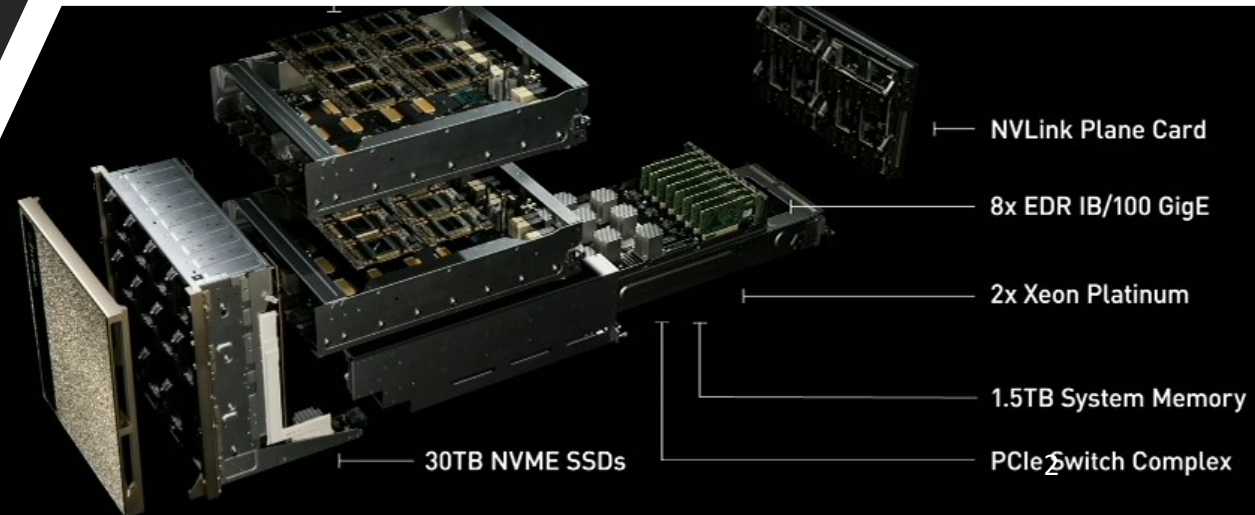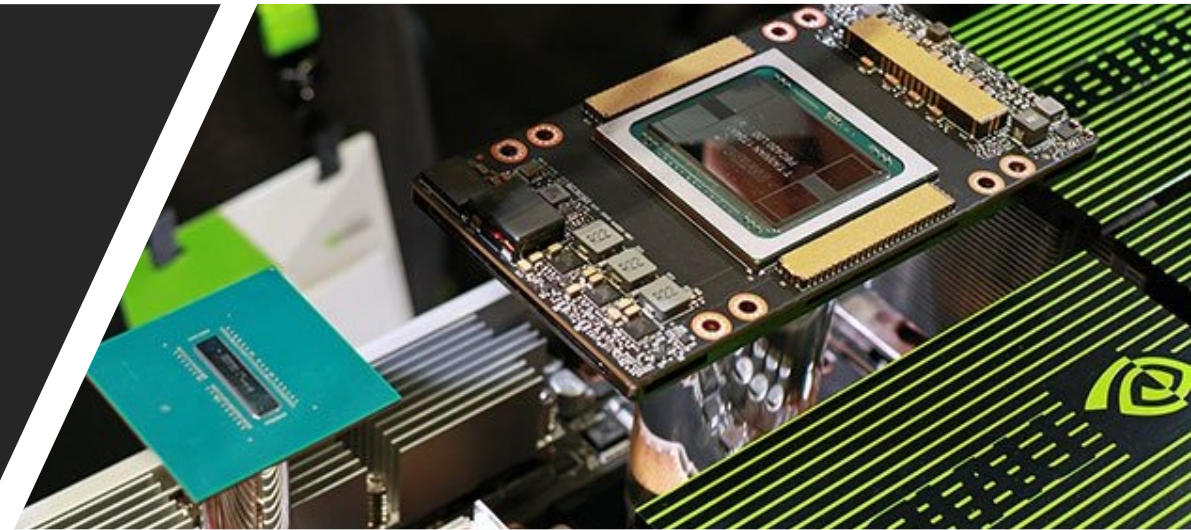# Parallel Architectures
# Parallel Algorithms
# CUDA

Chris Rossbach

cs378h

# Outline for Today

- Questions?
- Administrivia
  - pedagogical-* machines should be available
- Agenda
  - Parallel Algorithms
  - CUDA


- Acknowledgements: http://developer.download.nvidia.com/compute/develo pertrainingmaterials/presentations/cuda_language/Intro duction_to_CUDA_C.pptx



