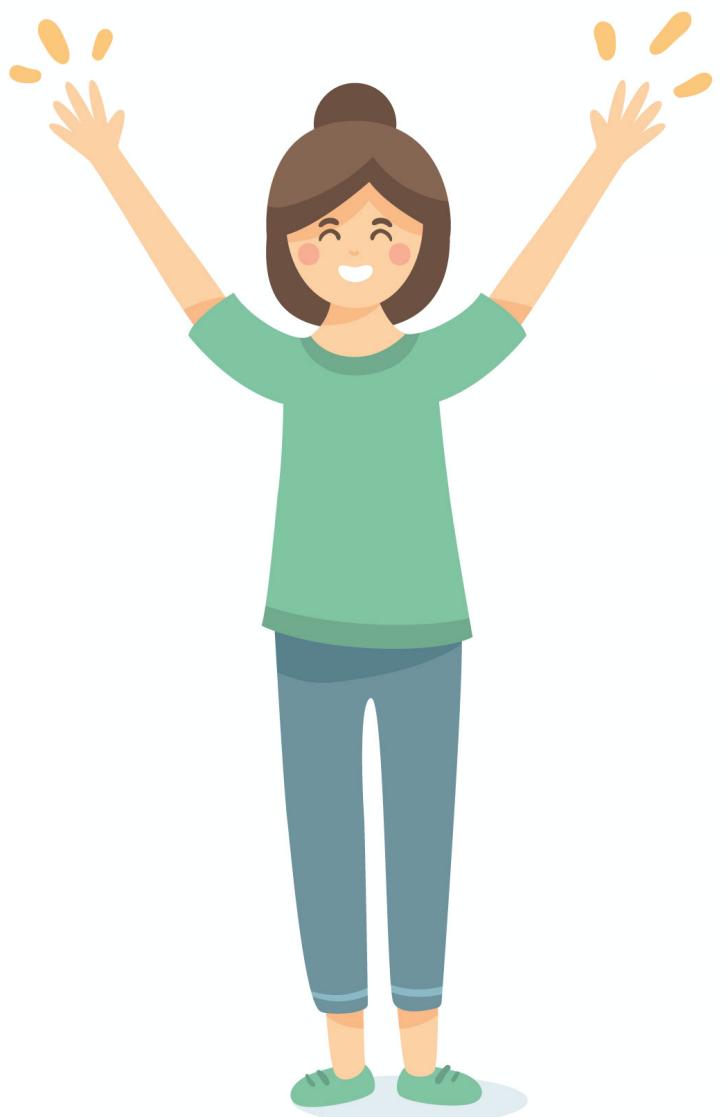
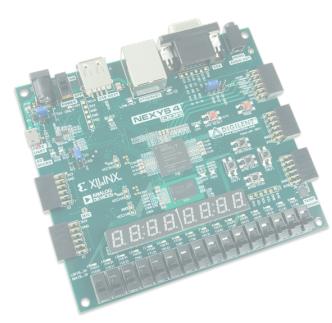
conditionals



**High-performance
hardware designs**



Ada's Journey



Super secret™ accelerator

Hardware

Super secret™ accelerator

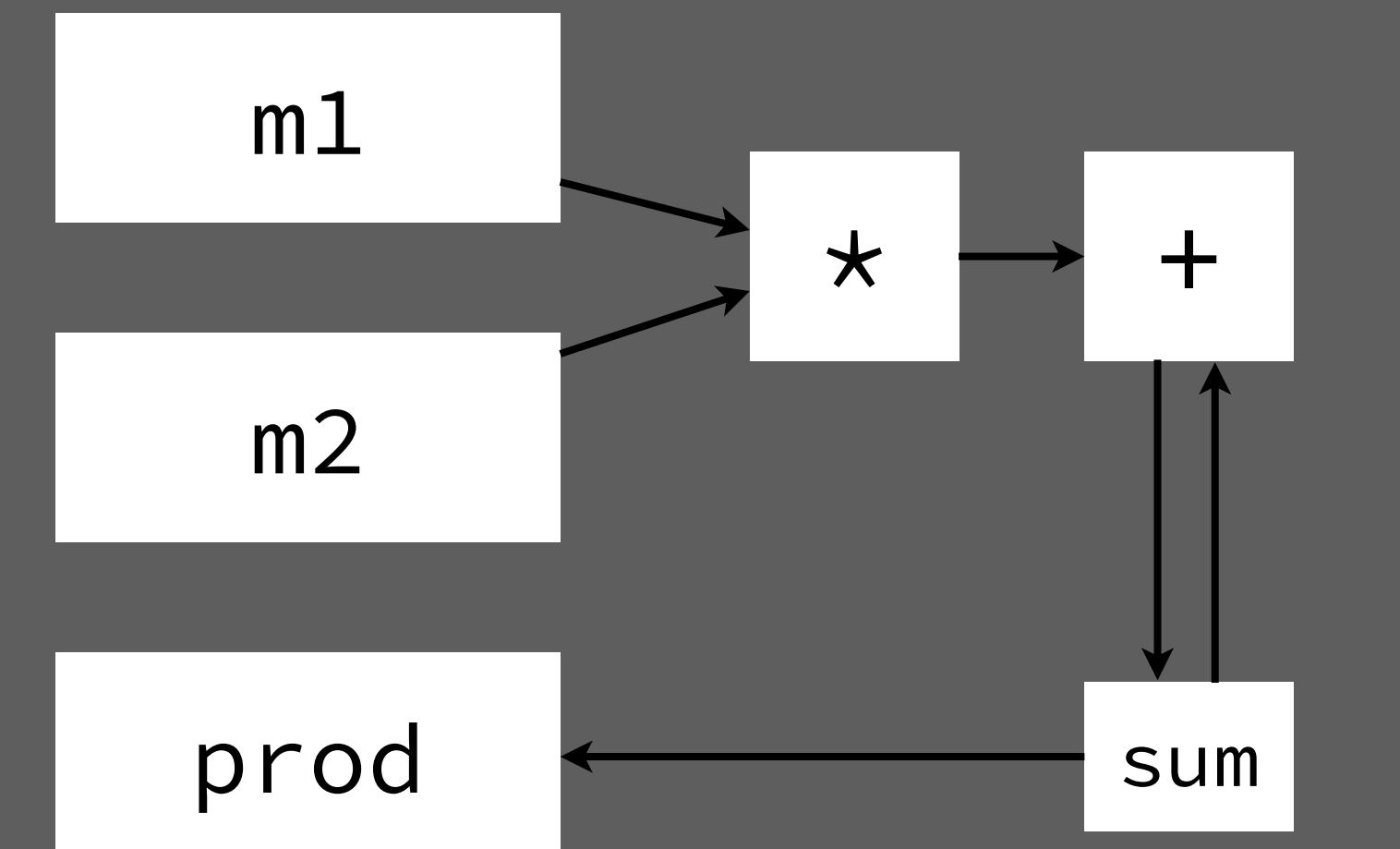
```
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

Hardware

Super secret™ accelerator

```
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

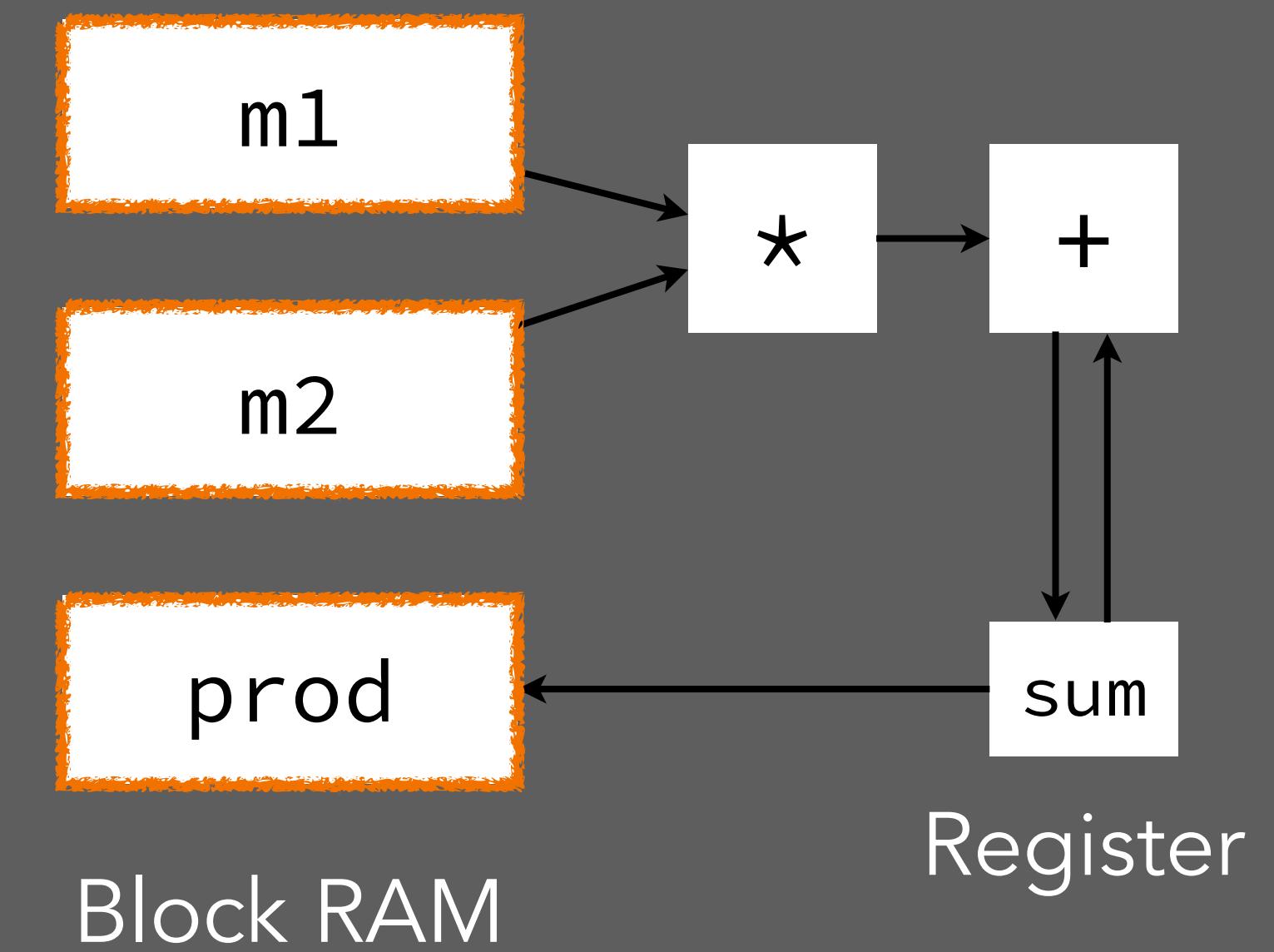
Hardware



Super secret™ accelerator

```
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

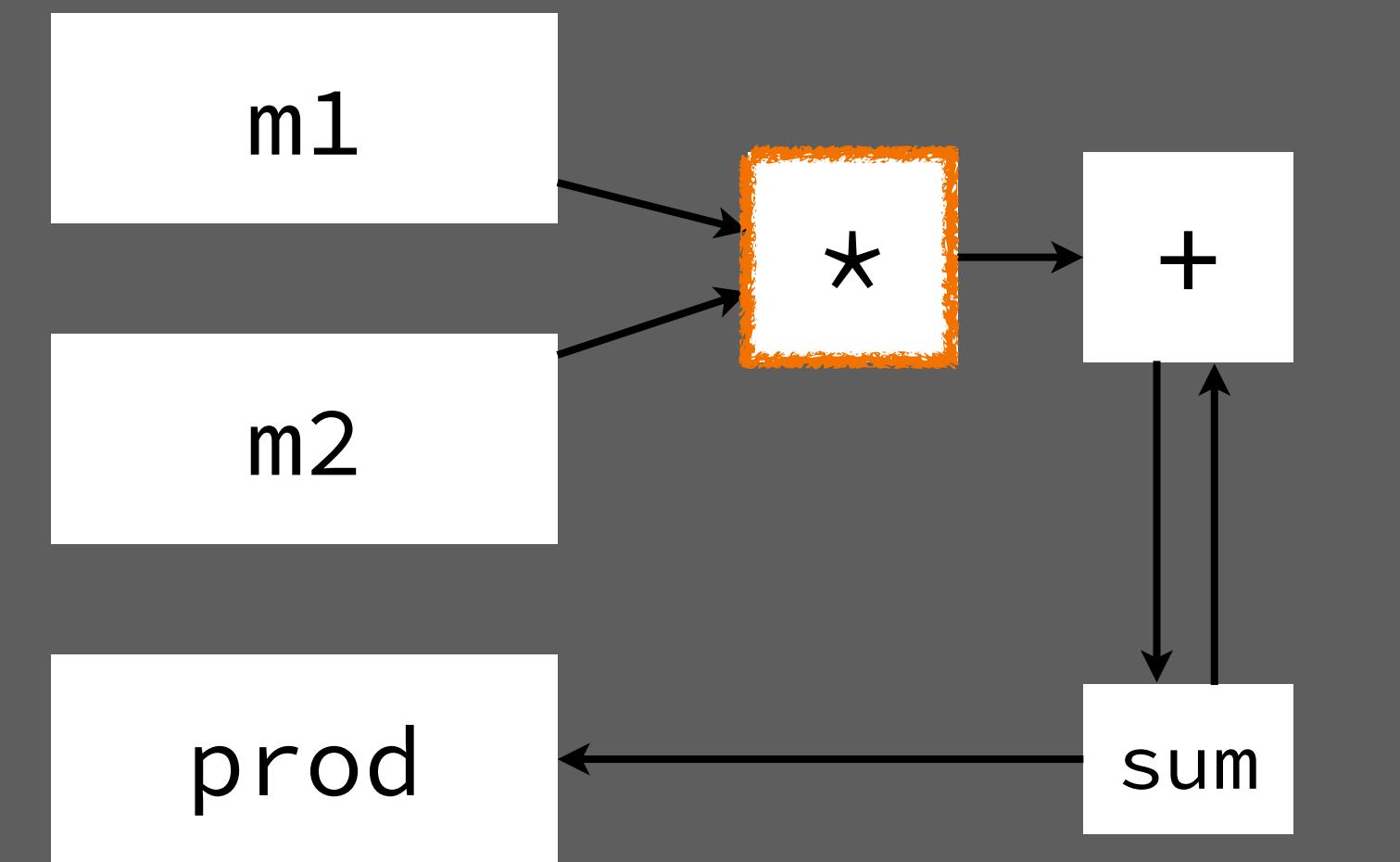
Hardware



Super secret™ accelerator

```
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

Hardware



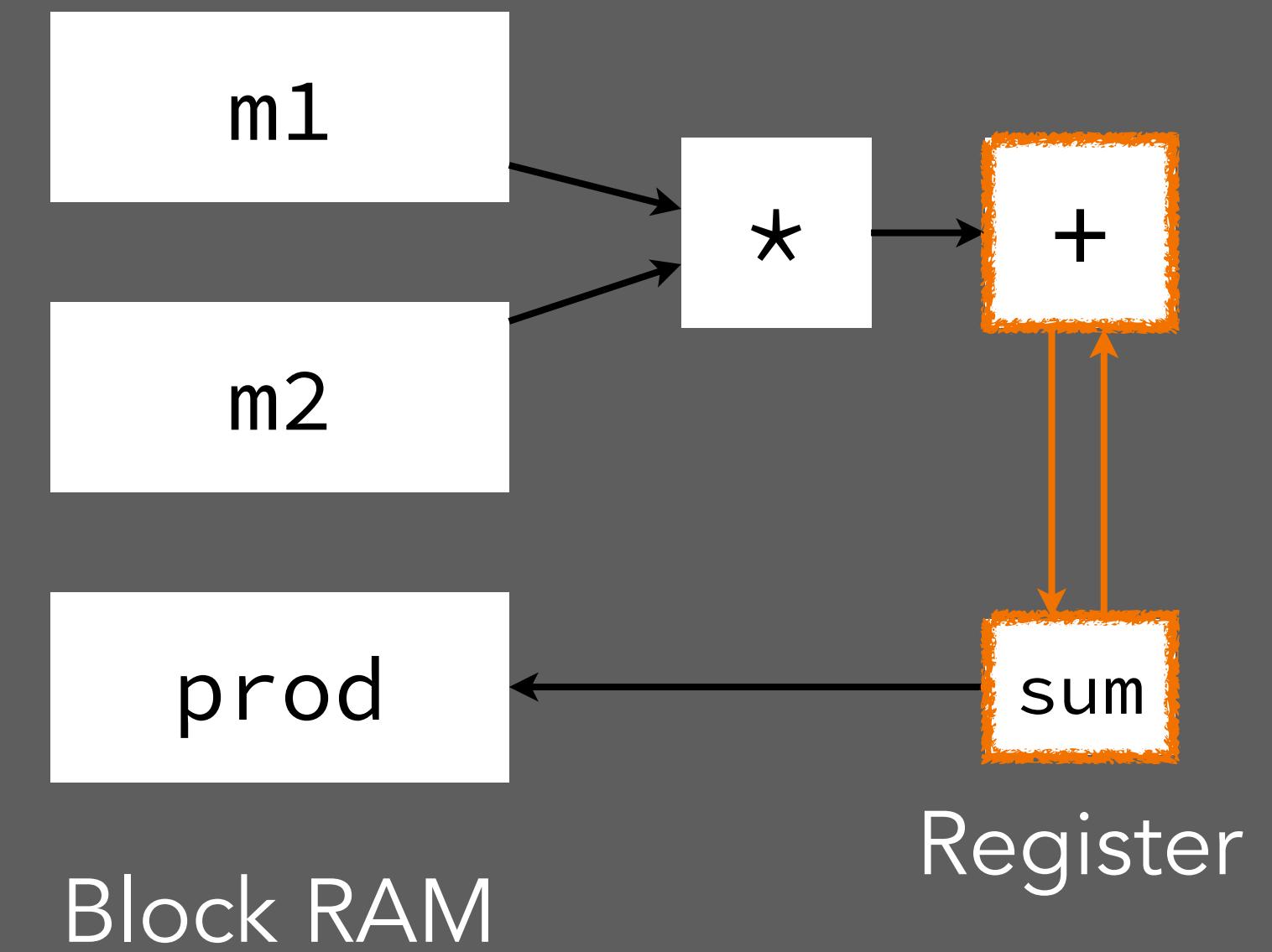
Block RAM

Register

Super secret™ accelerator

```
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

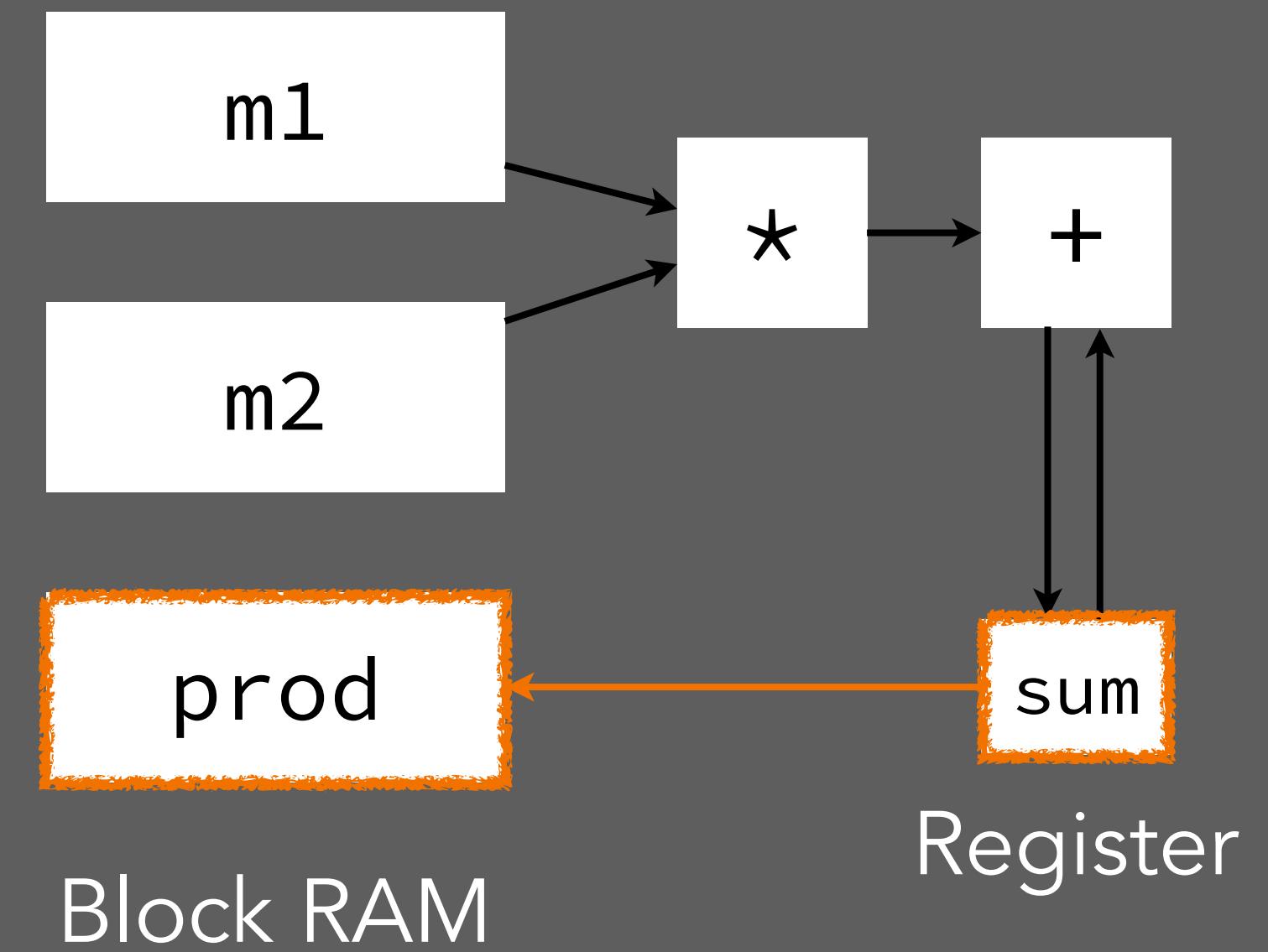
Hardware



Super secret™ accelerator

```
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

Hardware





reaching up to **35.9 effective TFLOPS** for a large GRU over hundreds of timesteps. This represents an approximate **two orders of magnitude** advantage over the Titan Xp. This is

Darwin: A Genomics Co-processor Provides up to **15,000×** acceleration on long read assembly

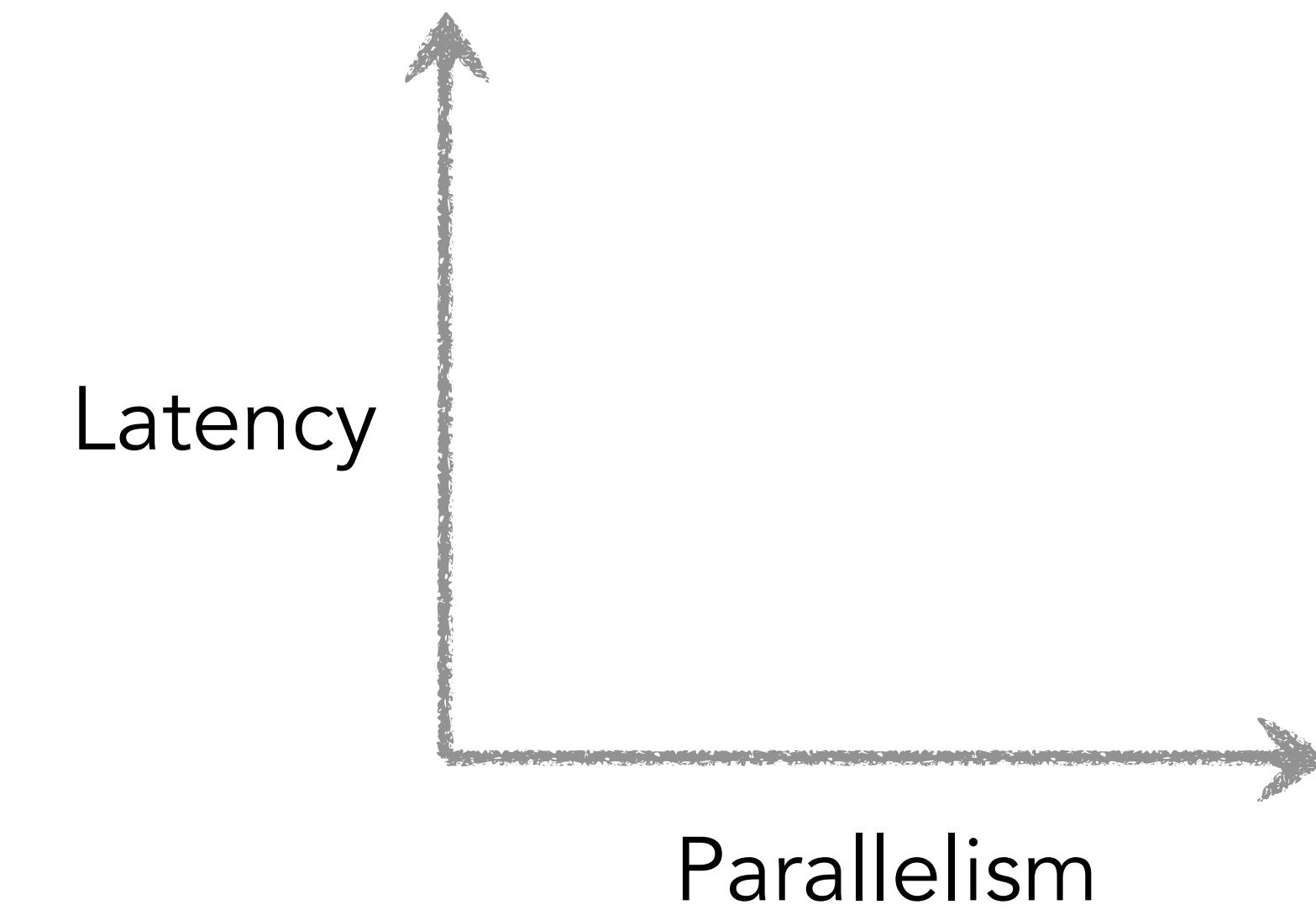
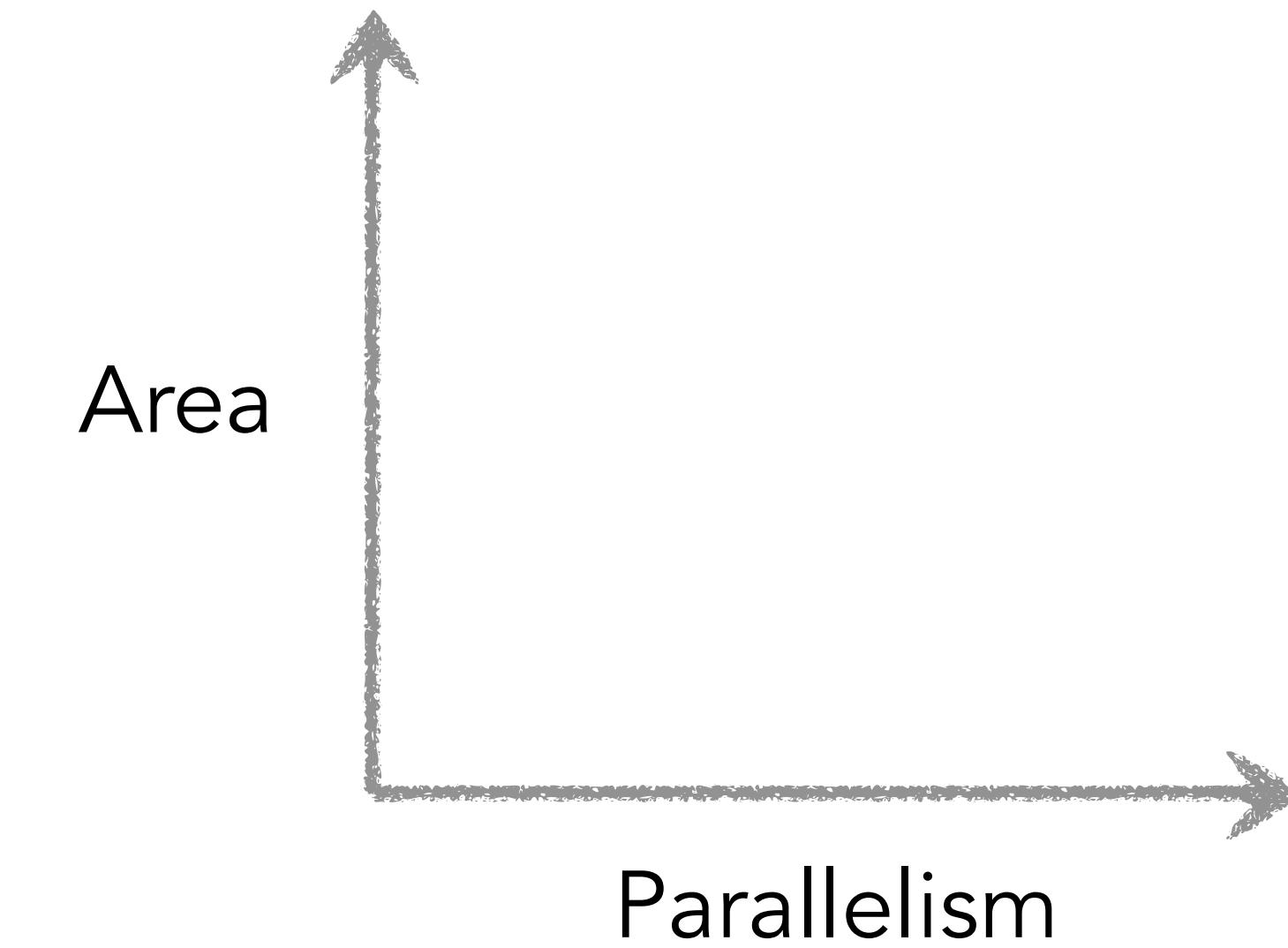
ranking candidate documents. Under high load, the large-scale reconfigurable fabric improves the ranking throughput of each server by a factor of 95% for a fixed latency distribution—or, while maintaining equivalent throughput, reduces the tail latency by 29%.



Super secret™
accelerator



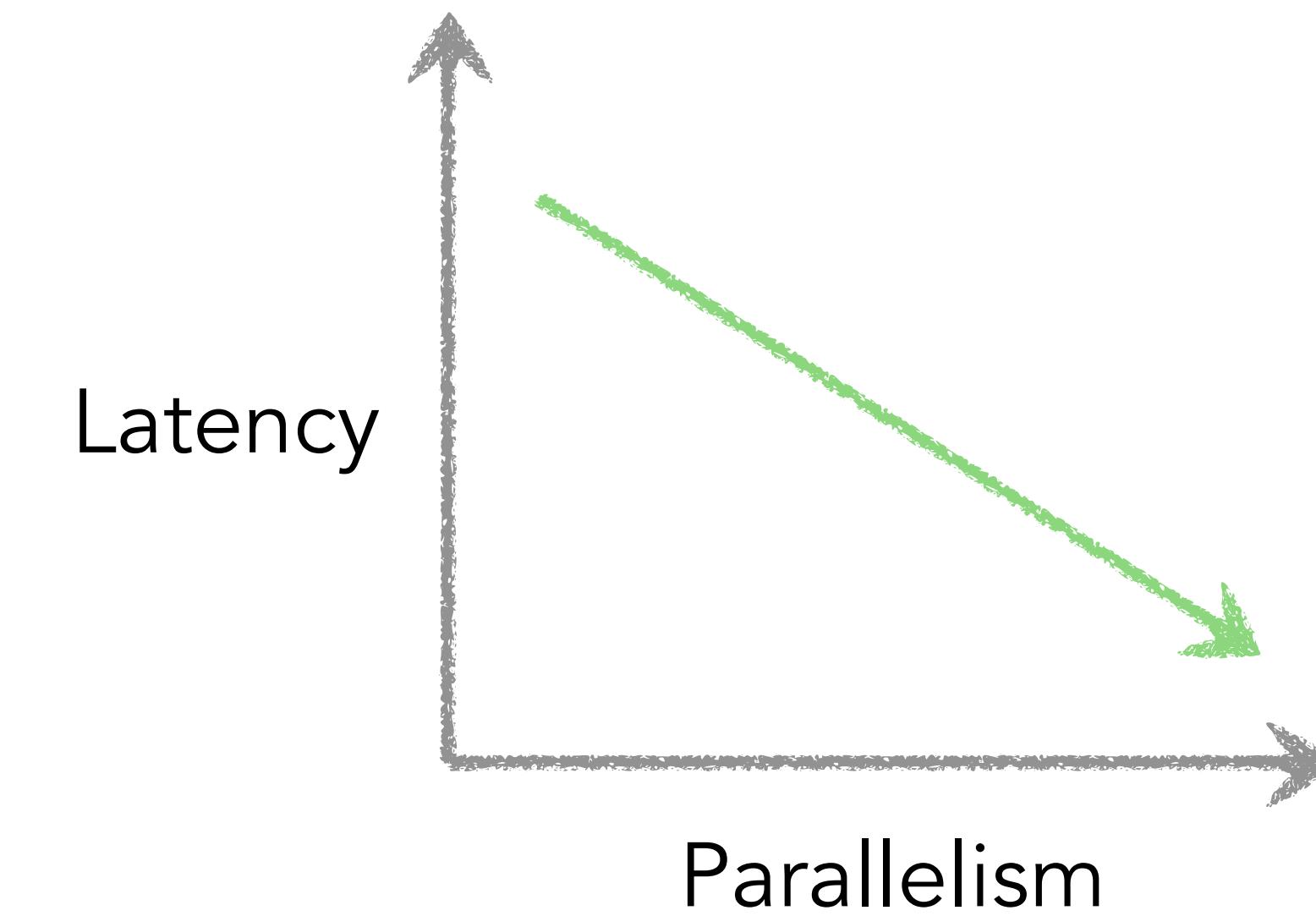
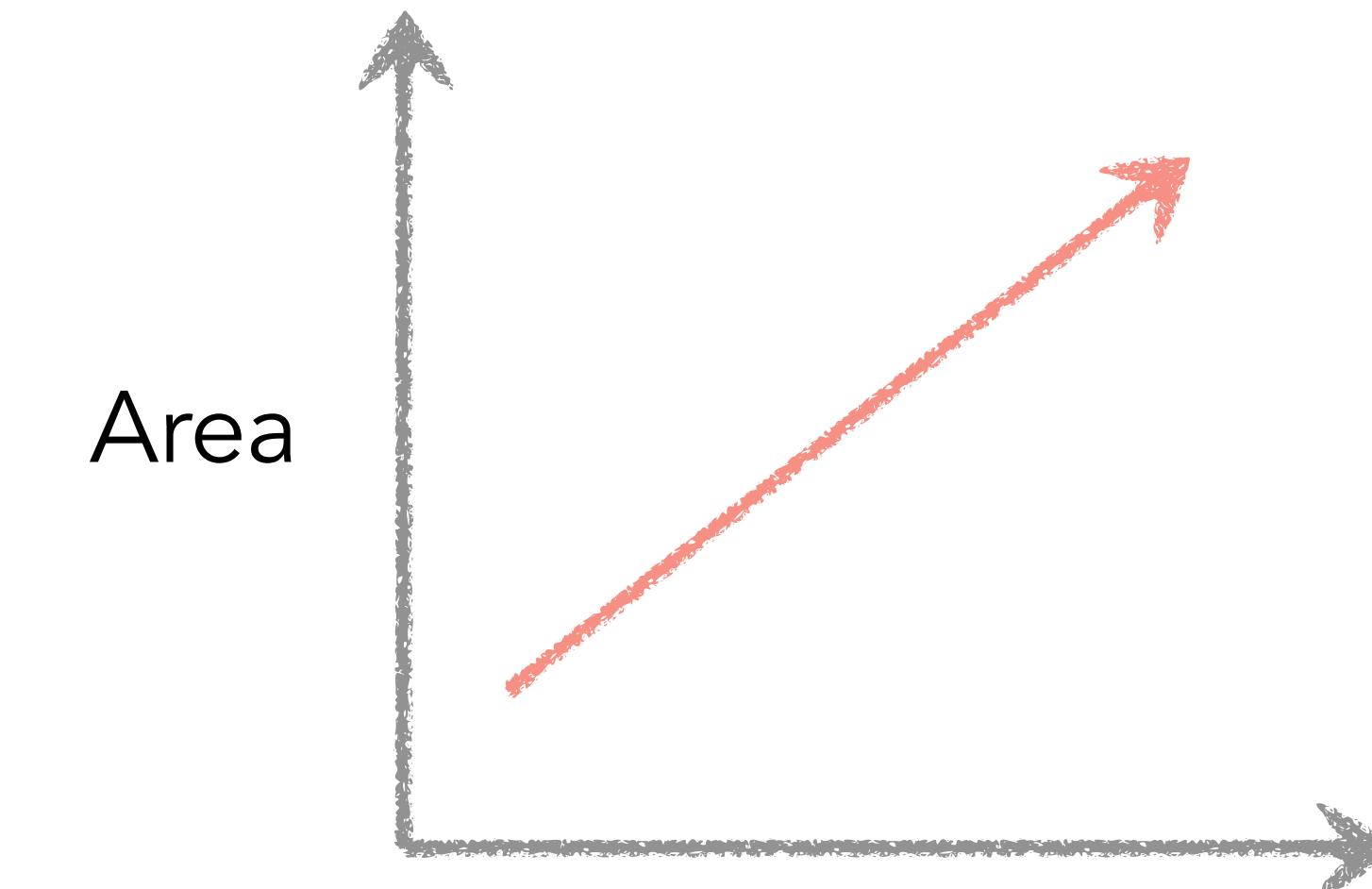
Resources on the FPGA

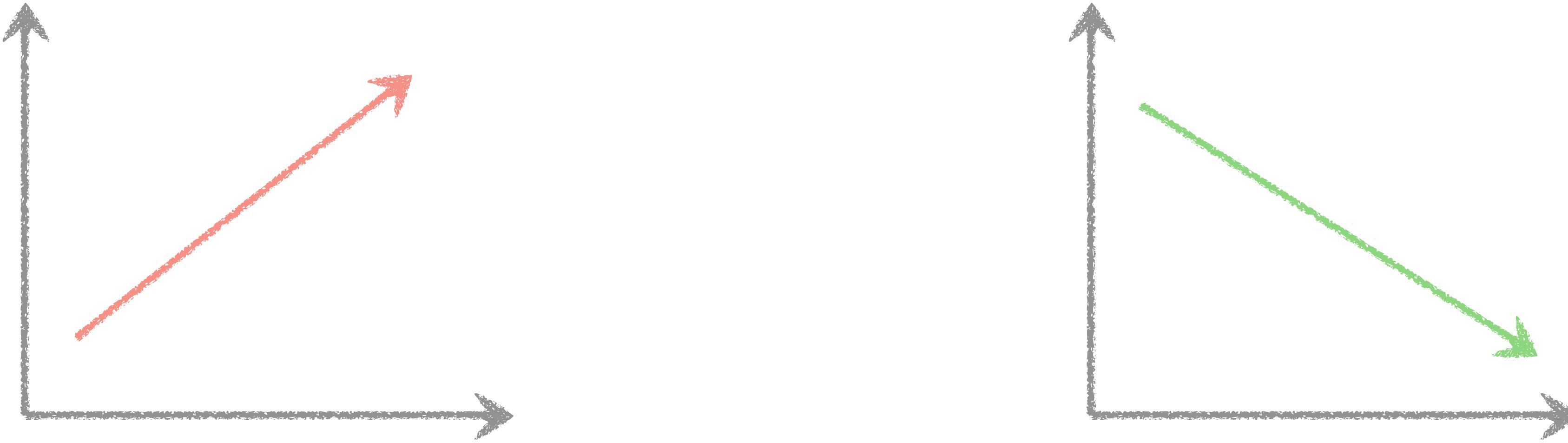


Super secret™
accelerator



Resources on the FPGA





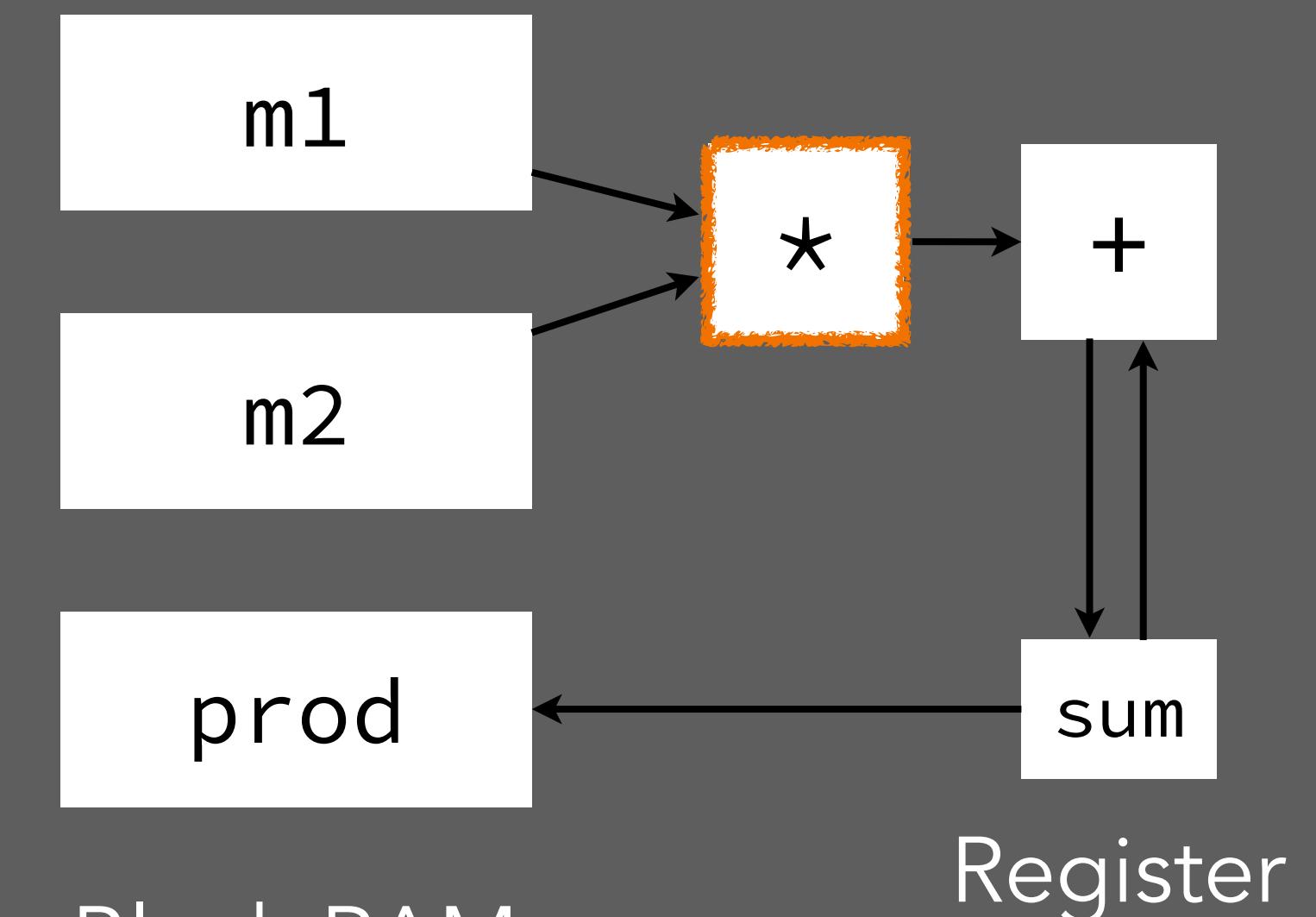
Area-Latency trade-offs

Create more processing elements to extract parallelism

Super secret™ accelerator

```
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

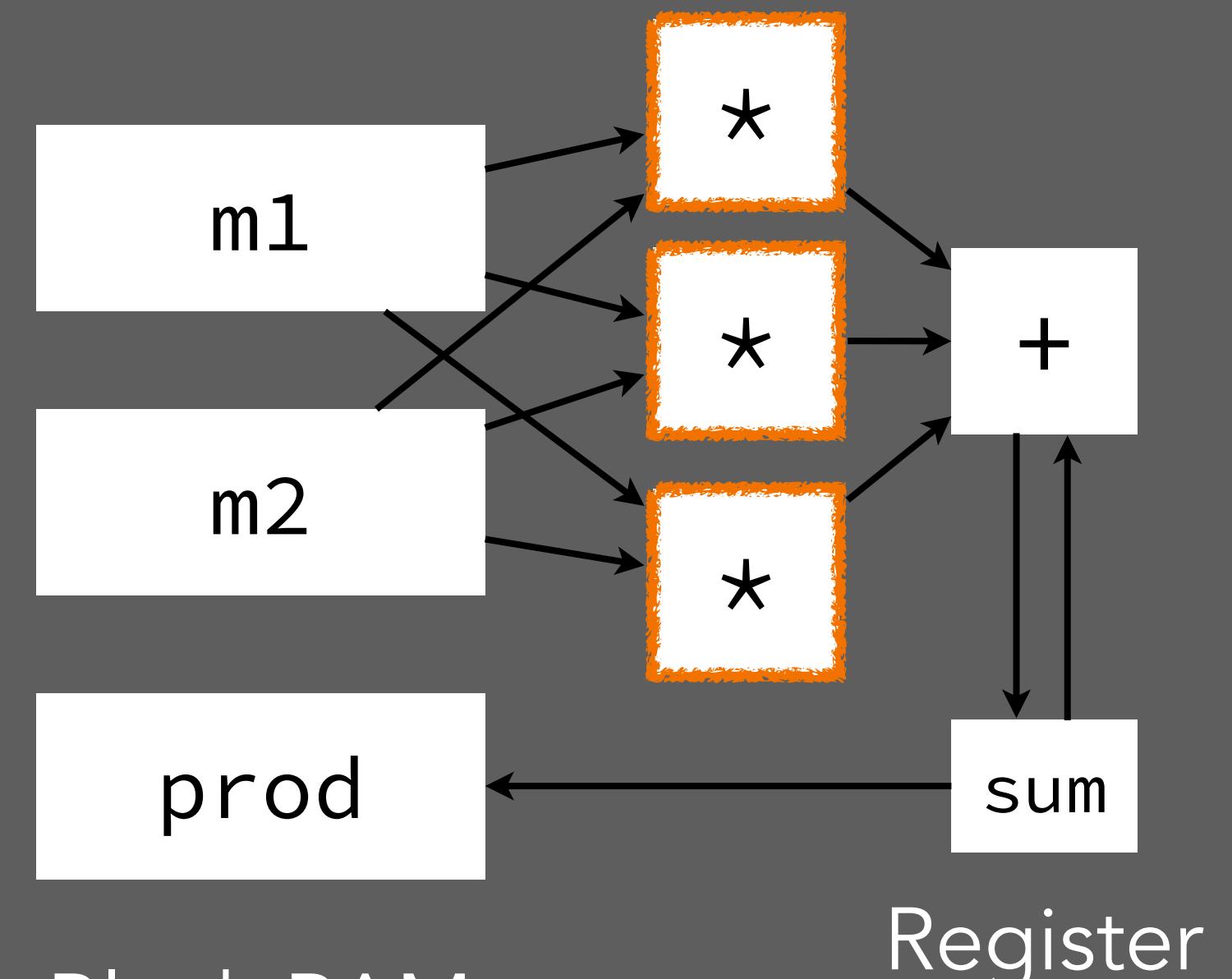
Hardware



Super secret™ accelerator

```
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

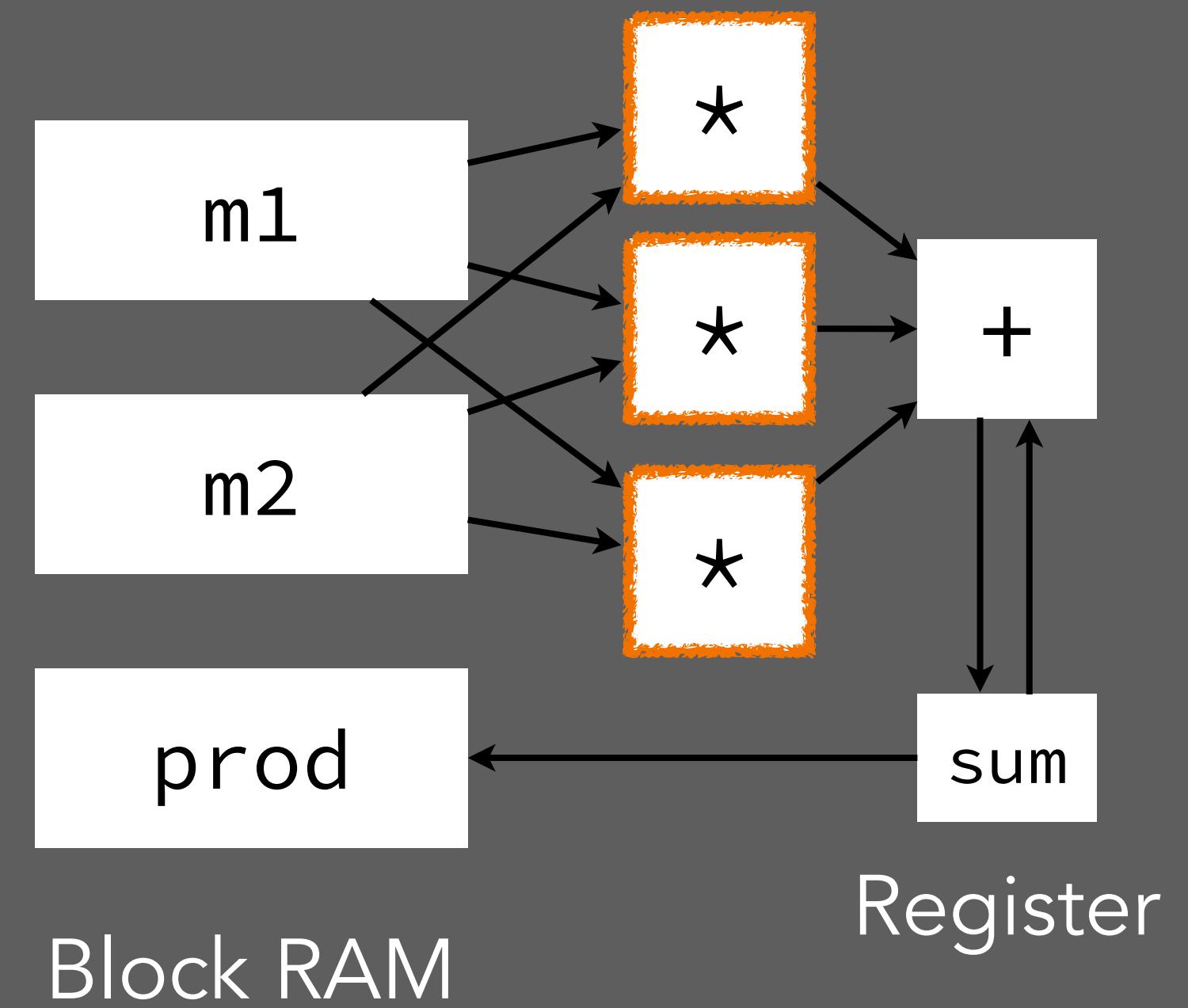
Hardware

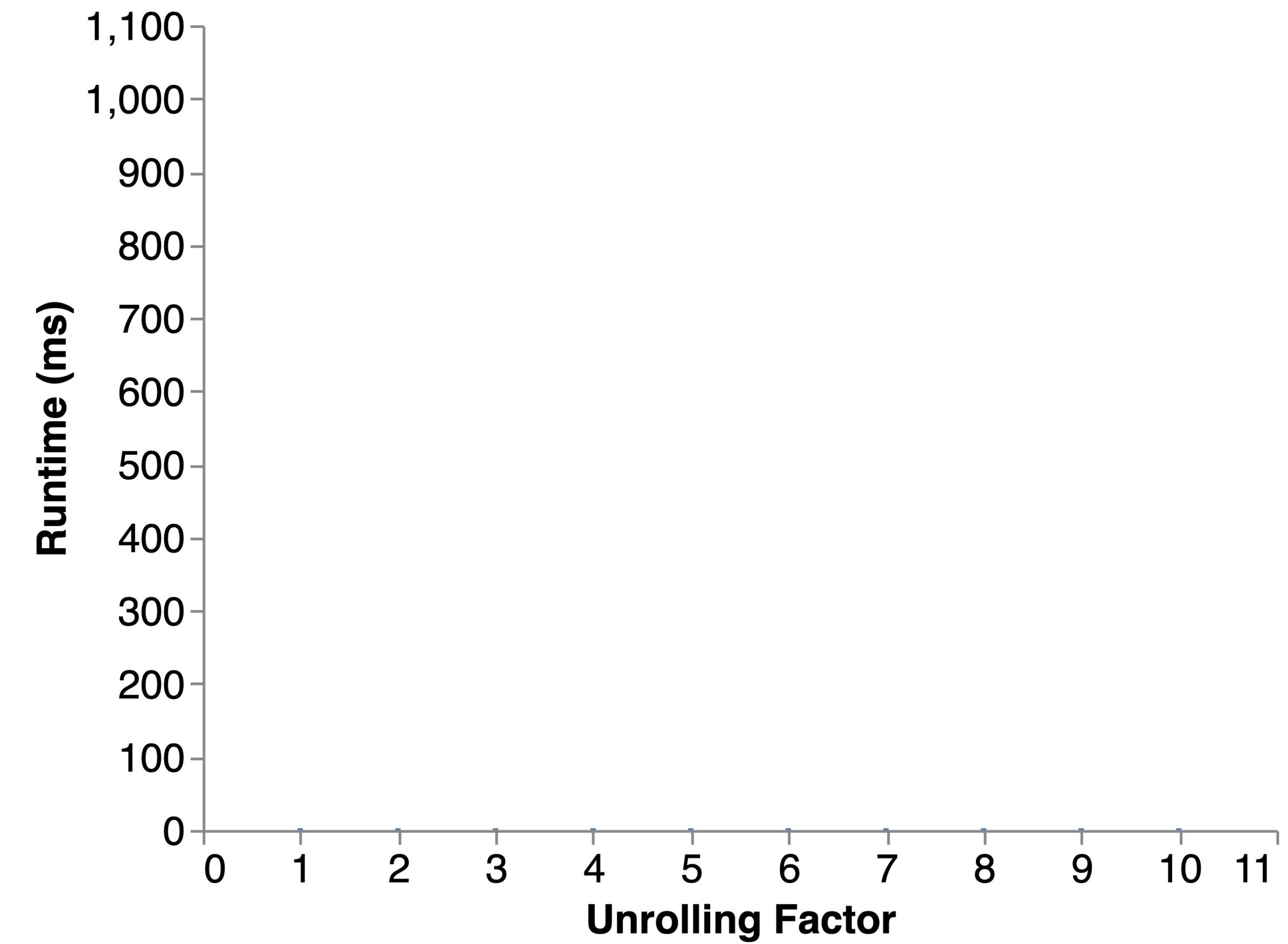
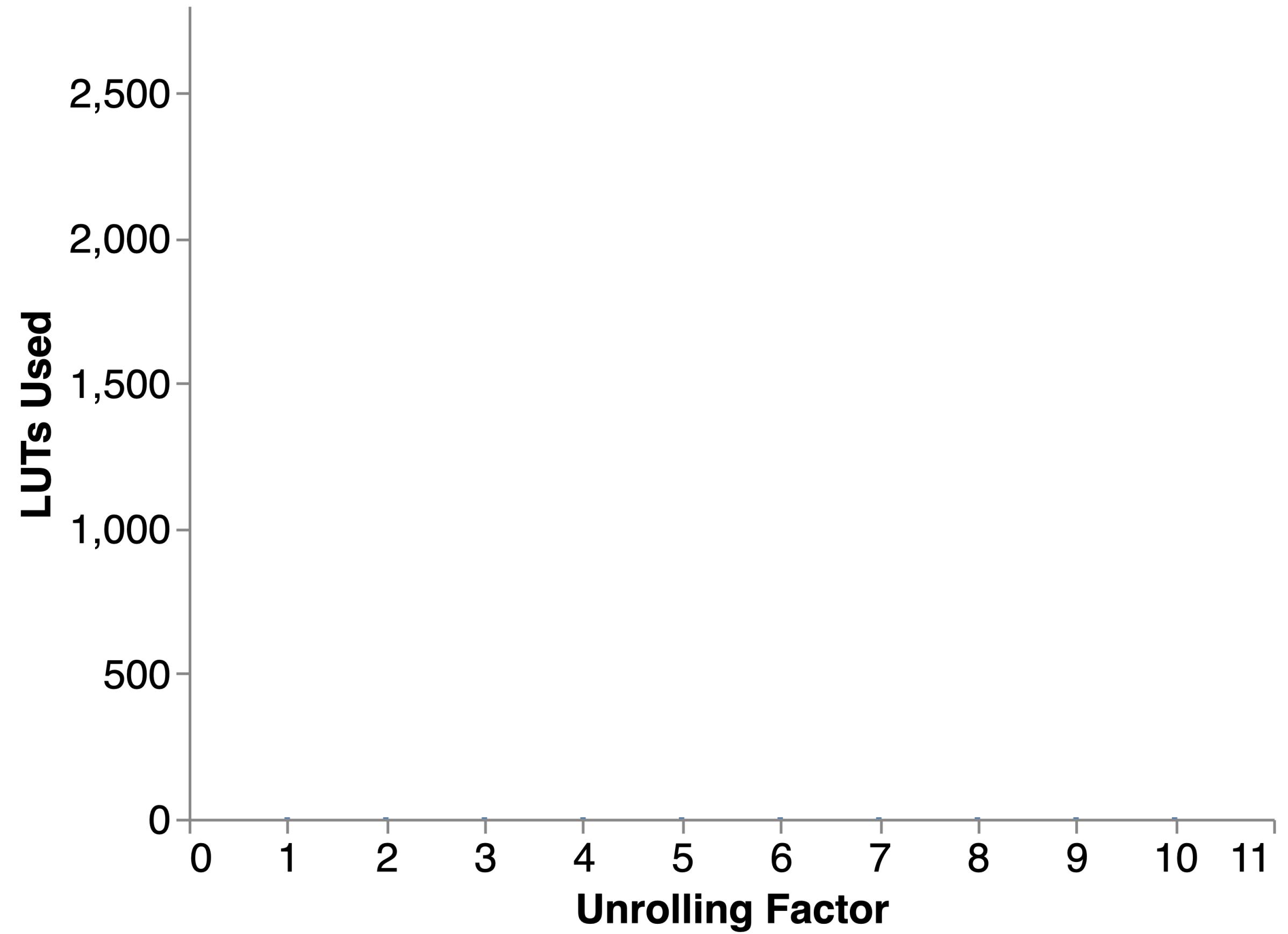


