

# Warp-Aware Trace Scheduling for GPUS

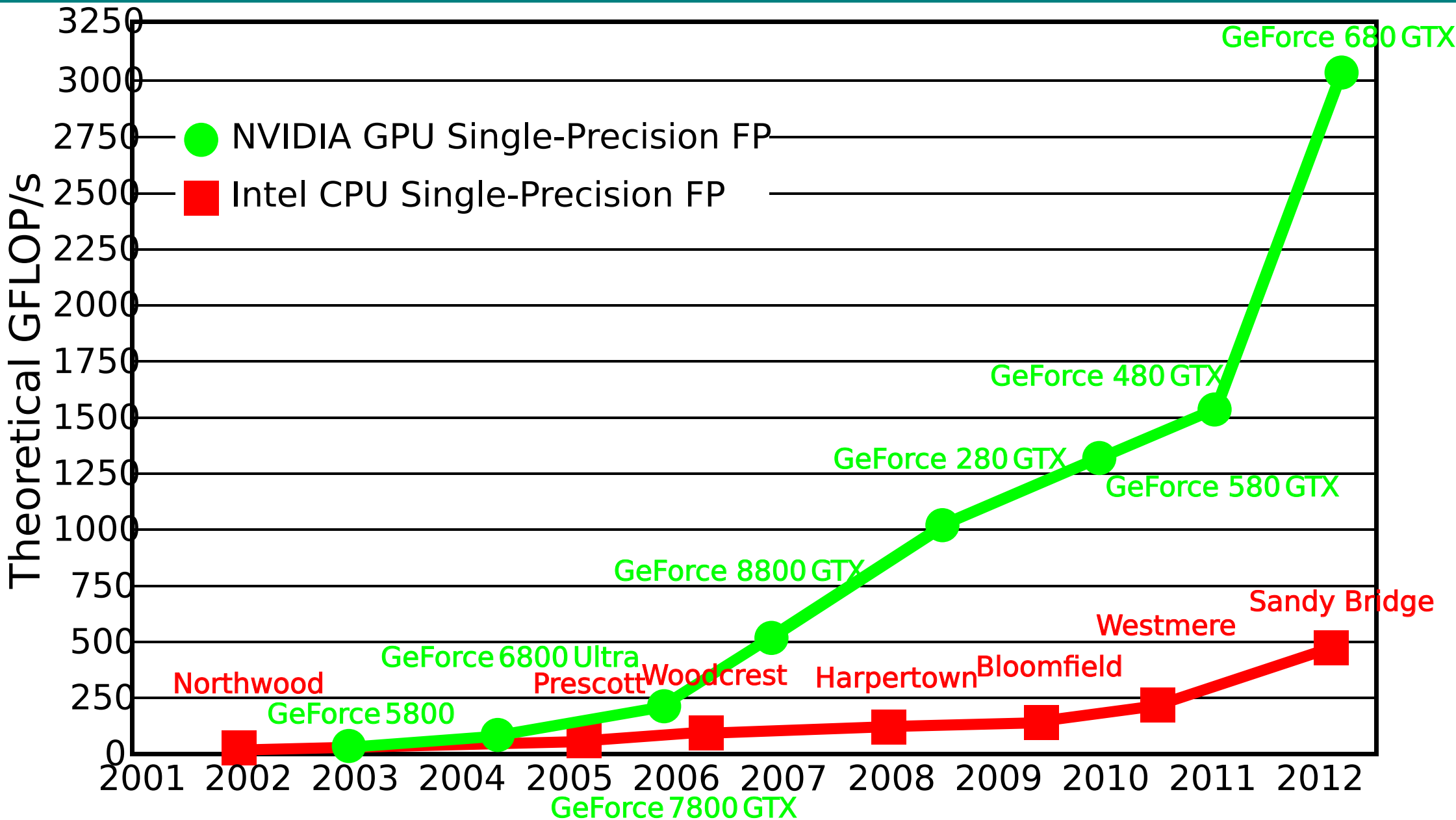
James Jablin (Brown)

Thomas Jablin (UIUC)

Onur Mutlu (CMU)

Maurice Herlihy (Brown)

# Historical Trends in GFLOPS: CPUs vs. GPUs



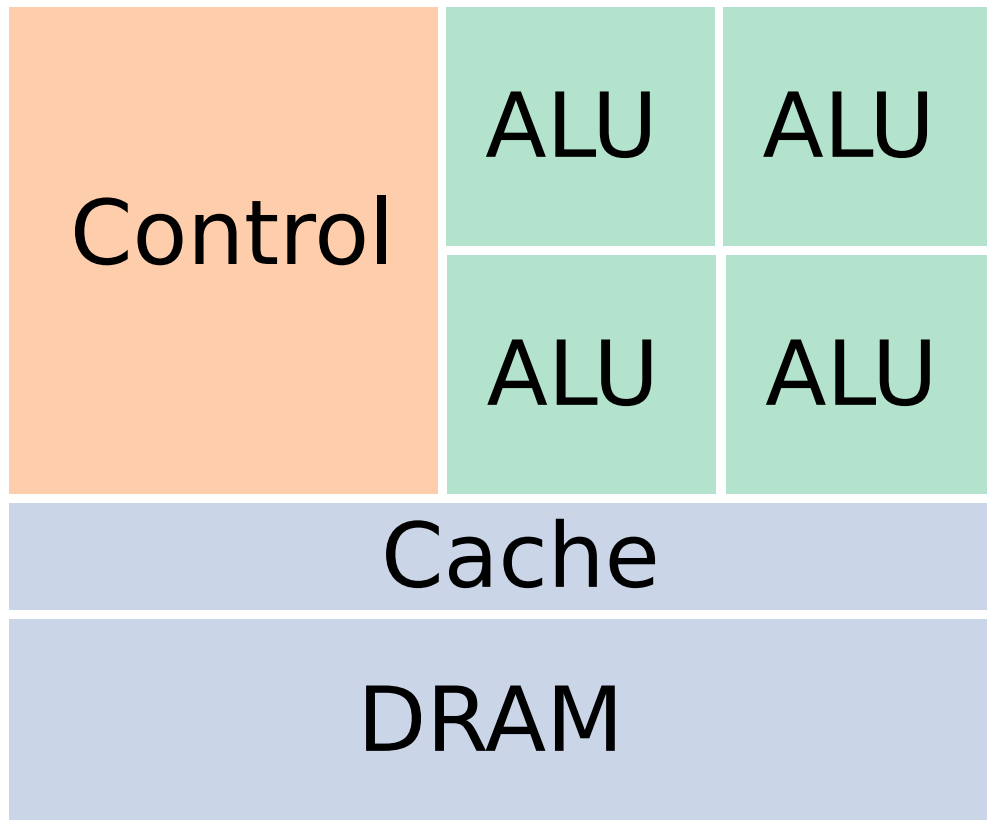
# Performance Pitfalls

Control flow can  
negatively affect performance.

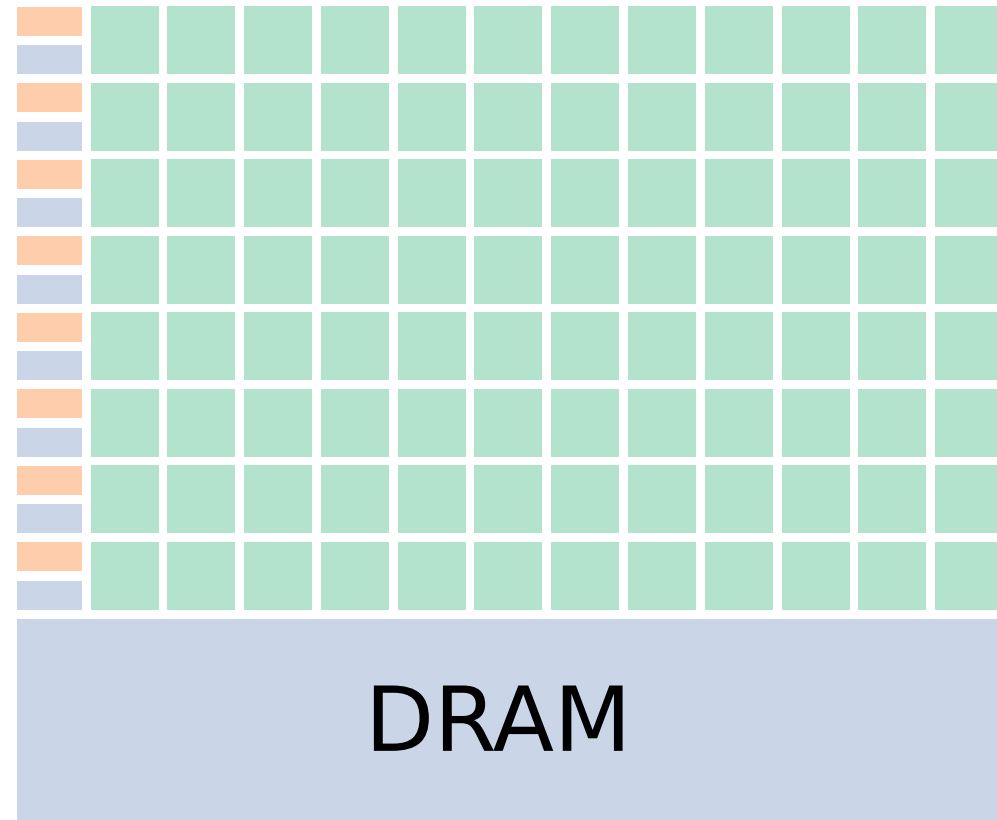
# Performance Pitfalls

Pipeline Stall - execution delay in an instruction pipeline to resolve a dependency

# Hardware: CPU versus GPU



CPU



GPU

# With Branch Prediction

# Without Branch Prediction

Clock Cycle

Clock Cycle

0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

Waiting  
Instructions

Pipeline Stages

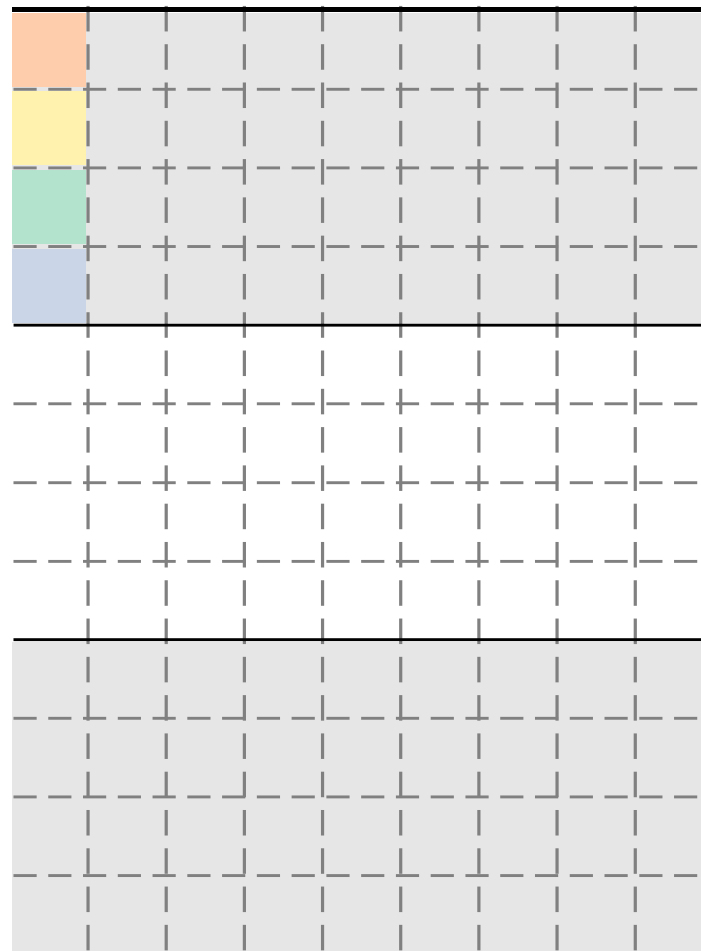
— Fetch

— Decode

— Execute

— Write

Completed  
Instructions



# With Branch Prediction

# Without Branch Prediction

Clock Cycle

Clock Cycle

0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

Waiting  
Instructions

Pipeline Stages

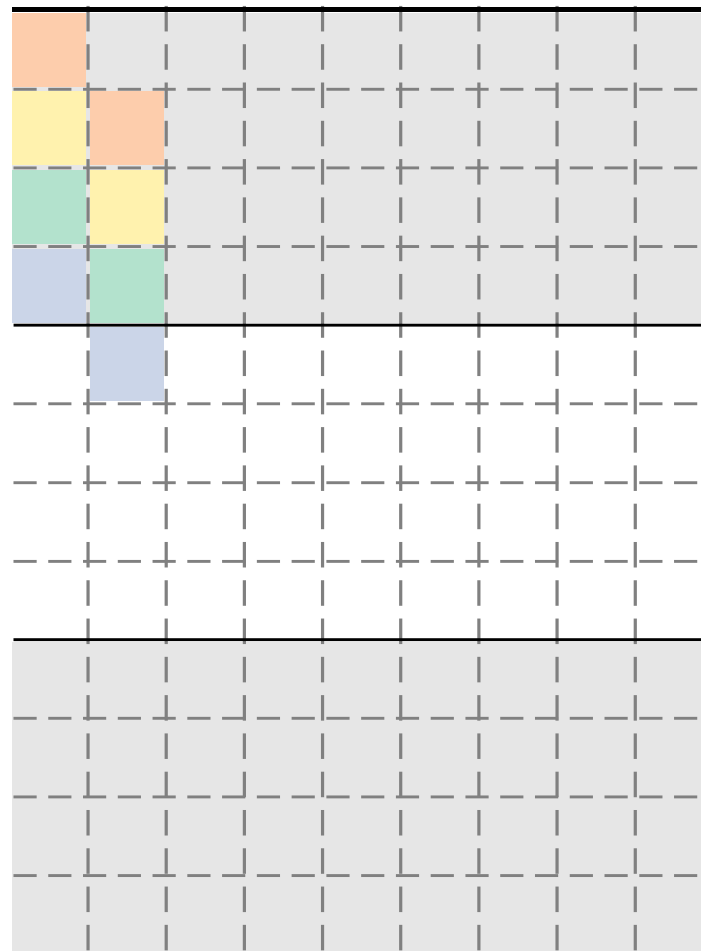
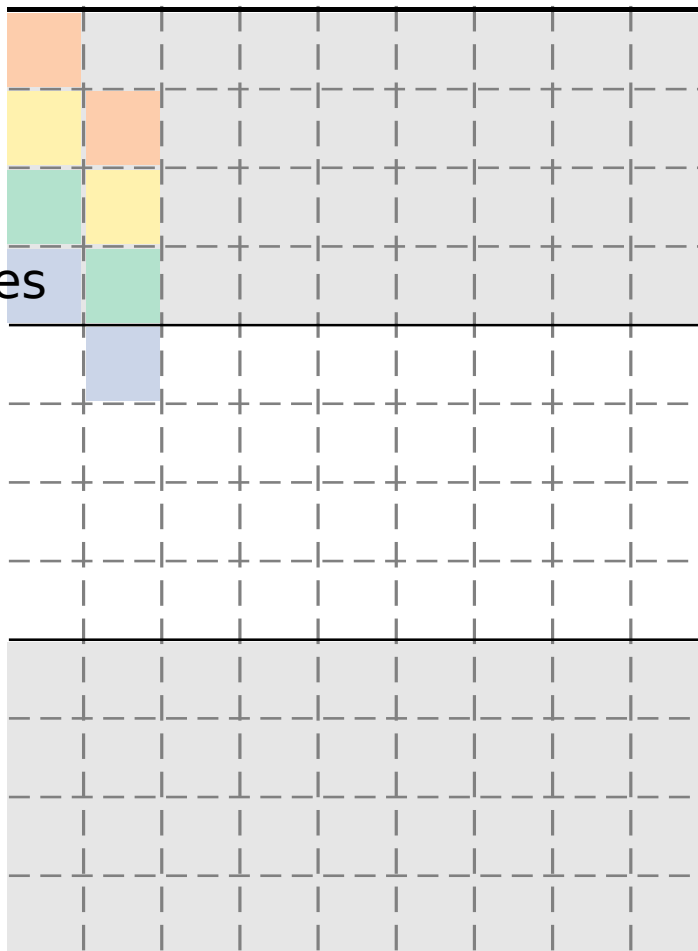
Fetch

Decode

Execute

Write

Completed  
Instructions







# With Branch Prediction

# Without Branch Prediction

Clock Cycle

Clock Cycle

0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

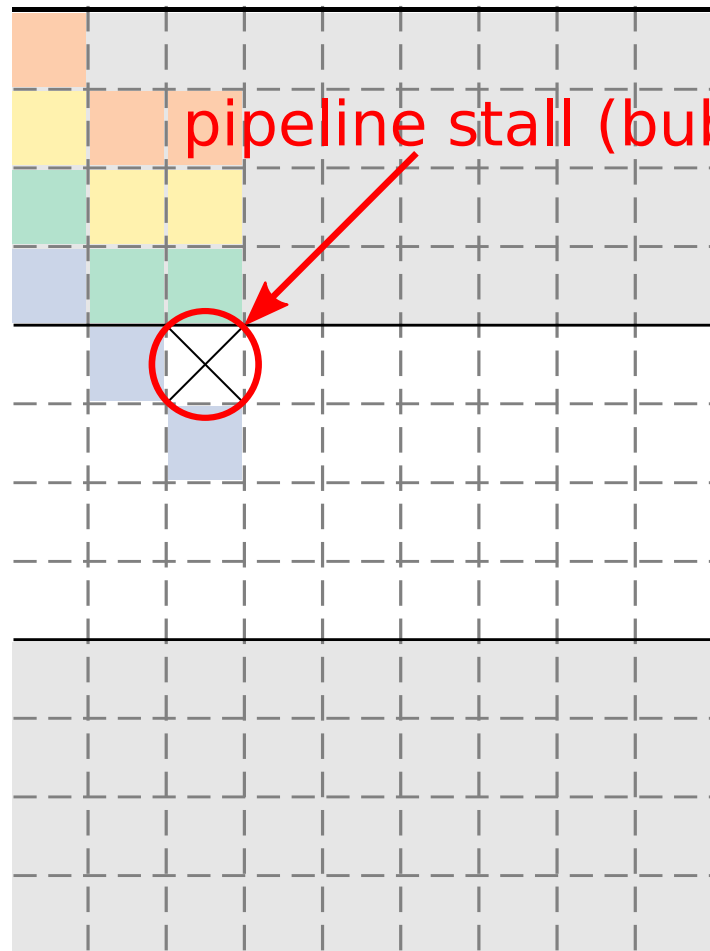
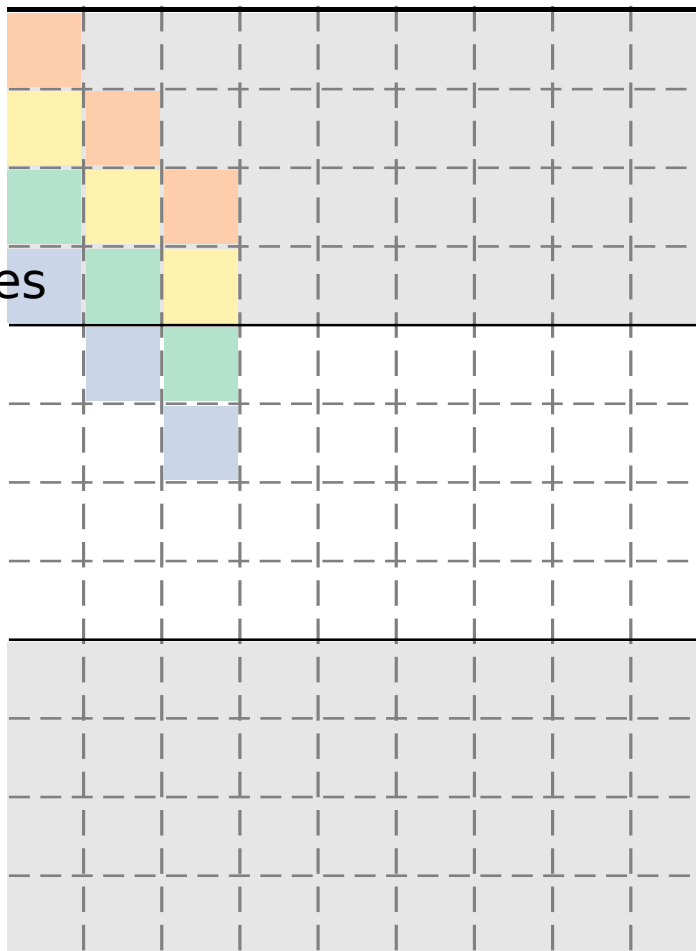
0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