NVLink Plane Card

8x EDR IB/100 GigE

2x Xeon Platinum

1.5TB System Memory

30TB NVME SSDs

PCIe Switch Complex
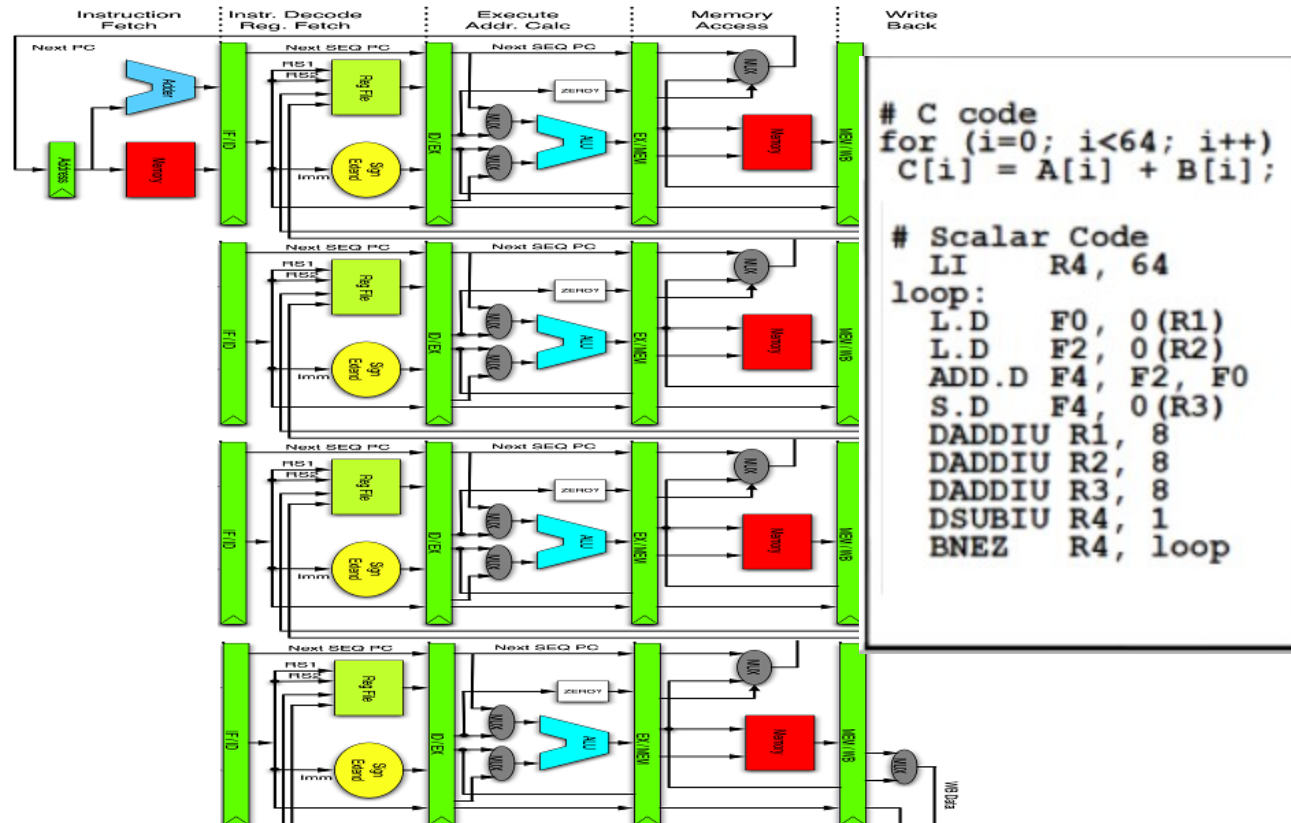
# Faux Quiz Questions

- What is a reduction? A prefix sum? Why are they hard to parallelize and what basic techniques can be used to parallelize them?

- Define flow dependence, output dependence, and anti-dependence: give an example of each. Why/how do compilers use them to detect loop-independent vs loop-carried dependences?

- What is the difference between a thread-block and a warp?

- How/Why must programmers copy data back and forth to a GPU?

- What is "shared memory" in CUDA? Describe a setting in which it might be useful.

- CUDA kernels have implicit barrier synchronization. Why is __syncthreads() necessary in light of this fact?

- How might one implement locks on a GPU?

- What ordering guarantees does a GPU provide across different hardware threads' access to a single memory location? To two disjoint locations?

- When is it safe for one GPU thread to wait (e.g. by spinning) for another?