Super secret™ accelerator

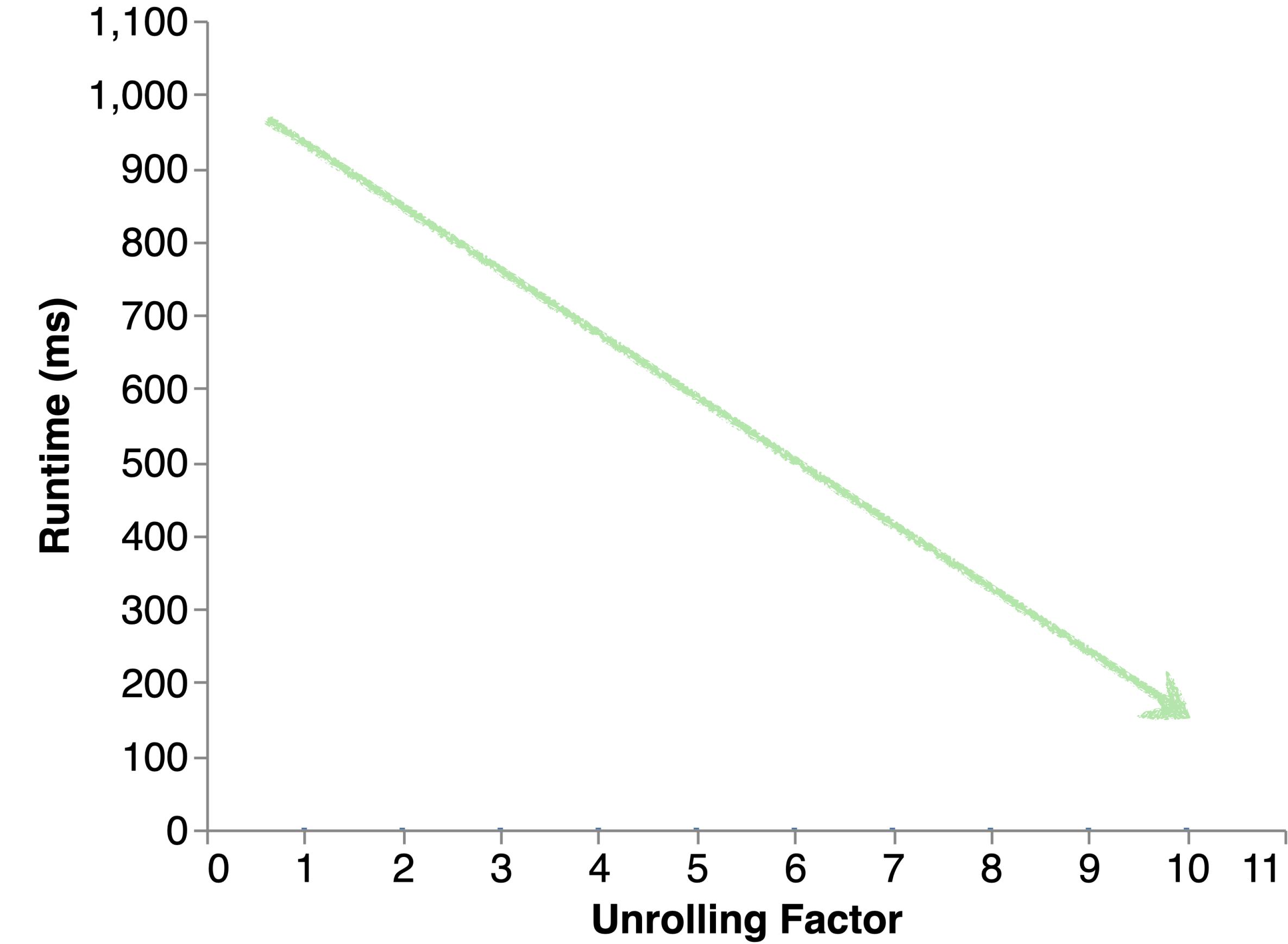
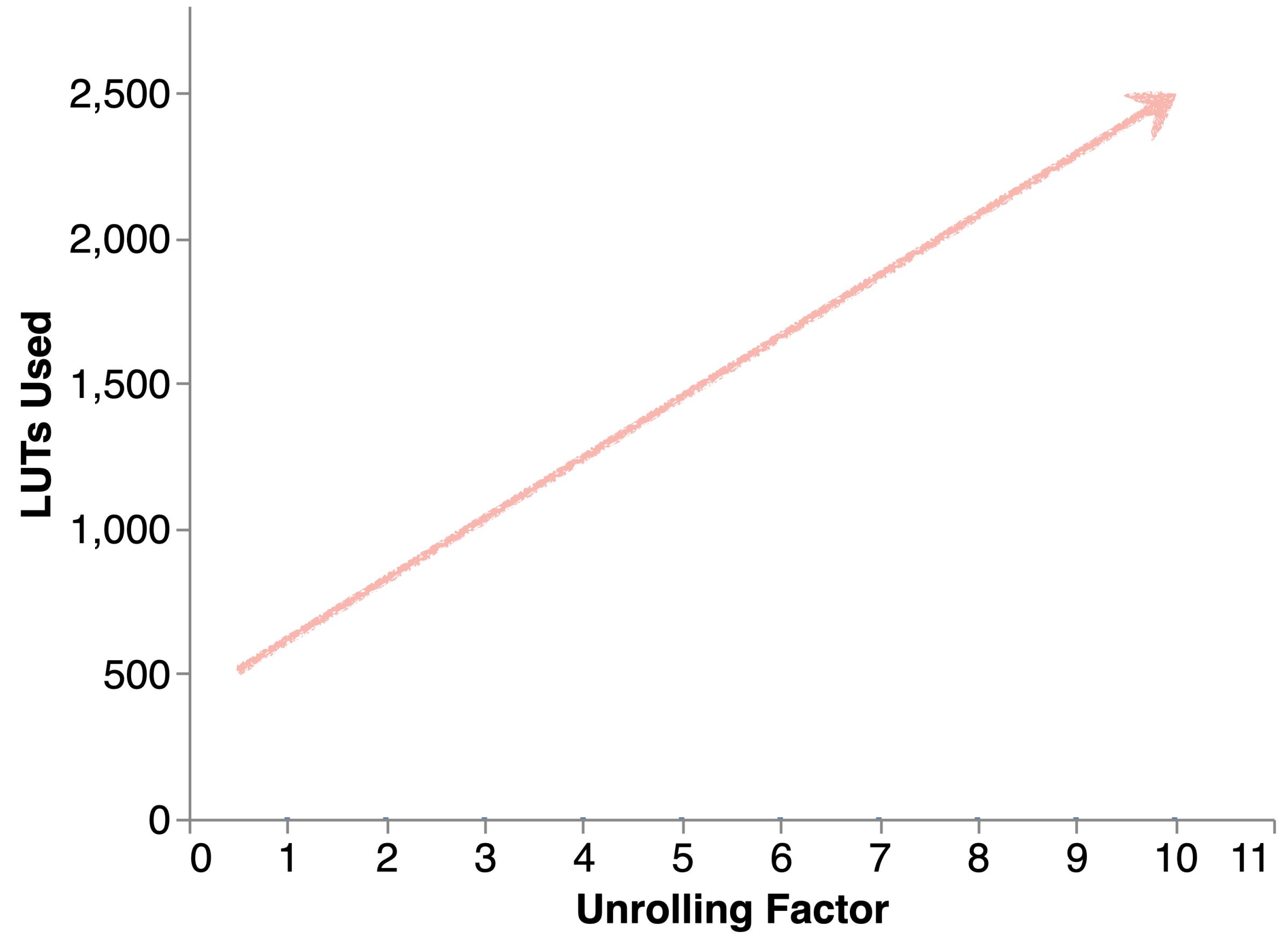
```
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            #pragma HLS UNROLL factor=3
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

Hardware

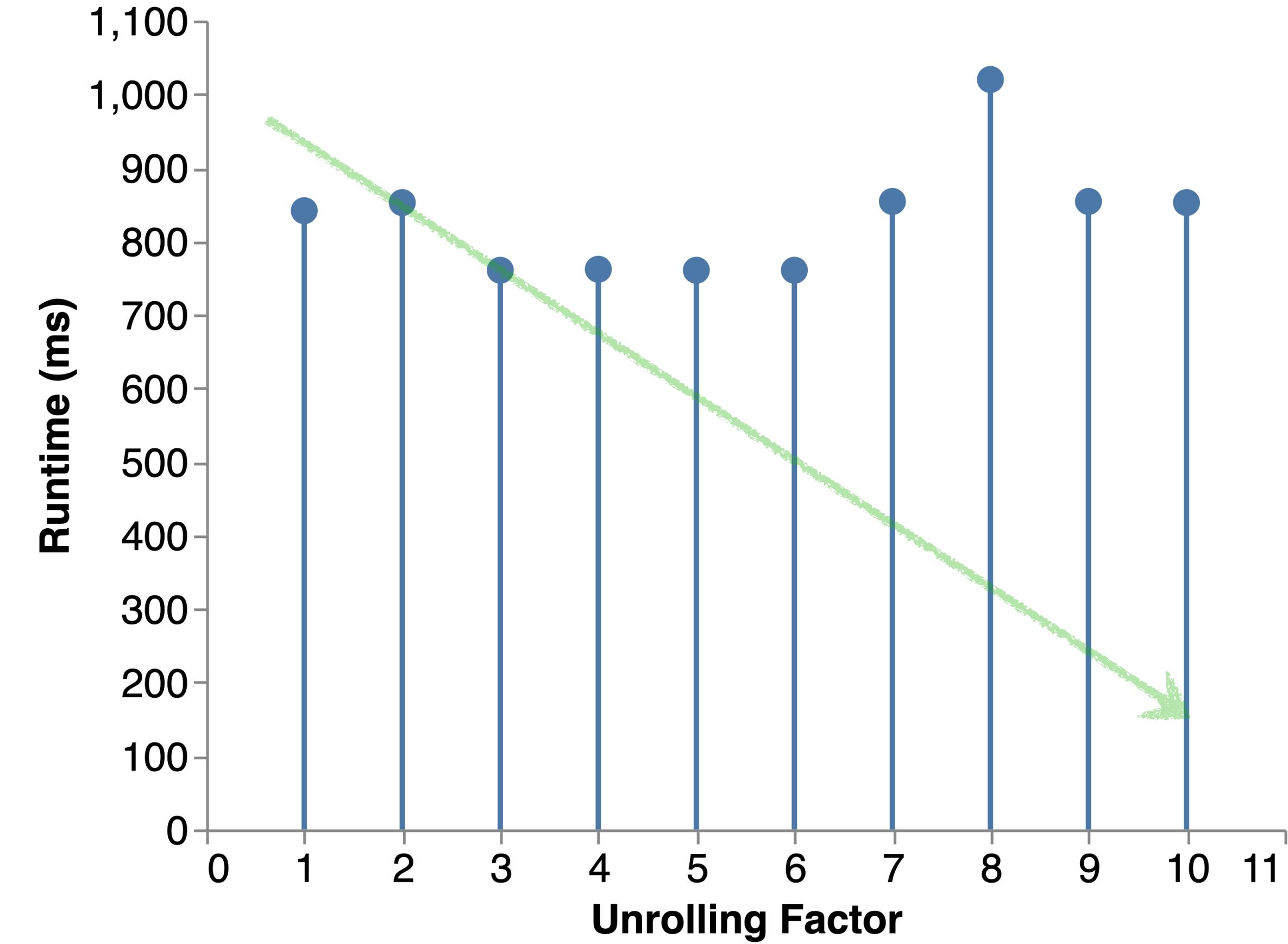
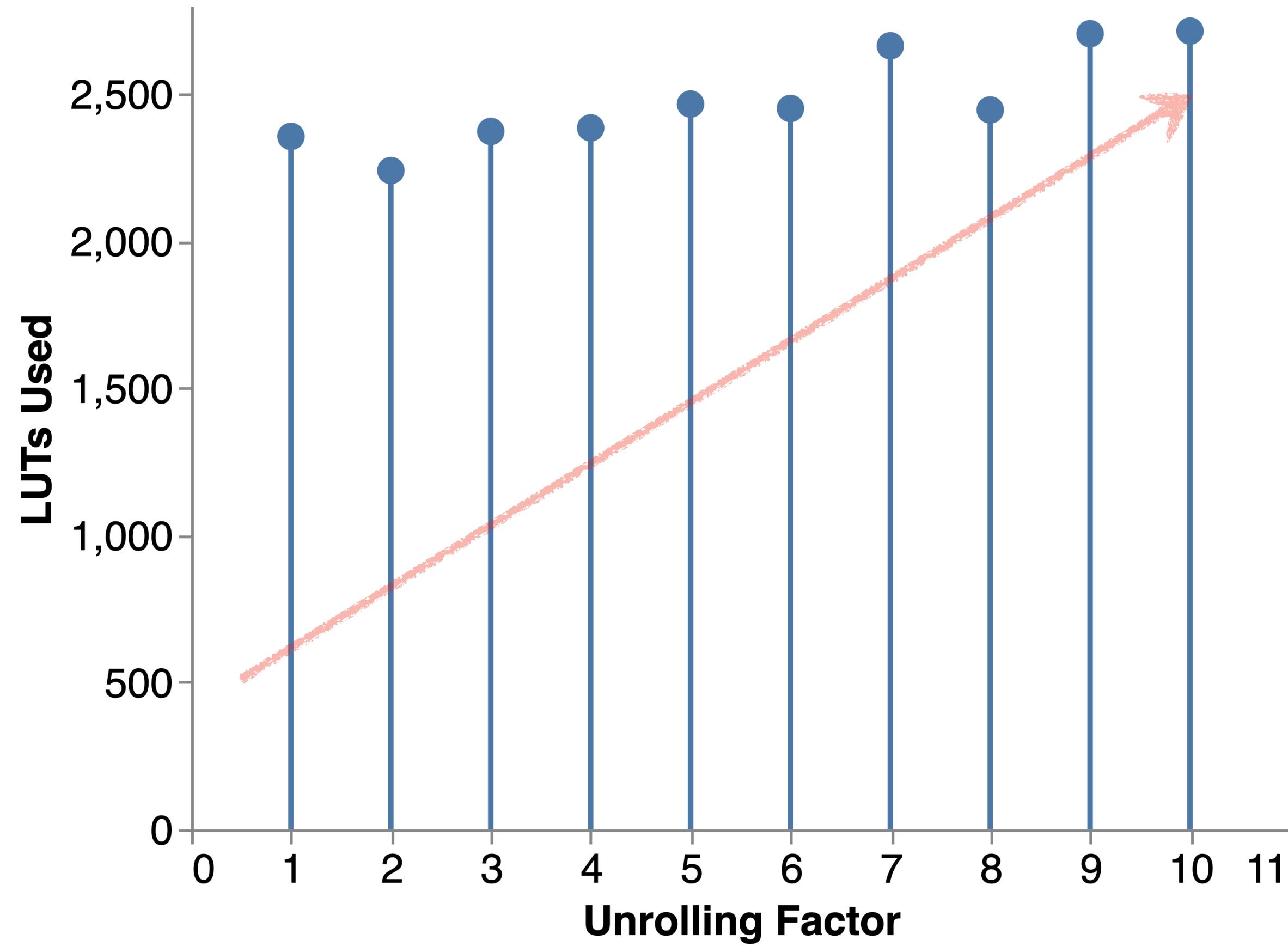




Area-Latency trade-offs



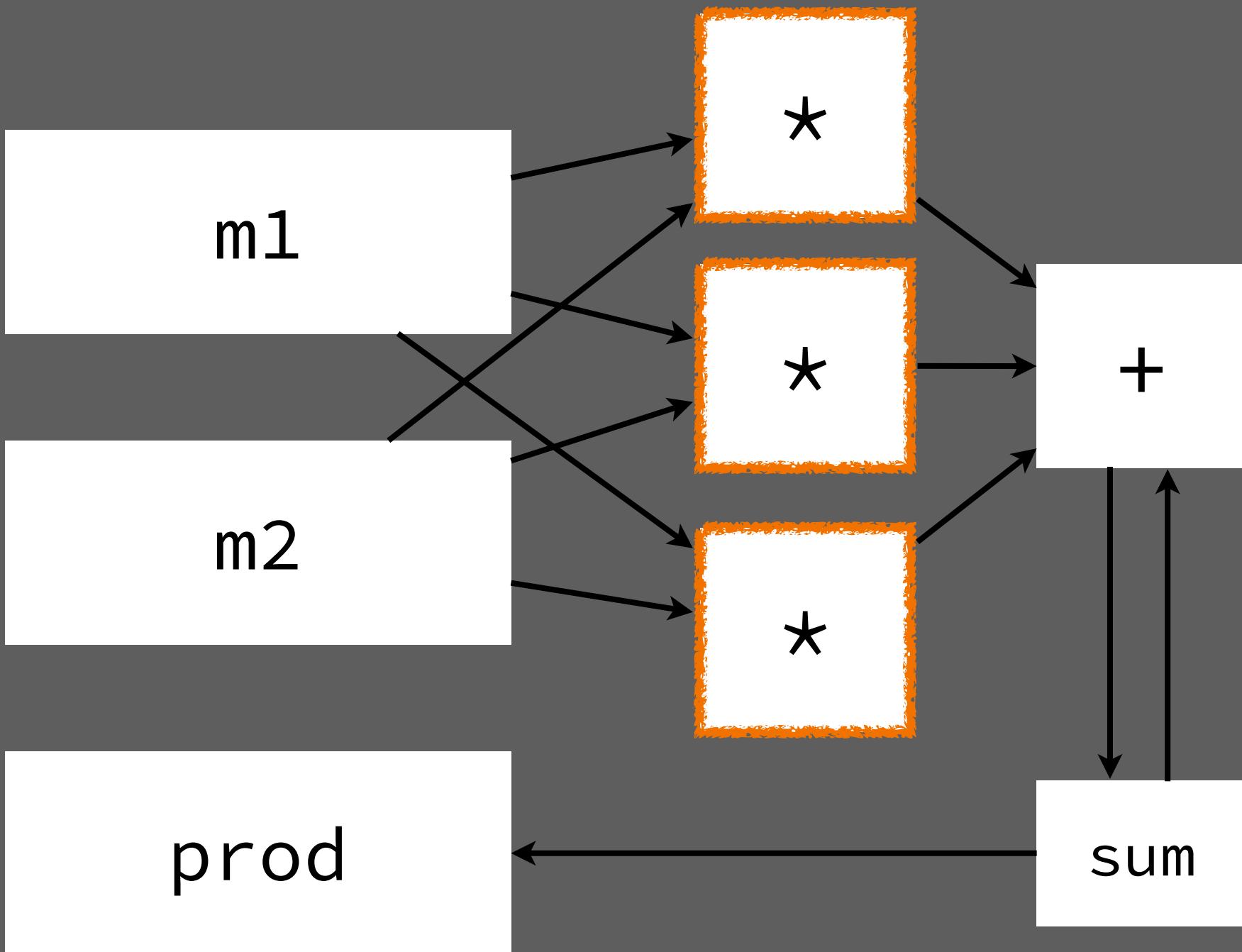
Area-Latency trade-offs



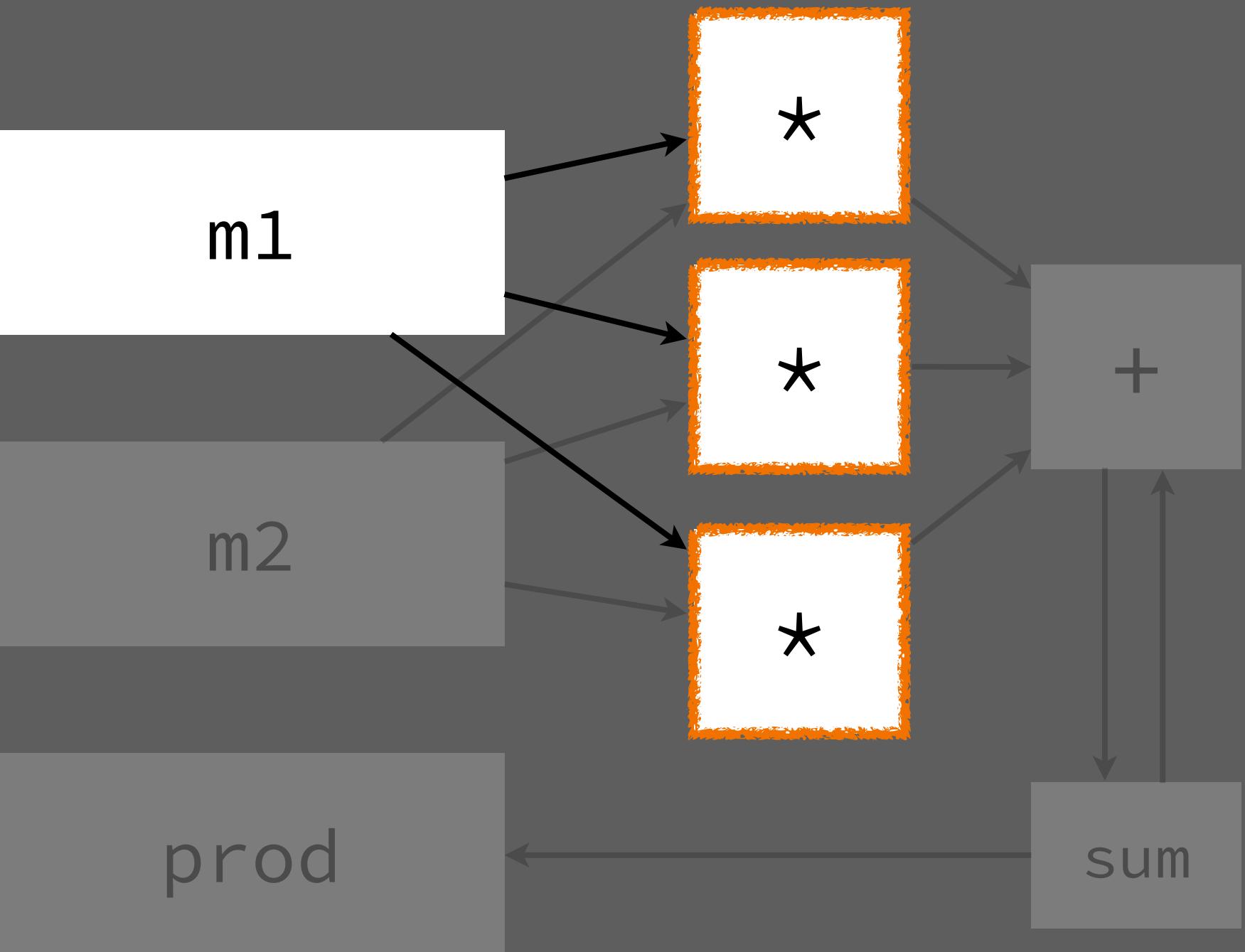
Area-Latency trade-offs



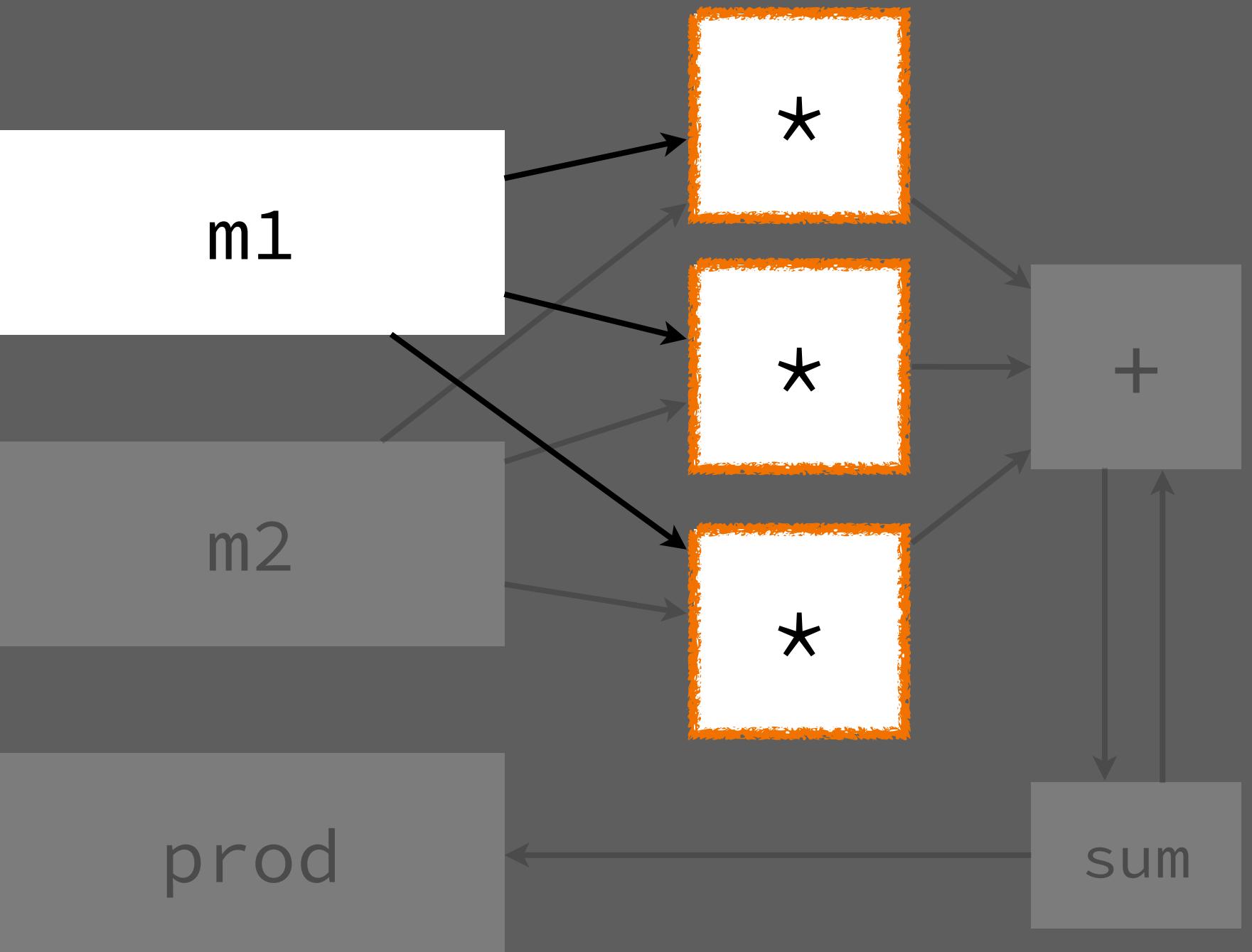
Hardware



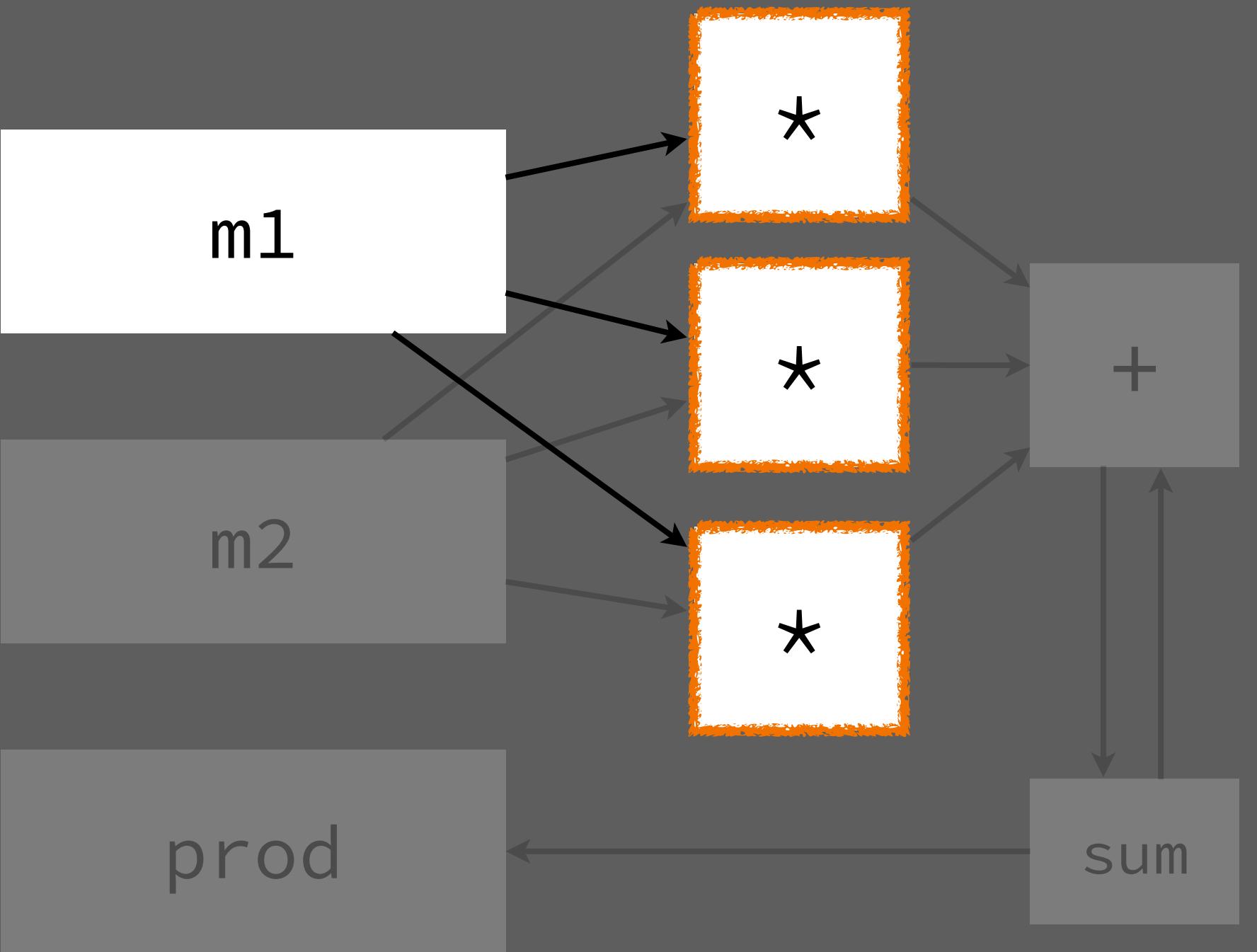
Hardware



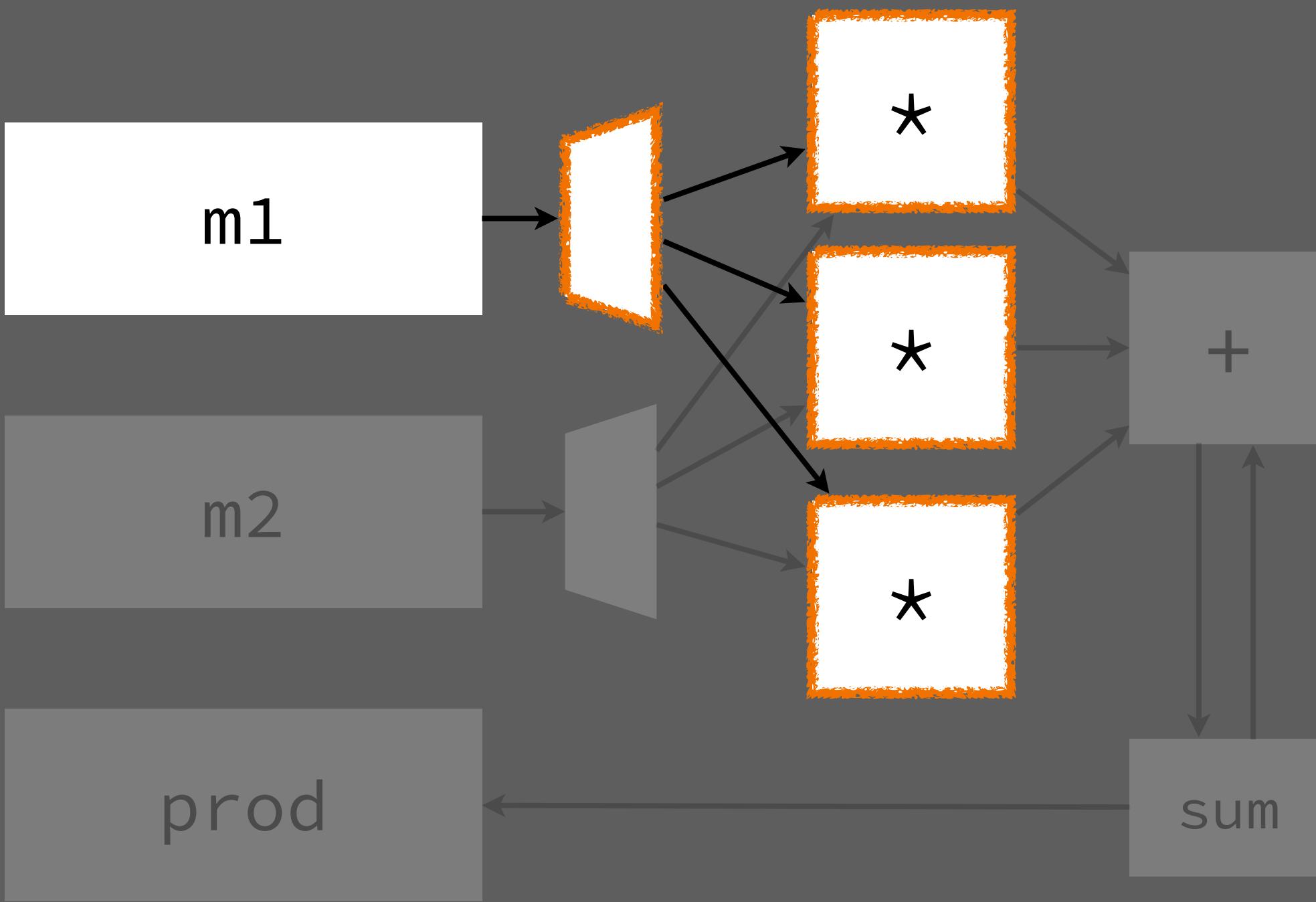
Hardware



Hardware

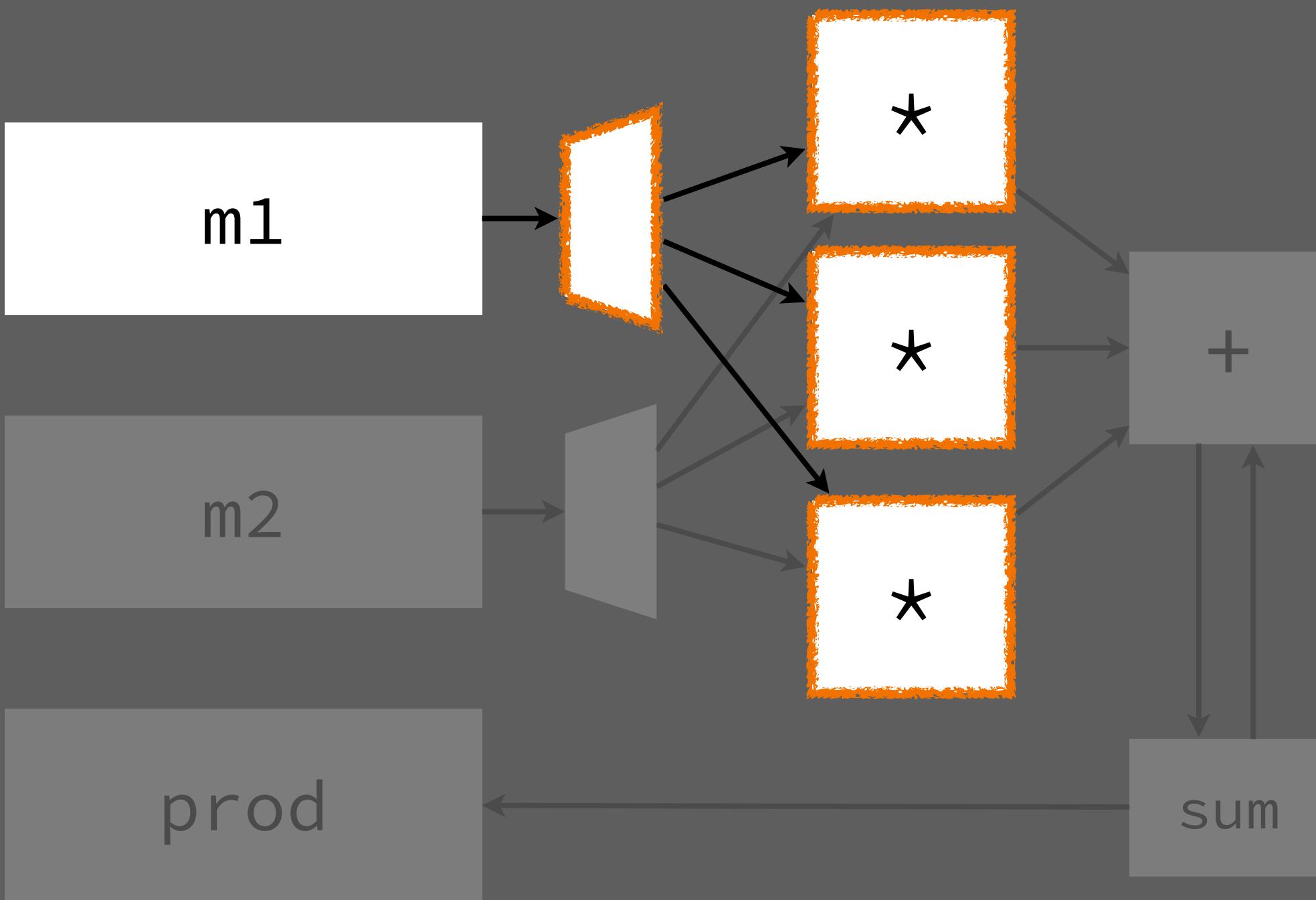


Hardware



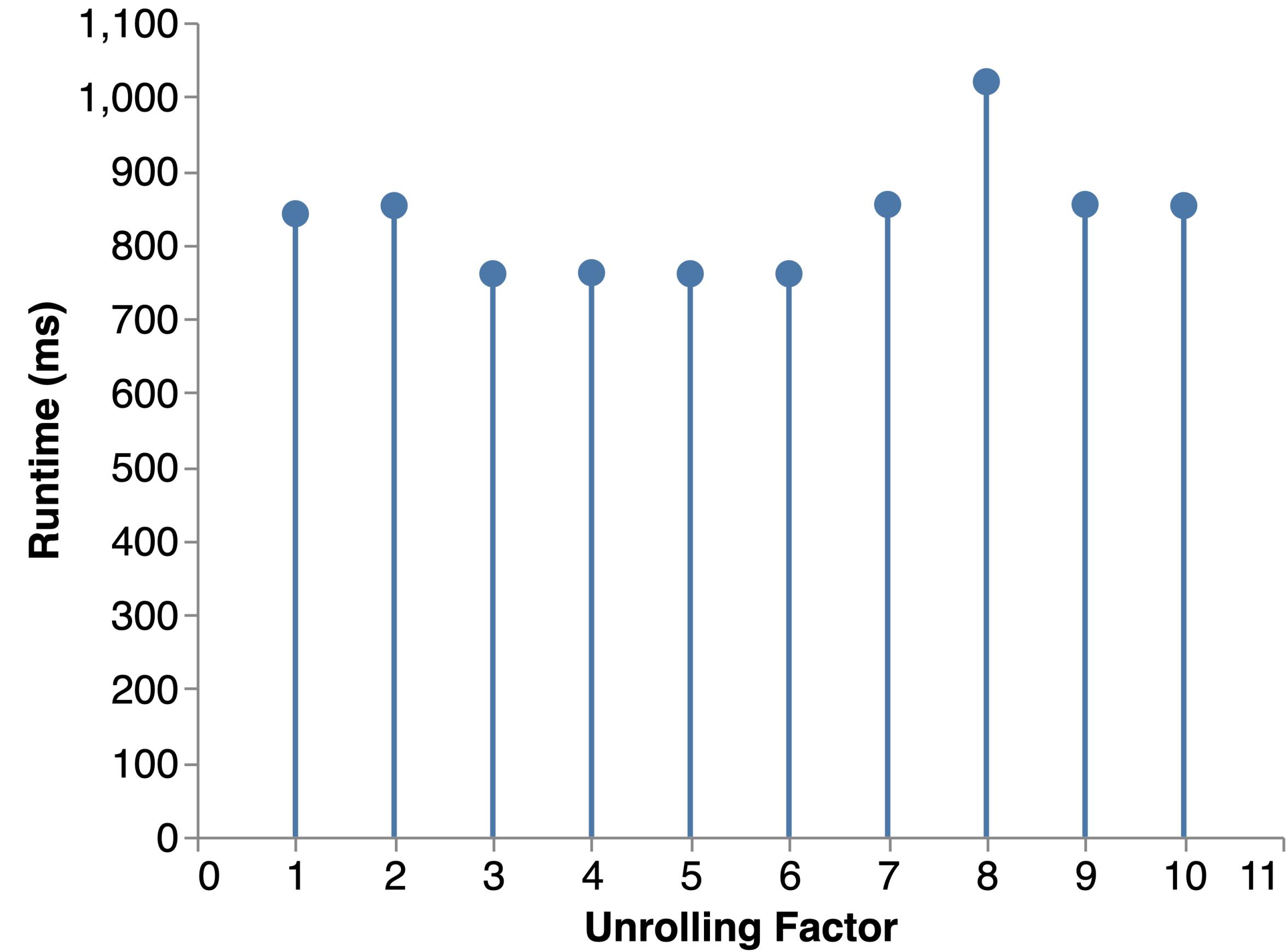
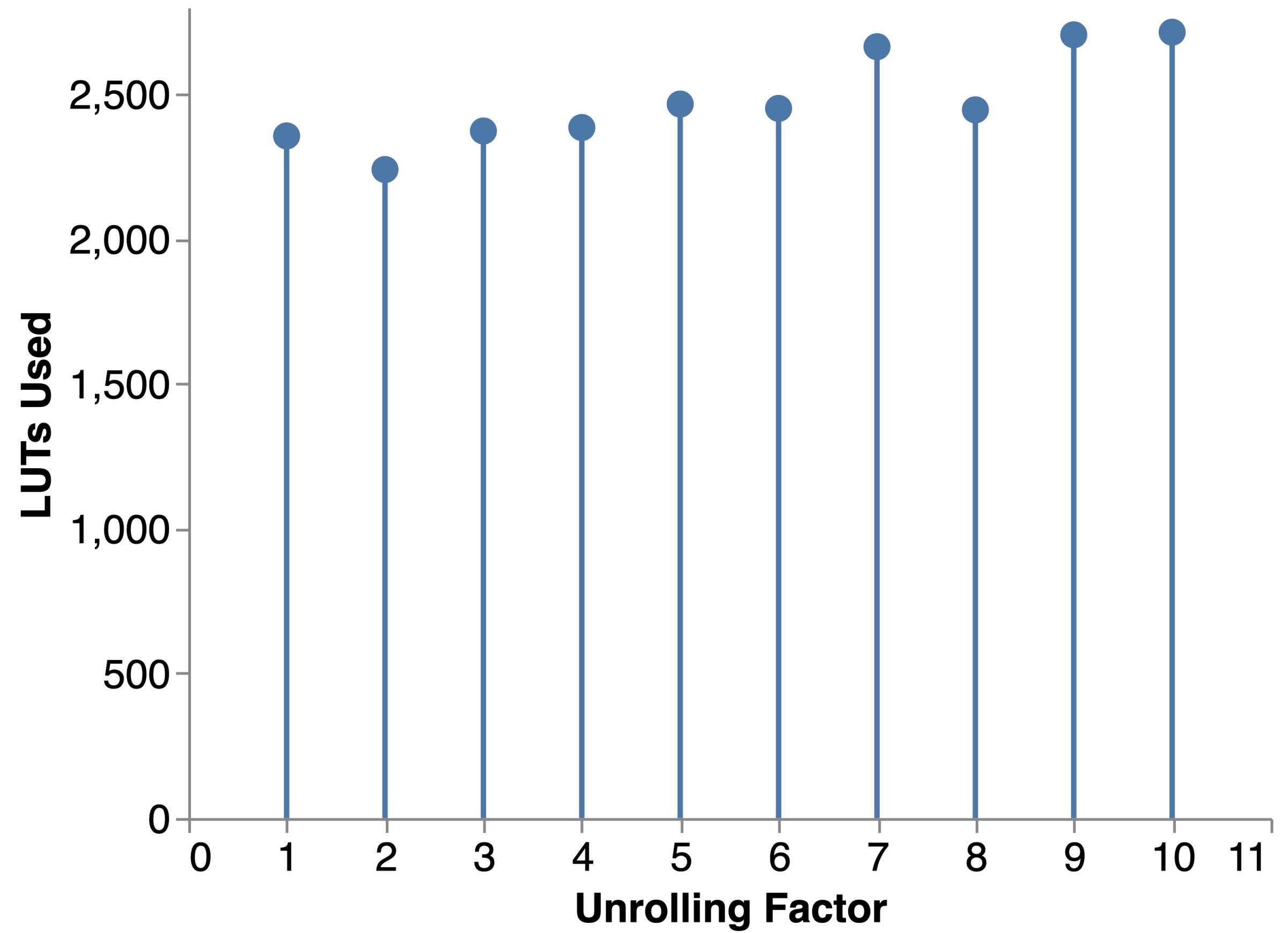
Multiplexers

Hardware



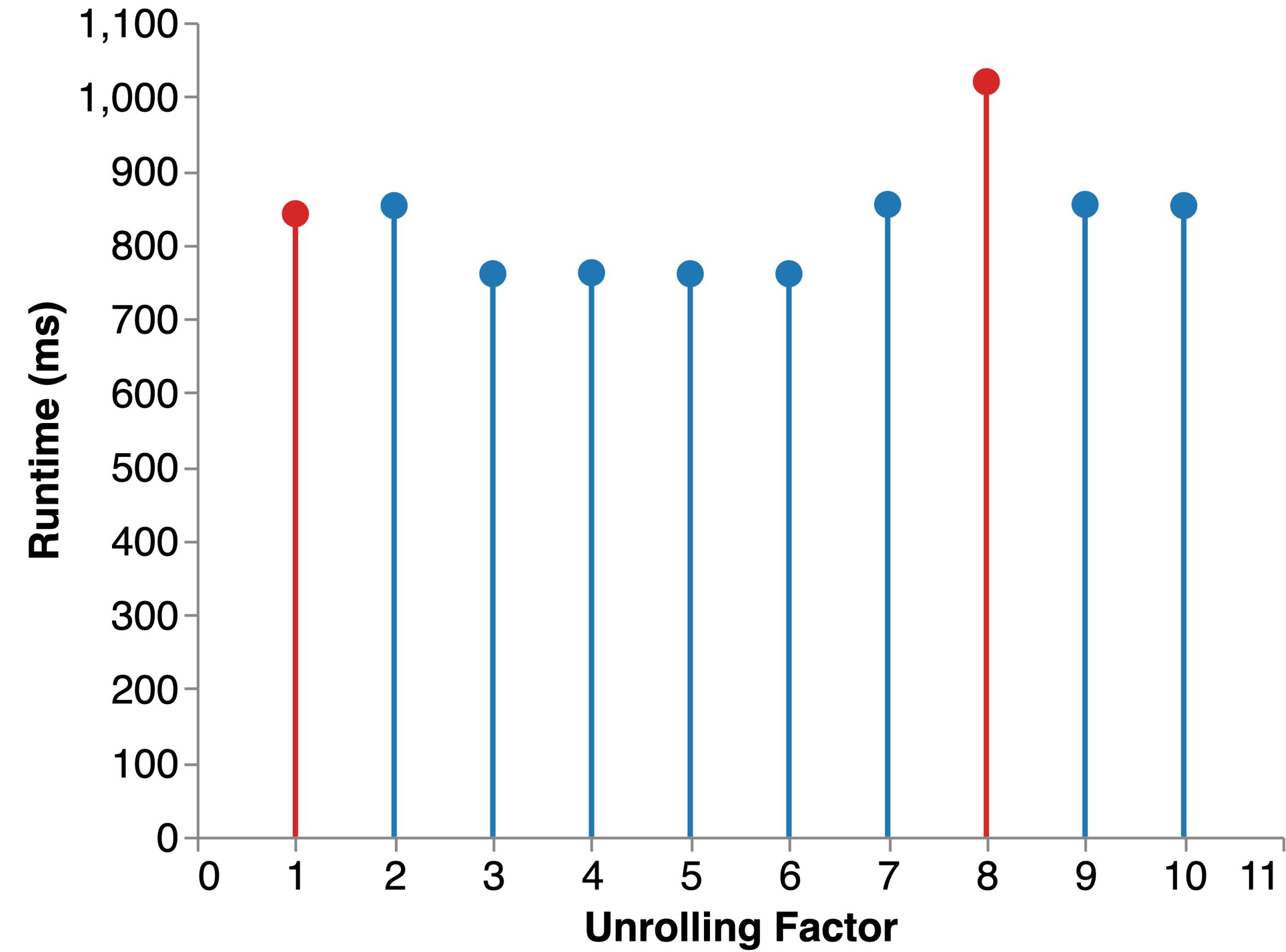
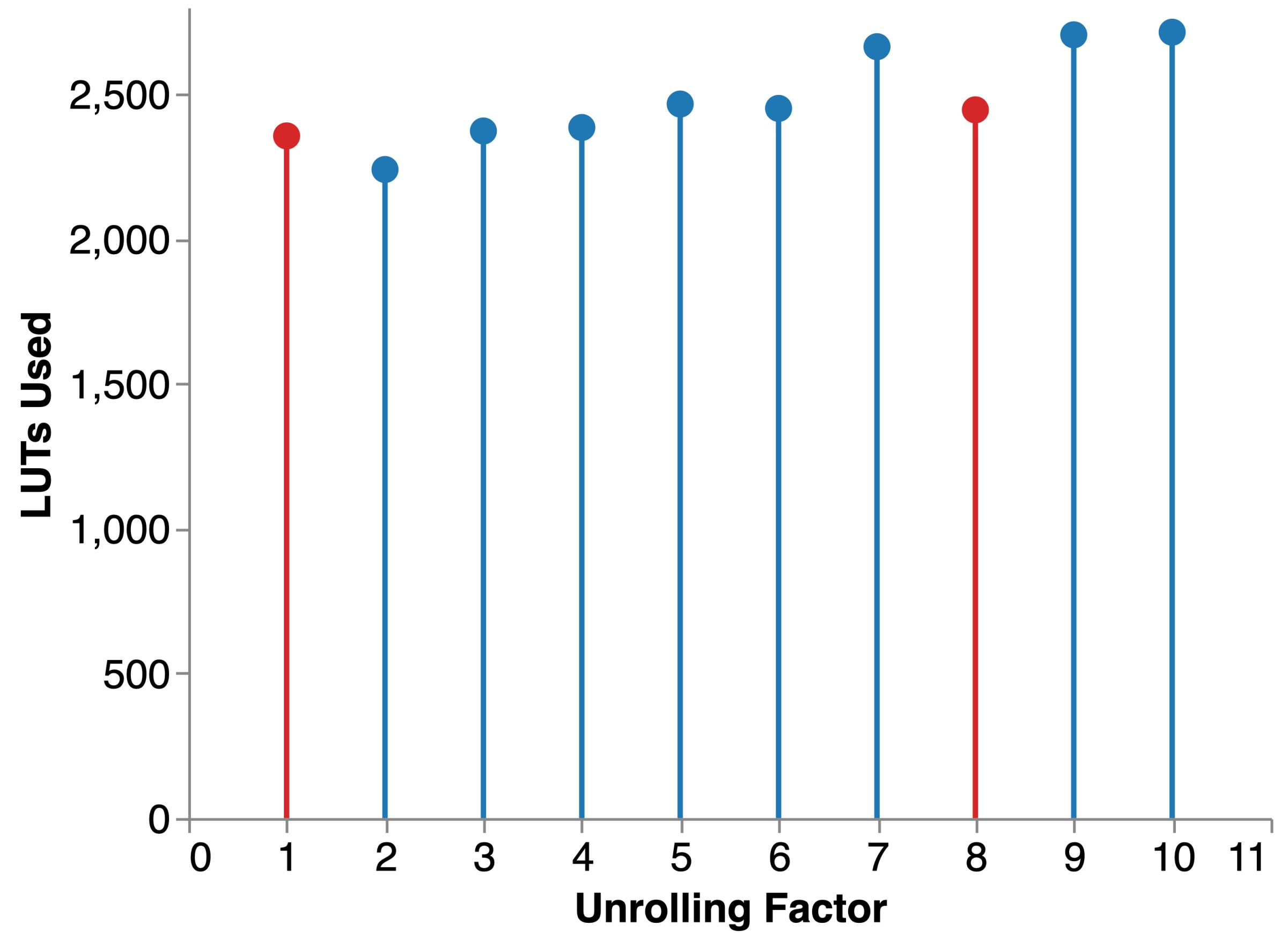
Multiplexers





Area-Latency trade-offs





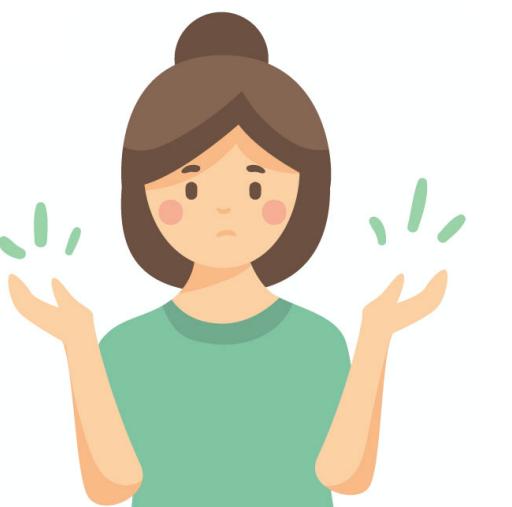
Area-Latency trade-offs



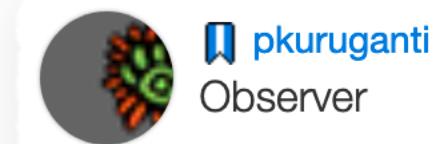
Un·pre·dict·a·ble

adjective

Behaving in a way that is not easily predicted.



2012.1 Vivado HLS - Mutually exclusive memory access is not implemented with MUX on addresses, and reports: "@W [SCHED-69] Unable to schedule 'load' operation on array 'x' due to limited resources (II = 1) "



pkuruganti
Observer

09-21-2017 02:38 PM

1,879 Views

Registered: 02-21-2013

SdAccle 2017.2 + KCU1500: [XOCC 204-69] Unable to schedule 'load/store' operation



tsoliman
Visitor

02-06-2019 04:47 AM - edited 02-06-2019 04:49 AM

1,973 Views

Registered: 02-06-2019

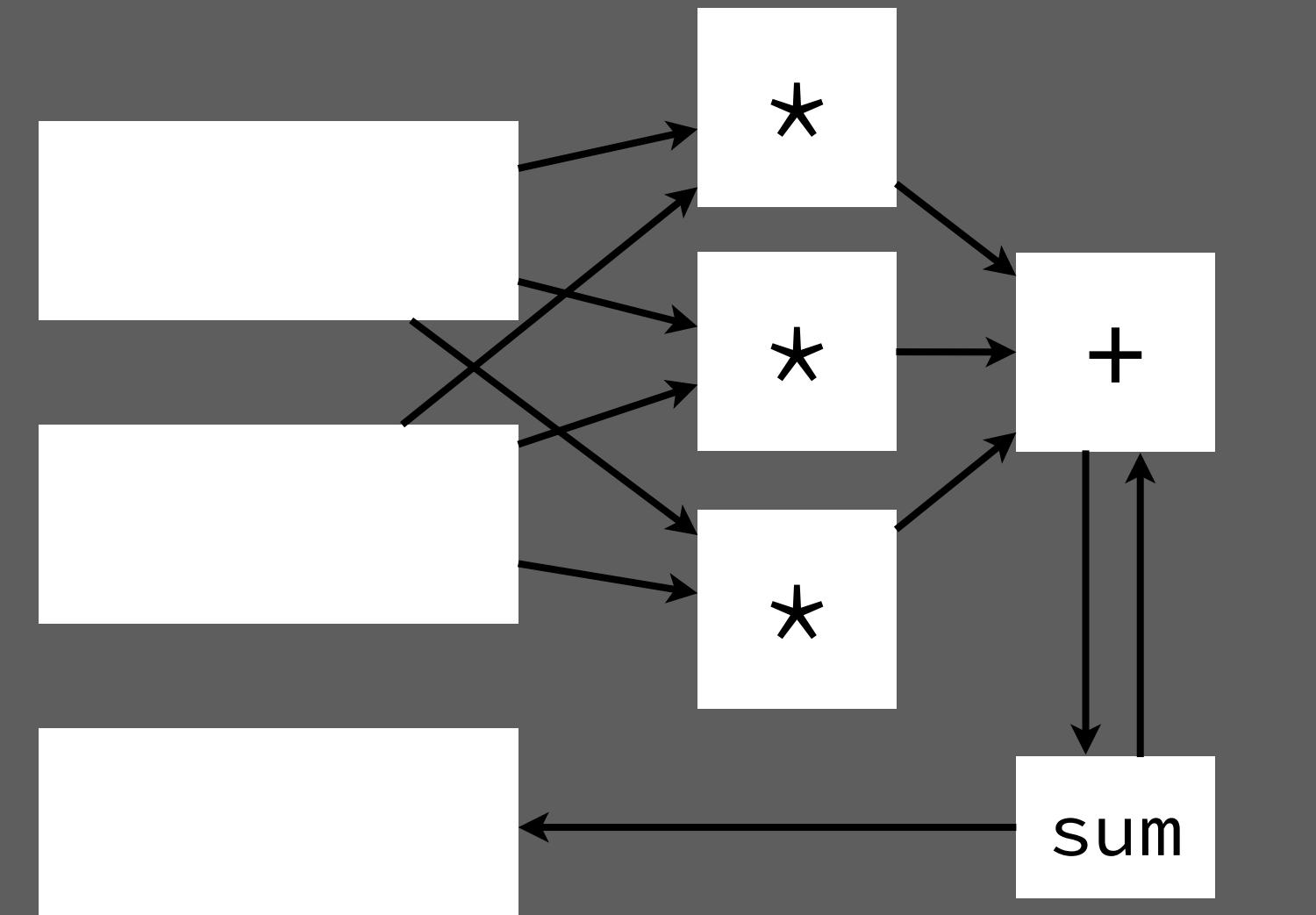
Unable to schedule 'load' operation



Super secret™ accelerator