Waiting  
Instructions

Pipeline Stages

- Fetch
- Decode
- Execute
- Write

Completed  
Instructions









# With Branch Prediction

# Without Branch Prediction

Clock Cycle

Clock Cycle

0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

Waiting  
Instructions

Pipeline Stages

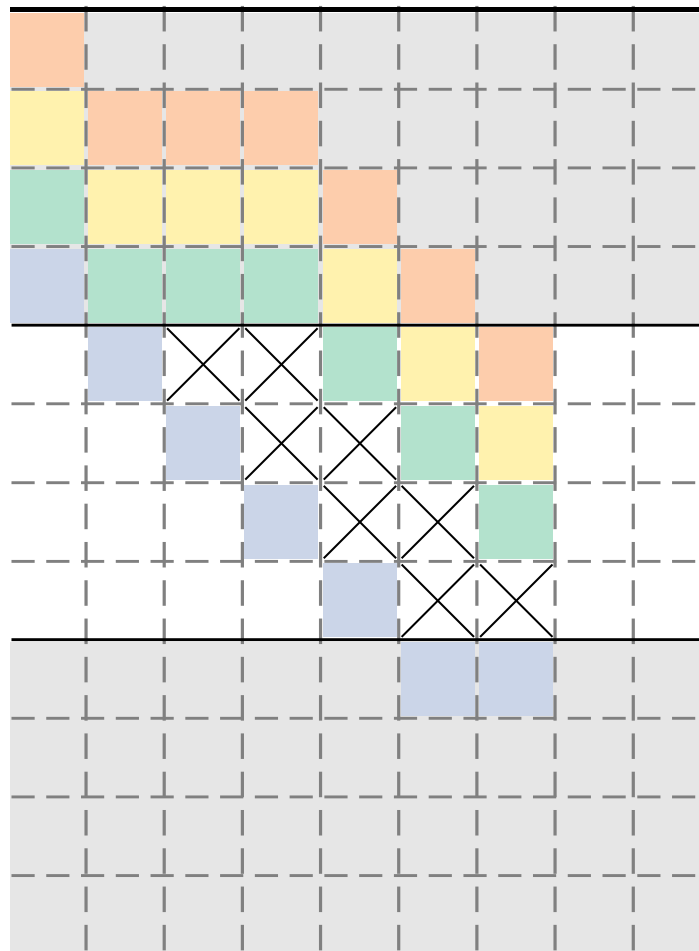
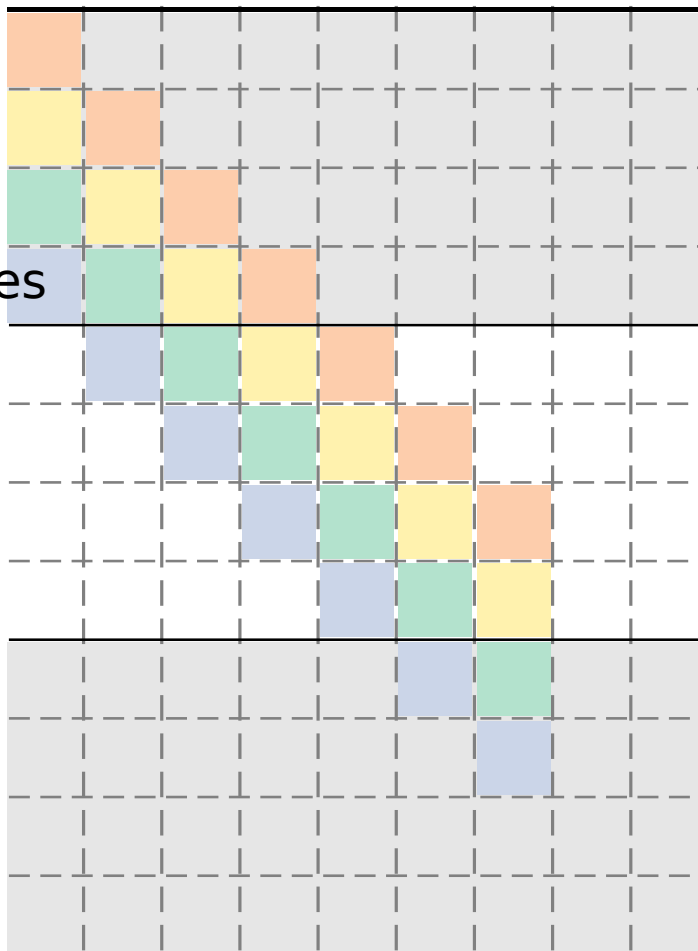
Fetch

Decode

Execute

Write

Completed  
Instructions



# With Branch Prediction

# Without Branch Prediction

Clock Cycle

Clock Cycle

0 1 2 3 4 5 6 7 8

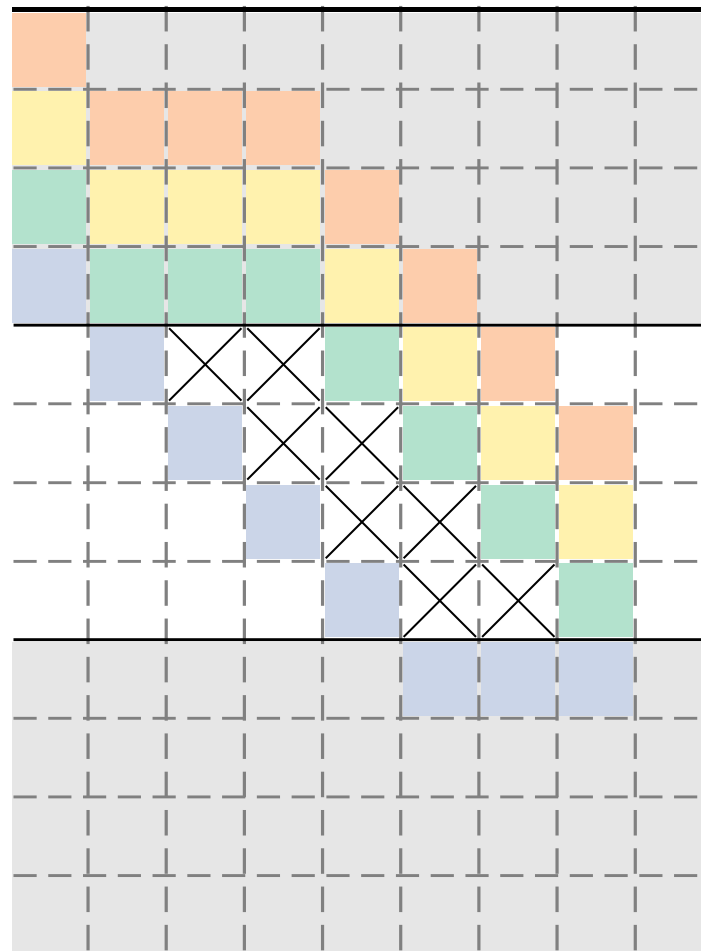
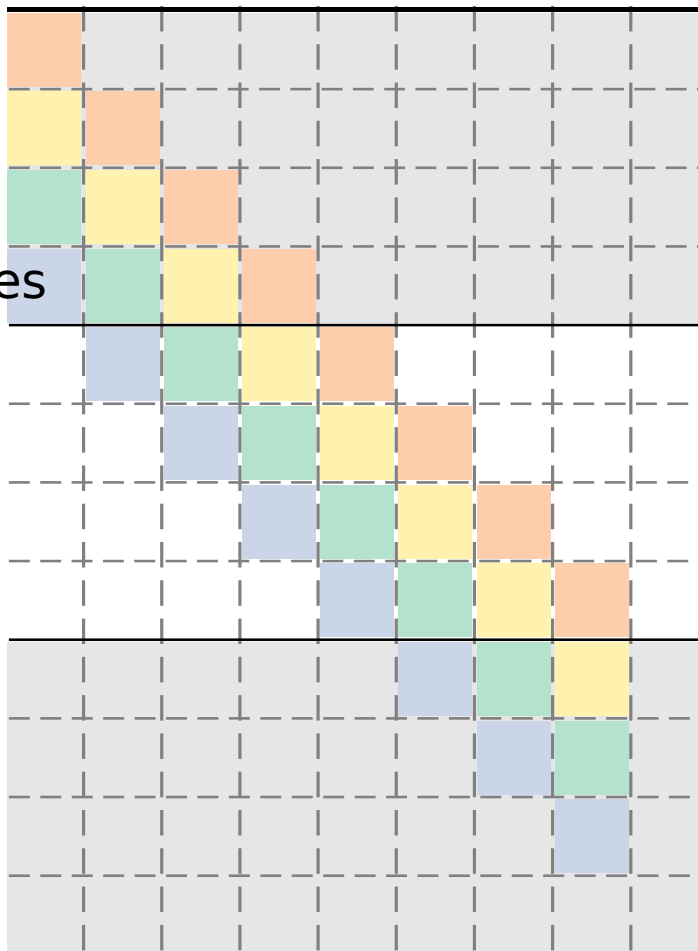
0 1 2 3 4 5 6 7 8

Waiting  
Instructions

Pipeline Stages

- Fetch
- Decode
- Execute
- Write

Completed  
Instructions



# With Branch Prediction

# Without Branch Prediction

Clock Cycle

Clock Cycle

0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

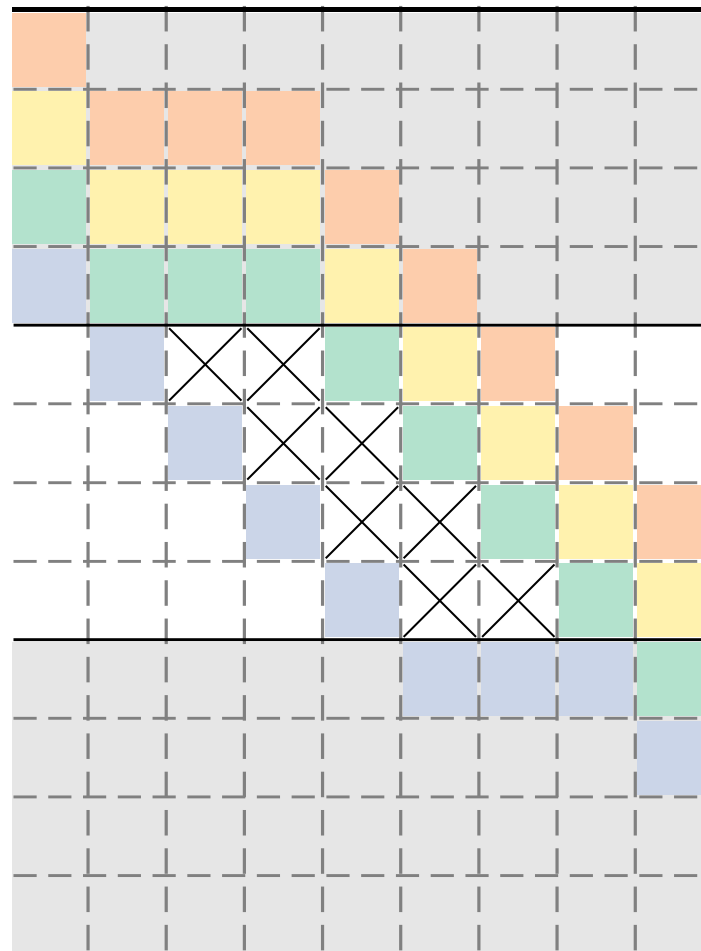
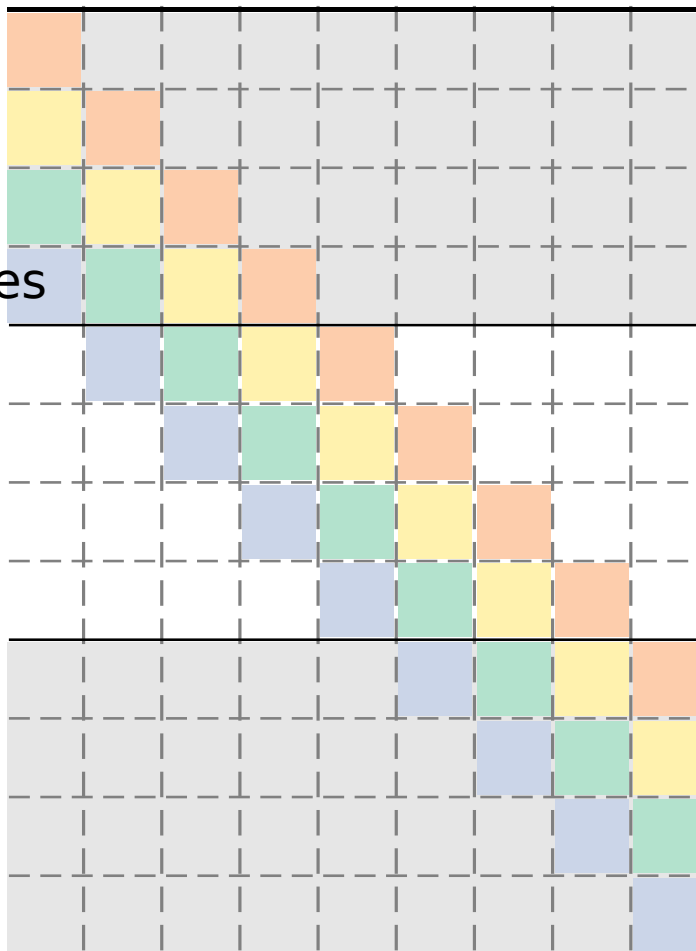
0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

Waiting  
Instructions

Pipeline Stages

- Fetch
- Decode
- Execute
- Write

Completed  
Instructions



# With Branch Prediction

# Without Branch Prediction

Clock Cycle

Clock Cycle

0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8

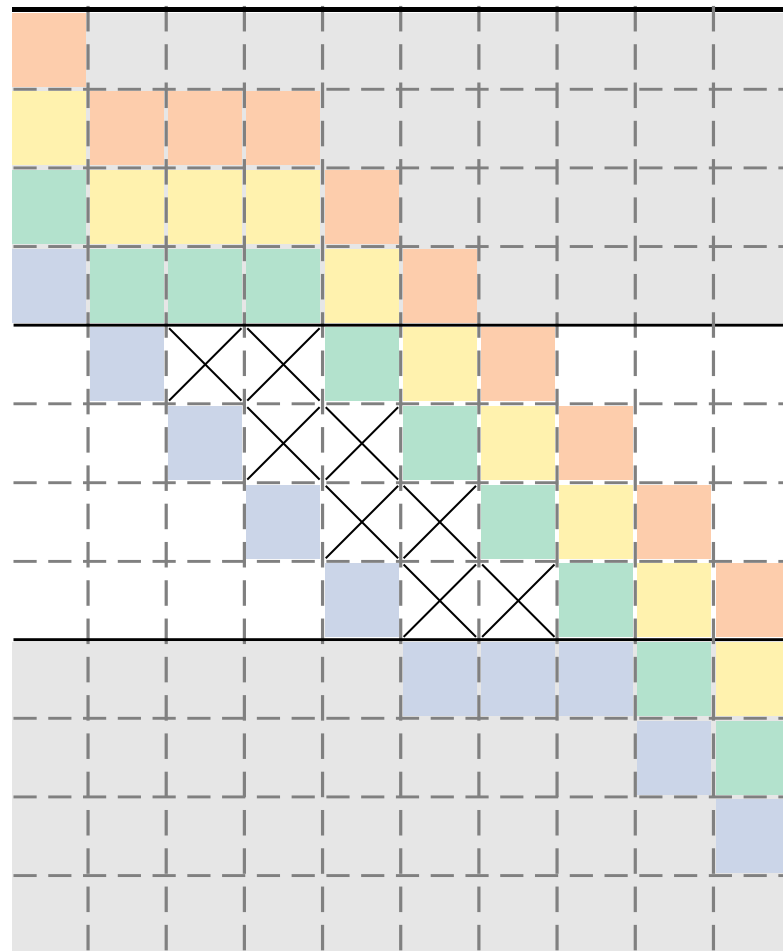
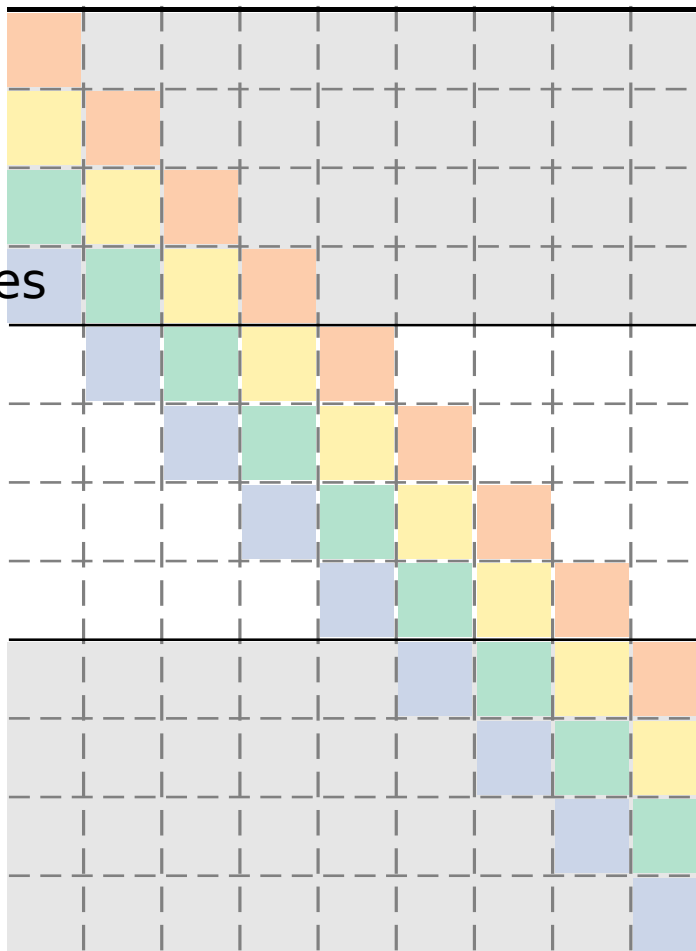
0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