# Review: what is a vector processor?

```
# C code
for (i=0; i<64; i++)
 C[i] = A[i] + B[i];

# Scalar Code
  LI      R4, 64
loop:
  L.D     F0, 0(R1)
  L.D     F2, 0(R2)
  ADD.D F4, F2, F0
  S.D     F4, 0(R3)
  DADDIU R1, 8
  DADDIU R2, 8
  DADDIU R3, 8
  DSUBIU R4, 1
  BNEZ    R4, loop
```
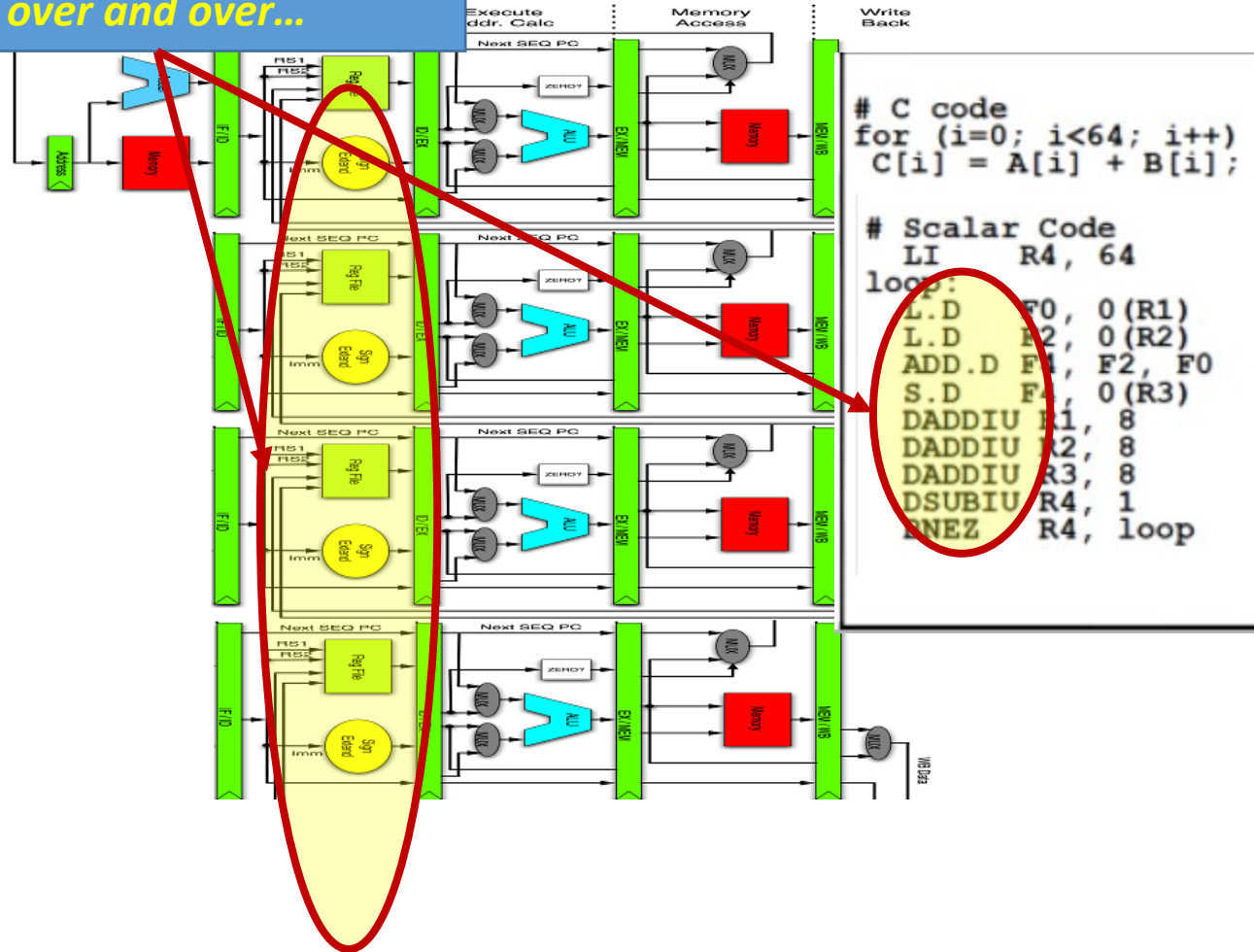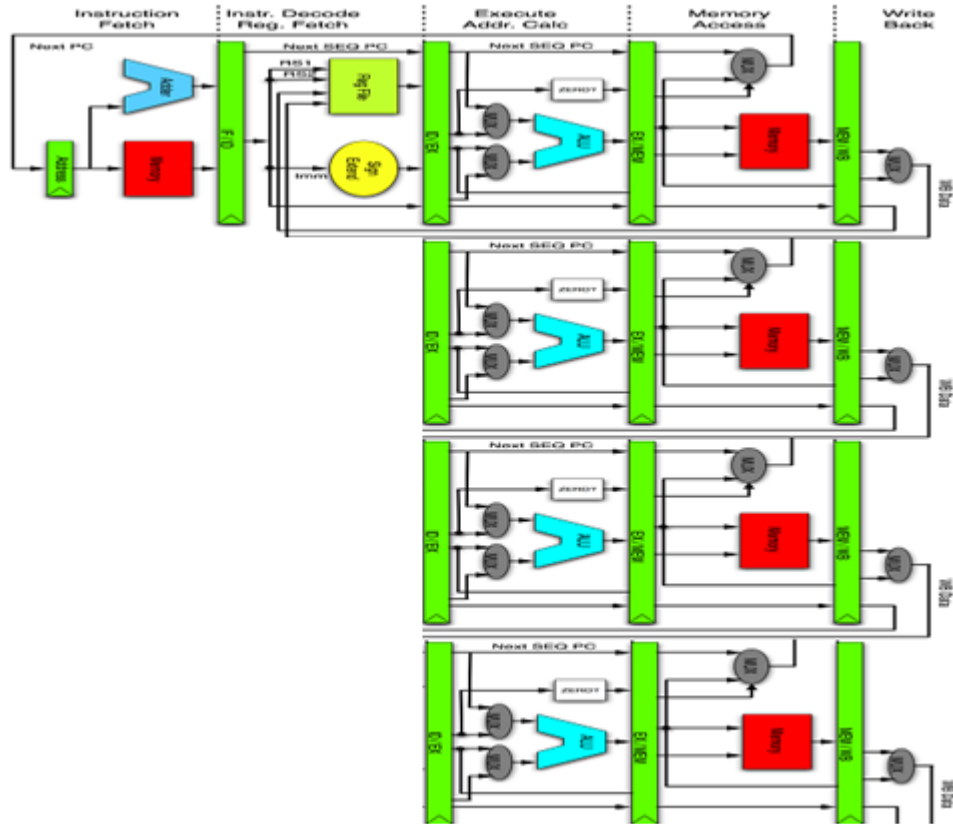
# Review: what is a vector processor?



```
# C code
for (i=0; i<64; i++)
  C[i] = A[i] + B[i];


# Scalar Code
    LI      R4, 64
loop:
    L.D     F0, 0(R1)
    L.D     F2, 0(R2)
    ADD.D   F4, F2, F0
    S.D     F4, 0(R3)
    DADDIU  R1, 8
    DADDIU  R2, 8
    DADDIU  R3, 8
    DSUBIU  R4, 1
    BNEZ    R4, loop
```

# Review: what is a vector processor?



Dont decode same instruction over and over…

```
# C code
for (i=0; i<64; i++)
  C[i] = A[i] + B[i];

# Scalar Code
    LI      R4, 64
loop:
    L.D     F0, 0(R1)
    L.D     F2, 0(R2)
    ADD.D   F4, F2, F0
    S.D     F4, 0(R3)
    DADDIU  R1, 8
    DADDIU  R2, 8
    DADDIU  R3, 8
    DSUBIU  R4, 1
    BNEZ    R4, loop
```

