# Automatic Pipelining for CGRA Applications

## Techniques, Analysis, and Future Directions

Jack Melchert

# Motivation

Problem: The applications running on our CGRA run at very low frequencies

| | |
|---|---|
| Gaussian Blur | 103 MHz |
| Harris Corner Detection | 26 MHz |
| Camera Pipeline | 17 Mhz |

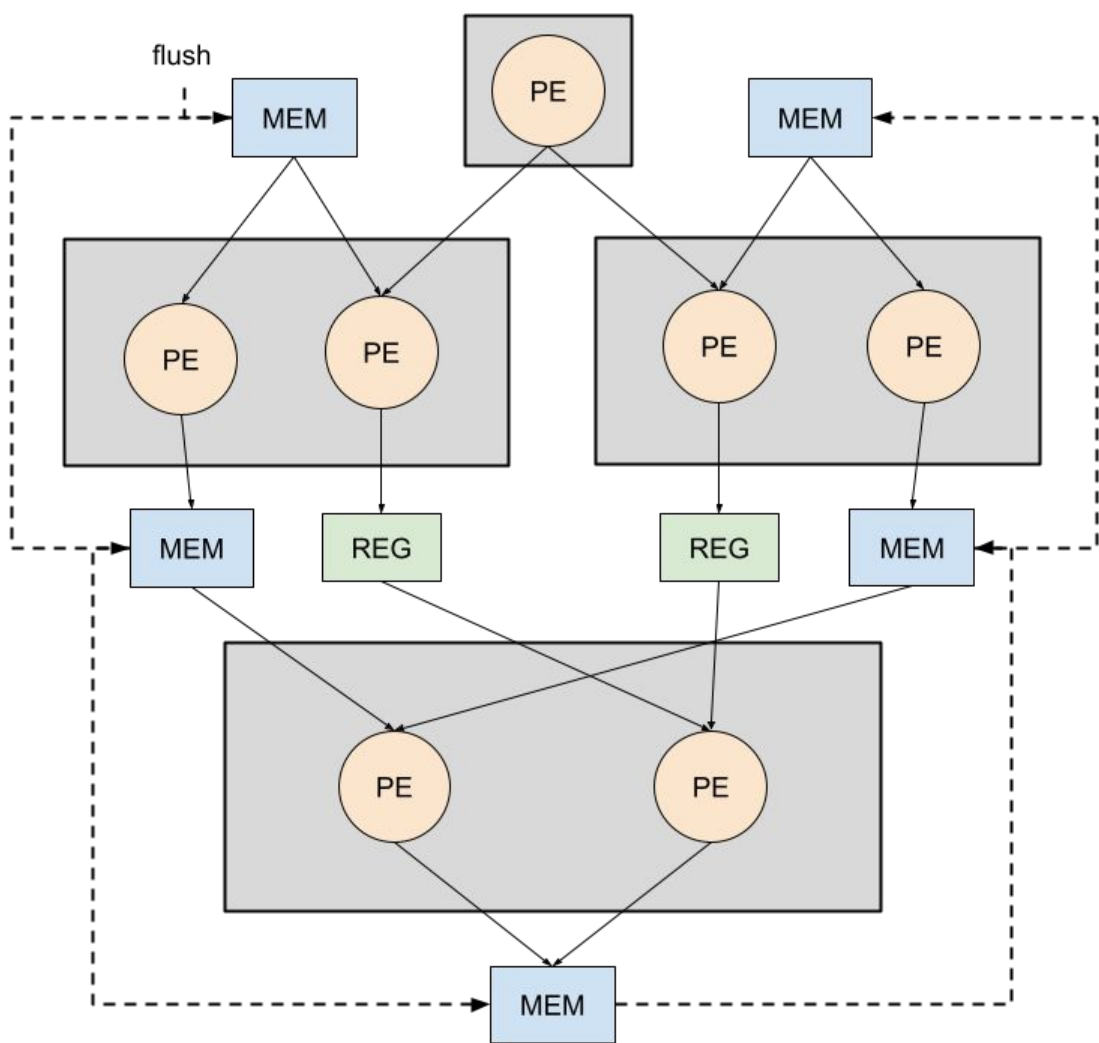The maximum frequency of the CGRA is 780 MHz at 1.1 V

