# CS 217 Lecture 5

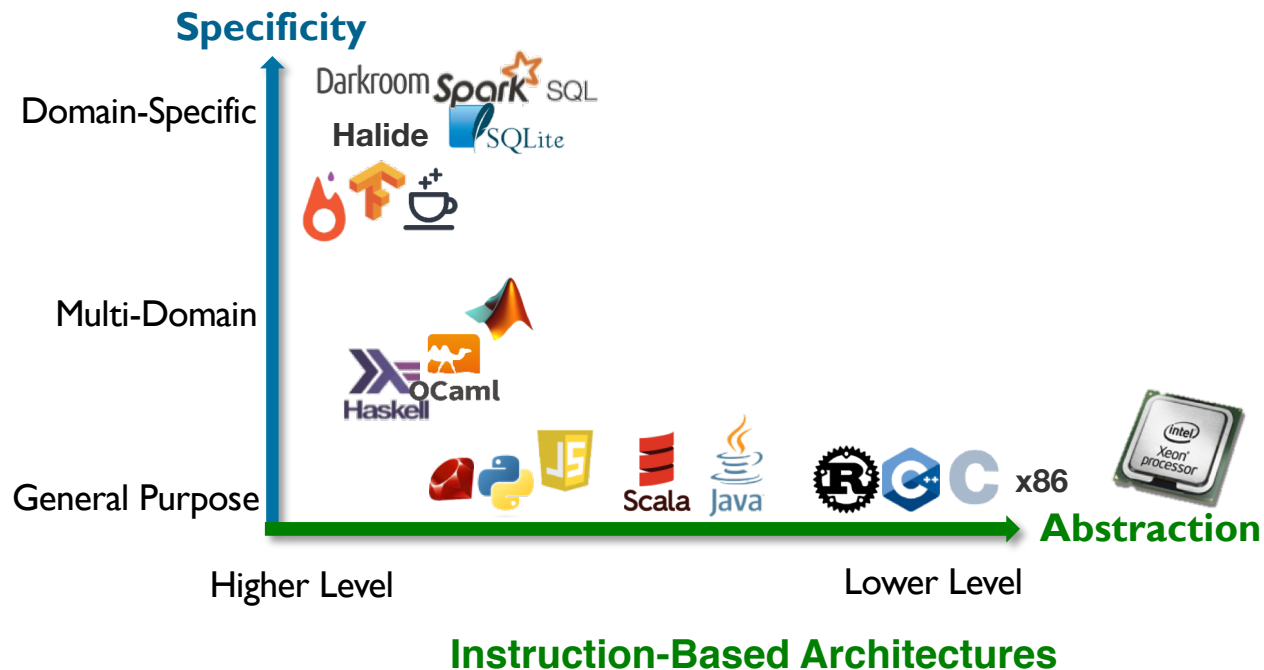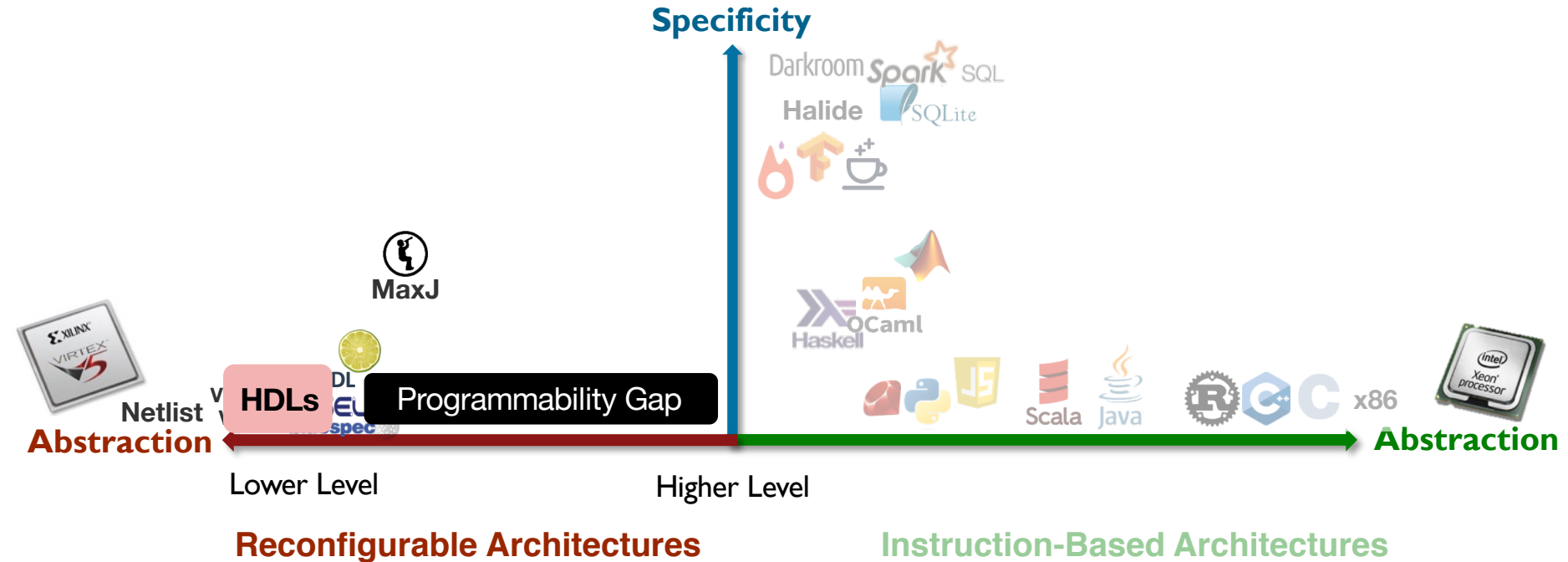## Introduction to Spatial

Fall 2018

# Considerations for Designing ML HW

- Pipelining
    - For better performance, keep as much of the HW as busy as possible at once
- Custom memory hierarchy
    - Lots of room to specialize memory hierarchy for specific access patterns
    - But custom memory hierarchy requires custom partitioning, allocation, etc.
- Finite physical compute and memory resources
    - Modern deep ML models require significant amount of compute and memory
- Huge parameter design spaces
    - Different machine learning applications have application-specific bottlenecks
    - The bottlenecks could be removed by choosing the right set of design parameters
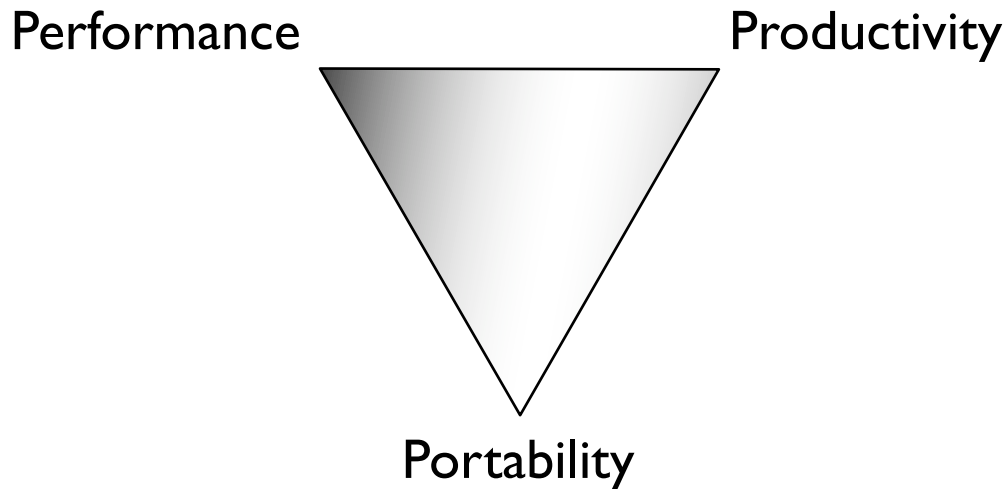
# Language Taxonomy

# Language Taxonomy



**Specificity**

Darkroom  Spark  SQL
Halide  SQLite

MaxJ

Haskell  OCaml

Scala  Java  x86

**HDLs**  Programmability Gap

Netlist

**Abstraction** ←————————————————→ **Abstraction**

Lower Level          Higher Level

**Reconfigurable Architectures**     **Instruction-Based Architectures**

# Language Requirements

1. Performant: generates efficient hardware
2. Productive: High-level language for "power" users
3. Portable: Target-generic source

Performance            Productivity

Portability

# Language Comparisons
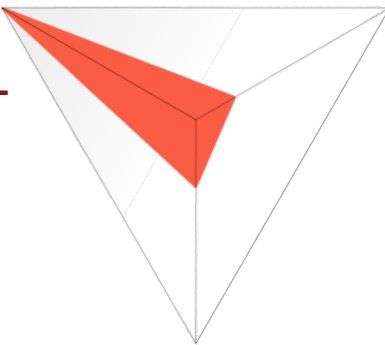
## Hardware Description Languages (HDLs)
e.g. Verilog, VHDL, Chisel, Bluespec
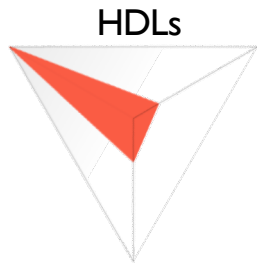
**Performance**

✓ Arbitrary RTL

**Productivity**

✗ No high-level abstractions

**Portability**

✗ Significant target-specific code
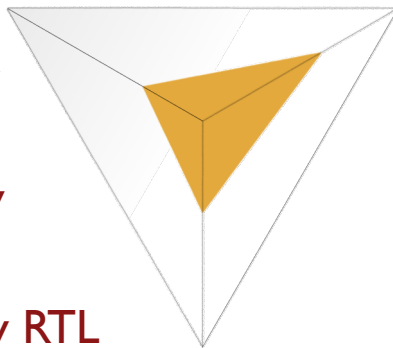
# Language Comparisons

HDLs

## High Level Synthesis Tools (HLS)
### e.g. SDAccel, OpenCL

**Performance**

✗ No memory hierarchy
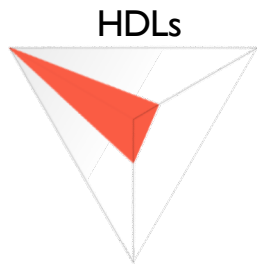
✗ No arbitrary pipelining

✗ No arbitrary RTL

**Productivity**

✓ Nested loops

✗ Difficult to tune

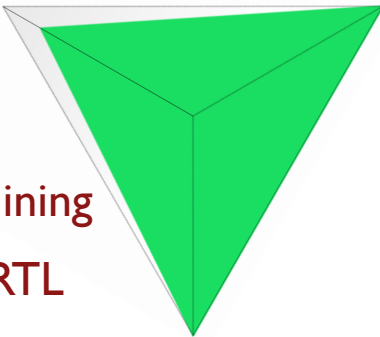✗ Banking, unrolling require pragmas

**Portability**

✓ Single vendor

# Language Comparisons



HDLs

HLS

Spatial

Performance

✓ Memory hierarchy

✓ Arbitrary pipelining

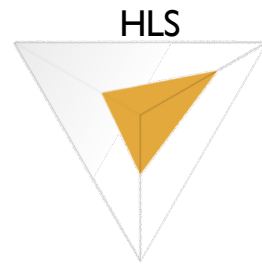✗ No arbitrary RTL

Productivity

✓ Nested loops

✓ Automatic memory banking/buffering

✓ Implicit design parameters (unrolling, banking, etc.)

✓ Automated design tuning

Portability

✓ Target-generic source

# Introducing Spatial

- Programming language to simplify accelerator design

    - Simple APIs to manage Host <-> Accelerator communication

    - Built-in constructs to express parallel datapaths,
      on-chip memories etc.