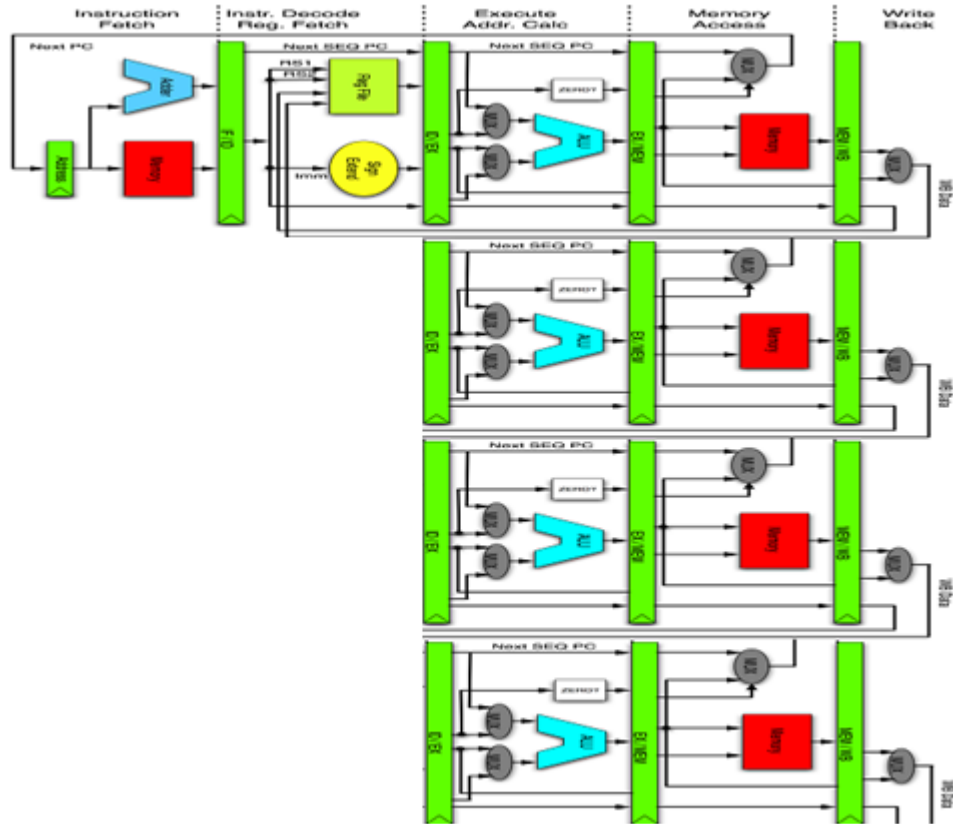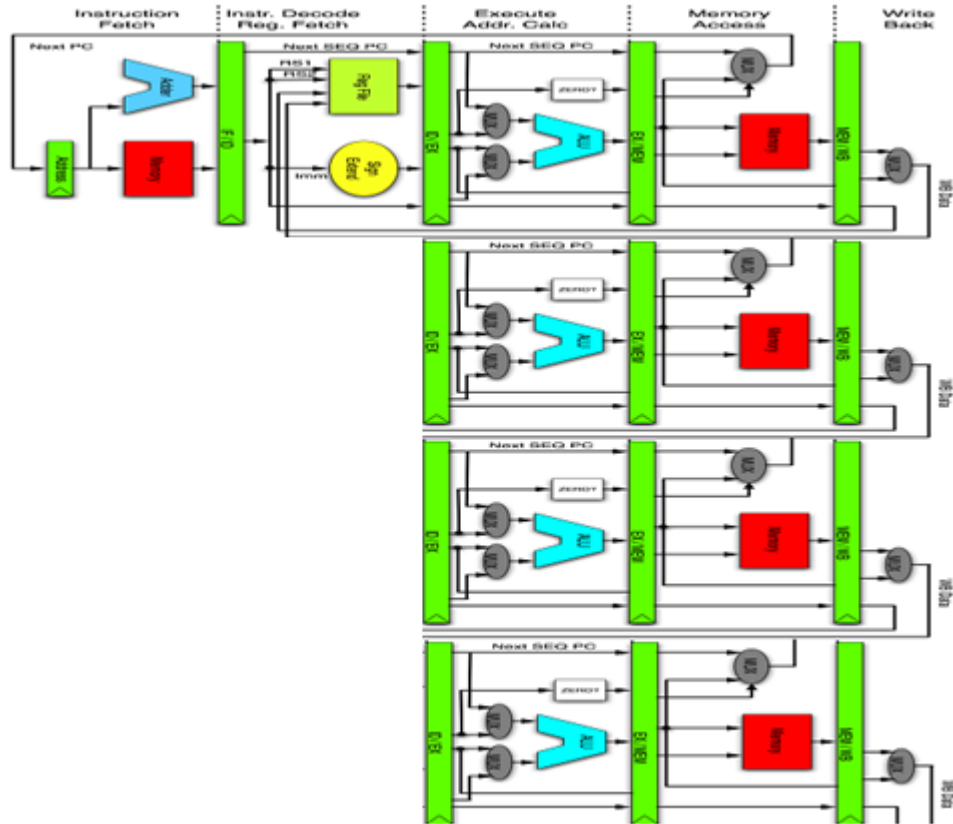# Review: what is a vector processor?