```
int m1[512][512];
#pragma PARTITION factor=3
int m2[512][512];
#pragma PARTITION factor=3
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            #pragma HLS UNROLL factor=3
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

Hardware



Block RAM

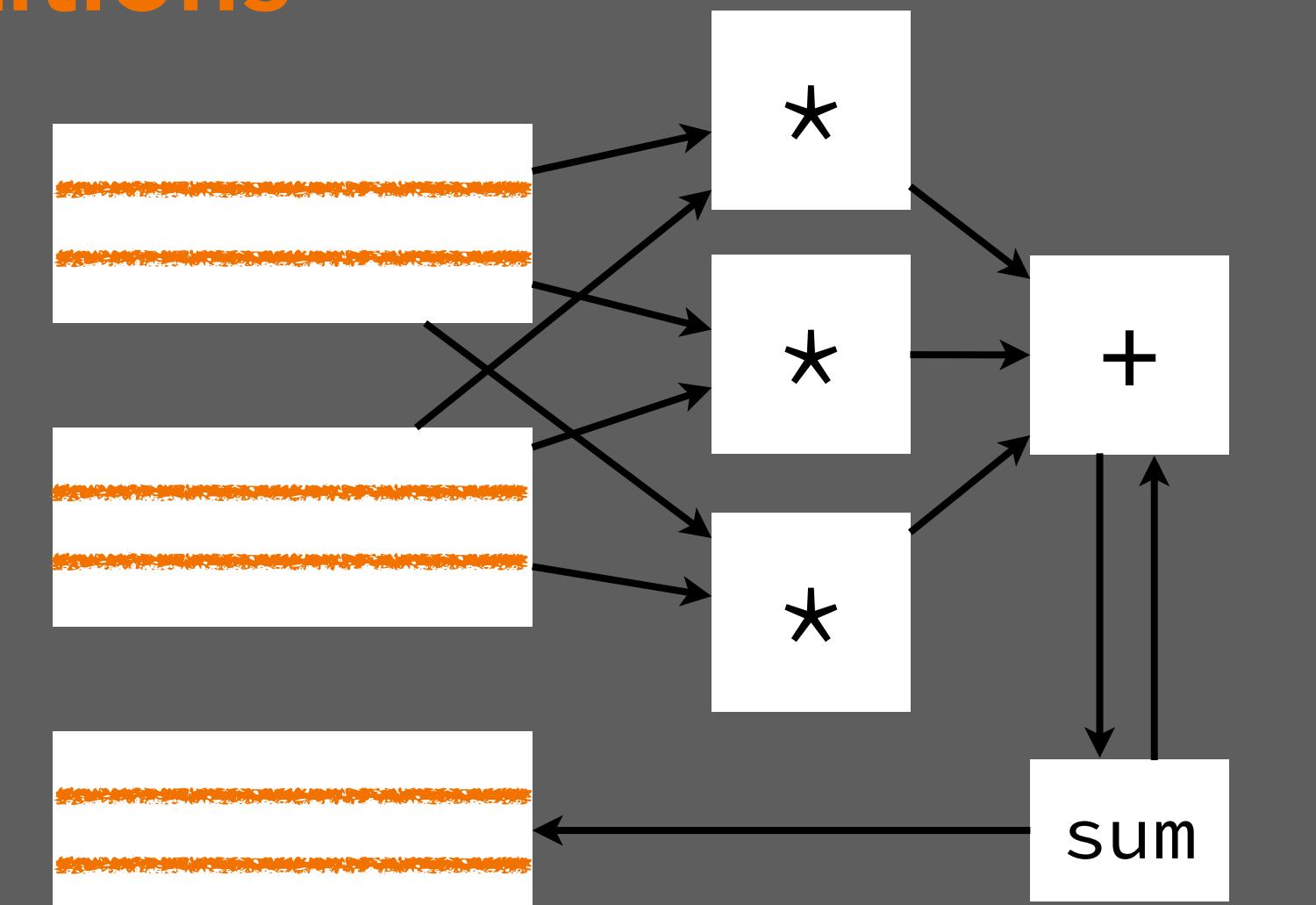
Register

Super secret™ accelerator

```
int m1[512][512];
#pragma PARTITION factor=3
int m2[512][512];
#pragma PARTITION factor=3
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            #pragma HLS UNROLL factor=3
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

Hardware

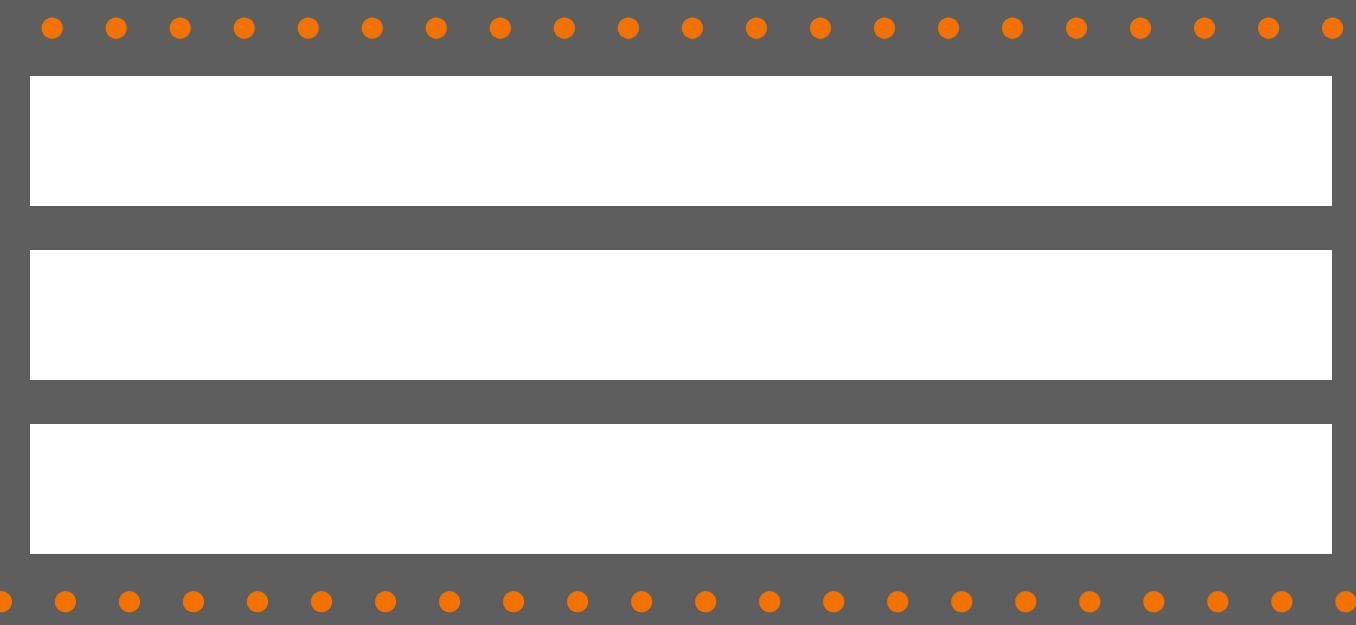
Memory Partitions



Block RAM

Register

Hardware



Logical
Memory

Hardware

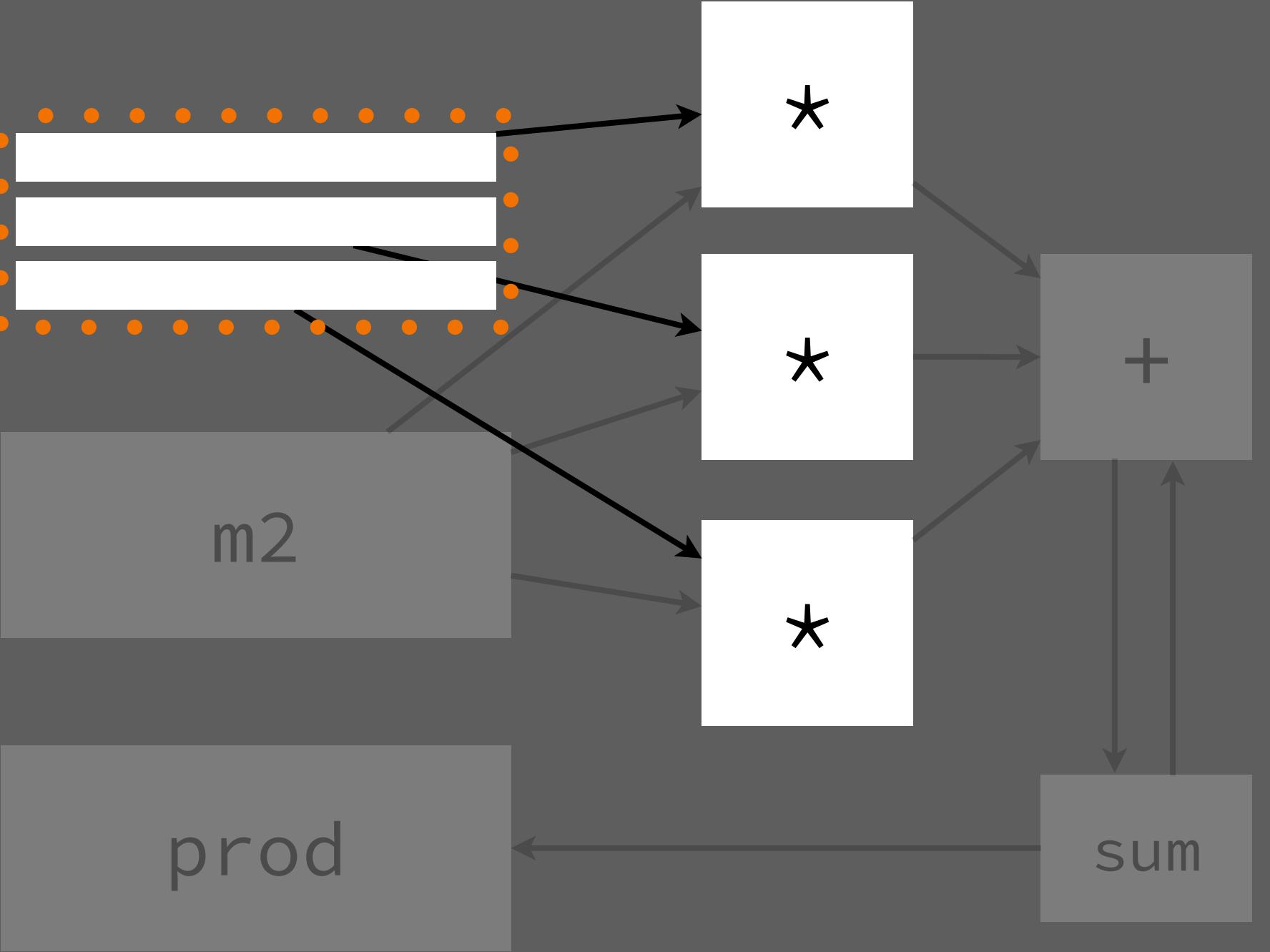
Physical
Partitions



Logical
Memory

Generated hardware

Memory
Partitions

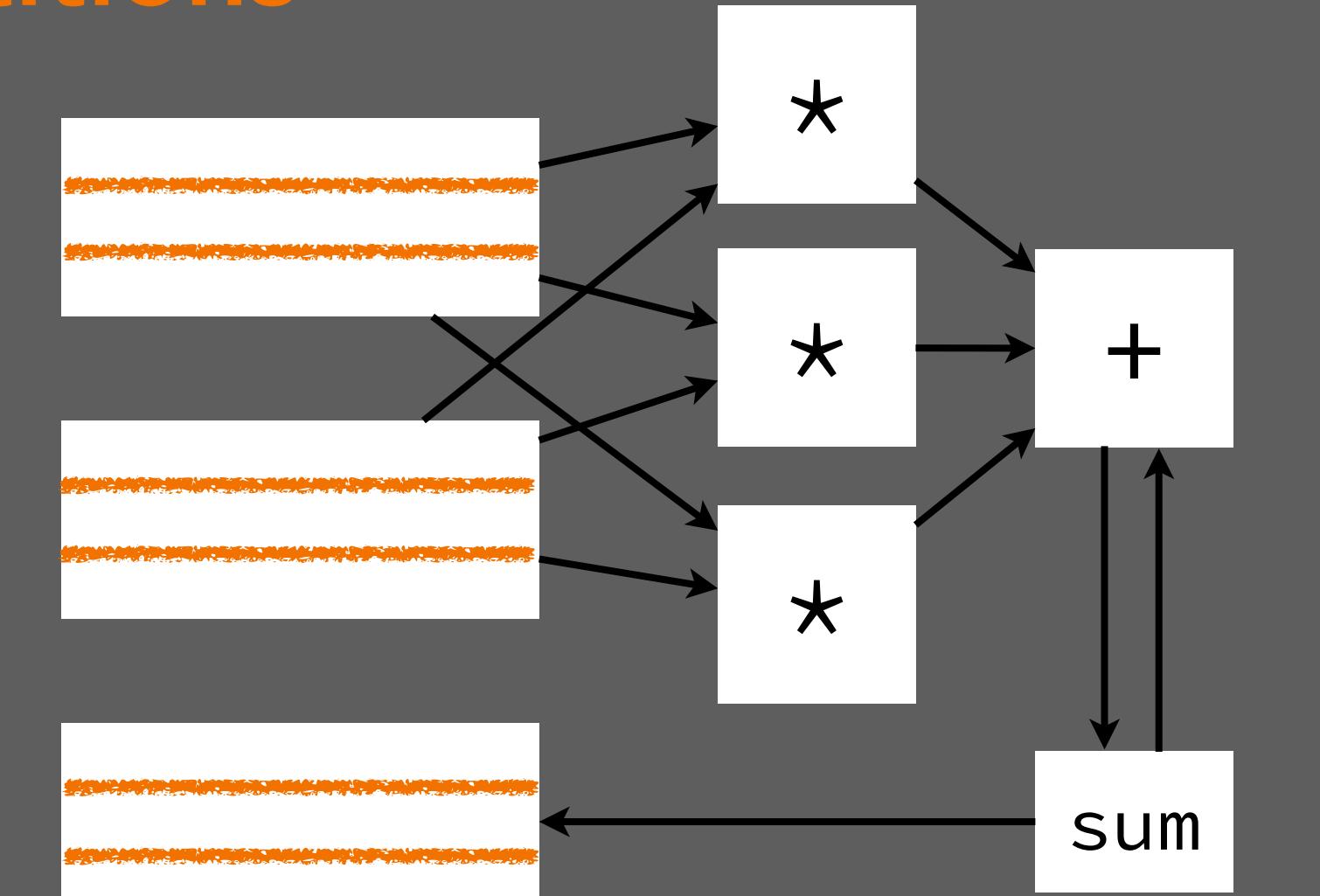


Super secret™ accelerator

```
int m1[512][512];
#pragma PARTITION factor=3
int m2[512][512];
#pragma PARTITION factor=3
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            #pragma HLS UNROLL factor=3
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

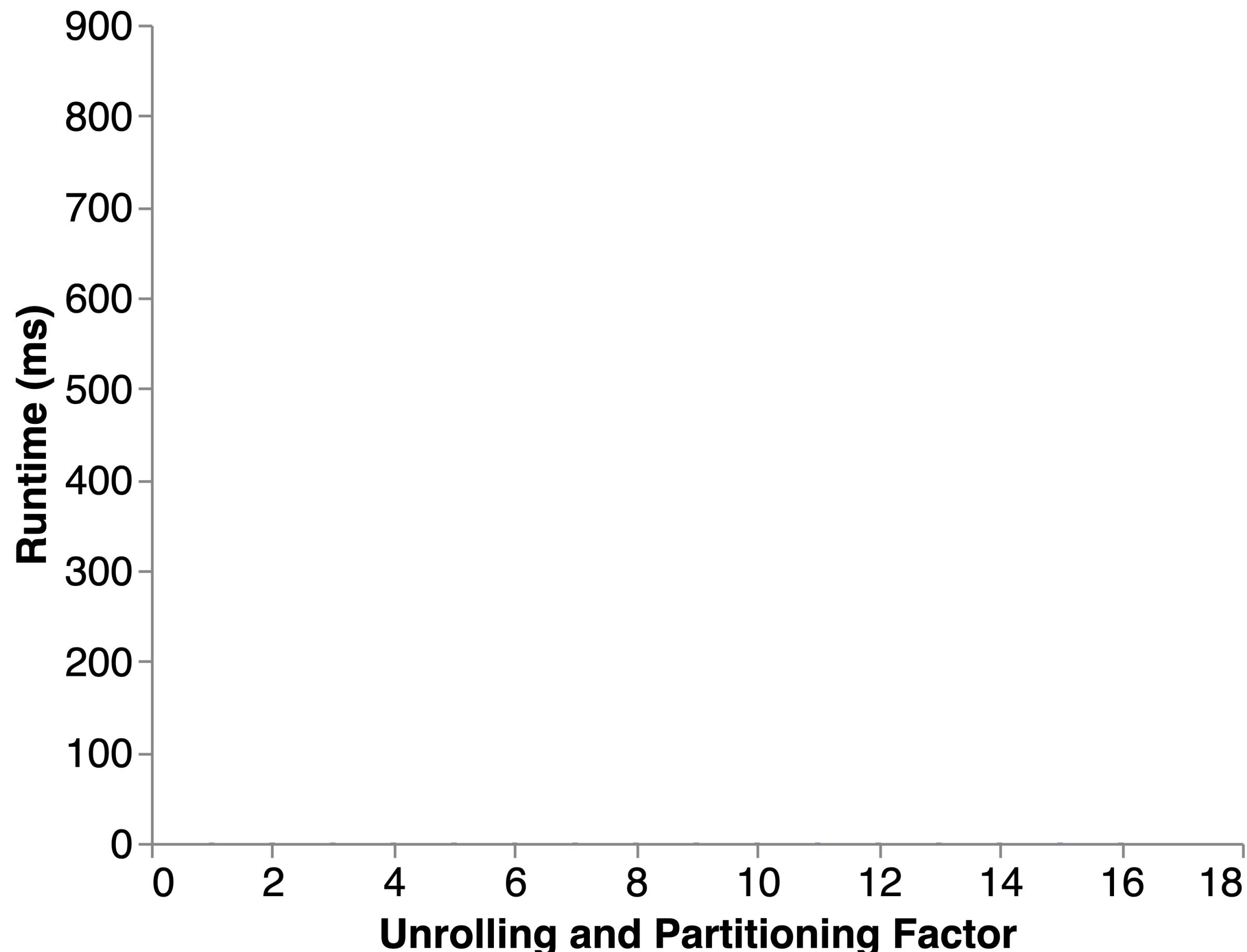
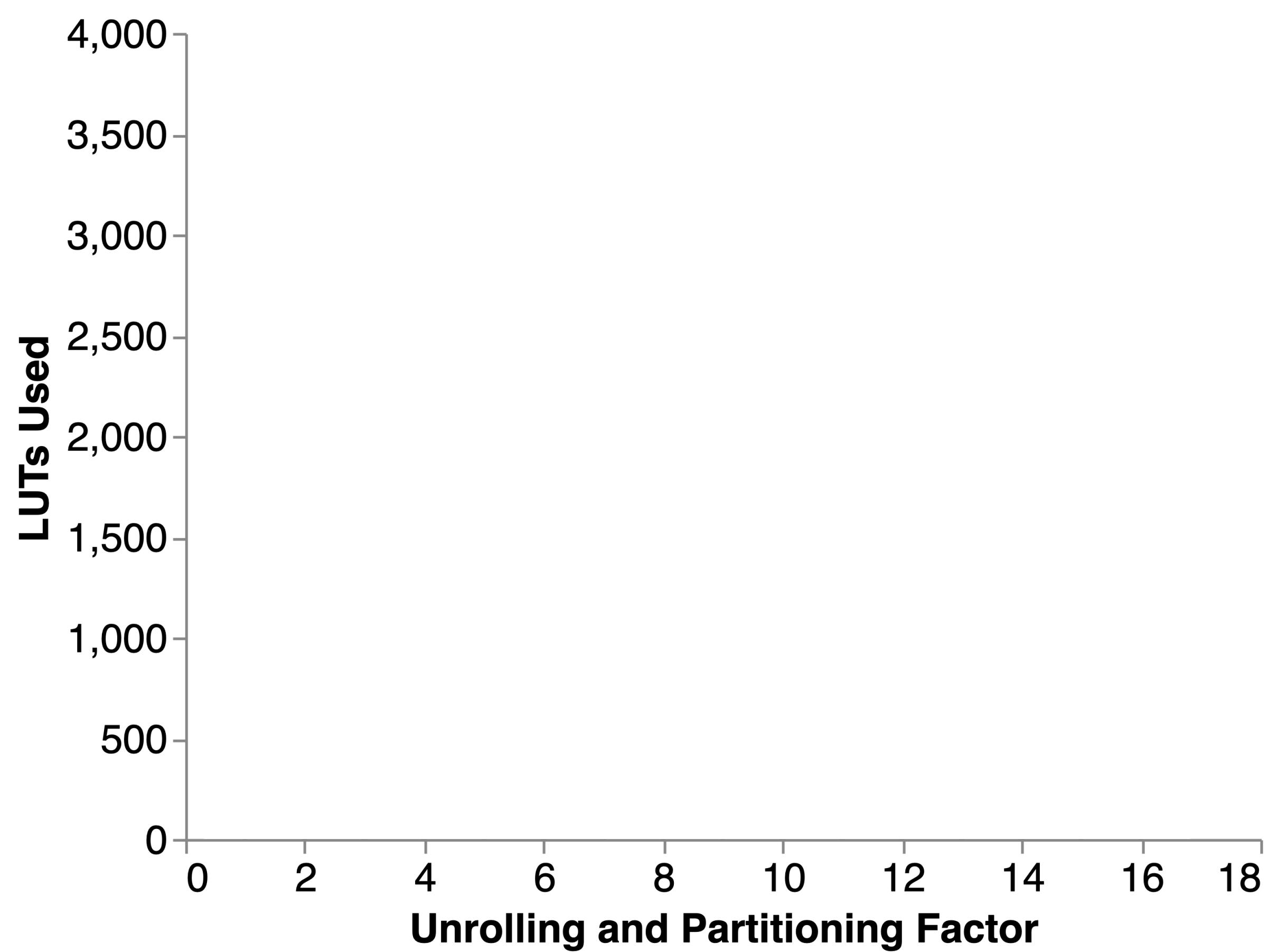
Hardware

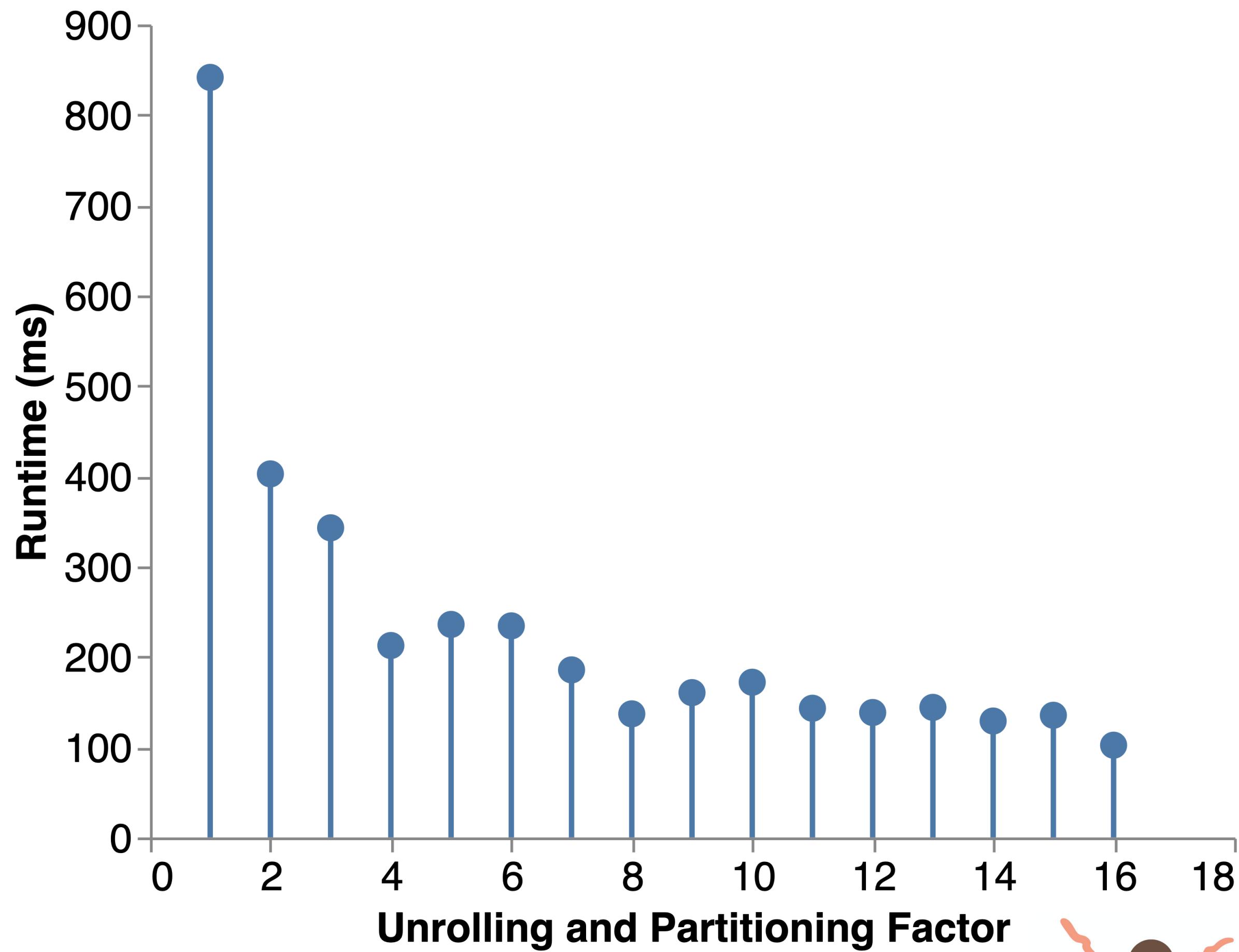
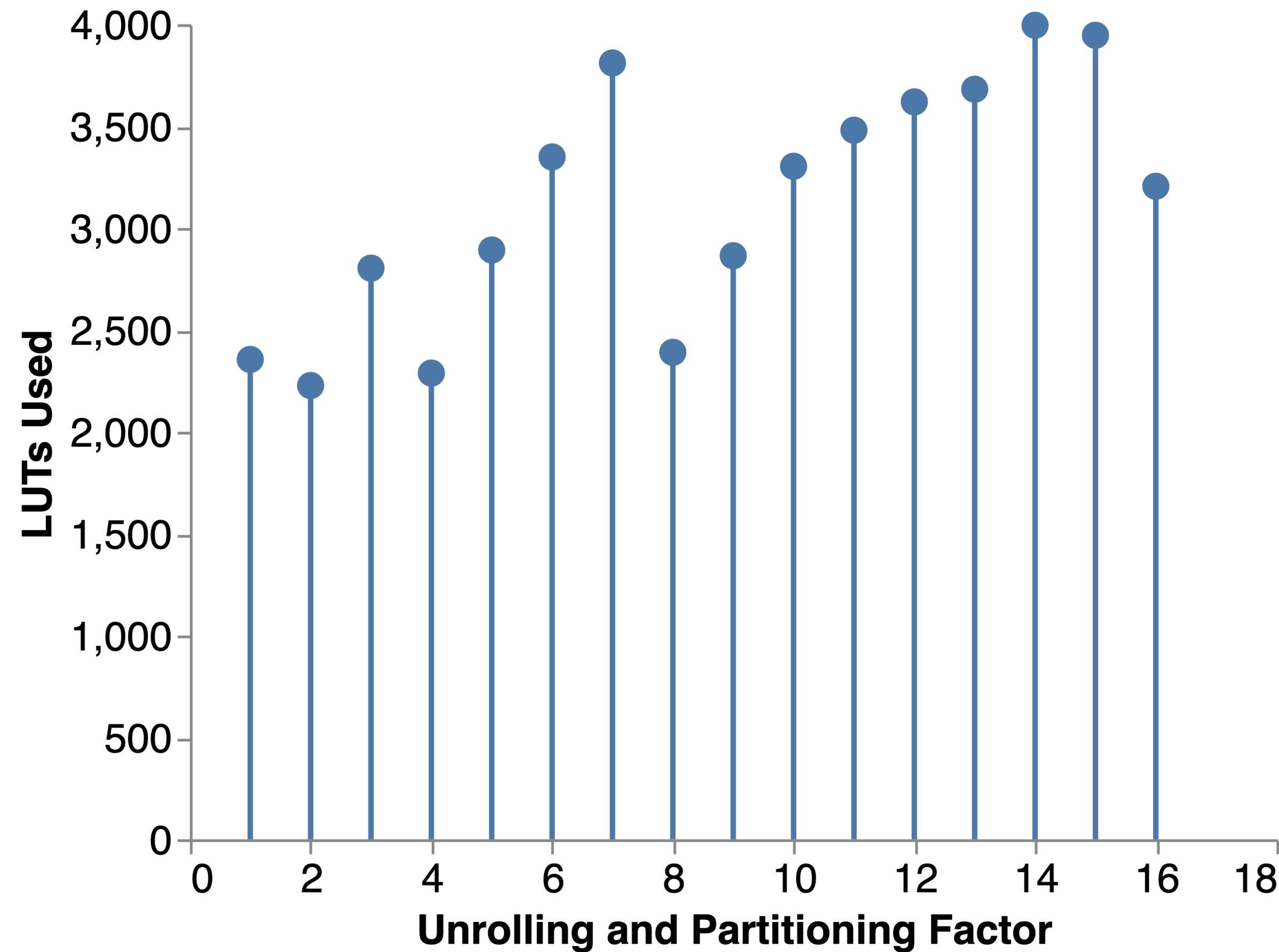
Memory Partitions

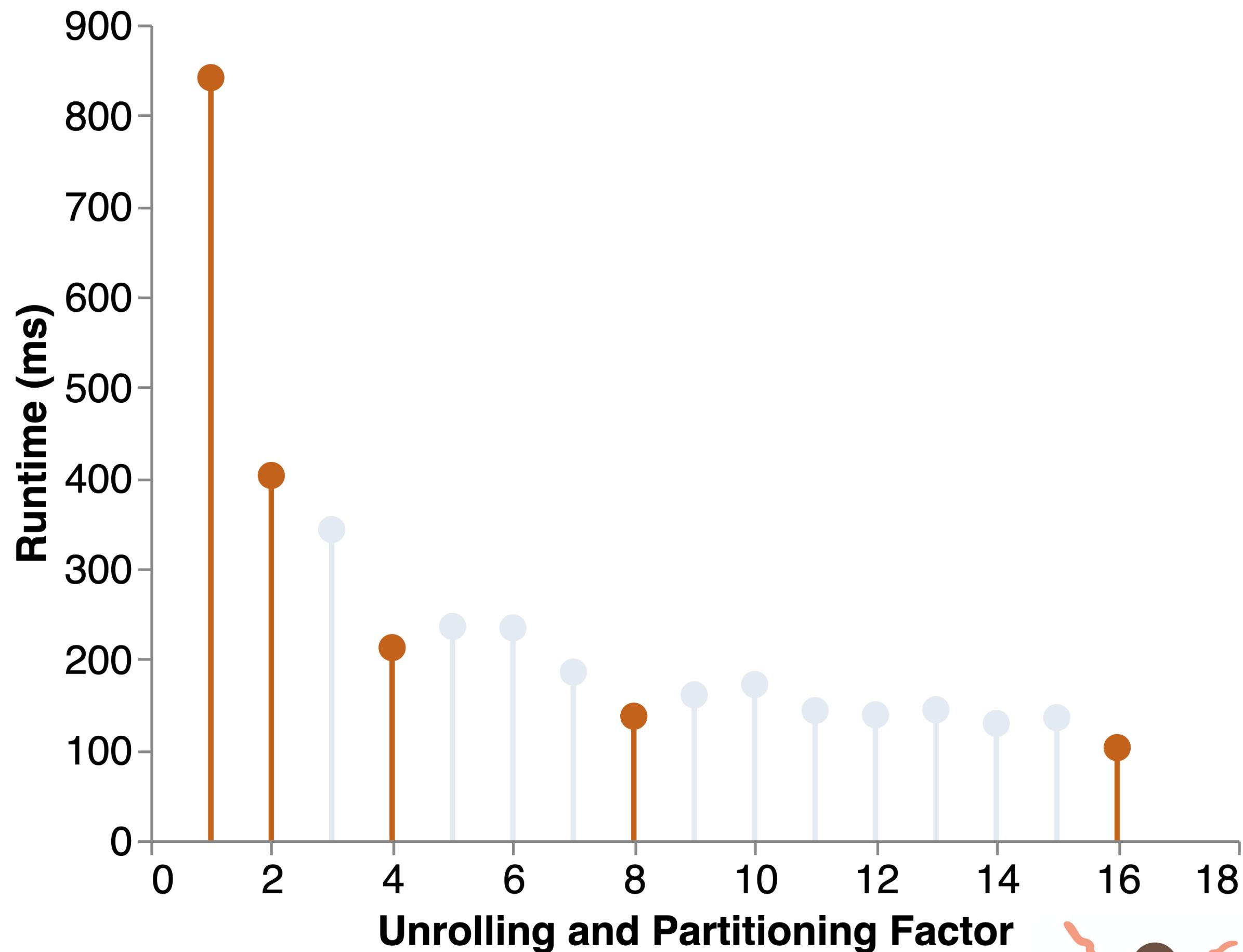
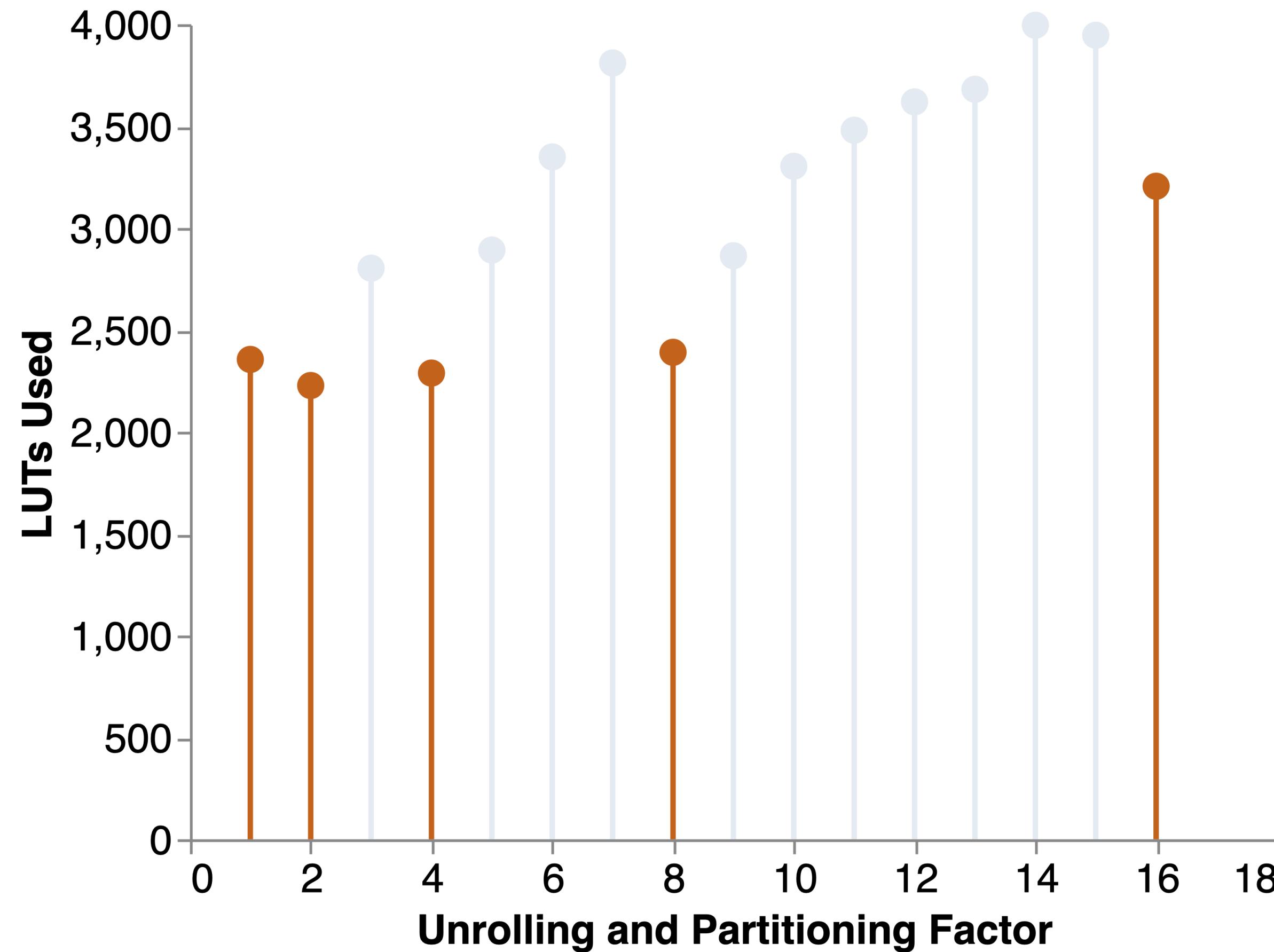