# Motivation

Solution: Pipeline applications

- Un-pipelined applications have critical paths that typically start or end at a memory tile or outside of the array
- Adding pipelining registers to these paths should dramatically increase the maximum frequency of the applications
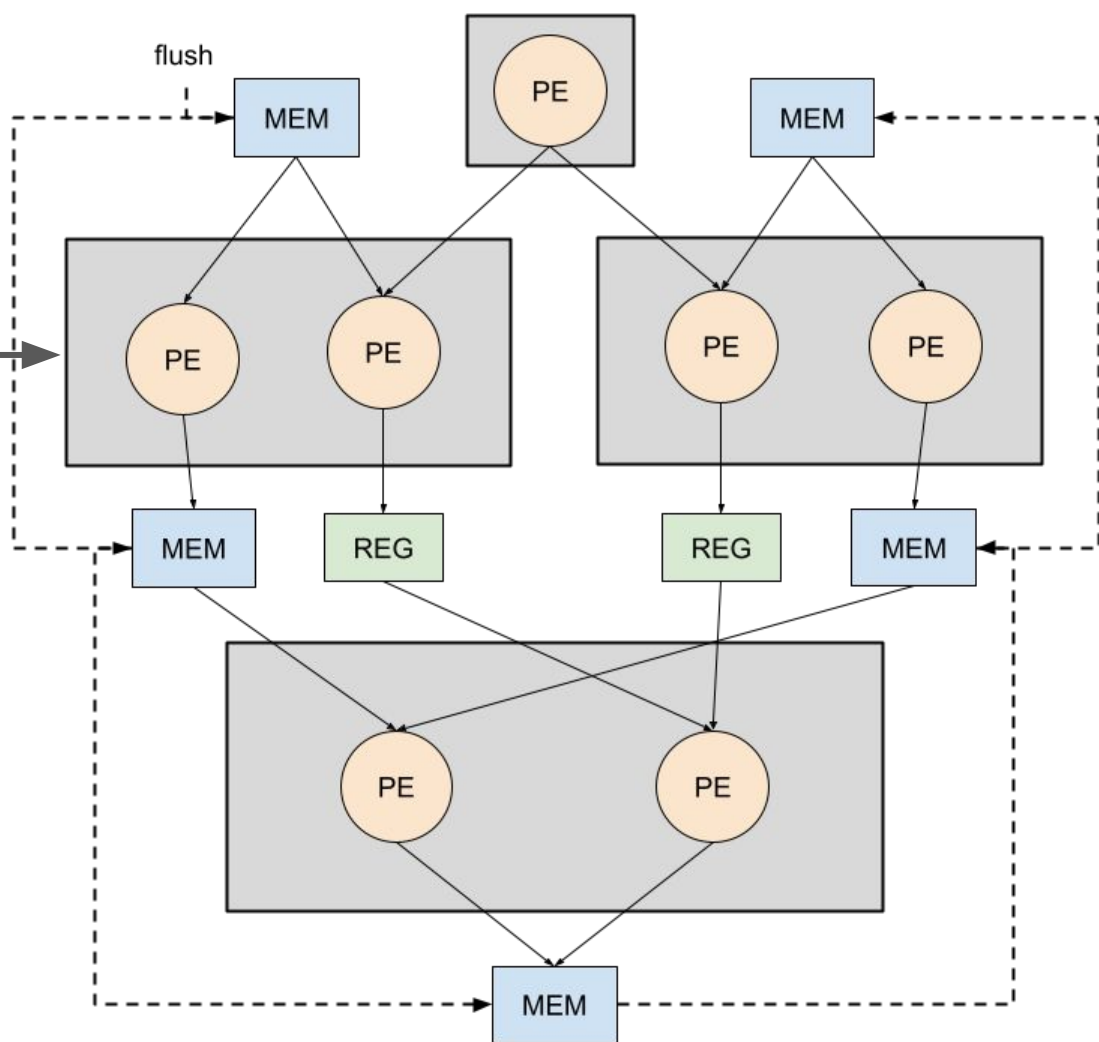
# Goals

- Understand and model the delays in the array