    - Automatic functional and cycle-accurate simulation

- Focus on "interesting stuff" aka accelerator datapath
  and control design

# Things Spatial (the Language) Has

✓ Nest-able control structures (like software)

✓ Arbitrary pipelining of loops

✓ Memory hierarchy (SRAM vs. DRAM)

✓ Specialized memory types (e.g. FIFO, RegFile)

✓ Streaming abstractions (StreamIn, StreamOut)

✓ High level host/Accelerator interfaces

# Things Spatial's Compiler Does

✓ Retiming (register insertion)

✓ Control scheduling

- Scheduling of inner loops (using dependencies) is done
- Scheduling of outer loops (also dependencies) soon!

✓ Multi-ported memory and buffer inference

- Infers multiple ports for multiple concurrent accesses
- Infers buffer depth for accesses across pipeline stages

✓ Some miscellaneous hardware optimizations

✓ Automated design tuning (updating soon)

# Things Spatial's Compiler Doesn't (Yet)

❌ Automatic loop/data tiling

- Splitting loops up
- Inference of DRAM/SRAM + transfers

❌ Automatic work allocation

- Breaking computation up into host/Accel regions

❌ Loop fusion

- Merging of loops which have element-wise dependencies
- Elimination of memories between such loops

# SPATIAL (AND SCALA) BASICS

# Hello Spatial!

```scala
1  import spatial.dsl._
2
3  @spatial
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]): Unit = {
8      val input = args(0).to[Int]
9      val in  = ArgIn[Int]
10     val out = ArgOut[Int]
11     setArg(in, input)
12     Accel {
13       out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

# **Spatial** is an *Embedded DSL* in **Scala**

```scala
1  import spatial.dsl._
2
3  @spatial
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]):
8      val input = args(0).to[Int]
9      val in  = ArgIn[Int]
10     val out = ArgOut[Int]
11     setArg(in, input)
12     Accel {
13       out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

An **embedded DSL** can be thought of as an automatically optimized **Scala** library.

# Scala Basics

```scala
1  import spatial.dsl._
2
3  @spatial
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]): Unit = {
8      val input = args(0).to[Int]
9      val in  = ArgIn[Int]
10     val out = ArgOut[Int]
11     setArg(in, input)
12     Accel {
13       out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

Semicolons are optional

# Import Statements

```
1 import spatial.dsl._
2
3
4 object HelloSpatial extends SpatialApp {
5
6
7
8
9
10     val out = ArgOut[Int]
11     setArg(in, input)
12     Accel {
13         out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

Same in every Spatial program
(Similar idea to **#include** in C,
Identical to **import** in Java, Python)

# Application Object Declaration

```
1  import spatial.dsl._
2
3  @spatial
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]): Unit = {
```

Spatial applications are always `objects`

```
10     val out = ArgOut[Int]
11     setArg(in, input)
12     Accel {
13       out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

# Application Object Declaration

```
1  import spatial.dsl._
2
3  @spatial
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String])
         .to[
         ]
10     val out = ArgOut[Int]
11     setArg(in, input)
12     Accel {
13       out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

Name of application

All Spatial applications inherit from ("extends") **SpatialApp**

# Spatial's Entry Function: "main()"

```
1  import spatial._
2
3  @spatial
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[      Spatial's entry function
8      val input = args(0).to[Int]
9      val in  = ArgIn[Int]
10     val out = ArgOut[Int]
11     setArg(in, input)
12     Accel {
13       out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

# Spatial's Entry Function: "main()"

```scala
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]): Unit = {
8      val input = args(0).to[Int]
9      val in   = ArgIn[Int]
         = ArgOut
              n, inp
              in + 
      }
15     println("Output: " + getArg(out))
16    }
17  }
18
```

Starts a function declaration

Function return type (Unit: same as void)

# Val Definitions

```
1  import spatial.dsl._
2
...
```

Declares an **immutable** value named "input" (value can't be modified later)

```
7
8    def main(args: Array[String]): Unit = {
9      val input = args(0).to[Int]
10     val in  = ArgIn[Int]
11     val out = ArgOut[Int]
12     setArg(in, input)
13     Accel {
14       out := in + 4
15     }
16     println("Output: " + getArg(out))
17   }
18 }
```

# Val Definitions

```
1  import spatial.dsl._
2
3
                                              {
   Value types are optional in Scala.

6
7
8    def main(args: Array[String]): Unit = {
9      val input: Int = args(0).to[Int]
10     val in  = ArgIn[Int]
11     val out = ArgOut[Int]
12     setArg(in, input)
13     Accel {
14       out := in + 4
15     }
16     println("Output: " + getArg(out))
17   }
18 }
```

# Val Definitions

```
1  import spatial.dsl.

6
7
8    def main(args: Array[String]): Unit = {
9      val input = args(0).to[Int]
10     val in  = ArgIn[Int]
11     val out = ArgOut[Int]
12     setArg(in, input)
13     Accel {
14       out := in + 4
15     }
16     println("Output: " + getArg(out))
17   }
18 }
```

Scala is statically typed (like C, Java)
Without the ": **Int**", the type of this
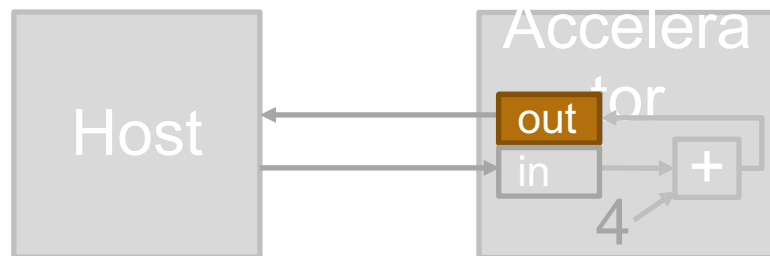value is **inferred** by the compiler.

# Method Calls

```
1  import spatial.dsl._
2
3
7
8    def main(args: Array[String]): Unit = {
9      val input = args(0).to[Int]
10     val in  = ArgIn[Int]
11     val out = ArgOut[Int]
12     setArg(in, input)
13     Accel {
14       out := in + 4
15     }
16     println("Output: " + getArg(out))
17   }
18 }
```

**Round brackets ( )** for value parameters
**Square brackets [ ]** are for **type** parameters

# Spatial Command-Line Arguments

```
1  import spatial.dsl._
2
3
4
5
6
7
8    def main(args: Array[String]): Unit = {
9      val input = args(0).to[Int]
10     val in  = ArgIn[Int]
11     val out = ArgOut[Int
12
13     A
14
15     }
16     println("Output: " + getArg(out))
17   }
18 }
```

Spatial app's command-line arguments

Conversion from **String** to **Int**

Host

```c
int main(int argc, char **argv) {
  int in = atoi(argv[1]);

  printf("Output: %d\n", out);
  return 0;
}
```

# Input Arguments (ArgIn)

```
1  import spatial.dsl._
2
3
4
5
6
7
8  def main(args: Array[String]): Unit = {
9    val input = args(0).to[Int]
10   val in  = ArgIn[Int]
11   val out = ArgOut[Int]
12   setArg(in, input)
13   Accel {
14     out := in + 4
15   }
16   println("Output: " + getArg(out))
17  }
18 }
```

Creates a new register to capture a scalar argument *from* the CPU



Host

Accelerator

out

in

+

4

# Output Arguments (ArgOut)

```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
7
8
9      val input = args(0).to[Int]
10     val in   = ArgIn[Int]
11     val out = ArgOut[Int]
12     setArg(in, input)
13     Accel {
14       out := in + 4
15     }
16     println("Output: " + getArg(out))
17   }
18 }
```

Creates a new scalar argument *to* the CPU *from* the Accelerator

Host

Accelerator

out

in

+

4

# Scalar Transfers (CPU → Accelerator)

# Accel Block

```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
11     setArg(in, input)
12     Accel {
13        out := in + 4
14     }
15     println("Output: " + getArg(out))
16  }
17 }
18
```

Defines an Accelerator computation scope. Everything in here goes on the Accelerator

Host

Accelerator

out

in

+

4

# Acceleratable Code

```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
7
8
9
10
11
12    Accel {
13      out := in + 4
14    }
15    println("Output: " + getArg(out))
16  }
17 }
18
```

The types of operations done in this scope are limited to **acceleratable (synthesizable)** Spatial

Host

Accelerator

out

in

+

4

# Accel Block

```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
...
12      Accel {
13        out := in + 4
14      }
15      println("Output: " + getArg(out))
16    }
17  }
18
```

**Accel** handles control signals for you.
It implicitly creates:
- a **start signal** (CPU → Accelerator)
- a **done signal** (Accelerator → CPU)

Host

done

start

Accelera tor

out

in

+

4

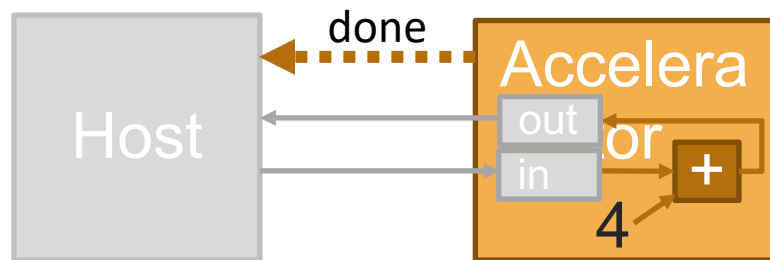# Implicit Register Reads



```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]): Unit = {
8
9
10
11     setArg(in, input)
12     Accel {
13        out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```
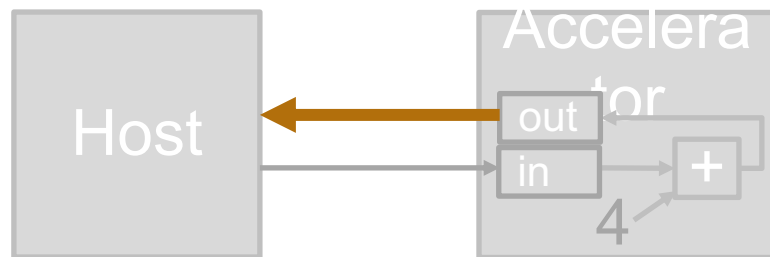
Implicitly creates a wire from the output of register (ArgIn) **in**

Host
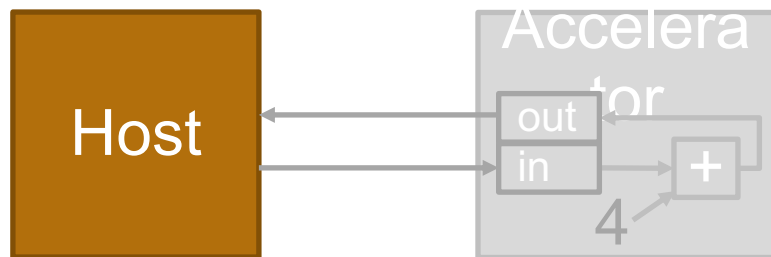
Accelerator

out

in

+

4

33

# Implicit Register Reads

```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]): Unit = {
8
9
10
11     setArg(in, input)
12     Accel {
13       out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

:= connects the value **in + 4**
to the input of the register **out**

Host

Accelerator

out

in

+

4

# Accel Block Scheduling

```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
7
8
9
10
11     setArg(in, input)
12     Accel {
13       out := in + 4
14     }
15     println("Output: " + getArg(out))
16  }
17 }
18
```

> **Accel** guarantees that Accelerator execution completes after all operations in this block complete

```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]): Unit = {
8      val input = args(0).to[Int]
9
10
11
12
13
14   }
15   println("Output: " + getArg(out))
16   }
17 }
18
```

Gets the value of the ArgOut **out** from the Accelerator back to the CPU

Host

Accelerator

out

in

+

4

# Printing in Spatial**

```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]): Unit = {
8      val input = args(0).to[Int]
9      val in   = ArgIn[Int]
10
11
12     Accel {
13        out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

Prints the output to the terminal

** Printing in Spatial isn't synthesizable, but it can be used in **host code** and in **debugging** (more later)



Host

Accelerator

out

in

+

4

```
int main(int argc, char **argv) {
  int in = atoi(argv[1]);
  …
  printf("Output: %d\n", out);
  return 0;
}
```

# Hello Spatial!

```
1  import spatial.dsl._
2
3
4  object HelloSpatial extends SpatialApp {
5
6
7    def main(args: Array[String]): Unit = {
8      val input = args(0).to[Int]
9      val in  = ArgIn[Int]
10     val out = ArgOut[Int]
11     setArg(in, input)
12     Accel {
13       out := in + 4
14     }
15     println("Output: " + getArg(out))
16   }
17 }
18
```

# SIMPLE TYPES

# Custom Types in Spatial

■ Now what if we want an ArgIn value that isn't an Int?