Block RAM

Register





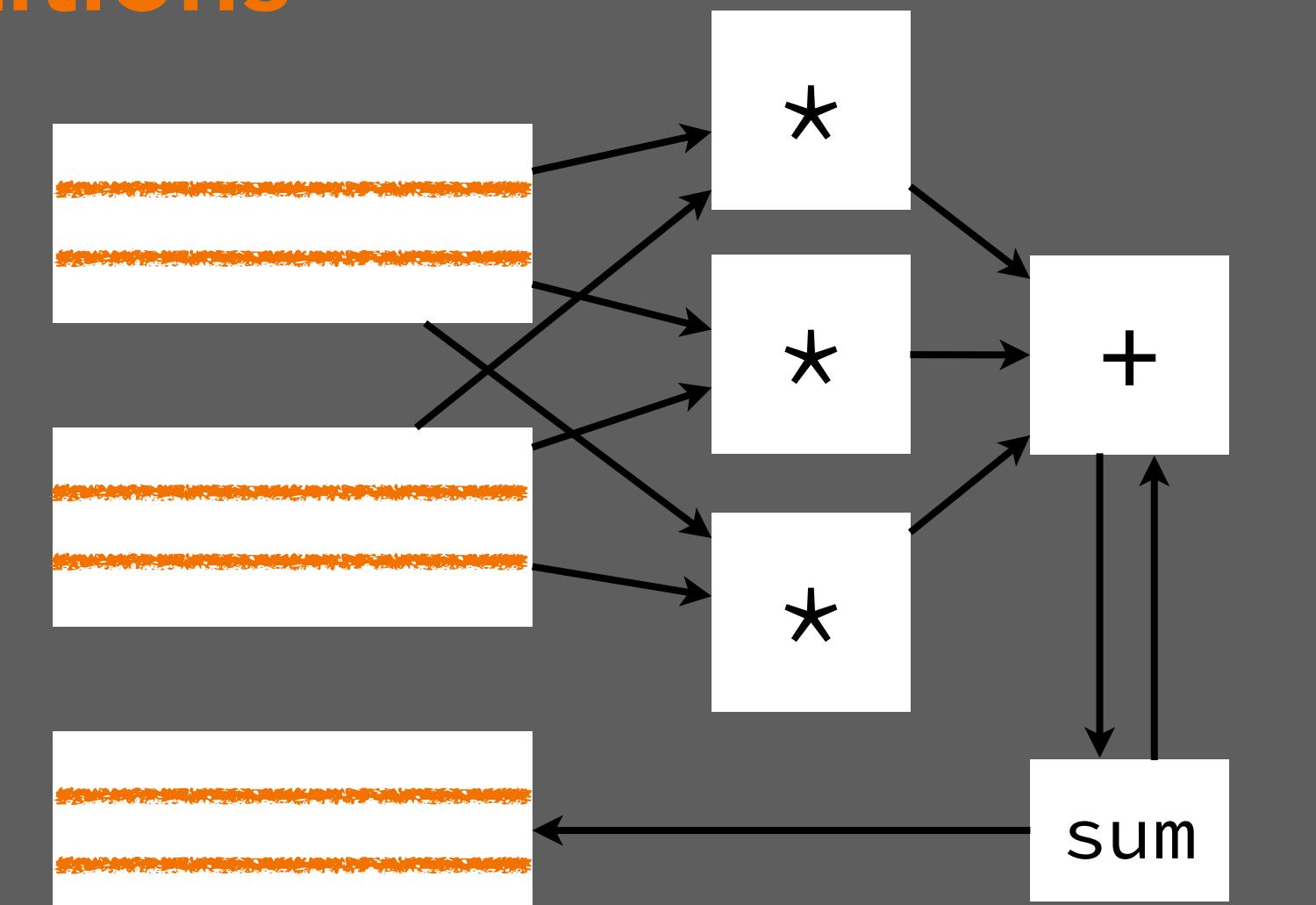


Super secret™ accelerator

```
int m1[512][512];
#pragma PARTITION factor=3
int m2[512][512];
#pragma PARTITION factor=3
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            #pragma HLS UNROLL factor=3
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

Hardware

Memory Partitions

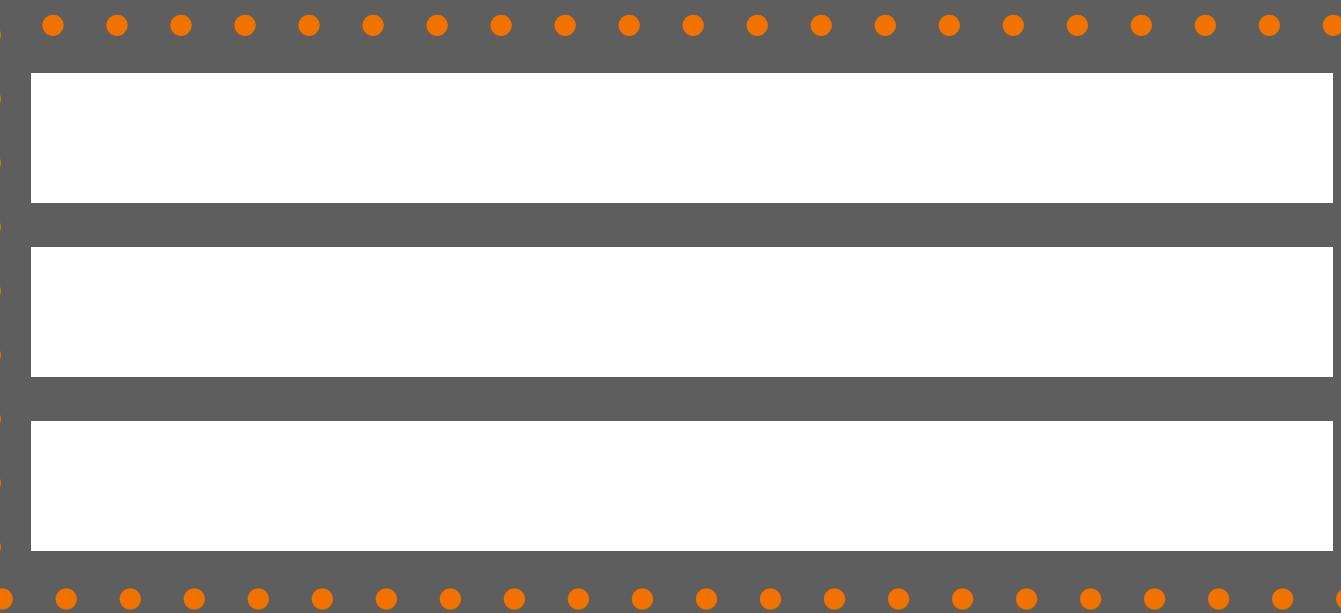


Block RAM

Register

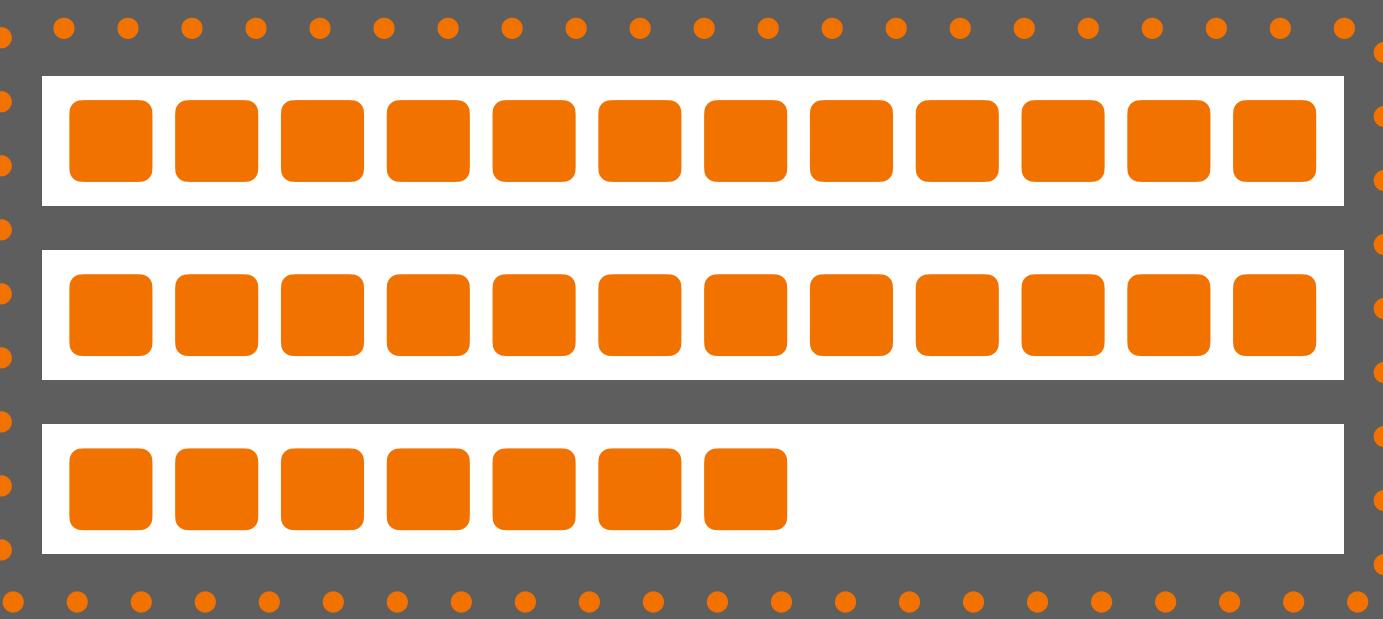
Hardware

512 % 3 ≠ 0



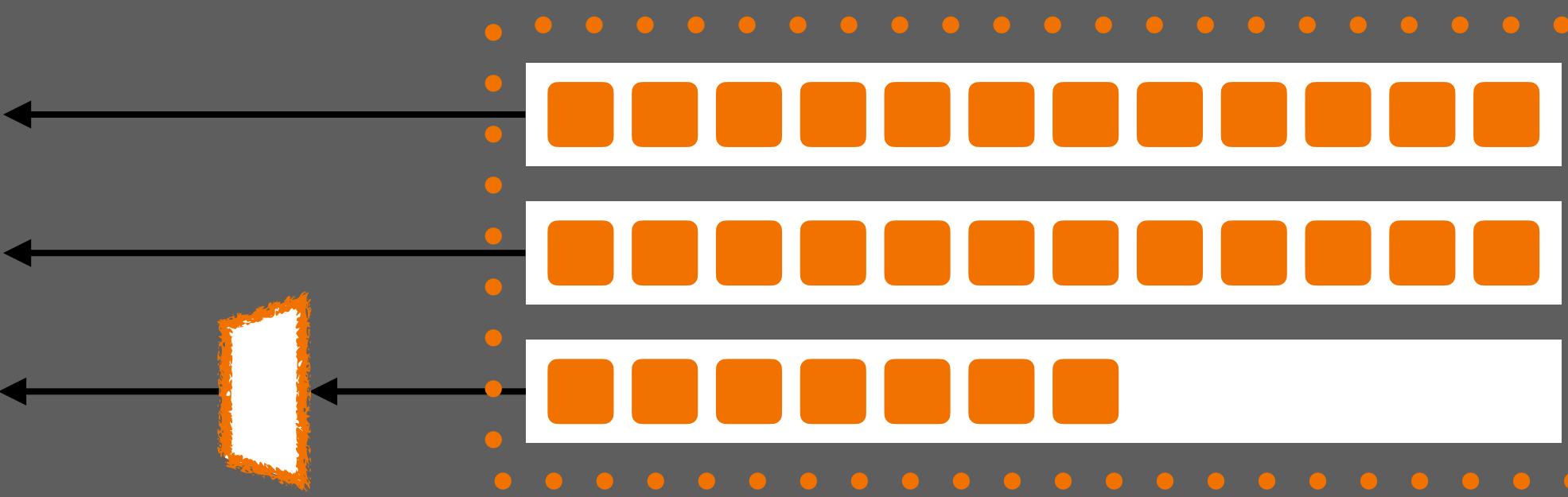
Hardware

512 % 3 ≠ 0

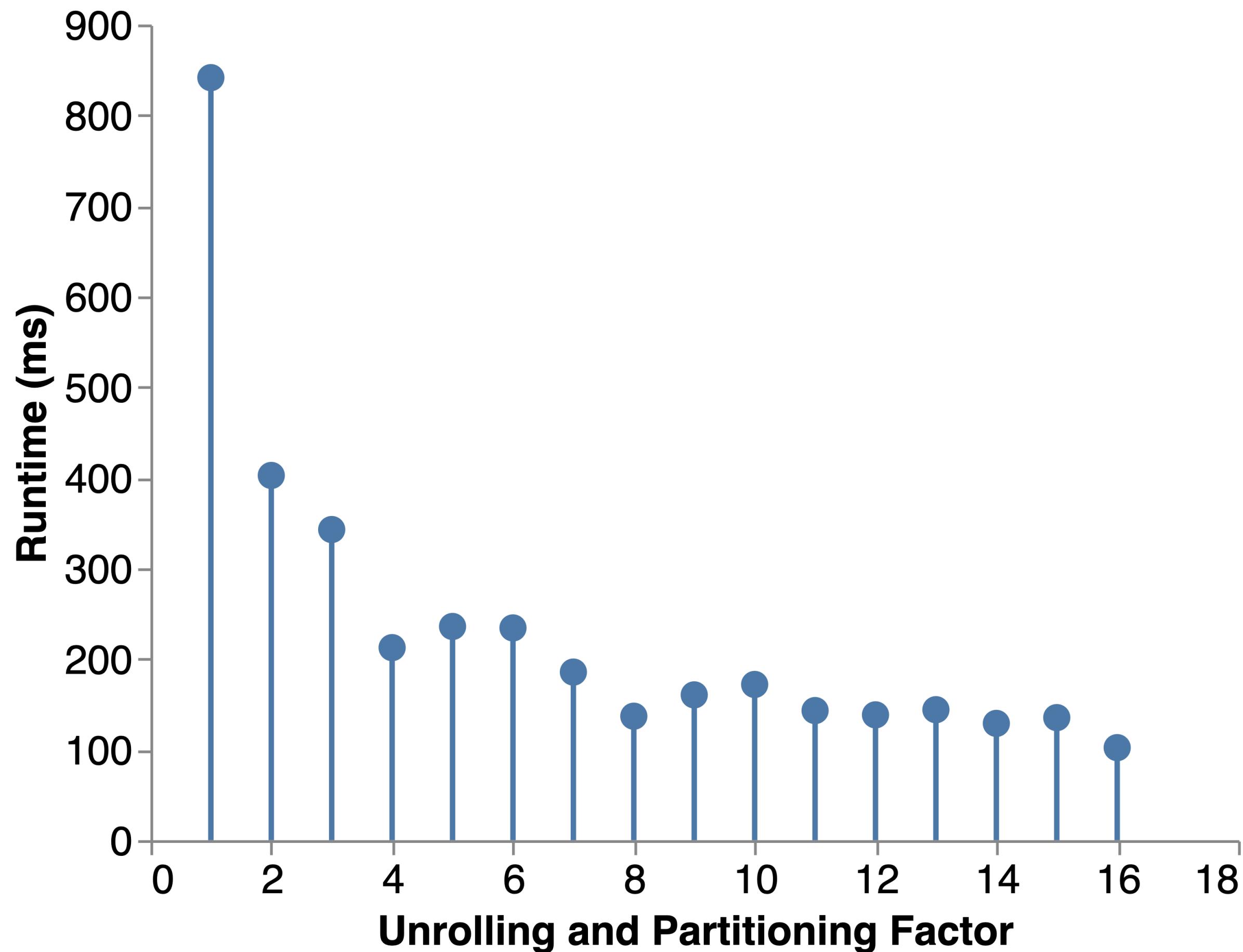
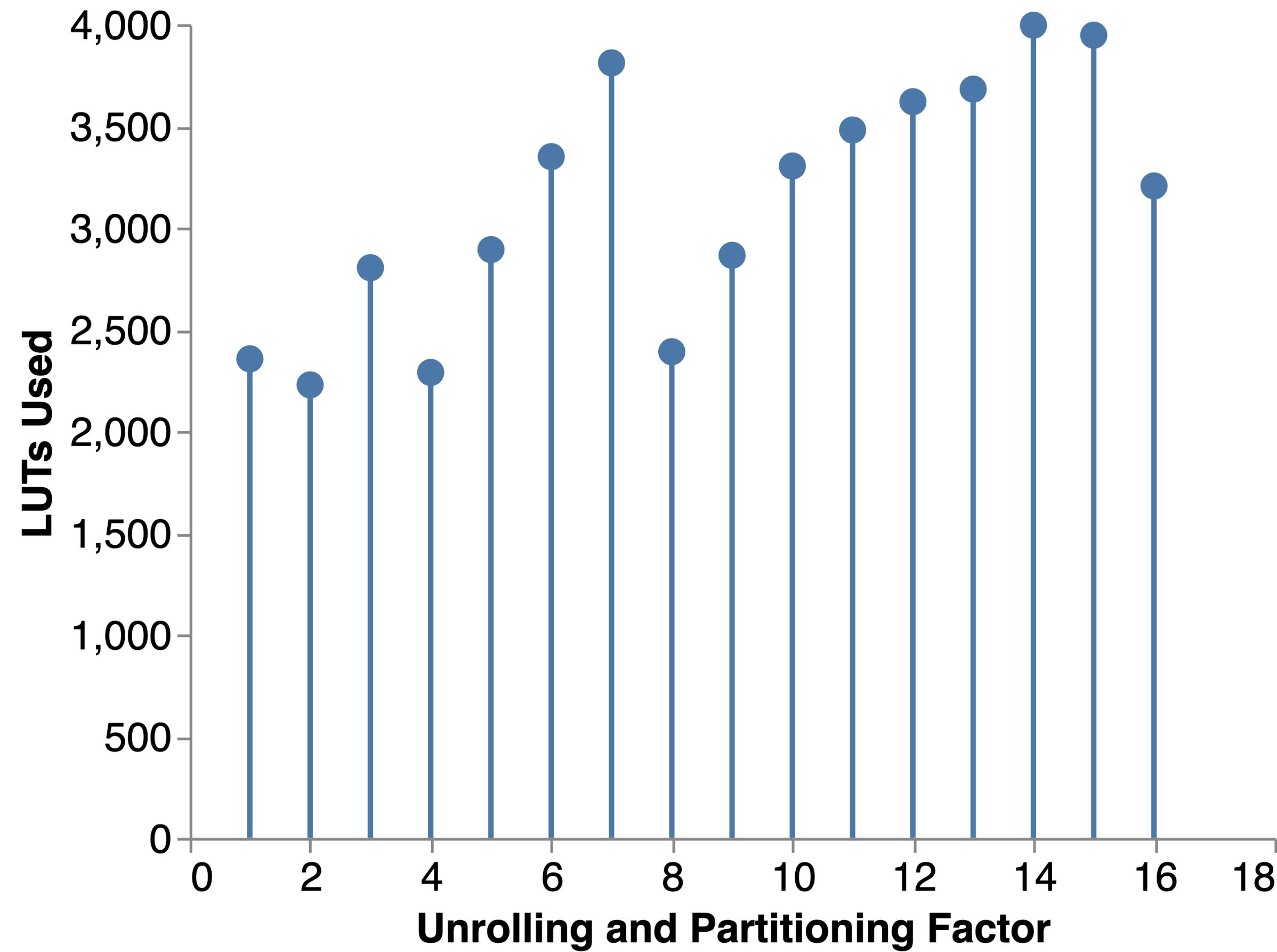


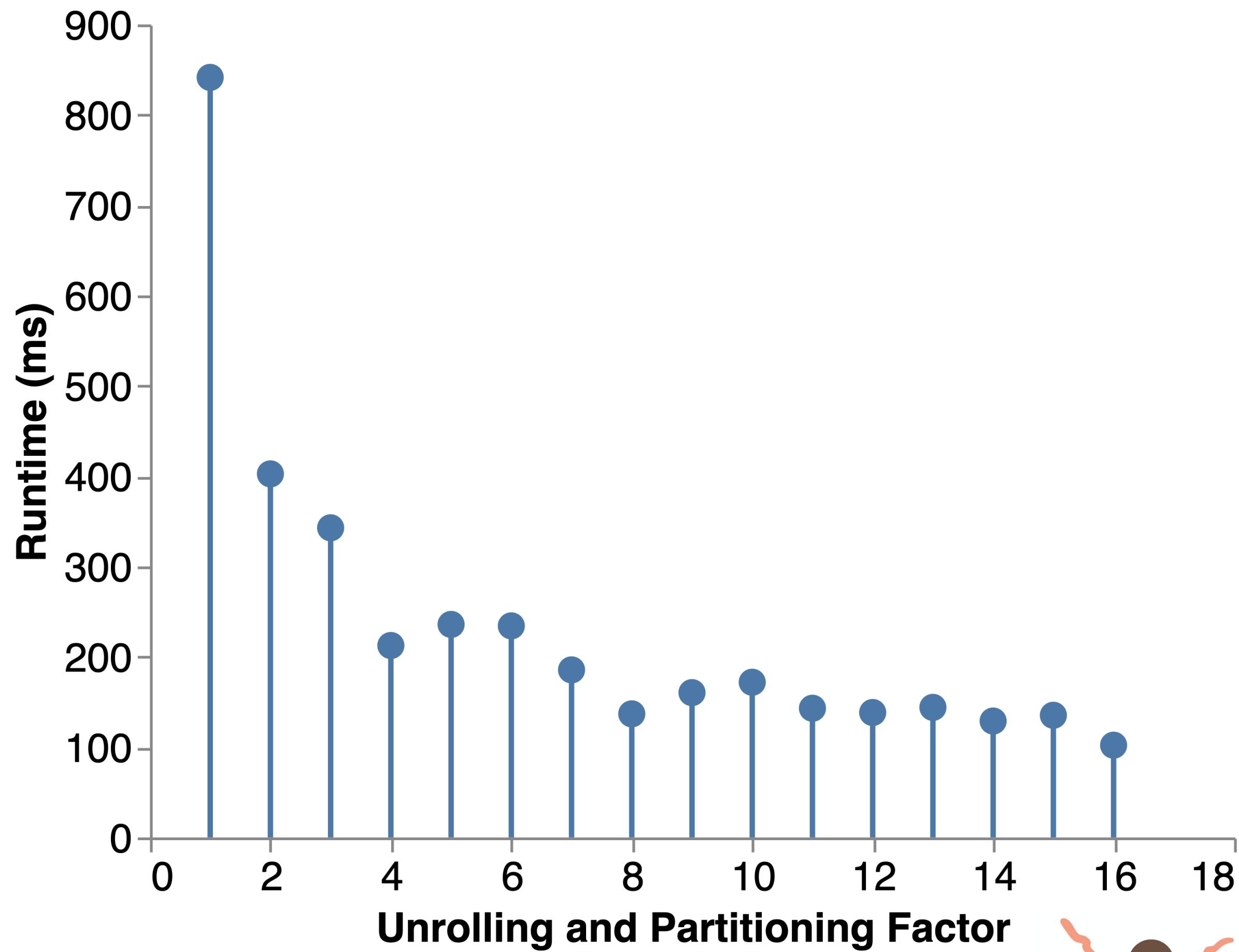
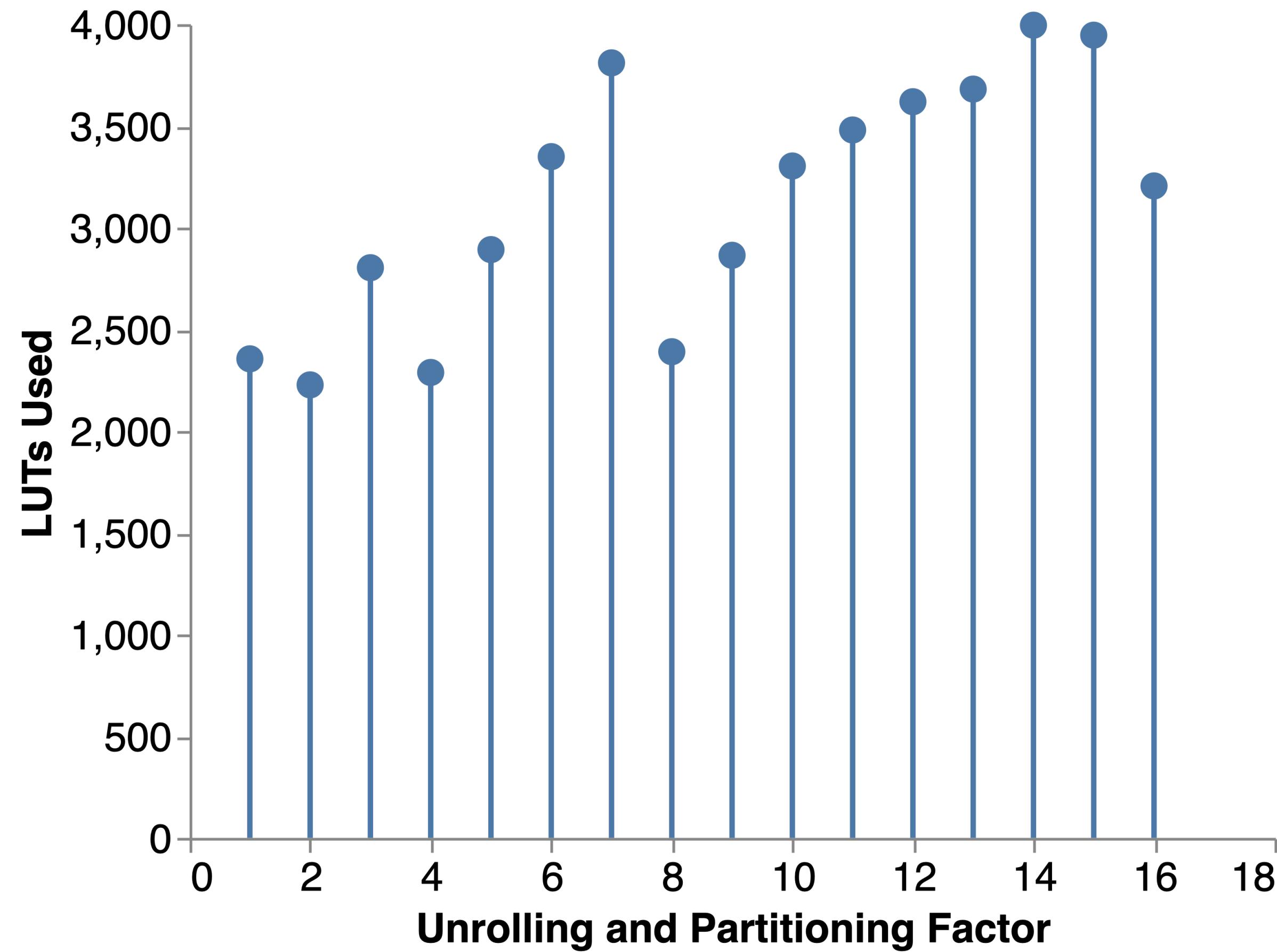
**Mismatched
Partition Sizes**

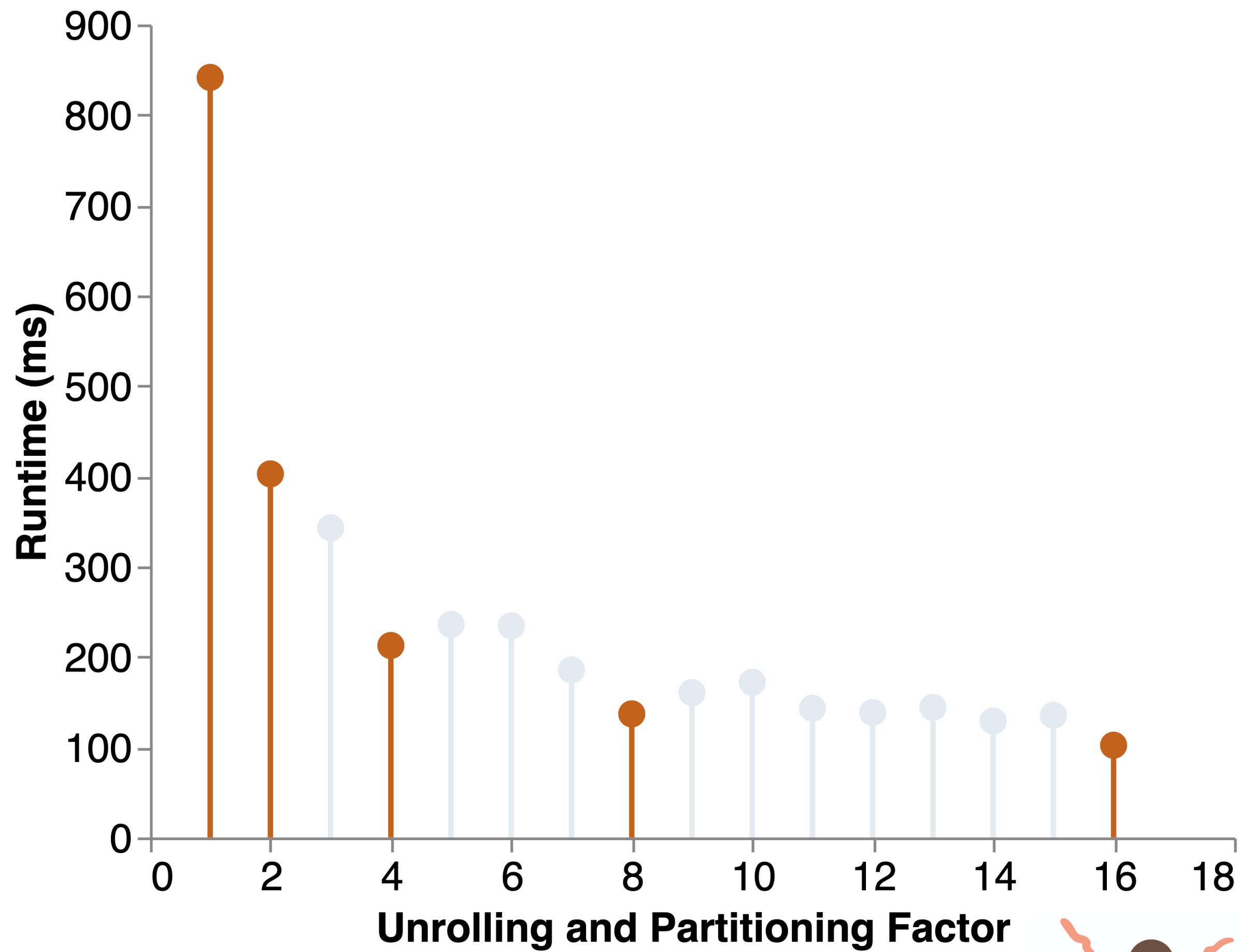
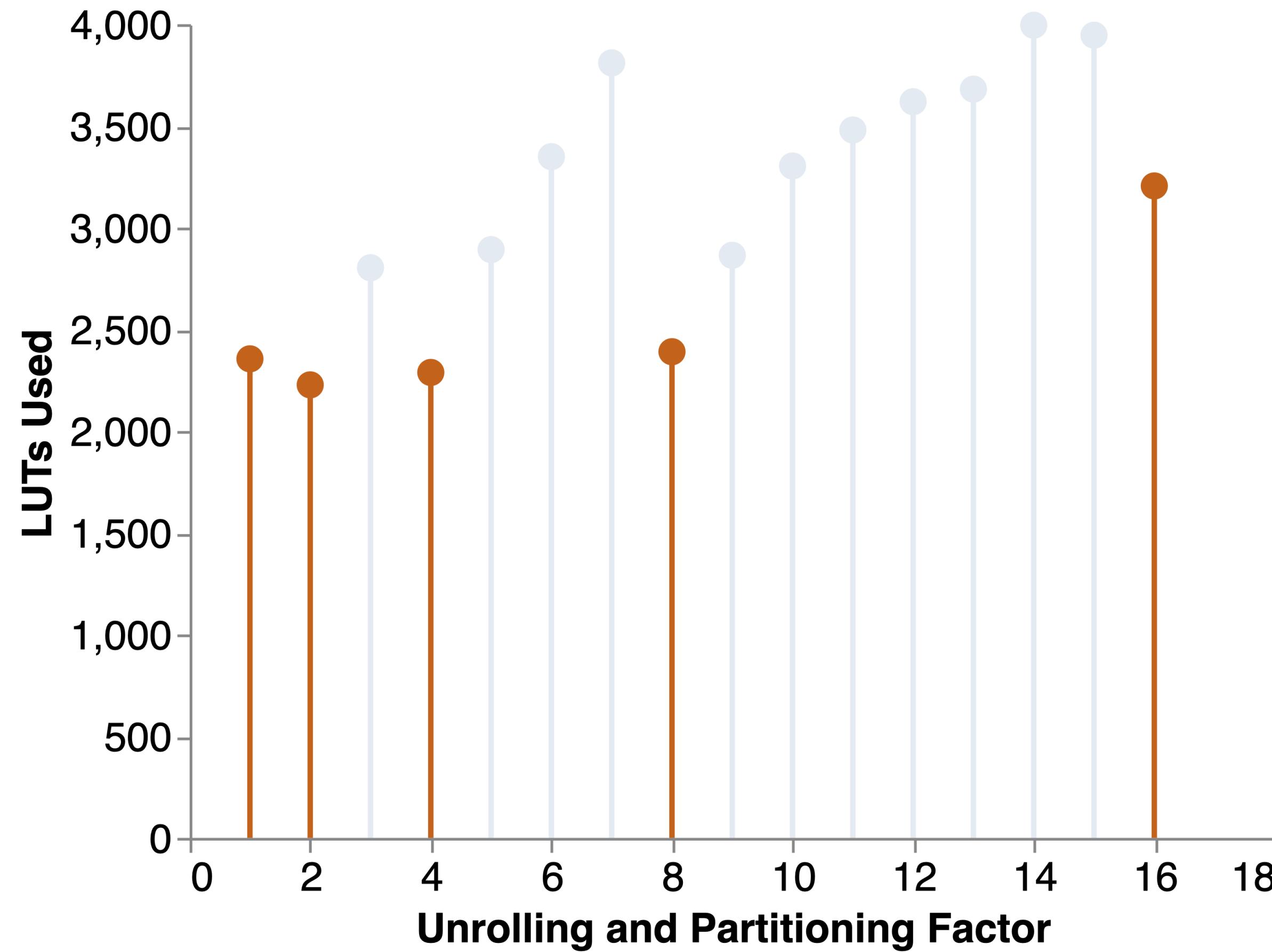
Hardware



Mismatched
Partition Sizes







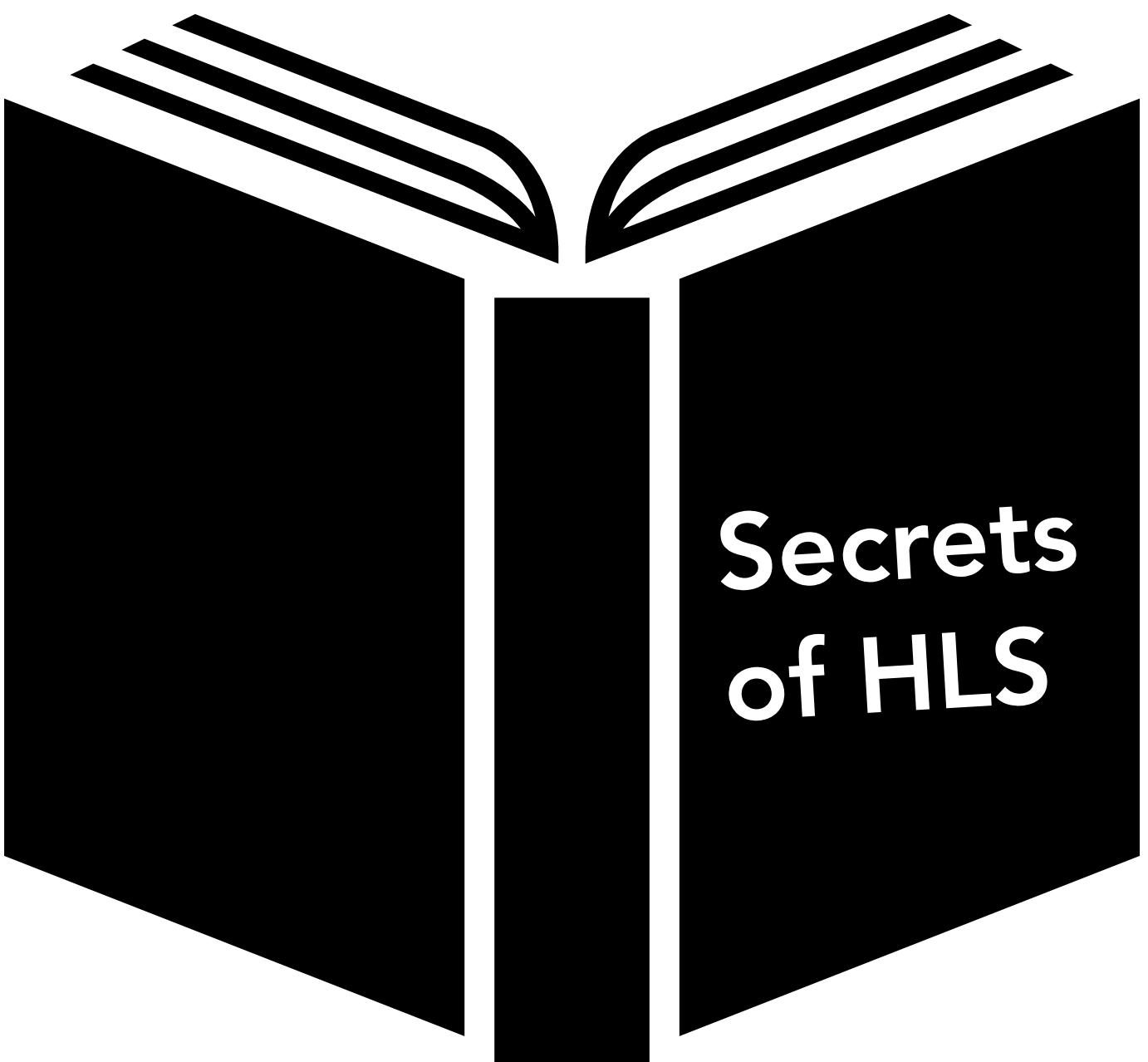
HLS Really Works!

HLS Really Works!

* * * * *

- * when you **unroll** designs
- * when **unrolling** and **partitioning** are aligned
- * when **partitioning** and **memory sizes** are aligned
- * when ports times **partitioning** is a factor of **unrolling**
- * when **memory accesses** are easily analyzable
- * when **reduction patterns** are easily analyzable

HLS Really Works!



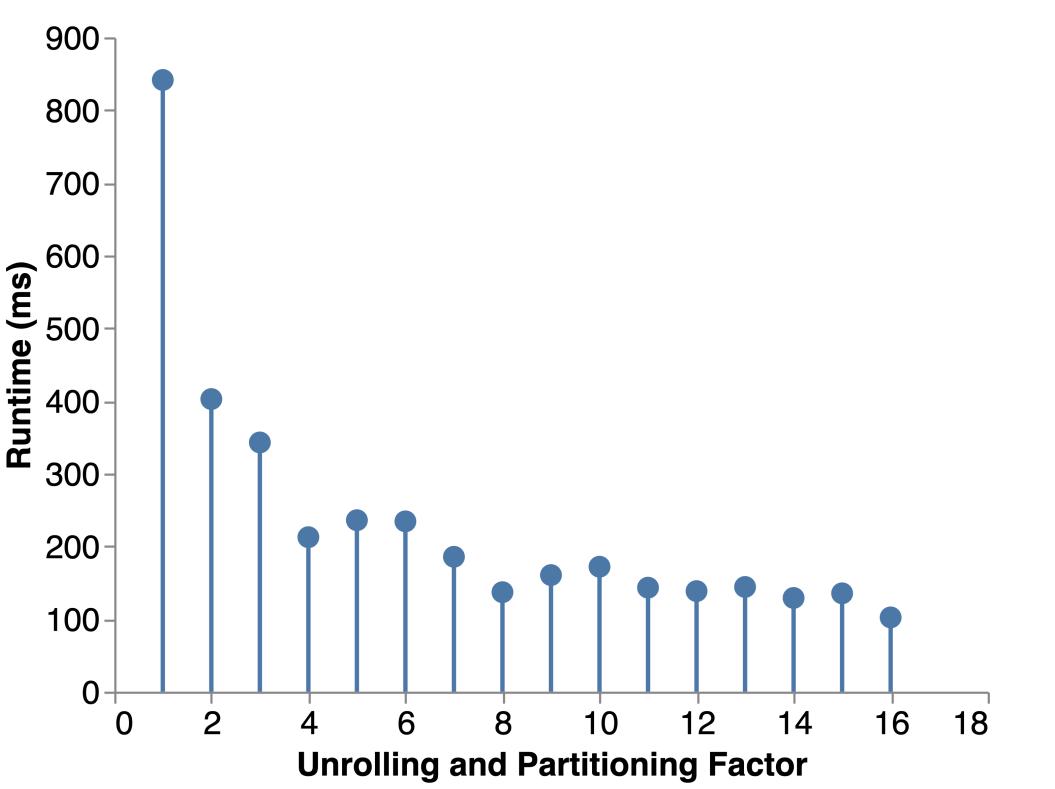
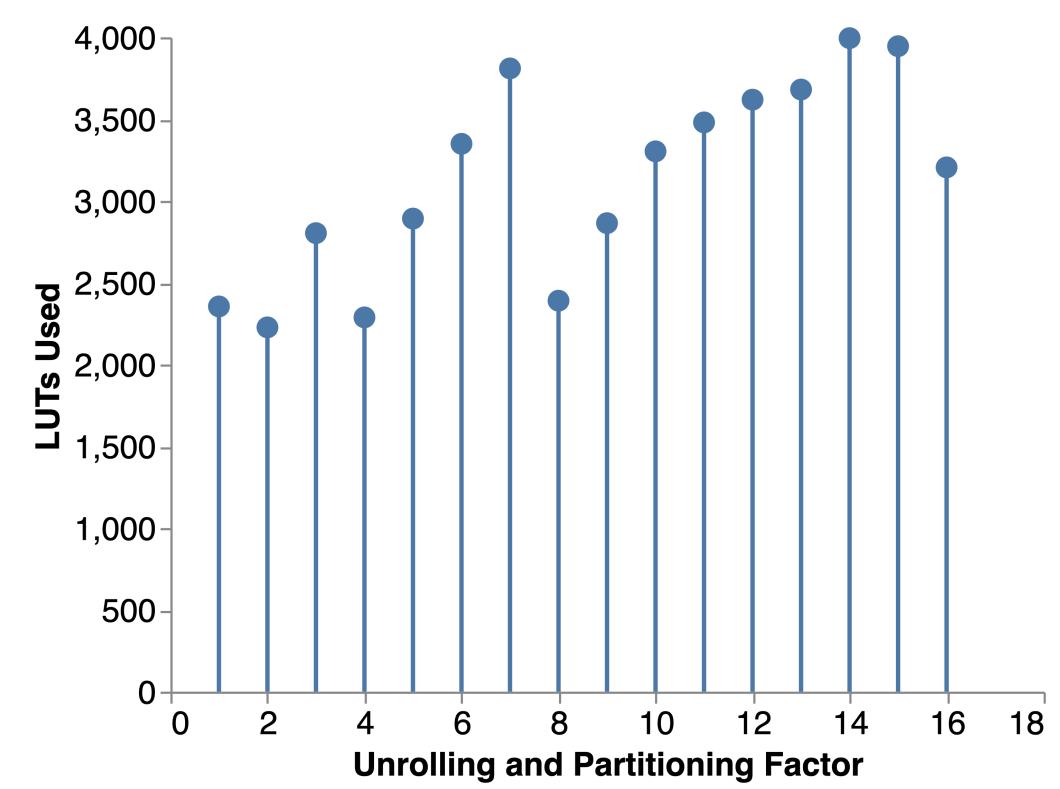
HLS Really Works!

* * * * *

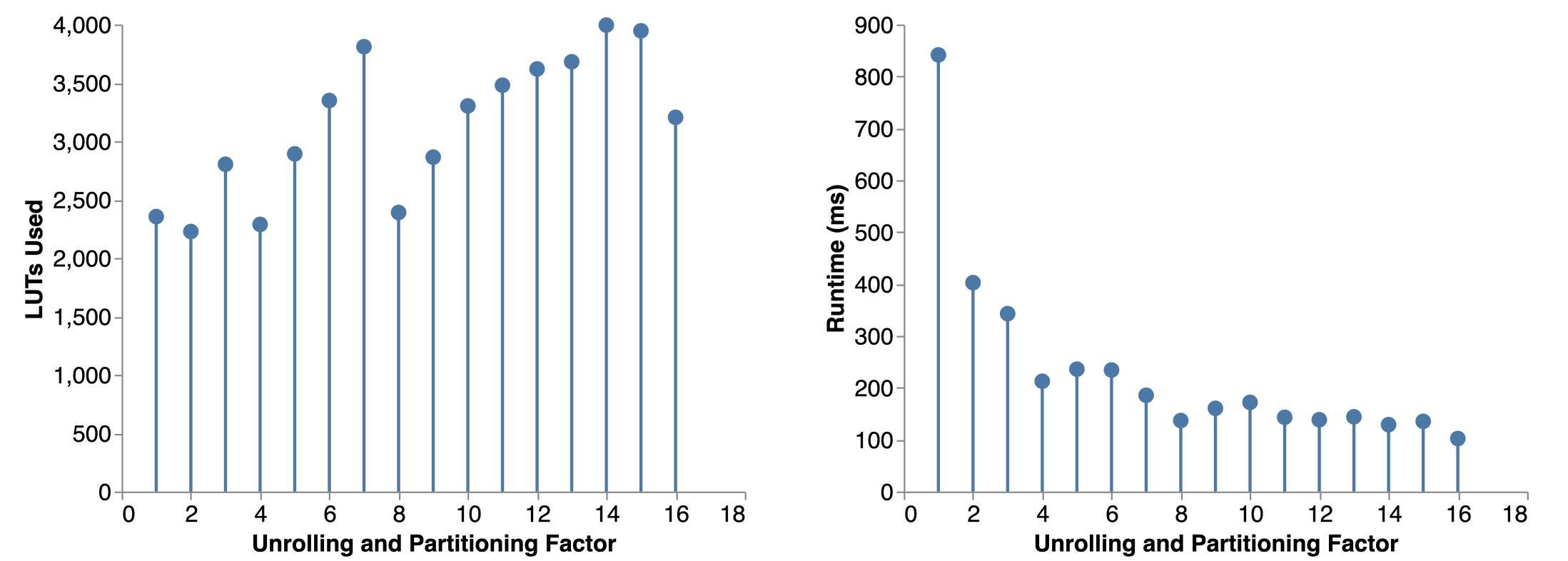
The language doesn't capture
timing and **resource constraints**.



Unpredictable

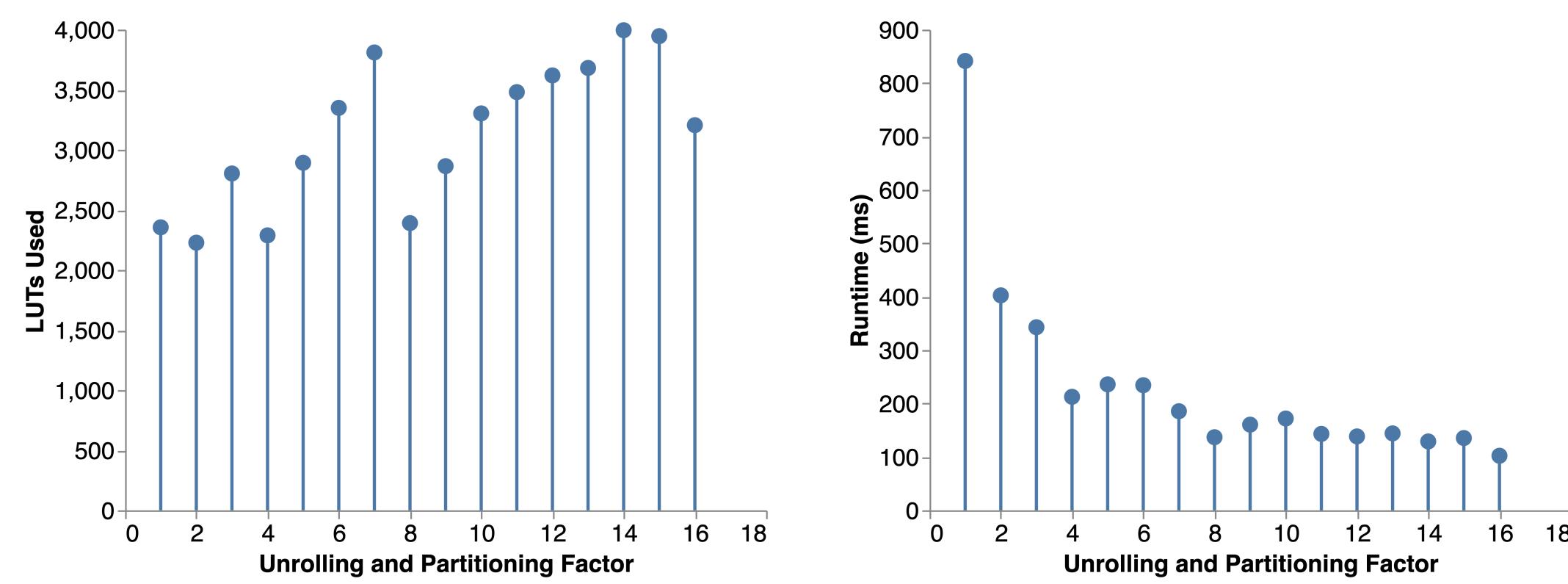


Unpredictable

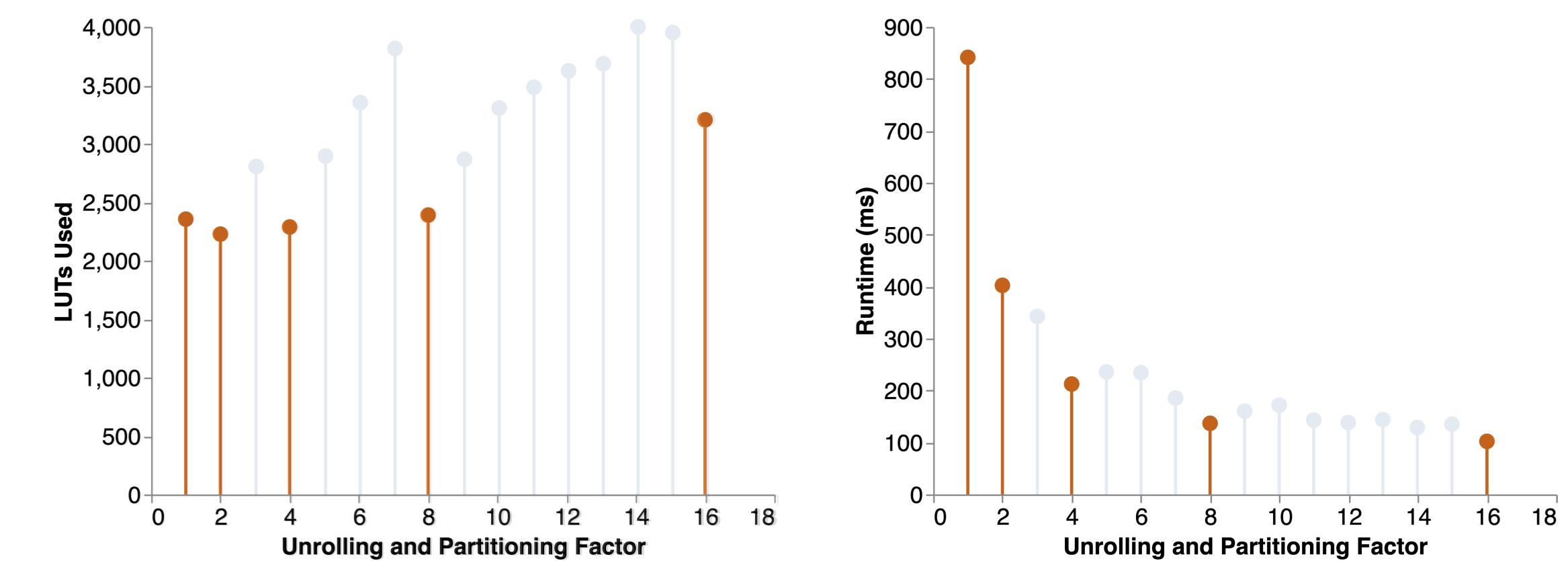


16 designs =
32 hours of
compilation

Unpredictable

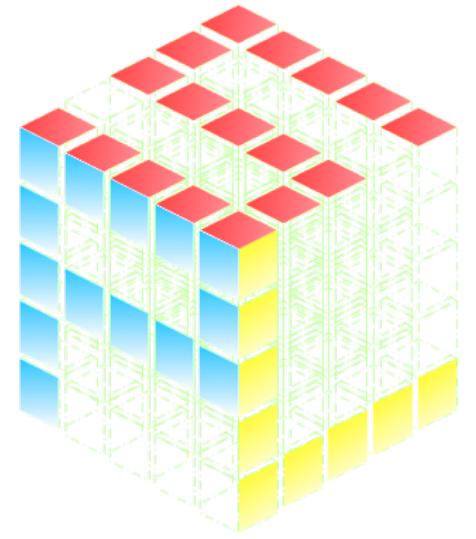


Predictable

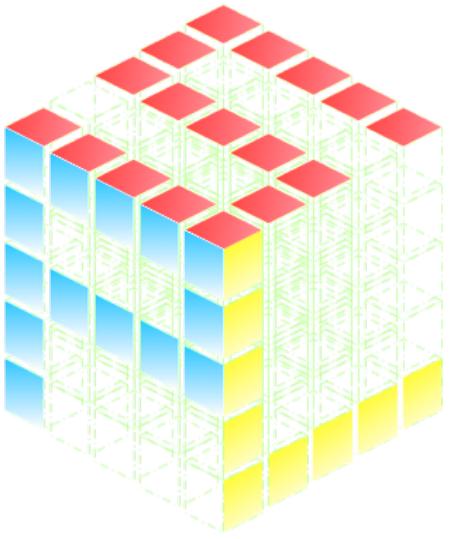


16 designs =
32 hours of
compilation

4 designs =
8 hours of
compilation



Spatial: A Language and Compiler for Application Accelerators

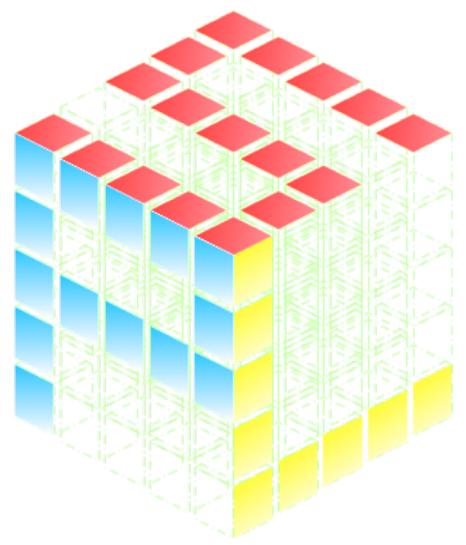


Spatial: A Language and Compiler for Application Accelerators

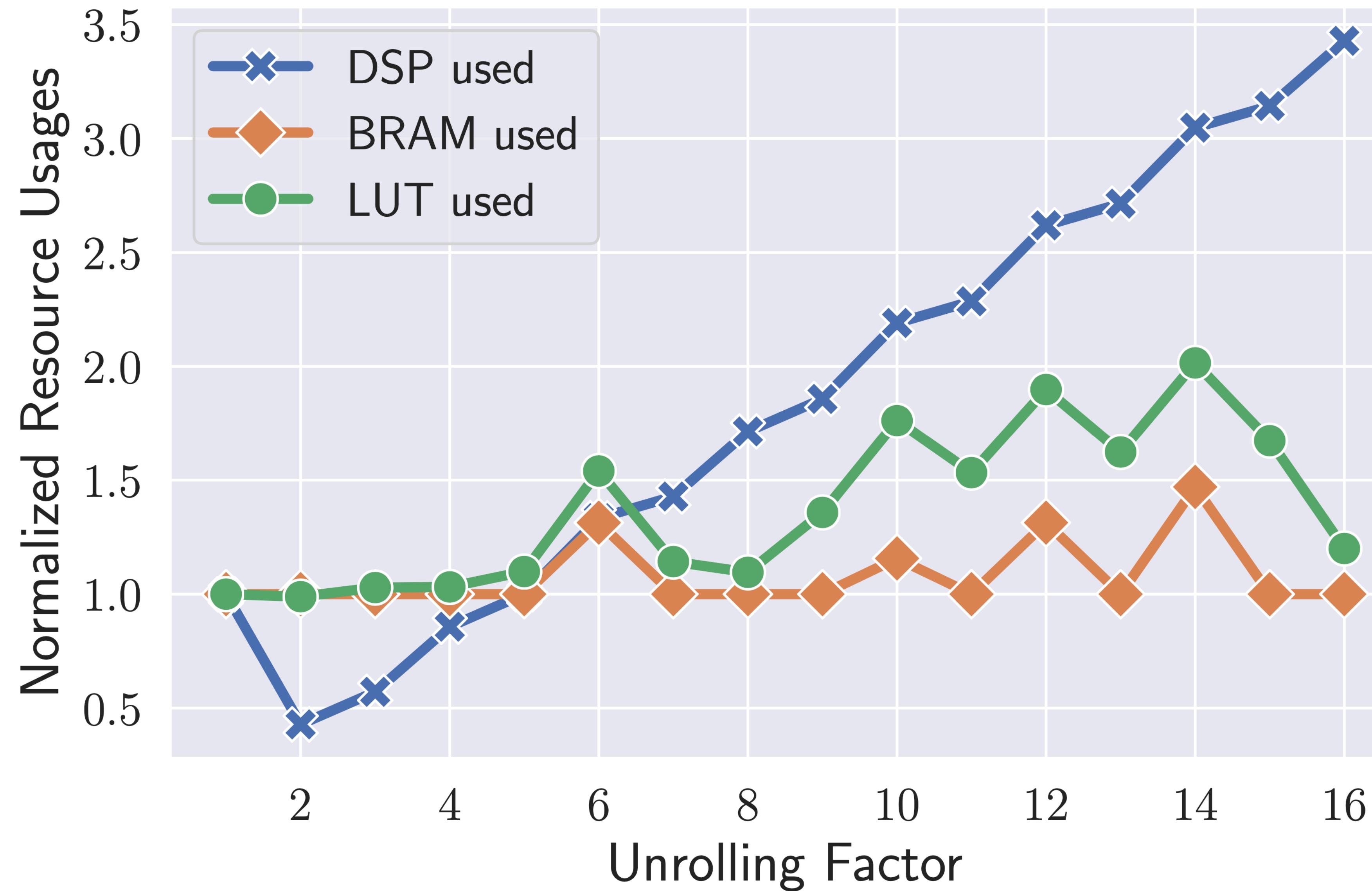
```
1 function GroupAccesses:
2   input:  $A \rightarrow$  set of reads or writes to  $m$ 
3    $G = \emptyset$  set of sets of compatible accesses
4   for all accesses  $a$  in  $A$ :
5     for all sets of accesses  $g$  in  $G$ :
6       if  $IComp(a, a')$  for all  $a'$  in  $g$  then
7         add  $a$  to  $g$ 
8         break
9       else add  $\{a\}$  to  $G$ 
10      return  $G$ 
11    end function
12
13 function ConfigureMemory:
14   input:  $A_r \rightarrow$  set of reads of  $m$ 
15   input:  $A_w \rightarrow$  set of writes to  $m$ 
16    $G_r = \text{GroupAccesses}(A_r)$ 
17    $G_w = \text{GroupAccesses}(A_w)$ 
18    $I = \emptyset$  set of memory instances
19   for all read sets  $R$  in  $G_r$ :
20      $I_r = \{R\}$ 
21      $I_w = \text{ReachingWrites}(G_w, I_r)$ 
22      $i = \text{BankAndBuffer}(I_r, I_w)$ 
23     for each  $inst$  in  $I$ :
24        $I'_r = \text{ReadSets}[inst] + R$ 
25        $I'_w = \text{ReachingWrites}(G_w, I'_r)$ 
26       if  $OComp(A_1, A_2) \ \forall A_1 \neq A_2 \in (G_w \cup I'_r)$  then:
27          $i' = \text{BankAndBuffer}(I'_r, I'_w)$ 
28         if  $\text{Cost}(i') < \text{Cost}(i) + \text{Cost}(inst)$  then:
29           remove  $inst$  from  $I$ 
30           add  $i'$  to  $I$ 
31           break
32         if  $i$  has not been merged then add  $i$  to  $I$ 
33       return  $I$ 
34     end function
```

Figure 7. Banking and buffering algorithm for calculating instances of on-chip memory m .

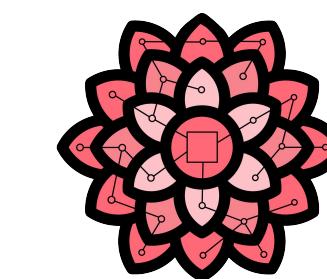
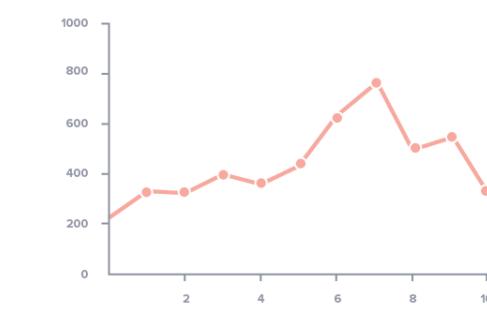
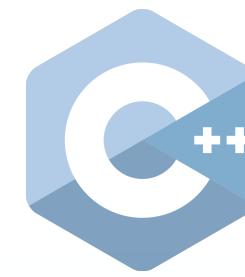
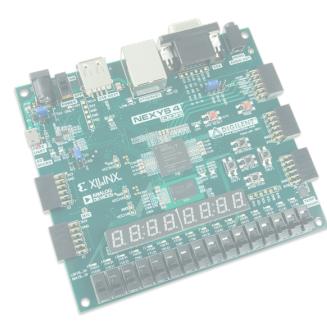
Automatic memory
partitioning



Spatial: A Language and Compiler for Application Accelerators

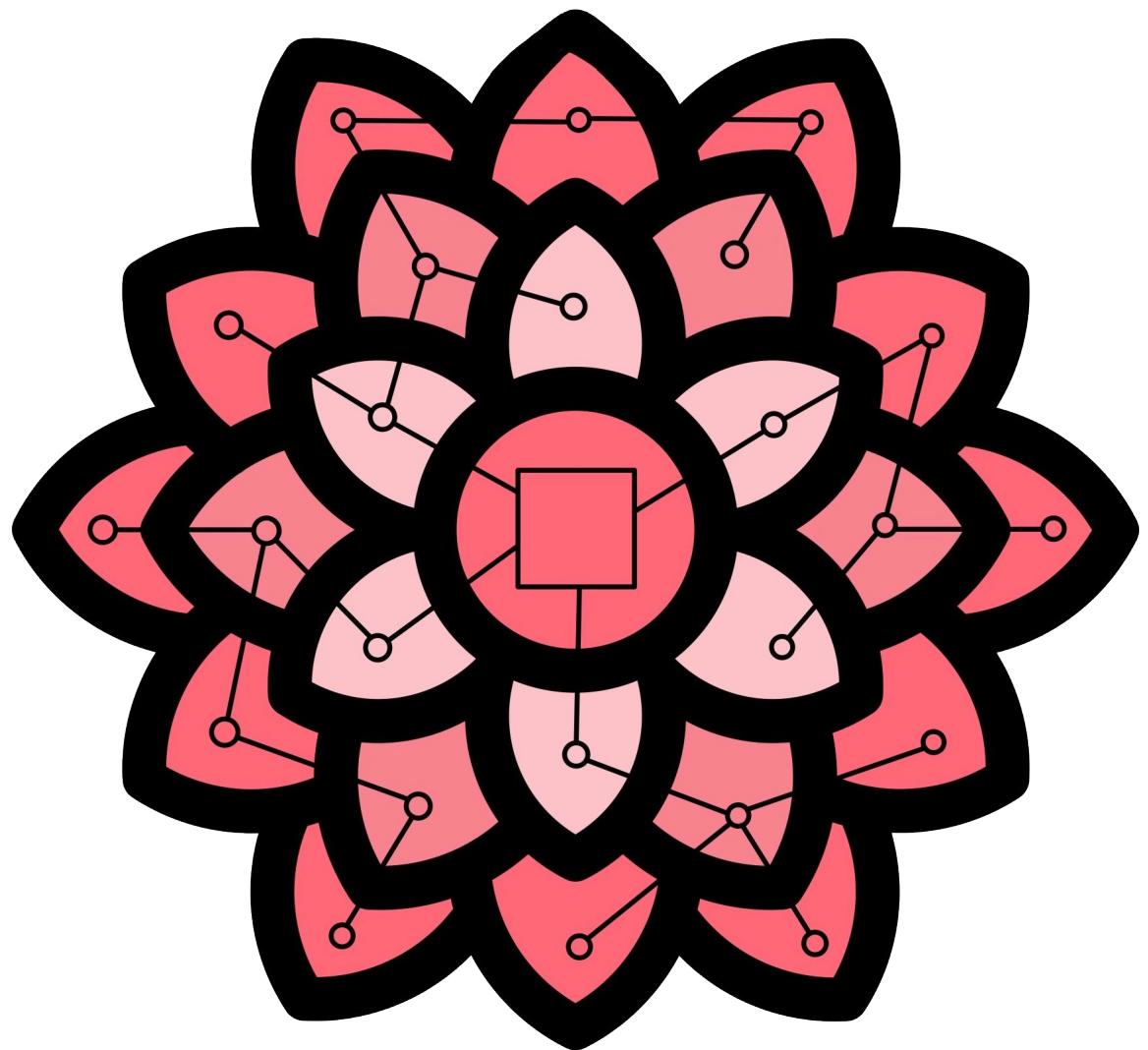


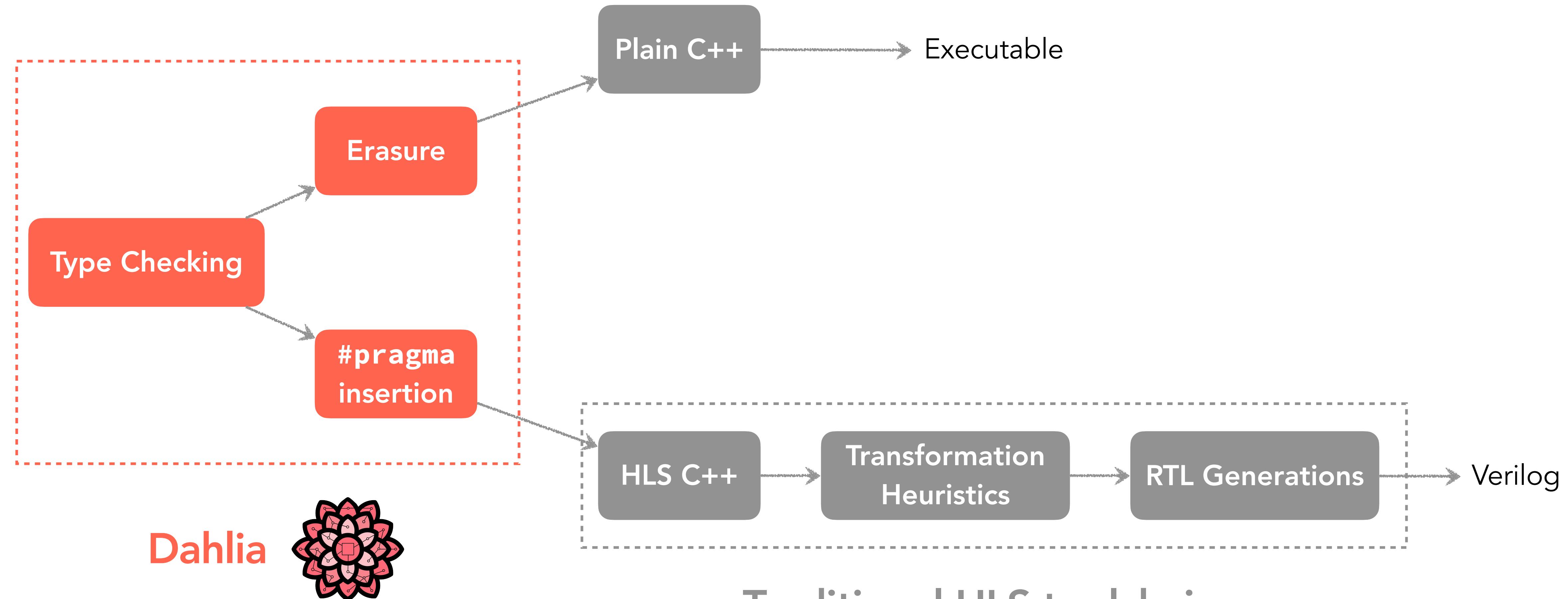
Ada's Journey

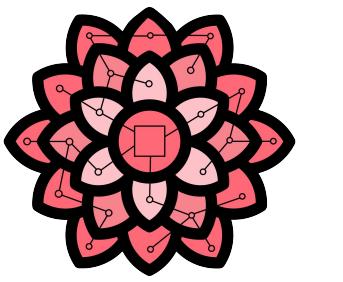


Dahlia

Predictable Accelerator Design

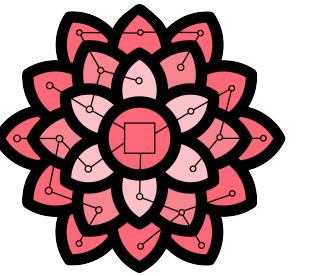






Dahlia

Hardware

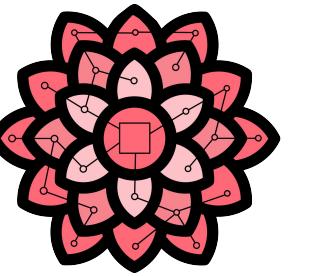


Dahlia

```
let m1: float[10];
```

```
let x = m1[0];
m1[1] := 1;
```

Hardware

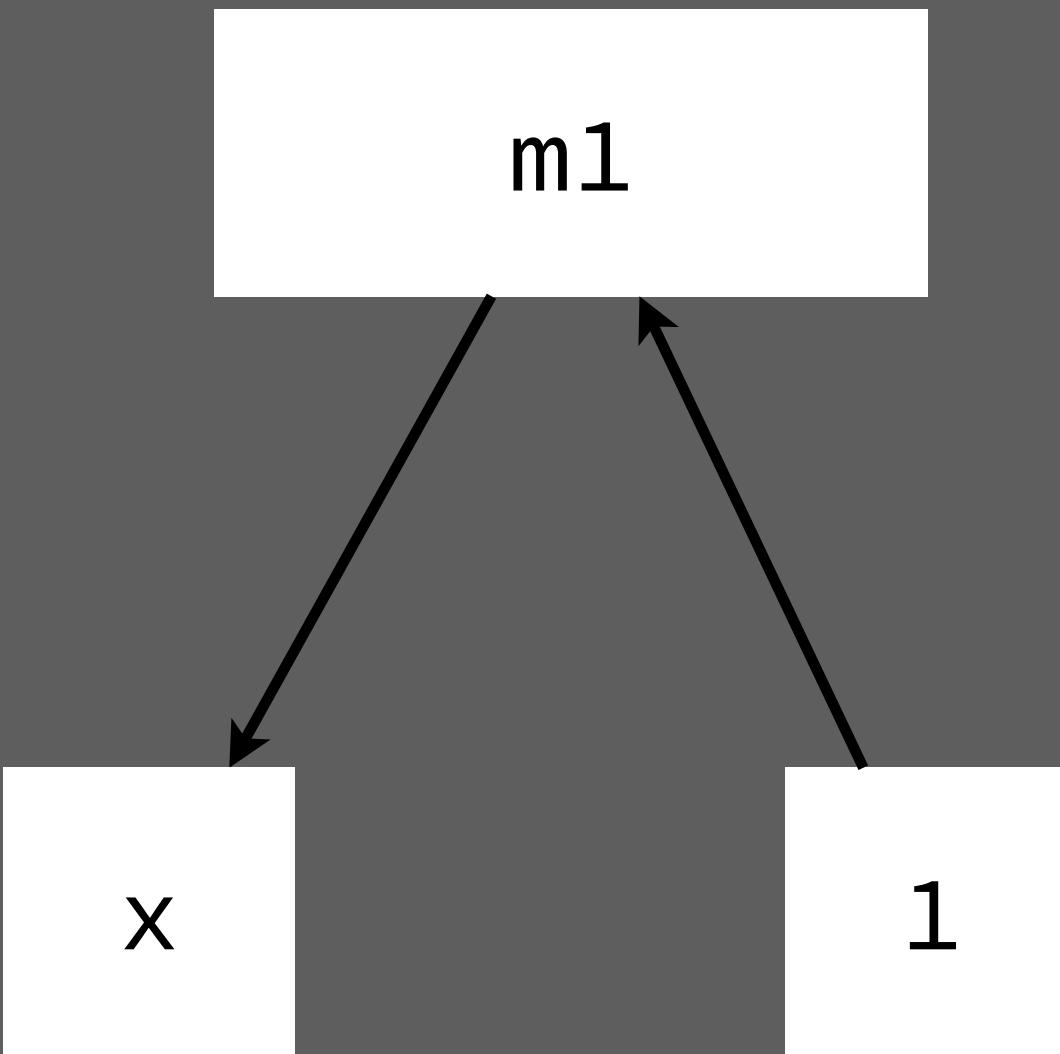


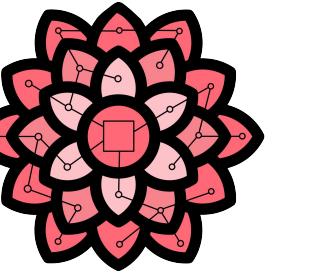
Dahlia