    - Measure all contributions to the critical path
    - Integrate that data into a timing analysis tool that will accurately model the critical path in an application
- Create an automated pipelining approach
    - Our applications and hardware change frequently, so an automated approach is required
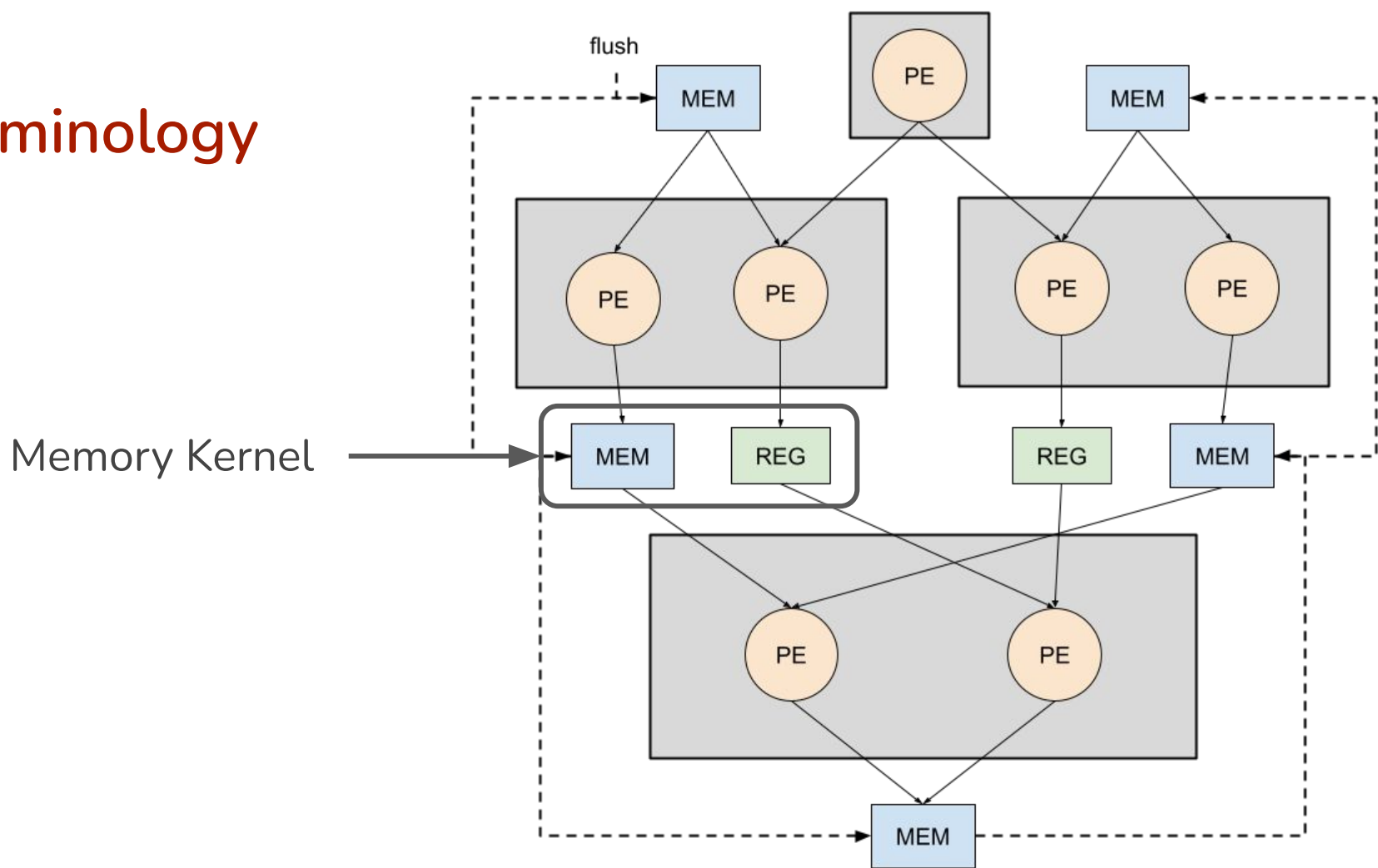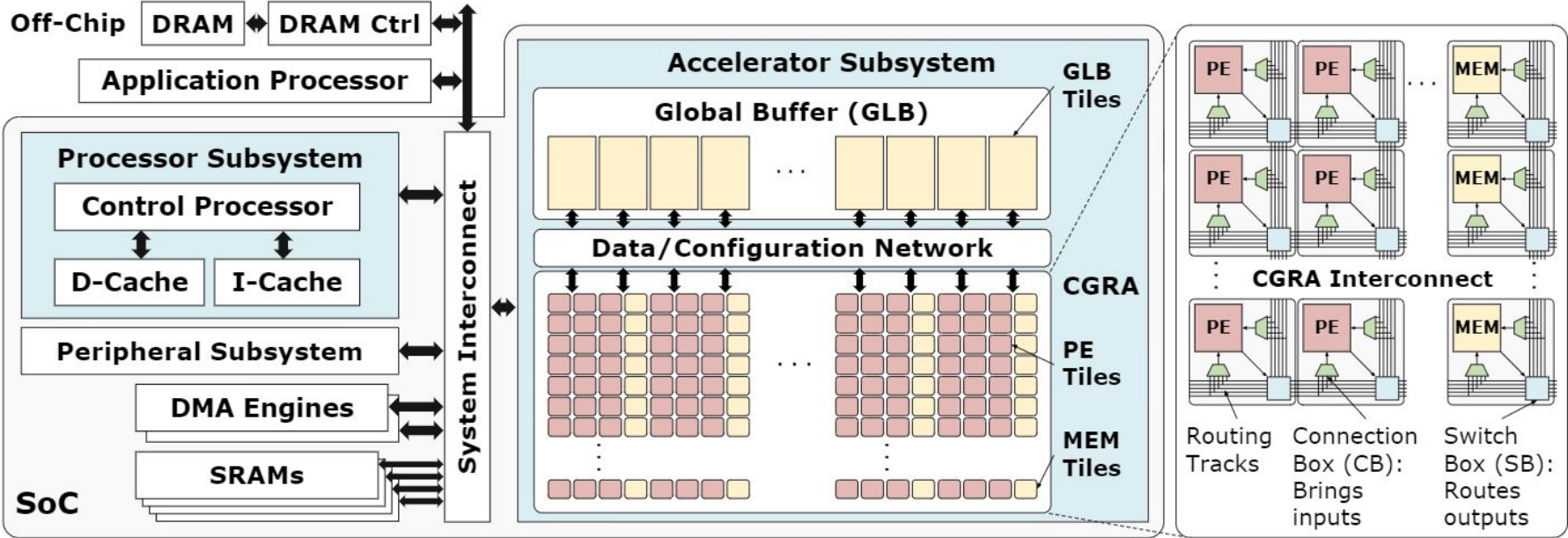
# Terminology

