# Garnet: The Next Generation CGRA Architecture

Priyanka Raina

May 10, 2019

# Challenges with Jade CGRA
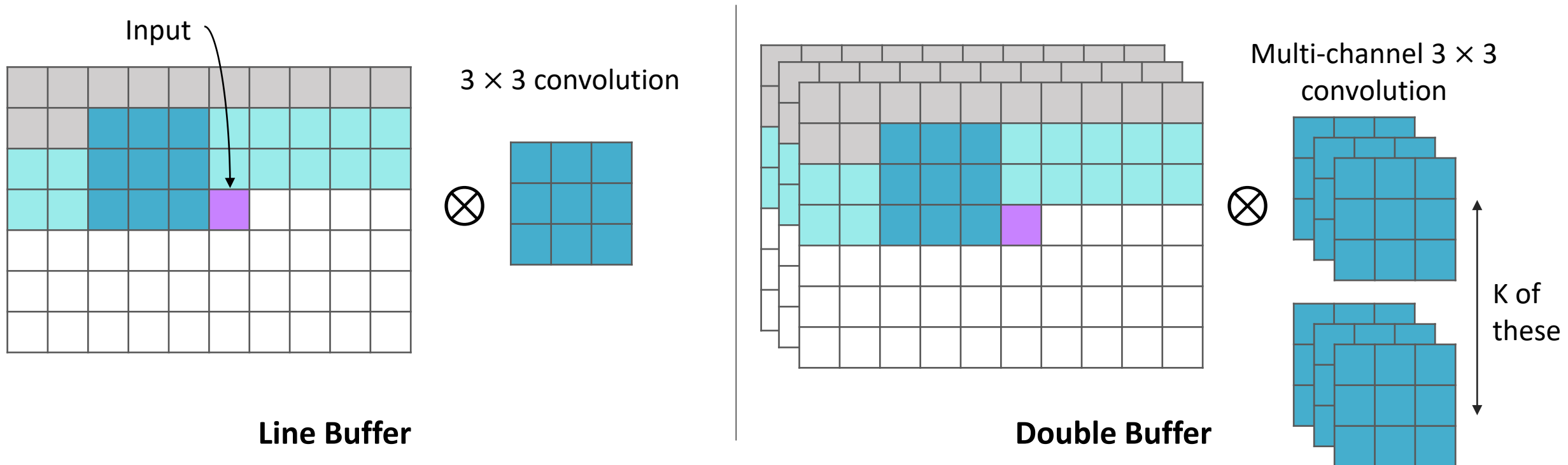
1. PE supports only simple integer operations
    - Porting applications that use floating point requires application expertise and manual effort

# Challenges with Jade CGRA

2. Memory supports only line buffered pipelines
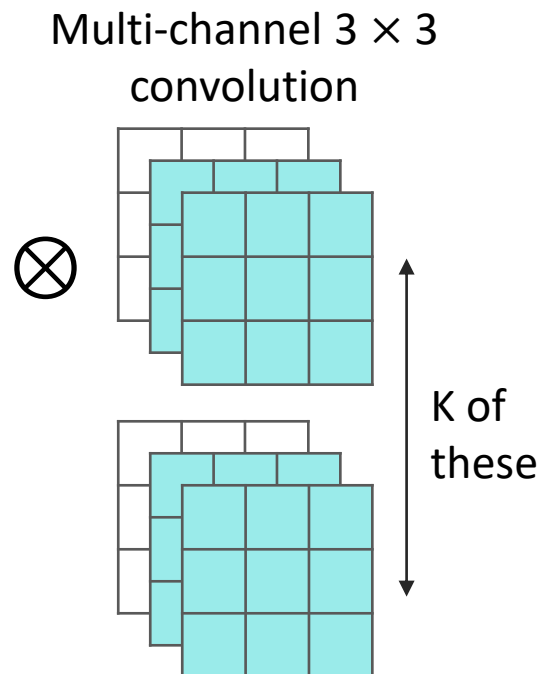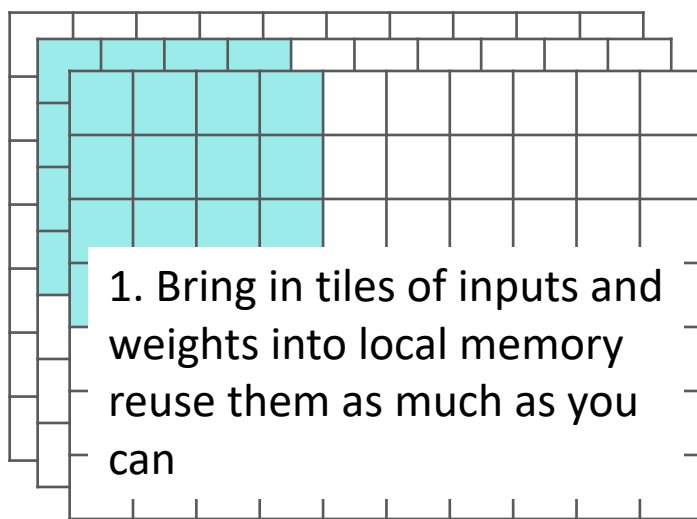   - Most new applications even for imaging and vision use neural networks
   - Need a memory hierarchy with double buffers for energy-efficiency



Input

$3 \times 3$ convolution

Line Buffer

Multi-channel $3 \times 3$ convolution

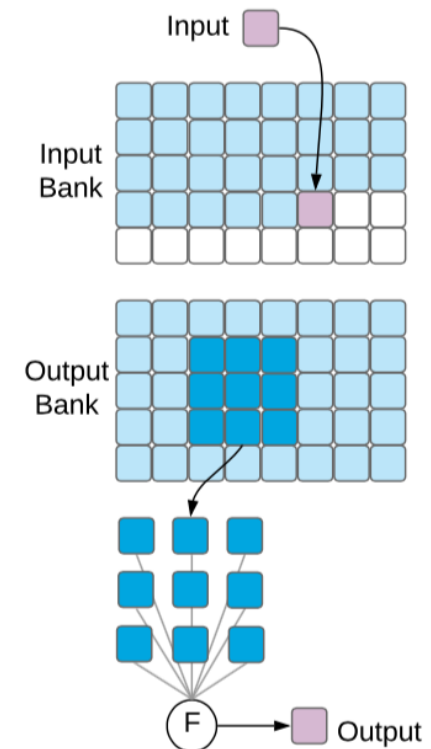K of these

Double Buffer

# Challenges with Jade CGRA

2. Memory supports only line buffered pipelines
   - Most new applications even for imaging and vision use neural networks
   - Need a memory hierarchy with double buffers for energy-efficiency

Multi-channel $3 \times 3$ convolution

1. Bring in tiles of inputs and weights into local memory reuse them as much as you can

$\otimes$

K of these

2. Use a double buffer to overlap compute on current tile and fetching of next tile

3. Use a hierarchy of double buffers to maximize energy-efficiency

Input

Input Bank

Output Bank

F → Output

# Challenges with Jade CGRA
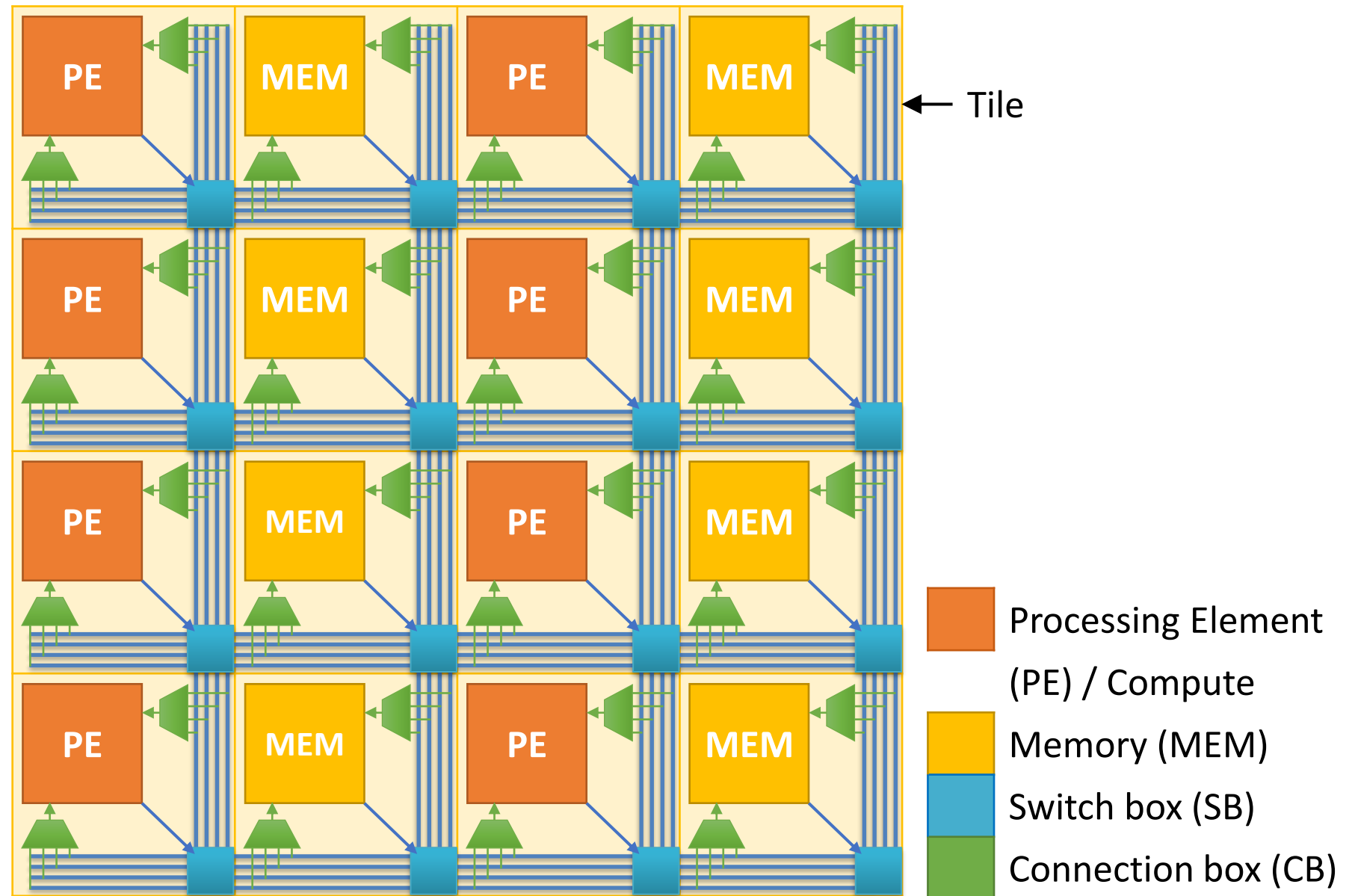
3. Configuration over JTAG is slow
   - Large applications with multiple kernels don't fit on the CGRA – need fast reconfiguration