■ Other options:

    ■ Custom fixed point types

    ■ Custom floating point types

    ■ Structs

    ■ Vectors

# Custom Types

```
1
2
3  val input = args(0).to[Int]
4  val in   = ArgIn[Int]
5
6  setArg(in, input)
7
8
9
10
11
12
13
14
15
16
17
18
```

# Custom Fixed Point Types

```
1  type Q8_8 = FixPt[FALSE,_8,_8]
2
3
4
```

_**N** = # of fraction bits
(N from 0 to 128)

**TRUE** = Signed
**FALSE** = Unsigned

_**N** = # of integer bits
(N from 1 to 128)

```
0b00000000.00000000
```

Integer bits    Fraction bits

# Custom Fixed Point Examples

```
1  type Q8_8 = FixPt[FALSE,_8,_8]
2
3  type UInt8 = FixPt[FALSE,_8,_0]
4
5  type LongLong = FixPt[TRUE,_128,_0]
```

0b00000000.00000000

Integer bits    Fraction bits

# Custom Fixed Point Types

```
1  type UInt8 = FixPt[FALSE,_8,_0]
2
3  val input = args(0).to[UInt8]
4  val in   = ArgIn[UInt8]
5
6  setArg(in, input)
7
8
9
10
11
12
13
14
15
16
17
18
```

# Custom Floating Point Types

```
1  type Float = FltPt[_23,_11]
2
3
4
```

**_N** = # of significand bits + 1
(N from 1 to 128)
**Includes sign bit!**

**_N** = # of exponent bits
(N from 0 to 128)

```
0 00000000 x 2^00000000
```

Sign bit          Significand bits          Exponent bits

45

# Custom Floating Point Types

```
1  type Half = FltPt[_11,_5]
2
3  val input = args(0).to[Half]
4  val in   = ArgIn[Half]
5
6  setArg(in, input)
```

# Predefined Type Aliases

```
1  type Char  = FixPt[TRUE,_8,_0]
2  type Short = FixPt[TRUE,_16,_0]
3  type Int   = FixPt[TRUE,_32,_0]
4  type Long  = FixPt[TRUE,_64,_0]
5
6  type Half   = FltPt[_11,_5]   // 754 Half
7  type Float  = FltPt[_24,_8]   // 754 Single
8  type Double = FltPt[_53,_11]  // 754 Double
9
10
11
12
13
14
15
16
17
18
```

# Note About Booleans

```
1
2
3  val input = args(0).to[Boolean]
4  val in   = ArgIn[Boolean]
5
6  setArg(in, input)
```

**Note**: For API purposes, Boolean is NOT the same as single bit fixed point number

Uses "false" and "true" rather than 0 and 1

# Custom Structs

```
1  @struct class MyStruct(
2    red:   Int,
3    green: Int,
4    blue:  Int
5  )
6
```

Declares a new Struct type with the given list of fields

# Custom Structs

```
1  @struct class MyStruct(
2    red:   Int,
3    green: Int,
4    blue:  Int
5  )
6
7  val bus = MyStruct(147, 192, 125)
8
9
10
11
12
13
14
15
16
17
18
```

147   125
  192

$\searrow 32$  $\searrow 32$  $\searrow 32$

$\searrow 96$

bus

Allocates an instance of the struct.
**Note**: NO *new* keyword used

In hardware, a struct instance is just
a concatenation of wires

# Custom Structs



```
1  @struct class MyStruct(
2    red:   Int,
3    green: Int,
4    blue:  Int
5  )
6
7  val bus = MyStruct(147, 192, 125)
8
9  val red = bus.red
10 val blue = bus.blue
11
12
13
14
15
16
17
18
```

Creates a reference to the struct field (equivalent to a bit slice)

51

# Custom Structs



```
1  @struct class MyStruct(
2    red:    Int,
3    green:  Int,
4    blue:   Int
5  )
6
7  val bus = MyStruct(147, 192, 125)
8
9  bus.blue = 45
10
11
12
13
14
15
16
17
18
```

**Note**: Allocated structs are immutable!
We can't write to them or change the contents!

# Nesting Structs

```
 1 @struct class RGB(
 2   red:    Int,
 3   green:  Int,
 4   blue:   Int
 5 )
 6
 7 @struct class RGBA(
 8   rgb:    RGB,
 9   alpha:  Int
10 )
11
12
13
14
15
16
17
18
```

# Registers of Custom Types

```
1  @struct class MyStruct(
2    red:    Int,
3    green:  Int,
4    blue:   Int
5  )
6
7  val in  = ArgIn[MyStruct]
8
9  in.red
10
11
12
13
14
15
16
17
18
```

Creates an ArgIn register which holds a value of type MyStruct

**Note**: Registers can hold structs as long as the fields are primitive values (FixPt, FltPt, Boolean) or other primitive-based structs

# Bit Casting

11.25

| 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|

bits: 11        6 5        0

data

720

| 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|

bits: 11            0

number

```
1  type UInt12 = FixPt[FALSE,_12,_0]
2  type Q6     = FixPt[FALSE,_6,_6]
3
4  val data = random[Q6]
5
6  val number = data.as[UInt12]
7
```

Generates a random 12-bit number, interpreted as a fixed point value with 6 integer and 6 fraction bits

Creates a view of these bits directly as an unsigned 12 bit integer

# Bit Casting



```
1  val data = random[Q6]
2
3  type UInt3 = FixPt[FALSE,_3,_0]
4  @struct class MyStruct(        ←   Bit-composed struct type
5    w0: UInt3,
6    w1: UInt3,
7    w2: UInt3,
8    w3: UInt3
9  )
10
11 val instance = vector.as[MyStruct]
12
```

Creates a view of these bits
directly as an instance of MyStruct

# All Bit Primitives Are Bit Vectors

```
 1 type UInt32 = FixPt[FALSE,_16,_0]
 2 type UInt16 = FixPt[FALSE,_16,_0]
 3 @struct class Split16(msByte: UInt16, lsByte: UInt16)
 4
 5 val a = 10.to[UInt32]
 6
 7 val bit3 = a(3)
 8
 9 val lsByte = a(17::2).as[UInt16]
10
11 val bits = a.as32b
12
13 val split = a.as[Split16]
14
15 val a_again = split.as[UInt32]
16
17
18
```

4th least significant bit of *a*

bit slice of *a*

vector view of all bits in *a*

**Split16** view of *a*

57

# OFF-CHIP MEMORIES

# Basic Machine Model

# DRAM

```
1  val data = DRAM[Int](192, 192)
2
3  Accel {
4    …
5  }
6
7
8
9
10
11
12
13
14
15
16
17
18
```

Declares a 192 x 192 (2D) block of memory in the Accel DRAM with words of type **Int**

**1 to 5 dimensions** are currently supported

Accelerator

Accel DRAM

data

# DRAM

```
1  val data = DRAM[Int](192, 192)
2
3  Accel {
4    …
5  }
6
7
8
9
10
11
12
13
14
15
16
17
18
```

Must be used **outside** of the *Accel* scope

Accelerator

Accel DRAM

data

```
 1  val N = ArgIn[Int]
 2  setArg(N, args(0).to[Int])
 3  val data = DRAM[Int](N)
 4
 5  val armData: Array[Int] = … //Info soon!
 6
 7  setMem(data, array)
 8
 9  Accel {
10    …
11  }
12
13
14
15
16
17
18
```

Copies data from CPU DRAM to Accel DRAM

| CPU (Host) | Accelerator |

| CPU DRAM | Accel DRAM |

# DRAM: Transfer from Accelerator to CPU

```
1  val N = ArgIn[Int]
2  setArg(N, args(0).to[Int])
3  val data = DRAM[Int](N)
4
5  val armData: Array[Int] = … //Info soon!
6
7  setMem(data, array)
8
9  Accel {
10    …
11 }
12
13 val outputData = getMem(data)
14
15
16
17
18
```

Copies data from Accel DRAM to CPU DRAM

CPU (Host)

Accelerator

CPU DRAM

Accel DRAM

# ON-CHIP MEMORIES

# Reg

```
1 val in  = ArgIn[Int]
2 val out = ArgOut[Int]
3 …
4 Accel {
5   val reg = Reg[Int]
6
7
8 }
9
10
11
12
13
14
15
16
17
18
```

Creates a **Reg** which holds a value of type **Int**

in

reg

out

# Reg Reset Value

```
1  val in  = ArgIn[Int]
2  val out = ArgOut[Int]
3  …
4  Accel {
5    val reg = Reg[Int](0)
6
7
8  }
9
10
11
12
13
14
15
16
17
18
```

Sets the reset value of this register to be *0*

**Note**: Reset values are currently restricted to constants

in

reg

out

# Reg Writing

```
1  val in  = ArgIn[Int]
2  val out = ArgOut[Int]
3  …
4  Accel {
5    val reg = Reg[Int](0)
6    reg := in
7
8  }
```

Creates a write to input of this register

# Reg Reading

```
1  val in  = ArgIn[Int]
2  val out = ArgOut[Int]
3  …
4  Accel {
5    val reg = Reg[Int](0)
6    reg := in
7    out := reg.value
8  }
9
10
11
12
13
14
15
16
17
18
```

Creates wires connected to the output of this register

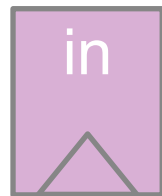**Note**: Register reads are normally implicit, but can be written explicitly

in → reg → out

# SRAM

```
1  val in   = ArgIn[Int]
2  val out  = ArgOut[Int]
3  …
4  Accel {
5    val sram = SRAM[Int](32, 32)
6
7
8  }
9
10
11
12
13
14
15
16
17
18
```

Creates an Accelerator **SRAM** (aka buffer, BRAM) of size 32 x 32 with values of type **Int**

**1 to 5 dimensions** are currently supported

in

sram

wr_addr
wr_data
wr_en          rd_data

rd_addr
rd_en

out

# SRAM

```
1  val in  = ArgIn[Int]
2  val out = ArgOut[Int]
3  …
4  Accel {
5    val sram = SRAM[Int](in.value, in.value)
6    val sram = SRAM[Int](32, 32)
7
8
9  }
10
11
12
13
14
15
16
17
18
```

SRAM dimensions **must** be statically known constants



sram

in

wr_addr
wr_data       rd_data
wr_en

rd_addr
rd_en

out

# SRAM Writes

```
1  val in  = ArgIn[Int]
2  val out = ArgOut[Int]
3  …
4  Accel {
5    val sram = SRAM[Int](32, 32)
6    sram(0, 0) = in.value
7
8  }
```