# Terminology



flush

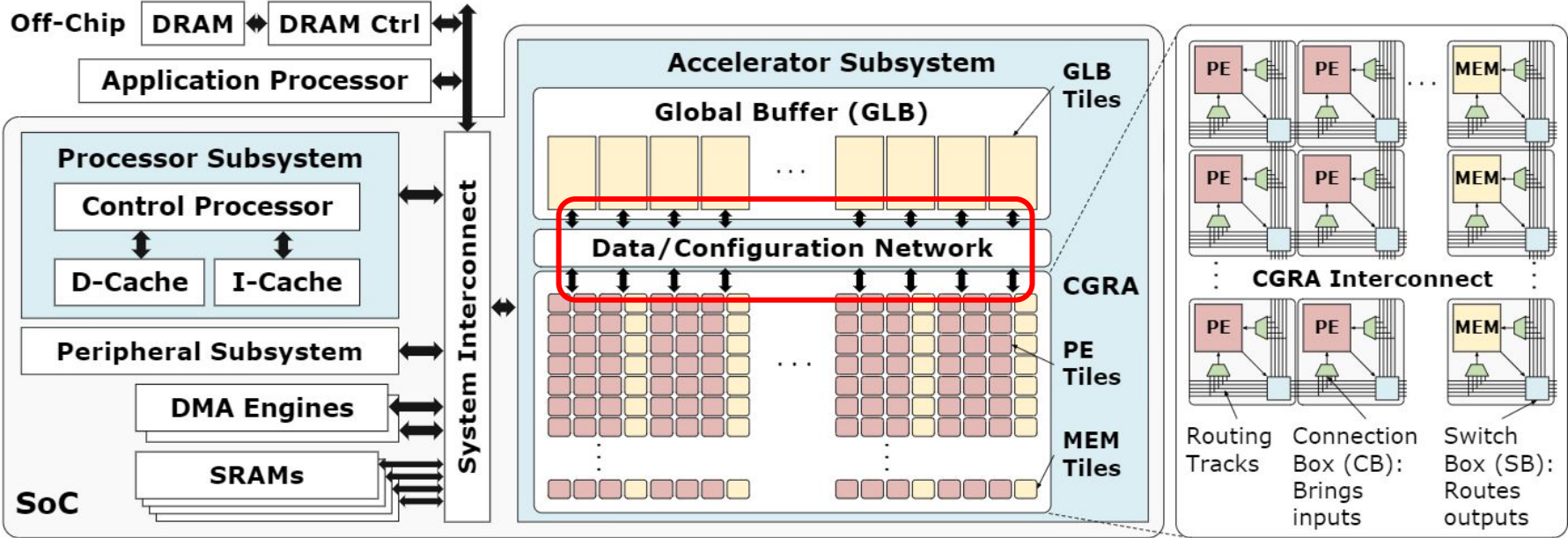Compute Kernel

# Terminology

flush



Memory Kernel

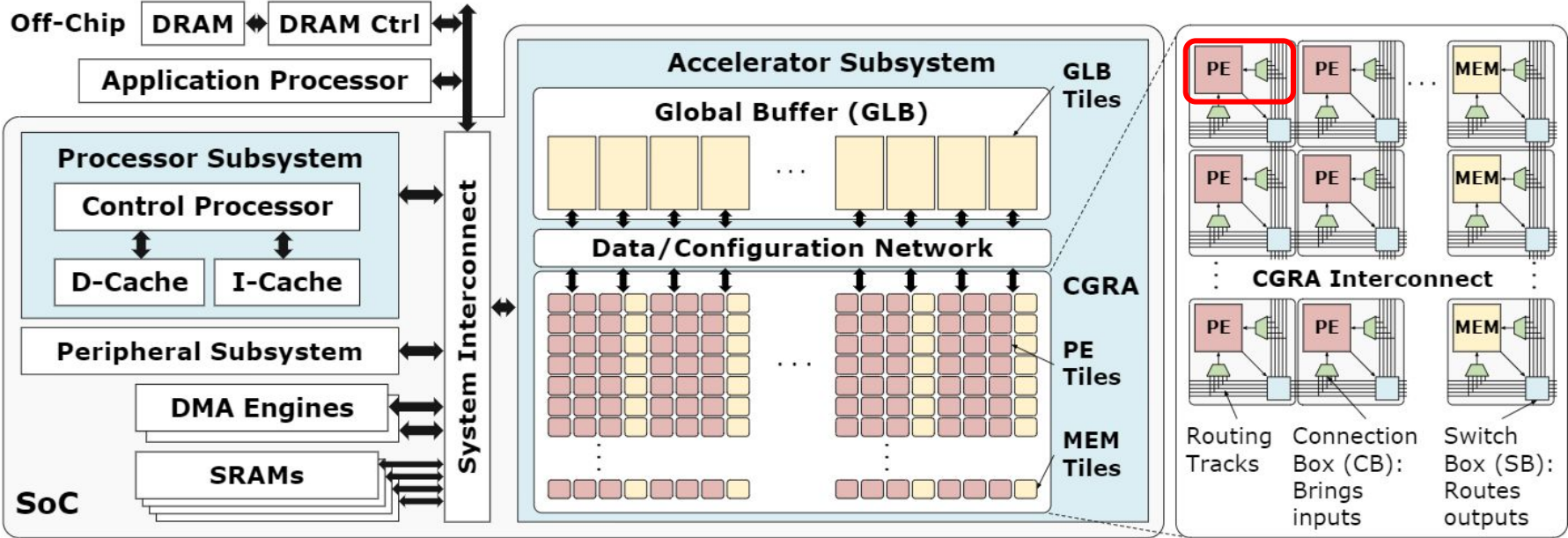# Contributions to Critical Path Delay

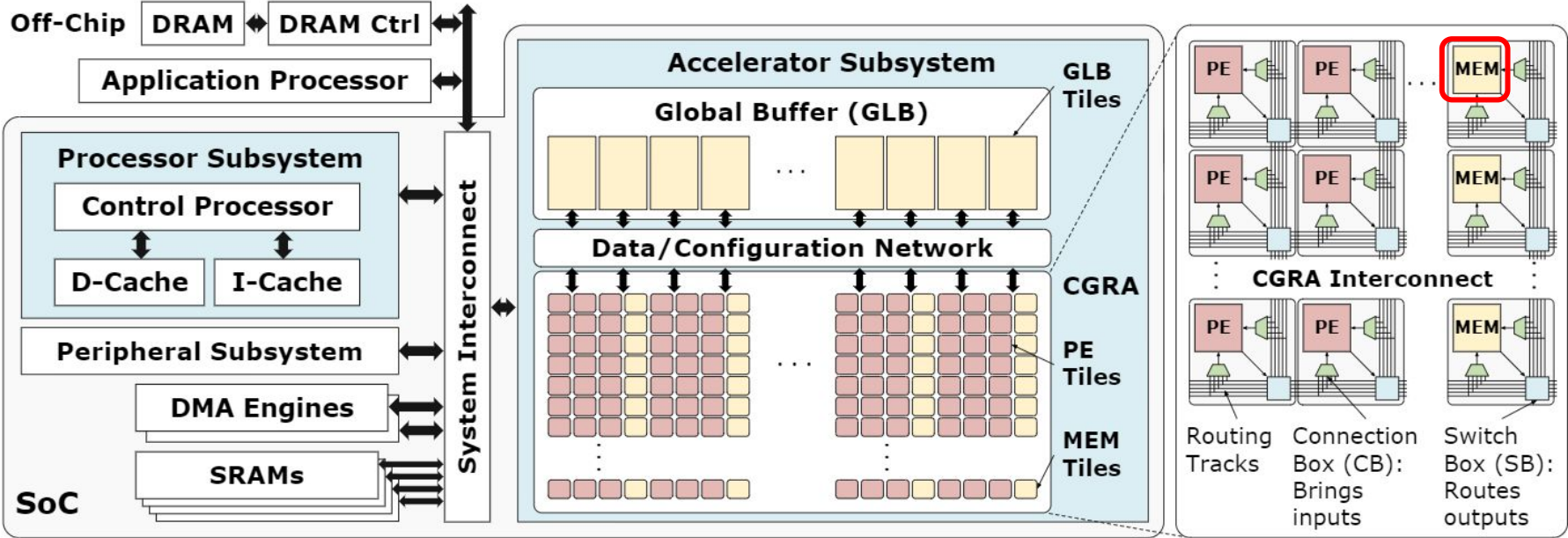# Contributions to Critical Path Delay



Global Buffer Delay: 1.1 ns