Waiting Instructions

Pipeline Stages

- Fetch
- Decode
- Execute
- Write

Completed Instructions







# Performance Pitfalls

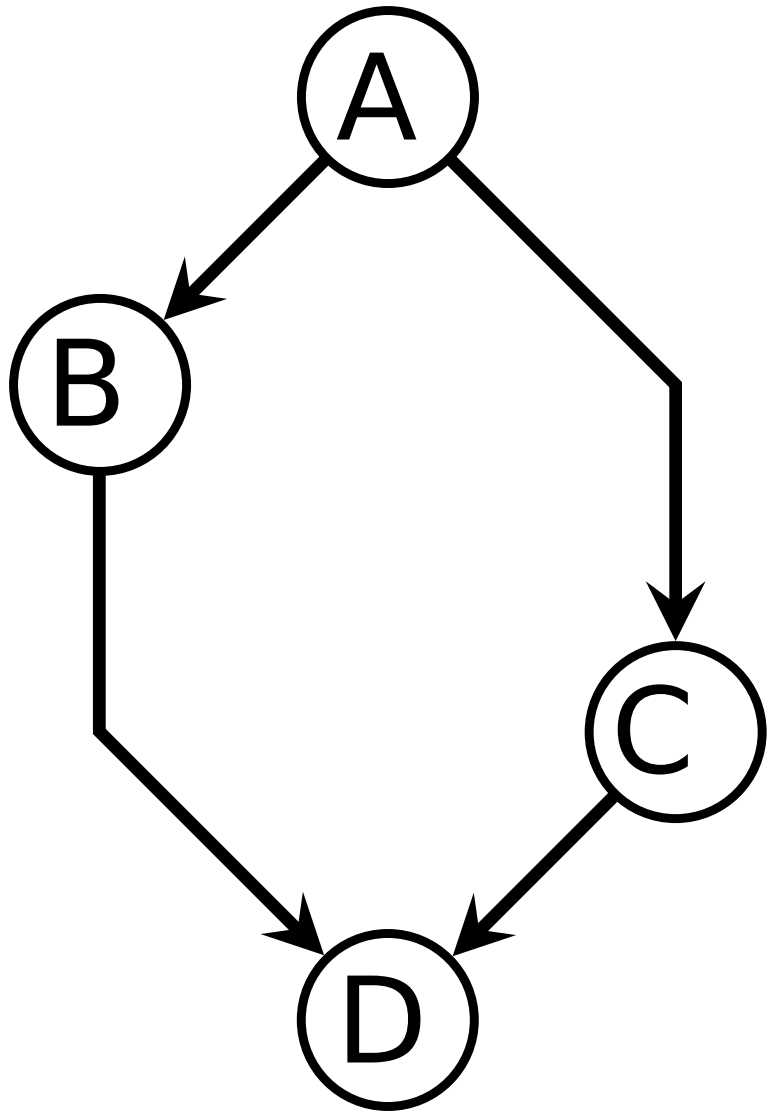
Pipeline Stall - execution delay in an instruction pipeline to resolve a dependency

# Performance Pitfalls

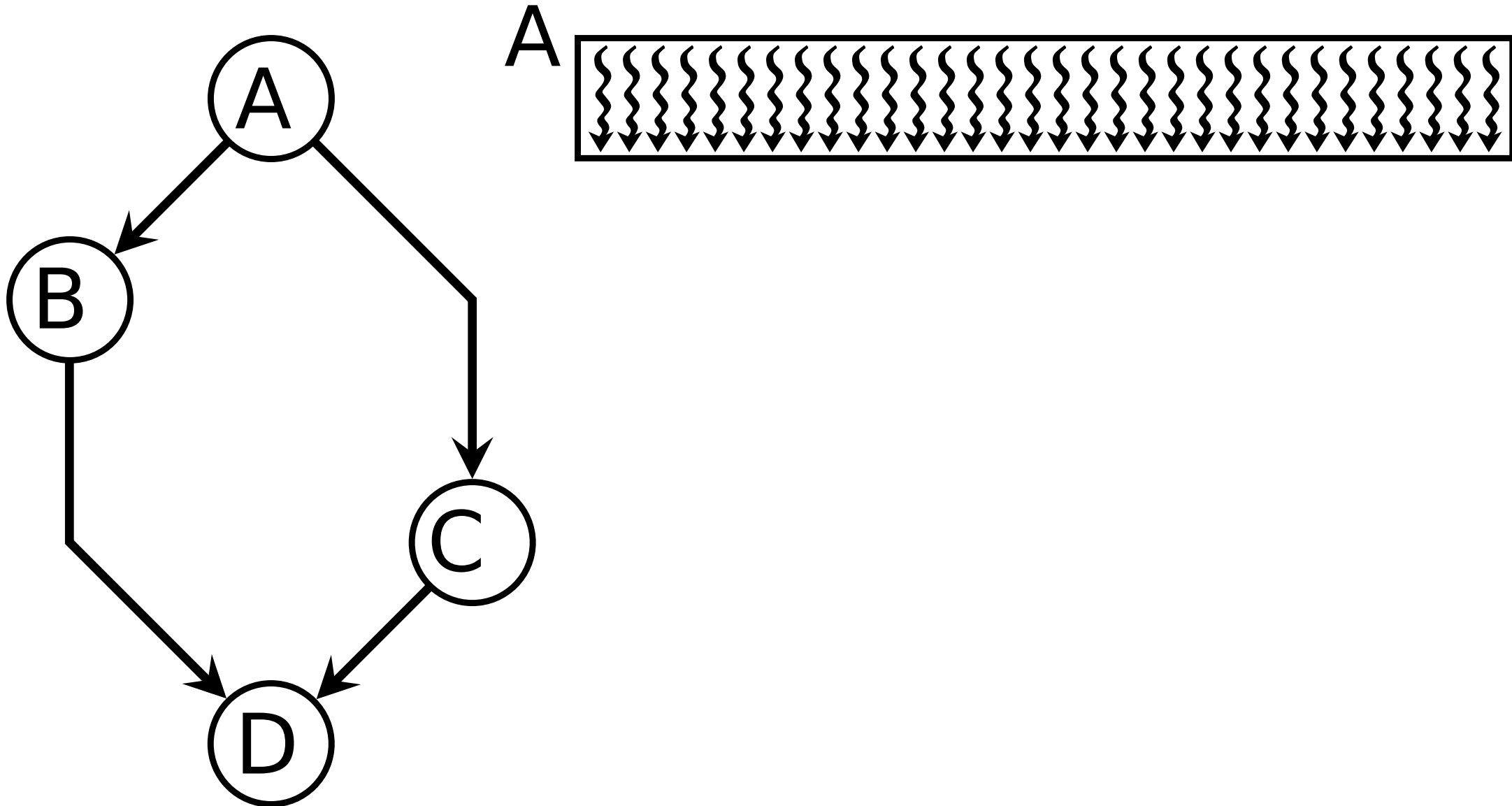
Pipeline Stall - execution delay in an instruction pipeline to resolve a dependency

Warp Divergence - threads within a warp take different paths and the different execution paths are serialized

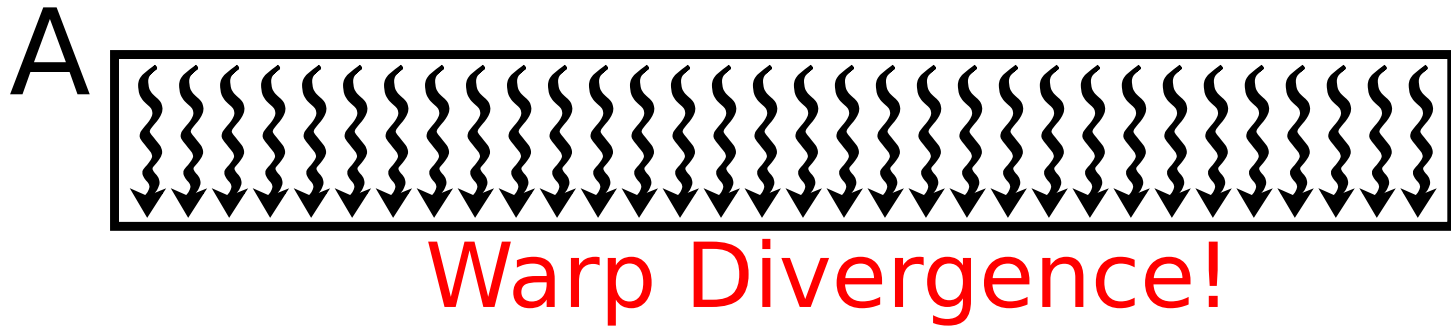
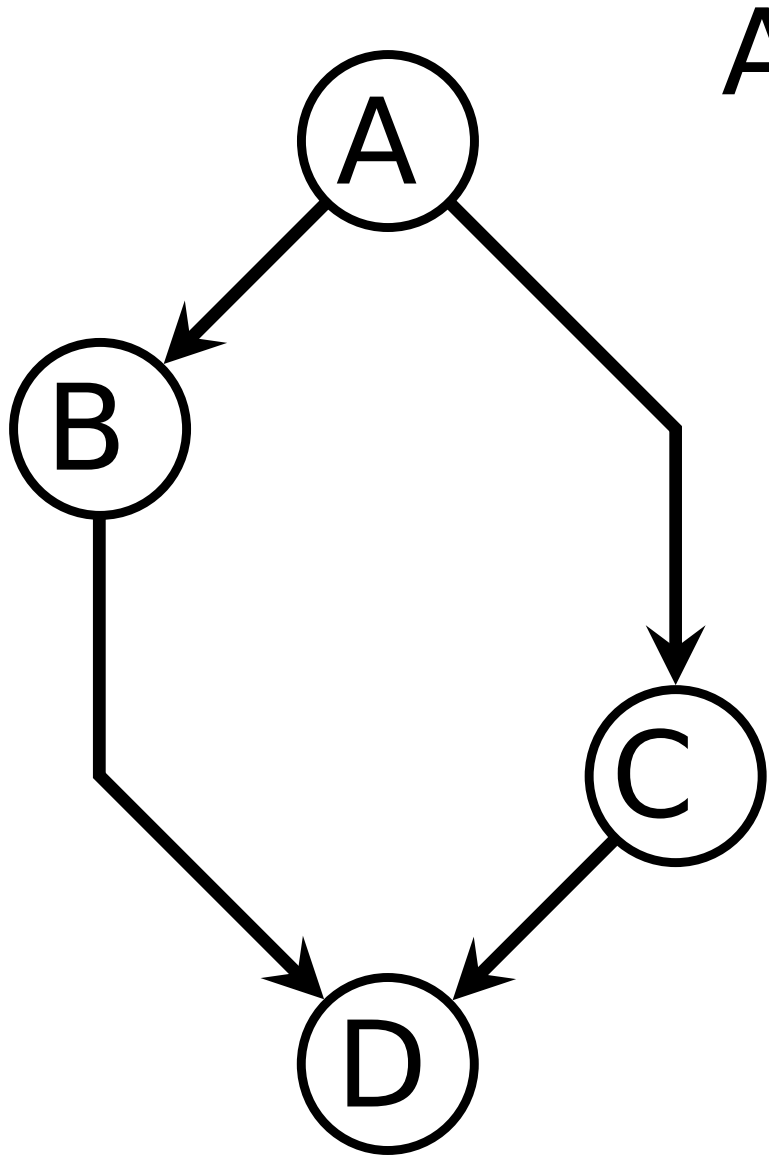
# Warp Divergence Example