```
let m1: float[10];
```

```
let x = m1[0];  
m1[1] := 1;
```

Hardware



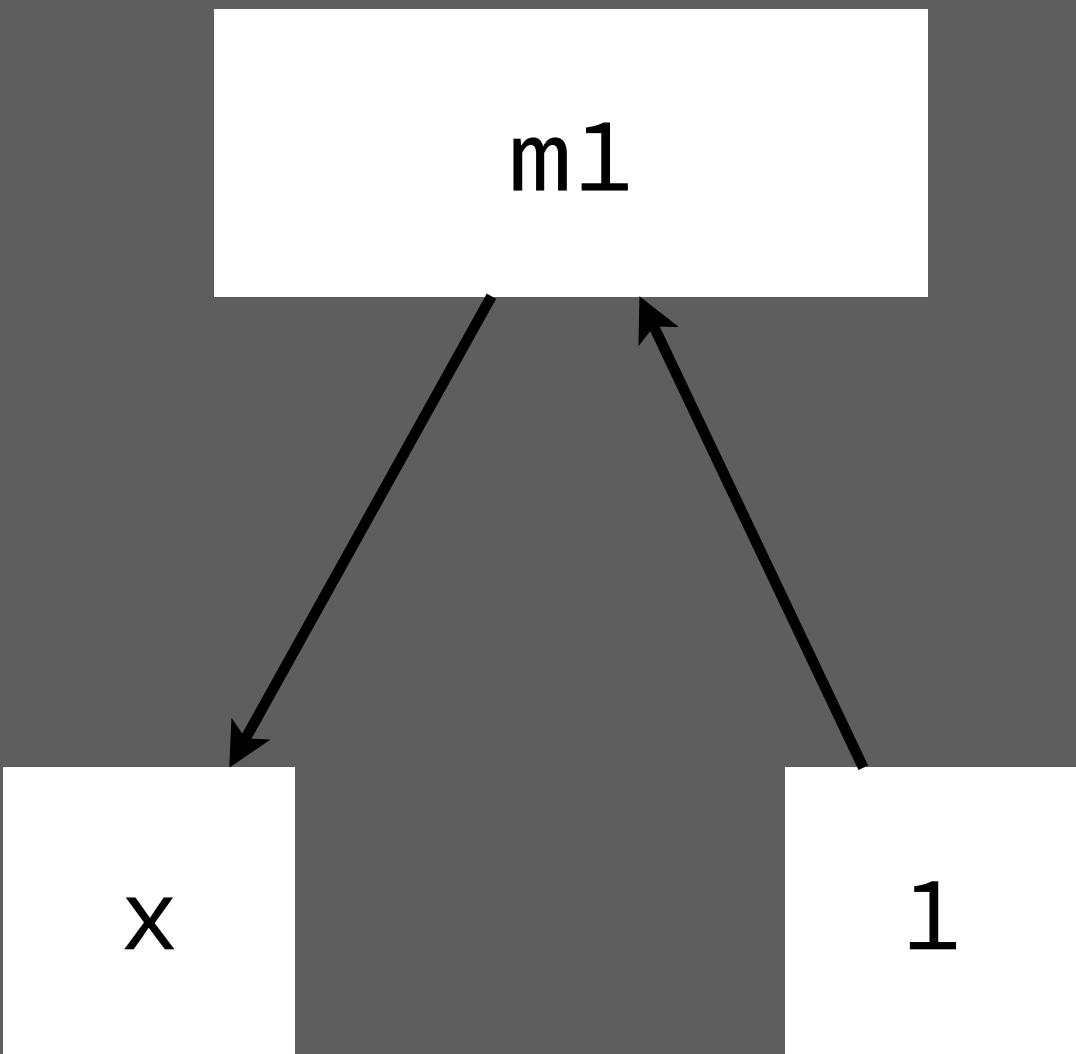


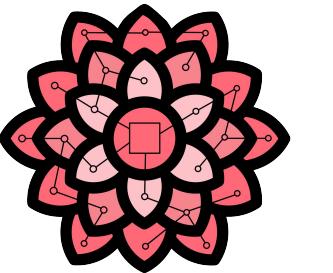
Dahlia

```
let m1: float[10];
```

```
let x = m1[0];  
m1[1] := 1;
```

Hardware





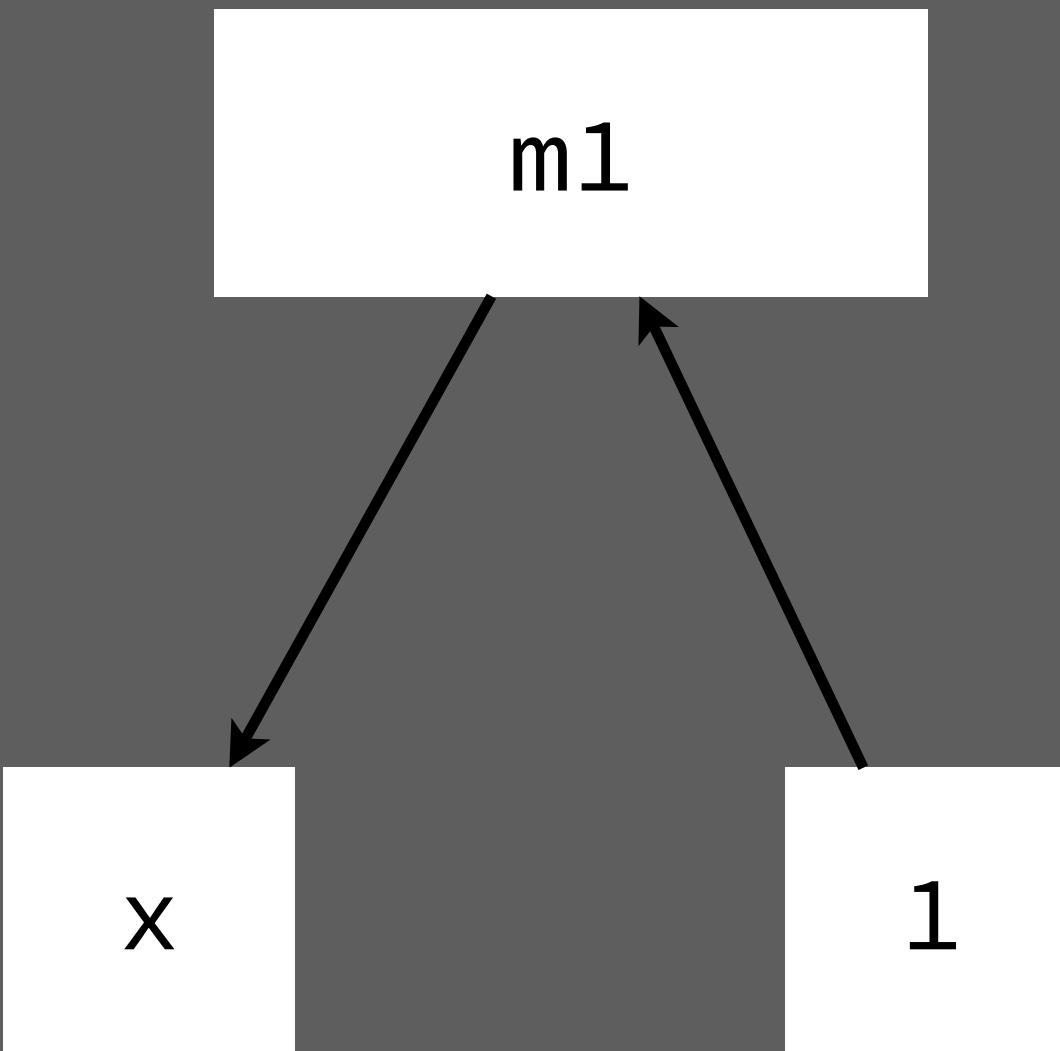
Dahlia

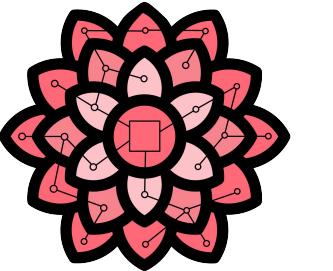
```
let m1: float[10];
```

```
let x = m1[0];  
m1[1] := 1;
```

Do these run
concurrently?

Hardware





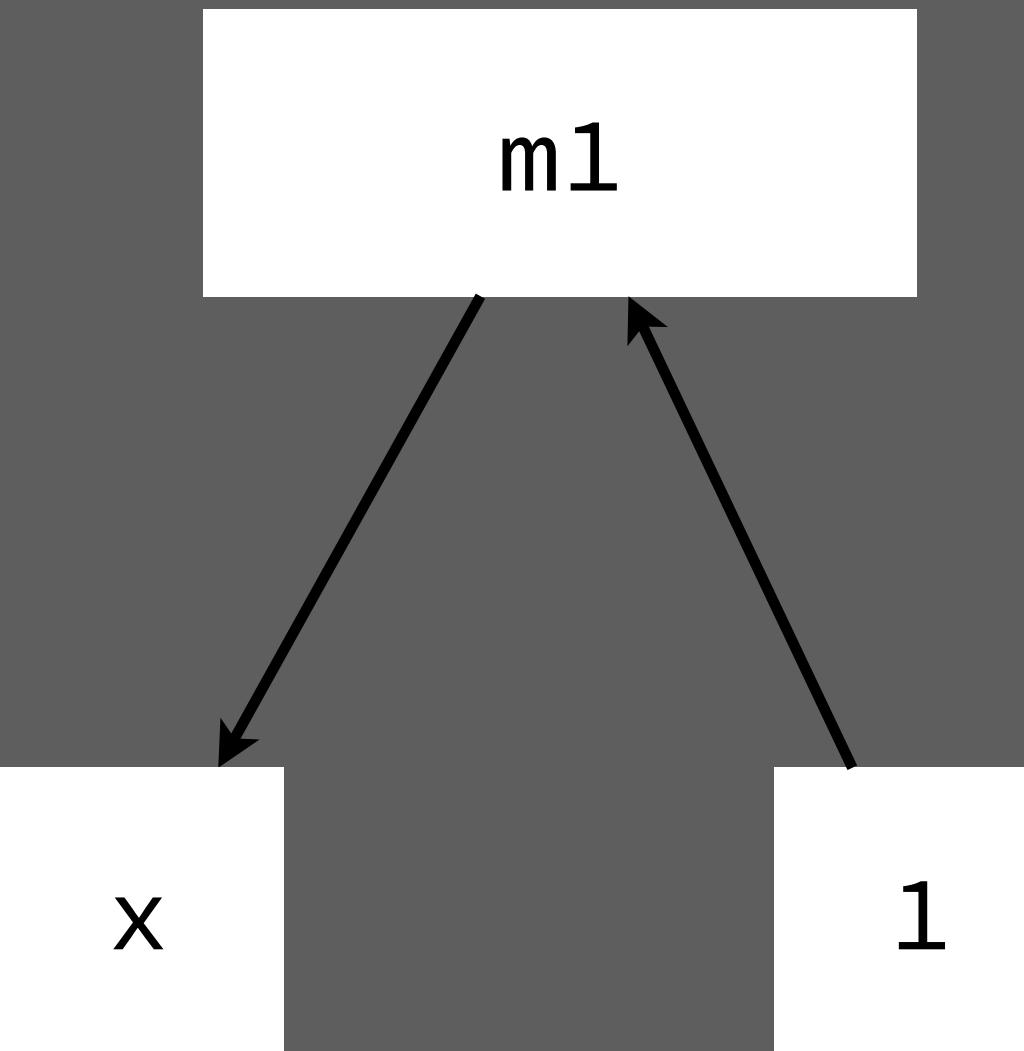
Dahlia

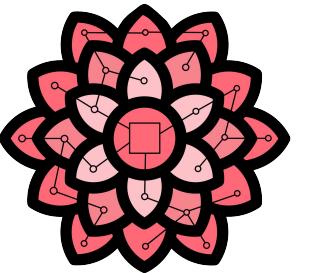
```
let m1: float[10];
```

```
let x = m1[0];  
m1[1] := 1;
```

Do these run
concurrently?

Hardware





Dahlia

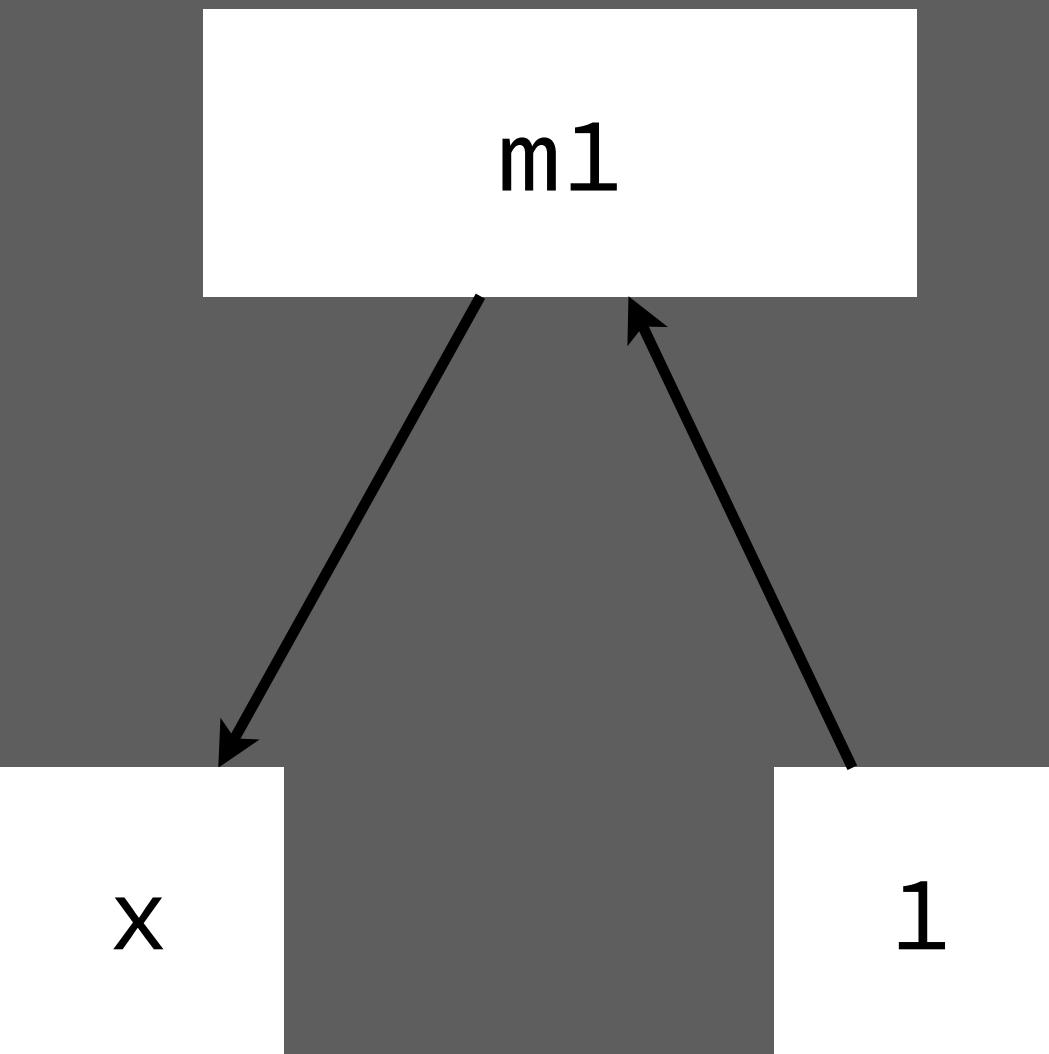
```
let m1: float[10];
```

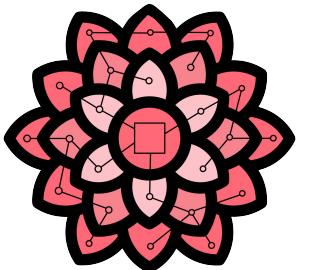
```
let x = m1[0];  
m1[1] := 1;
```

Do these run
concurrently?

Error: Affine resource 'm1'
already used in this context.

Hardware



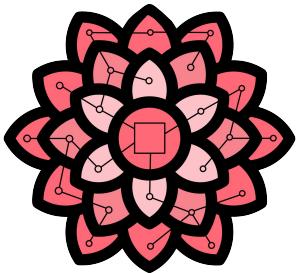


Dahlia

Affine Type System

```
let m1: float[10];
```

```
let x = m1[0];
m1[1] := 1;
```



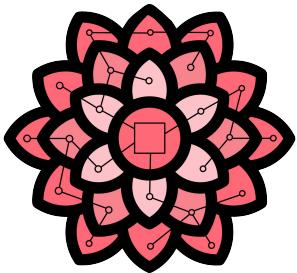
Dahlia

Affine Type System

```
let m1: float[10];
```

```
let x = m1[0];  
m1[1] := 1;
```

At most one **use** of a
variable



Dahlia

Affine Type System

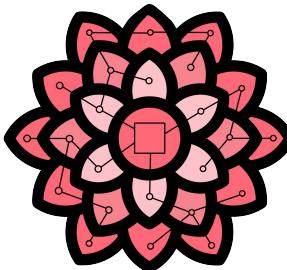
Declaration



```
let m1: float[10];
```

```
let x = m1[0];  
m1[1] := 1;
```

At most one **use** of a
variable



Dahlia

Affine Type System

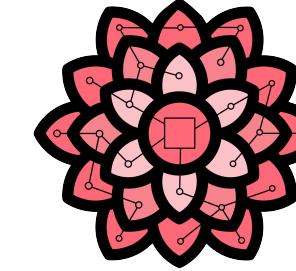
Declaration

```
let m1: float[10];
```

```
let x = m1[0];  
m1[1] := 1;
```

Use

At most one **use** of a variable



Dahlia

Affine Type System

Declaration

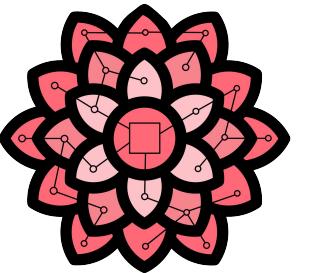
```
let m1: float[10];
```

```
let x = m1[0];  
m1[1] := 1;
```

Use

At most one **use** of a variable

Error: Affine resource 'm1'
already used in this context.



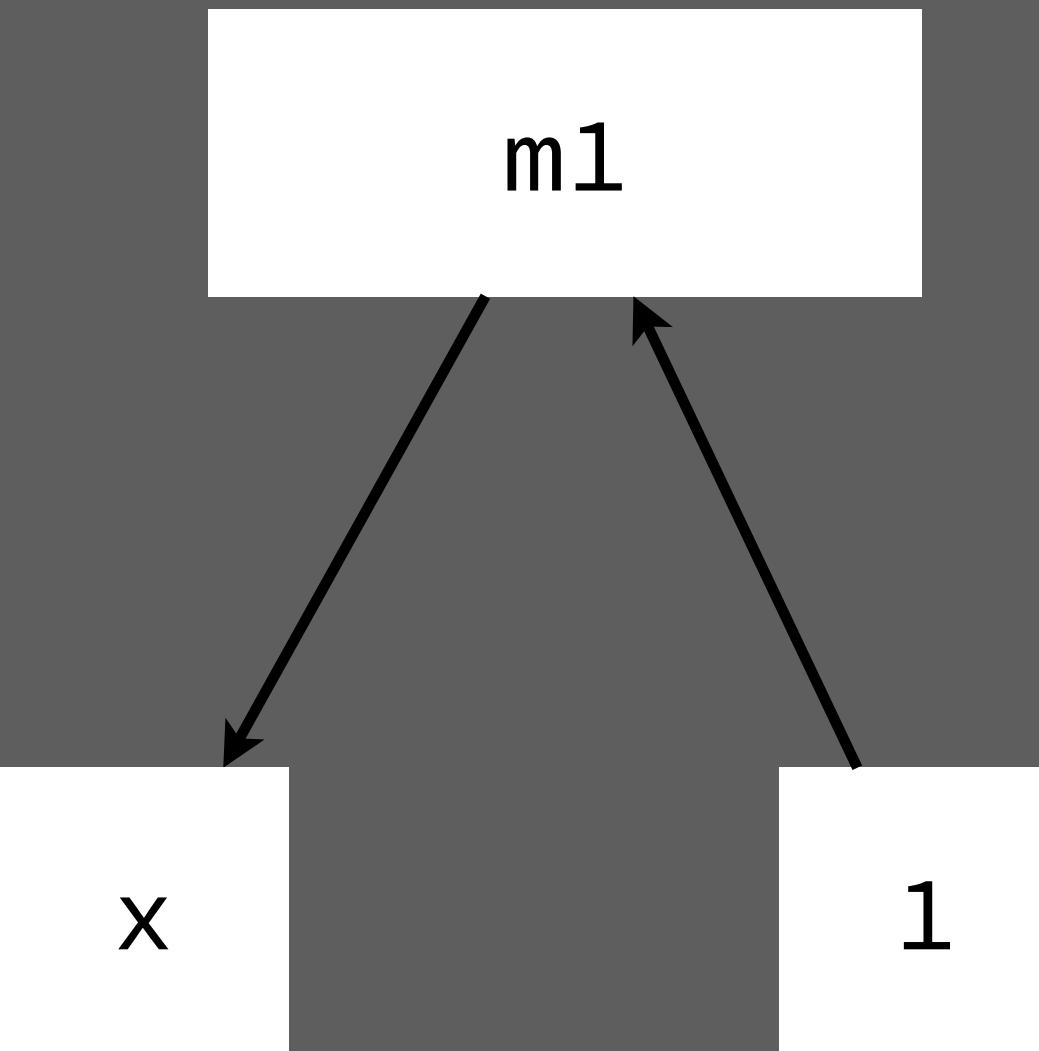
Dahlia

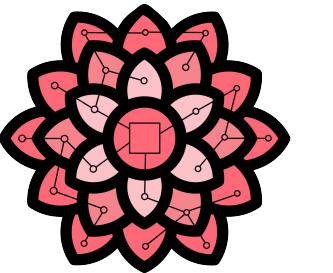
```
let m1: float[10];
```

```
let x = m1[0];  
m1[1] := 1;
```

Do these run
concurrently?

Hardware





Dahlia

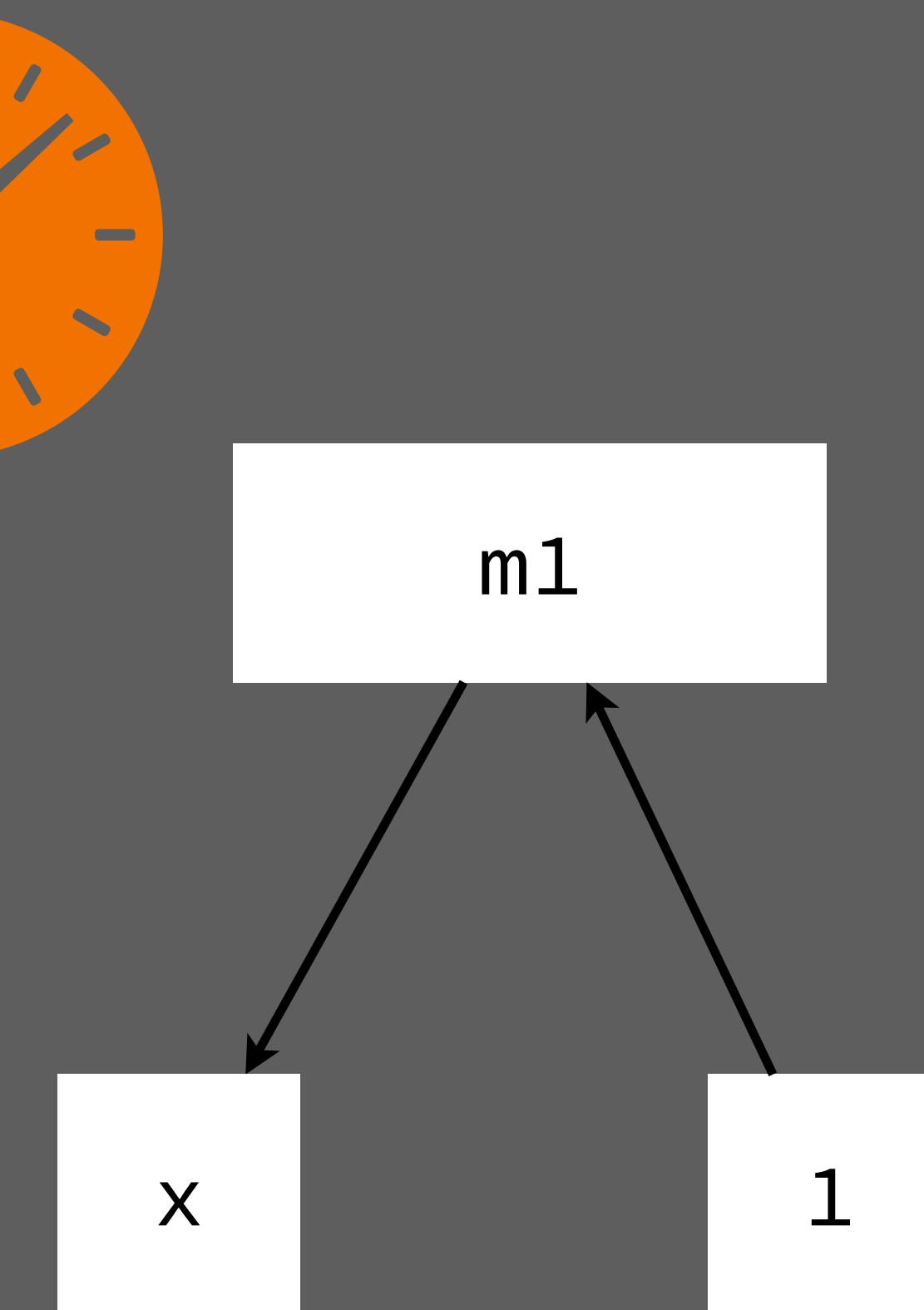
```
let m1: float[10];
```

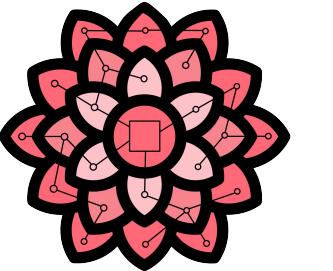
```
let x = m1[0];  
m1[1] := 1;
```

Unordered Composition

Do these run
concurrently?

Hardware





Dahlia

```
let m1: float[10];
```

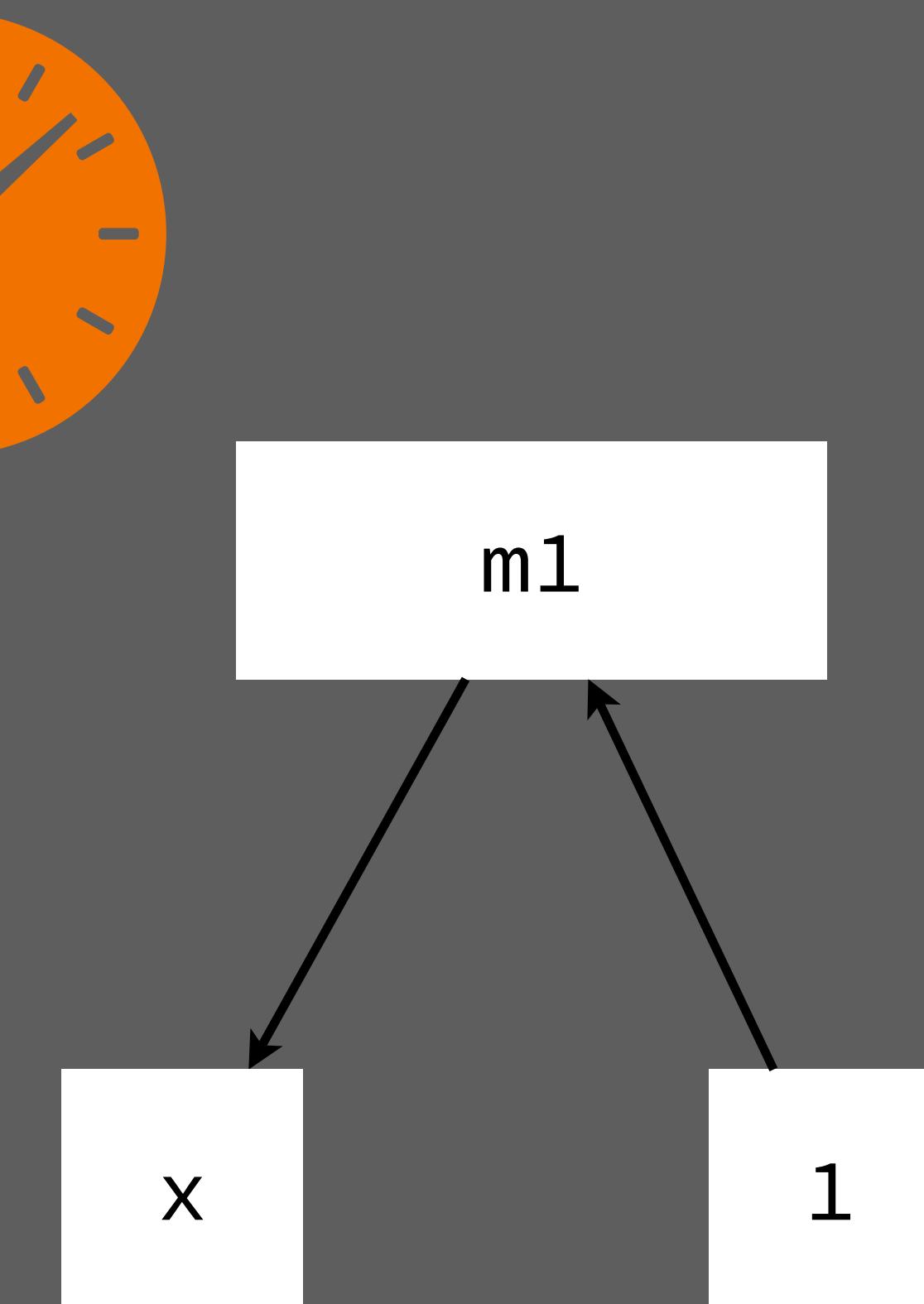
```
let x = m1[0];  
m1[1] := 1;
```

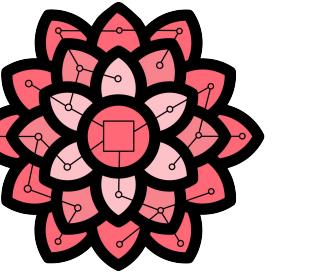
Do these run
concurrently?

Unordered Composition

- Run in parallel respecting data flow.
- Consume affine resources.

Hardware





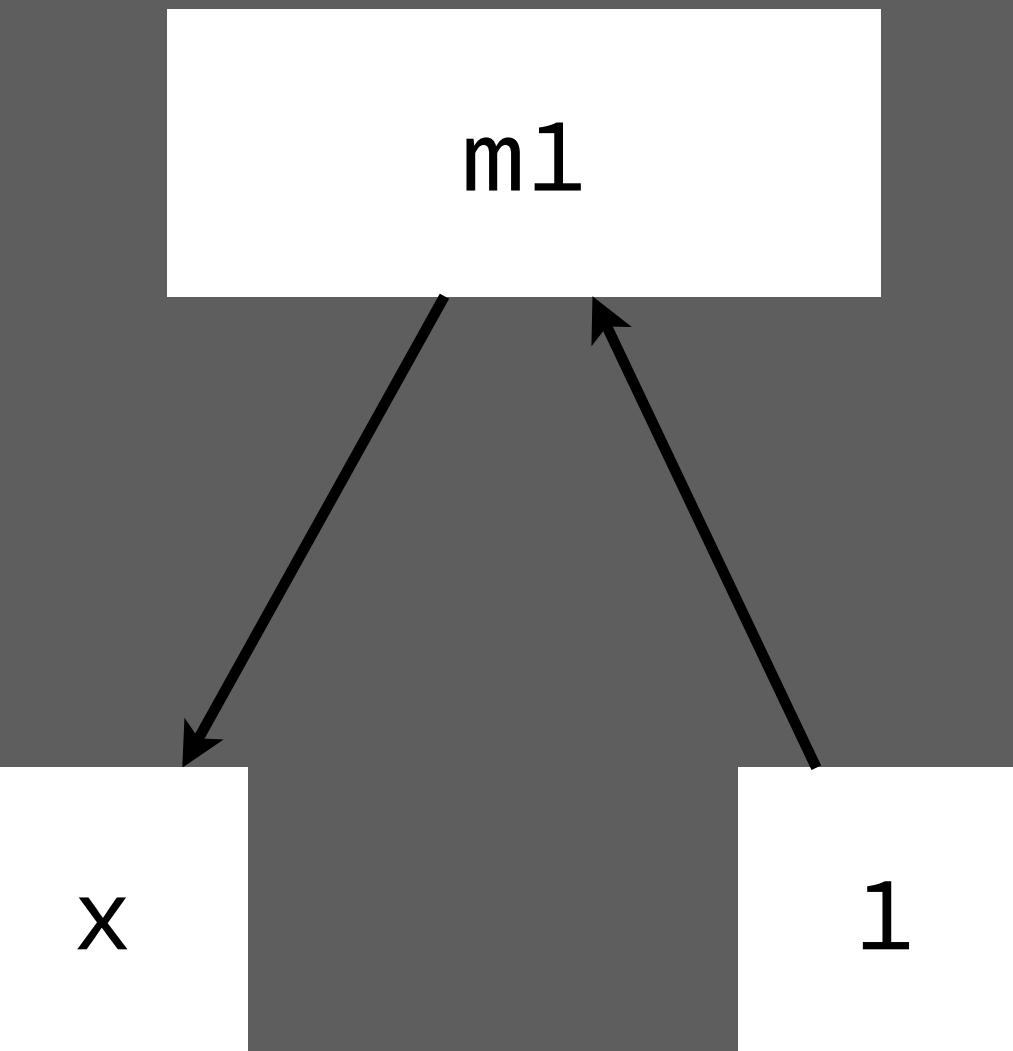
Dahlia

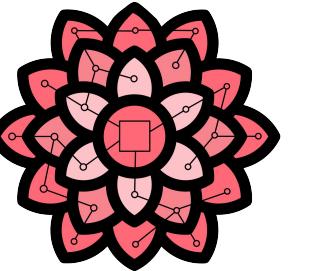
```
let m1: float[10];
```

```
let x = m1[0];
```

```
m1[1] := 1;
```

Hardware





Dahlia

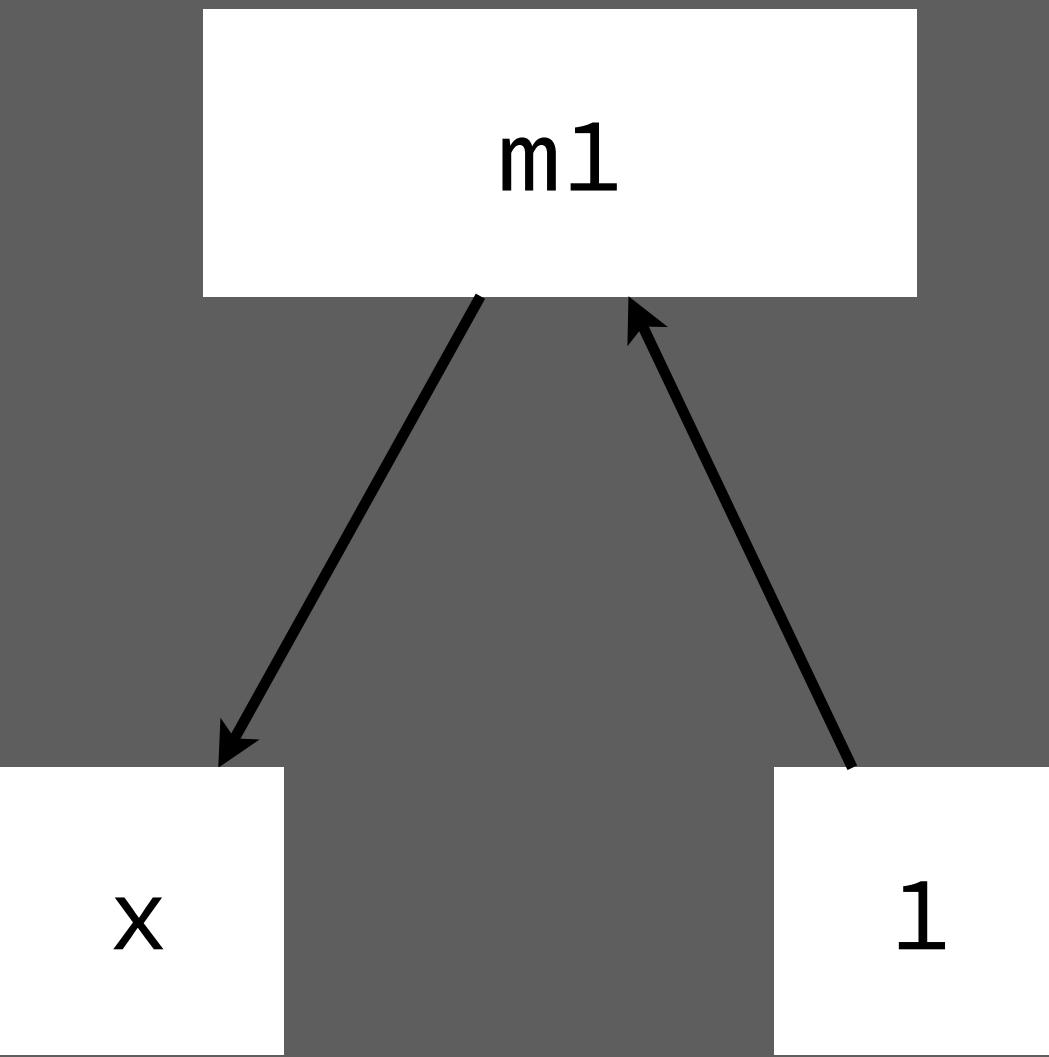
```
let m1: float[10];
```

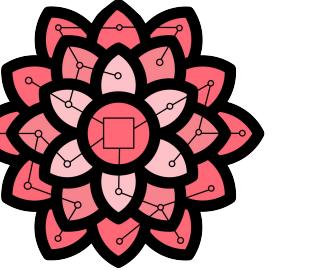
```
let x = m1[0];
```

```
m1[1] := 1;
```

Ordered Composition

Hardware





Dahlia

```
let m1: float[10];
```

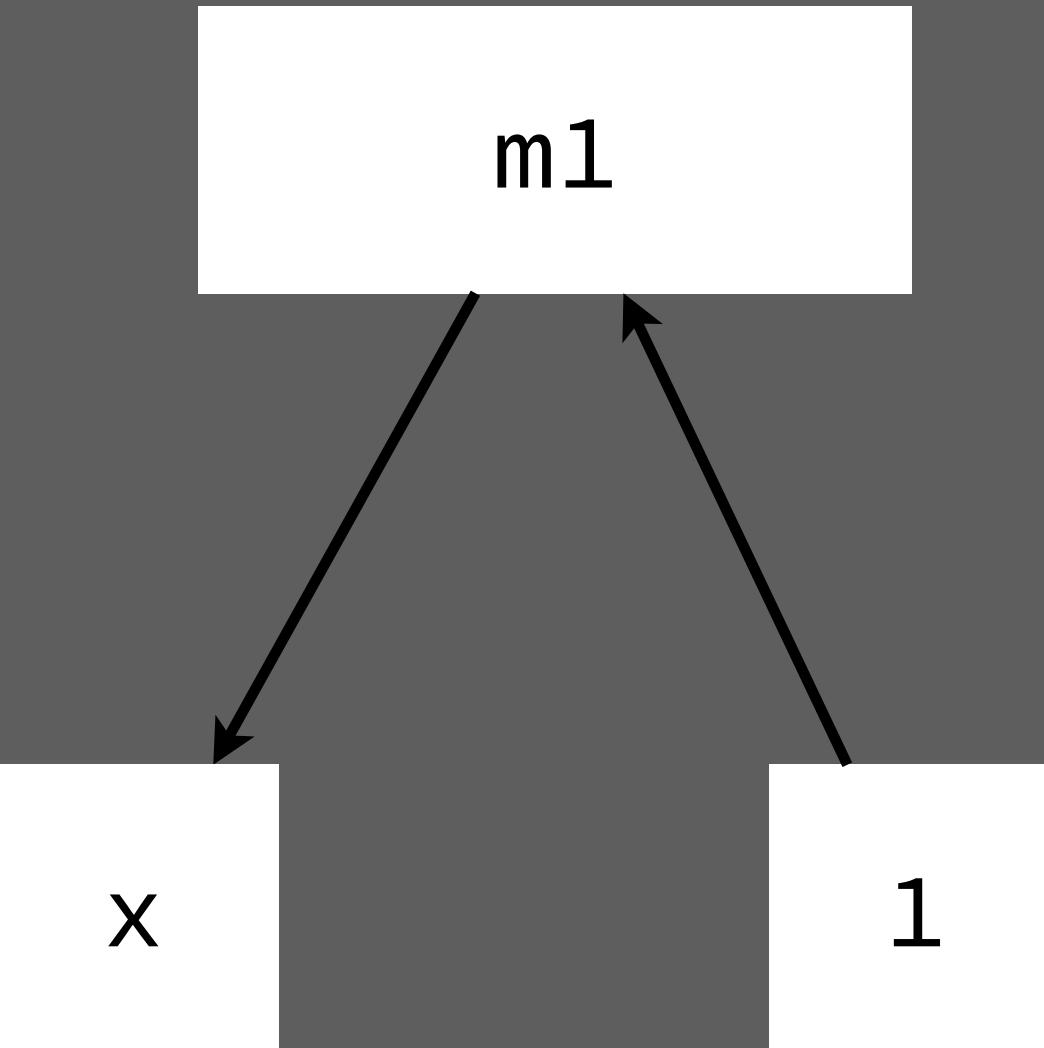
```
let x = m1[0];
```

```
m1[1] := 1;
```

Ordered Composition

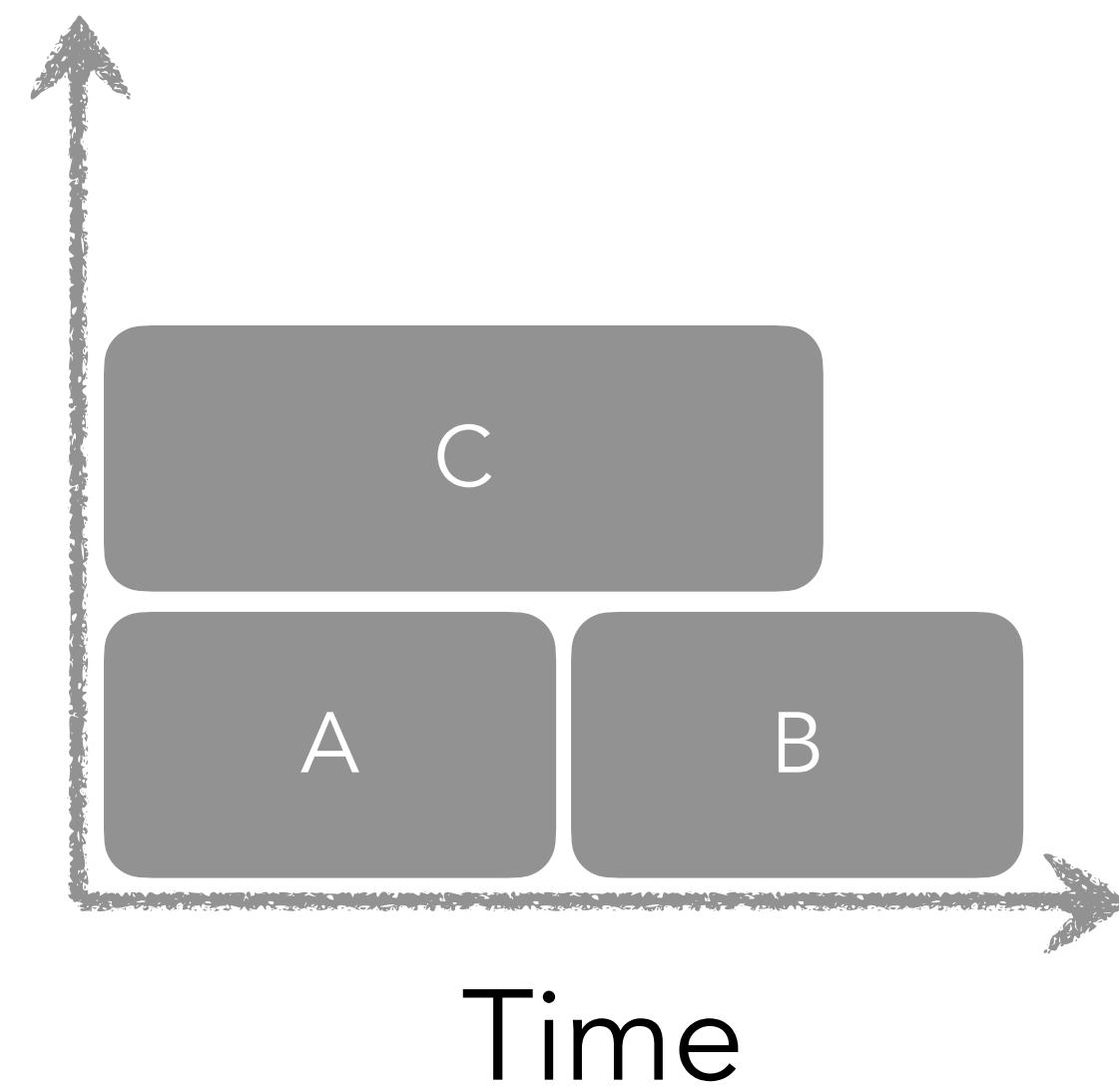
- Run **sequentially**.
- Renew **affine** resources.

Hardware



{ A --- B } ;

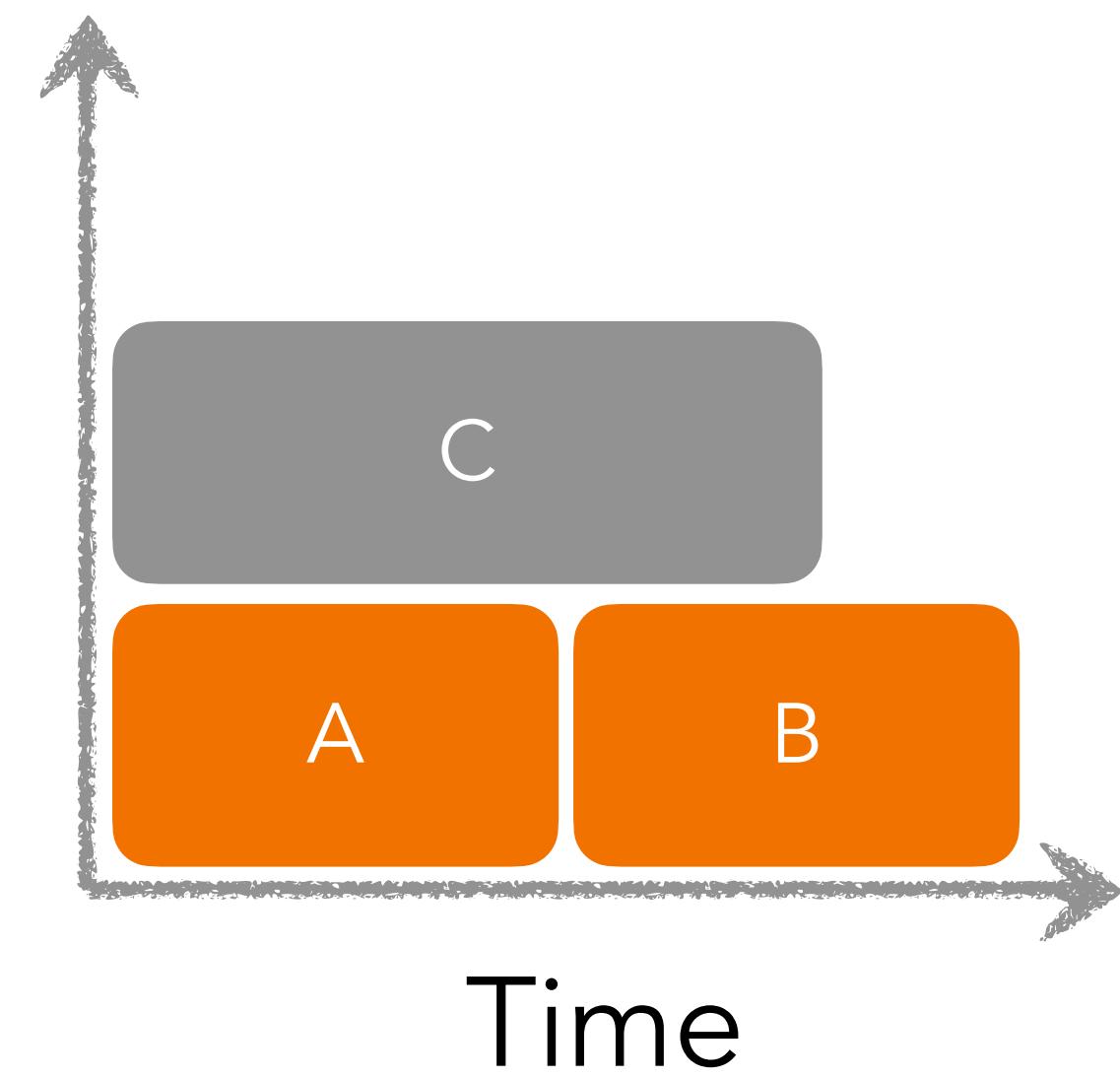
C



Temporally exclusive use of resources

{ A --- B } ;

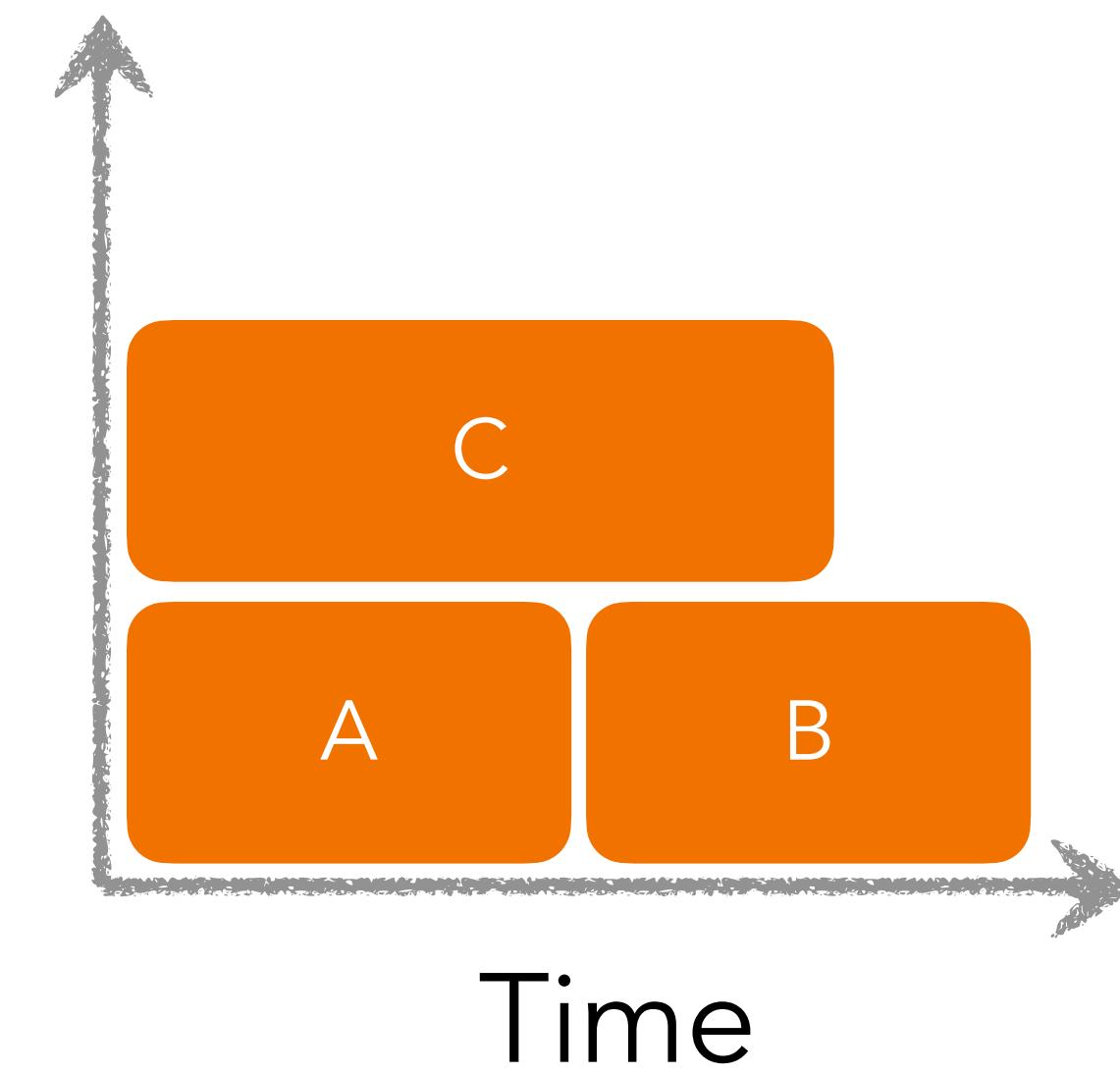
C



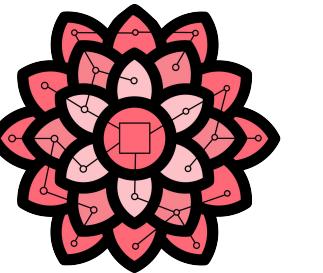
Temporally exclusive use of resources

{ A --- B } ;

C



Temporally exclusive use of resources

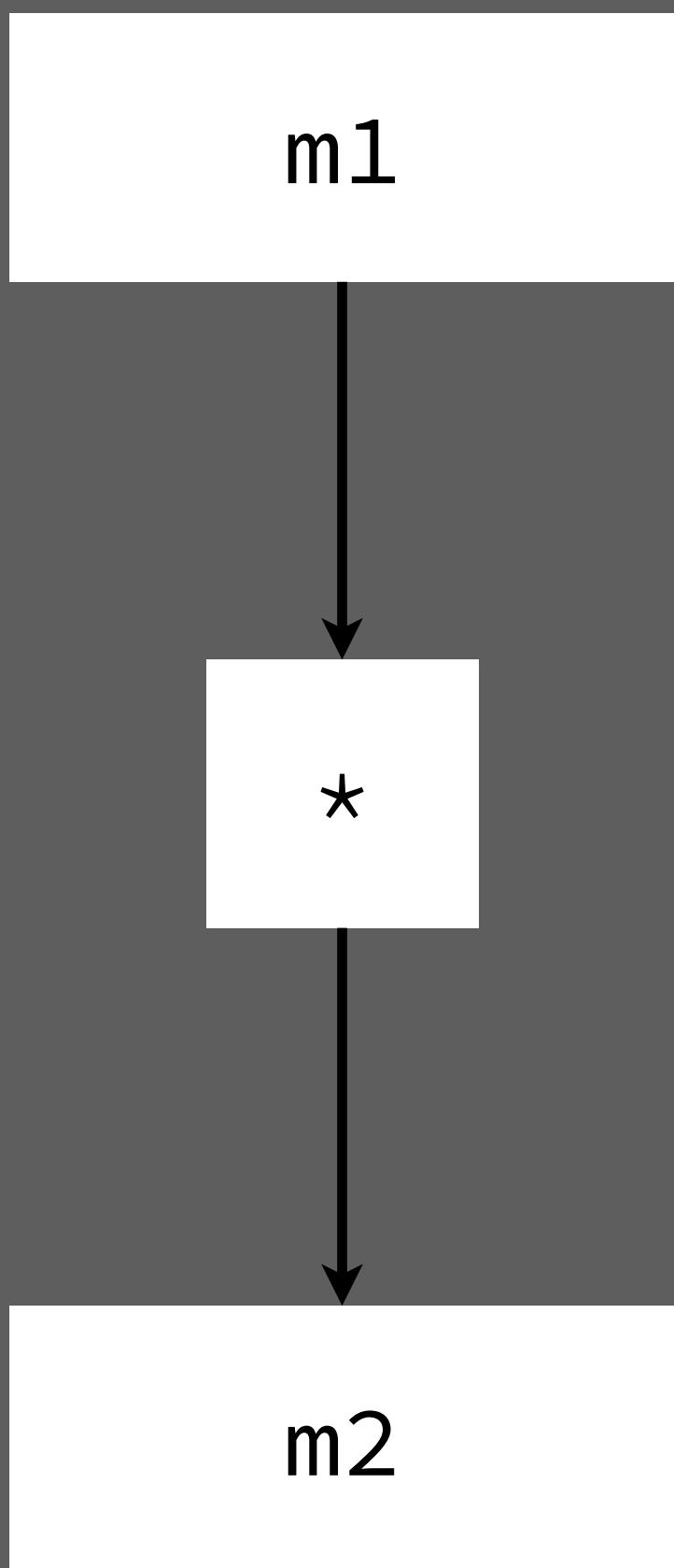


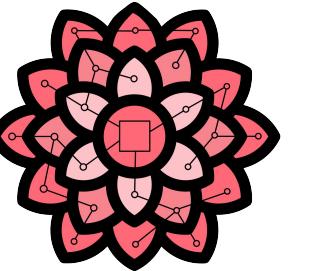
Dahlia