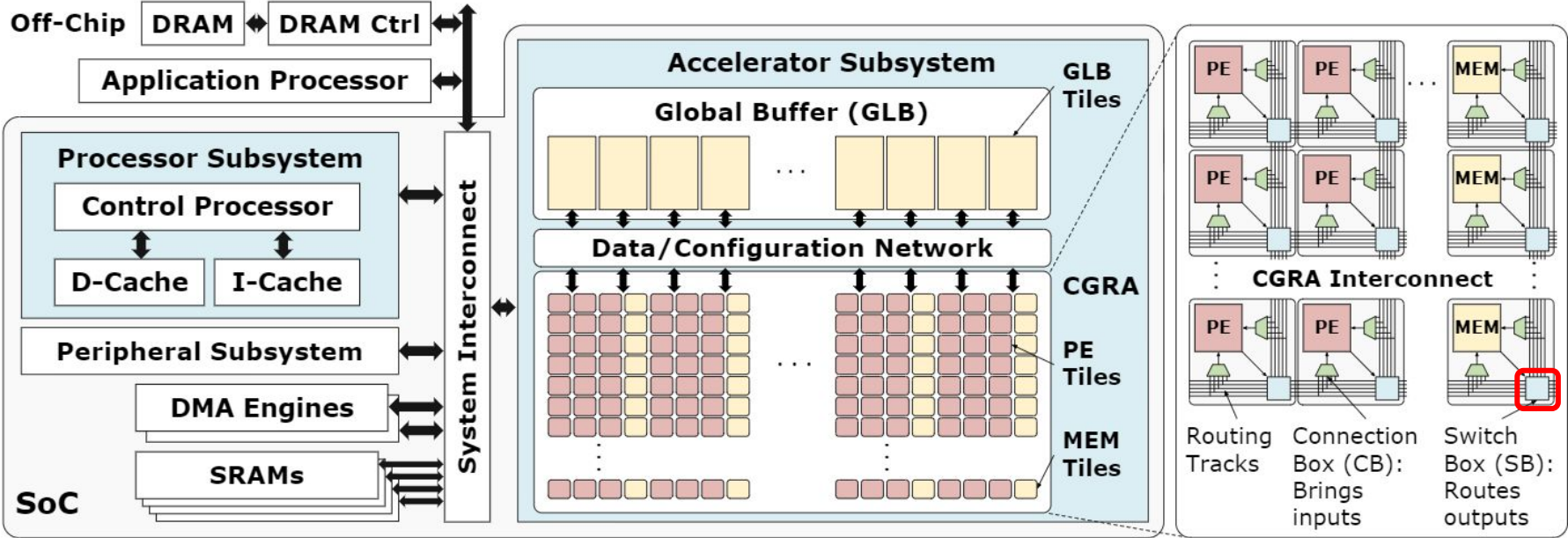
# Contributions to Critical Path Delay



PE Delay: 0.48-0.70 ns

# Contributions to Critical Path Delay



MEM Delay: 0 ns
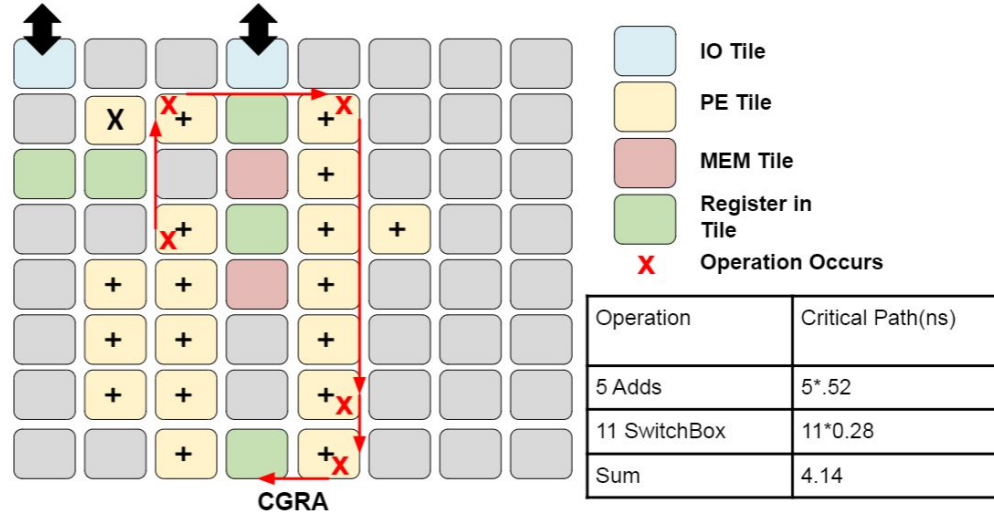
# Contributions to Critical Path Delay



SB Delay: 0.14 ns

# Critical Path Model

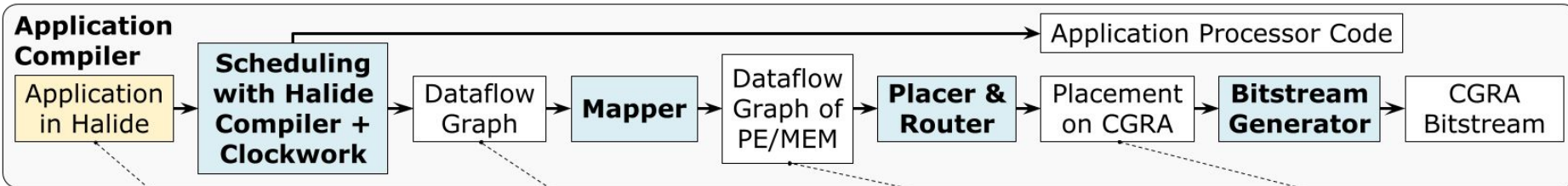Contributions:

- Global Buffer: 1.1 ns
- Switch Boxes: 0.14 ns
- PEs: 0.48-0.7 ns

After about 3 hops in the interconnect, a critical path with 1 PE will be dominated by the interconnect



| | IO Tile |
| | PE Tile |
| | MEM Tile |
| | Register in Tile |
| X | Operation Occurs |

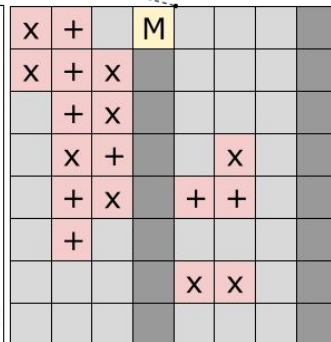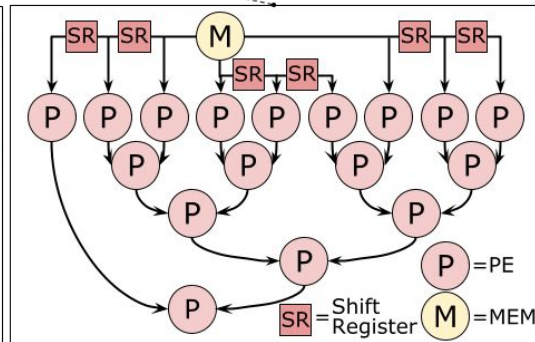| Operation | Critical Path(ns) |
| --- | --- |
| 5 Adds | 5*.52 |
| 11 SwitchBox | 11*0.28 |
| Sum | 4.14 |

CGRA

# Critical Path Model



```
conv(x, y) += kernel(r.x, r.y) *
              input(x+r.x, y+r.y);
conv.in().compute_root();
conv.in()
    .tile(x,y,xo,yo,xi,yi,64,64)
    .hw_accelerate(xi, xo);
conv.update()
    .unroll(r.y, 3)
    .unroll(r.x, 3);
conv.compute_at(conv.in(), xo);
input.stream_to_accelerator();
```
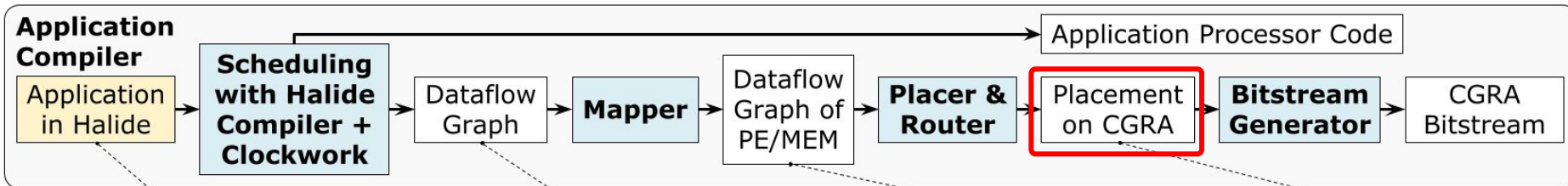
Example: 3x3 Convolution

# Critical Path Model



Example: 3x3 Convolution

We must construct the delay model using the output of the place and route tool

# Critical Path Model - Evaluation

| Application | Modeled (MHz) | Measured (MHz) | % Difference |
|---|---|---|---|
| Gaussian v1 | 158 | 200 | -21.0% |
| Gaussian v2 | 295 | 280 | 5.4% |
| Gaussian v3 | 515 | 420 | 22.6% |
| Gaussian v4 | 793 | 600 | 32.2% |
| Harris v1 | 30 | 25 | 20.0% |
| Harris v2 | 137 | 160 | -14.4% |
| Harris v3 | 335 | 300 | 11.7% |
| Harris v4 | 373 | 360 | 3.6% |

# Critical Path Model - Evaluation

| Application | Modeled (ns) | Measured (ns) | % Difference |
|---|---|---|---|
| Gaussian v1 | 6.33 | 5.00 | 26.6% |
| Gaussian v2 | 3.39 | 3.57 | -5.1% |
| Gaussian v3 | 1.94 | 2.38 | -18.4% |
| Gaussian v4 | 1.26 | 1.67 | -24.3% |
| Harris v1 | 33.33 | 40.00 | -16.7% |
| Harris v2 | 7.30 | 6.25 | 16.8% |
| Harris v3 | 2.99 | 3.33 | -10.4% |
| Harris v4 | 2.68 | 2.78 | -3.5% |

# Pipelining Techniques and Analysis

1. Compute Pipelining

2. Broadcast Pipelining

3. Placement Algorithm Tweaks

4. Post-PnR Pipelining

5. Register File Pipelining

# Compute Pipelining

- At compute mapping, we know how many PEs we will use and how they are connected
- Have all information needed to do branch delay matching
  - Ensures all paths from one memory to another are the same number of cycles
- All registers are added to the compute graph
  - Need to be packed into the PEs or placed onto the routing fabric

# Compute Pipelining

|  |  | Baseline | Compute Pipelining |
|---|---|---|---|
| Harris | Clk Freq (MHz) | 26 | 147 |
| Harris | Resource Utilization (PE/Mem/Reg) | 91/6/30 | 91/6/66 |
| Camera Pipeline | Clk Freq (MHz) | 17 | 32 |
| Camera Pipeline | Resource Utilization (PE/Mem/Reg) | 281/34/8 | 281/34/138 |
| Gaussian | Clk Freq (MHz) | 103 | 164 |
| Gaussian | Resource Utilization (PE/Mem/Reg) | 160/8/136 | 160/8/160 |

# Broadcast Signal Pipelining

# Broadcast Signal Pipelining

# Broadcast Signal Pipelining

| | | Baseline | Compute Pipelining | Broadcast Pipelining |
|---|---|---|---|---|
| Harris | Clk Freq (MHz) | 26 | 147 | 370 |
| | Resource Utilization (PE/Mem/Reg) | 91/6/30 | 91/6/66 | 91/6/86 |
| Camera Pipeline | Clk Freq (MHz) | 17 | 32 | 84 |
| | Resource Utilization (PE/Mem/Reg) | 281/34/8 | 281/34/138 | 281/34/206 |
| Gaussian | Clk Freq (MHz) | 103 | 164 | 291 |
| | Resource Utilization (PE/Mem/Reg) | 160/8/136 | 160/8/160 | 160/8/211 |