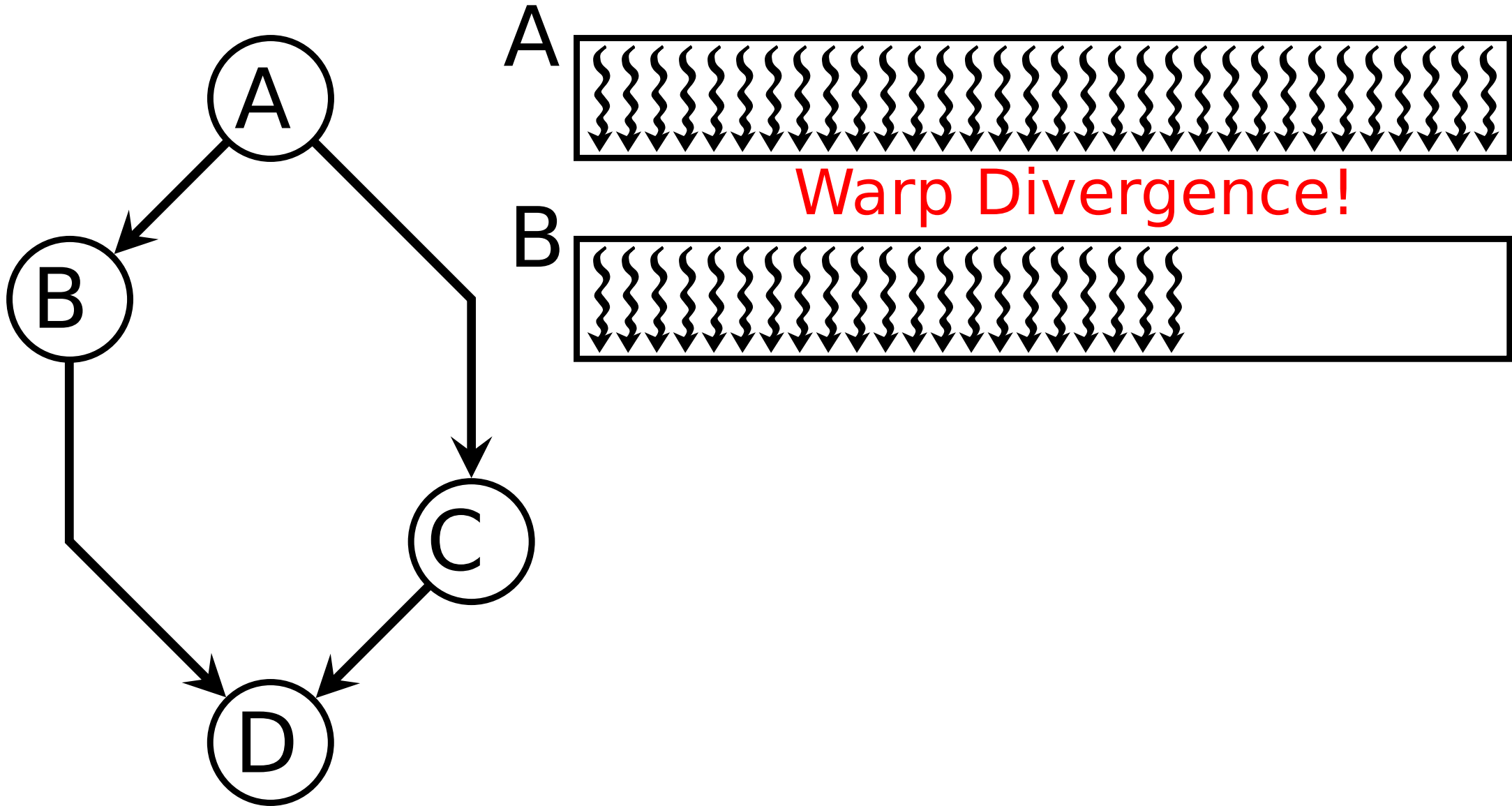
# Warp Divergence Example



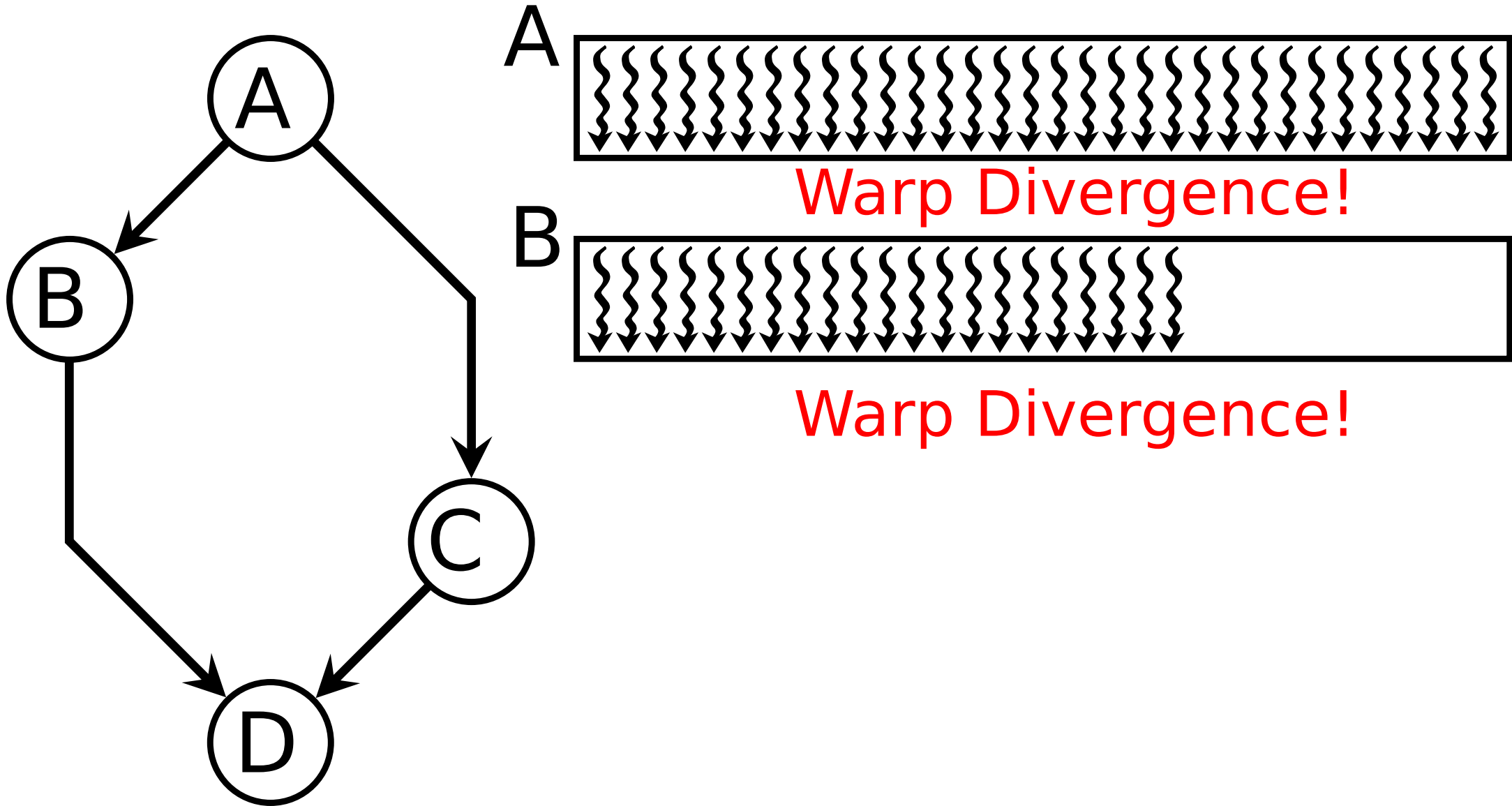
# Warp Divergence Example



# Warp Divergence Example

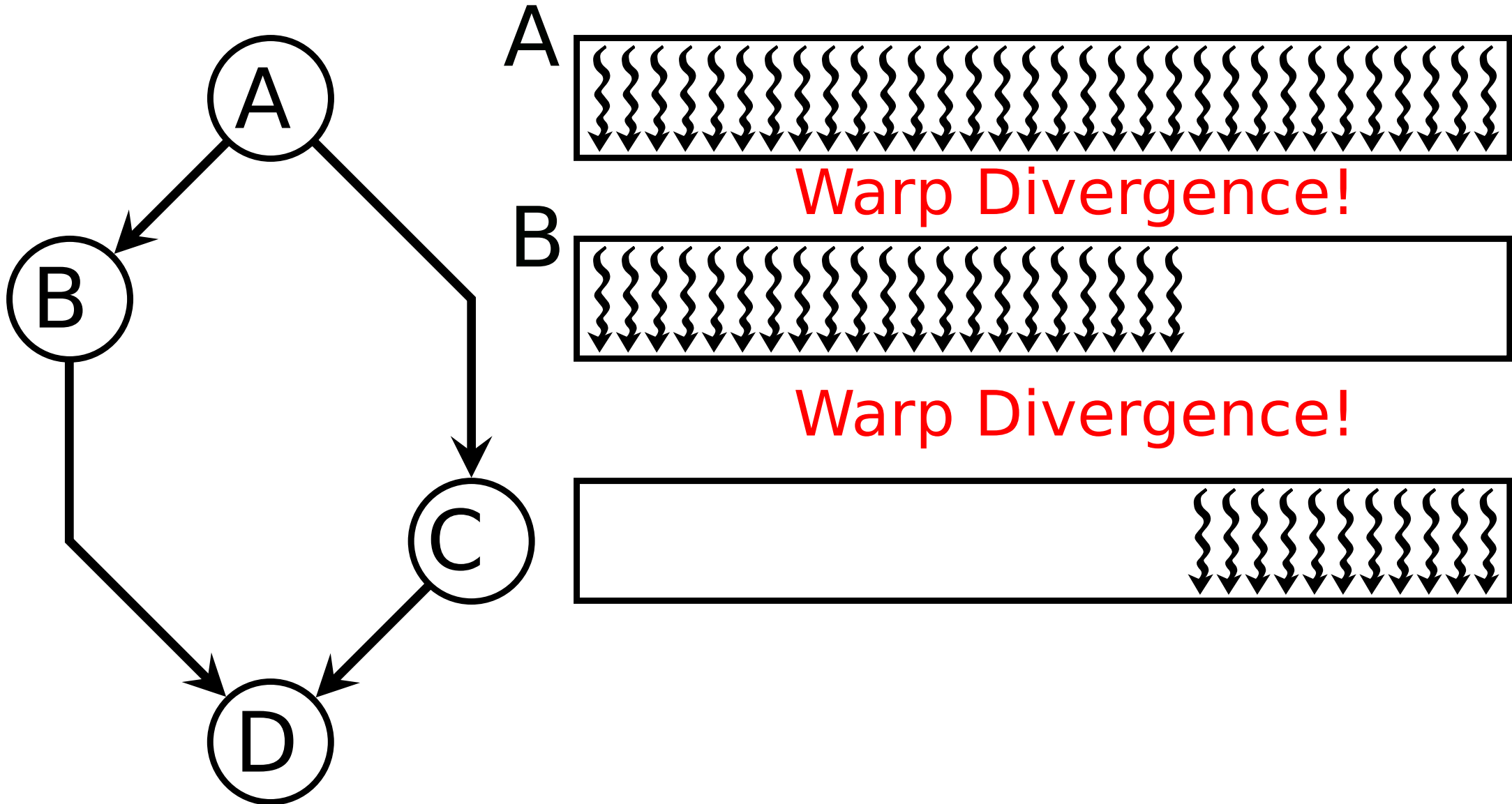


# Warp Divergence Example

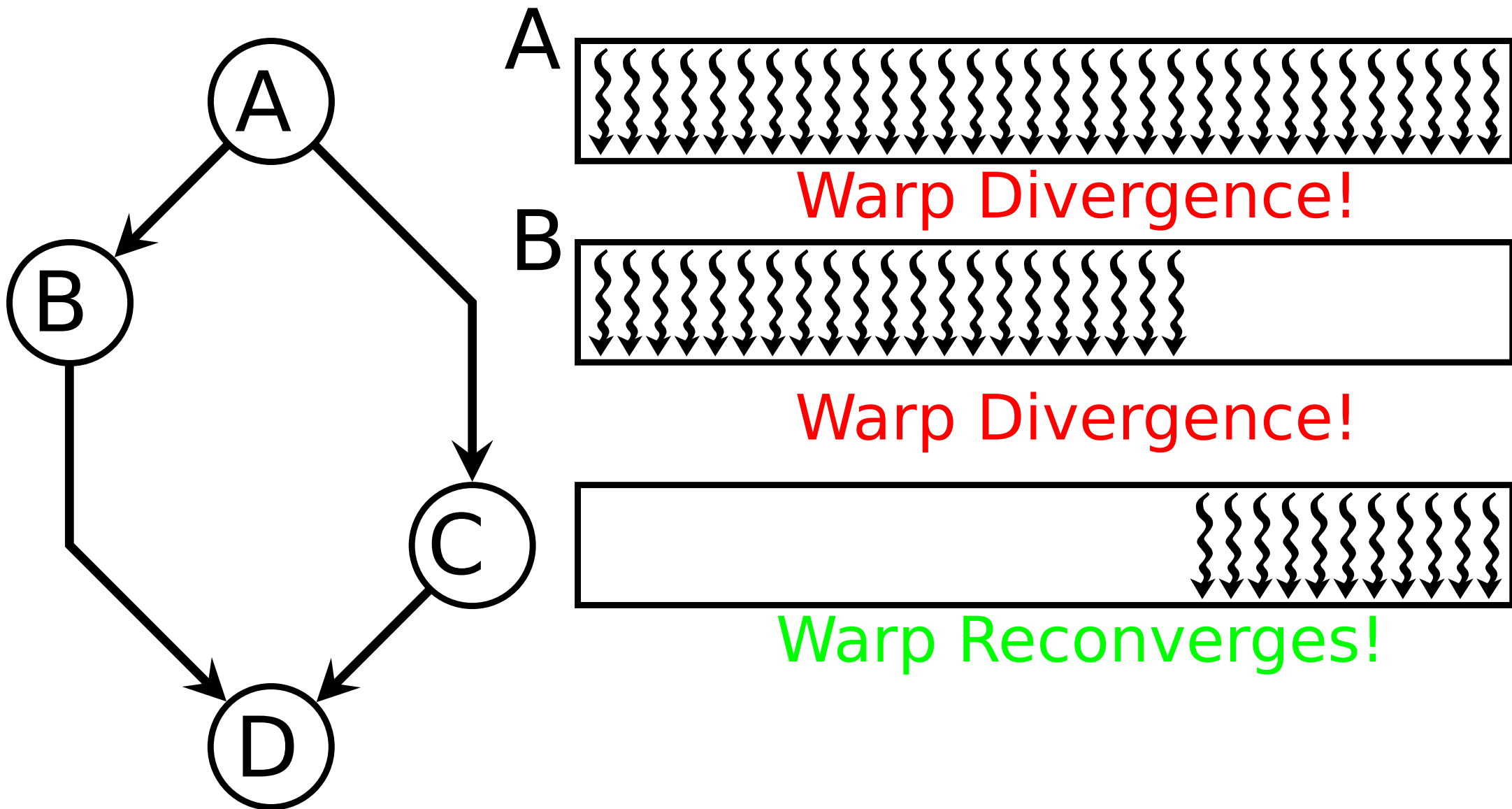




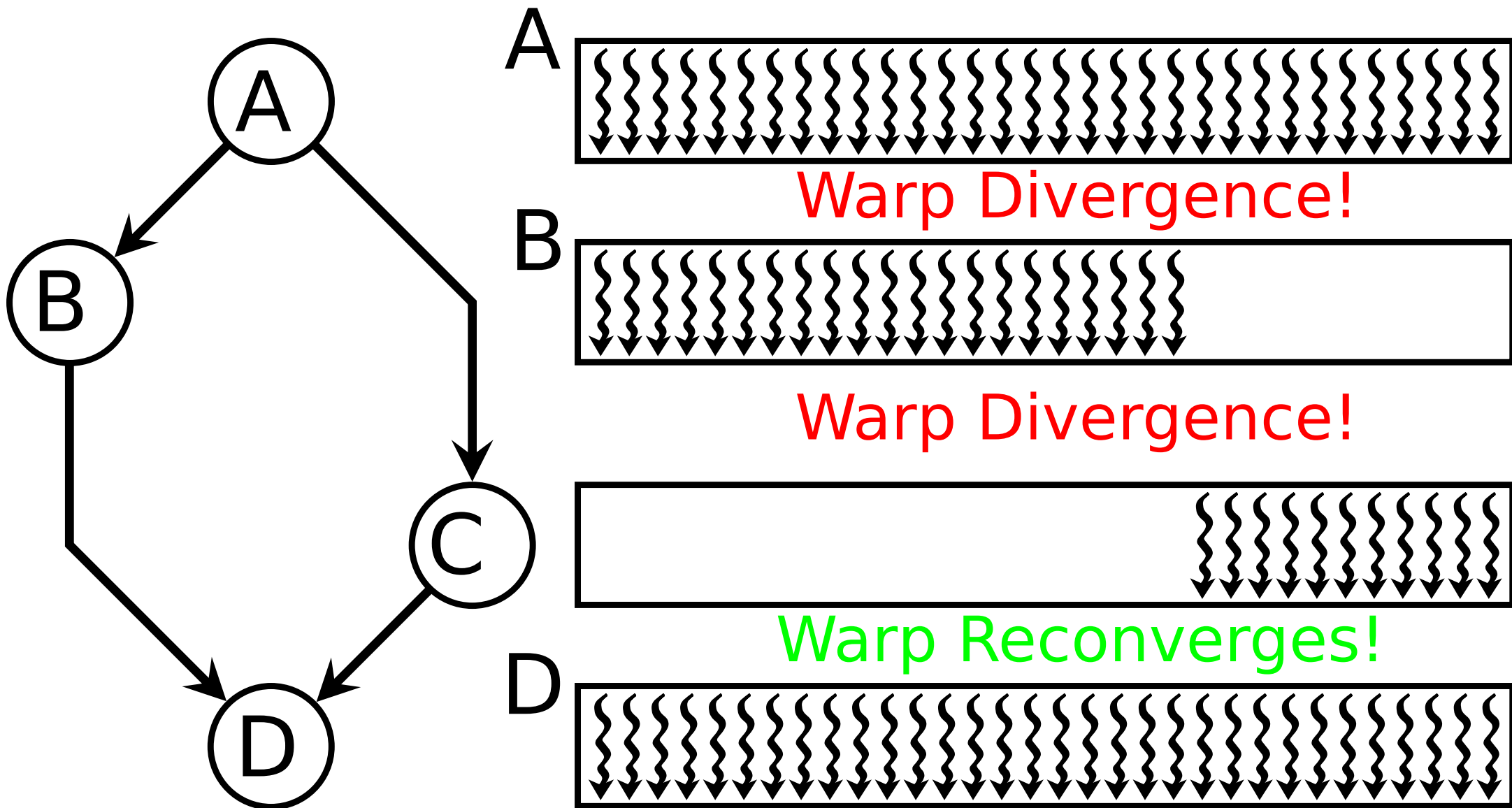
# Warp Divergence Example



# Warp Divergence Example



# Warp Divergence Example



# Warp-Aware Trace Scheduling

Schedule instructions across basic block boundaries to expose additional ILP...

# Warp-Aware Trace Scheduling

Schedule instructions across basic block boundaries to expose additional ILP...

while managing and  
optimizing warp divergence.

# Origins: Microcode Trace Scheduling

...generalizing local and disparate vertical-to-horizontal microcode compaction

Step	Description

# Origins: Microcode Trace Scheduling

...generalizing local and disparate vertical-to-horizontal microcode compaction

Step	Description
1. Trace Selection	

# Origins: Microcode Trace Scheduling

...generalizing local and disparate vertical-to-horizontal microcode compaction

Step	Description
1. Trace Selection	
2. Trace Formation	



# Origins: Microcode Trace Scheduling

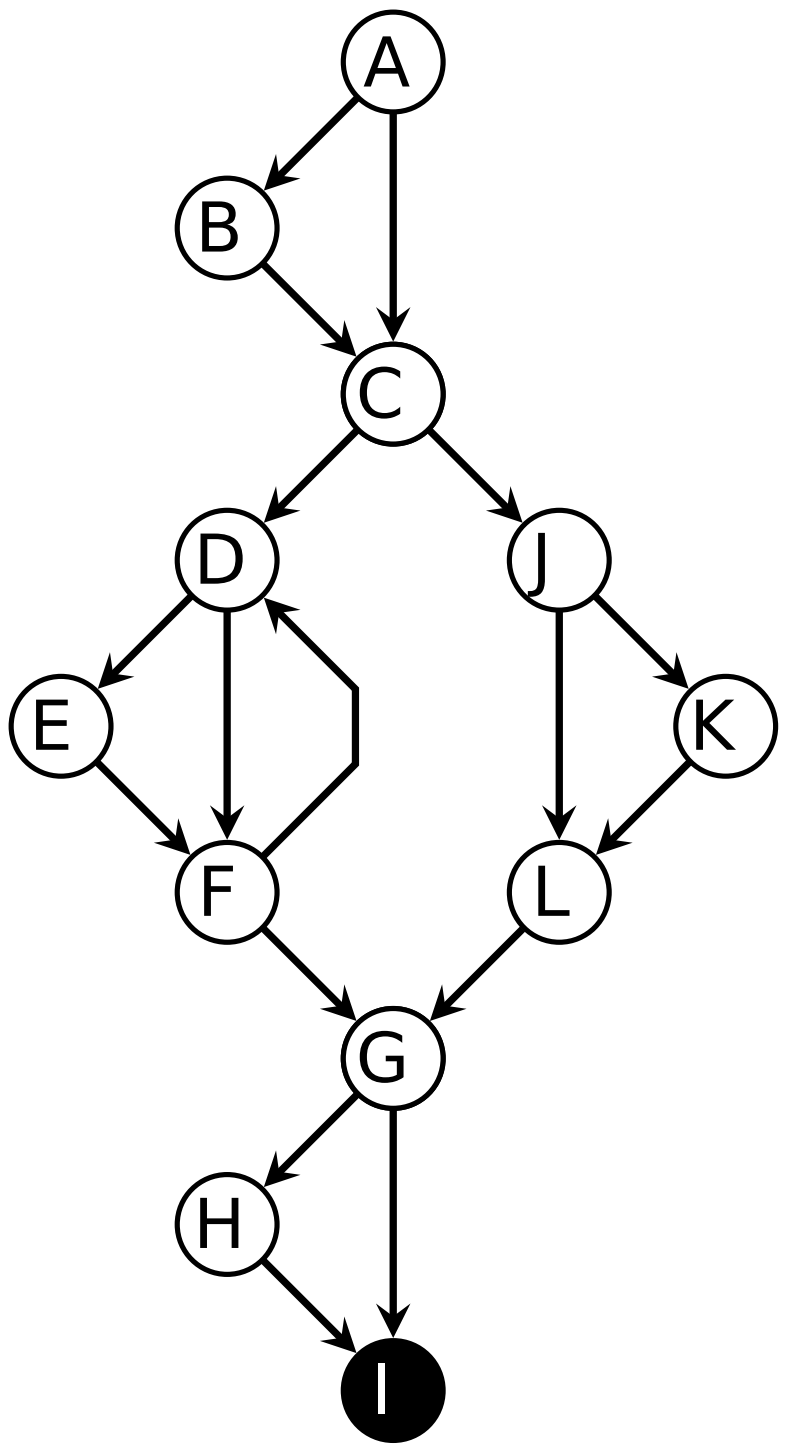
...generalizing local and disparate vertical-to-horizontal microcode compaction

Step	Description
1. Trace Selection	
2. Trace Formation	
3. Local Scheduling	

# Origins: Microcode Trace Scheduling

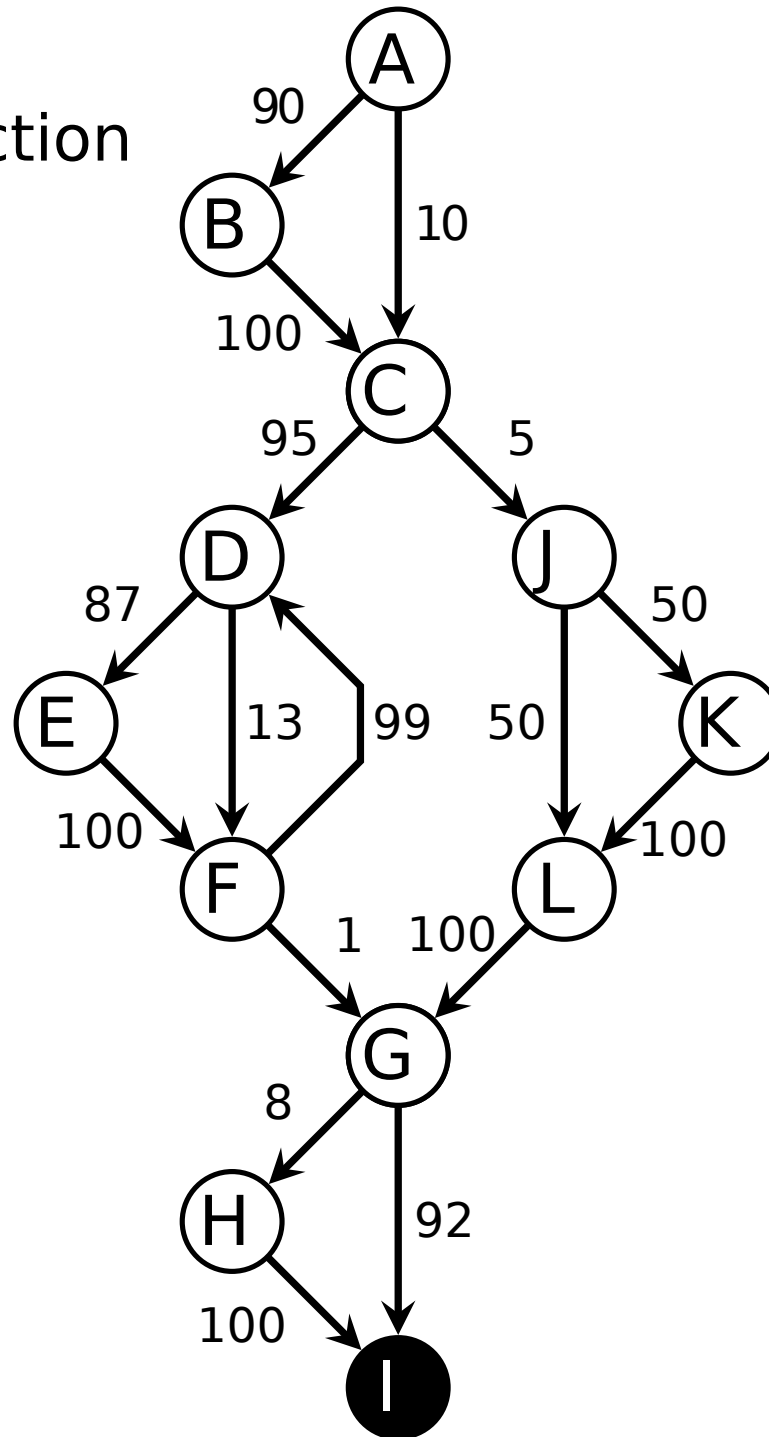
...generalizing local and disparate vertical-to-horizontal microcode compaction

Step	Description
1. Trace Selection	partition basic blocks into regions
2. Trace Formation	facilitate local scheduling, potentially adding nodes and edges
3. Local Scheduling	schedule instructions within each region