Creates an SRAM write port with data *in.value,* address *0* which is enabled by *Accel*

# SRAM Reads



```
1 val in  = ArgIn[Int]
2 val out = ArgOut[Int]
3 …
4 Accel {
5   val sram = SRAM[Int](32, 32)
6   sram(0, 0) = in.value
7   out := sram(0,0)
8 }
9
10
11
12
13
14
15
16
17
18
```
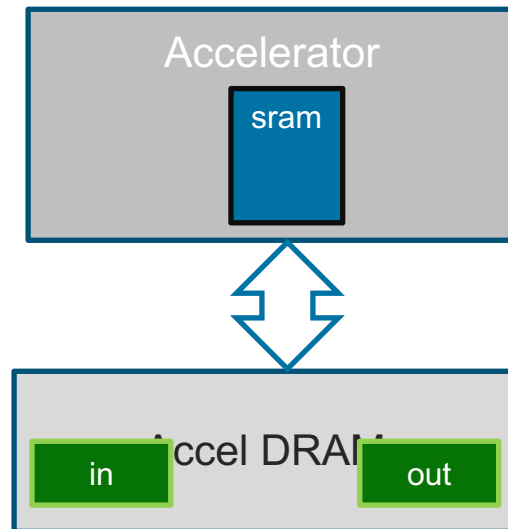
Creates an SRAM read port with address *0* which is enabled by *Accel*

# SRAM: Interfacing with DRAM?

```
1 val in  = DRAM[Int](32)
2 val out = DRAM[Int](32)
3 …
4 Accel {
5   val sram = SRAM[Int](16)
6
7
8 }
9
10
11
12
13
14
15
16
17
18
```

How can we copy data to/from Accel DRAM?

Accelerator

sram

Accel DRAM

in

out

# SRAM: Dense Loading from DRAM

```
1  val in  = DRAM[Int](32)
2  val out = DRAM[Int](32)
3  …
4  Accel {
5    val sram = SRAM[Int](16)
6    sram load in(0::16)
7
8  }
```

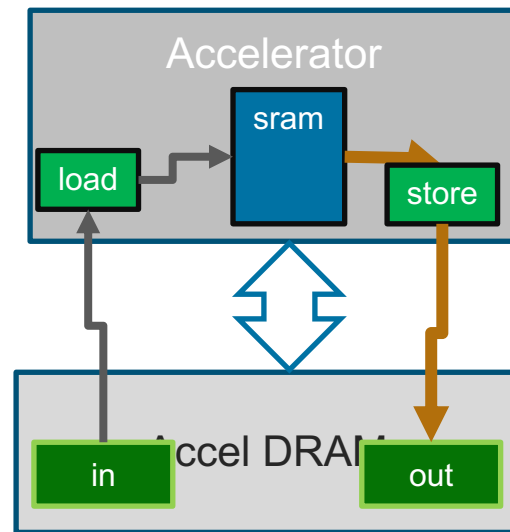Creates logic which loads data within *in,* address range **0 until 16** (exclusive), to *sram*

**Note**: The address range can be omitted if SRAM and DRAM are the same size

# SRAM: Dense Storing to DRAM

```
1  val in  = DRAM[Int](32)
2  val out = DRAM[Int](32)
3  …
4  Accel {
5    val sram = SRAM[Int](16)
6    sram load in(0::16)
7    out(0::16) store sram
8  }
9
10
11
12
13
14
15
16
17
18
```

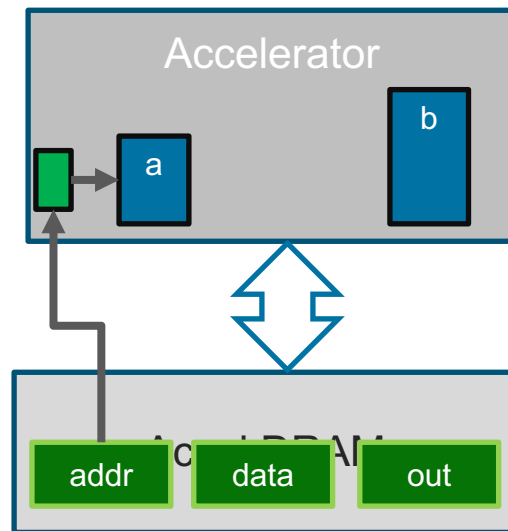Creates logic which stores contents of *sram* to *out's* address range **0 until 16** (exclusive)

# SRAM: Gather from DRAM

```
 1 val data = DRAM[Int](32)
 2 val addr = DRAM[Int](32)
 3 val out  = DRAM[Int](32)
 4 …
 5 Accel {
 6   val a = SRAM[Int](16)
 7   val b = SRAM[Int](16)
 8   a load addr(0::16) // Addresses
 9
10 }
11
12
13
14
15
16
17
18
```
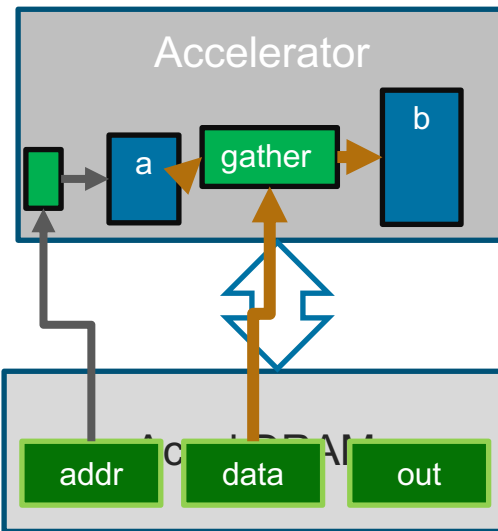
**Equivalent C:**
```
for (i=0; i<16; i++) {
  b[i] = data[a[i]]
}
```



Accelerator

b

a

Accel DRAM

addr    data    out

# SRAM: Gather from DRAM

```
 1 val data = DRAM[Int](32)
 2 val addr = DRAM[Int](32)
 3 val out  = DRAM[Int](32)
 4 …
 5 Accel {
 6   val a = SRAM[Int](16)
 7   val b = SRAM[Int](16)
 8   a load addr(0::16) // Addresses
 9   b gather data(a)
10 }
11
12
13
14
15
16
17
18
```

Creates logic which gathers elements in *data* at addresses in *a* into *b*



Accelerator

b

gather

a

Accel DRAM

addr    data    out

```
for (i=0; i<16; i++) {
  b[i] = data[a[i]]
}
```

# SRAM: Gather from DRAM

```
 1  val data = DRAM[Int](32)
 2  val addr = DRAM[Int](32)
 3  val out  = DRAM[Int](32)
 4  …
 5  Accel {
 6    val a = SRAM[Int](16)
 7    val b = SRAM[Int](16)
 8    a load addr(0::16) // Addresses
 9    b gather data(a, 10)
10  }
```

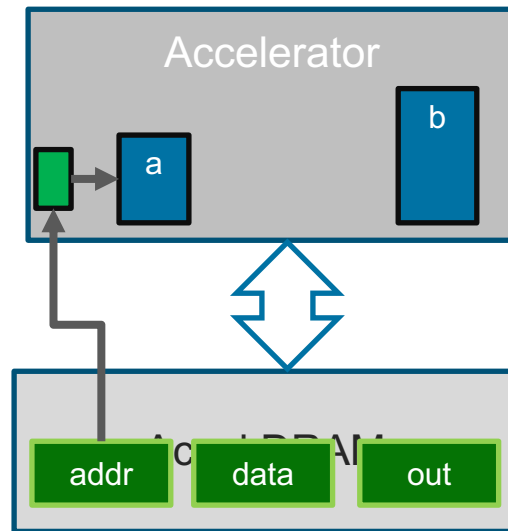Uses the first 10 elements in a only



Accelerator

b
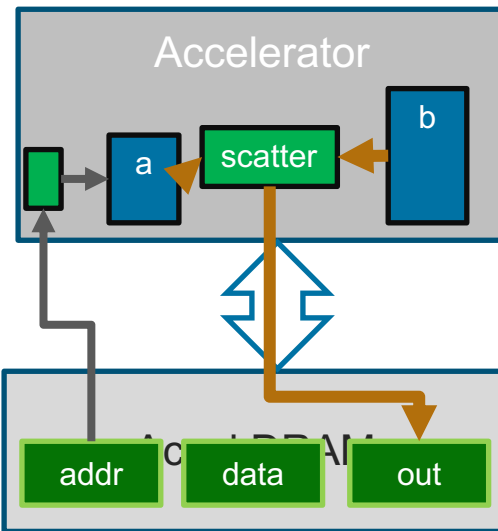
a   gather

addr   data   out

```
for (i=0; i<10; i++) {
  b[i] = data[a[i]]
}
```

# SRAM: Scatter to DRAM

```
 1 val data = DRAM[Int](32)
 2 val addr = DRAM[Int](32)
 3 val out  = DRAM[Int](32)
 4 …
 5 Accel {
 6    val a = SRAM[Int](16)
 7    val b = SRAM[Int](16)
 8    a load addr(0::16) // Addresses
 9
10 }
11
12
13
14
15
16
17
18
```

**Equivalent C:**
```
for (i=0; i<16; i++) {
   out[a[i]] = b[i]
}
```



Accelerator

b

a

Accel DRAM

addr    data    out

# SRAM: Scatter to DRAM

```
 1 val data = DRAM[Int](32)
 2 val addr = DRAM[Int](32)
 3 val out  = DRAM[Int](32)
 4 …
 5 Accel {
 6   val a = SRAM[Int](16)
 7   val b = SRAM[Int](16)
 8   a load addr(0::16) // Addresses
 9   out(a) scatter b
10 }
11
12
13
14
15
16
17
18
```
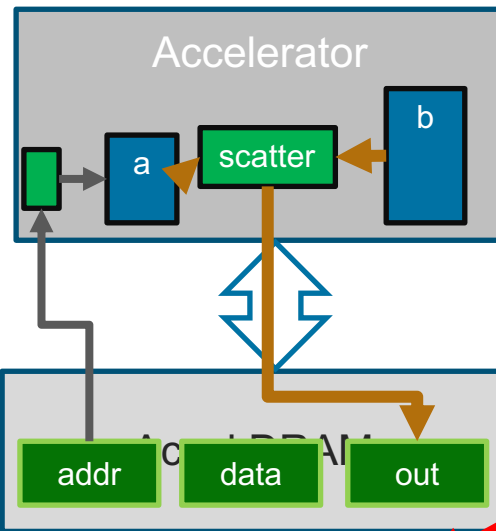
Creates logic which scatters elements in *b* into *out* at addresses in *a*



Accelerator

b

scatter

a

Accel DRAM

addr    data    out

```
for (i=0; i<16; i++) {
  out[a[i]] = b[i]
}
```

# SRAM: Scatter to DRAM

```
 1 val data = DRAM[Int](32)
 2 val addr = DRAM[Int](32)
 3 val out  = DRAM[Int](32)
 4 …
 5 Accel {
 6   val a = SRAM[Int](16)
 7   val b = SRAM[Int](16)
 8   a load addr(0::16) // Addresses
 9   out(a, 10) scatter b
10 }
11
12
13
14
15
16
17
18
```
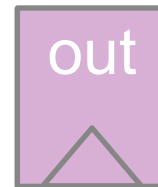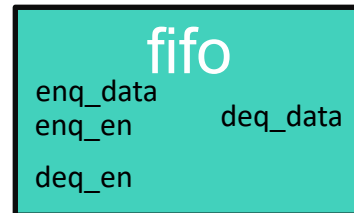


```
for (i=0; i<10; i++) {
  out[a[i]] = b[i]
}
```