# Placement Algorithm Tweaks

- The cost metric for placement is total wirelength
  - Generally leads to good routability

- For pipelining we care about max wirelength, not total wirelength
  - Using max wirelength alone as a cost metric for placement leads to routability issues
  - Instead we can create a new cost metric that incorporates both maximum and total wirelength
  - Total wirelength$^n$ for n = [1...10] works well



half-perimeter wirelength = x+y

# Placement Algorithm Exponent Parameter Sweep

# Placement Algorithm Tweaks

|  |  | **Baseline** | **Compute Pipelining** | **Broadcast Pipelining** | **Placement Tweaks** |
|---|---|---|---|---|---|
| Harris | Clk Freq (MHz) | 26 | 147 | 370 | 444 |
|  | Resource Utilization (PE/Mem/Reg) | 91/6/30 | 91/6/66 | 91/6/86 | 91/6/86 |
| Camera Pipeline | Clk Freq (MHz) | 17 | 32 | 84 | 231 |
|  | Resource Utilization (PE/Mem/Reg) | 281/34/8 | 281/34/138 | 281/34/206 | 281/34/206 |
| Gaussian | Clk Freq (MHz) | 103 | 164 | 291 | 417 |
|  | Resource Utilization (PE/Mem/Reg) | 160/8/136 | 160/8/160 | 160/8/211 | 160/8/211 |

# Post Place and Route Pipelining

- Iteratively break the critical path determined by the STA tool
- Re-analyze and determine new critical path
- Continue until there are no more registers available to use on the interconnect
  - Adding pipelining registers to a placement result is not always possible because branch delay matching is required

# Post Place and Route Pipelining

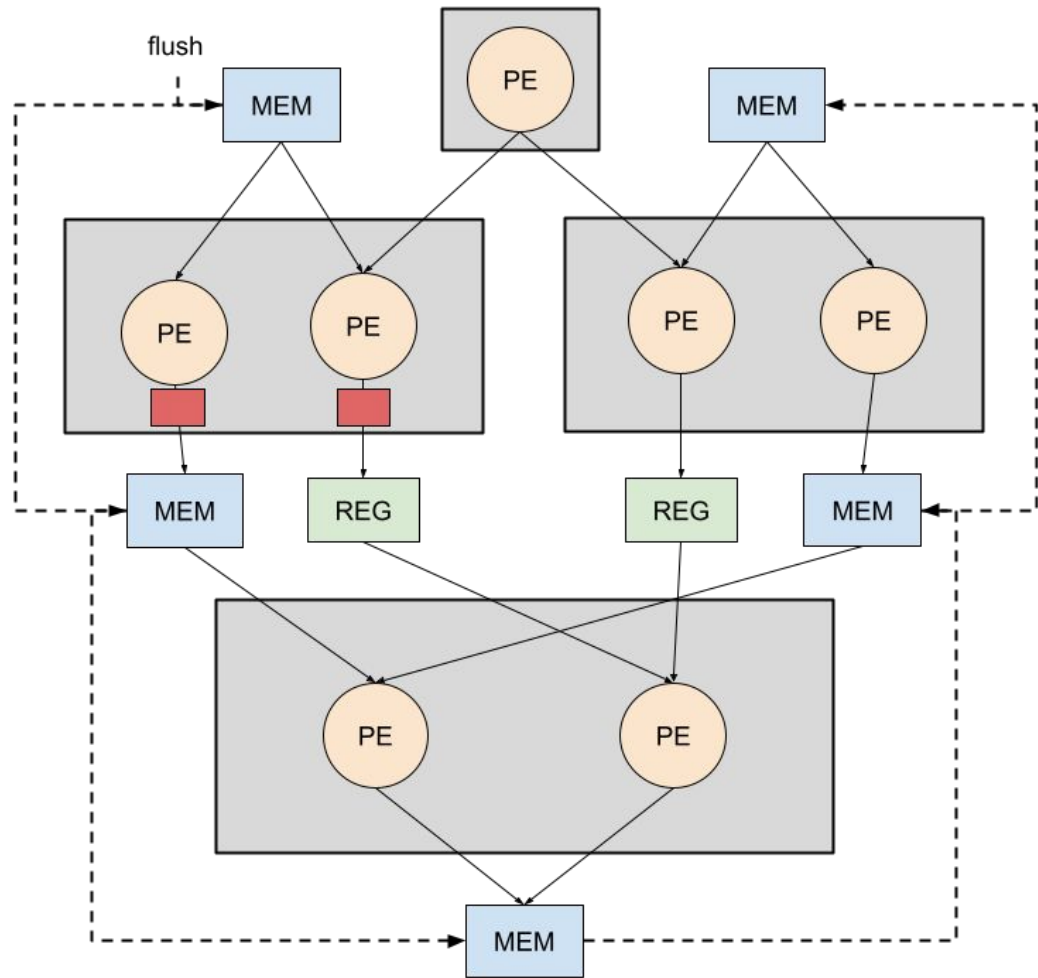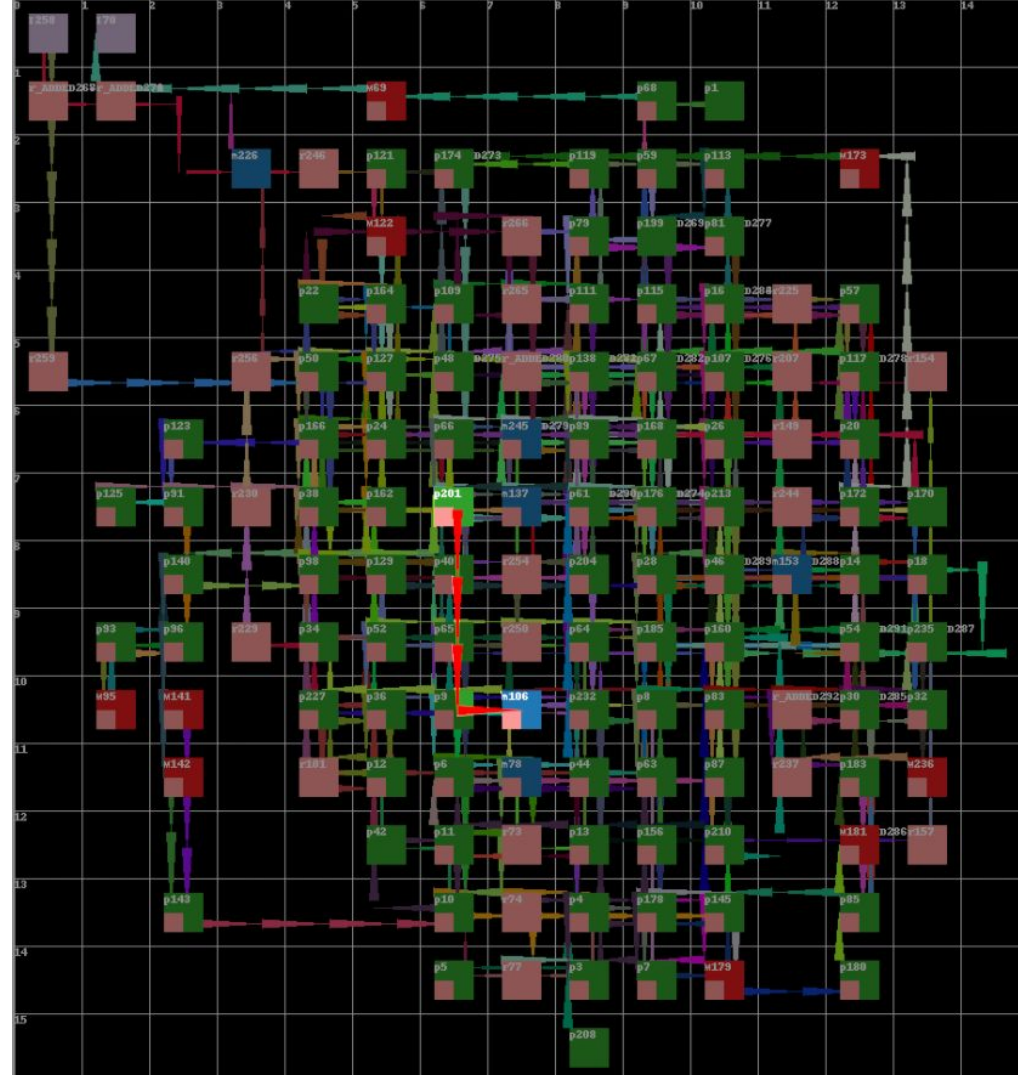| | | Baseline | Compute Pipelining | Broadcast Pipelining | Placement Tweaks | PnR Pipelining |
|---|---|---|---|---|---|---|
| Harris | Clk Freq (MHz) | 26 | 147 | 370 | 444 | 706 |
| | Resource Utilization (PE/Mem/Reg) | 91/6/30 | 91/6/66 | 91/6/86 | 91/6/86 | 91/6/102 |
| Camera Pipeline | Clk Freq (MHz) | 17 | 32 | 84 | 231 | 231 |
| | Resource Utilization (PE/Mem/Reg) | 281/34/8 | 281/34/138 | 281/34/206 | 281/34/206 | 281/34/206 |
| Gaussian | Clk Freq (MHz) | 103 | 164 | 291 | 417 | 417 |
| | Resource Utilization (PE/Mem/Reg) | 160/8/136 | 160/8/160 | 160/8/211 | 160/8/211 | 160/8/211 |