# Key (application driven!) architectural changes in Garnet CGRA

1. Support for Bfloat16, and for executing complex operations like divide using multiple PEs

2. Addition of a global buffer to create a memory hierarchy for efficiently executing neural networks, and double buffer support in all memories

3. Fast reconfiguration support using global buffer and control processor

4. Addition of configurable power domains

# CGRA

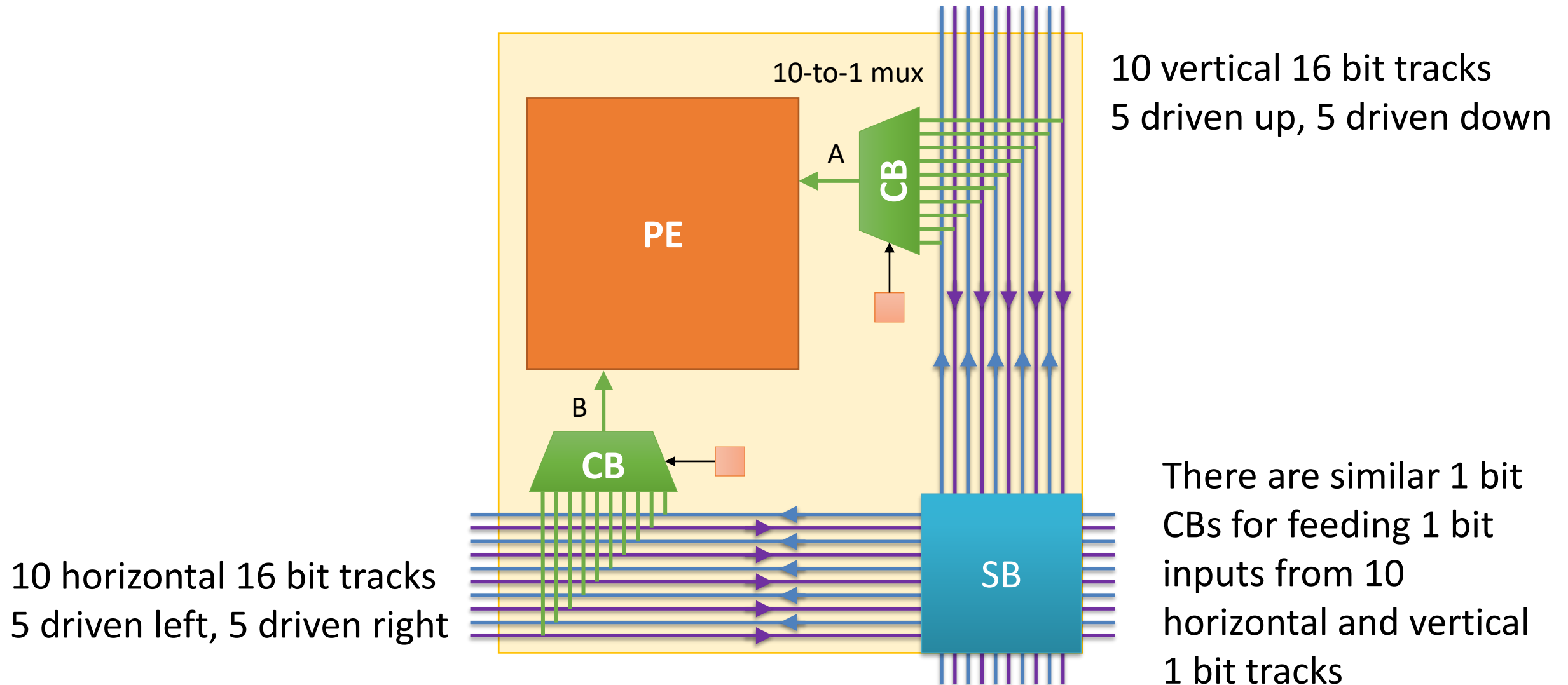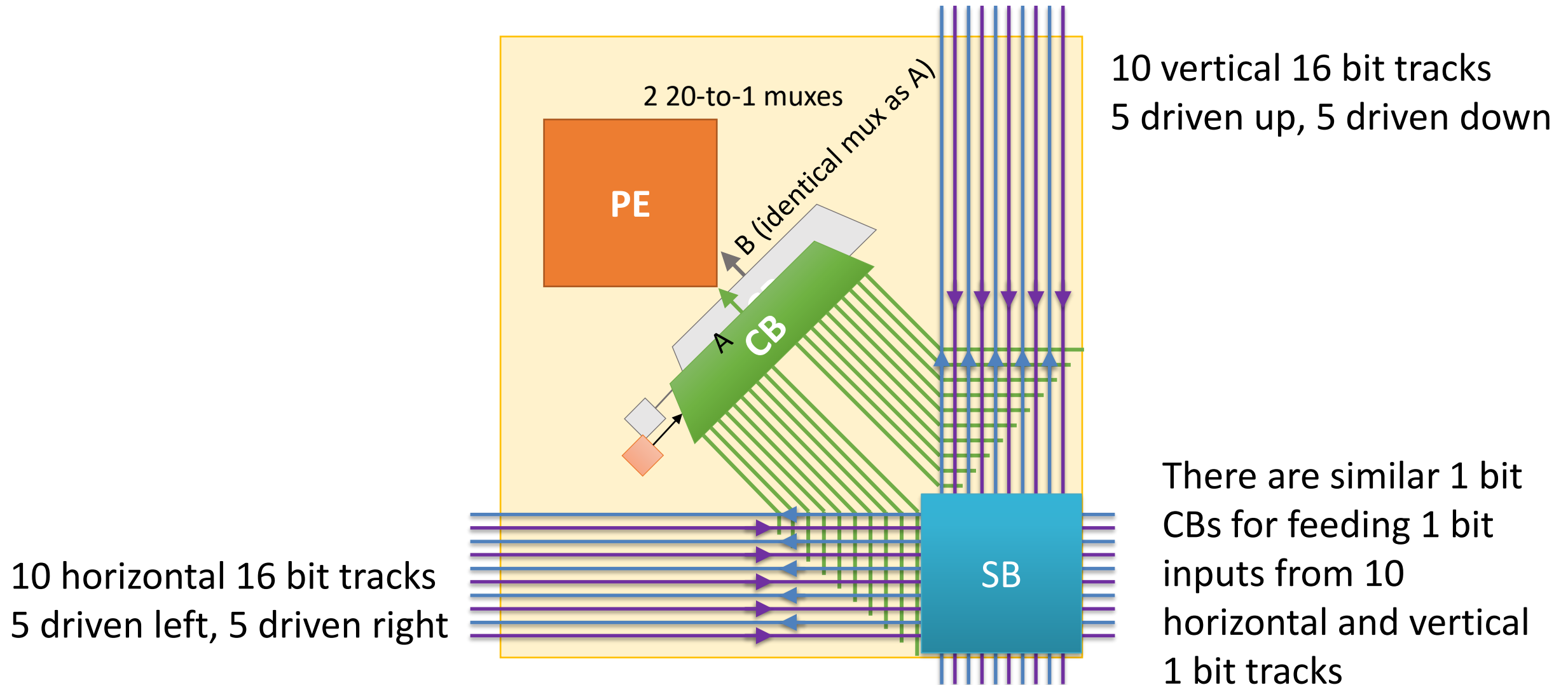Coarse Grained Reconfigurable Array (CGRA)