# Review: what is a vector processor?



```
# C code
for (i=0; i<64; i++)
  C[i] = A[i] + B[i];
```

```
# Scalar Code
    LI      R4, 64
loop:
    L.D     F0, 0(R1)
    L.D     F2, 0(R2)
    ADD.D   F4, F2, F0
    S.D     F4, 0(R3)
    DADDIU R1, 8
    DADDIU R2, 8
    DADDIU R3, 8
    DSUBIU R4, 1
    BNEZ    R4, loop
```

```
# Vector Code
    LI      VLR, 64
    LV      V1, R1
    LV      V2, R2
    ADDV.D V3, V1, V2
    SV      V3, R3
```
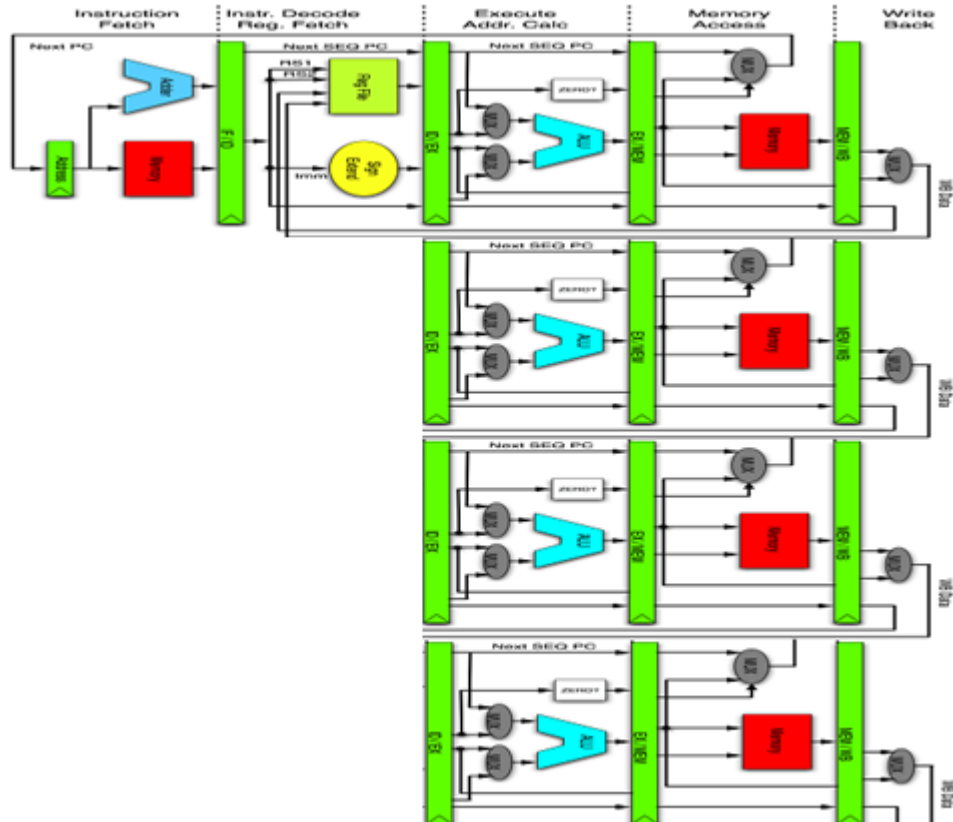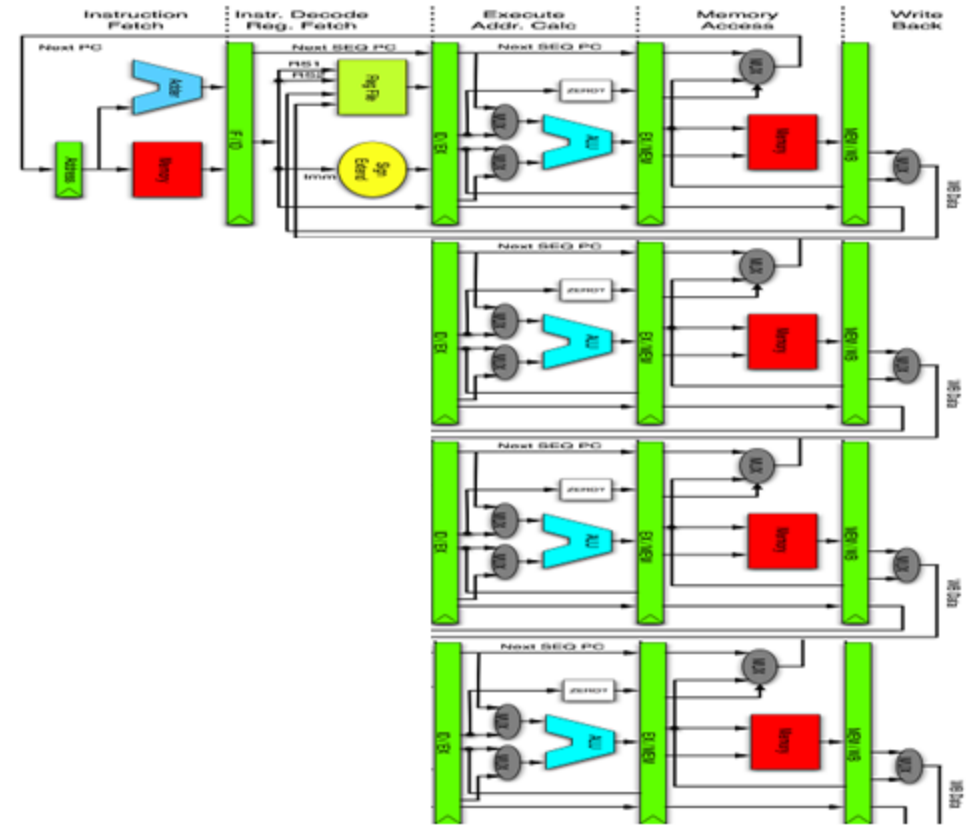
# Review: what is a vector processor?