```
let m1: float[12];
let m2: float[12];

for (let i = 0 .. 12) {
    m2[i] = m1[i] * 2;
}
```

Hardware



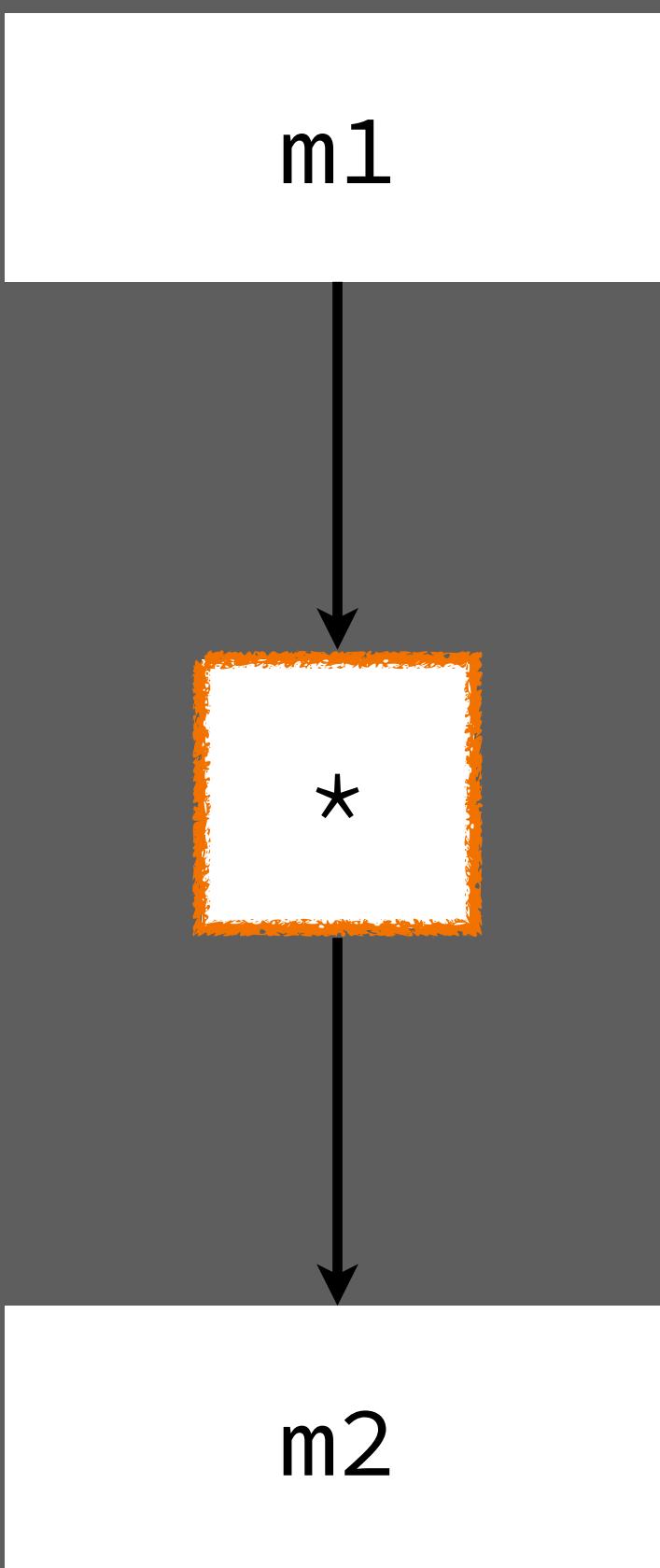


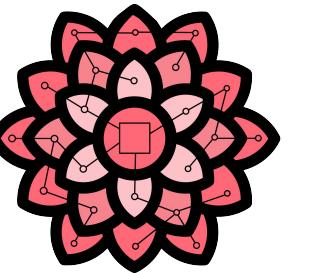
Dahlia

```
let m1: float[12];
let m2: float[12];

for (let i = 0 .. 12) unroll 3 {
    m2[i] = m1[i] * 2;
}
```

Hardware



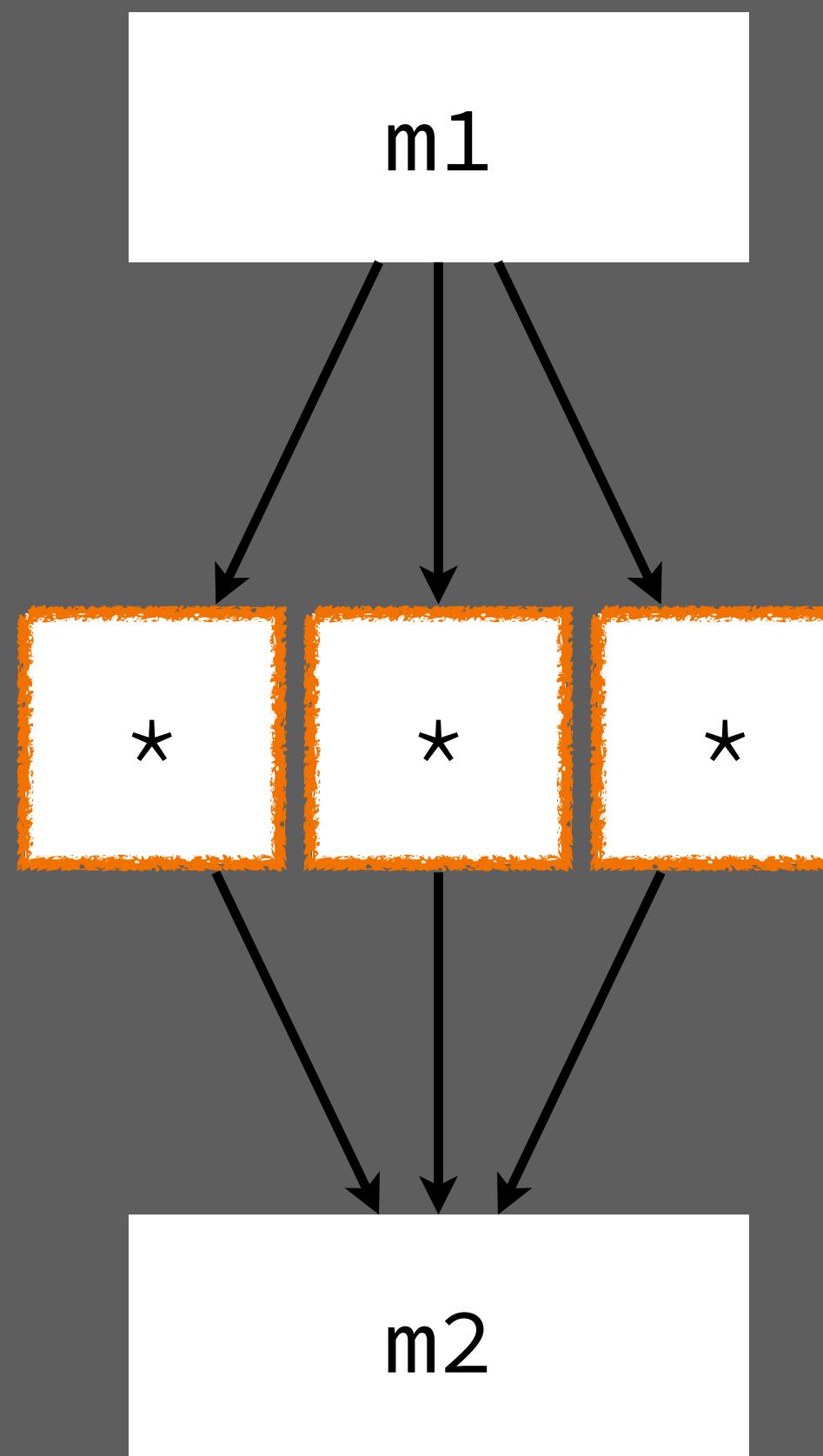


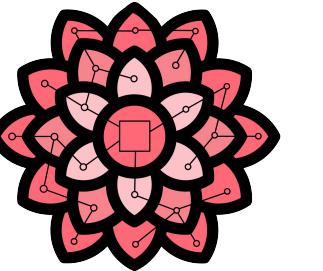
Dahlia

```
let m1: float[12];
let m2: float[12];

for (let i = 0 .. 12) unroll 3 {
    m2[i] = m1[i] * 2;
}
```

Hardware





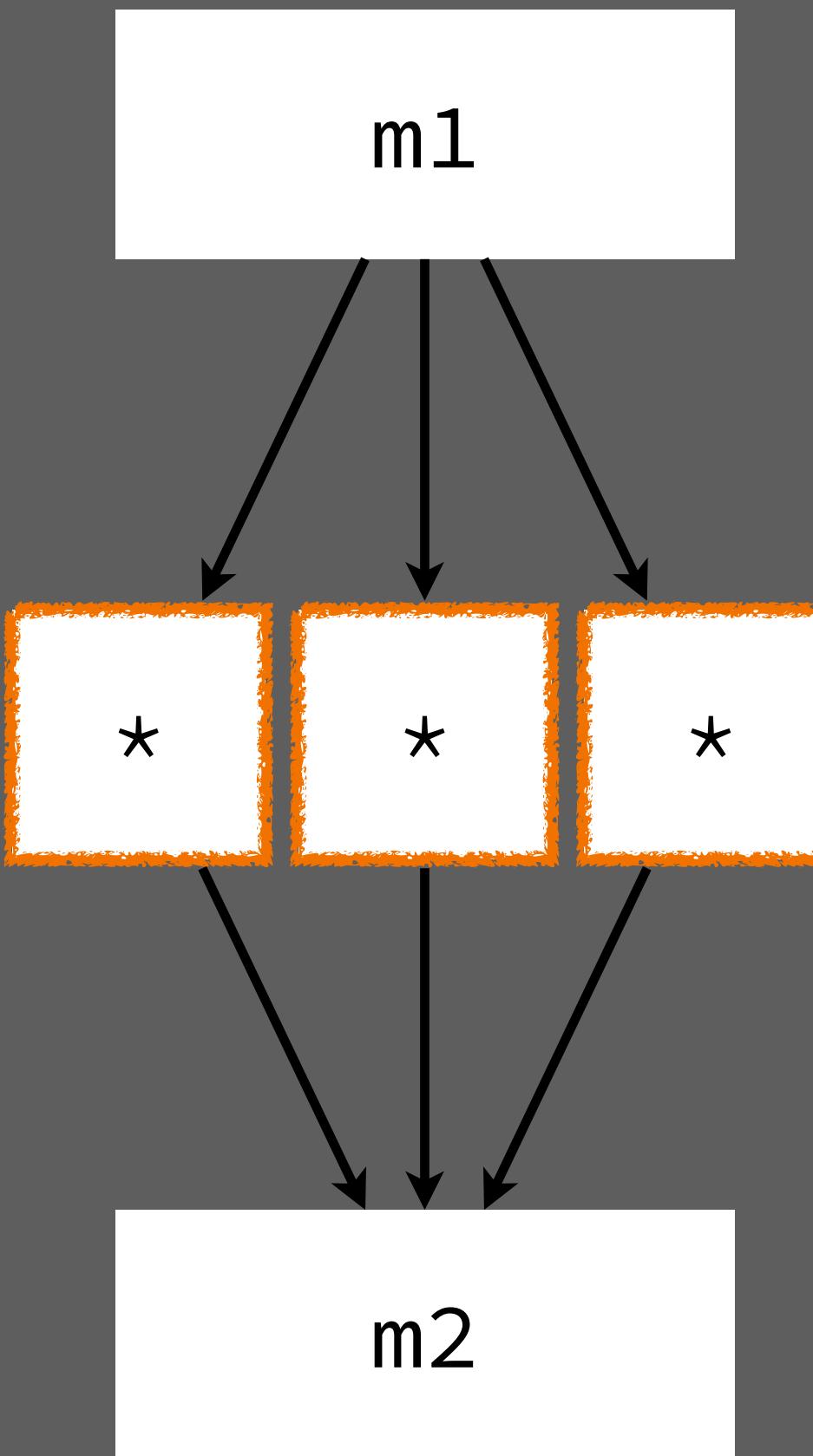
Dahlia

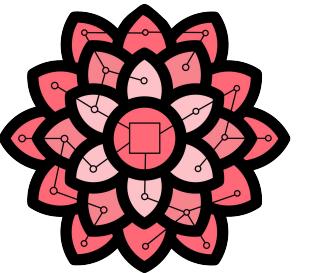
```
let m1: float[12];
let m2: float[12];

for (let i = 0 .. 12) unroll 3 {
    m2[i] = m1[i] * 2;
}
```

Error: Affine resource 'm1'
already used in this context.

Hardware





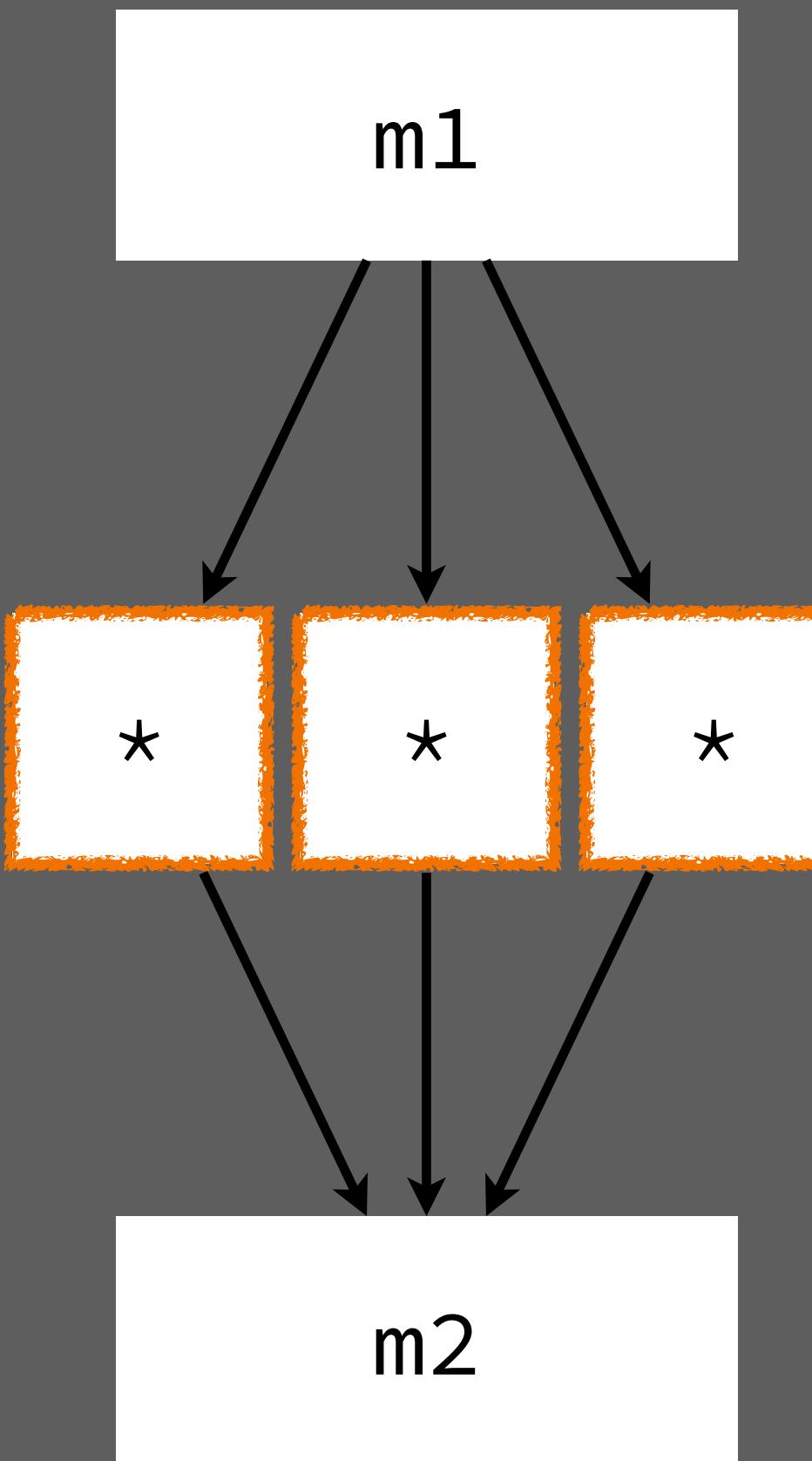
Dahlia

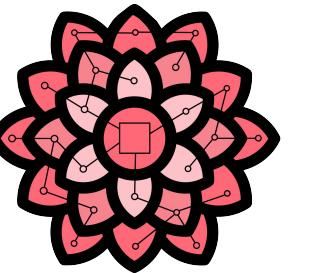
```
let m1: float[12];
let m2: float[12];

for (let i = 0 .. 4) {
    m2[3*i+0] = m1[3*i+0] * 2;
    m2[3*i+1] = m1[3*i+1] * 2;
    m2[3*i+2] = m1[3*i+2] * 2;
}
```

Error: Affine resource 'm1'
already used in this context.

Hardware



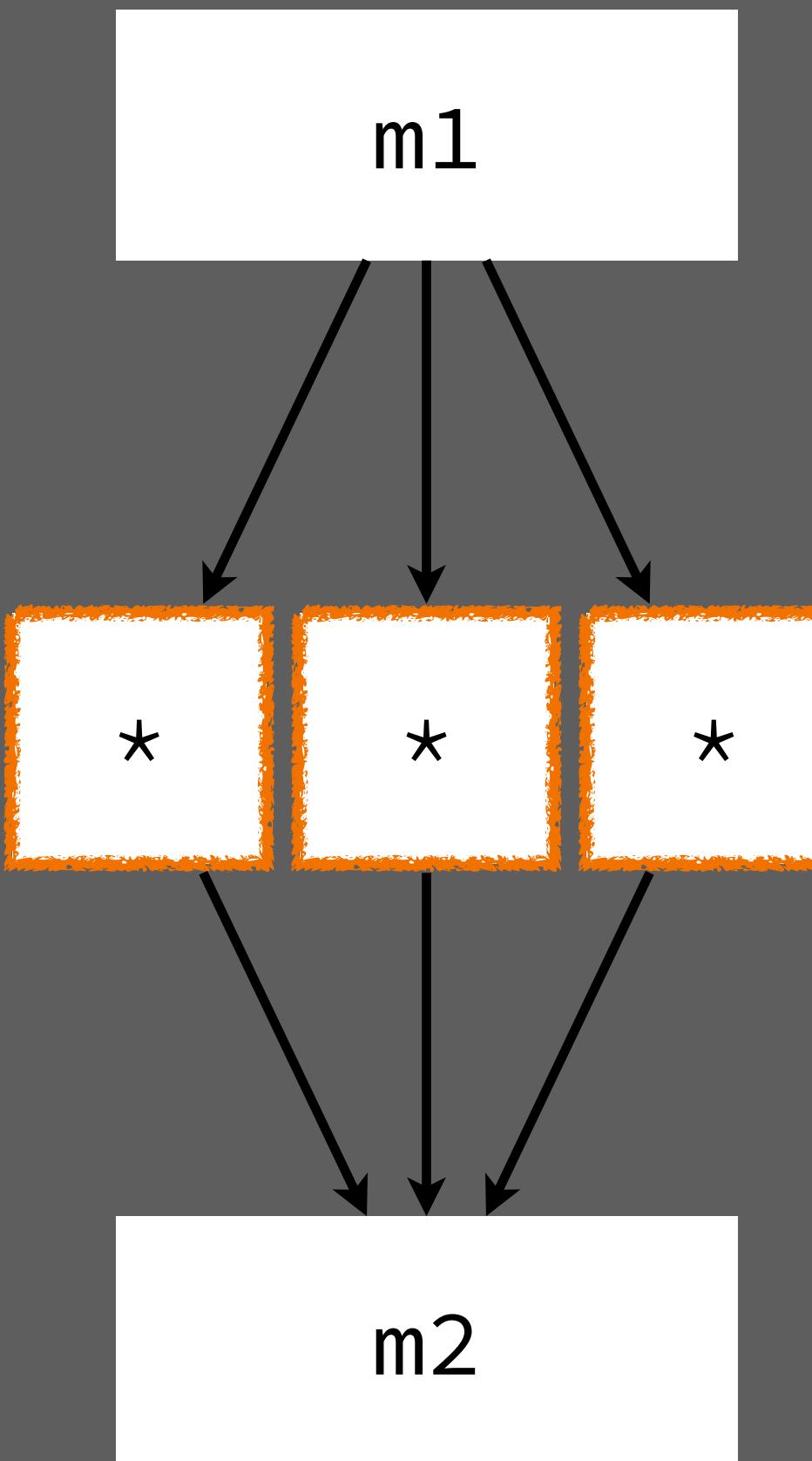


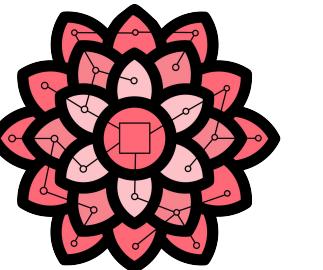
Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 12) unroll 3 {
    m2[i] = m1[i] * 2;
}
```

Hardware



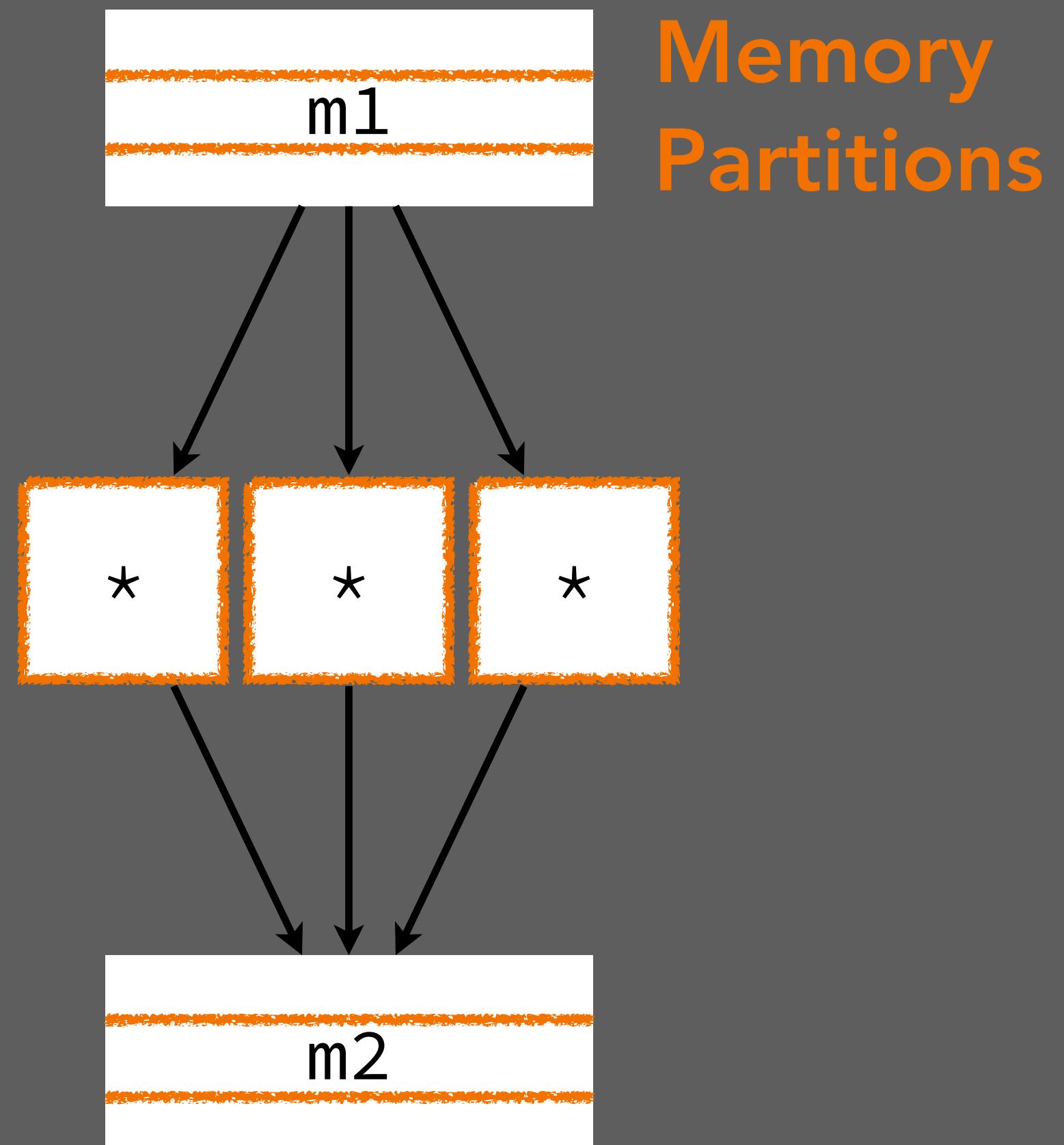


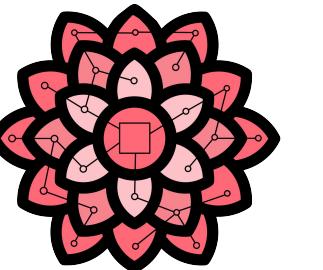
Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 12) unroll 3 {
    m2[i] = m1[i] * 2;
}
```

Hardware





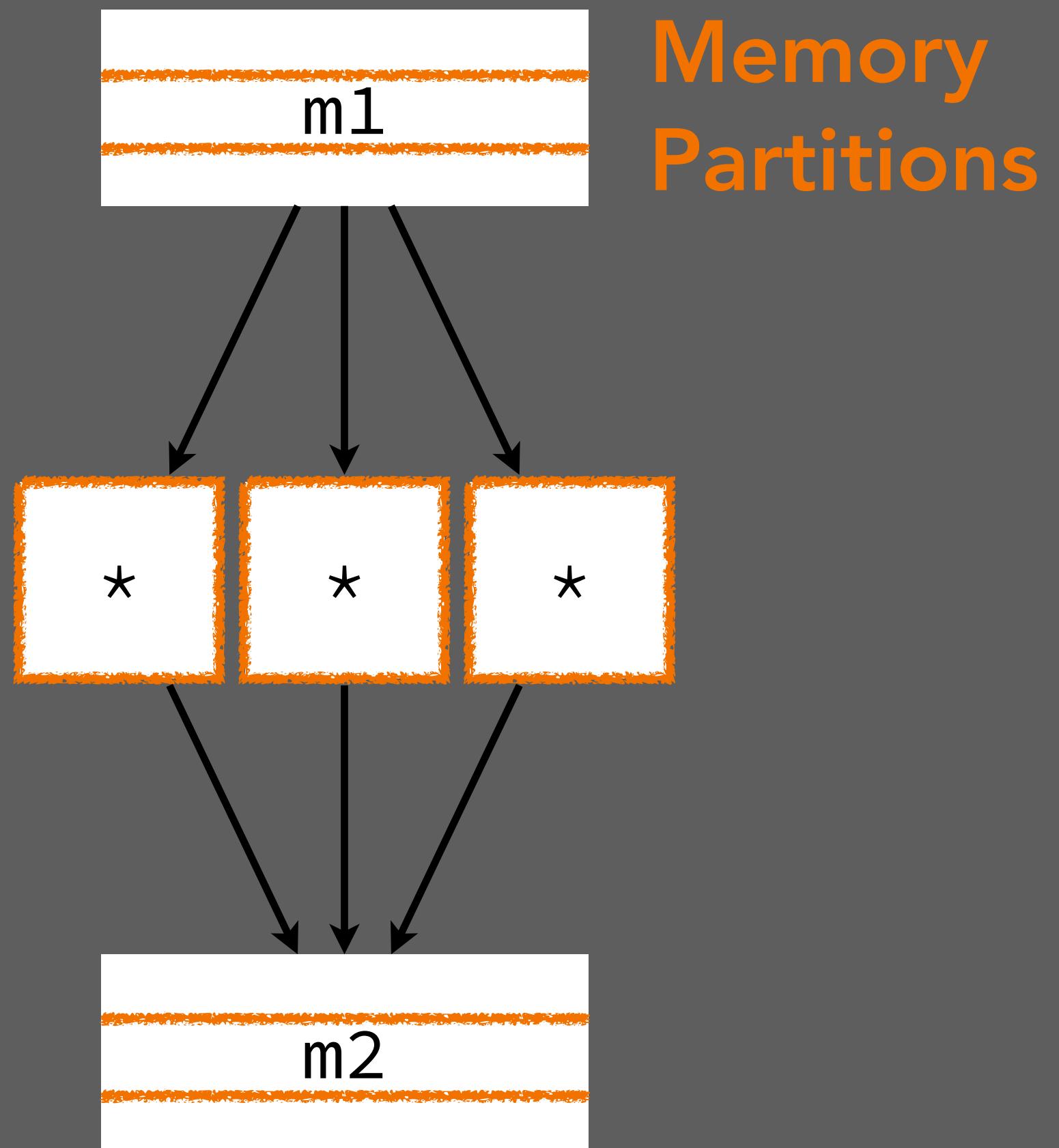
Dahlia

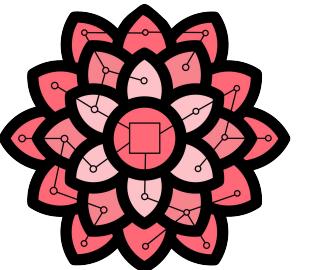
```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 12) unroll 3 {
    m2[i] = m1[i] * 2;
}
```

OK: Dahlia guarantees that parallel
accesses use disjoint partitions.

Hardware





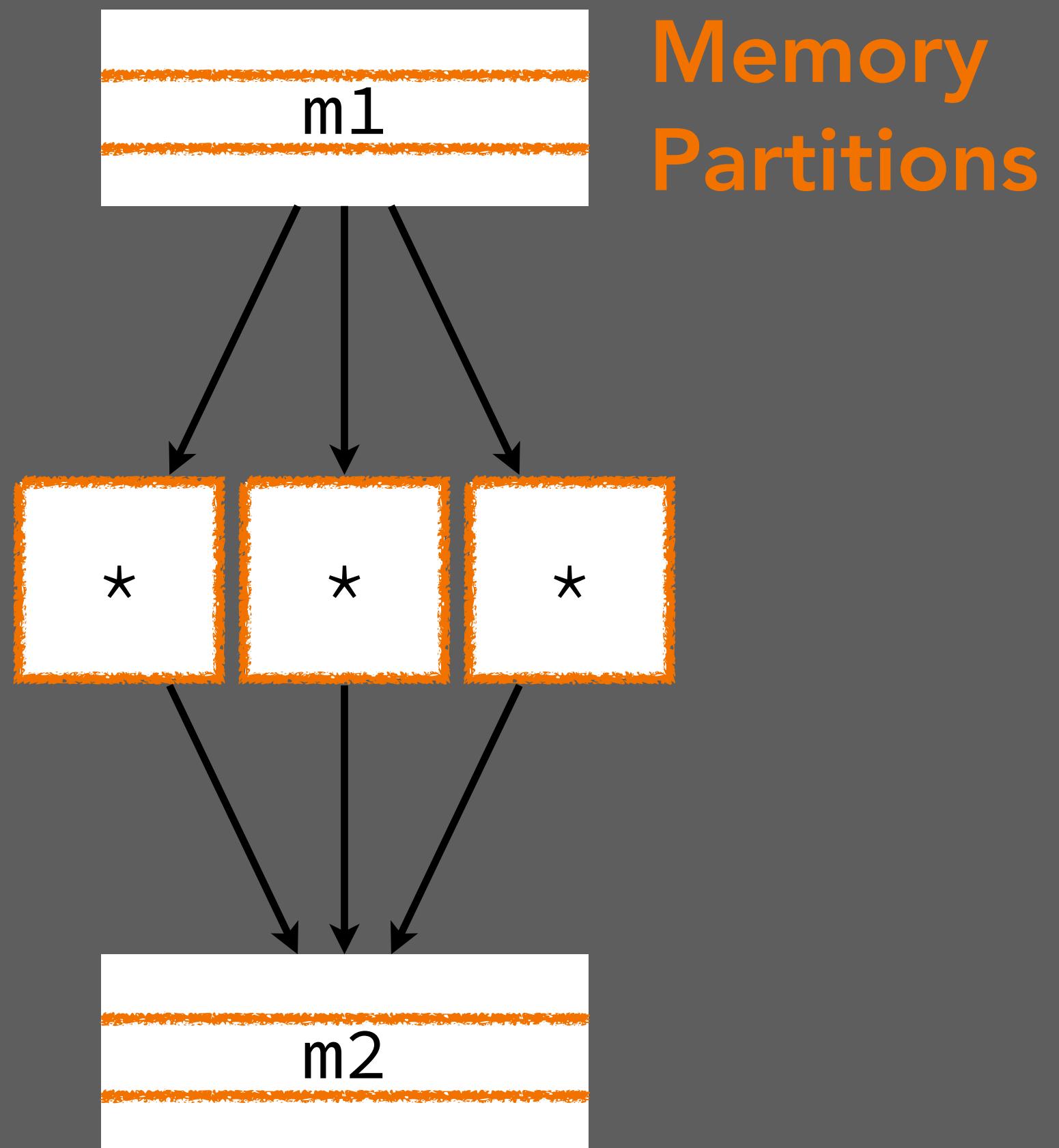
Dahlia

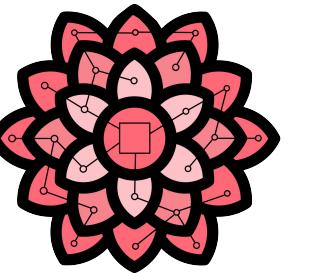
```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 4) {
    m2[3*i+0] = m1[3*i+0] * 2;
    m2[3*i+1] = m1[3*i+1] * 2;
    m2[3*i+2] = m1[3*i+2] * 2;
}
```

OK: Dahlia guarantees that parallel
accesses use disjoint partitions.

Hardware



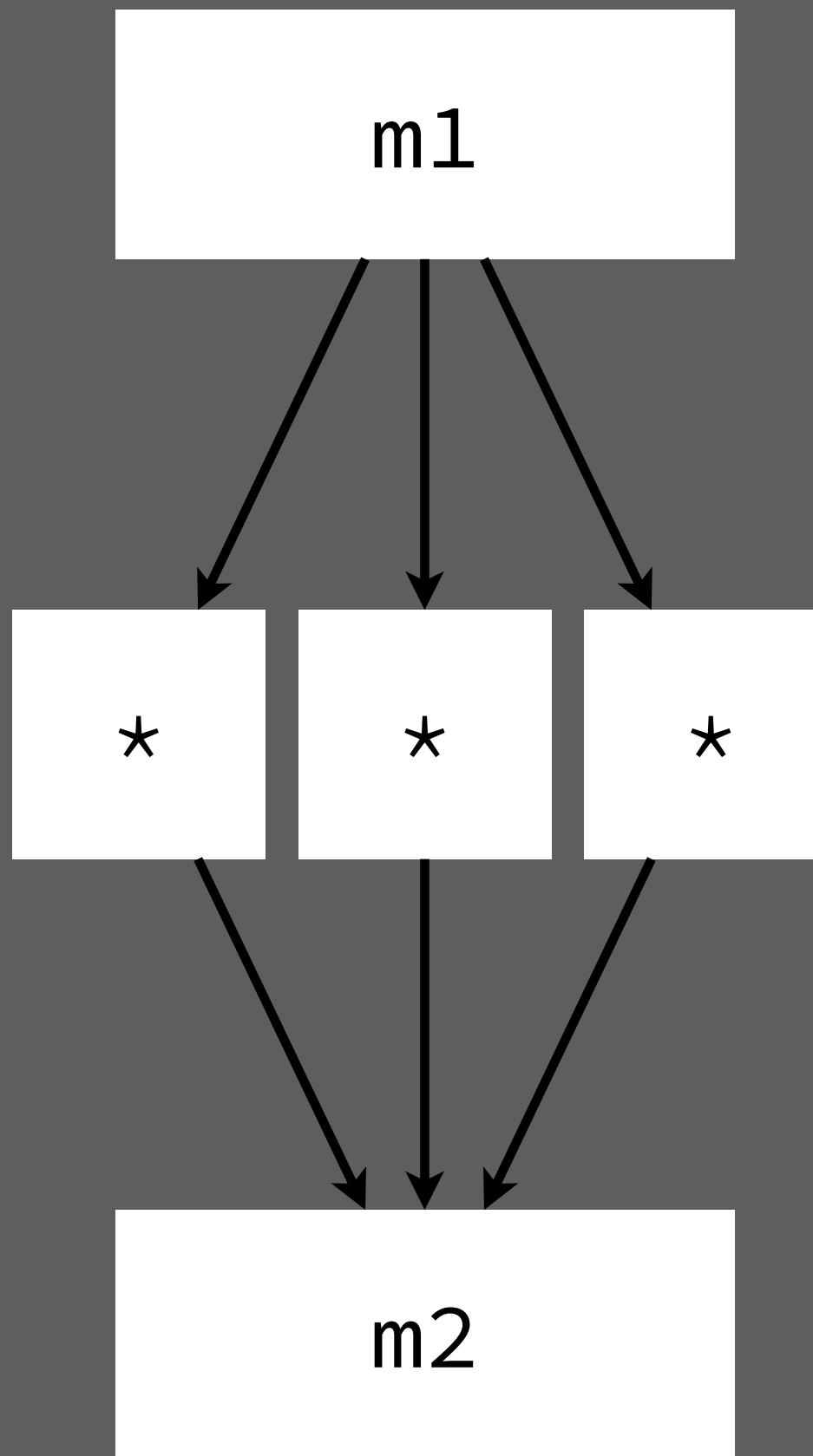


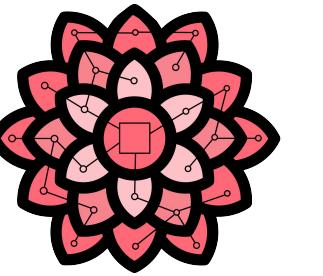
Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 4) {
  for (let j = 0 .. 3) {
    m2[i] = m1[3*i+j] * 2;
  }
}
```

Hardware



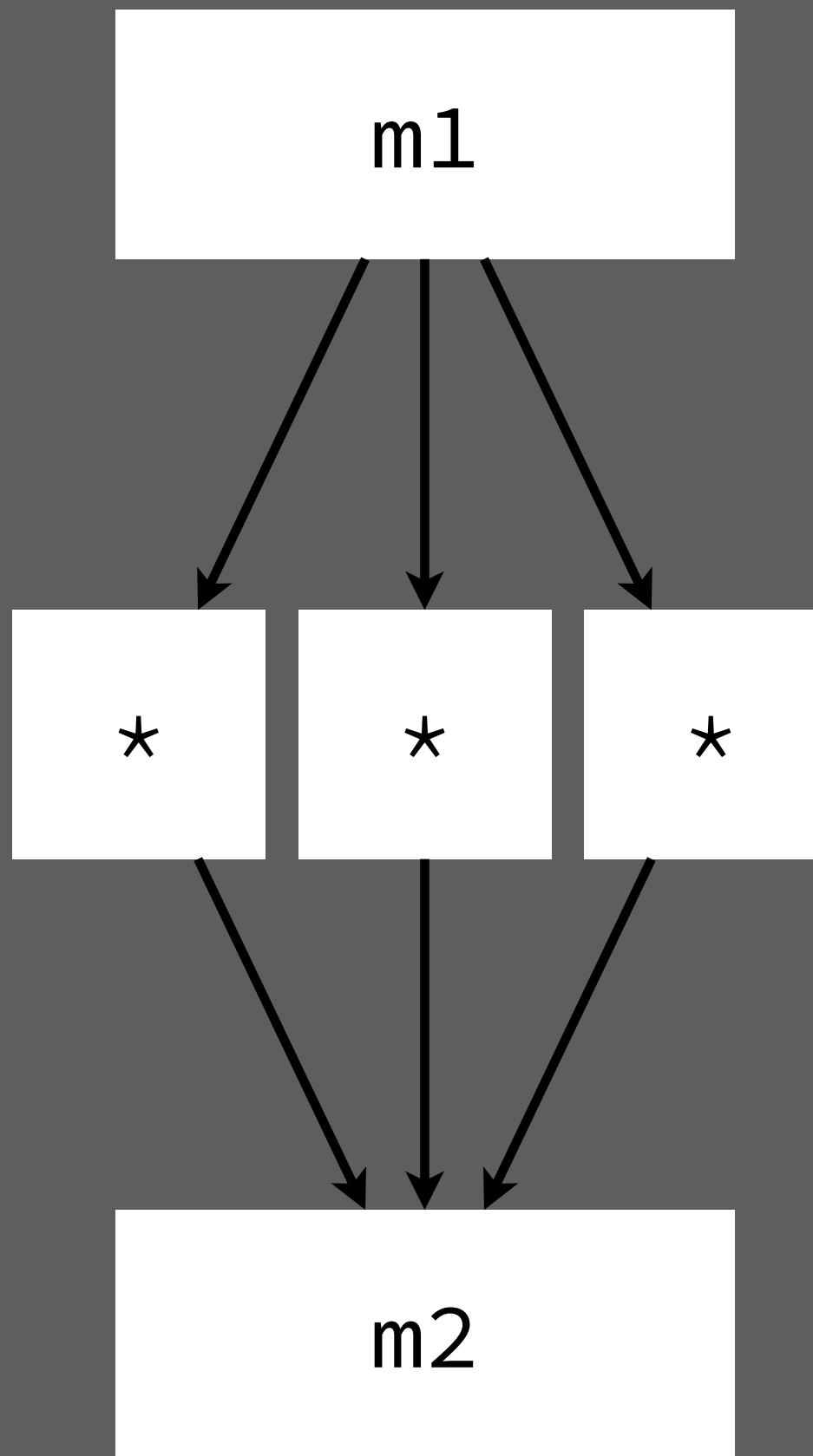


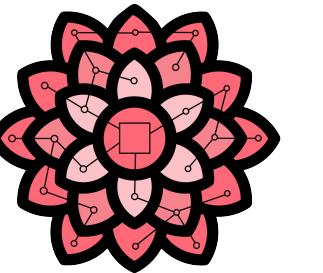
Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 4) {
  for (let j = 0 .. 3) unroll 3 {
    m2[i] = m1[3*i+j] * 2;
  }
}
```

Hardware



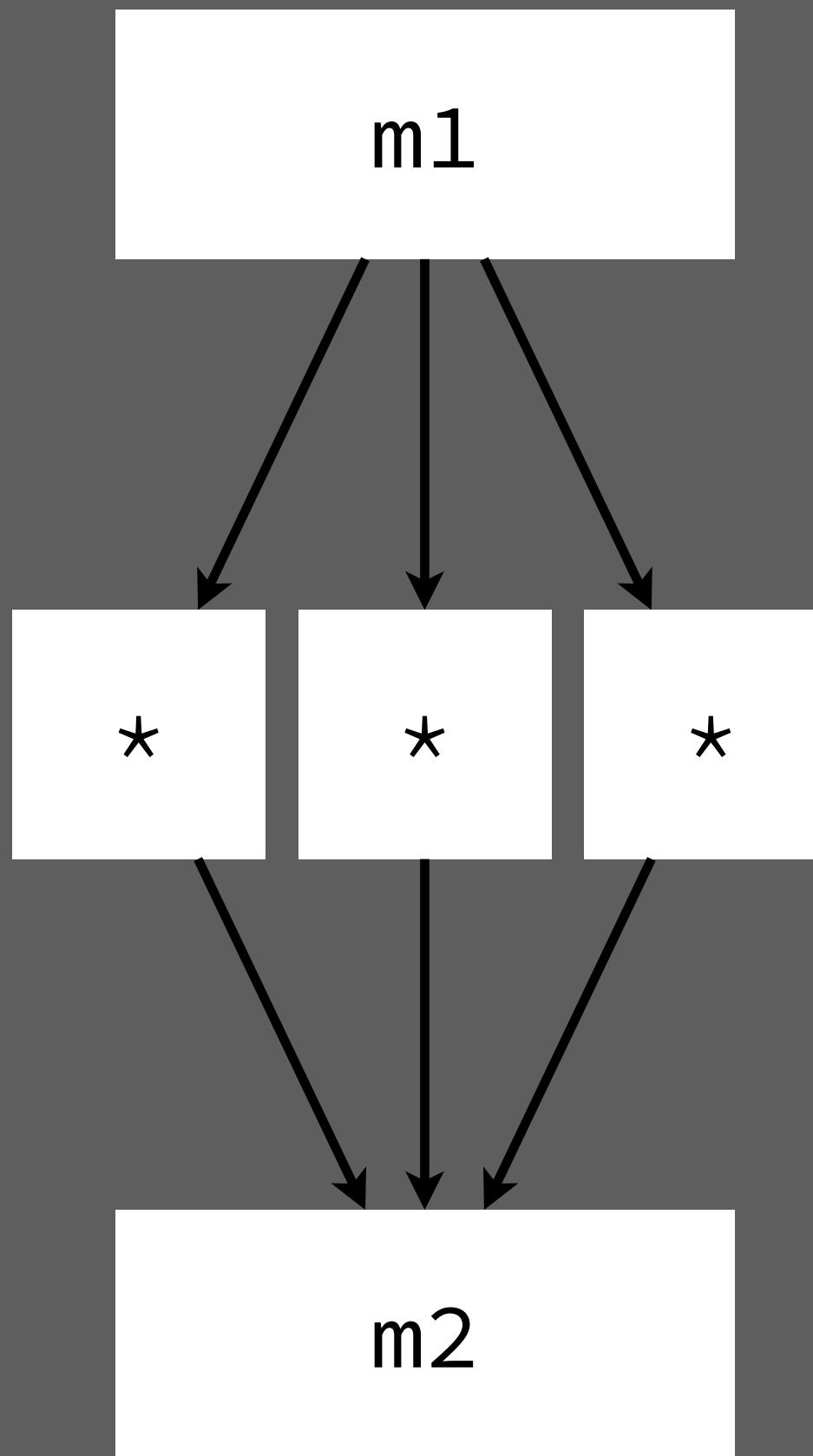


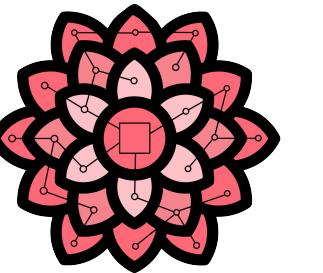
Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 4) {
  for (let j = 0 .. 3) unroll 3 {
    m2[i] = m1[f(i, j)] * 2;
  }
}
```

Hardware





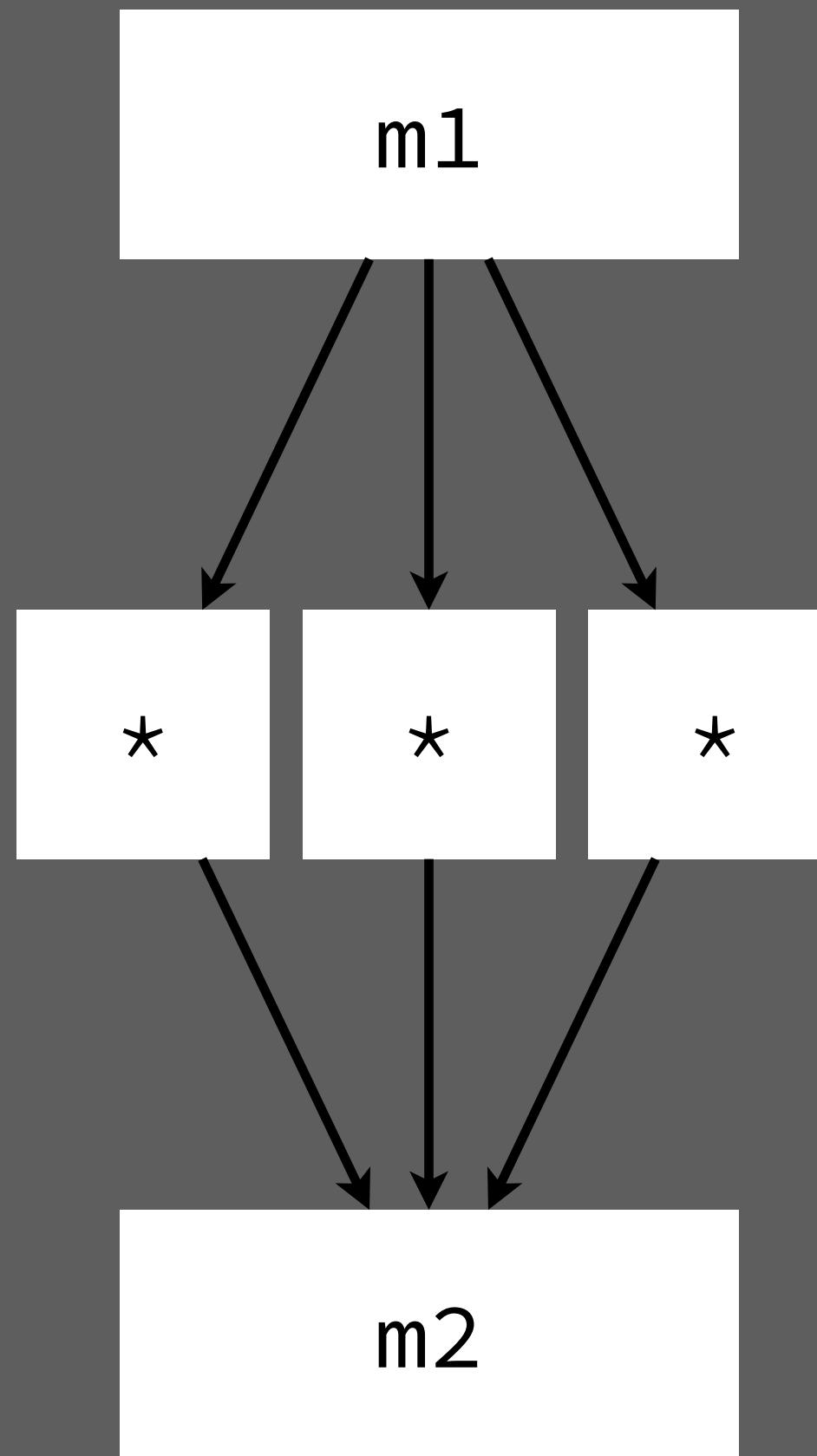
Dahlia

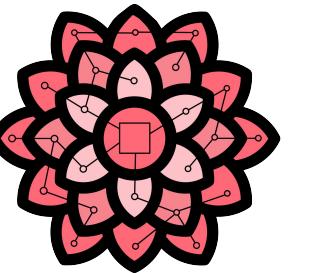
```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 4) {
  for (let j = 0 .. 3) unroll 3 {
    m2[i] = m1[f(i, j)] * 2;
  }
}
```

Parallelizing access patterns requires
unpredictable analyses

Hardware



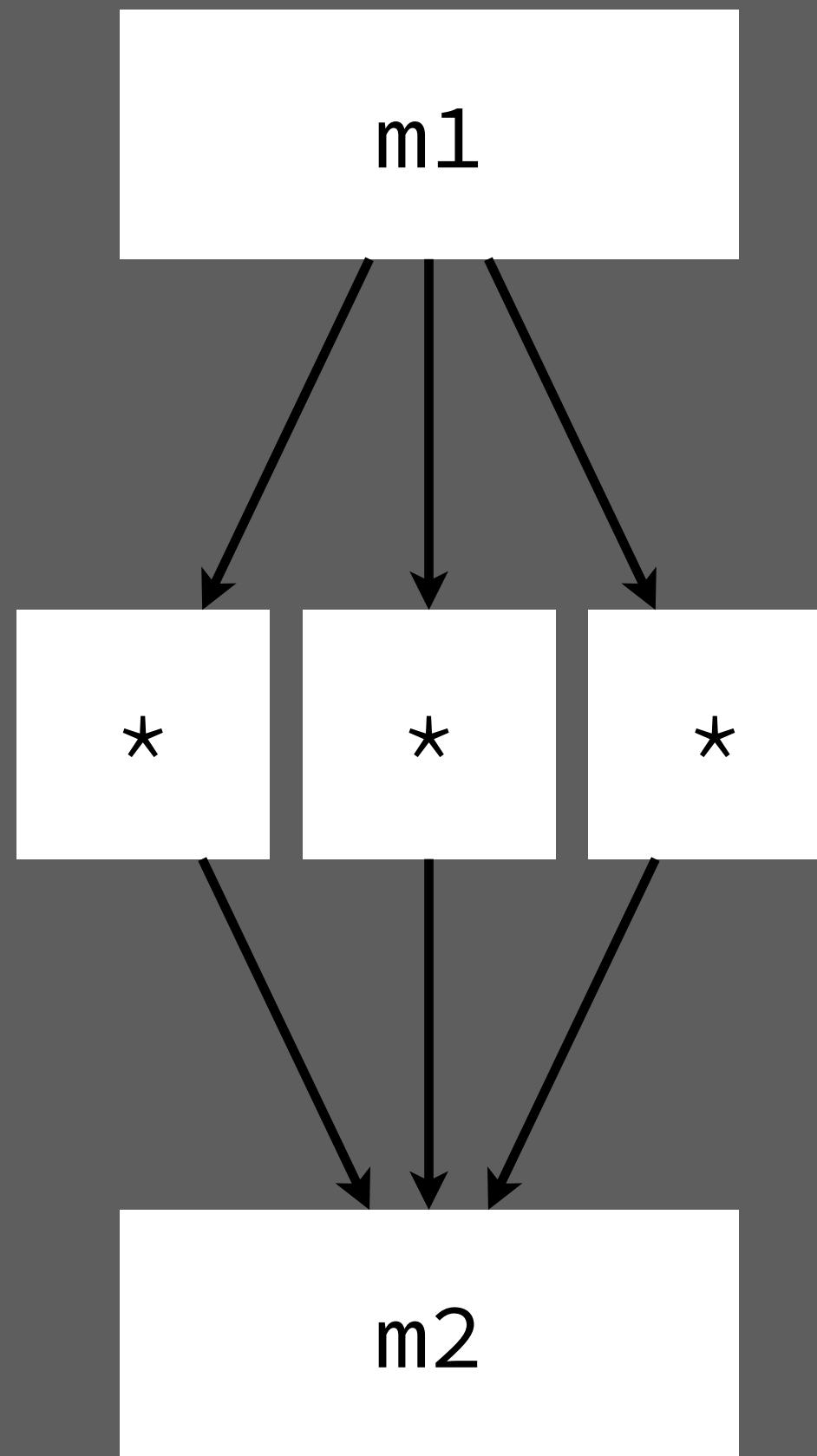


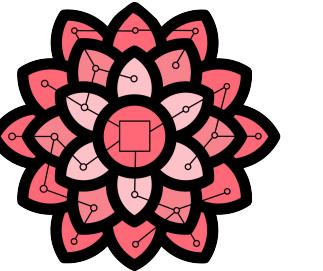
Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 4) {
  for (let j = 0 .. 3) unroll 3 {
    m2[i] = m1[3*i+j] * 2;
  }
}
```

Hardware





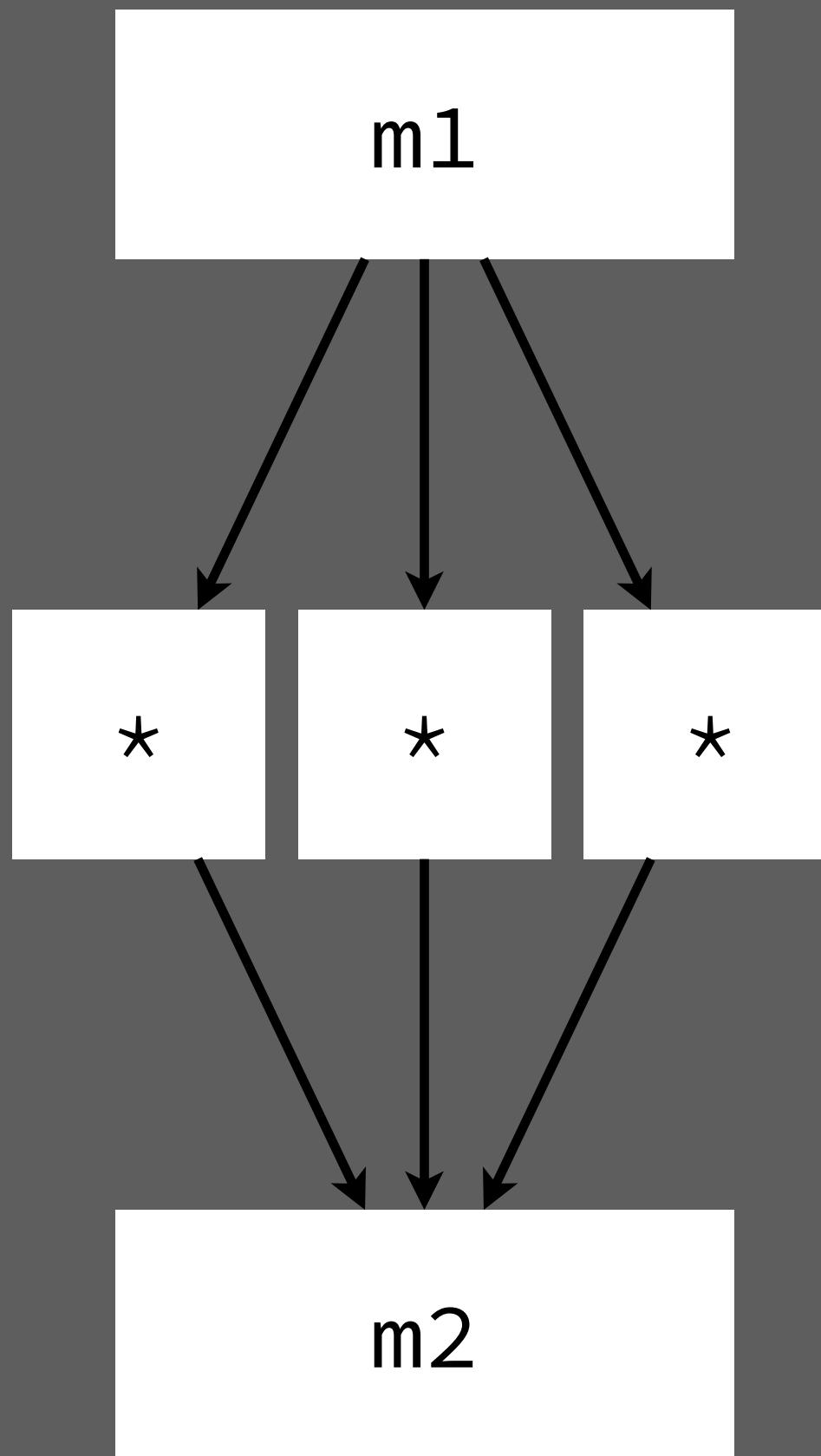
Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 4) {
  for (let j = 0 .. 3) unroll 3 {
    m2[i] = m1[3*i+j] * 2;
  }
}
```

Error: Cannot parallelize
dynamic access pattern.