← Tile

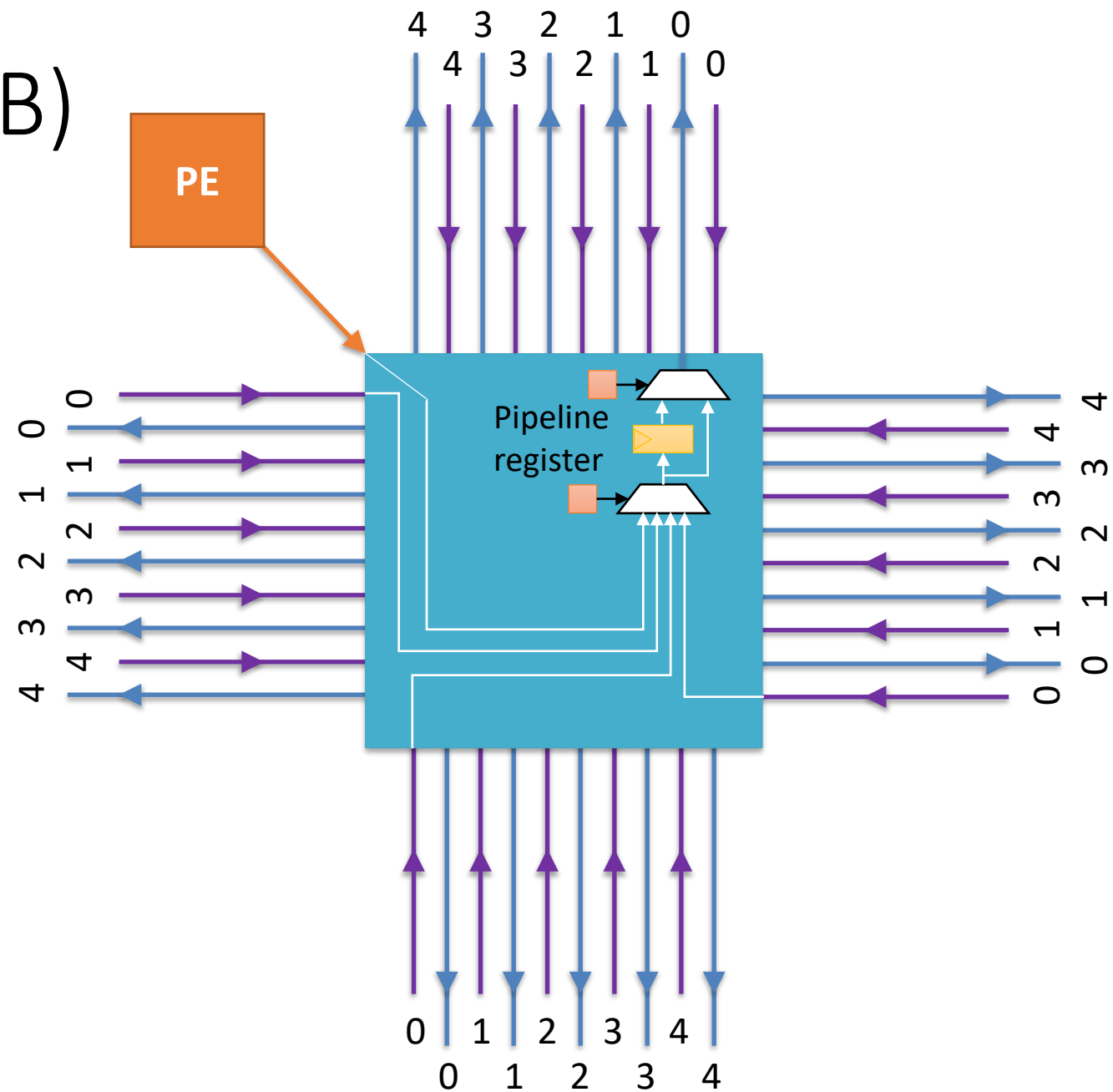Processing Element (PE) / Compute

Memory (MEM)

Switch box (SB)

Connection box (CB)

# Interconnect

Keyi Zhang

# Connection Box (CB) - Jade



10-to-1 mux

A

CB

PE

B

CB

SB

10 vertical 16 bit tracks
5 driven up, 5 driven down

There are similar 1 bit
CBs for feeding 1 bit
inputs from 10
horizontal and vertical
1 bit tracks

10 horizontal 16 bit tracks
5 driven left, 5 driven right

# Connection Box (CB) - Garnet



2 20-to-1 muxes

B (identical mux as A)

A

PE

CB

SB

10 vertical 16 bit tracks
5 driven up, 5 driven down

There are similar 1 bit CBs for feeding 1 bit inputs from 10 horizontal and vertical 1 bit tracks

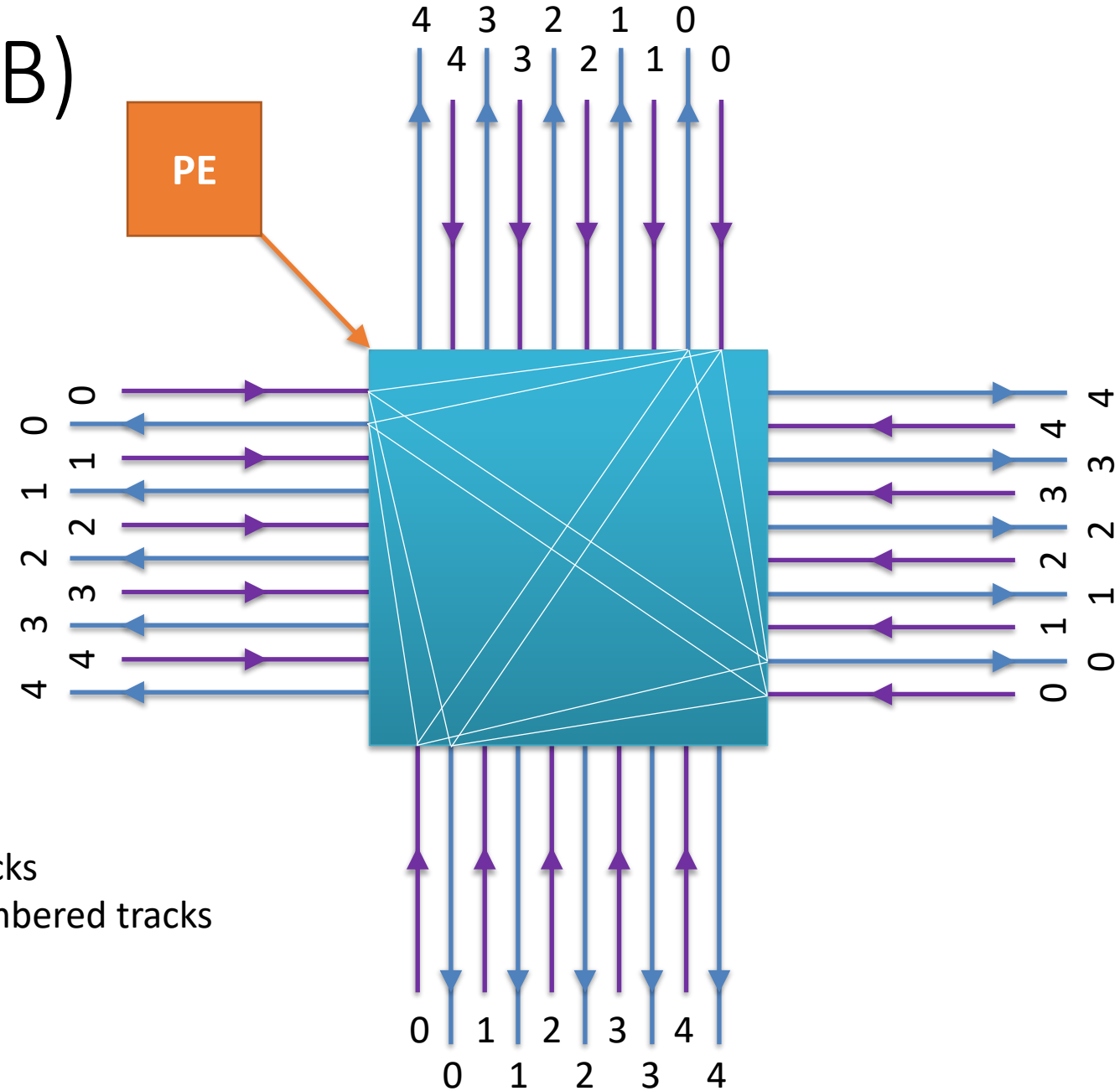10 horizontal 16 bit tracks
5 driven left, 5 driven right
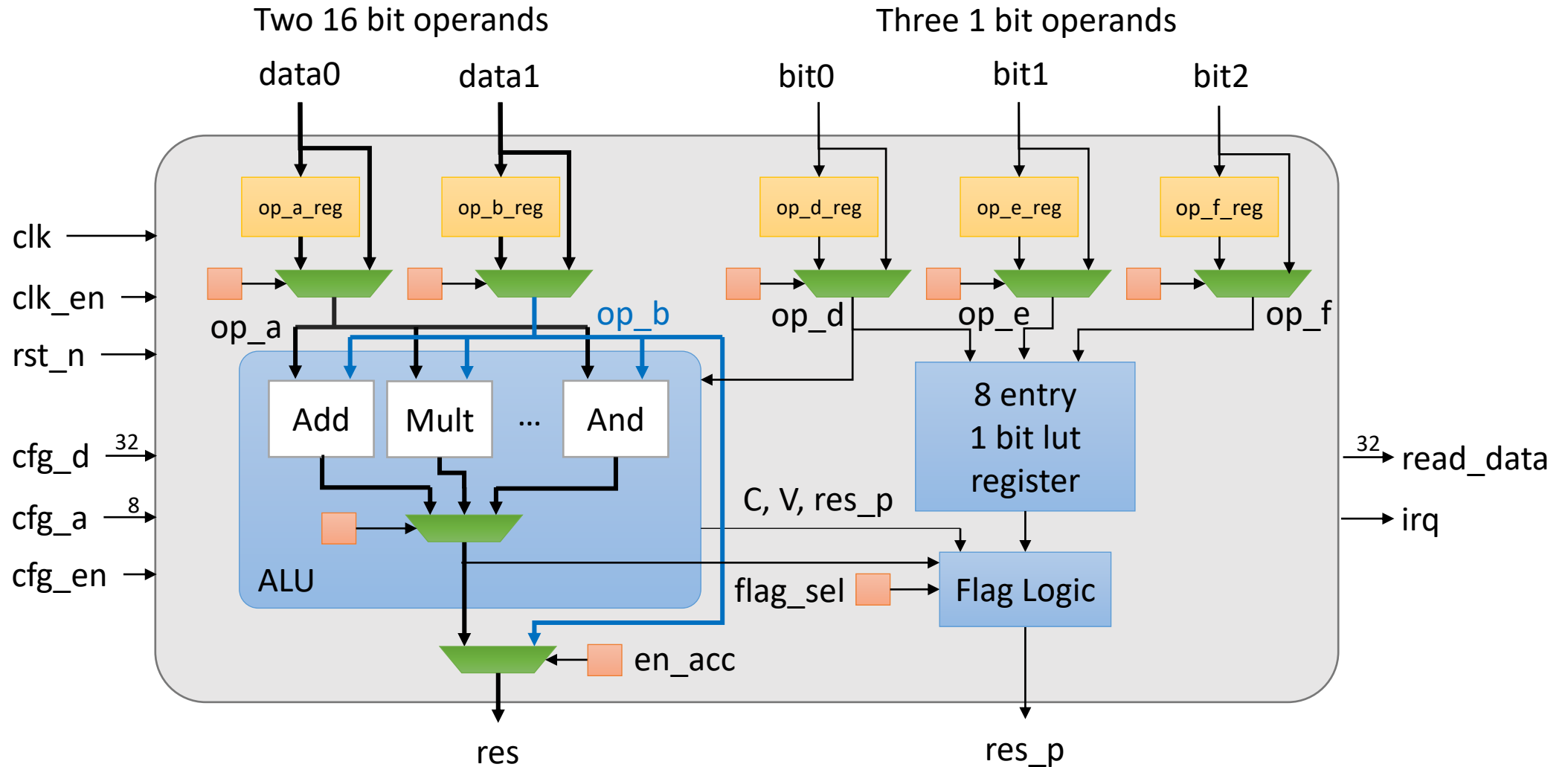
# Switch Box (SB)