# Annotate CFG

- dynamic profiling
- static branch prediction



# Annotate CFG

- dynamic profiling
- static branch prediction

while there are unvisited nodes

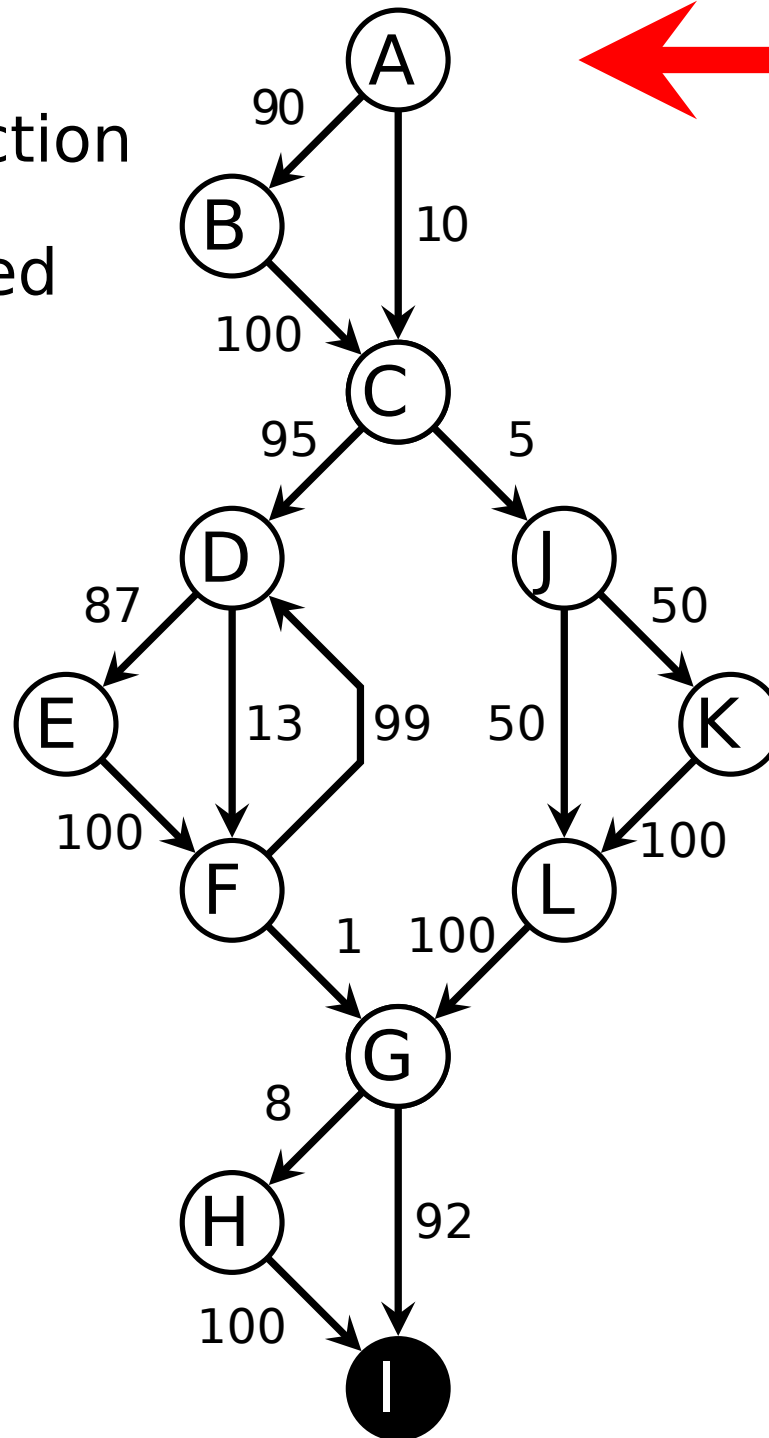
loop

Find the next unvisited node, with highest edge weight

Add node to trace

end loop

end while



# Annotate CFG

- dynamic profiling
- static branch prediction

while there are unvisited nodes

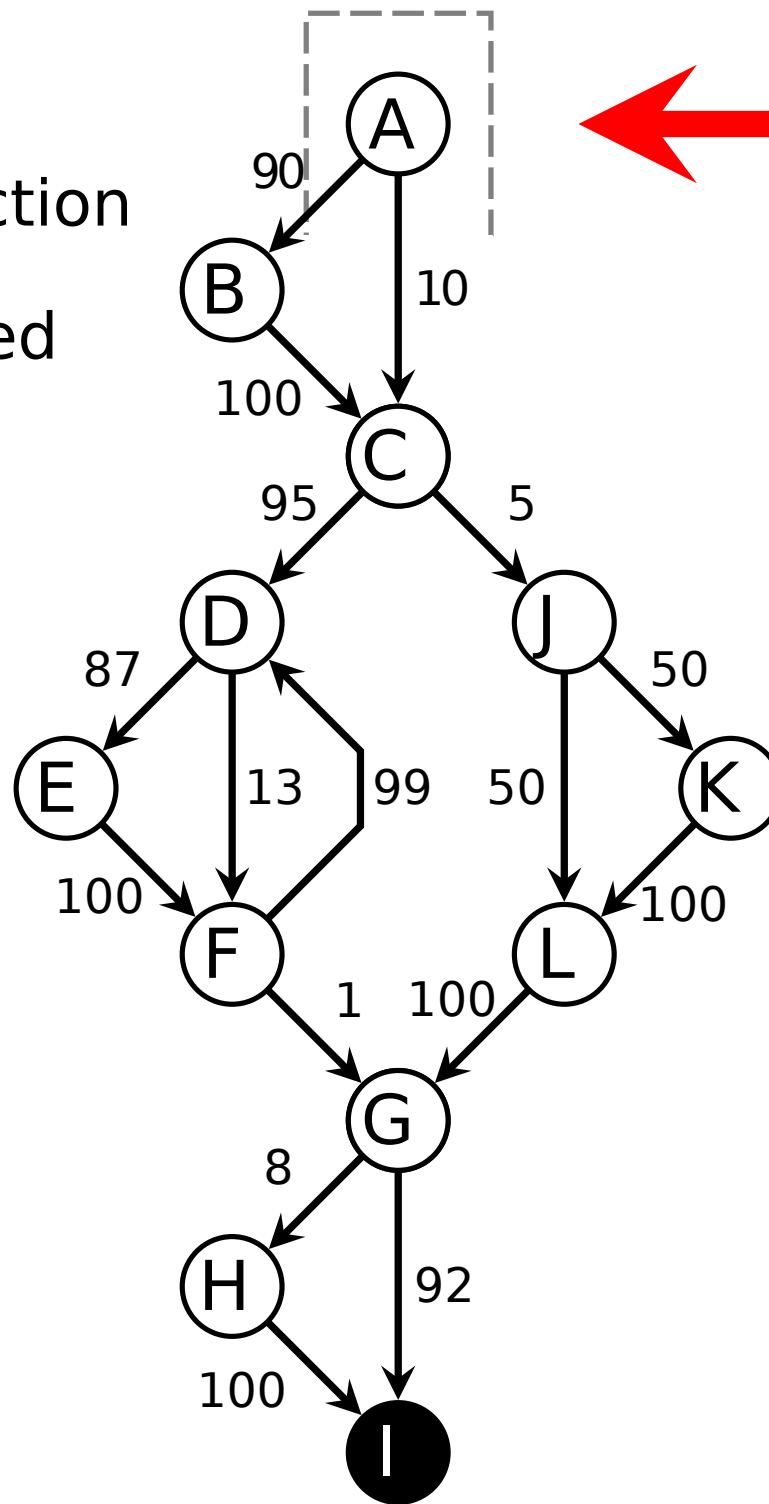
loop

Find the next unvisited node, with highest edge weight

Add node to trace

end loop

end while



# Annotate CFG

- dynamic profiling
- static branch prediction

while there are unvisited nodes

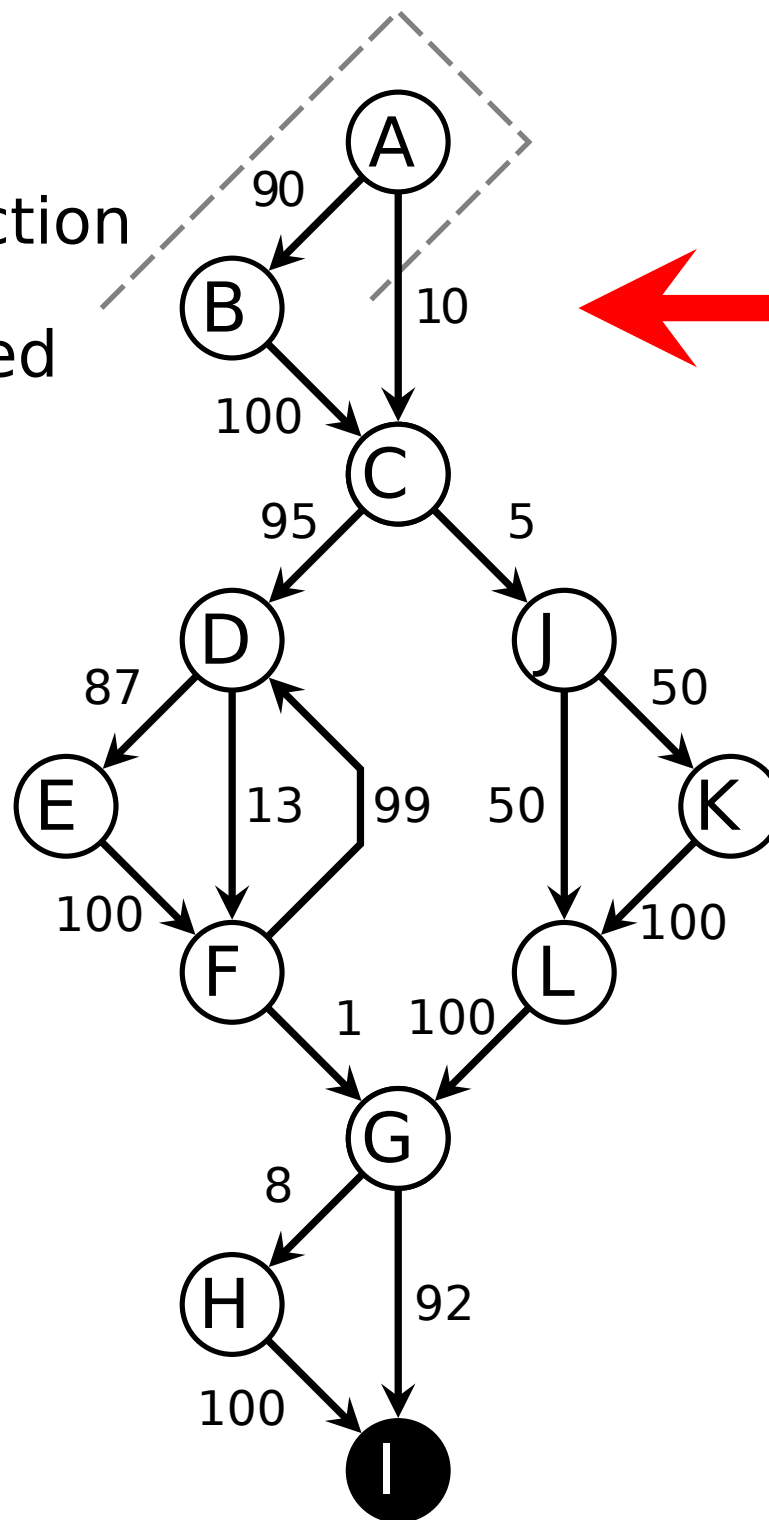
loop

Find the next unvisited node, with highest edge weight

Add node to trace

end loop

end while



# Annotate CFG

- dynamic profiling
- static branch prediction

while there are unvisited nodes

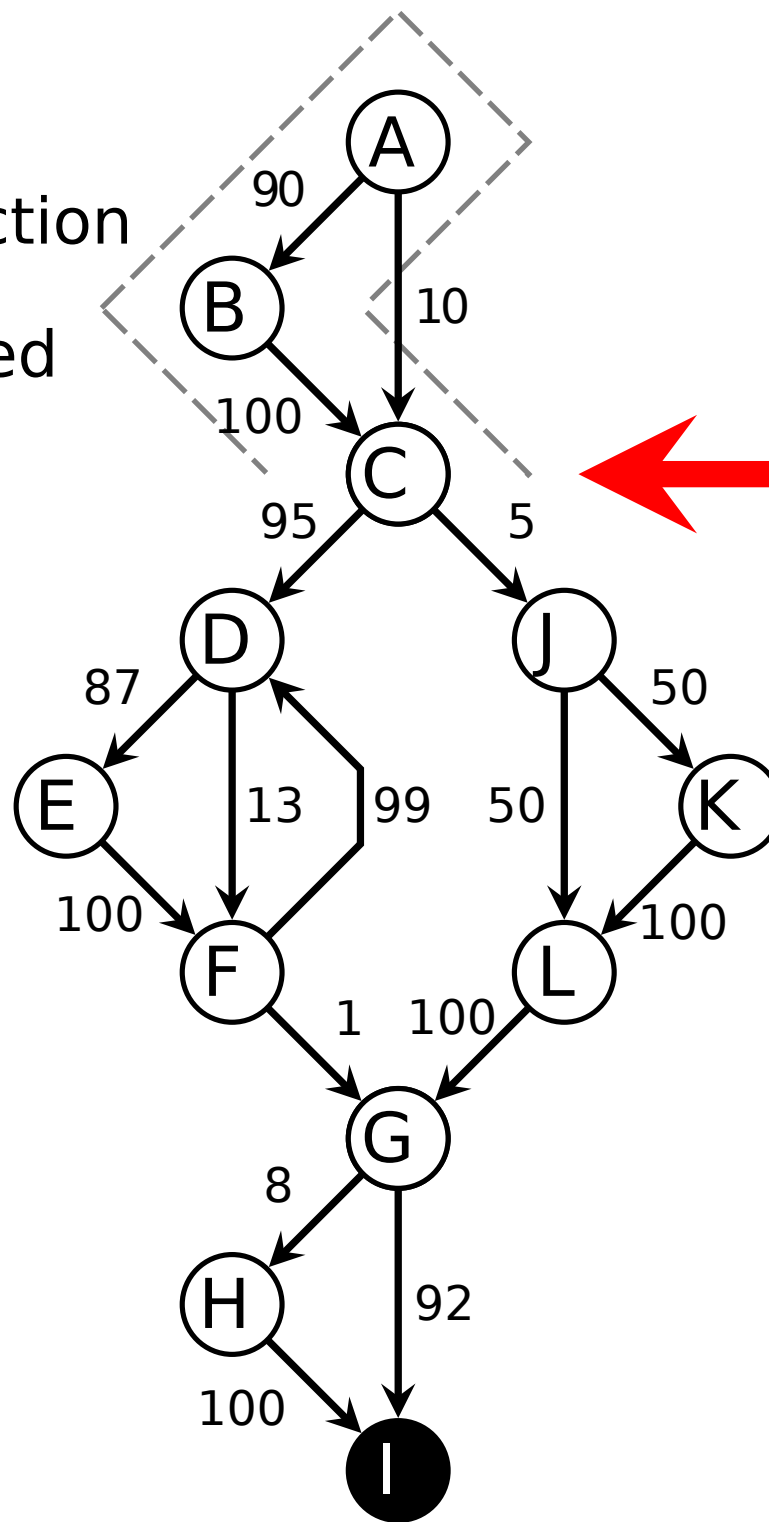
loop