# FIFO

```
1  val in  = ArgIn[Int]
2  val out = ArgOut[Int]
3  …
4  Accel {
5    val fifo = FIFO[Int](16)
6
7
8  }
9
10
11
12
13
14
15
16
17
18
```

FIFO with depth 16

**Note**: Depth must be statically known

in

fifo
enq_data
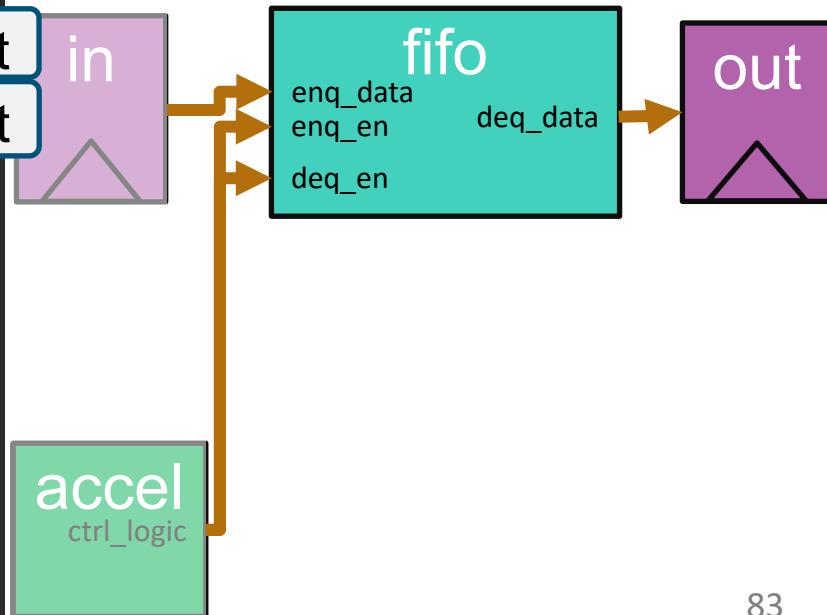enq_en          deq_data
deq_en

out

accel
ctrl_logic

# FIFO: Enqueueing / Dequeueing

```
1  val in  = ArgIn[Int]
2  val out = ArgOut[Int]
3  …
4  Accel {
5    val fifo = FIFO[Int](16)
6    a.enq(scalarIn)
7    scalarOut := a.deq()
8  }
9
10
11
12
13
14
15
16
17
18
```

FIFO enqueue port

FIFO dequeue port
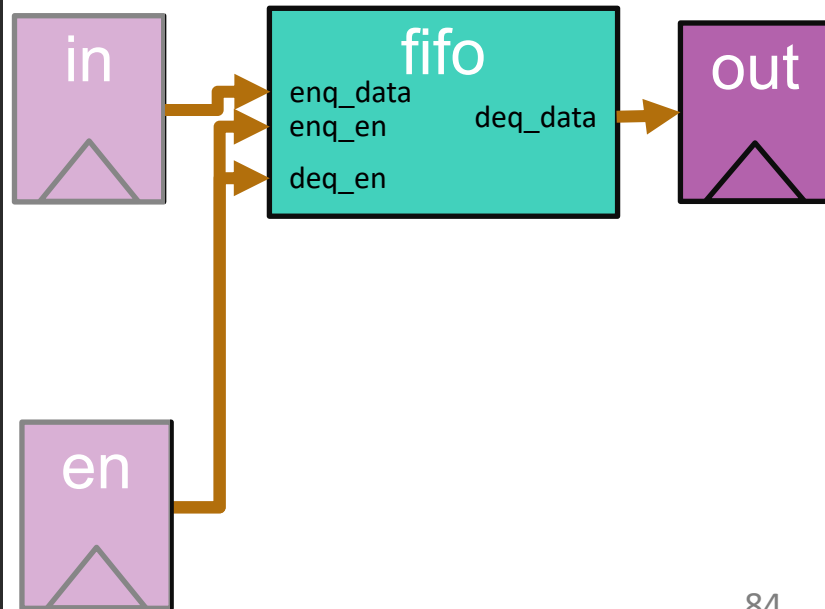
in

fifo

enq_data
enq_en

deq_en

deq_data

out

accel

ctrl_logic

# FIFO: Enabled Enqueuing/Dequeueing

```
1 val in  = ArgIn[Int]
2 val en  = ArgIn[Bool]
3 val out = ArgOut[Int]
4 …
5 Accel {
6   val fifo = FIFO[Int](16)
7   a.enq(scalarIn, en)
8   scalarOut := a.deq(en)
9 }
10
11
12
13
14
15
16
17
18
```

Can also set data-dependent enables for enqueue/dequeue e.g. for data filtering

in

en

fifo

enq_data
enq_en

deq_data

deq_en

out

# FIFO: Transfers to/from DRAM

```
1  val data = DRAM[Int](32)
2  val out  = DRAM[Int](32)
3  …
4  Accel {
5    val fifo = FIFO[Int](32)
6    fifo load data
7    out store fifo
8  }
9
10
11
12
13
14
15
16
17
18
```
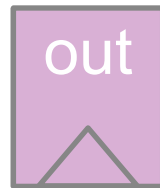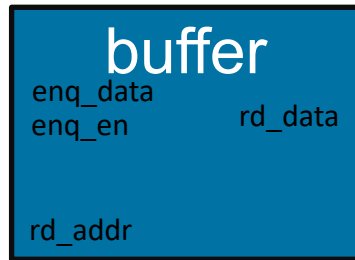
Load from DRAM

Store to DRAM

# LineBuffer

```
1 val scalarIn  = ArgIn[Int]
2 val scalarOut = ArgOut[Int]
3 …
4 Accel {
5   val buffer = LineBuffer[Int](3, 1024)
6
7
8 }
9
10
11
12
13
14
15
16
17
18
```

**LineBuffer** with 3 rows, each with 1024 columns

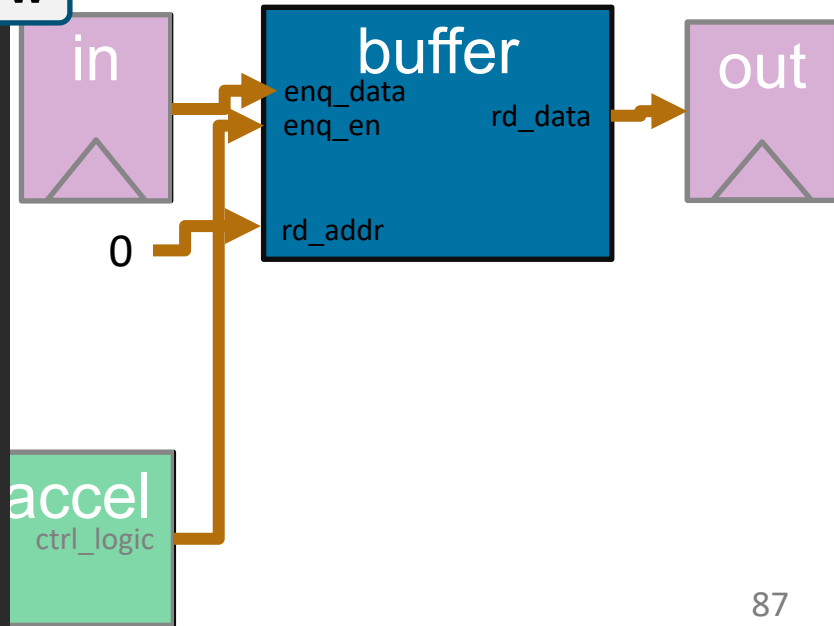**Note**: Only 2-dimensional buffers currently supported. Dimensions must be statically known

in

buffer
enq_data
enq_en          rd_data

rd_addr

out

accel
ctrl_logic

```
1  val in  = ArgIn[Int]
2  val out = ArgOut[Int]
3  …
4  Accel {
5    val buffer = Lin
6    buffer.enq(in)
7    out := buffer(0, 0)
8  }
9
10
11
12
13
14
15
16
17
18
```

Enqueue to **current row**

**Addressed** linebuffer read

**Note**: **LineBuffer** contains internal logic to increment **current row** when *#columns* elements have been stored
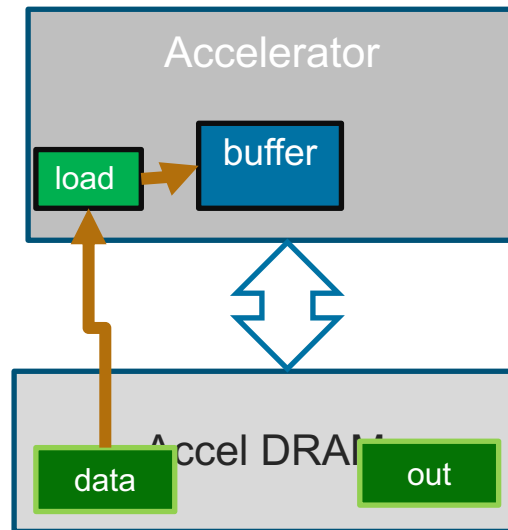
in

buffer

enq_data
enq_en

rd_data

rd_addr

0

out

accel

ctrl_logic

87

# LineBuffer: Loading from DRAM

```
1  val data = DRAM[Int](100,1024)
2  val out  = DRAM[Int](32)
3  …
4  Accel {
5     val buffer = LineBuffer[Int](3,1024)
6     buffer load data(0, 0::1024)
7
8  }
```

Load *data* row 0, columns 0 until 1024 to **current row** of buffer

**Note**: Storing to **DRAM** from **LineBuffer** is currently unsupported

# A Note About Ports

- Spatial compiler makes best effort to minimize amount of resources needed to implement logical memories

- However, writing and reading from the same memory many times can be expensive!

- Be aware of how many times you read/write a given memory, and try to minimize the number of concurrent reads

# CONTROLLERS
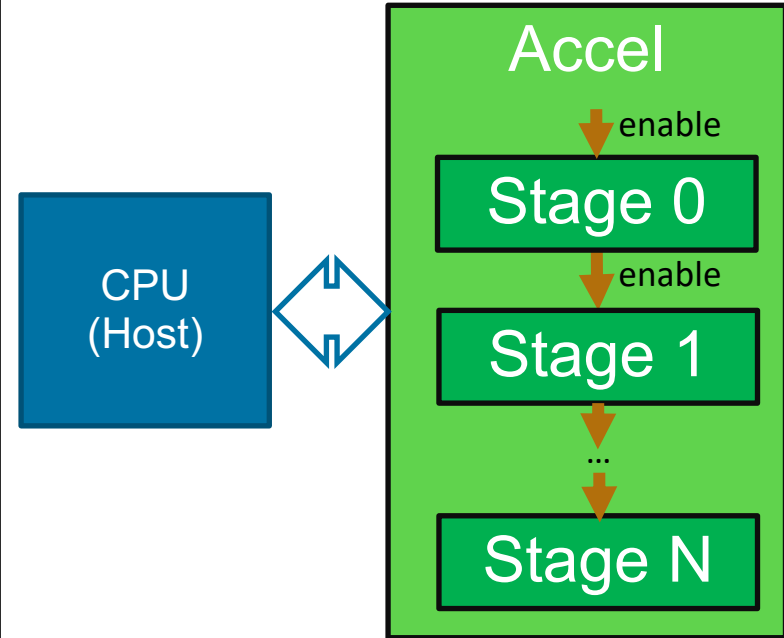
# Accel

```
1  Accel {
2    …
3
4  }
```

Receives **start** from CPU
Executes stages sequentially
Sends **done** to CPU

A **stage** is either:
- one primitive operation
- one controller

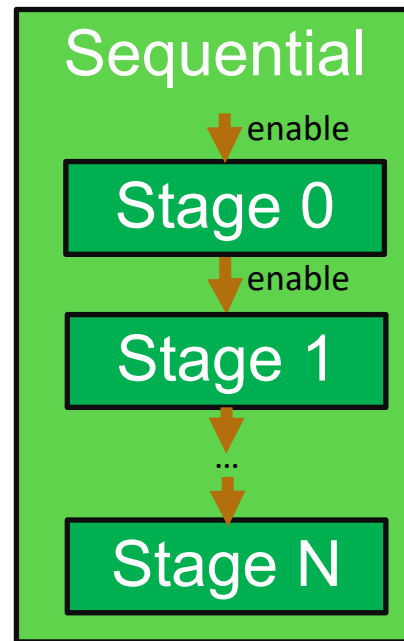**Note**: May not be nested in any other controller