Implementation:
- Instruction fetch control logic shared
- Same instruction stream executed on
- Multiple pipelines
- Multiple different operands in parallel

```
# C code
for (i=0; i<64; i++)
  C[i] = A[i] + B[i];
```

```
# Scalar Code
    LI     R4, 64
loop:
    L.D    F0, 0(R1)
    L.D    F2, 0(R2)
    ADD.D  F4, F2, F0
    S.D    F4, 0(R3)
    DADDIU R1, 8
    DADDIU R2, 8
    DADDIU R3, 8
    DSUBIU R4, 1
    BNEZ   R4, loop
```

```
# Vector Code
    LI     VLR, 64
    LV     V1, R1
    LV     V2, R2
    ADDV.D V3, V1, V2
    SV     V3, R3
```
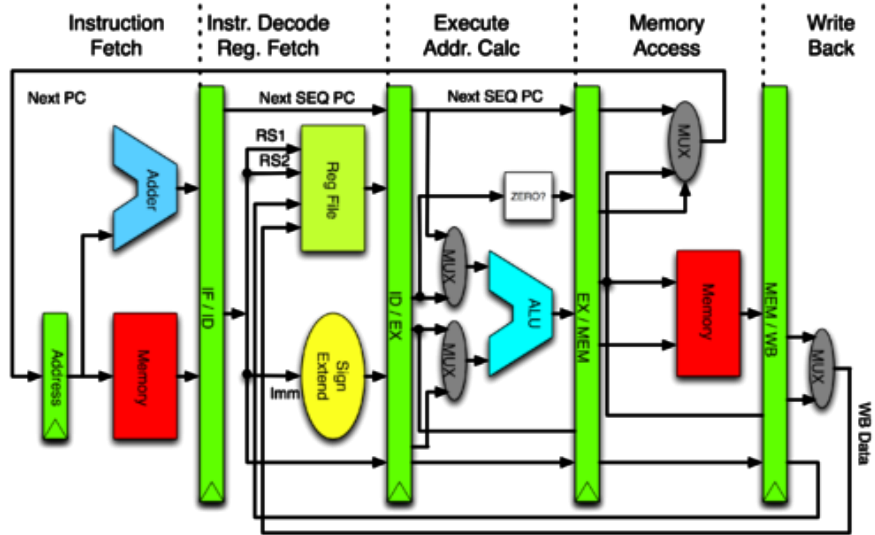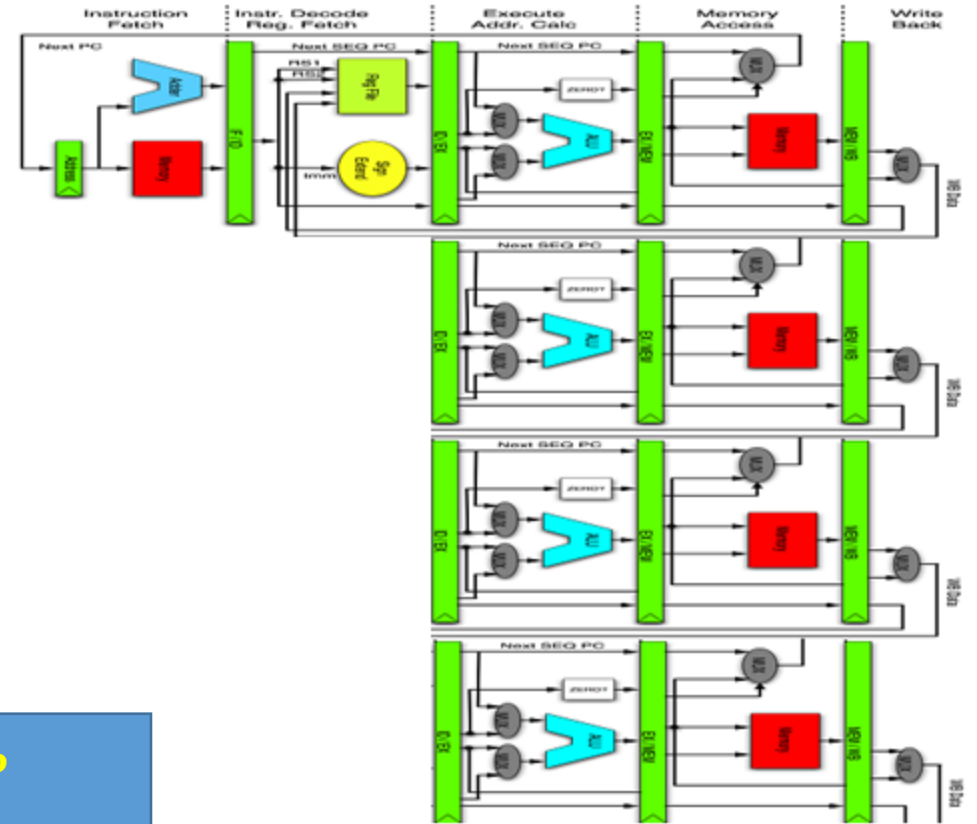
# Review: what is a vector proc



Imp

- In
- Sa
- M
- Multiple different operands in parallel

**Scalar Registers**
r15
r0

**Vector Registers**
v15
v0 [0] [1] [2] [VLRMAX-1]
[63], [127], [255], …

Vector Length Register VLR

**Vector Arithmetic Instructions**
ADDV v3, v1, v2

v1
v2
v3
[0] [1] [VLR-1]

**Vector Load and Store Instructions**
LV v1, r1, r2

Vector Register
v1

Base, r1   Stride, r2   Memory

```
# C code
for (i=0; i<64; i++)
 C[i] = A[i] + B[i];

            # Scalar Code
            LI      R4, 64
         loop:
            L.D     F0, 0(R1)
            L.D     F2, 0(R2)
            ADD.D   F4, F2, F0
            S.D     F4, 0(R3)
            DADDIU  R1, 8
            DADDIU  R2, 8
            DADDIU  R3, 8
            DSUBIU  R4, 1
            BNEZ    R4, loop

# Vector Code
LI      VLR, 64
LV      V1, R1
LV      V2, R2
ADDV.D  V3, V1, V2
SV      V3, R3
```