# Switch Box (SB)

Figure shows connectivity for '0' tracks
Similar connectivity for all other numbered tracks
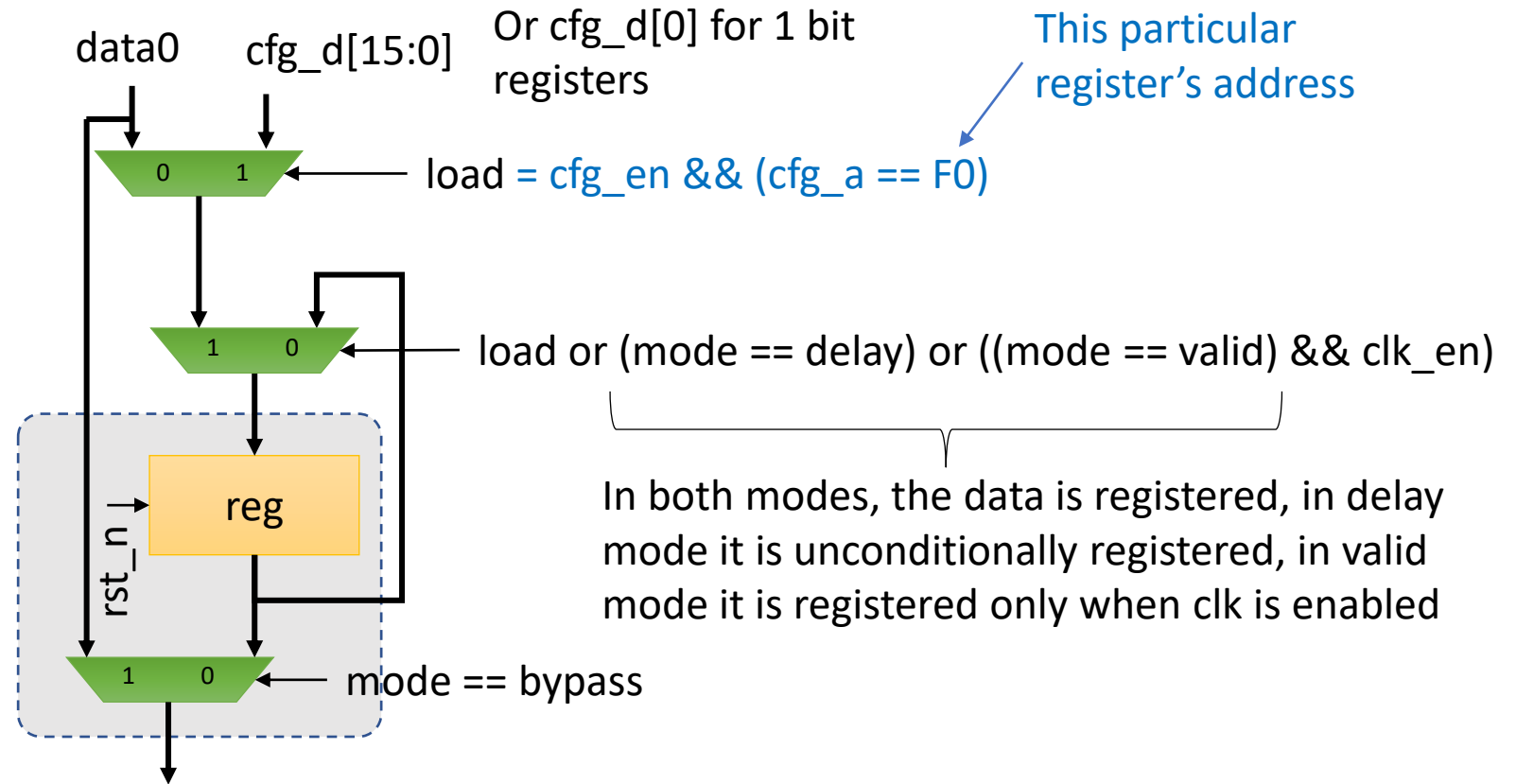PE output goes to all muxes

# Processing Element
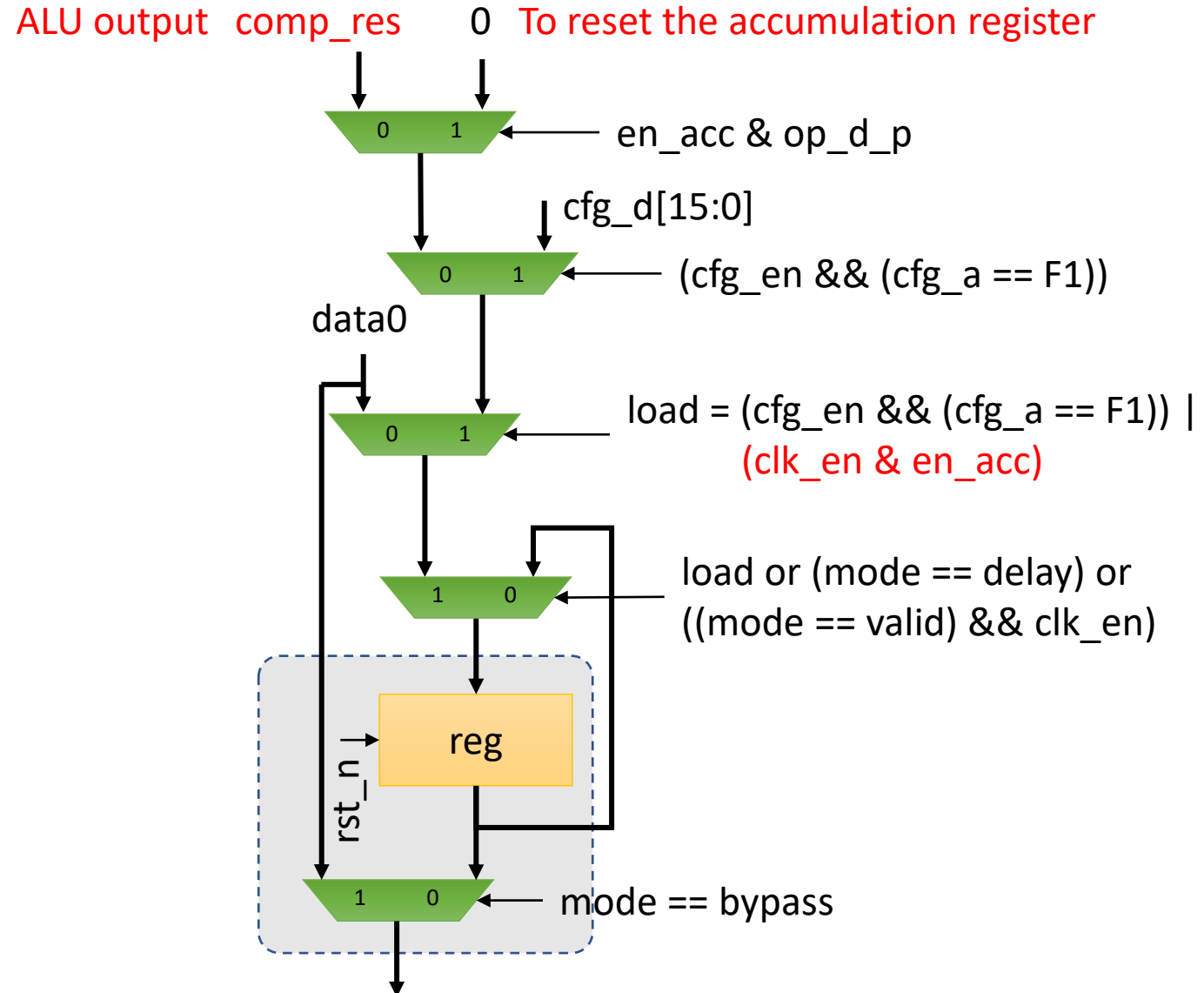
Nikhil Bhagdikar

# Processing Element (PE) in Diablo

# Register Modes

data0    cfg_d[15:0]    Or cfg_d[0] for 1 bit
registers

This particular
register's address

0    1    ←    load = cfg_en && (cfg_a == F0)

1    0    ←    load or (mode == delay) or ((mode == valid) && clk_en)