CPU (Host)

Accel

enable

Stage 0

enable

Stage 1

…

Stage N

# Sequential

```
1  Accel {
2    …
3    Sequential {
4      …
5
6    }
7    …
8  }
```

Executes stages sequentially

# Parallel

```
1   Accel {
2     …
3     Parallel {
4       …
5
6     }
7     …
8   }
9
10
11
12
13
14
15
16
17
18
```
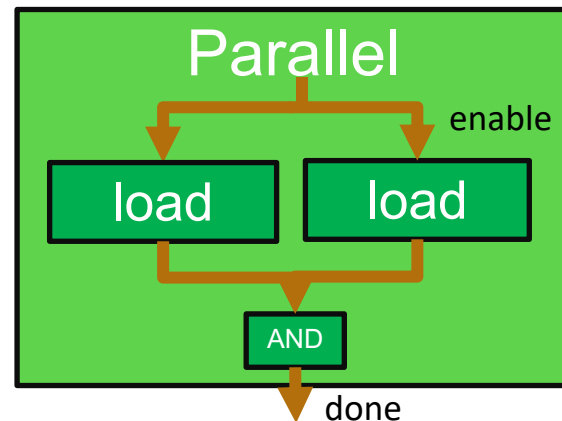
Executes stages in parallel
Completes when all stages finish

**Note**: Spatial will soon infer control structures for parallel execution automatically
But for now, use **Parallel** when you want to guarantee parallel execution



93

# Parallel: Example

```
1  val dataA = DRAM[Int](1024)
2  val dataB = DRAM[Int](1024)
3
4  Accel {
5    val a = SRAM[Int](16)
6    val b = SRAM[Int](16)
7    Parallel {
8      a load dataA(0::16)
9      b load dataB(0::16)
10   }
11   …
12 }
```
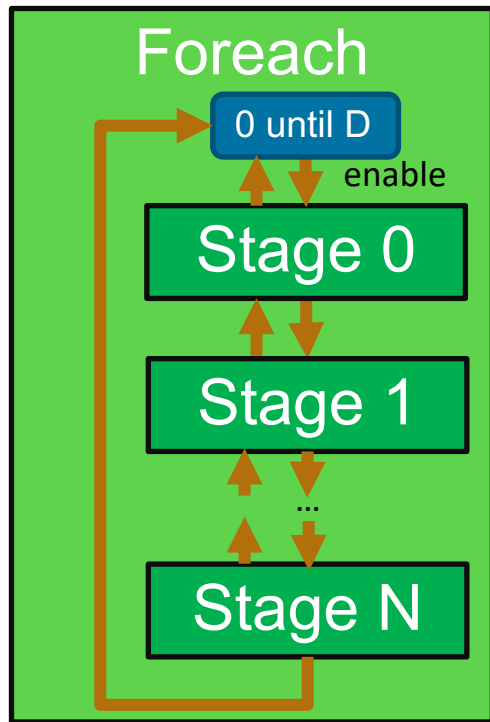
# Foreach

```
1  val D = ArgIn[Int]
2  Accel {
3    Foreach(0 until D) {i =>
4      …
5
6    }
7    …
8  }
9
10
11
12
13
14
15
16
17
18
```

Loop iterator

Executes each stage in a pipelined fashion, repeating for *D* iterations

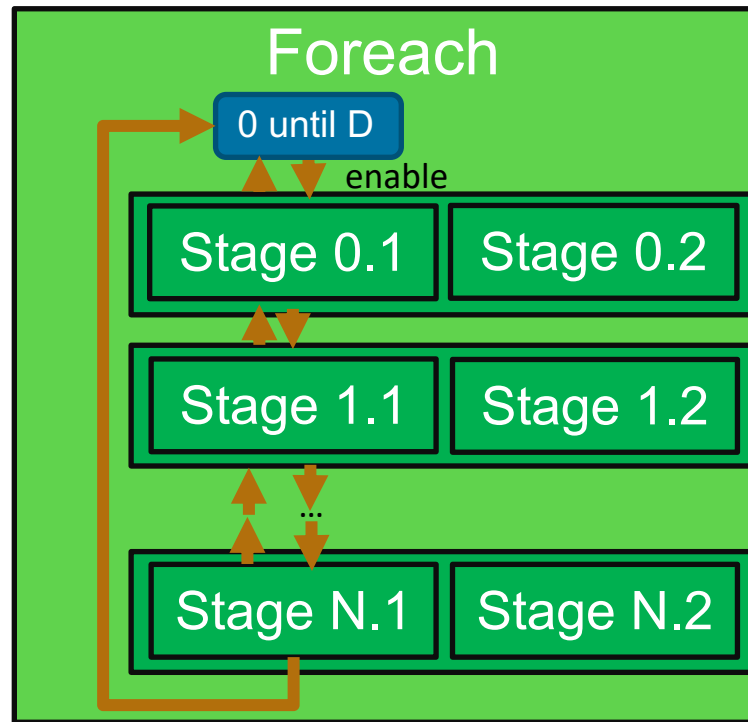Spatial handles memory buffering and stall signals!

Foreach

0 until D

enable

Stage 0

Stage 1

…

Stage N

# Foreach: Parallelization

```
1  val D = ArgIn[Int]
2  Accel {
3    Foreach(0 until D par 2) {i =>
4      …
5
6    }
7    …
8  }
9
10
11
12
13
14
15
16
17
18
```

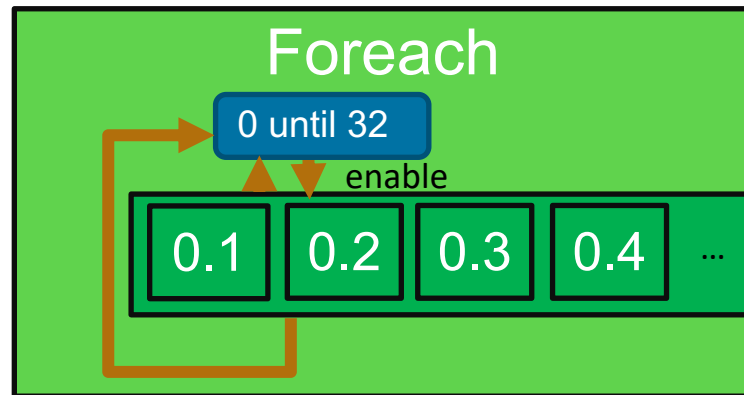Parallelizes the pipeline by duplicating the body

Spatial also handles memory banking for you!

# Foreach: Example

```
 1 val data = DRAM[Int](32)
 2 Accel {
 3   val input  = SRAM[Int](32)
 4   val output = SRAM[Int](32)
 5   input load data
 6   Foreach(0 until 32 par 16) {i =>
 7     output(i) = input(i) * 2
 8   }
 9   data store output
10 }
```
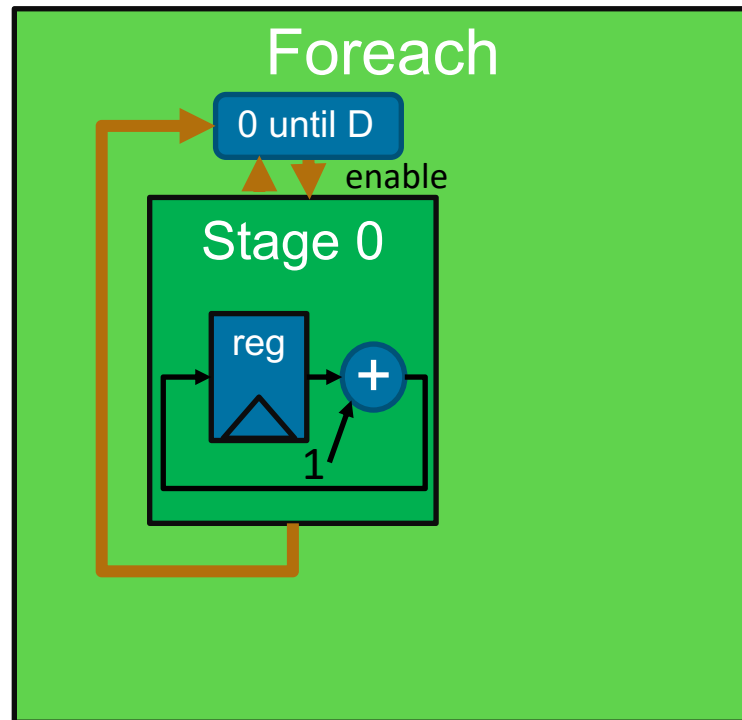
Multiply every element by 2, store back to DRAM



Foreach

0 until 32

enable

0.1  0.2  0.3  0.4  …

# Foreach: Illegal Parallelization Cases

```
1  val D = ArgIn[Int]
2  Accel {
3    val reg = Reg[Int](0)
4    Foreach(0 until D par 2) {i =>
5      reg := reg + 1
6    }
7    …
8  }
9
```

It's **unsafe** to parallelize pipelines with loop-carry dependencies!
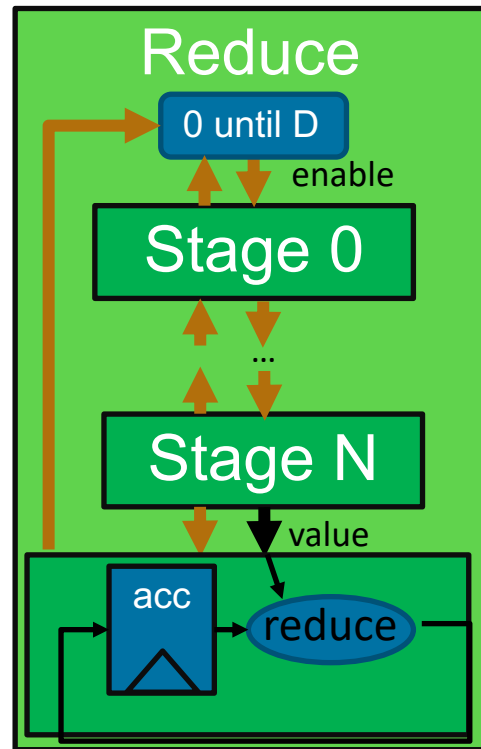
# Reduce



Zero value OR accumulator

Loop iterator

```
…
Reduce(acc)(0 until D){i =>


}{(a,b) => reduce(a,b) }
…
}
```

Value function (aka map)

reduce(a,b)

Executes each stage in a pipelined fashion, repeating for *N* iterations.
Reduces the result of the value function into an accumulator

Reduce

0 until D

enable

Stage 0

...