Find the next unvisited node, with highest edge weight

Add node to trace

end loop

end while





# Annotate CFG

- dynamic profiling
- static branch prediction

while there are unvisited nodes

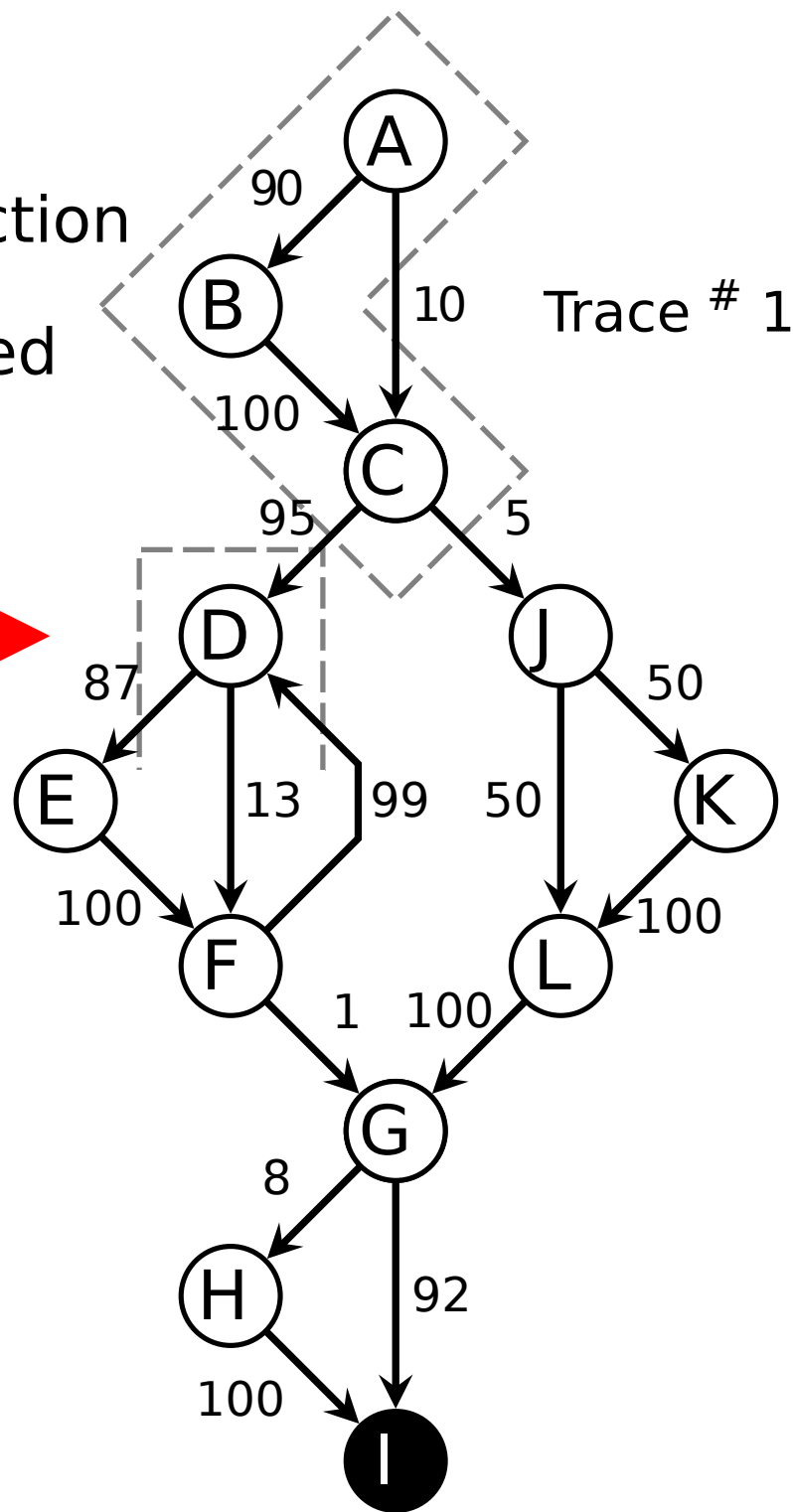
loop

Find the next unvisited node, with highest edge weight

Add node to trace

end loop

end while



# Annotate CFG

- dynamic profiling
- static branch prediction

while there are unvisited nodes

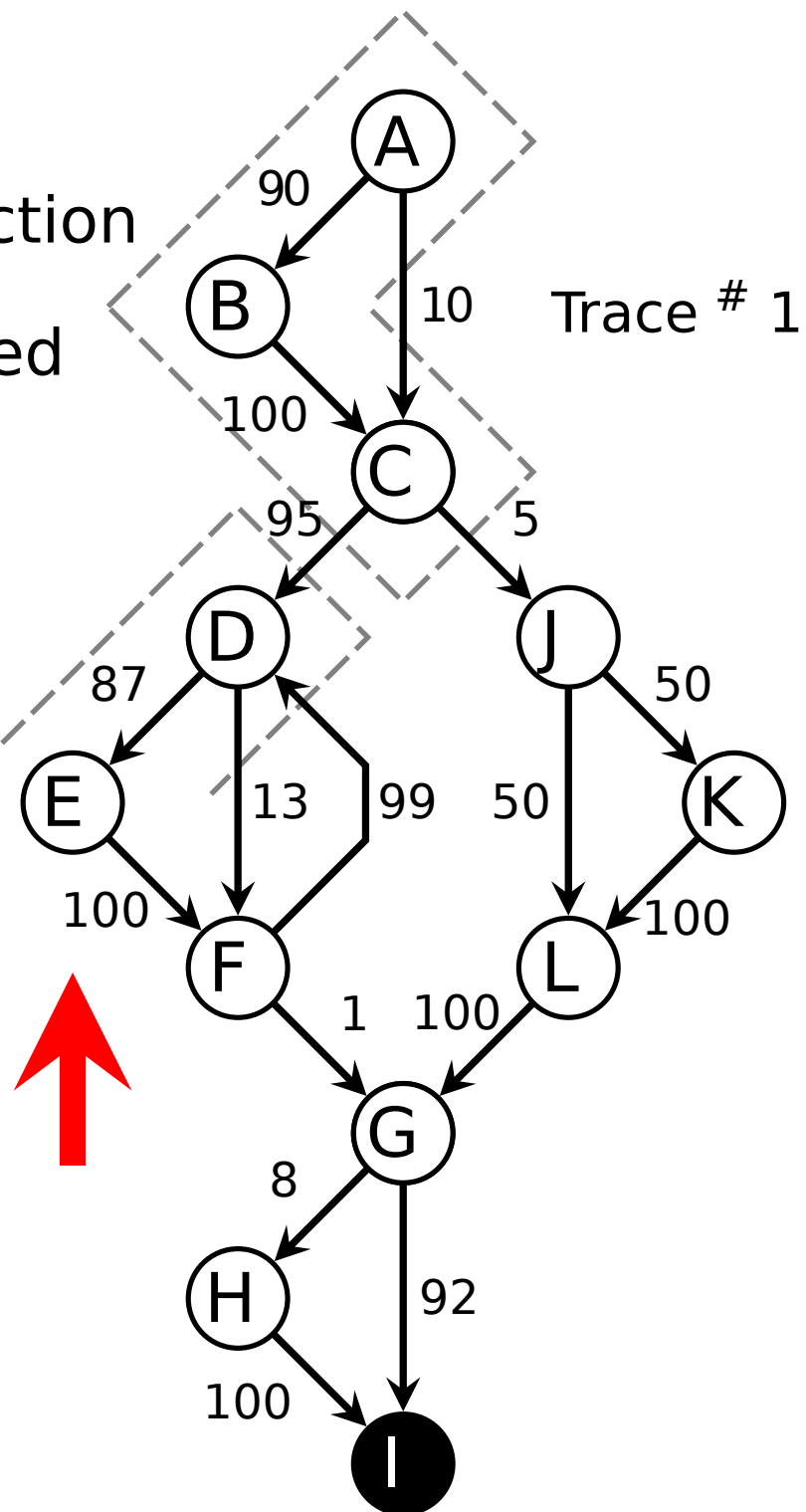
loop

Find the next unvisited node, with highest edge weight

Add node to trace

end loop

end while



# Annotate CFG

- dynamic profiling
- static branch prediction

while there are unvisited nodes

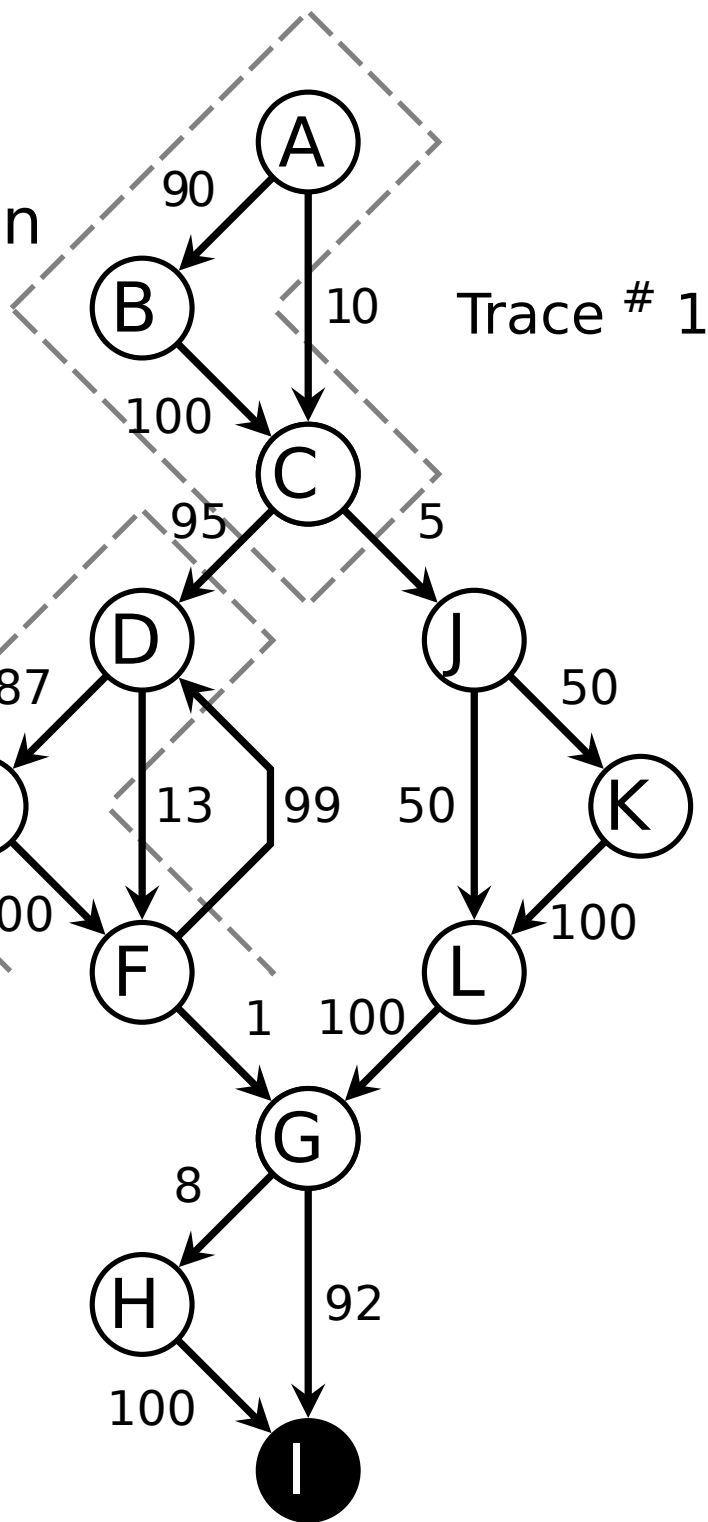
loop

Find the next unvisited node, with highest edge weight

Add node to trace

end loop

end while



# Annotate CFG

- dynamic profiling
- static branch prediction

while there are unvisited nodes

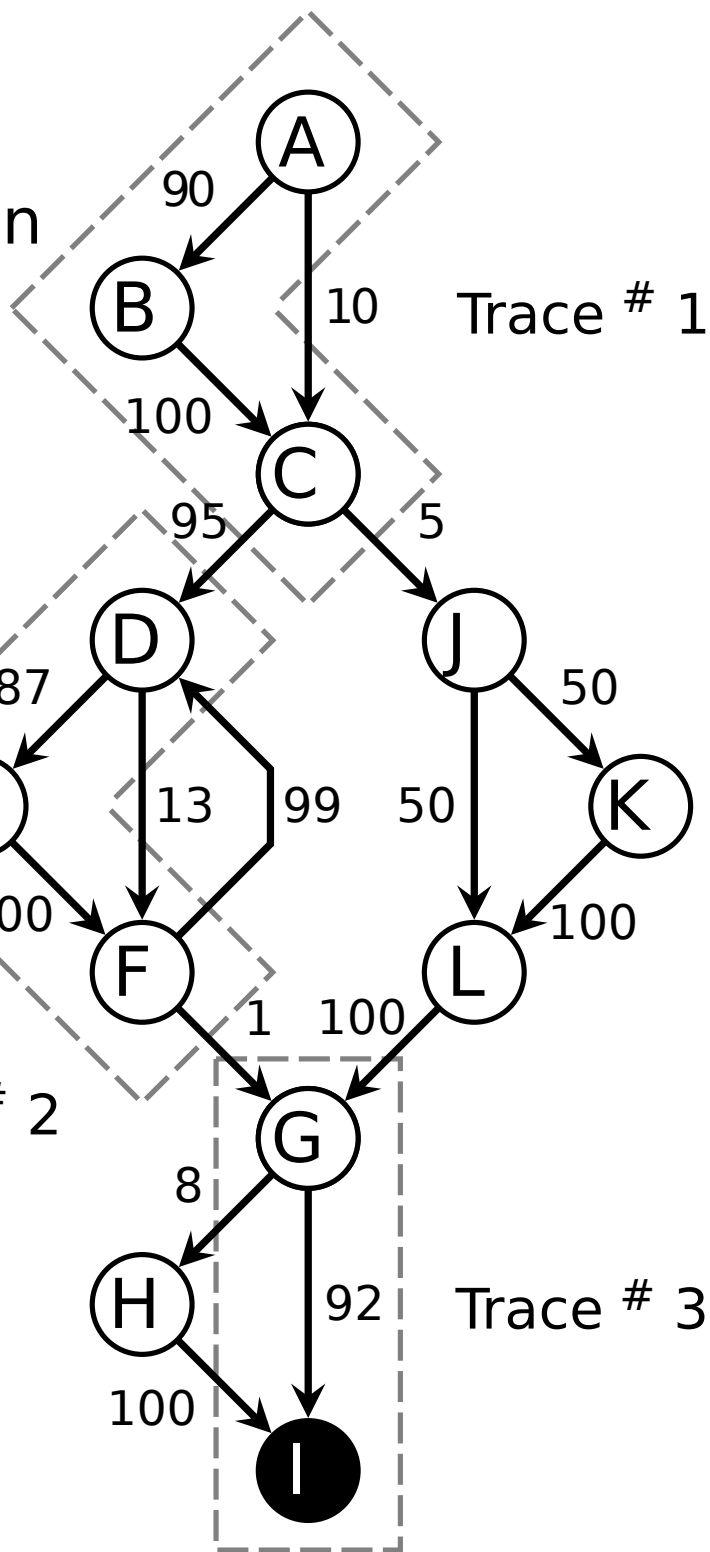