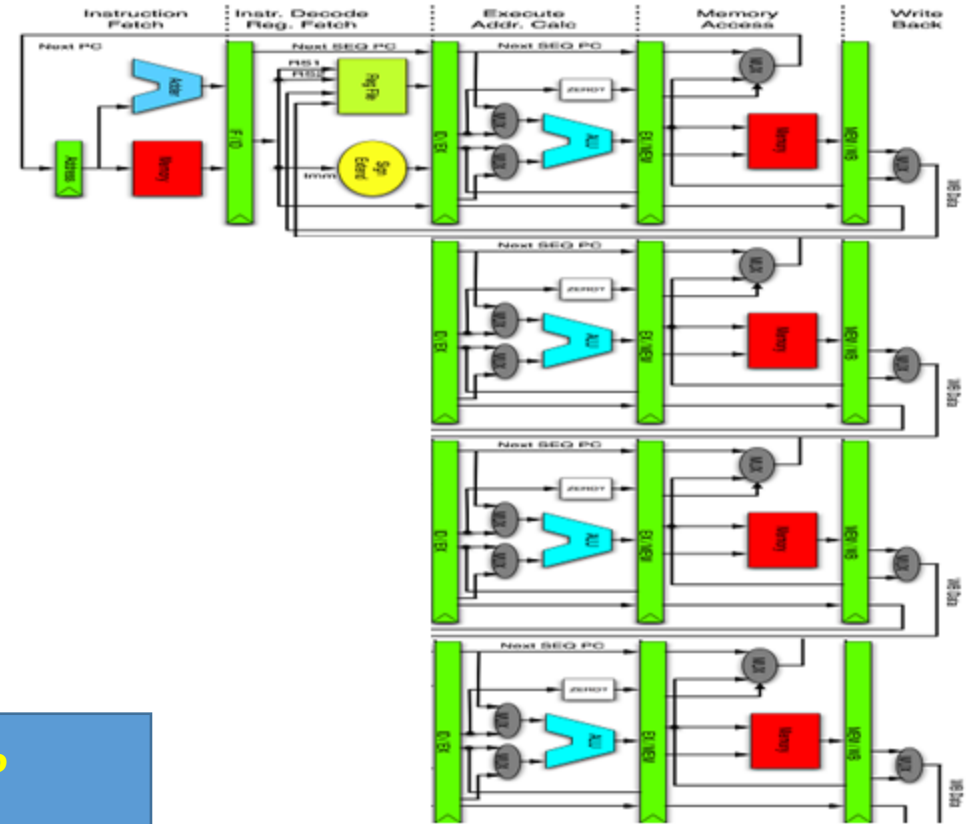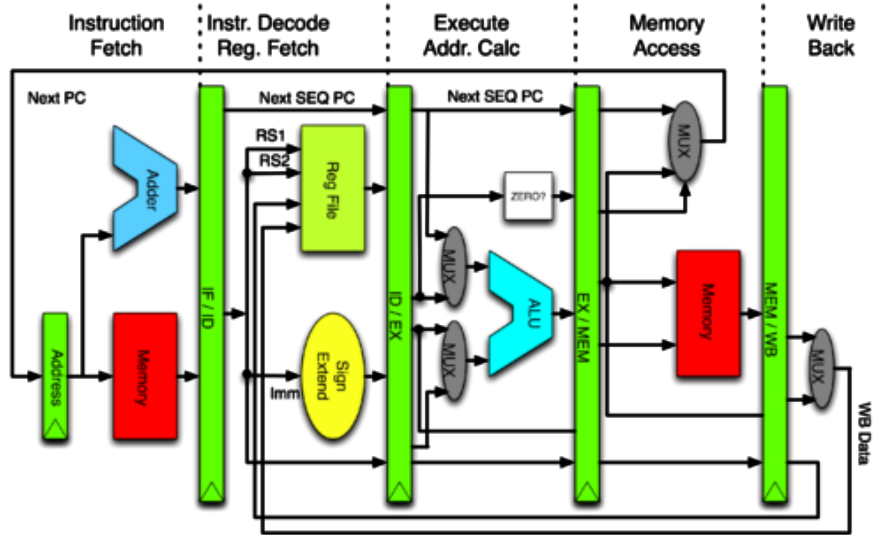# When does vector processing help?

# When does vector processing help?



What are the potential bottlenecks here?
When can it improve throughput?

# When does vector processing help?



*What are the potential bottlenecks here?*
*When can it improve throughput?*

*Only helps if memory can keep the pipeline busy!*

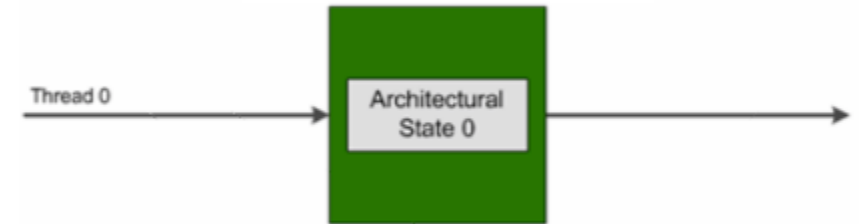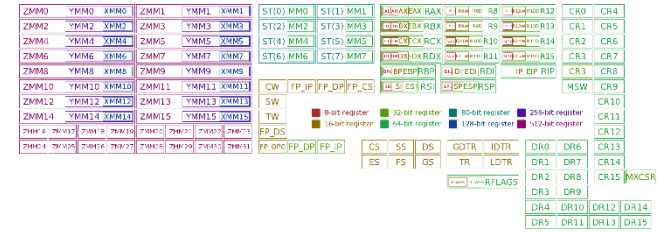# Hardware multi-threading