In both modes, the data is registered, in delay
mode it is unconditionally registered, in valid
mode it is registered only when clk is enabled

reg

rst_n

1    0    ←    mode == bypass

Previous slide shows only this
dotted portion
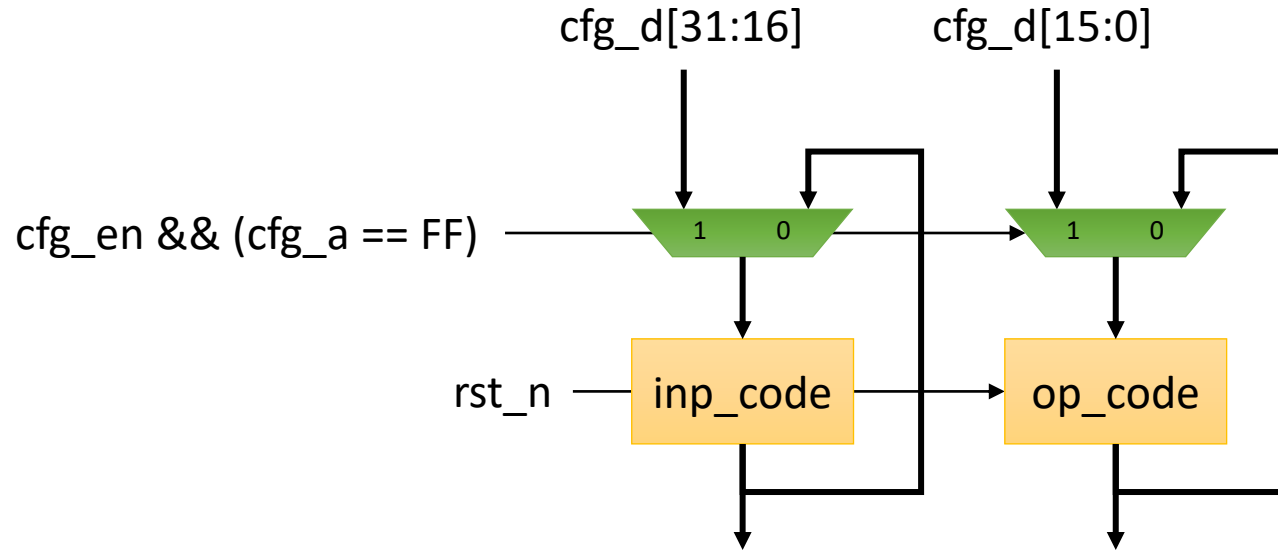In reality, the register is more
complex because of debug
functionality

# Accumulation Operation

- For b, load signal can also be clk_en & en_acc (bit in the op code that enables accumulation)

- The data going into the register can also be 0/output of the PE after accumulation

ALU output   comp_res   0   To reset the accumulation register

en_acc & op_d_p

cfg_d[15:0]

(cfg_en && (cfg_a == F1))

data0

load = (cfg_en && (cfg_a == F1)) | (clk_en & en_acc)

load or (mode == delay) or ((mode == valid) && clk_en)

reg

rst_n

mode == bypass

# Configuration Registers in Diablo

# Flag Logic

**ALU produces**

C = carry is generated
V = overflow
res 16b
res_p 1b

Z = res == 0
N = res is negative

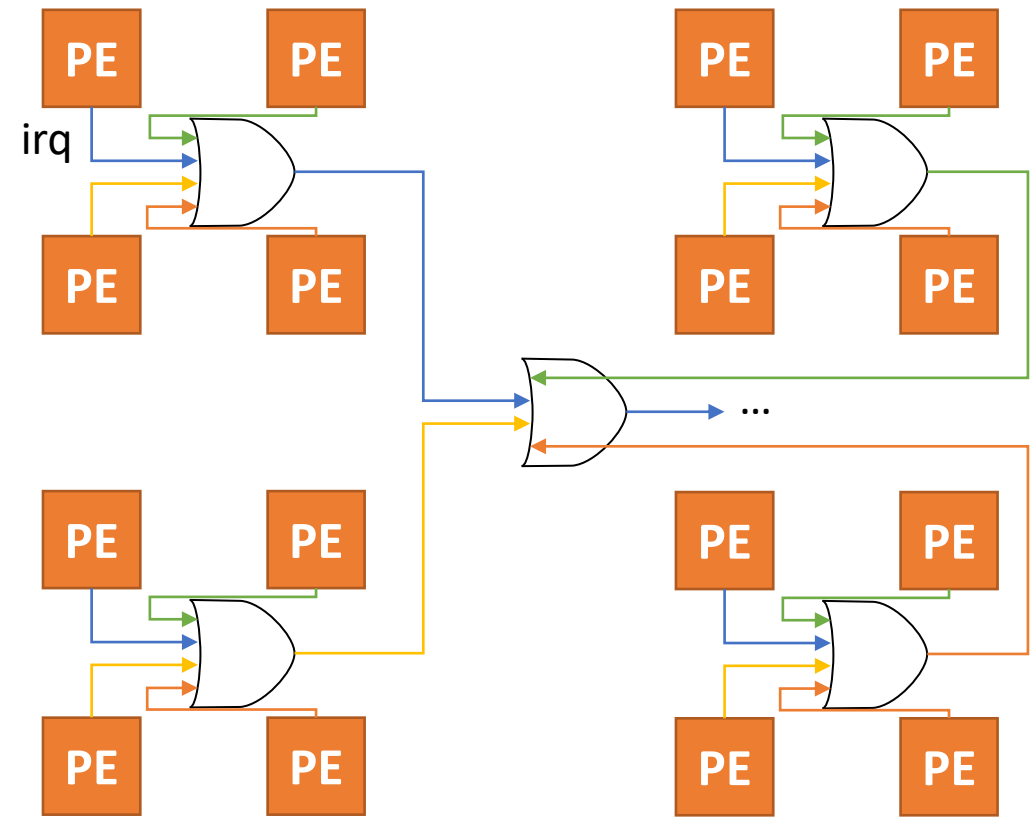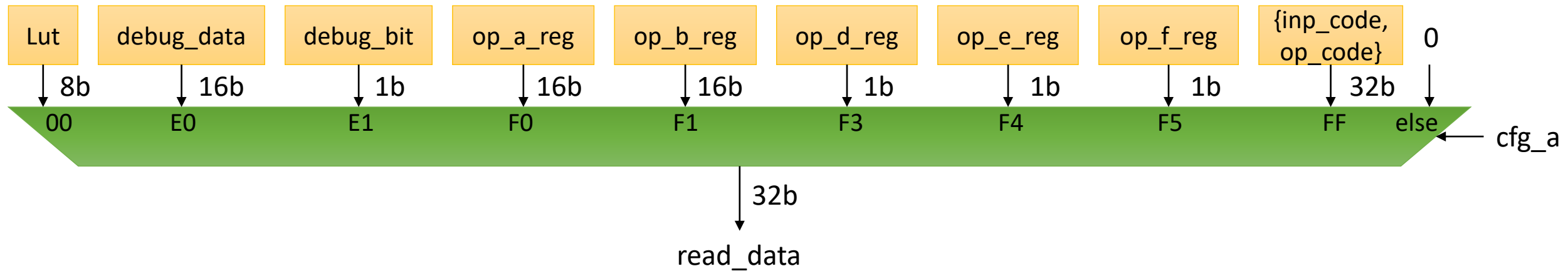| Z | ~Z | C | ~C | N | ~N | V | ~V | C&~Z | ~C\|Z | N==V | N!=V | ~Z\|(N==V) | Z\|(N!=V) | LUT | ALU res_p |
|---|----|----|-----|----|-----|----|-----|------|-------|------|------|-----------|-----------|-----|-----------|
| EQ | NE | CS | CC | MI | PL | VS | VC | HI | LS | GE | LT | GT | LE | LUT | ALU |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

← flag_sel

1b

res_p

# Breakpoint in Diablo

# Global Signal Flow

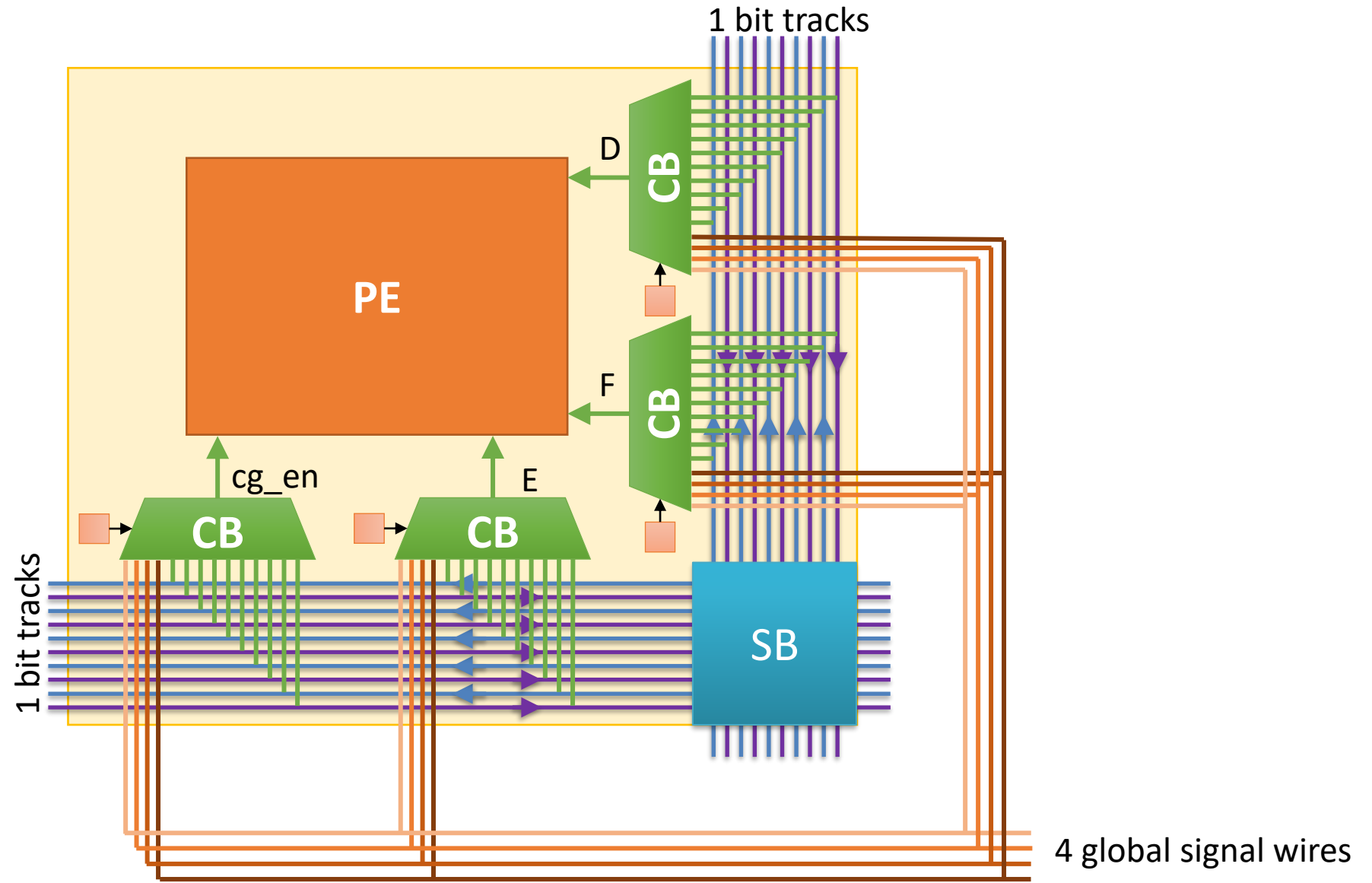- irq (interrupt request) signal is connected to global signal tiles (GST) in a hierarchical tree fashion