Hardware



Memory Views

Memory **Views**

```
let m1: float[12 bank 3];
```

Memory Views

```
let m1: float[12 bank 3];
```

```
view v1 = suffix m1[by 2*i];
```

Memory Views

```
let m1: float[12 bank 3];
```

```
view v1 = suffix m1[by 2*i];
```

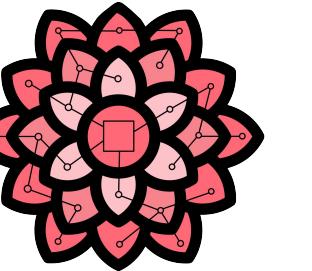
- Hardware cost of indexing logic

Memory Views

```
let m1: float[12 bank 3];
```

```
view v1 = suffix m1[by 2*i];
```

- Hardware cost of indexing logic
- Proof that accesses can be parallelized



Dahlia

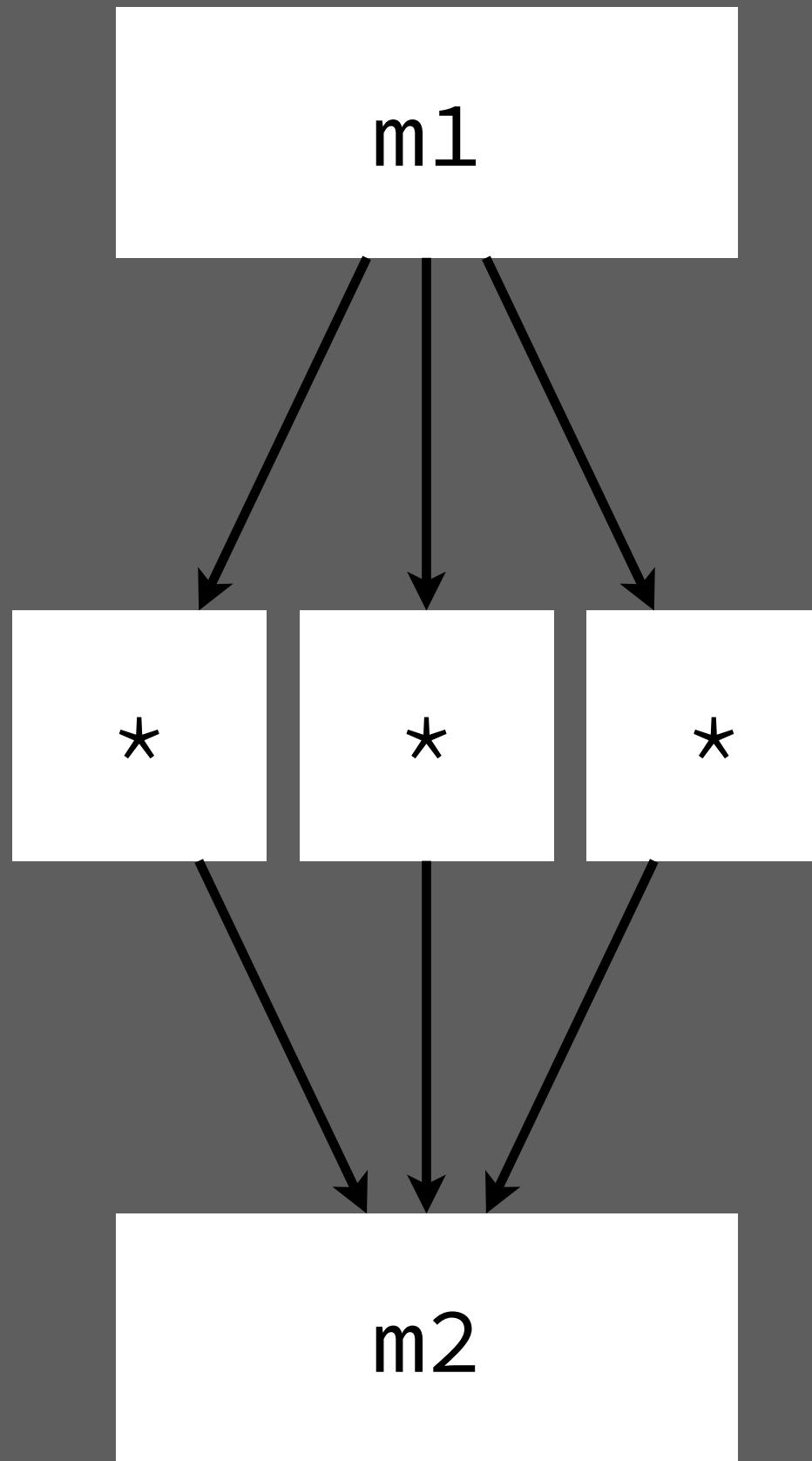
```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

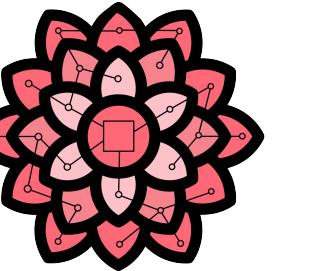
for (let i = 0 .. 4) {

    view v1 = suffix m1[by 3*i];

    for (let j = 0 .. 3) unroll 3 {
        m2[i] = v1[j] * 2; // m1[3*i+j]
    }
}
```

Hardware





Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

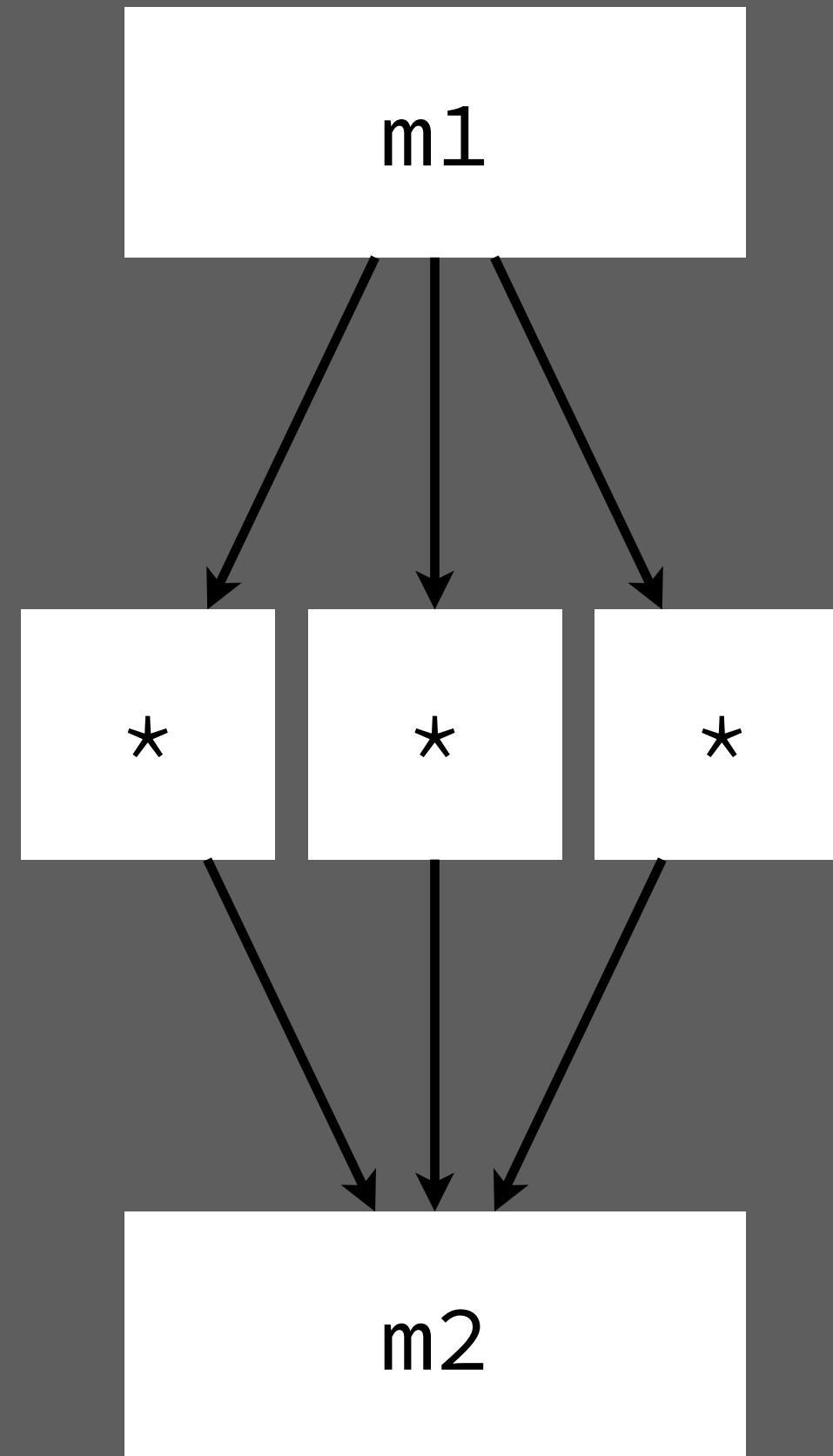
for (let i = 0 .. 4) {

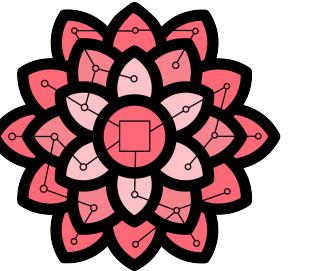
    view v1 = suffix m1[by 3*i];

    for (let j = 0 .. 3) unroll 3 {
        m2[i] = v1[j] * 2; // m1[3*i+j]
    }
}
```

Suffix View

Hardware





Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 4) {

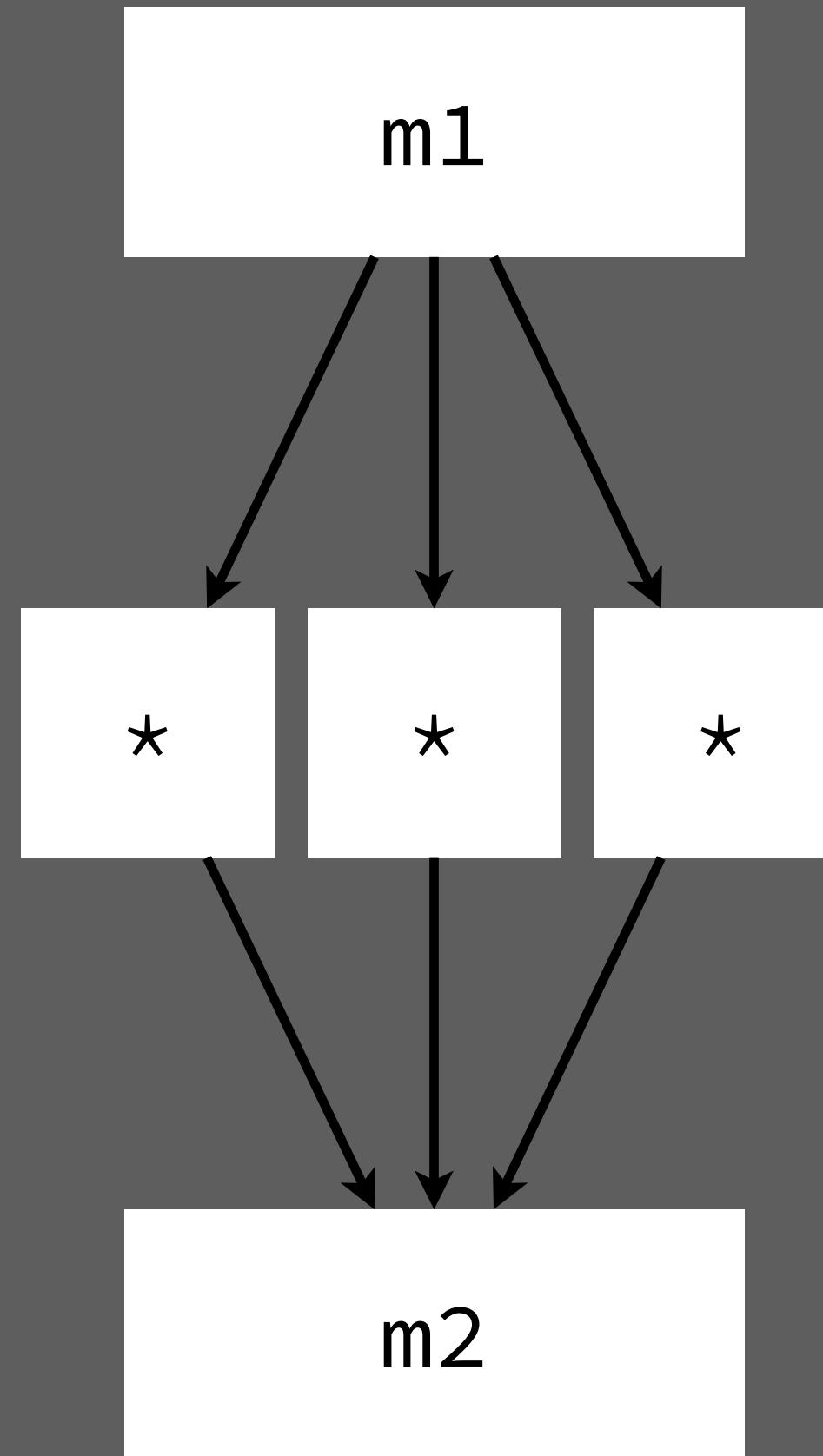
    view v1 = suffix m1[by 3*i];

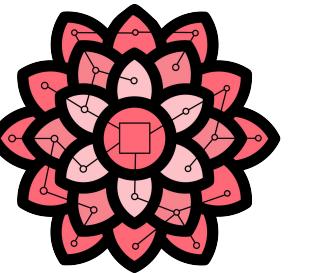
    for (let j = 0 .. 3) unroll 3 {
        m2[i] = v1[j] * 2; // m1[3*i+j]
    }
}
```

Suffix View

- Offset by banking factor

Hardware





Dahlia

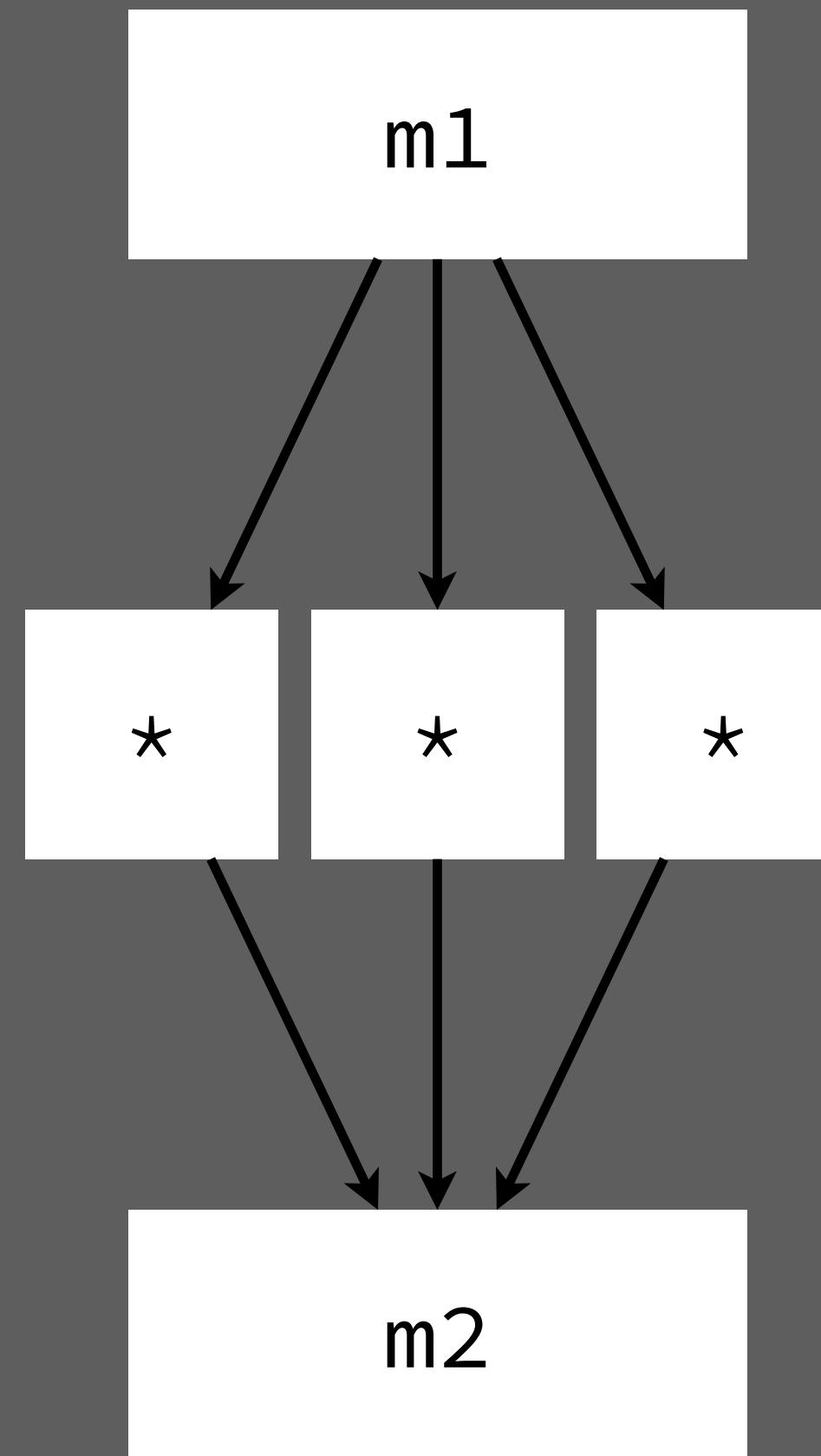
```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

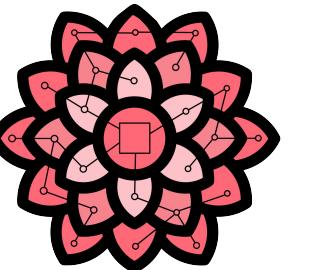
for (let i = 0 .. 4) {

    view v1 = shift m1[by f(i)];

    for (let j = 0 .. 3) unroll 3 {
        m2[i] = v1[j] * 2; // m1[f(i)+j]
    }
}
```

Hardware





Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

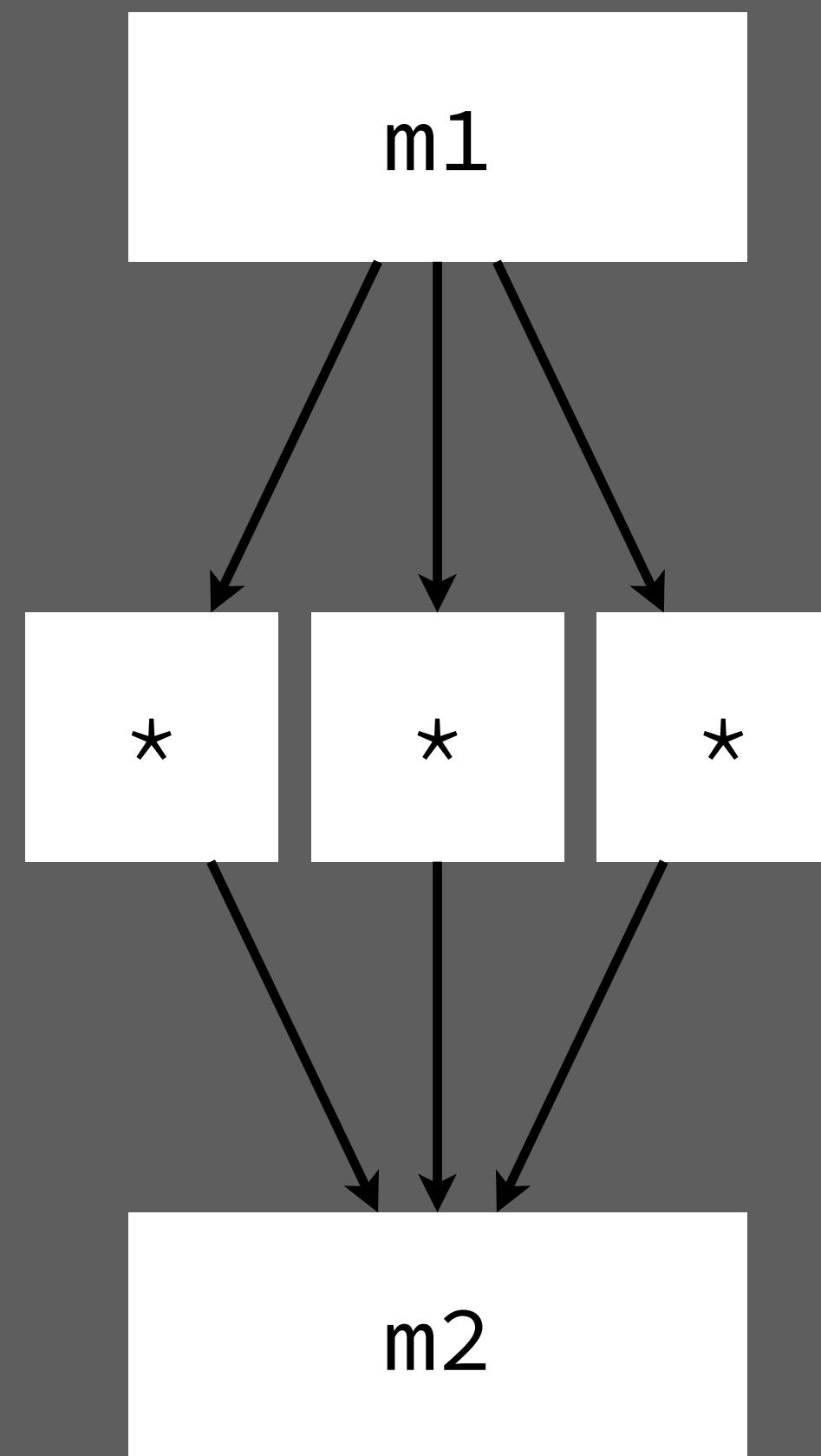
for (let i = 0 .. 4) {

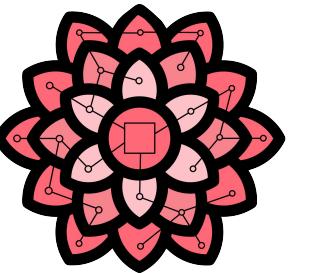
    view v1 = shift m1[by f(i)];

    for (let j = 0 .. 3) unroll 3 {
        m2[i] = v1[j] * 2; // m1[f(i)+j]
    }
}
```

Shift View

Hardware





Dahlia

```
let m1: float[12 bank 3];
let m2: float[12 bank 3];

for (let i = 0 .. 4) {

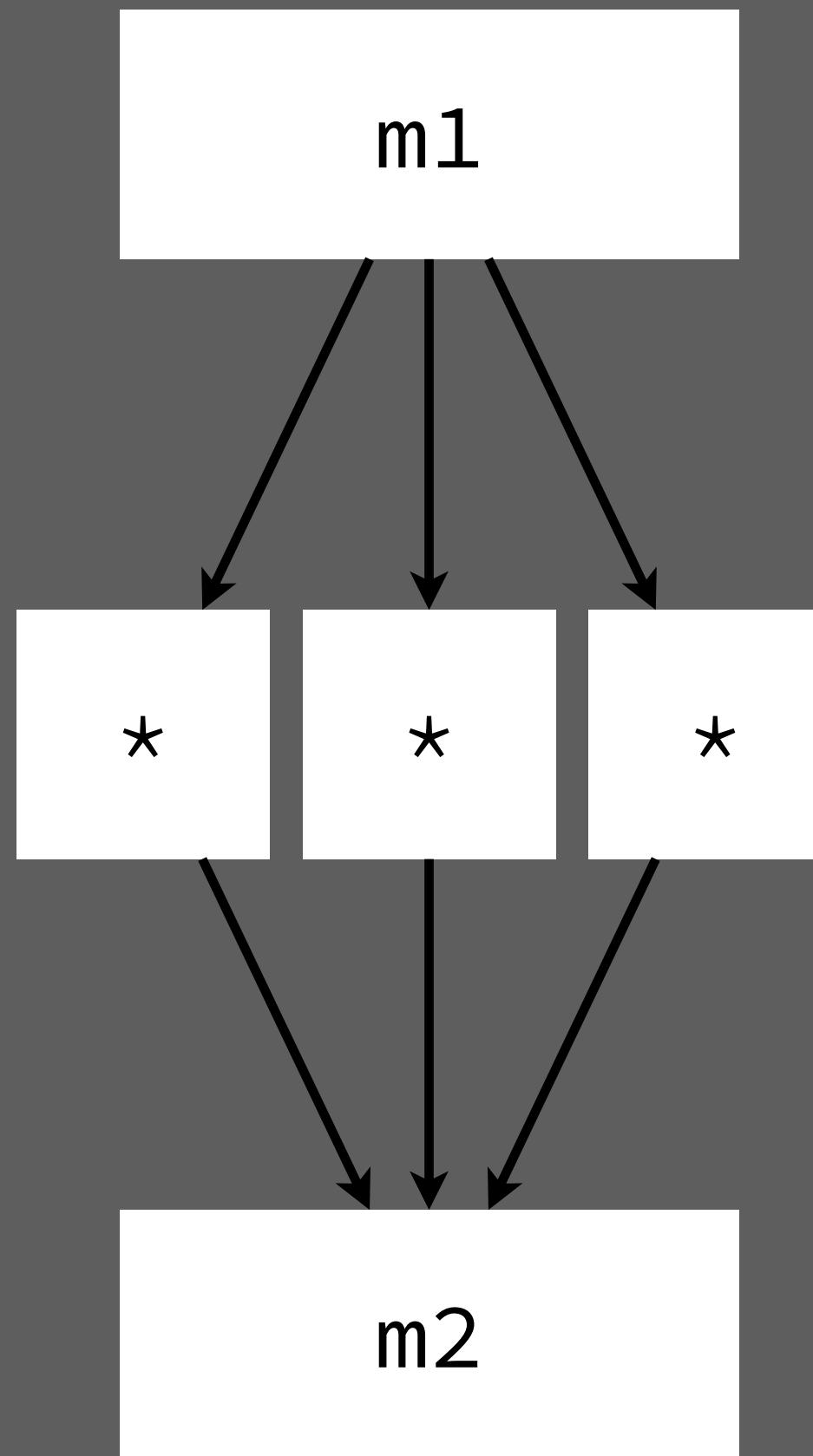
    view v1 = shift m1[by f(i)];