Stage N

value

acc

reduce

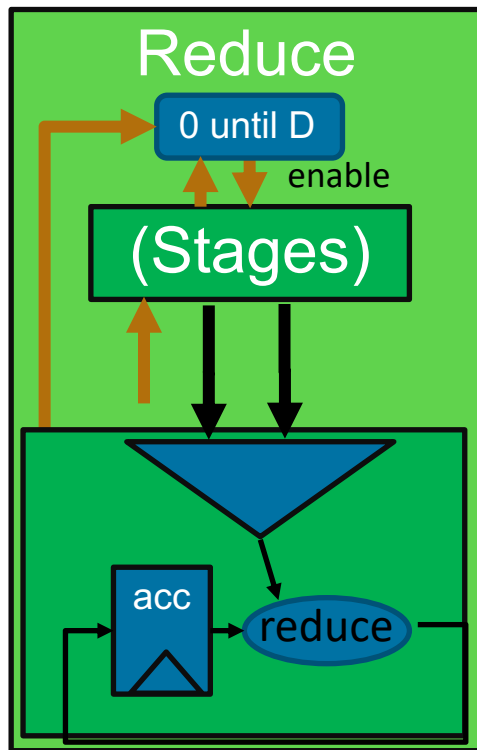# Reduce: Parallelization

```
1  val D = ArgIn[Int]
2  Accel {
3    …
4    Reduce(acc)(0 until D par 2){i =>
5      valueFunction(i)
6    }{(a,b) => reduce(a,b) }
7    …
8  }
9
10
11
12
13
14
15
16
17
18
```

Parallelize!

Value function is parallelized like Foreach
Reduction is parallelized using a tree



Reduce

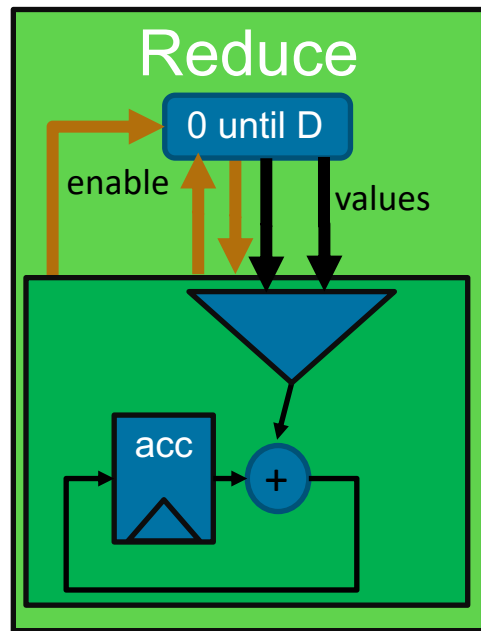0 until D

enable

(Stages)

acc
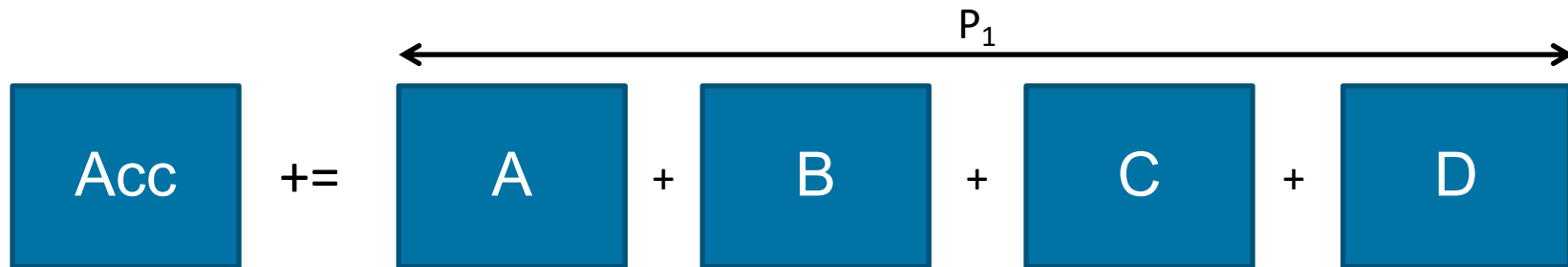
reduce

# Reduce: Example

```
1  val D = ArgIn[Int]
2  val out = ArgOut[Int]
3  Accel {
4    val acc = Reg[Int](0)
5    Reduce(acc)(0 until D par 16){i =>
6      i
7    }{(a,b) => a + b }
8    out := acc
9  }
10
11
12
13
14
15
16
17
18
```

Sum the values 0 until D,
adding 16 values in parallel

*accum* contains the sum
after the controller ends



Reduce

0 until D

enable          values

acc

+

# Accumulation

$$\text{Acc} \mathrel{+}= A + B + C + D$$
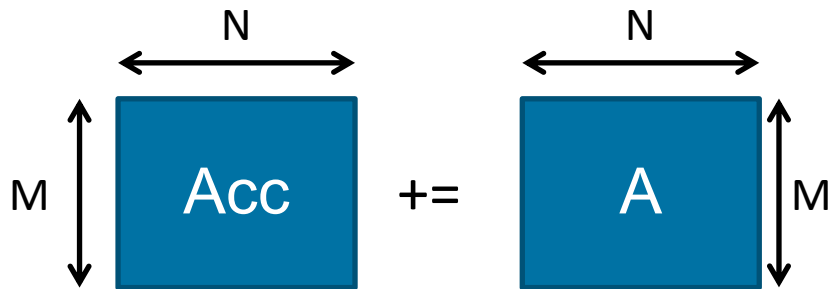
with bracket $P_1$ spanning $A + B + C + D$
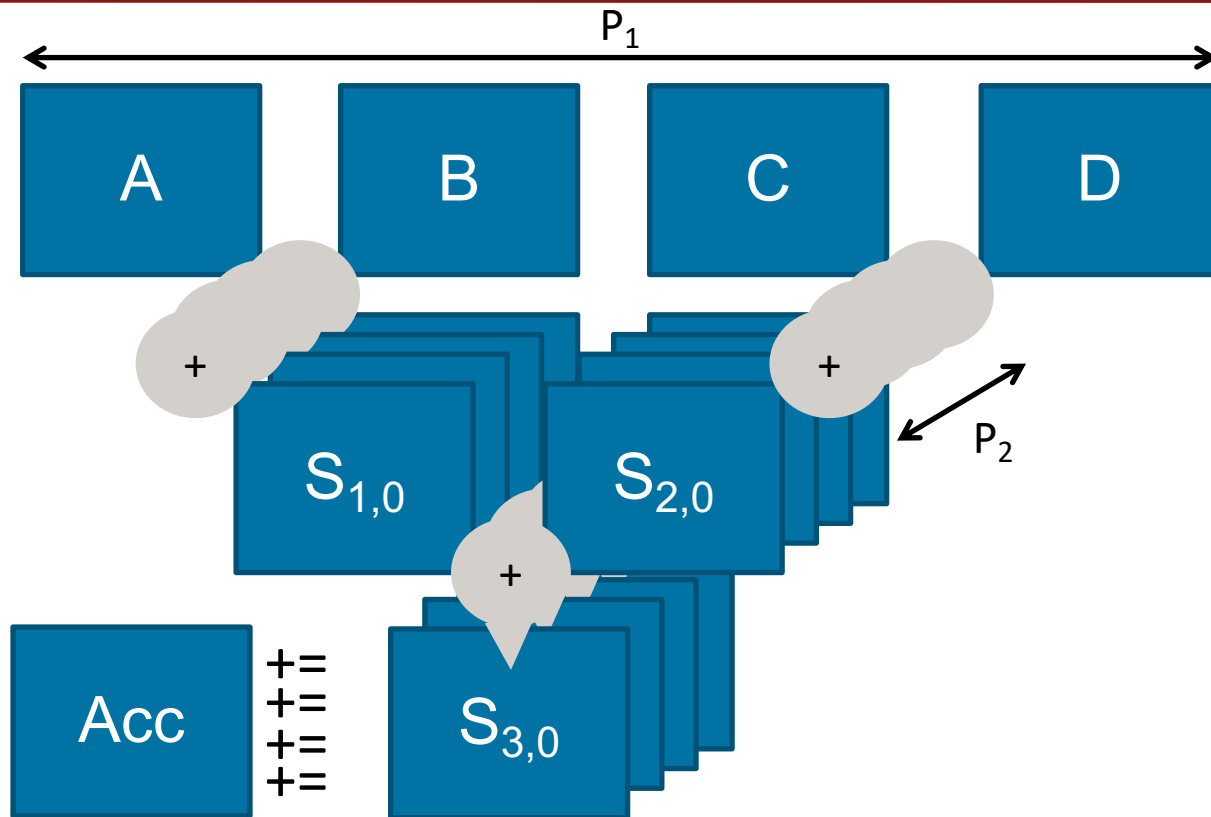
# Reduce (Visualization)
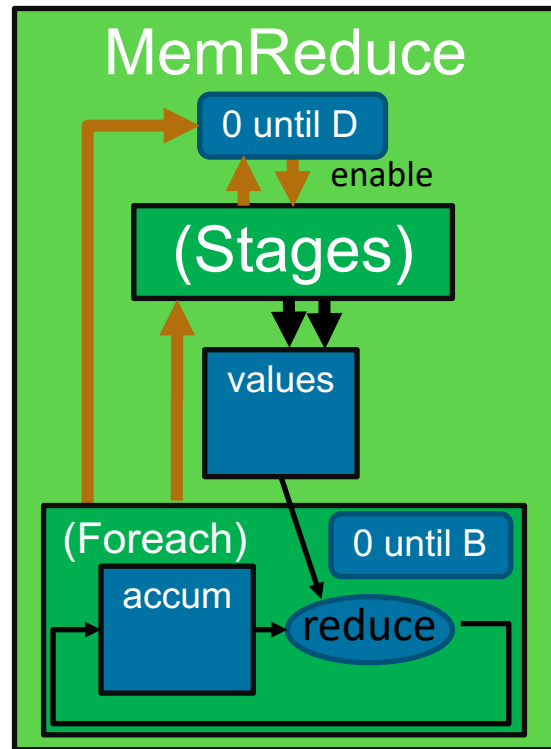
# MemReduce

# MemReduce

# MemReduce

```
1  val D = ArgIn[Int]
2  Accel {
3    val accum = SRAM[Int](B)
4    MemReduce(accum)(0 until D){i =>
5      val values = SRAM[Int](B)
6
7
8      values
9    }{(a,b) =>          }
10
11
12
13
14
15
16
17
18
```

Loop iterator

Value function (aka map)

reduce(a,b)

Executes each stage in a pipelined fashion, repeating for *N* iterations.
**Value function** populates an **SRAM**
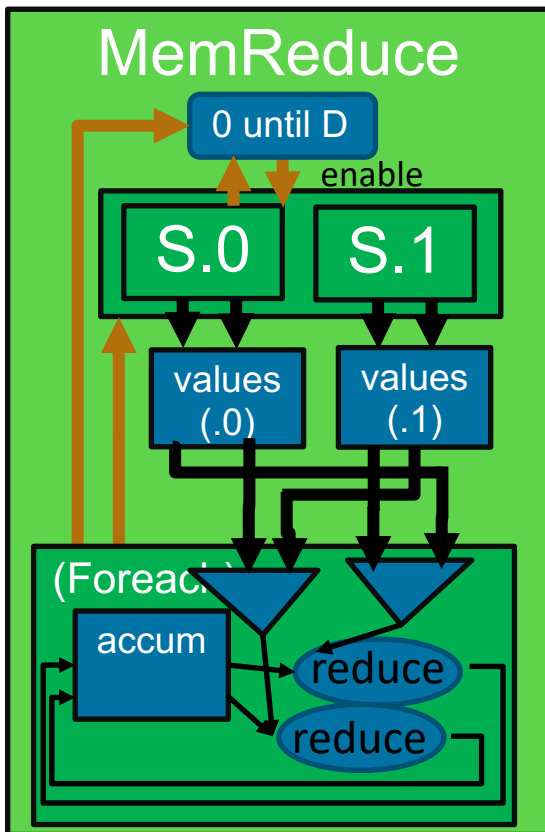**Reduce** says how to combine an element from *value* into *accum*



MemReduce

0 until D

enable

(Stages)

values

(Foreach)

accum

0 until B

reduce

106

# MemReduce: Parallelization



```
1  val D = Arg
2  Accel {
3    val accum = SRAM[Int](B)
4    MemReduce(accum par 2)(0 until D par 2){i =>
5      val values = SRAM[Int](B)
6      valueFunction(values, i)
7      values
8    }{(a,b) => reduce(a,b) }
9    …
10 }
```

Parallelize!