loop

Find the next unvisited node, with highest edge weight

Add node to trace

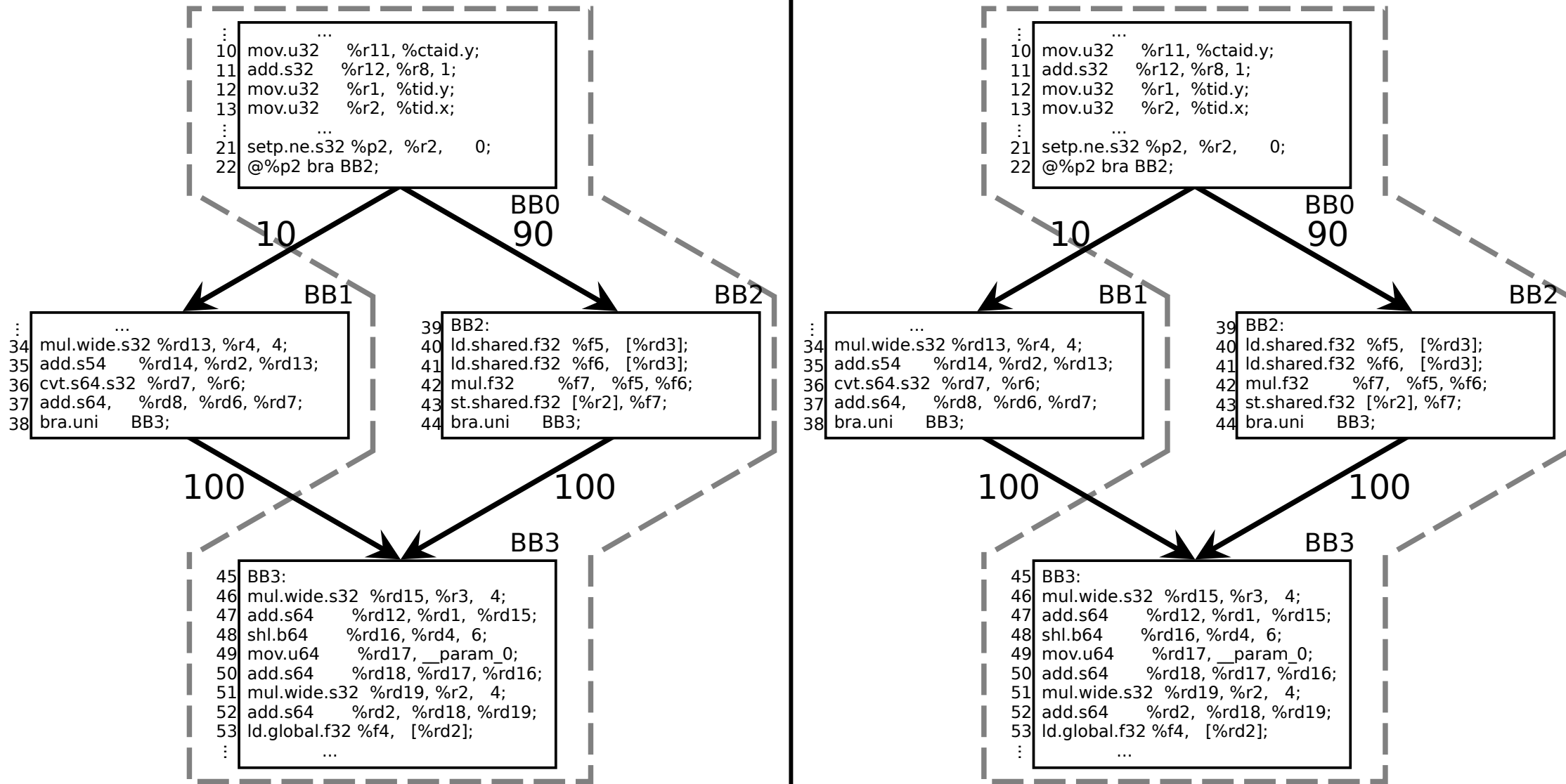
end loop

end while



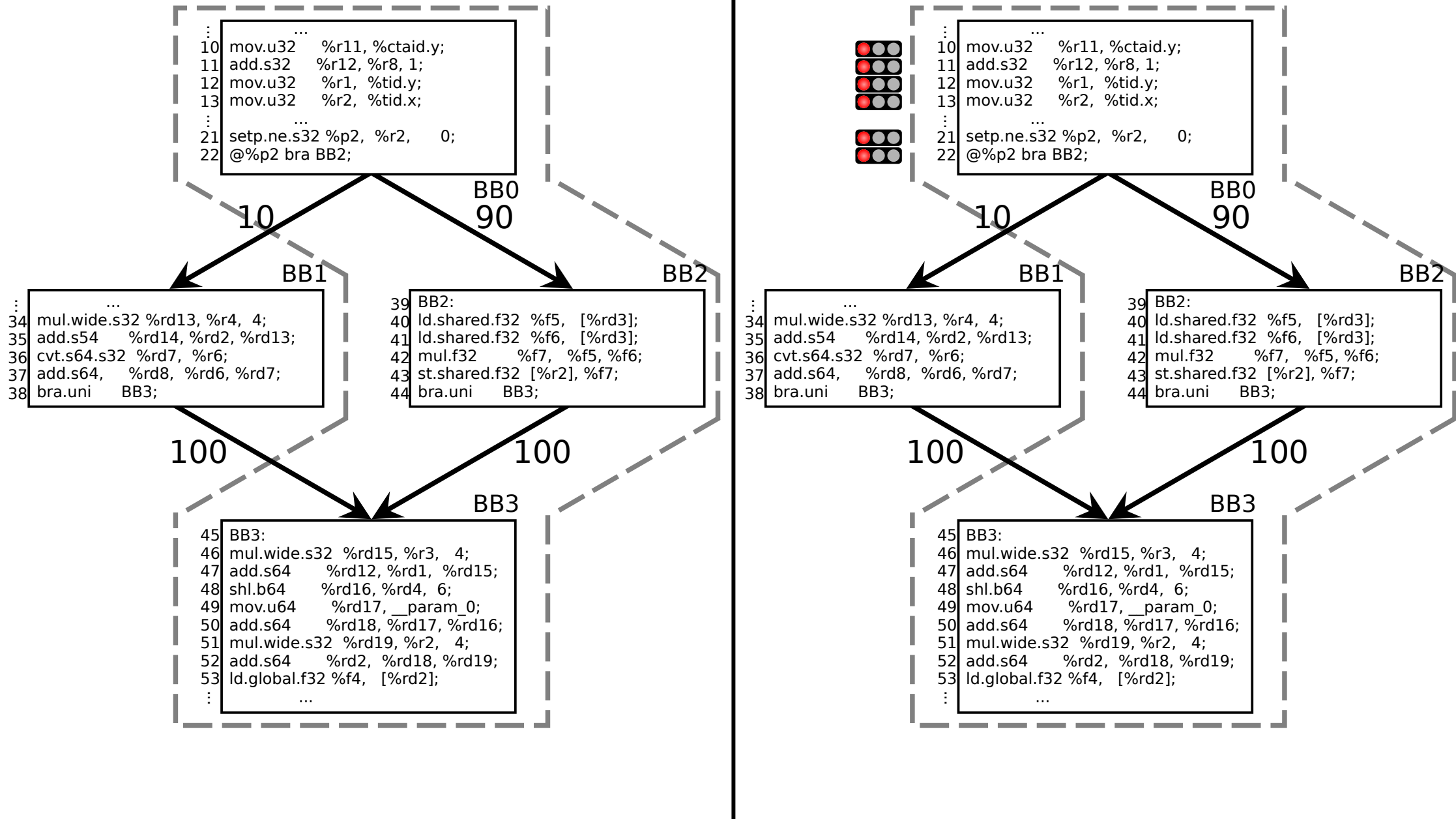
# Before

# After



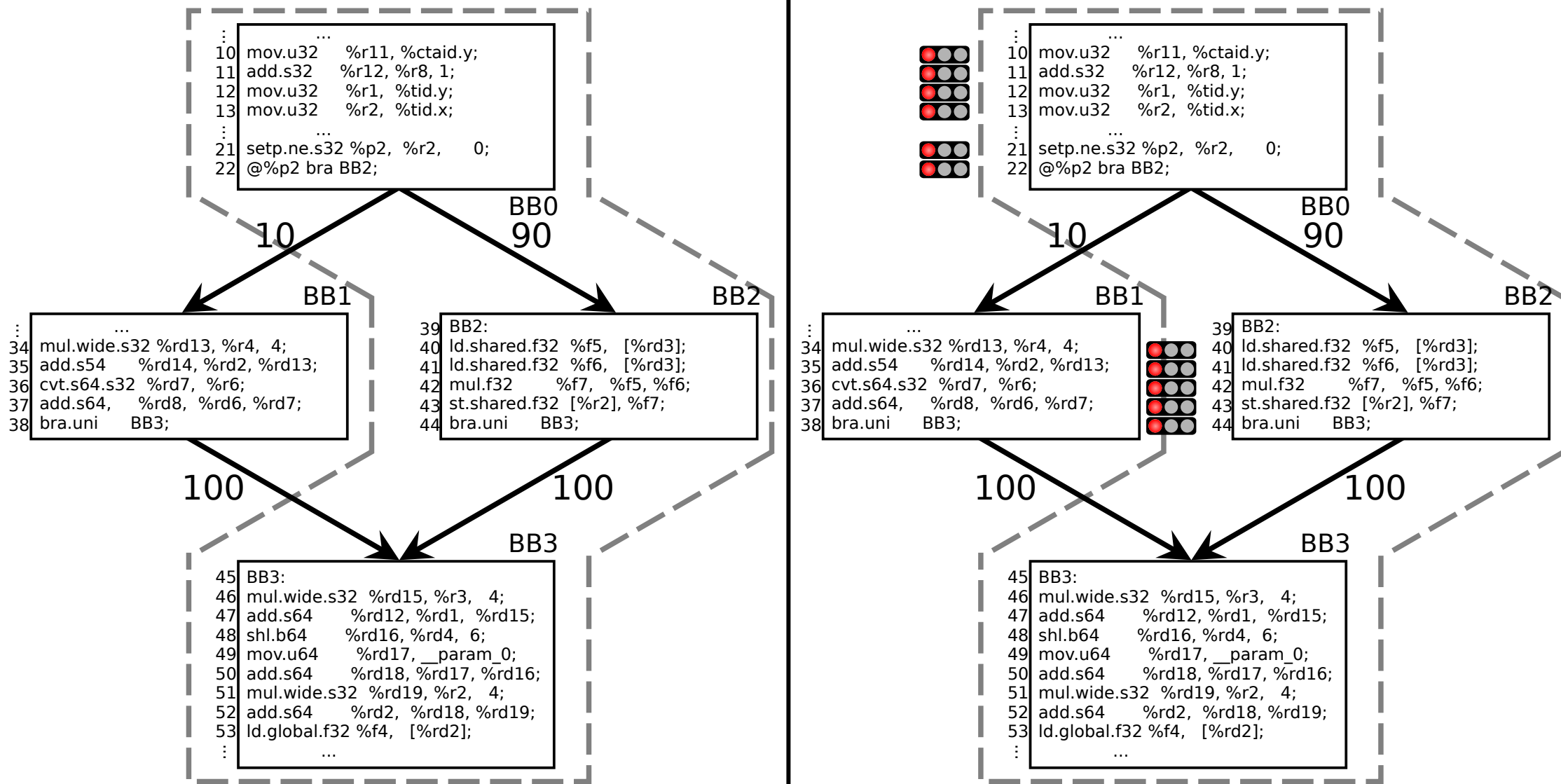
# Before

# After



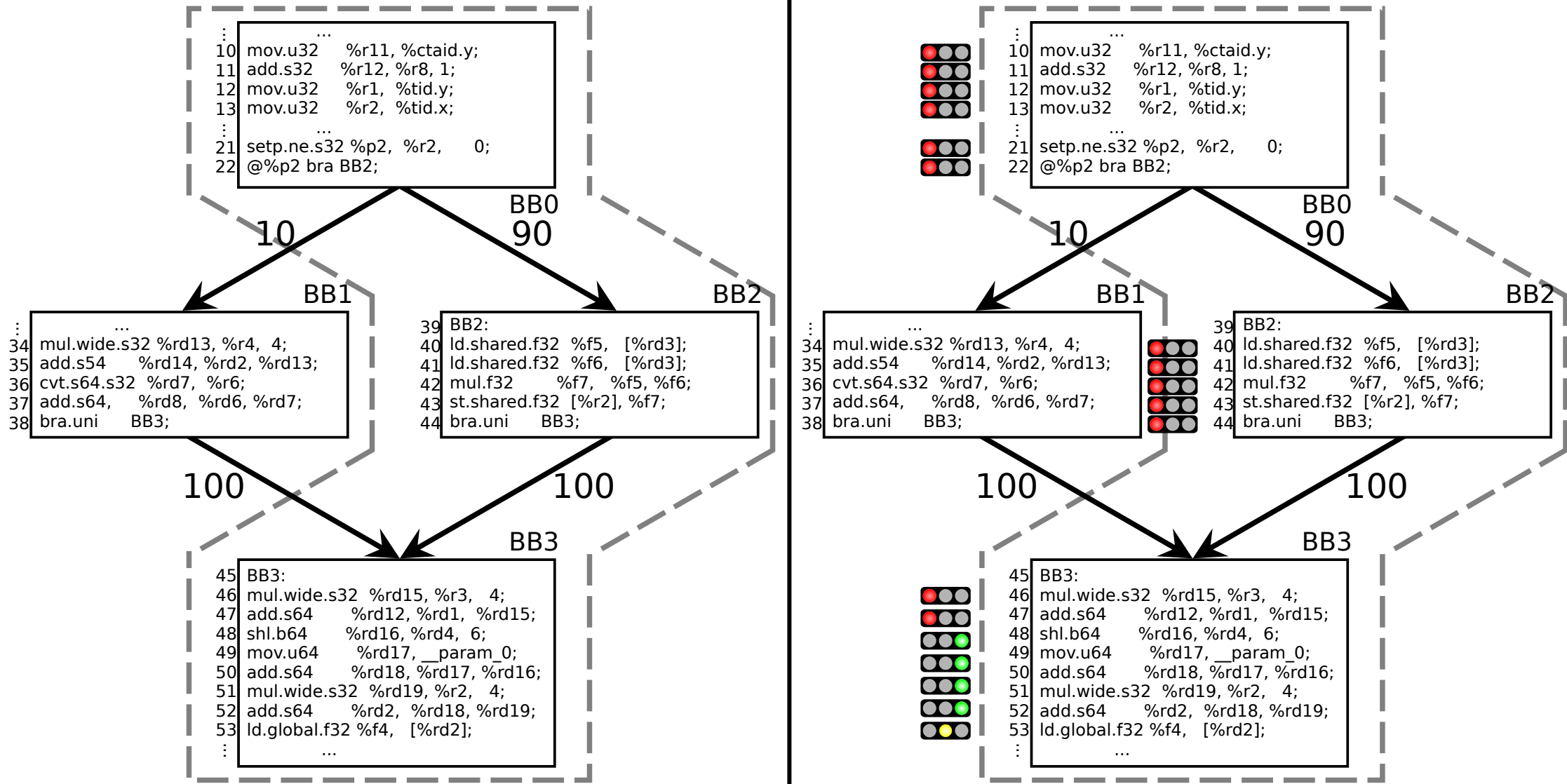
# Before

# After



# Before

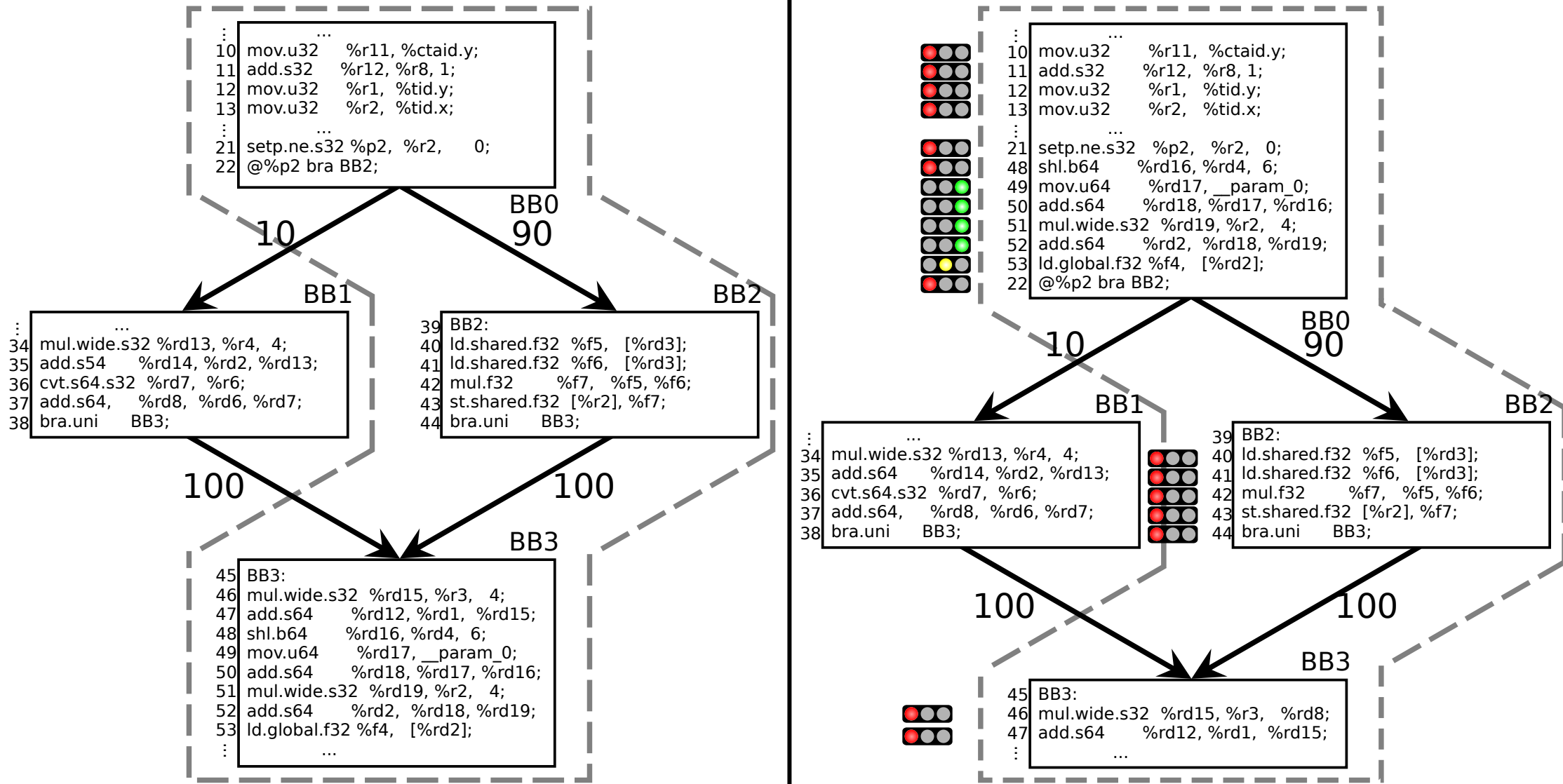
# After





# Before

# After

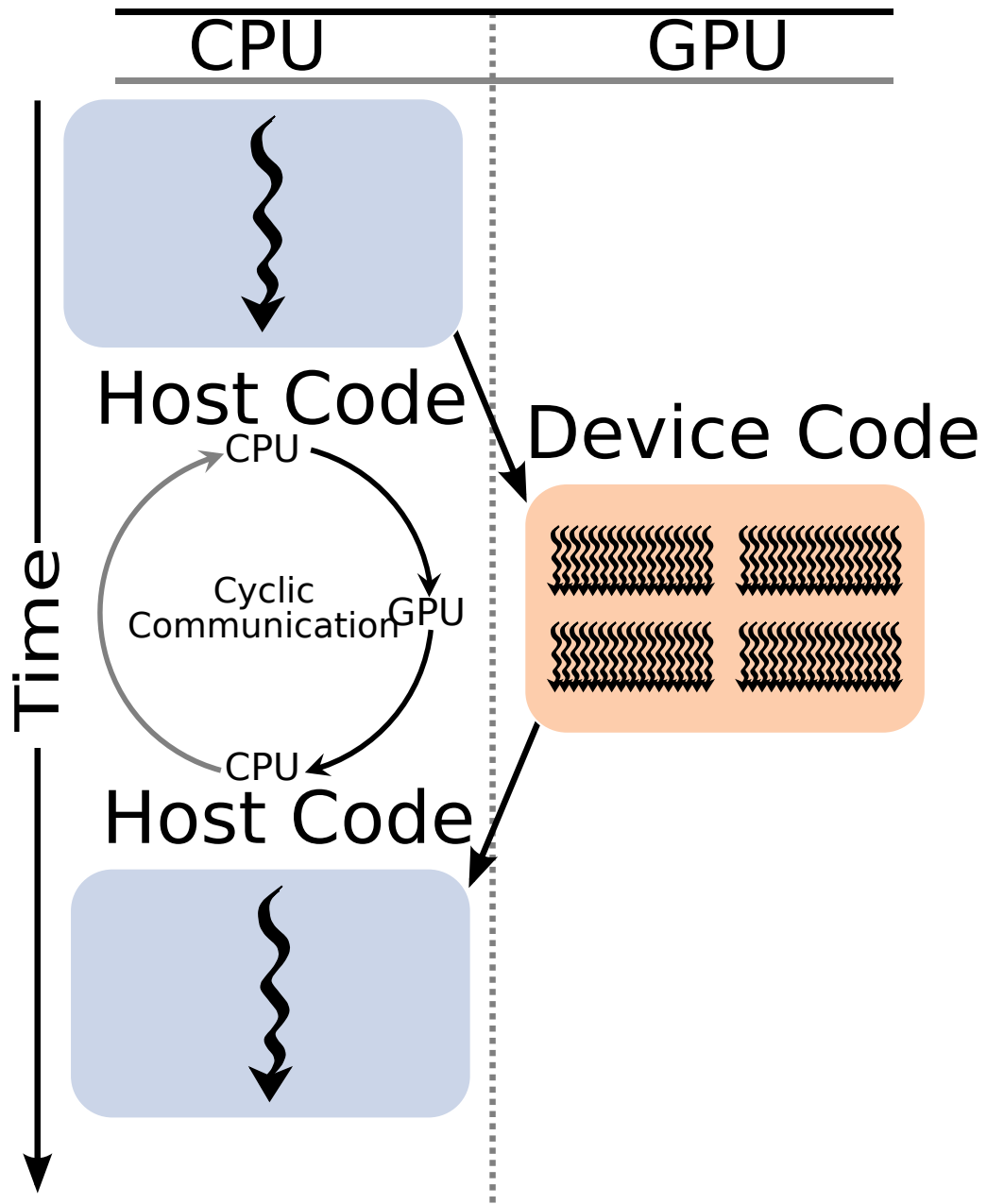




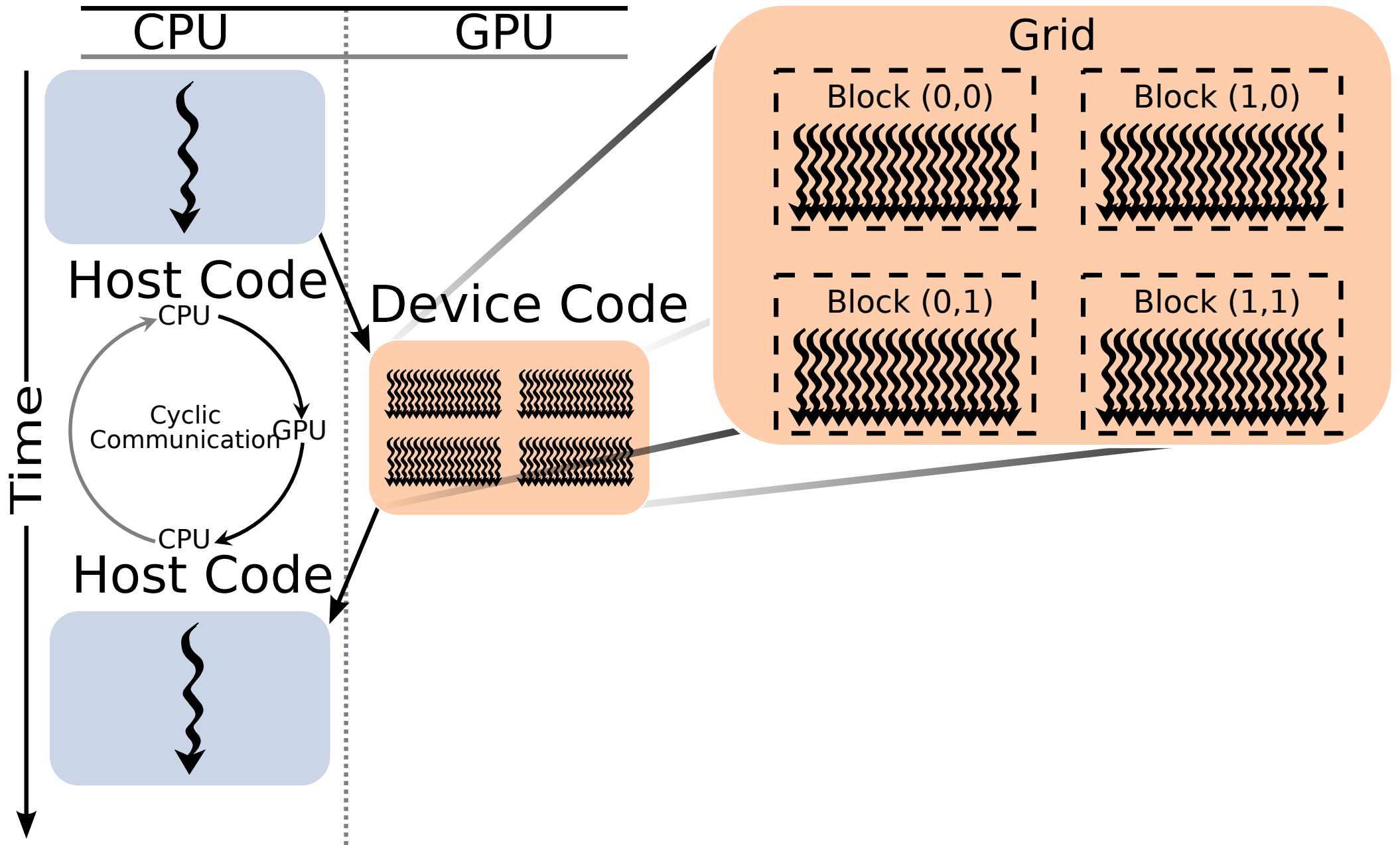
# Backup Slides

	Restricted [6]	General [6]	Boosting [36]	Deviant (GPU)
Scheduling Restrictions	Legal and Safe	Legal	none	excludes texture, shared and constant memory operations and all store instructions
Hardware Support	none	non-trapping instructions	shadow register file, shadow store buffer, and support for re-executing instructions	none
Exception Handling for Speculative Instructions	prohibited	ignored	supported	absent