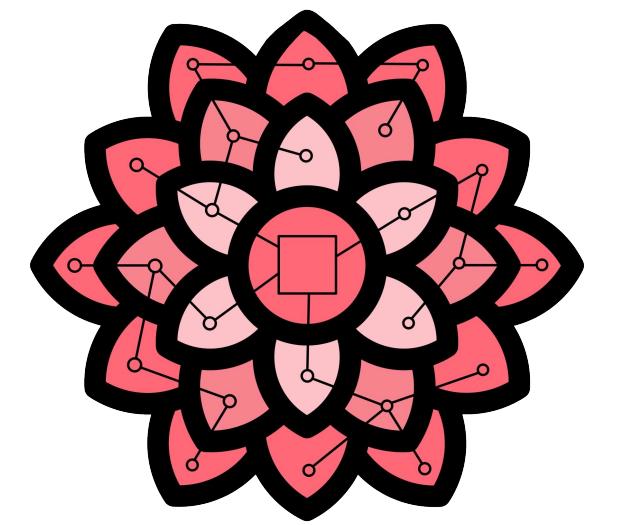
    for (let j = 0 .. 3) unroll 3 {
        m2[i] = v1[j] * 2; // m1[f(i)+j]
    }
}
```

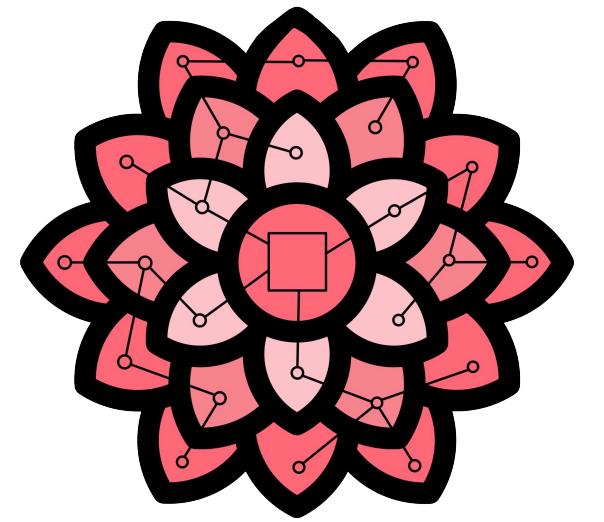
Shift View

- Arbitrary offset

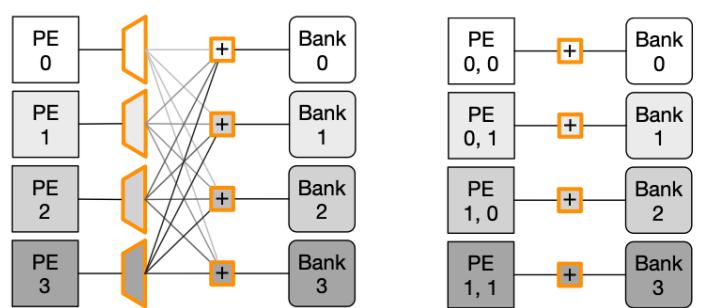
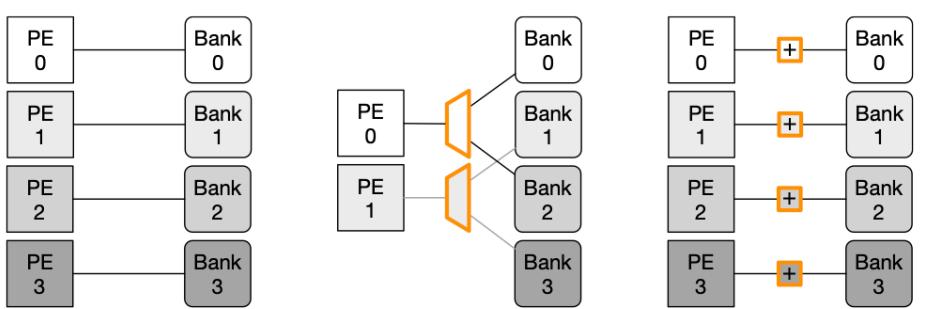
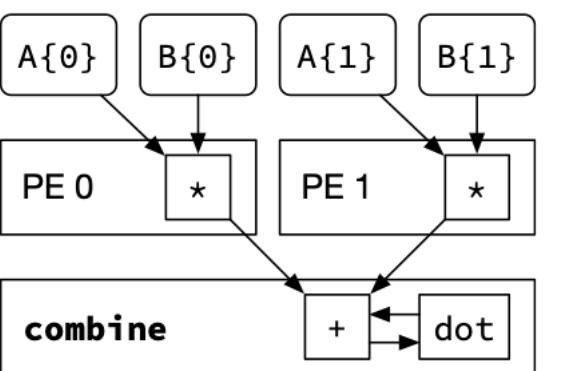
Hardware

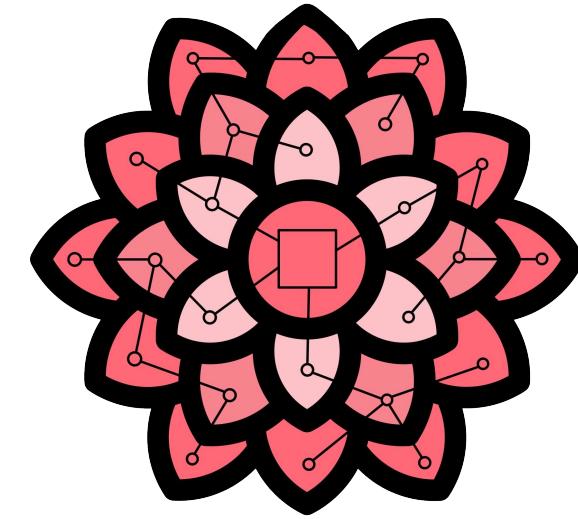






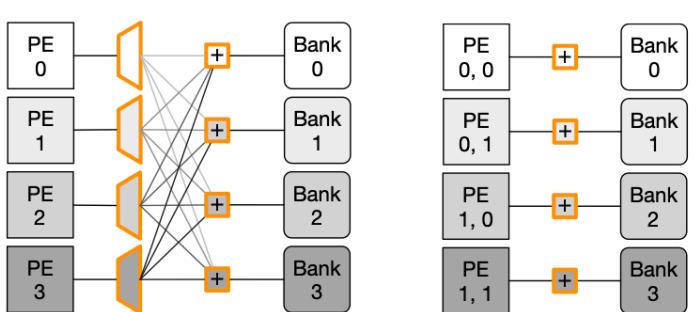
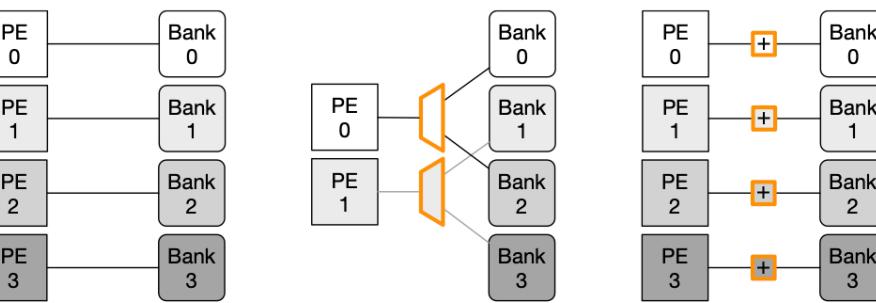
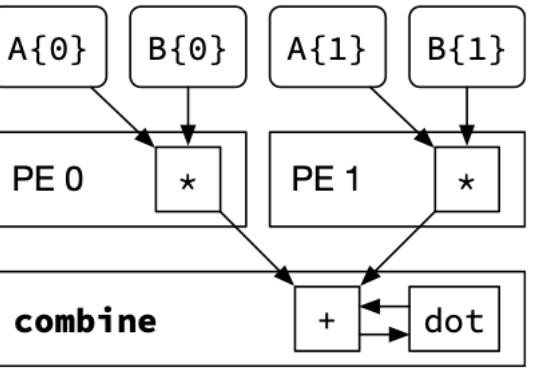
Language Features





Language Features

Formalism



$x \in \text{variables}$ $a \in \text{memories}$ $n \in \text{numbers}$

$b ::= \text{true} \mid \text{false}$ $v ::= n \mid b$
 $e ::= v \mid \text{bop } e_1 \ e_2 \mid x \mid a[e]$
 $c ::= e \mid \text{let } x = e \mid c_1 — c_2 \mid c_1 ; c_2 \mid \text{if } x \ c_1 \ c_2 \mid$
 $\text{while } x \ c \mid x := e \mid a[e_1] := e_2 \mid \text{skip}$
 $\tau ::= \text{bit}\langle n \rangle \mid \text{float} \mid \text{bool} \mid \text{mem } \tau[n_1]$

$$\frac{\sigma_1, \rho_1, c \rightarrow \sigma', \rho', c'}{\sigma, \rho, e_1 \rightarrow \sigma', \rho', e'_1}$$

$$\frac{\sigma, \rho, e_1 \rightarrow \sigma', \rho', e'_1}{\sigma, \rho, a[e_1] := e_2 \rightarrow \sigma', \rho', a[e'_1] := e_2}$$

$$\frac{\sigma, \rho, e \rightarrow \sigma', \rho', e'}{\sigma, \rho, a[n] := e \rightarrow \sigma', \rho', a[n] := e'}$$

$$\frac{a \notin \rho}{\sigma, \rho, a[n] := v \rightarrow \sigma[a[n] \mapsto v], \rho \cup \{a\}, \text{skip}}$$

$$\frac{\sigma, \rho, e \rightarrow \sigma', \rho', e'}{\sigma, \rho, \text{let } x = e \rightarrow \sigma', \rho', \text{let } x = e'}$$

$$\frac{}{\sigma, \rho, \text{let } x = v \rightarrow \sigma[x \mapsto v], \rho, \text{skip}}$$

$$\frac{\sigma, \rho, c_1 \rightarrow \sigma', \rho', c'_1}{\sigma, \rho, c_1 ; c_2 \rightarrow \sigma', \rho', c'_1 ; c_2}$$

$$\frac{}{\sigma, \rho, \text{skip} ; c_2 \rightarrow \sigma, \rho, c}$$

$$\frac{}{\sigma, \rho, c_1 — c_2 \rightarrow \sigma, \rho, c_1 \xrightarrow{\rho} c_2}$$

$$\frac{\sigma, \rho, c_1 \rightarrow \sigma', \rho', c'_1}{\sigma, \rho, c_1 \xrightarrow{\rho''} c_2 \rightarrow \sigma', \rho', c'_1 \xrightarrow{\rho''} c_2}$$

$$\frac{\sigma, \rho'', c_2 \rightarrow \sigma', \rho''', c'_2}{\sigma, \rho, \text{skip} \xrightarrow{\rho''} c_2 \rightarrow \sigma', \rho, \text{skip} \xrightarrow{\rho'''} c'_2}$$

$$\frac{\sigma, \rho, \text{skip} \xrightarrow{\rho''} \text{skip} \rightarrow \sigma, \rho \cup \rho'', \text{skip}}{\sigma, \rho, \text{skip} \xrightarrow{\rho''} \text{skip} \rightarrow \sigma, \rho \cup \rho'', \text{skip}}$$

$$\frac{\sigma_1, \rho_1, e \Downarrow \sigma_2, \rho_2, v}{\sigma_1, \rho_1, e_1 \Downarrow \sigma_2, \rho_2, v_1} \quad \frac{\sigma_2, \rho_2, e_2 \Downarrow \sigma_3, \rho_3, v_2}{\sigma_2, \rho_2, e_1 \Downarrow \sigma_3, \rho_3, v_3} \quad \frac{v_3 = v_1 \text{ bop } v_2}{\sigma_1, \rho_1, \text{bop } e_1 \ e_2 \Downarrow \sigma_3, \rho_3, v_3}$$

$$\frac{\sigma(x) = v}{\sigma, \rho, x \Downarrow \sigma, \rho, v}$$

$$\frac{\sigma_1, \rho_1, e \Downarrow \sigma_2, \rho_2, n \quad \sigma_2(a)(n) = v}{\sigma_1, \rho_1, a[e] \Downarrow \sigma_2, \rho_2 \cup \{a\}, v}$$

$$\frac{\sigma_1, \rho_1, c \Downarrow \sigma_2, \rho_2}{\sigma_1, \rho_1, \text{let } x = e \Downarrow \sigma_2[x \mapsto v], \rho_2}$$

$$\frac{\sigma_1, \rho_1, e \Downarrow \sigma_2, \rho_2, v}{\sigma_1, \rho_1, \text{let } x = e \Downarrow \sigma_2[x \mapsto v], \rho_2}$$

$$\frac{\sigma_1, \rho_1, c_1 \Downarrow \sigma_2, \rho_2 \quad \sigma_2, \rho_1, c_2 \Downarrow \sigma_3, \rho_3}{\sigma_1, \rho_1, c_1 — c_2 \Downarrow \sigma_3, \rho_2 \cup \rho_3}$$

$$\frac{\sigma_1, \rho_1, c_1 \Downarrow \sigma_2, \rho_2 \quad \sigma_2, \rho, c_2 \Downarrow \sigma_3, \rho_3}{\sigma_1, \rho_1, c_1 \xrightarrow{\rho} c_2 \Downarrow \sigma_3, \rho_2 \cup \rho_3}$$

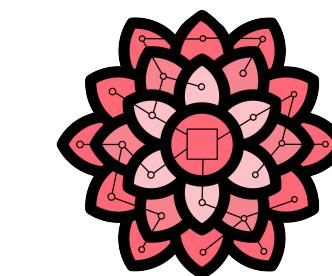
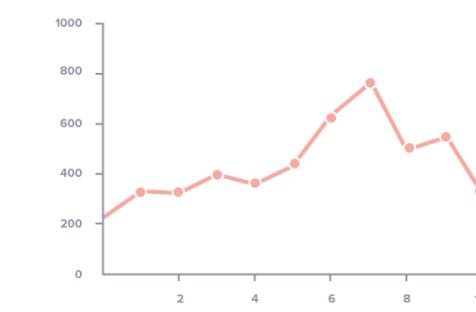
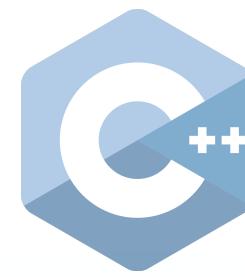
$$\frac{\sigma_1, \rho_1, c_1 \Downarrow \sigma_2, \rho_2 \quad \sigma_2, \rho_2, c_2 \Downarrow \sigma_3, \rho_3}{\sigma_1, \rho_1, c_1 ; c_2 \Downarrow \sigma_3, \rho_3}$$

$$\frac{\sigma_1, \rho_1, e_1 \Downarrow \sigma_2, \rho_2, \text{true} \quad \sigma_2, \rho_2, c_1 \Downarrow \sigma_3, \rho_3}{\sigma_1, \rho_1, \text{if } x \ c_1 \ c_2 \Downarrow \sigma_3, \rho_3}$$

$$\frac{\sigma_1, \rho_1, e_1 \Downarrow \sigma_2, \rho_2, \text{false} \quad \sigma_2, \rho_2, c_2 \Downarrow \sigma_3, \rho_3}{\sigma_1, \rho_1, \text{if } x \ c_1 \ c_2 \Downarrow \sigma_3, \rho_3}$$

$$\frac{\sigma_1, \rho_1, e_1 \Downarrow \sigma_2, \rho_2, \text{true} \quad \sigma_2, \rho_2, c \Downarrow \sigma_3, \rho_3}{\sigma_1, \rho_1, \text{while } x \ c \Downarrow \sigma_3, \rho_3} \quad \frac{\sigma_1, \rho_1, e_1 \Downarrow \sigma_2, \rho_2, \text{false} \quad \sigma_1, \rho_1, \text{while } x \ c \Downarrow \sigma_2, \rho_2}{\sigma_1, \rho_1, \text{while } x \ c \Downarrow \sigma_2, \rho_2}$$

Ada's Journey



Type Systems

Well-typed programs
cannot **go wrong**

Type Systems

Well-typed programs
make **predictable trade-offs**

Design Space Exploration

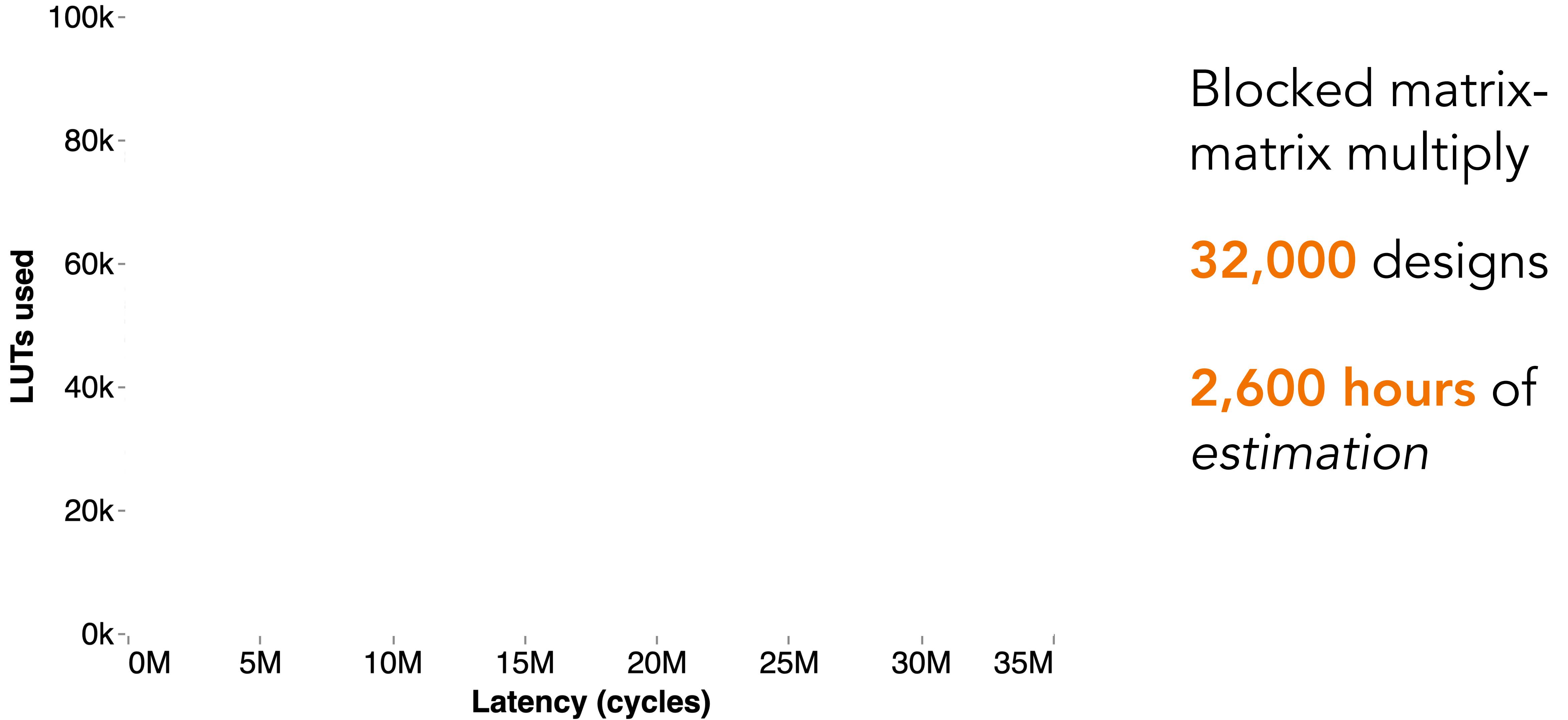
Design Space Exploration

Blocked matrix-matrix multiply

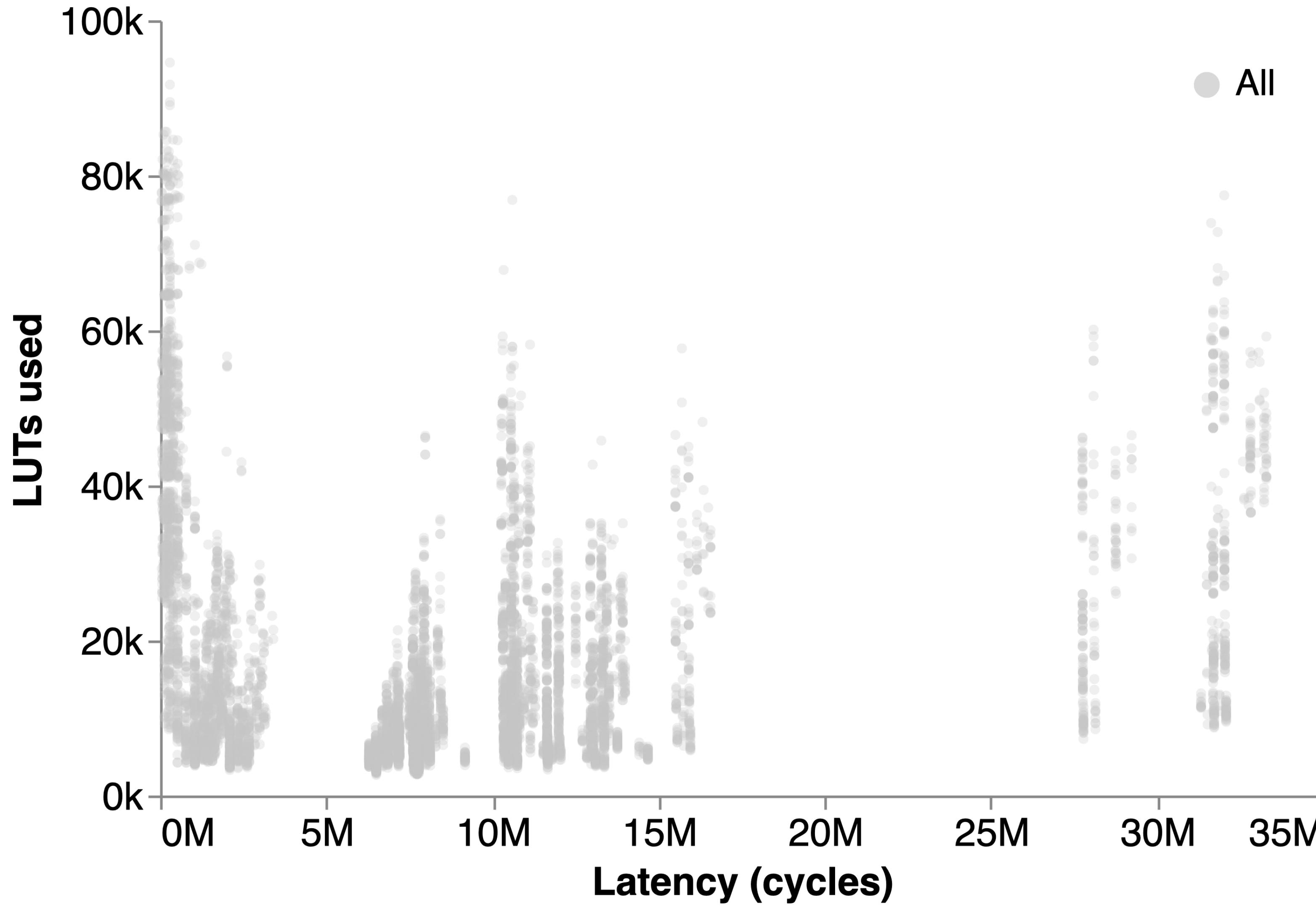
32,000 designs

2,600 hours of estimation

Design Space Exploration



Design Space Exploration

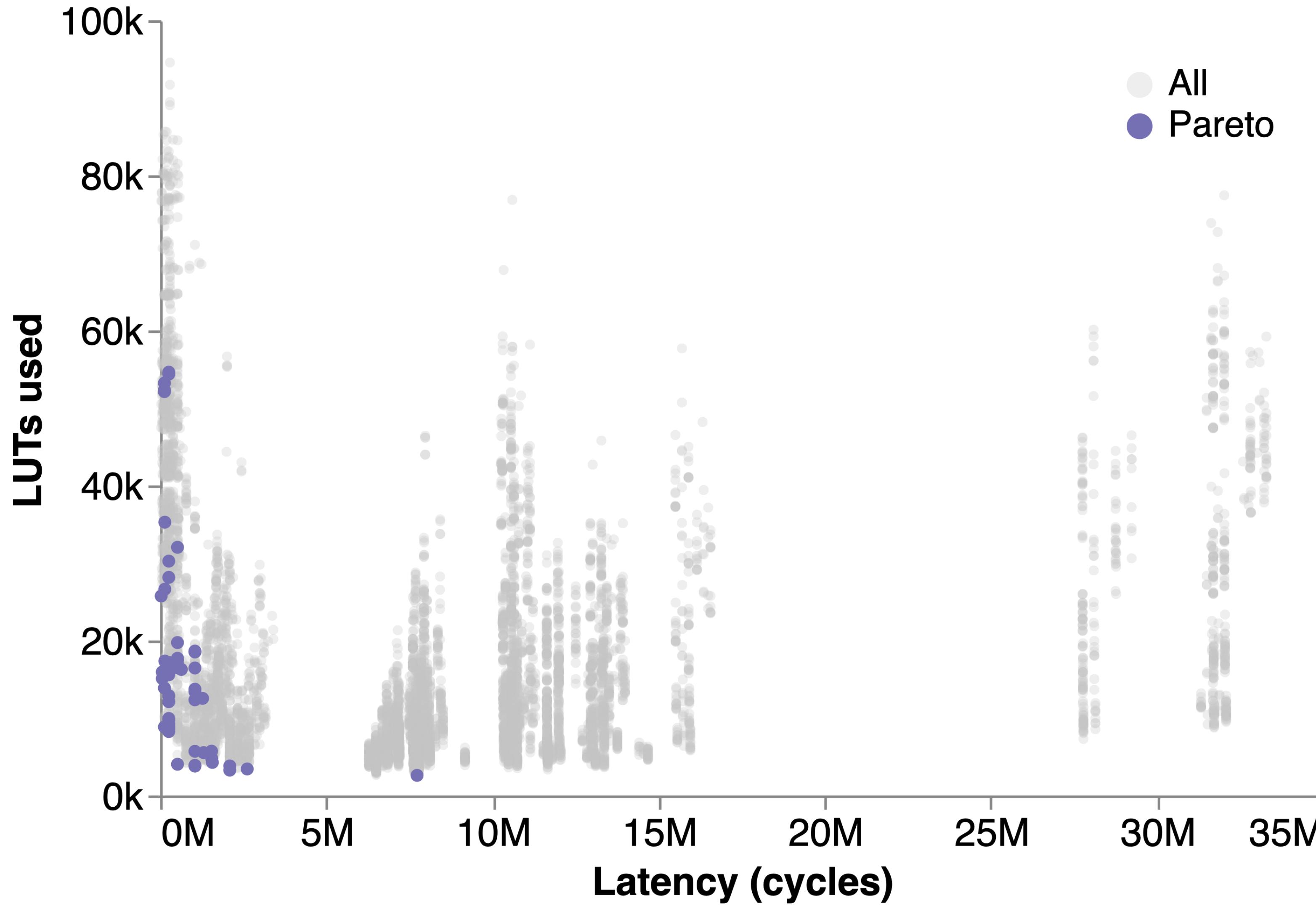


Blocked matrix-matrix multiply

32,000 designs

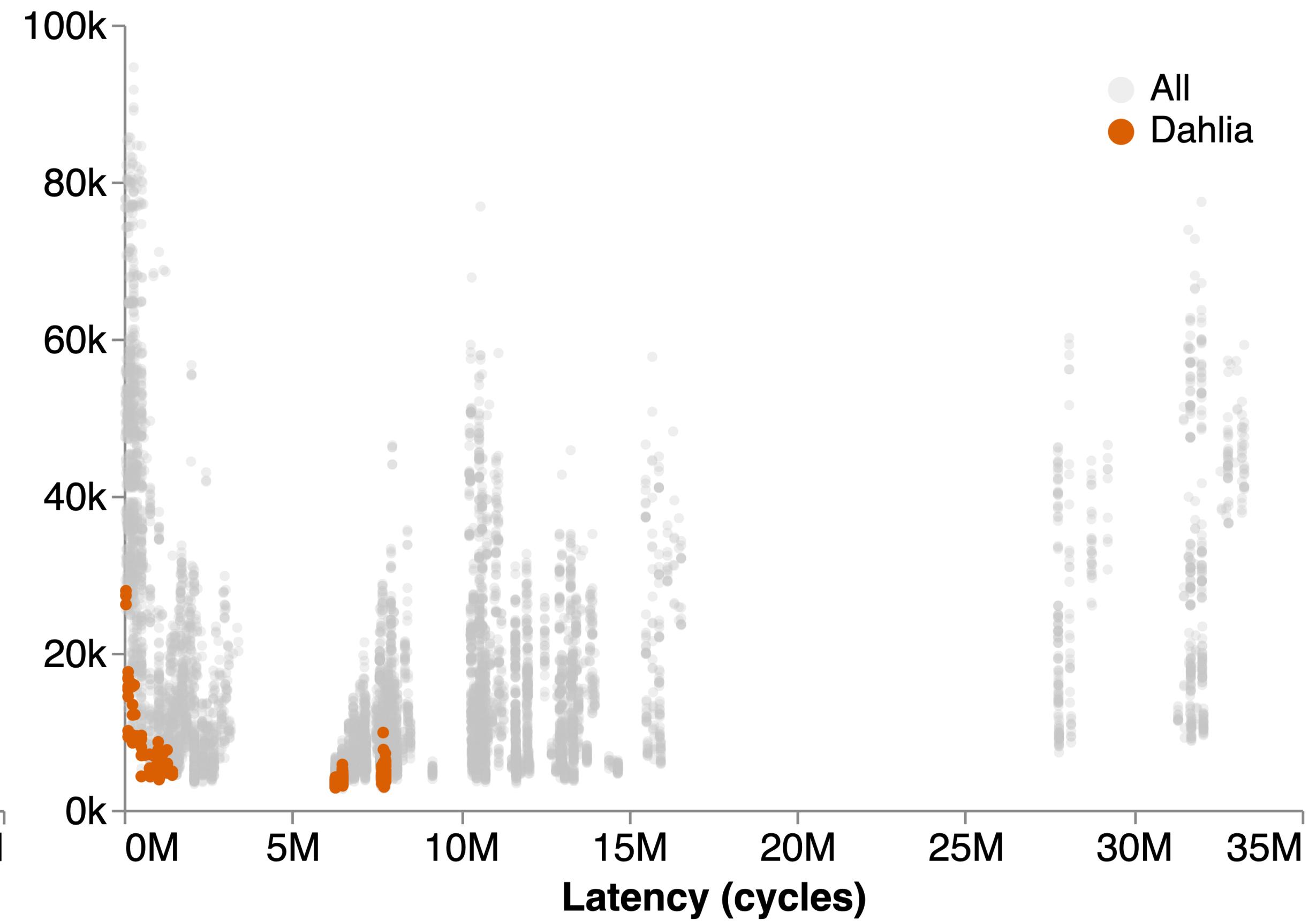
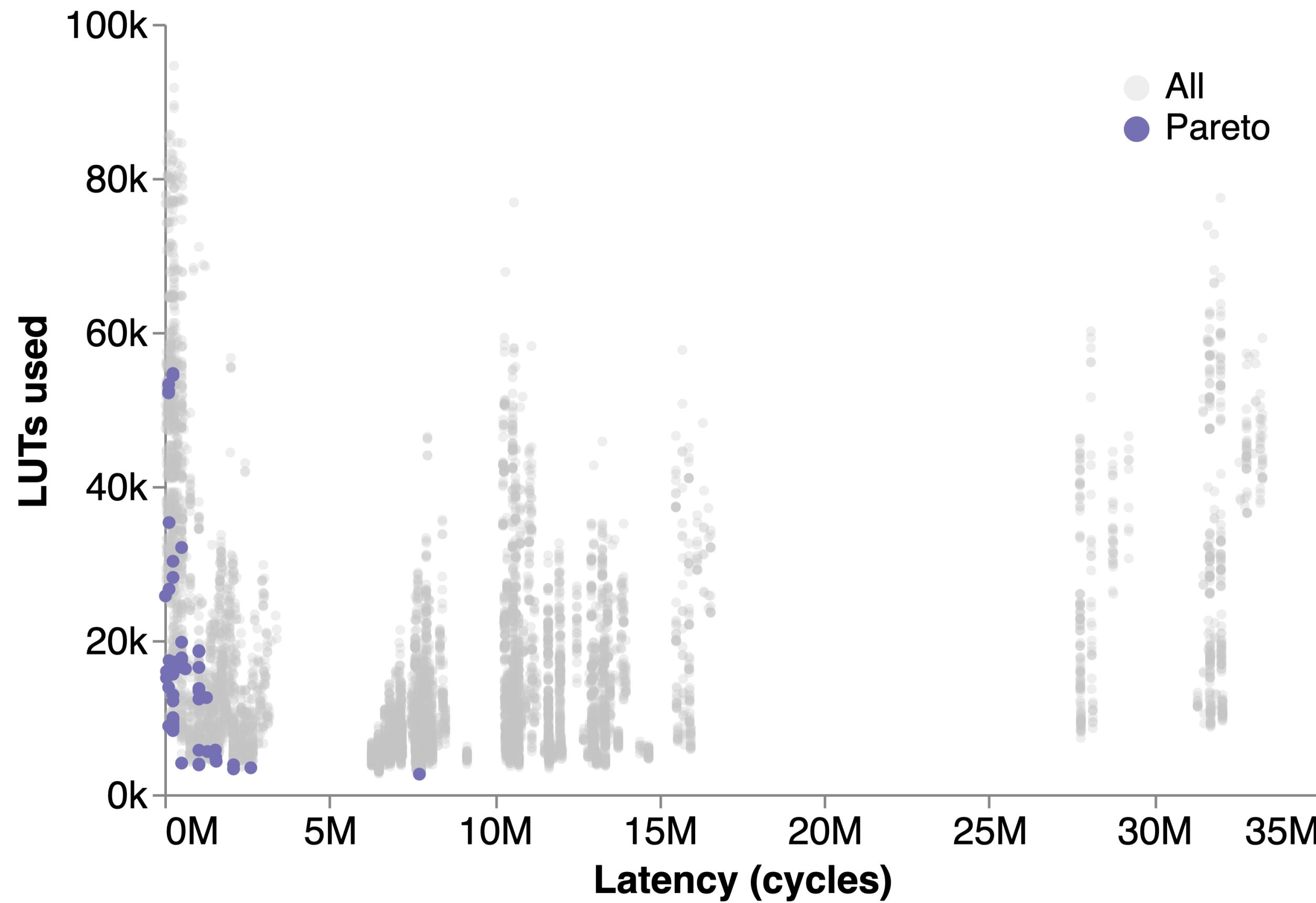
2,600 hours of estimation

Design Space Exploration

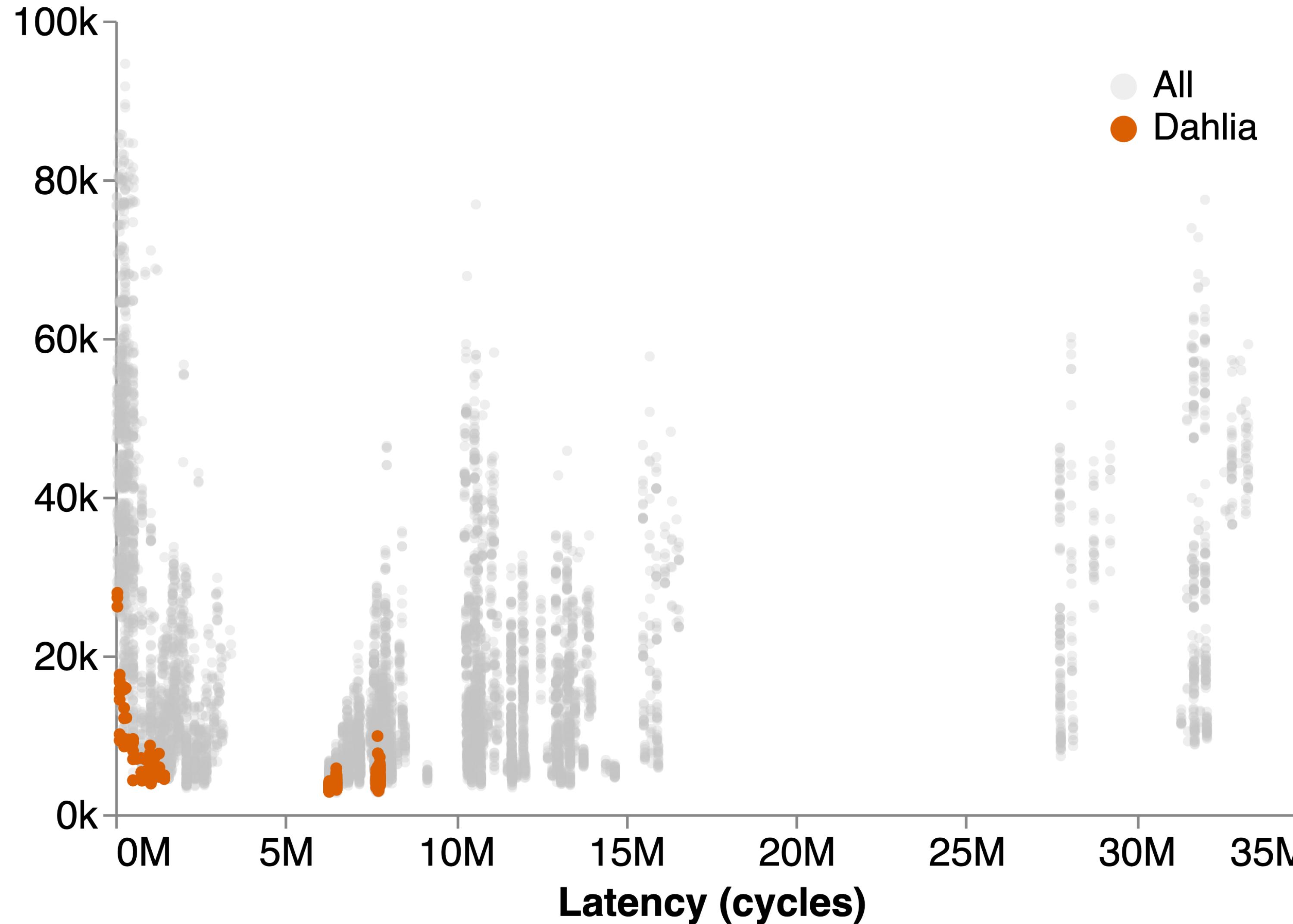


Pareto optimal w.r.t
FPGA resources and
runtime

Design Space Exploration

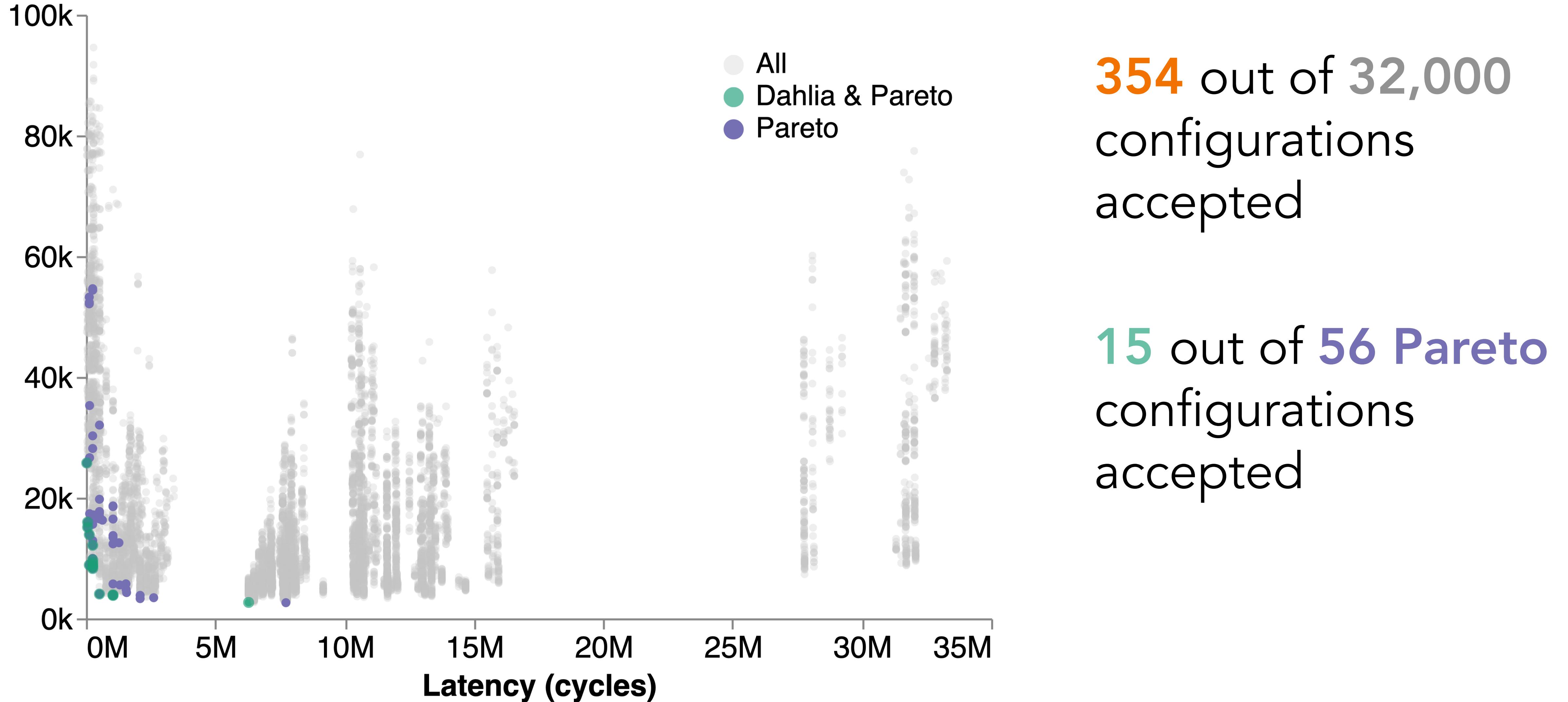


Design Space Exploration

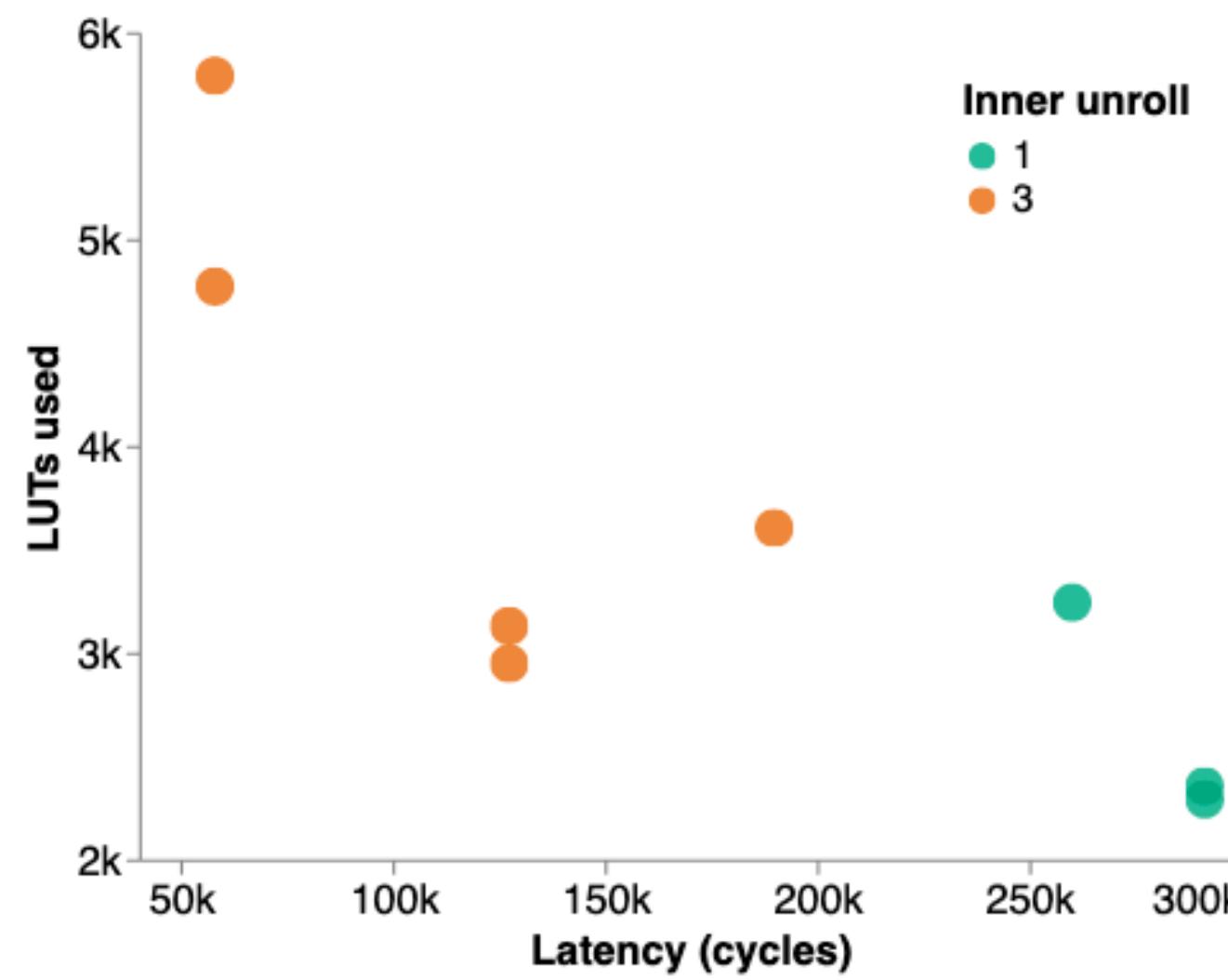


354 out of **32,000**
configurations
accepted

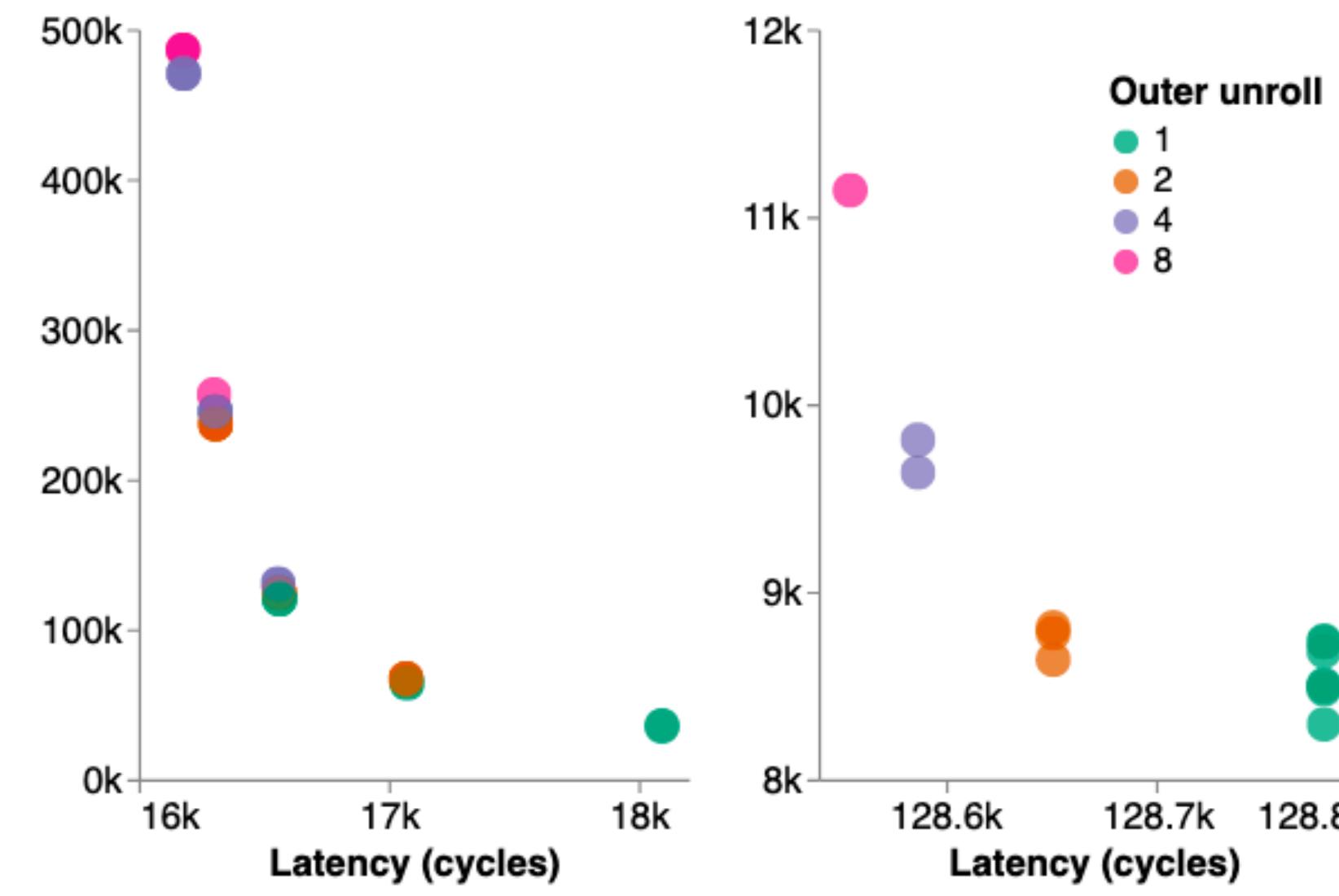
Design Space Exploration



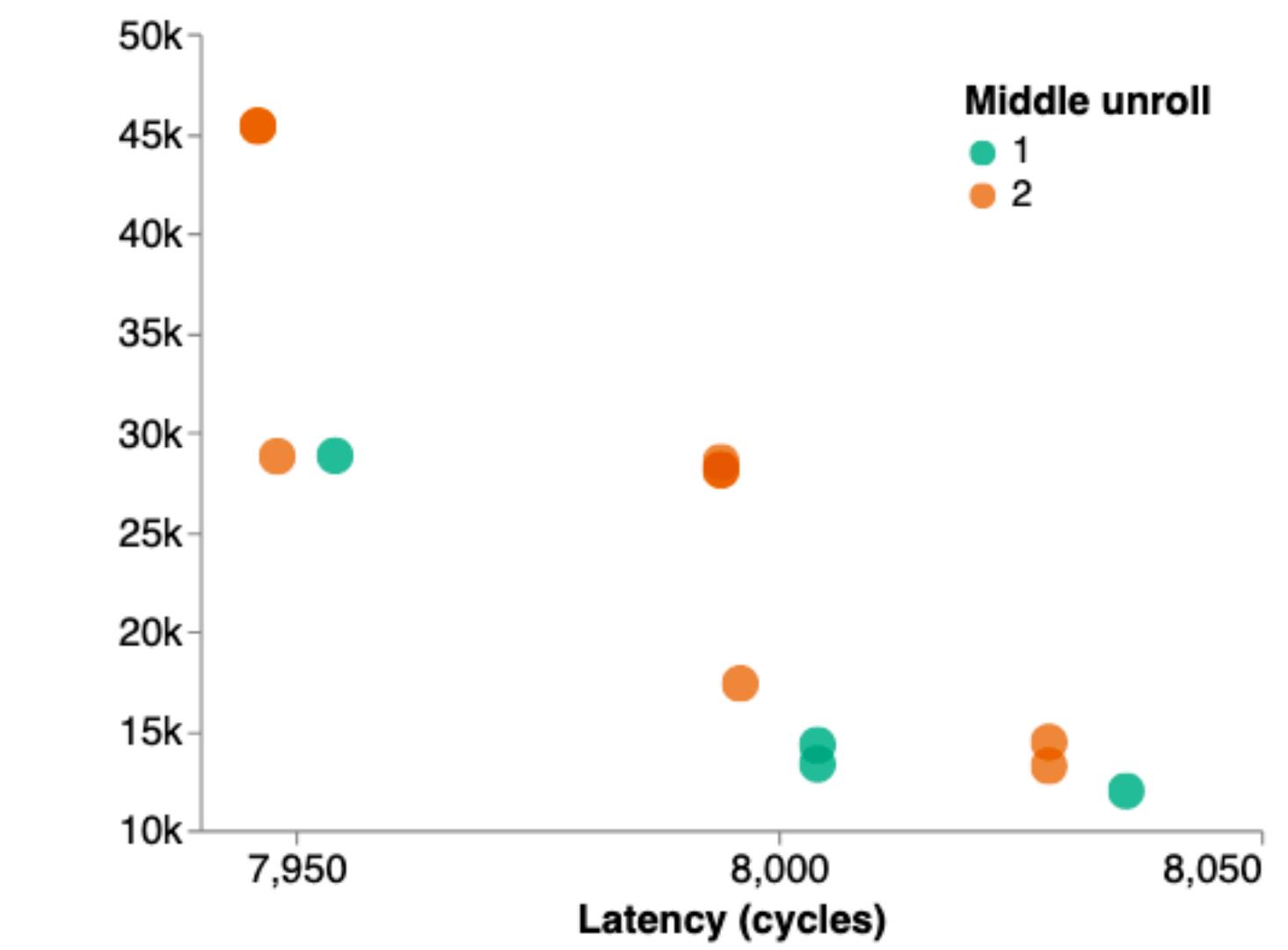
Qualitative Studies



(a) stencil2d with inner unroll.



(b) md-knn with outer unroll.



(c) md-grid with middle unroll.

Qualitative Studies

```
for (r=0; r<row_size; r++)
    for (c=0; c<col_size; c++)
        for (k1=0; k1<3; k1++)
            for (k2=0; k2<3; k2++)
                mul = filter[k1*3 + k2] *
                      orig[(r+k1)*col_size + c+k2]
```

2D-Stencil

Qualitative Studies

```
for (r=0; r<row_size; r++)  
    for (c=0; c<col_size; c++)  
        for (k1=0; k1<3; k1++)  
            for (k2=0; k2<3; k2++)  
                mul = filter[k1*3 + k2] *  
                      orig[(r+k1)*col_size + c+k2]
```

Unroll?

2D-Stencil

Qualitative Studies

```
for (r=0; r<row_size; r++)  
  for (c=0; c<col_size; c++)  
    for (k1=0; k1<3; k1++)  
      for (k2=0; k2<3; k2++)  
        mul = filter[k1*3 + k2] *  
              orig[(r+k1)*col_size + c+k2]
```

2D-Stencil

Unroll?

Partition?

Qualitative Studies

```
for (let row = 0..126)
  for (let col = 0..62)
    view window = shift orig[by row][by col];
    for (let k1 = 0..3)
      for (let k2 = 0..3)
        mul = filter[k1][k2] * window[k1][k2]
```

2D-Stencil

Qualitative Studies

```
for (let row = 0..126)
  for (let col = 0..62)
    view window = shift orig[by row][by col];
    for (let k1 = 0..3)
      for (let k2 = 0..3)
        mul = filter[k1][k2] * window[k1][k2]
```

The diagram illustrates a code transformation process. It starts with a large orange curved arrow pointing from the outermost loops (`row` and `col`) down towards the inner loops (`k1` and `k2`). From the inner loops, two smaller orange arrows point leftwards, indicating the flow of control or data from the inner loop body back up to the loop headers. To the right of the code, the text "Unroll?" is written in orange, with a large orange curved arrow originating from the inner loops and pointing directly at it.

2D-Stencil

Qualitative Studies

The diagram shows a loop nest with several annotations:

- Two orange arrows point from the outermost loops (`for (let row = 0..126)` and `for (let col = 0..62)`) to a large orange circle containing two gray X marks.
- A single orange arrow points from the innermost loop (`for (let k1 = 0..3)`) to another orange circle.
- A single orange arrow points from the innermost loop (`for (let k2 = 0..3)`) to a third orange circle.
- The word "Unroll?" is written in orange next to the first circle.

```
for (let row = 0..126)
  for (let col = 0..62)
    view window = shift orig[by row][by col];
    for (let k1 = 0..3)
      for (let k2 = 0..3)
        mul = filter[k1][k2] * window[k1][k2]
```

2D-Stencil

Qualitative Studies

```
for (let row = 0..126)
  for (let col = 0..62)
    view window = shift orig[by row][by col];
    for (let k1 = 0..3) unroll 3
      for (let k2 = 0..3) unroll 3
        mul = filter[k1][k2] * window[k1][k2]
```

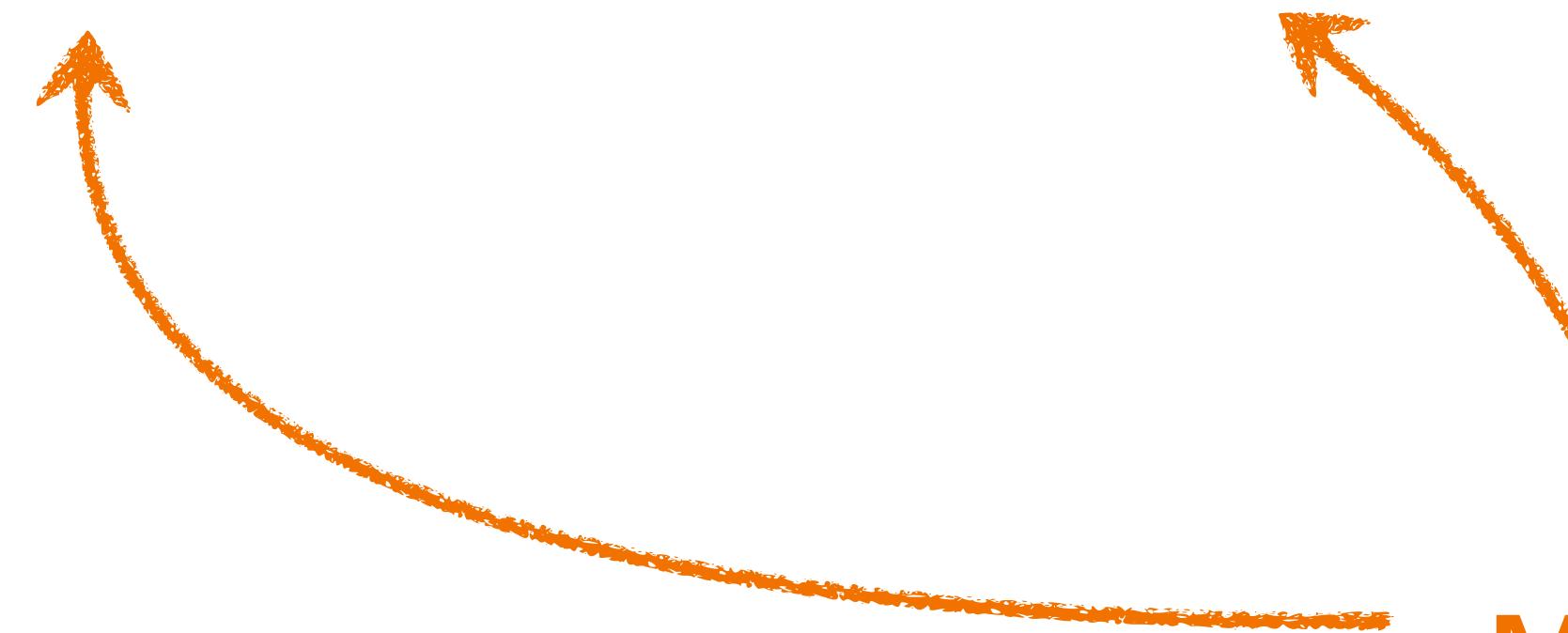
2D-Stencil

Qualitative Studies

```
for (let row = 0..126)
  for (let col = 0..62)
    view window = shift orig[by row][by col];
    for (let k1 = 0..3) unroll 3
      for (let k2 = 0..3) unroll 3
        mul = filter[k1][k2] * window[k1][k2]
```

2D-Stencil

Must partition



The Future

The Future

Resource Polymorphism

```
let m1: float[12 bank M1];
let m2: float[12 bank M2];

for (let i = 0 .. 12) unroll N {
    m2[i] = m1[i] * 2;
}

for (let i = 0 .. 12) unroll K {
    sum += m2[i];
}
```

The Future

Resource Polymorphism

```
let m1: float[12 bank M1];
let m2: float[12 bank M2];

for (let i = 0 .. 12) unroll N {
    m2[i] = m1[i] * 2;
}

for (let i = 0 .. 12) unroll K {
    sum += m2[i];
}
```

The diagram illustrates resource polymorphism. It shows two memory banks, M_1 and M_2 , represented by orange ovals at the top. Two arrows point from these ovals to the `m1` and `m2` variables in the code. Below the code, two more arrows point from the `N` and `K` parameters in the `unroll` annotations to the corresponding loops. This indicates that the compiler is generating code that can handle different numbers of iterations (N and K) by utilizing different memory banks (M_1 and M_2) and unrolling the loops accordingly.

The Future

Resource Polymorphism

```
let m1: float[12 bank M1];
let m2: float[12 bank M2];

for (let i = 0 .. 12) unroll N {
    m2[i] = m1[i] * 2;
}

for (let i = 0 .. 12) unroll K {
    sum += m2[i];
}
```

M₁ = { ... }
M₂ = { ... }
N = { ... }
K = { ... }

The Future

Modularity

The Future

Modularity

```
def dot(m1: float[10], m2: float[10]) {  
    let sum = 0;  
    for (let i = 0 .. 10) {  
        sum += m1[i] * m2[i];  
    }  
    return sum;  
}  
  
dot(A, B);
```

The Future

Modularity

```
def dot(m1: float[10 bank 5], m2: float[10 bank 5]) {  
    let sum = 0;  
    for (let i = 0 .. 10) unroll 5 {  
        sum += m1[i] * m2[i];  
    }  
    return sum;  
}
```

dot(A, B); // A, B need exactly 5 banks

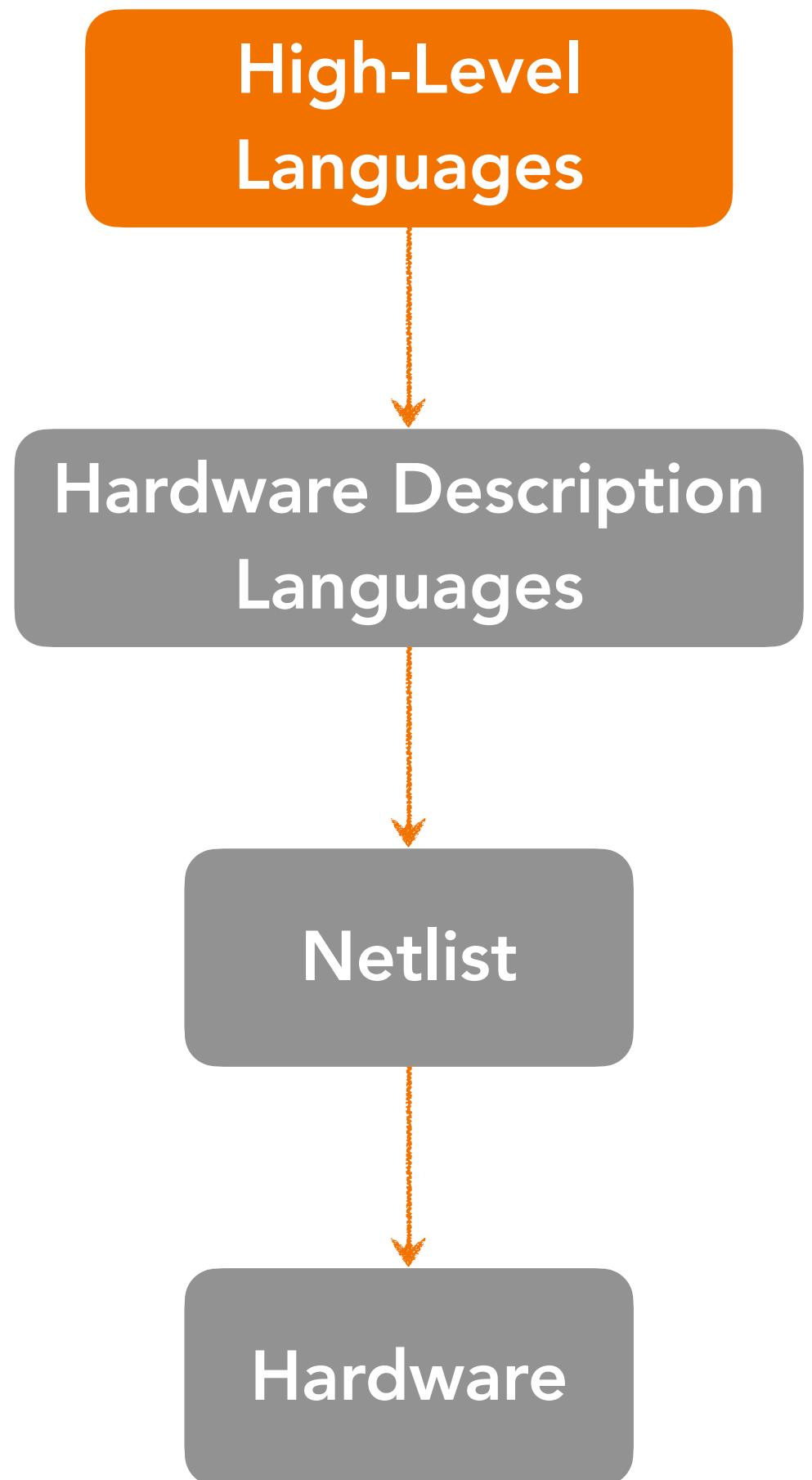
The Future

Modularity

```
def dot(m1: float[10 bank 5], m2: float[10 bank 5]) {  
    let sum = 0;  
    for (let i = 0 .. 10) unroll 5 {  
        sum += m1[i] * m2[i];  
    }  
    return sum;  
}  
  
for (let i = 0..2) unroll 2 {  
    dot(A, B); // How many banks? How many copies of dot?  
}
```

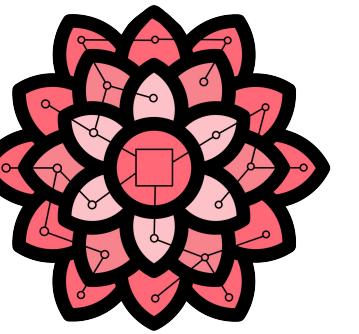
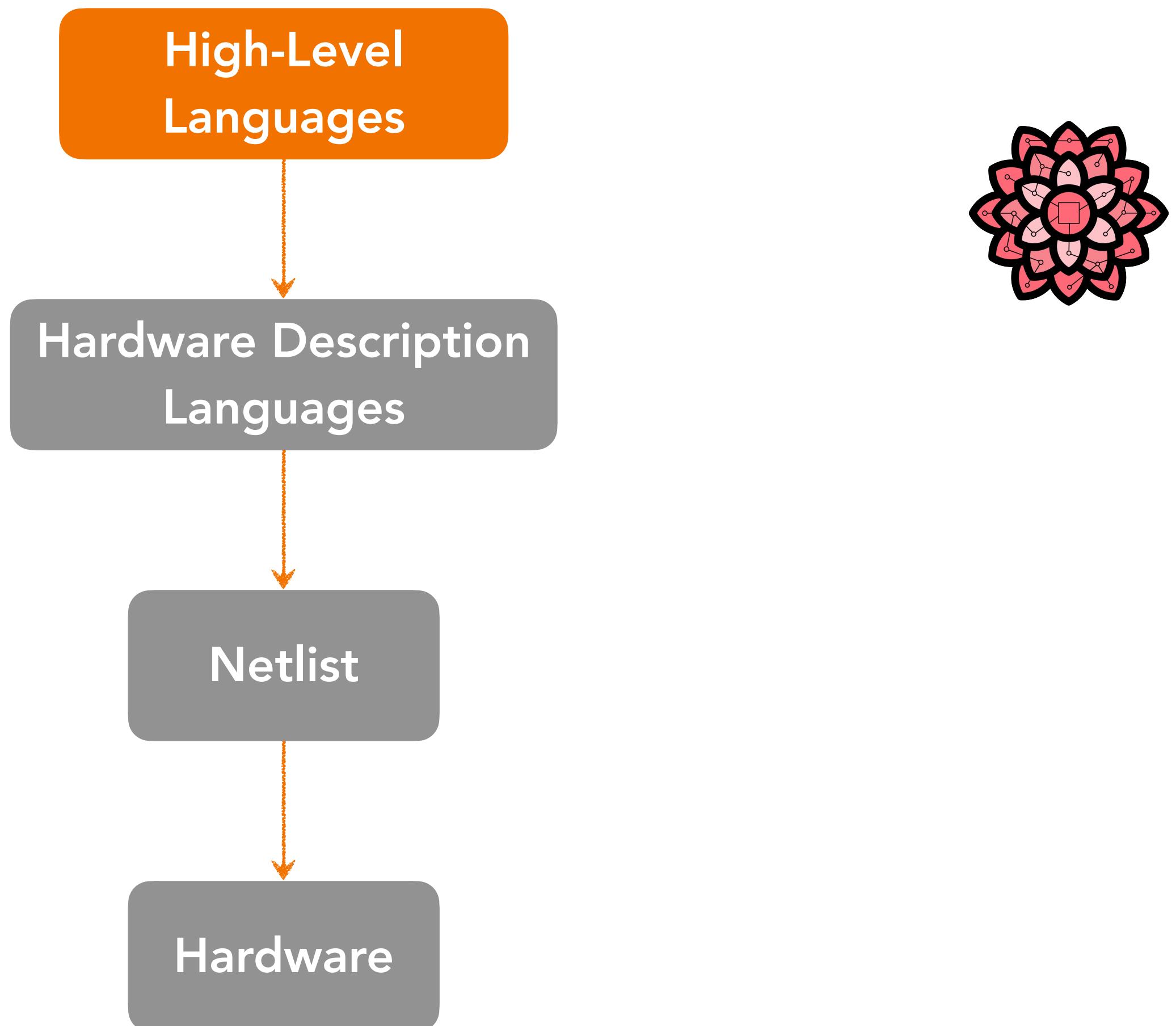
The Future

A Predictable Stack



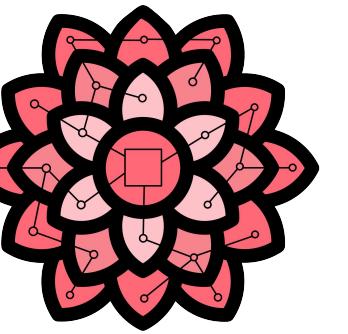
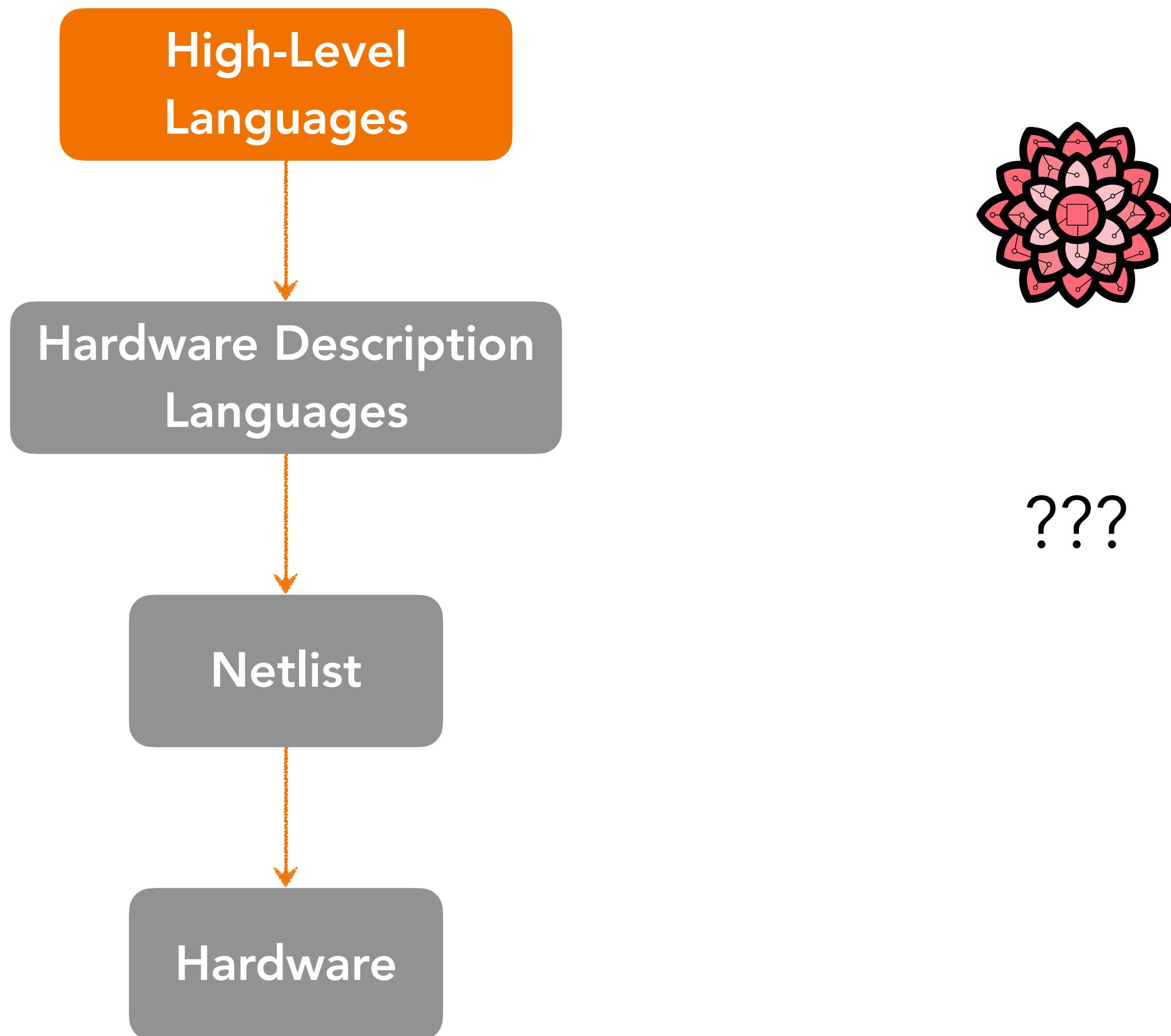
The Future

A Predictable Stack



The Future

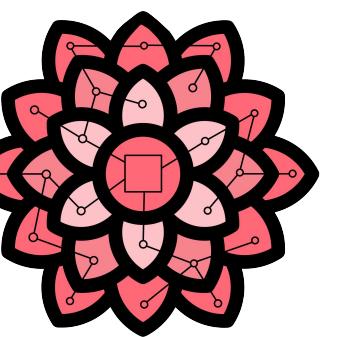
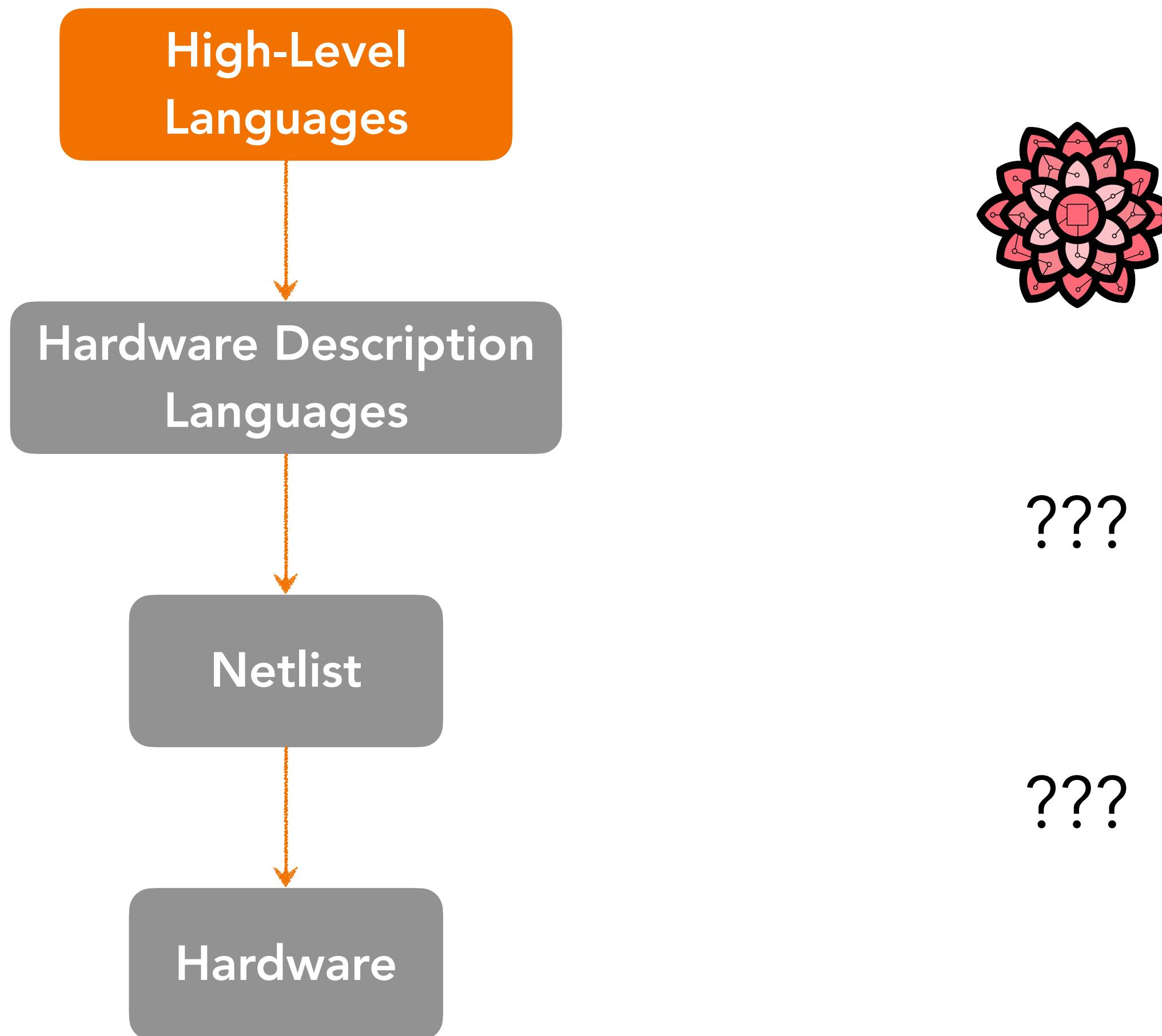
A Predictable Stack



???

The Future

A Predictable Stack



???

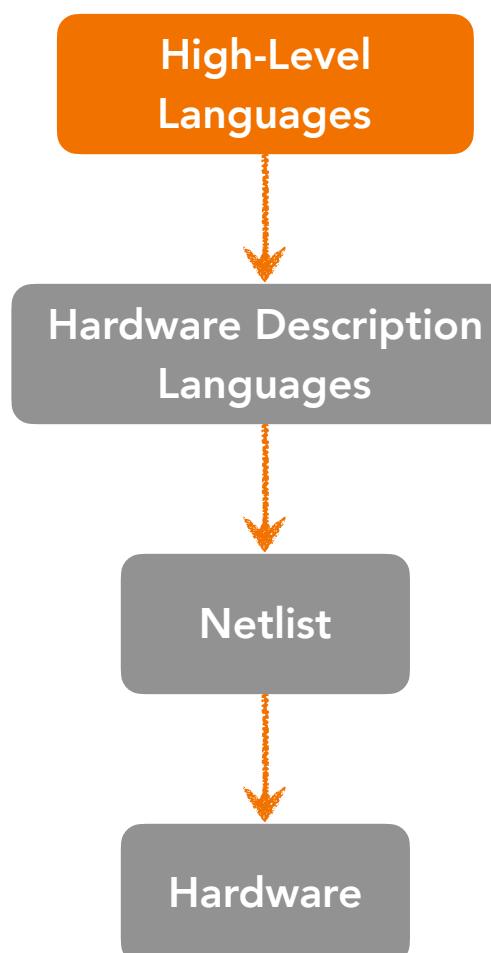
???

The Future

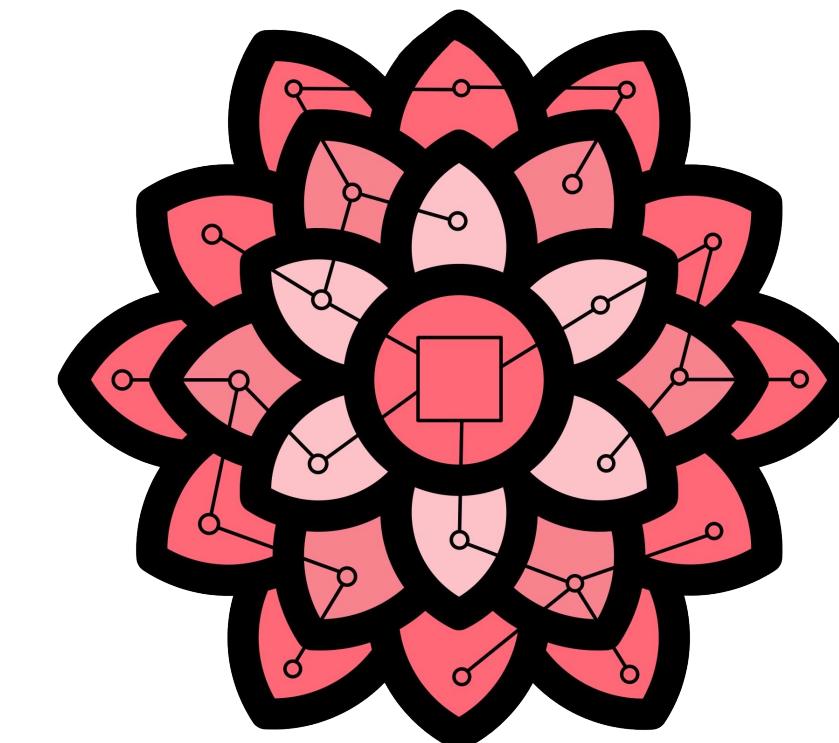
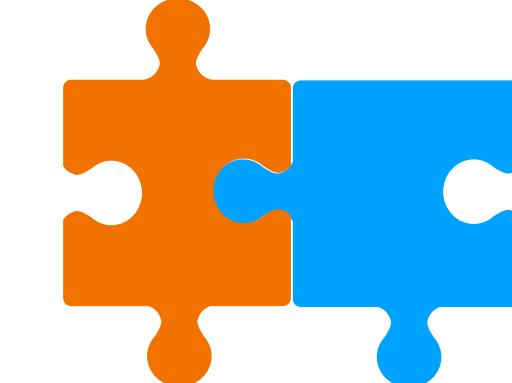
Resource Polymorphism

```
let m1: float[12 bank M];  
let m2: float[12 bank M];  
  
for (let i = 0 .. 12) unroll N {  
    m2[i] = m1[i] * 2;  
}
```

A Predictable Stack



Modularity



Predictability from
languages to LUTs

capra.cs.cornell.edu/dahlia