# Reading Diablo registers in debug

| Lut | debug_data | debug_bit | op_a_reg | op_b_reg | op_d_reg | op_e_reg | op_f_reg | {inp_code, op_code} | 0 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 8b | 16b | 1b | 16b | 16b | 1b | 1b | 1b | 32b | |
| 00 | E0 | E1 | F0 | F1 | F3 | F4 | F5 | FF | else |

cfg_a

32b

read_data

# Clock Gating in Diablo

- Each tile has 4 1 bit global signals coming in

- These go into all three 1 bit connection boxes

- There is another connection box that outputs the cg_en (clock gating enable) signal, reverse of  which is the clk_en signal sent to the PE core
  - This connection box chooses this cg_en signal from 10 of the routing tracks on side 1 and the 4 global signals
  - Block diagram is on next slide

# 1 Bit Connection Boxes in Diablo

1 bit tracks

D

PE

CB

F

CB

cg_en

E

CB

CB

1 bit tracks

SB

4 global signal wires

# Improvements in Lassen over Diablo

- Lassen supports 16 bit Bfloat type
  - 1 b sign, 8 b exponent, 7 bit mantissa
  - Same range as Float32

- Bfloat add and multiply natively

- Micro-instructions to implement Bfloat transcendentals (divide, exp, power, log, sin, cos)

- Optimization passes (currently manual) to improve efficiency

- Proposed improvements : 16-bit LUT for masking operations

- Need review: Are the operations in lassen sufficient for implementing target applications

# Diablo vs Lassen

| Diablo Ops | New Ops in Lassen |
|---|---|
| Add | FPAdd |
| Sub | FPMult |
| Abs | |
| GTE_Max | FGetMant |
| LTE_Min | FAddIExp |
| Sel | FSubExp |
| Mult0 | FCnvExp2F |
| Mult1 | FGetFInt |
| Mult2 | FGetFFrac |
| SHR | |
| SHL | |
| Or | |
| And | |
| XOr | |

Micro-instructions

1. Allows easy porting of applications!
2. Enables new applications that need a large range like DNN training or complex ops like divide

**Other FP instructions requested by Jeff:**
**sub,**
**le lt, ge, gt, neq, neq,**
**neg, flr, ceil, abs, min, max, sqr**

# Transcendentals Example - Divide

```
divident = Data(0x4288)  #68

divisor  = Data(0x4020)  #2.5

quotient = Data(0x41D9)  #68/2.5 = 27.2


mant          = pe_get_mant(asm.fgetmant(), divisor, Data(0))

lookup_result = mem_lut.div_lut(mant)

scaled_result = pe_scale_res(asm.fsubexp(), lookup_result, divisor)

quotient.     = pe_mult(asm.fp_mult(), scaled_result, divident)
```

# Bfloat LUTs

```python
class tlut:
    def div_lut(self, index):
        return {
        0   : 0x3f80,
        1   : 0x3f7e,
        2   : 0x3f7c,

         …
        126 : 0x3f01,
        127 : 0x3f00
        } [index]

    def ln_lut(self, index):
```

# Memory

Max Strange

# Memory Core in Jade

- There is a mux that chooses which one bit control signal goes out
  - 0, 1, 2, valid, full, almost_full
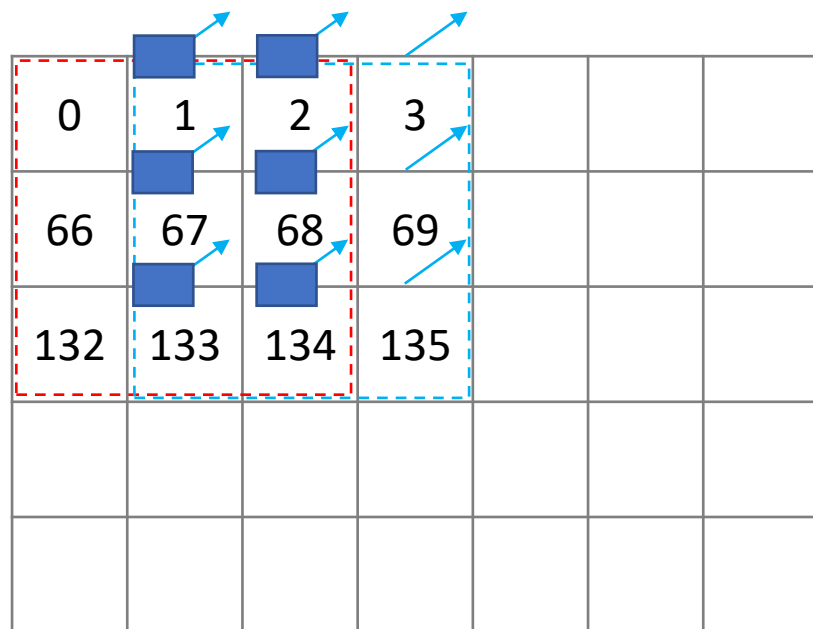- Two banks both 16 bit wide and 512 entry
- There is a read and a write phase

# Memory

9 output ports

# Global Buffer

Taeyoung Kong

# Global Buffer – Main Functions

- Highly banked buffer for **data reuse**
- Buffer to store bitstream for **parallel reconfiguration**

- Design Philosophy: Make global buffer simple and flexible

# Global Buffer – Components

**Global Buffer**

AXI – bank Interconnect

**Bank Array**

...

| 0 | 1 | 2 | 3 | 28 | 29 | 30 | 31 |

32 SRAM banks
4 MB total
64 bit wide

8 I/O controller > bank interconnect

I/O controller_0    I/O controller_1    ...    I/O controller_6    I/O controller_7

# Global Buffer – SoC interface

- CGRA appears as a Memory to the SoC

To/from SoC

wdata (64b)   waddr   wen   raddr   ren   rdata (64b)

Address Decoders

| 0 | 1 | 2 |  | 31 |

- SoC provides a wrapper that converts between this simple interface into AXI interface which has block transfers
  - Block = 256 words * 64 bits/word
- Pipelined to support high clock frequency

# Global Buffer – I/O controller



**Global Buffer**

Weights

Inputs

0 1 2 3 4 28 29 30 31

...

I/O controller_0 ... I/O controller_7

**Configurable in modes, number of types, capacity for different data**

din addr addr dout

0 1 2 3 4 ... 28 29 30 31

16 x 32 CGRA

# Global Buffer – I/O controller

- To CGRA, we send out data and valid
- From CGRA, we get data and valid



<FIFO mode>
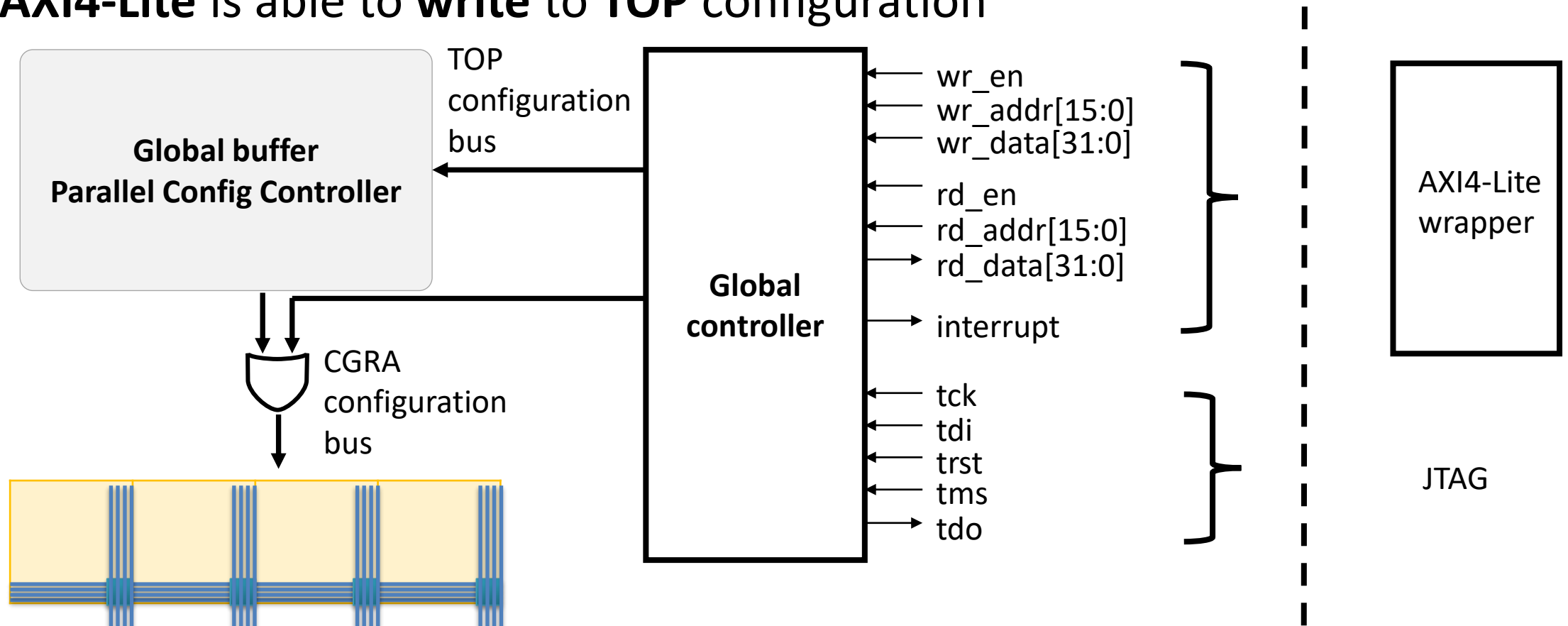
<SRAM mode>

# Global Buffer – Parallel Reconfiguration



32 SRAM Banks

Decode Logic
Generates rden, addr

Global Address

8 config controllers

Configuration Bus
{config_write,
32b config_addr,
32b config_data}

Configuration Bus
from Global
Controller

32 CGRA Columns

# Global Controller

Taeyoung Kong & Alex Carsello

# Configuration – AXI4-Lite & JTAG

- **JTAG** is able to **write/read** to **CGRA** configuration and **TOP** configuration
- **AXI4-Lite** is able to **write** to **TOP** configuration