Parallelize!

Can parallelize production of values
AND reduction of values

MemReduce

0 until D

enable

S.0    S.1

values (.0)    values (.1)
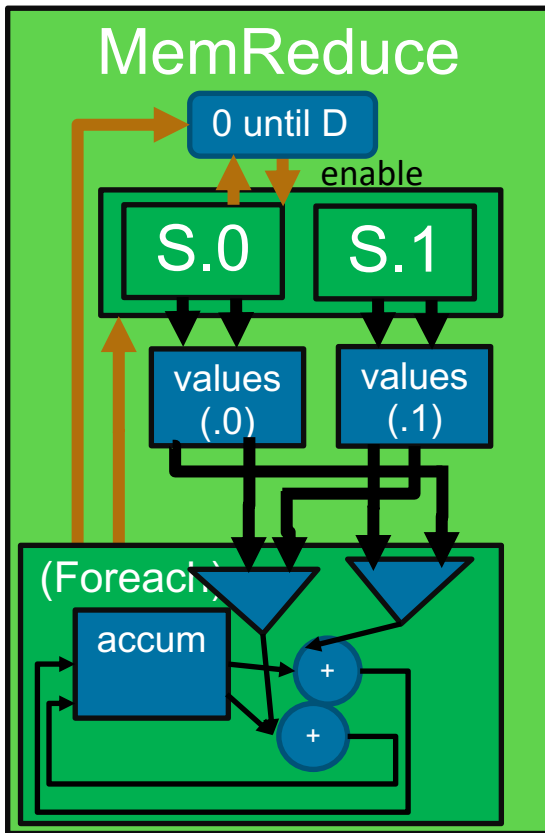
(Foreach)

accum    reduce

reduce

# MemReduce: Example

```
1  val data = DRAM[Int](D)
2  val out  = DRAM[Int](16)
3  Accel {
4    val accum = SRAM[Int](16)
5    MemReduce(accum par 2)(D by 16 par 2){i =>
6     val values = SRAM[Int](16)
7     values load data(i::i+16)
8     values
9    }{(a,b) => a + b }
10
11   out store accum
12 }
```
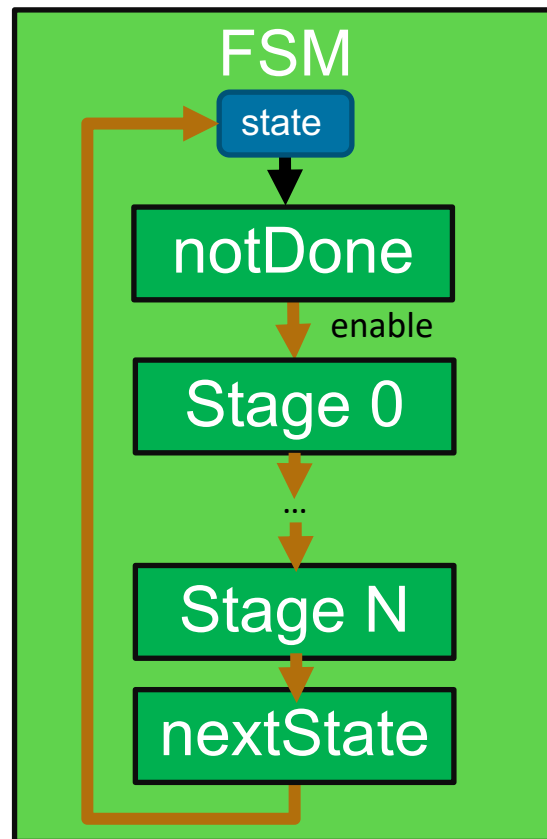
0 until D, strided by 16

Chunks up *data* into 16 element blocks and combines the blocks using element-wise addition

# FSM (State Machine)

```
1  Accel {
2    …
3    FSM(init){state => notDone(state) }{state =>
4
5      action(state)
6
7    }{state => nextState(state) }
8    …
9  }
```

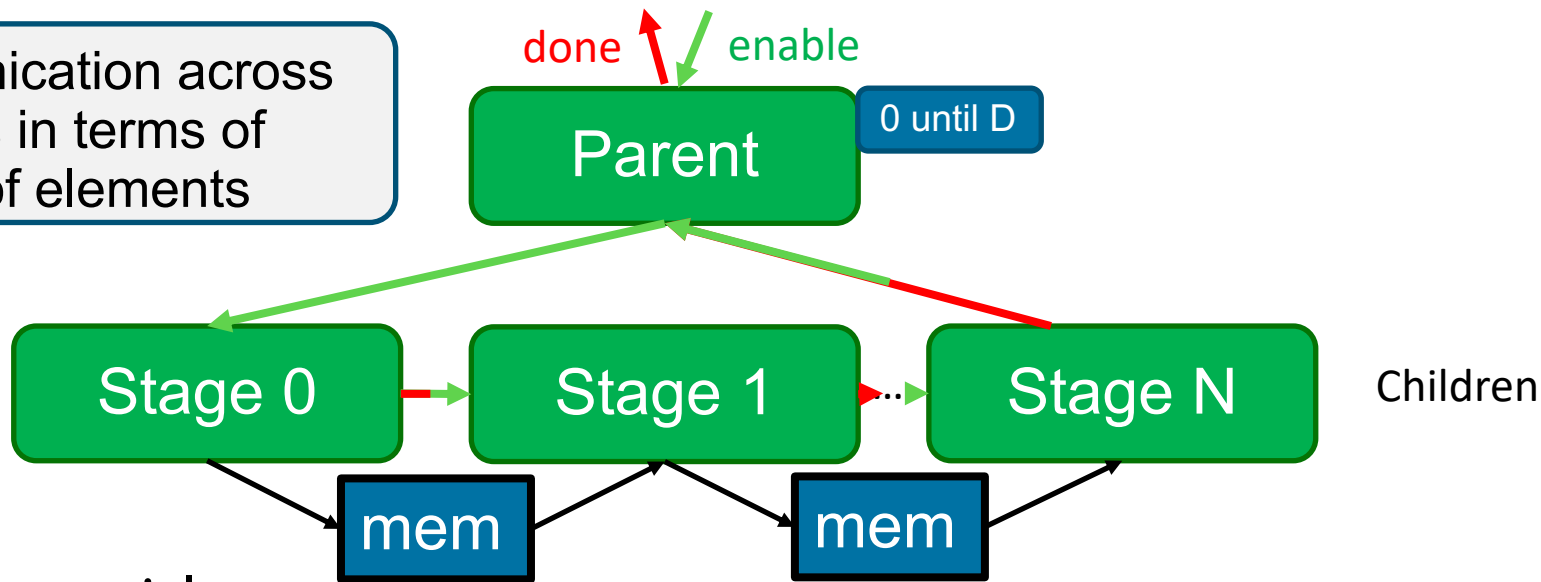Can also give an explicit initial state

(Otherwise initial state is zero)

# Controller Tags

```
1   Accel {
2     …
3     Sequential.Foreach(0 until D){i =>
4       …
5     }
6
7     Sequential.Reduce(0)(0 until D){i =>
8       …
9     }
10
```

**Sequential** can be added as a tag on looping controllers to change execution of stages from pipelined to purely sequential

# Sequential Execution

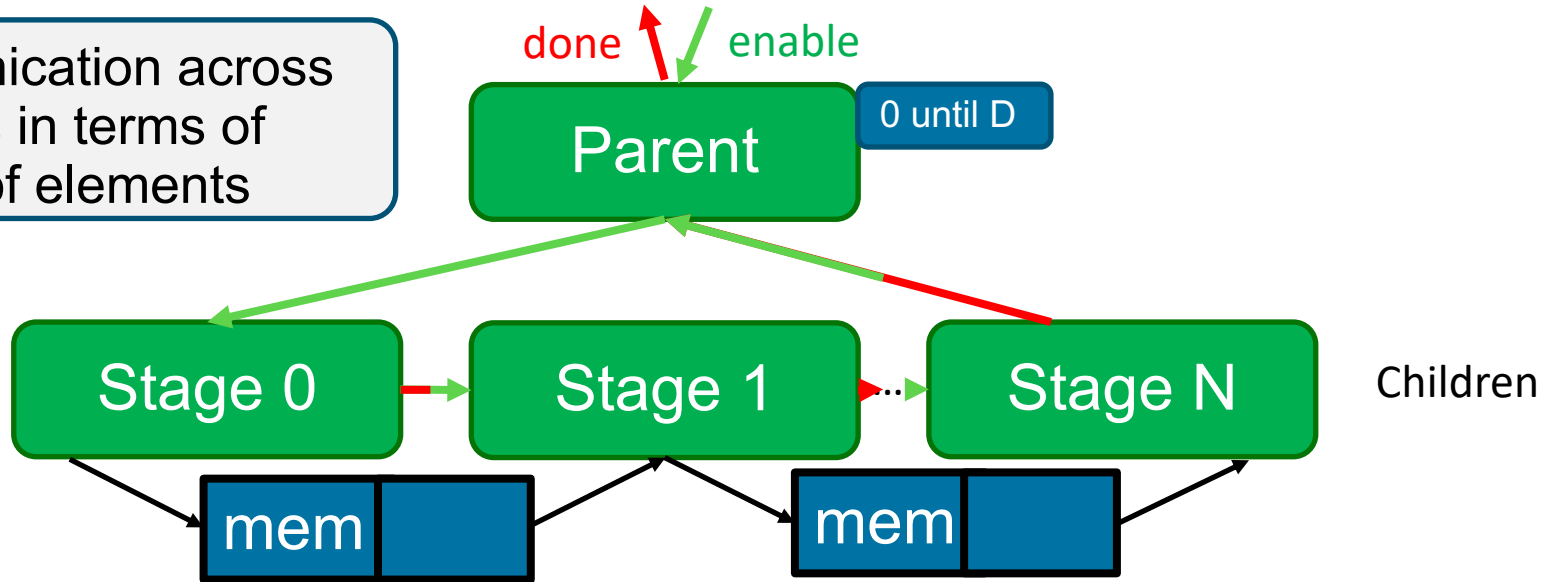Communication across stages is in terms of **blocks** of elements

done    enable

Parent    0 until D

Stage 0    Stage 1    ...    Stage N    Children

mem    mem

- Sequential
  - Parent enables Stage 0 when enabled, as long as counter < D
  - Stage K is enabled when Stage K−1 completes (K > 0)

# Pipelined Execution

Communication across stages is in terms of **blocks** of elements

done    enable

Parent    0 until D

Stage 0    Stage 1    ...    Stage N    Children

mem    mem

- Pipelined
  - Parent enables Stage 0 TWICE when enabled, as long as counter < D
  - Stage K is enabled when Stage K−1 completes (K > 0)