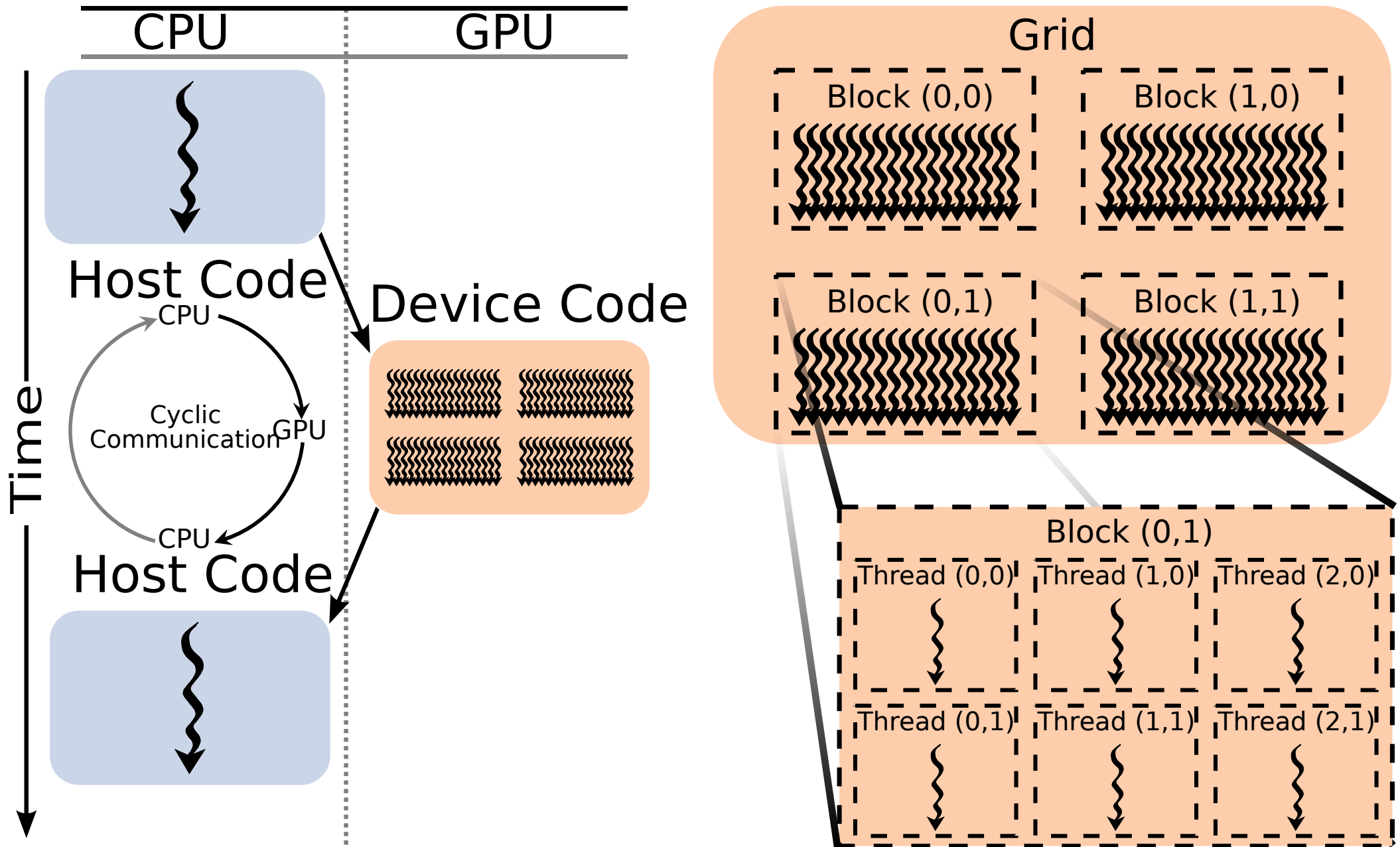
# GPU Programming Model



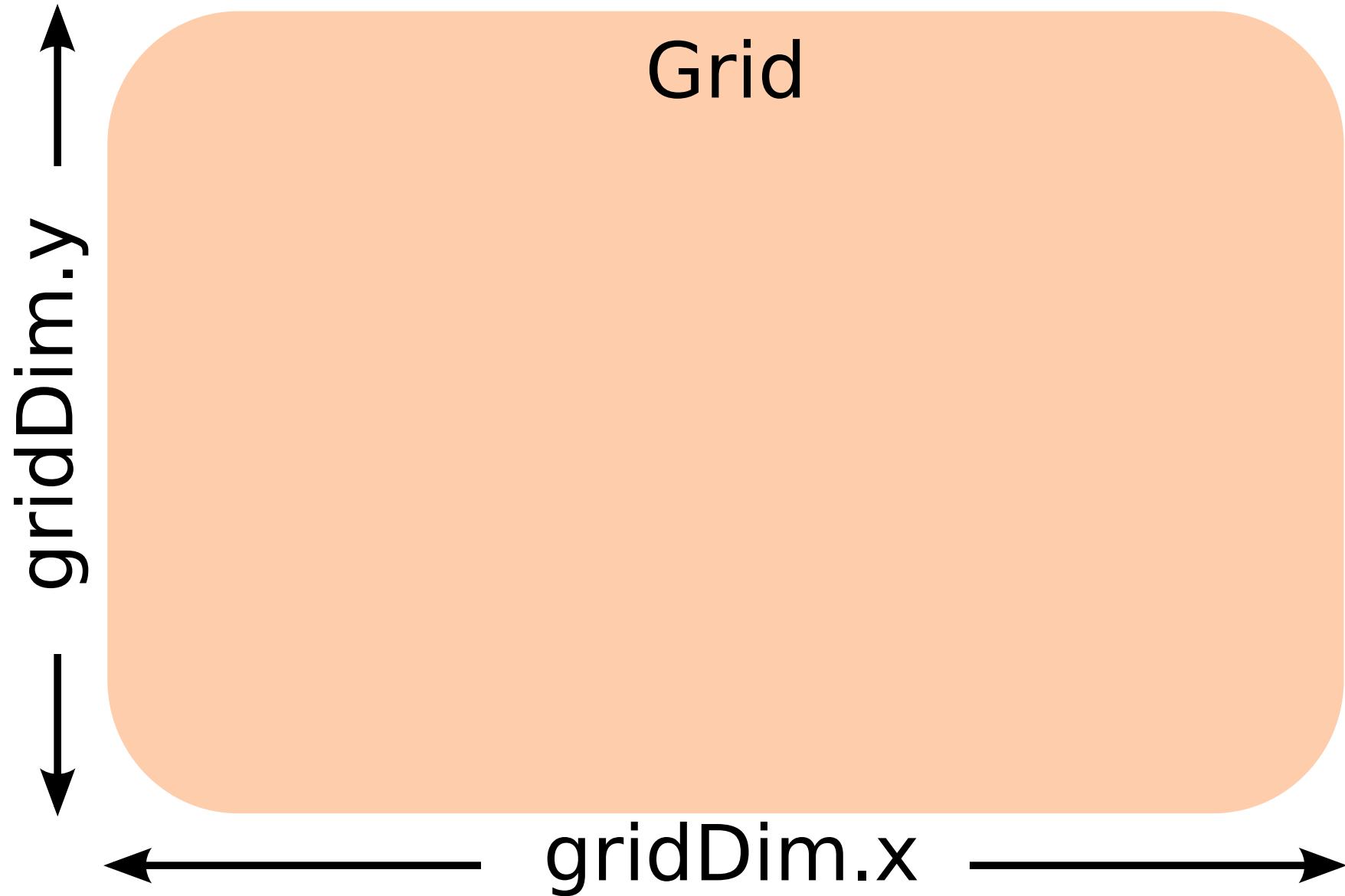
# GPU Programming Model



# GPU Programming Model

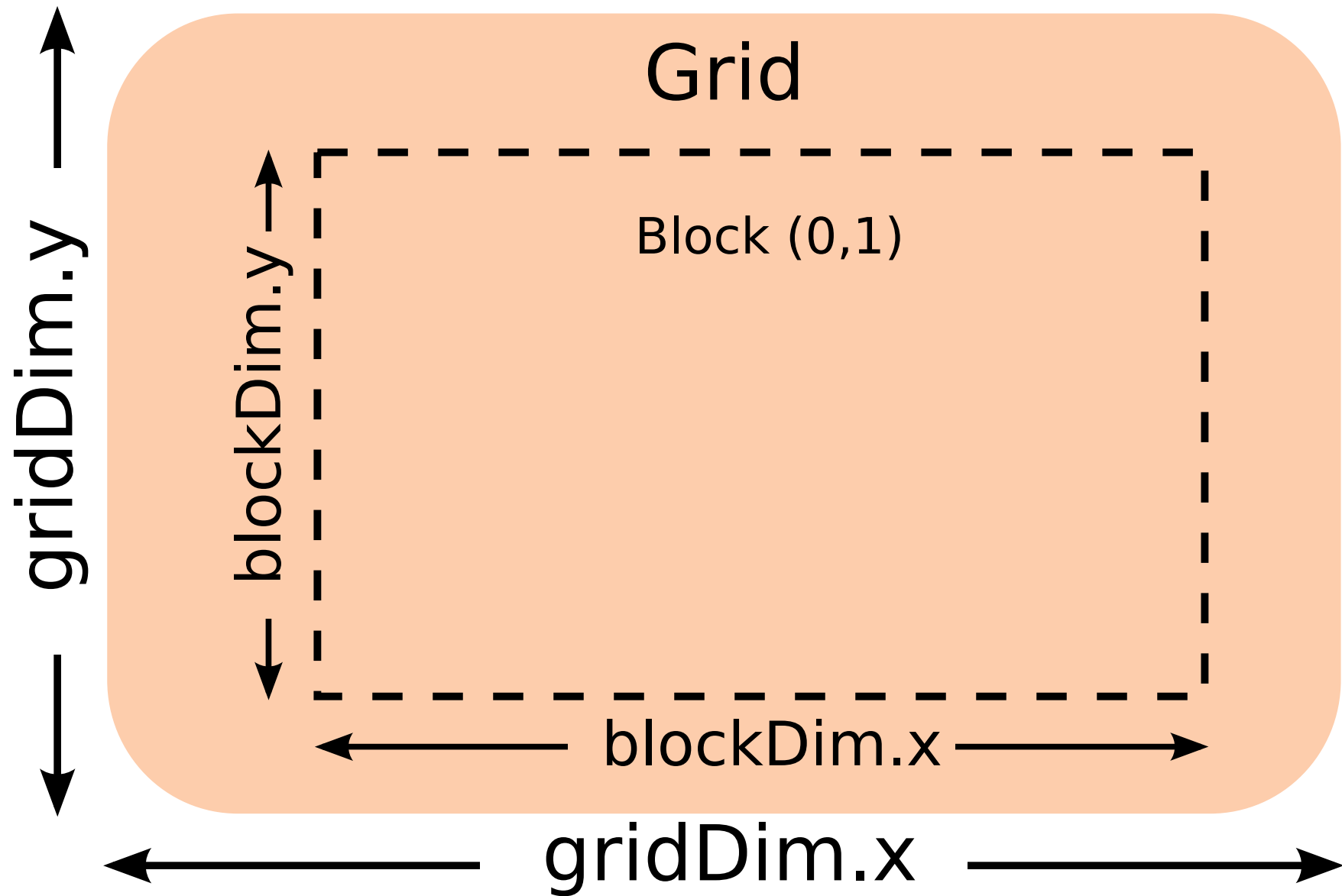


# Characterizing the Grid...

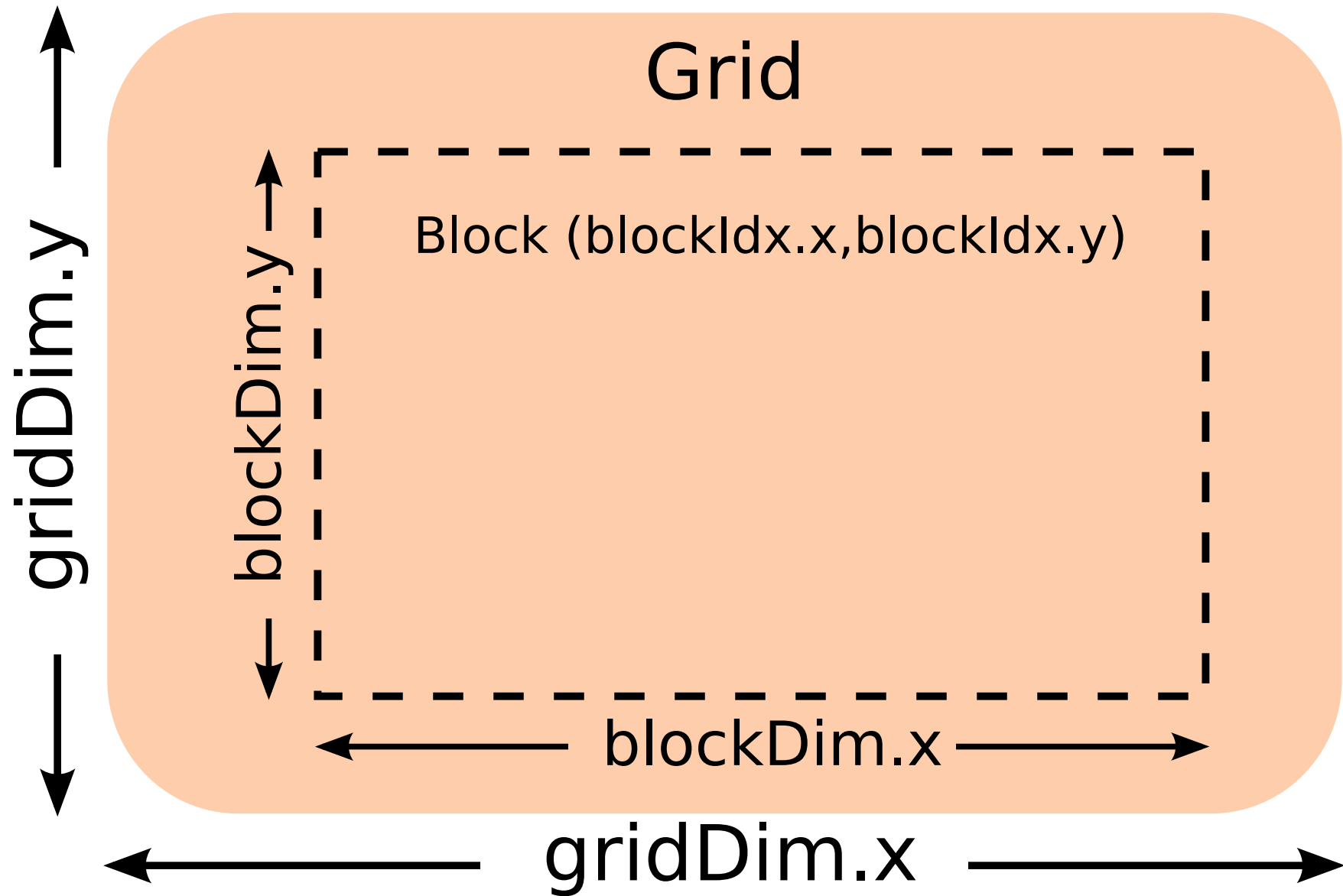




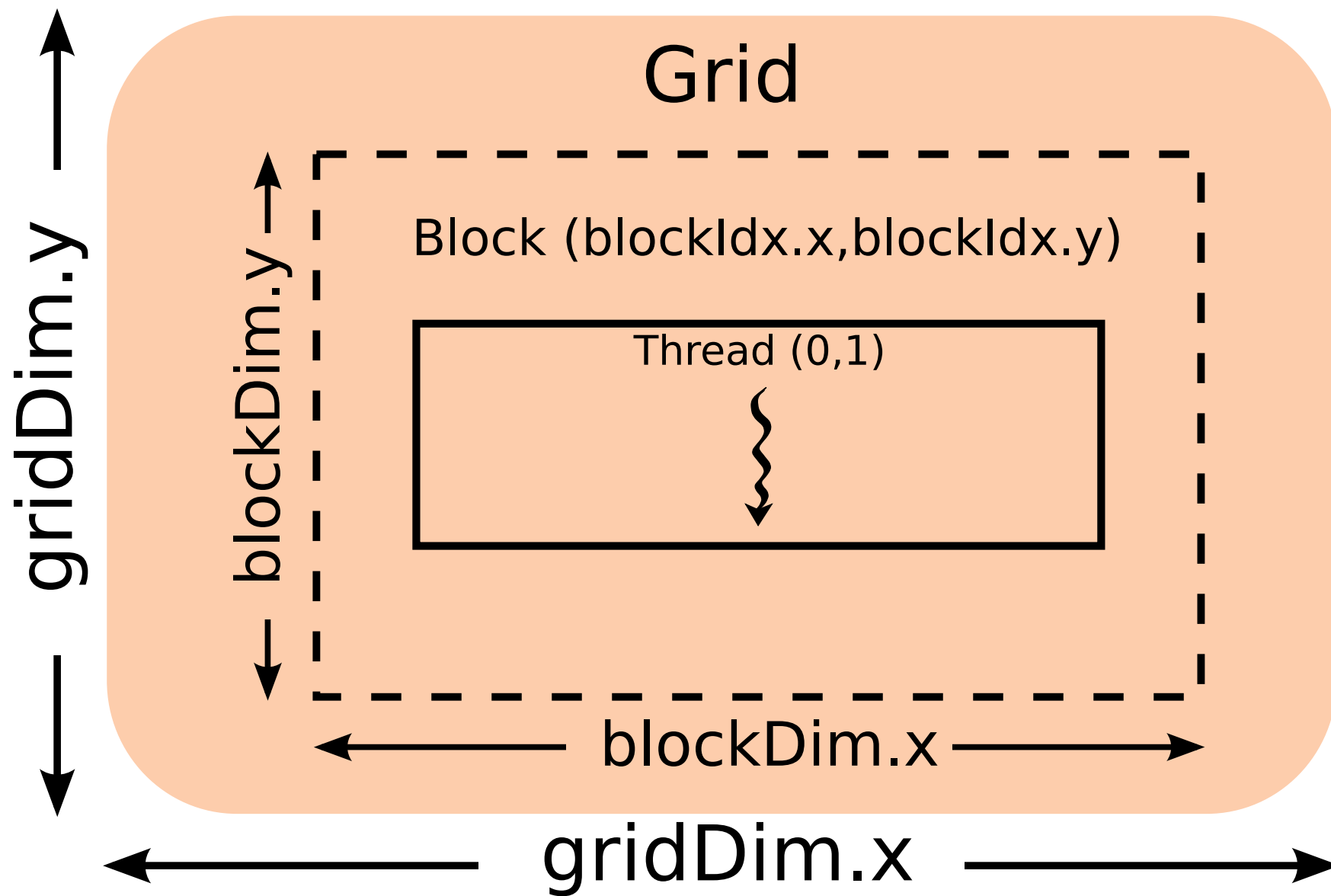
# Characterizing the Grid, Blocks...



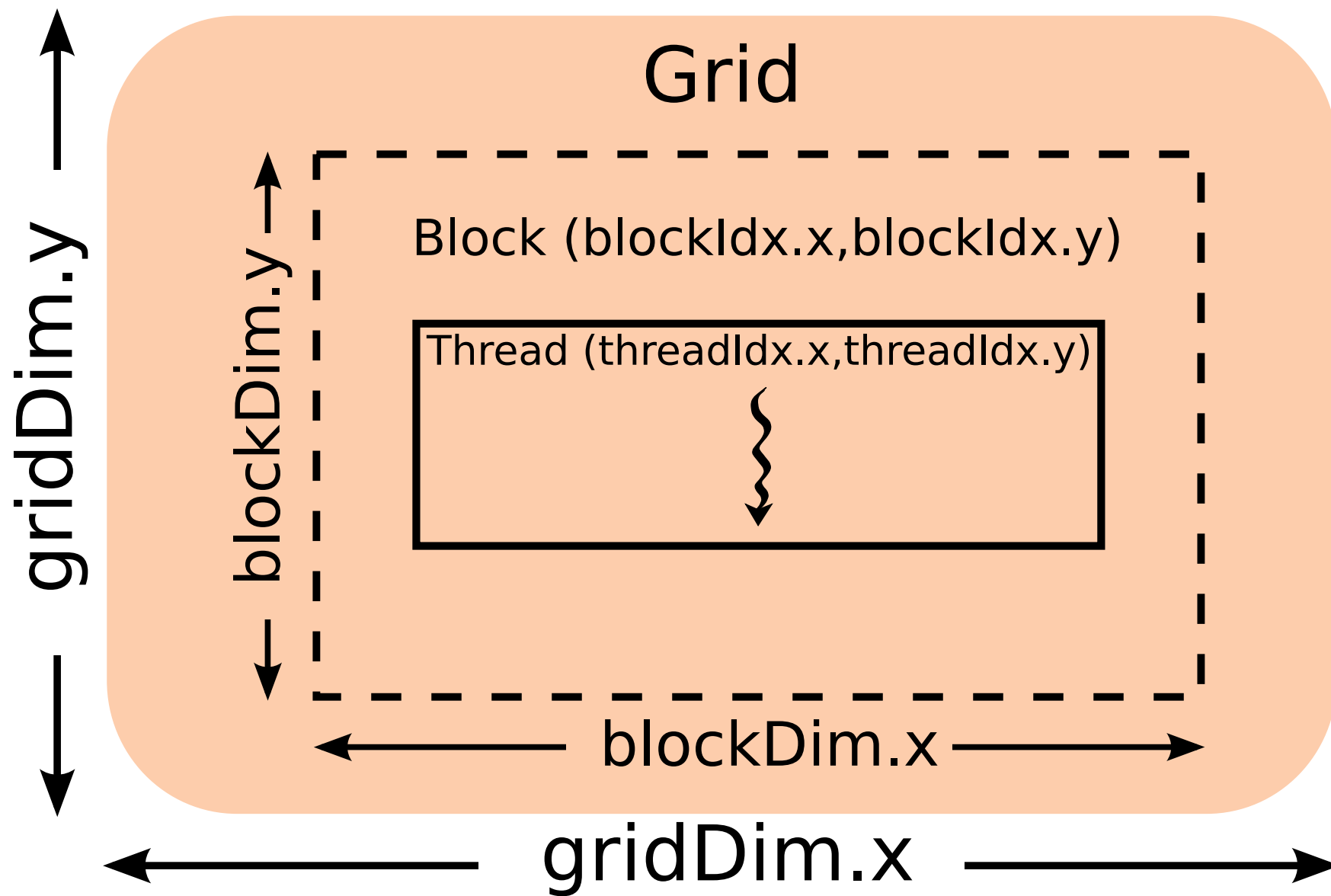
# Characterizing the Grid, Blocks...



# Characterizing the Grid, Blocks, and Threads



# Characterizing the Grid, Blocks, and Threads



# Warp Divergence Examples

Assuming one block of 128 threads...

Example	Divergence?
if (threadIdx.x < 32) { }	

# Warp Divergence Examples

Assuming one block of 128 threads...

Example	Divergence?
<code>if (threadIdx.x &lt; 32) { }</code>	NO
<code>if (threadIdx.x &gt; 15) { }</code>	

# Warp Divergence Examples

Assuming one block of 128 threads...

Example	Divergence?
<code>if (threadIdx.x &lt; 32) { }</code>	NO
<code>if (threadIdx.x &gt; 15) { }</code>	YES
<code>if (threadIdx.x &gt; 65) { }</code>	

# Warp Divergence Examples

Assuming one block of 128 threads...

Example	Divergence?
<code>if (threadIdx.x &lt; 32) { }</code>	NO
<code>if (threadIdx.x &gt; 15) { }</code>	YES
<code>if (threadIdx.x &gt; 65) { }</code>	YES
<code>if (BlockIdx.x &gt; 1) { }</code>	



# Warp Divergence Examples

Assuming one block of 128 threads...

Example	Divergence?
<code>if (threadIdx.x &lt; 32) { }</code>	NO
<code>if (threadIdx.x &gt; 15) { }</code>	YES
<code>if (threadIdx.x &gt; 65) { }</code>	YES
<code>if (blockIdx.x &gt; 1) { }</code>	NO