# CGRA control – start & done

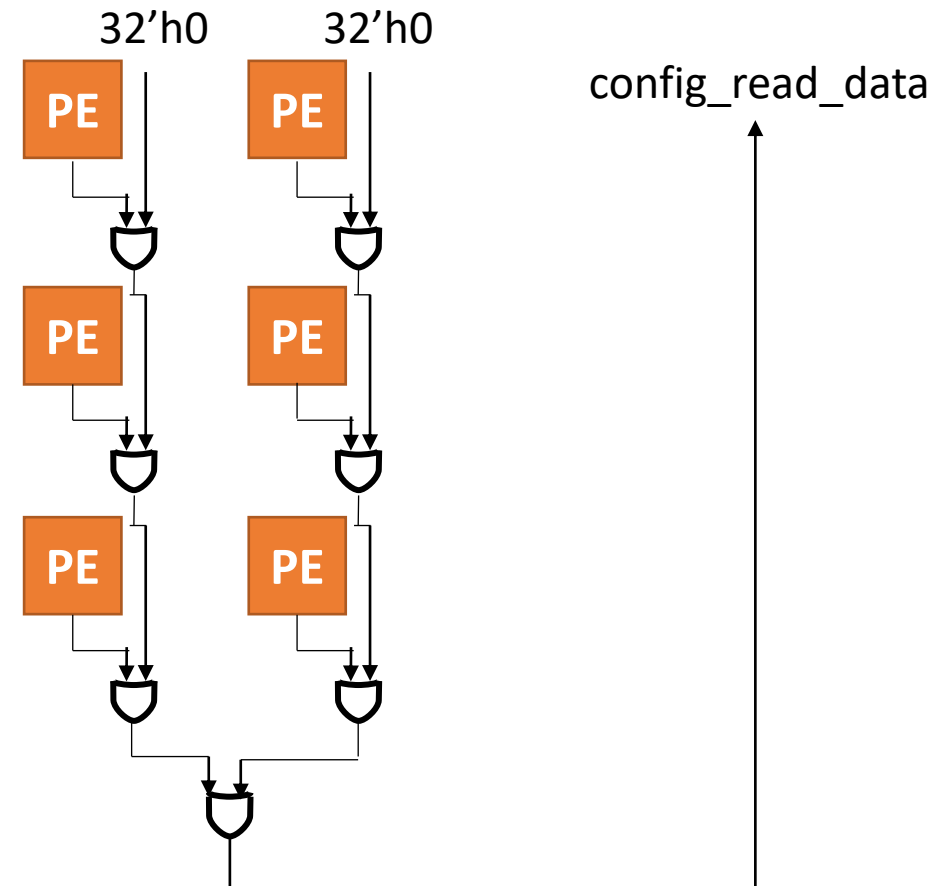- **Block level** of application is controlled by **cgra_start** and **cgra_done**

# Global Signal Flow

# Debug Support by JTAG

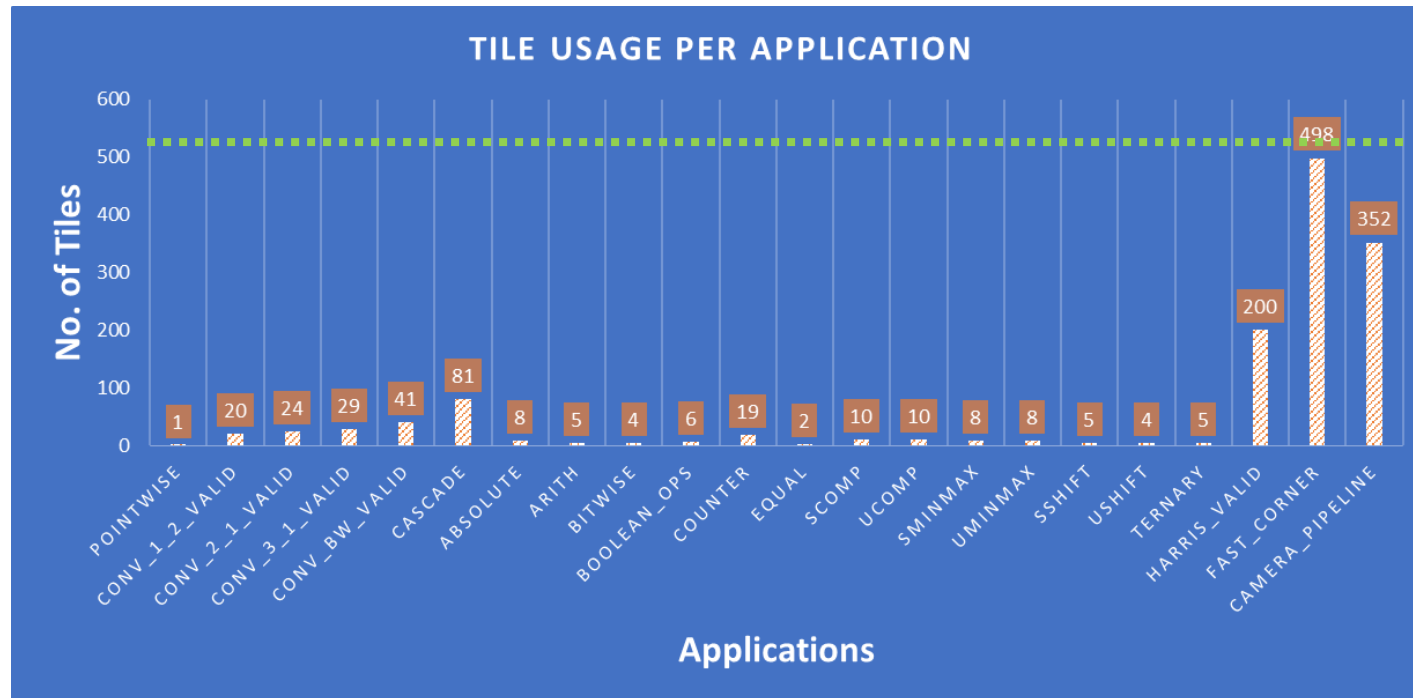- JTAG can do everything!

- Debug CGRA fabric
  - Write/read CGRA-configuration registers
  - Write/read SRAM in CGRA fabric
  - 4bit stall signal
    - Every tile has configuration registers to select a stall signal from 4bits

- Debug CGRA fabric + global buffer
  - Write/read TOP-configuration registers
  - Write/read SRAM in global buffer
  - Write/read CGRA control register (cgra_start)
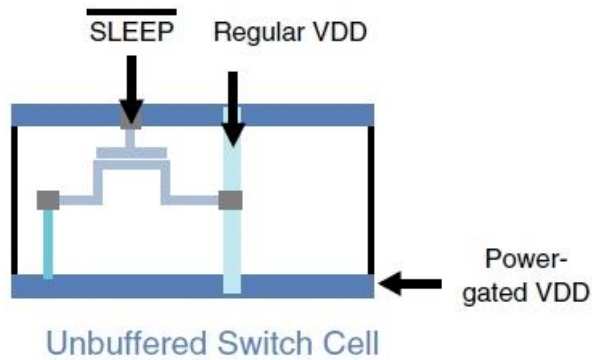
# Power Domains

Ankita Nayak

# Motivation

- Adoption of 'Design for PPA' mindset

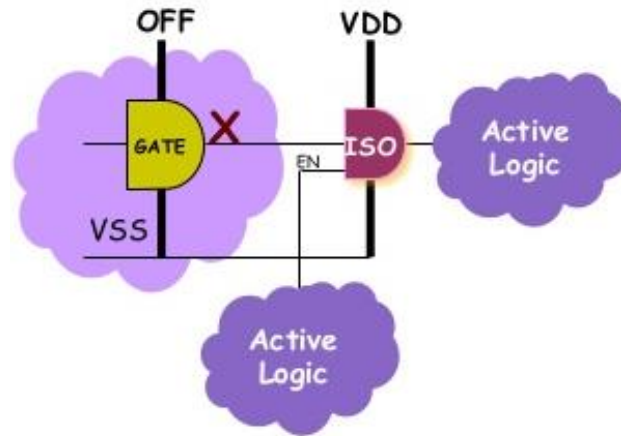- Save power by turning off tiles in the fabric not actively used
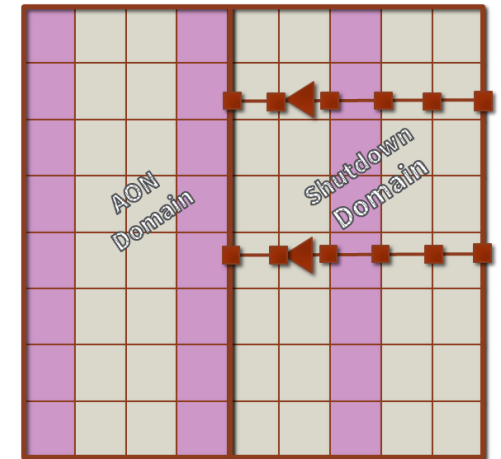
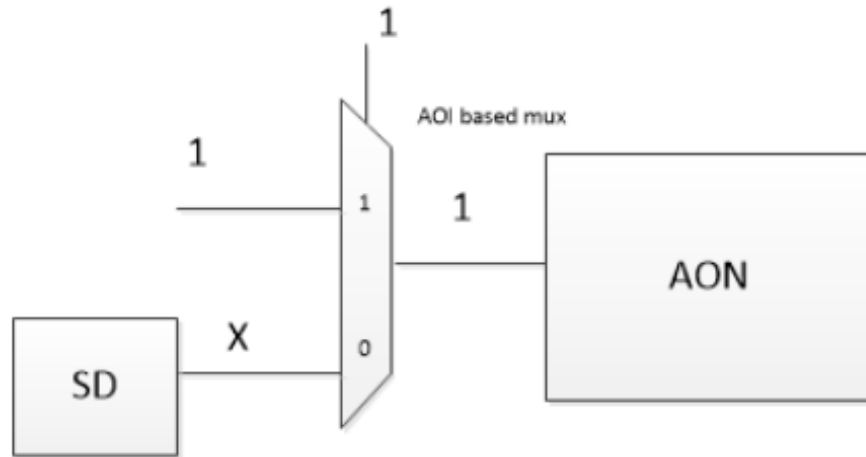# Architectural Features