# Register File Pipelining

- We have a register file in every PE in the CGRA
- If we write and read every cycle to the same address, we can use this reg file as a pipeline register
- If we write and read to offset incrementing addresses, we can use it as a variable length pipeline register chain

# Register File Pipelining

| | | Baseline | Compute Pipelining | Broadcast Pipelining | Placement Tweaks | PnR Pipelining | With Pipelining Ponds |
|---|---|---|---|---|---|---|---|
| Harris | Clk Freq (MHz) | 26 | 147 | 370 | 444 | 706 | 706 |
| | Resource Utilization (PE/Mem/Reg) | 91/6/30 | 91/6/66 | 91/6/86 | 91/6/86 | 91/6/102 | 91/6/84 |
| Camera Pipeline | Clk Freq (MHz) | 17 | 32 | 84 | 231 | 231 | 264 |
| | Resource Utilization (PE/Mem/Reg) | 281/34/8 | 281/34/138 | 281/34/206 | 281/34/206 | 281/34/206 | 281/68/76 |
| Gaussian | Clk Freq (MHz) | 103 | 164 | 291 | 417 | 417 | 471 |
| | Resource Utilization (PE/Mem/Reg) | 160/8/136 | 160/8/160 | 160/8/211 | 160/8/211 | 160/8/211 | 160/8/187 |

# Future Directions - Hardware Changes

- Using ready-valid signaling would make branch delay matching easier
  - Rarely need to add registers to branch delay match
  - FIFOs at tile inputs can be used to buffer data that arrives early

- Optimize global buffer to array path
  - Current path is by far the longest in the array
  - We would never be able to run applications at 1 GHz

# Future Directions - Compiler Changes

- Make the scheduler more adaptable to different delays
    - Our pipelining tools are currently over constrained
    - The limiting factor when adding registers post-pnr is branch delay matching

- Improve the place and route tool
    - There is too much variability in the max wirelength from run to run
    - In the current PnR tool, routability as the primary concern
    - We can add more emphasis on reducing maximum wirelength
    - Fix the global placement algorithm and add congestion estimation