1. Power Switch Network Architecture



SLEEP    Regular VDD

Power-gated VDD

Unbuffered Switch Cell

2. Boundary protection between the active and shutdown tiles



OFF      VDD

GATE      X      ISO      Active Logic

EN

VSS

Active Logic

3. FT to connect IO near shutdown logic to the active part & for global signals



AON Domain      Shutdown Domain
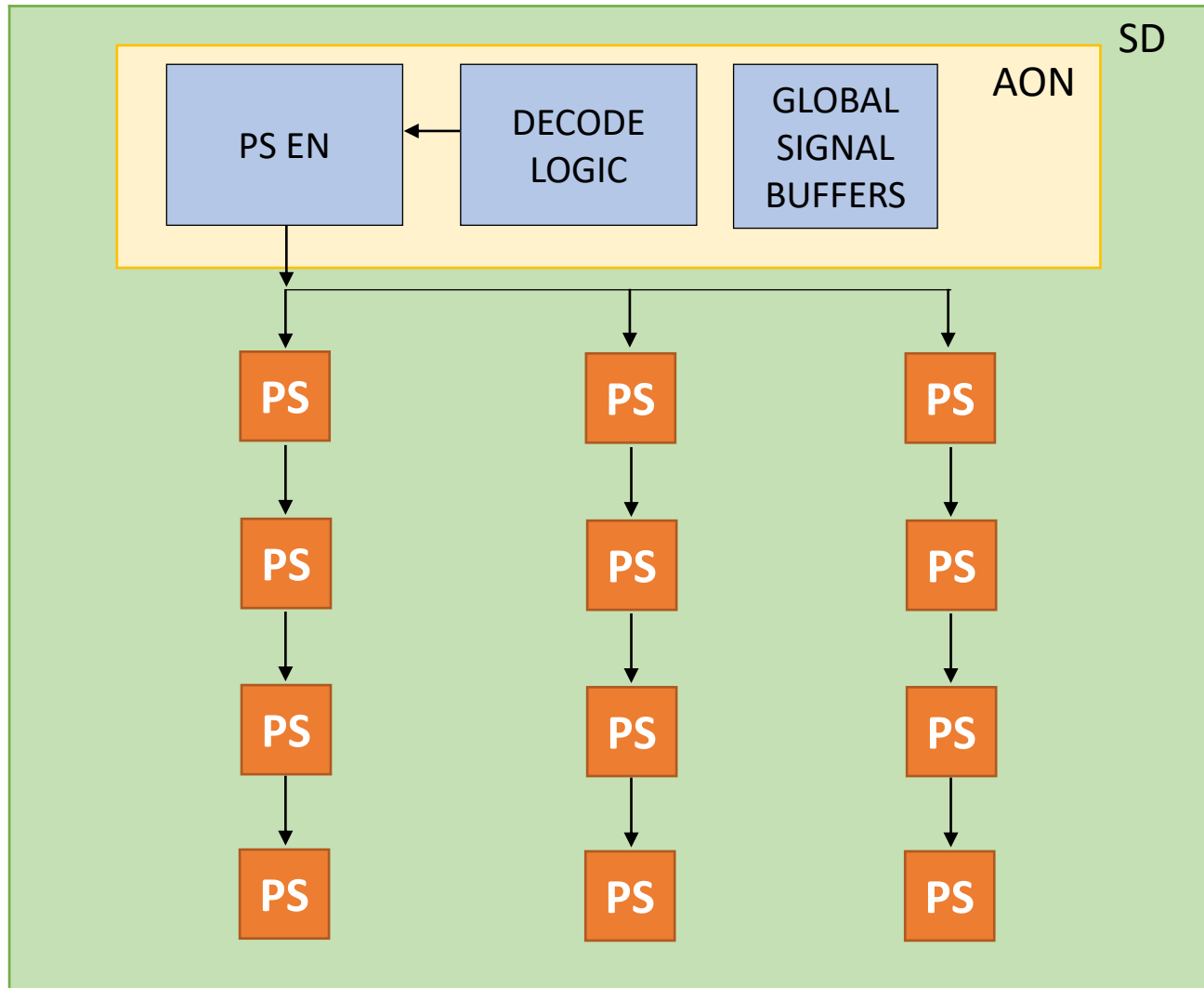
# ON-OFF Boundary Protection



1. Redesign the muxes in the design to allow for inherent clamping mechanism

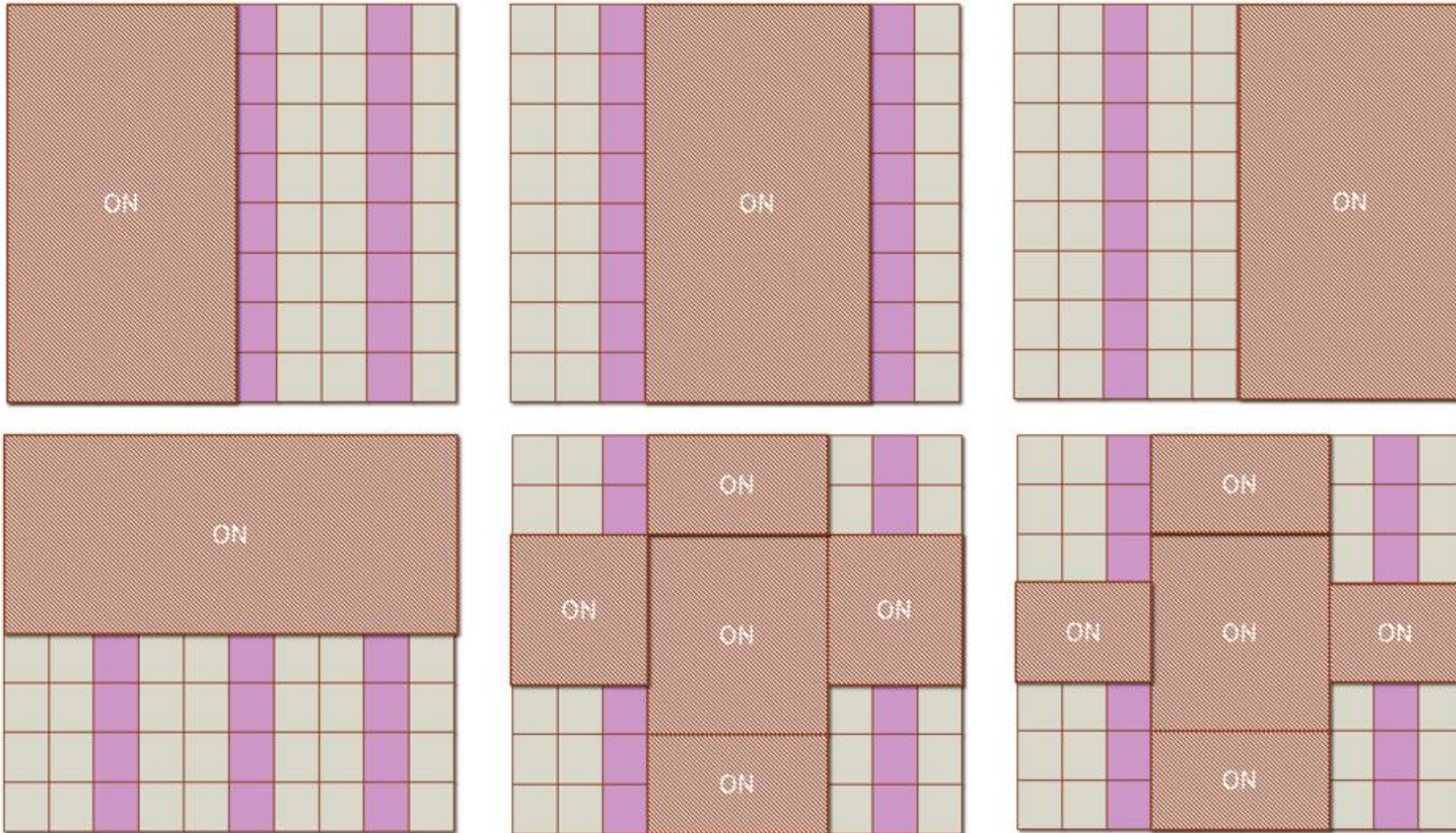2. Use the existing muxes in the design for boundary protection

AON : Always ON
SD    : Shutdown

# Tile Architecture



- Configuration register and its decode logic for enabling the power switches is kept ON

- Buffers and logic on global signals are kept ON

- All other logic is shutdown when the tile is turned OFF

# Top Level – Different sections can be turned ON/OFF by P&R

# Summary of Garnet architectural features

1. Support for Bfloat16, and for executing complex operations like divide using multiple PEs

2. Addition of a global buffer to create a memory hierarchy for efficiently executing neural networks

3. Fast reconfiguration support using global buffer and control processor

4. Addition of configurable power domains