# Hardware multi-threading

- Address memory bottleneck

# Hardware multi-threading

- Address memory bottleneck

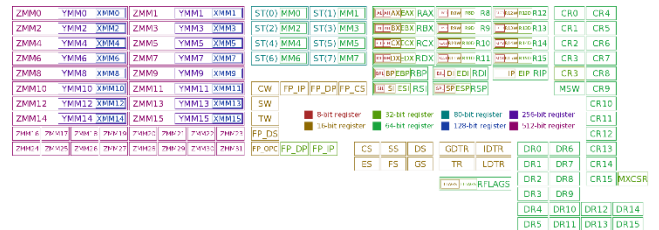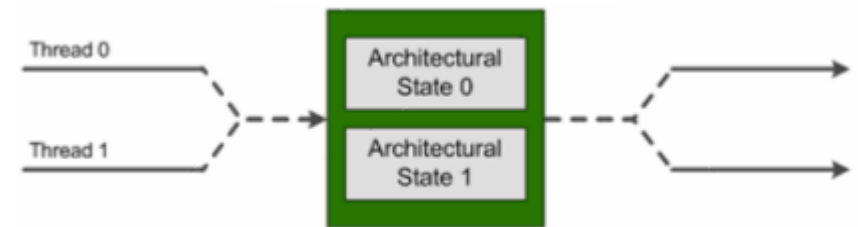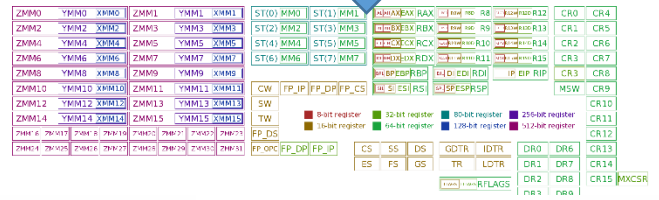- Share exec unit across
  - Instruction streams
  - Switch on stalls

# Hardware multi-threading



- Address memory bottleneck
- Share exec unit across
  - Instruction streams
  - Switch on stalls



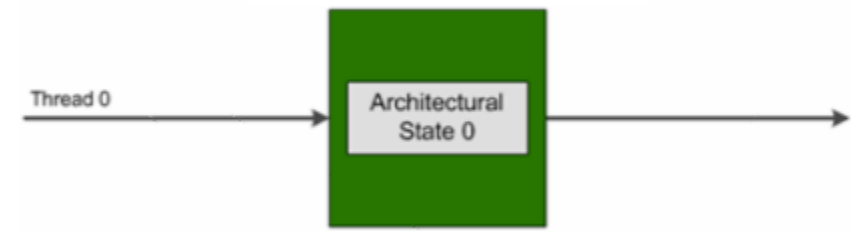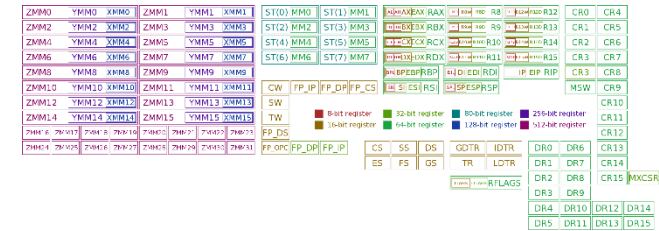Thread 0 → Architectural State 0 →

# Hardware multi-threading

- Address memory bottleneck

- Share exec unit across
  - Instruction streams
  - Switch on stalls

# Hardware multi-threading

- Address memory bottleneck

- Share exec unit across
  - Instruction streams
  - Switch on stalls

- Looks like multiple cores to the OS

# Hardware multi-threading

- Address memory bottleneck

- Share exec unit across
  - Instruction streams
  - Switch on stalls

- Looks like multiple cores to the OS
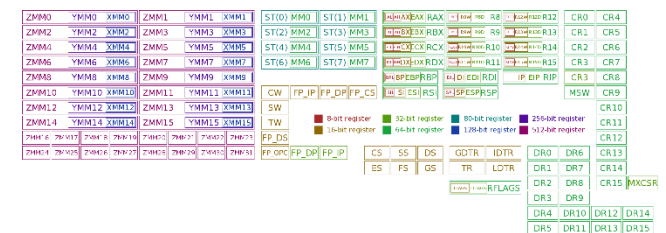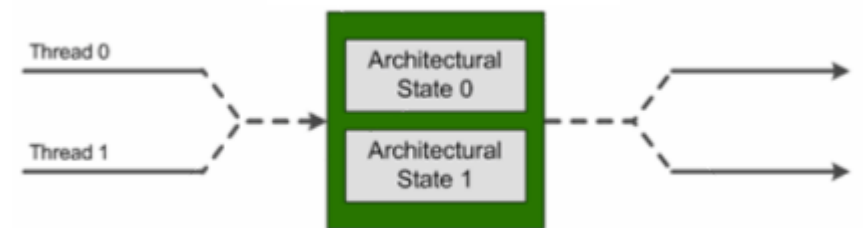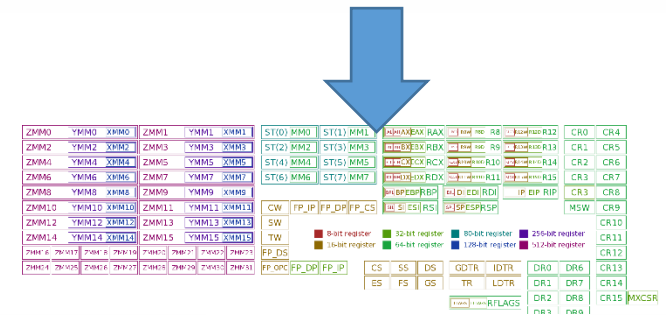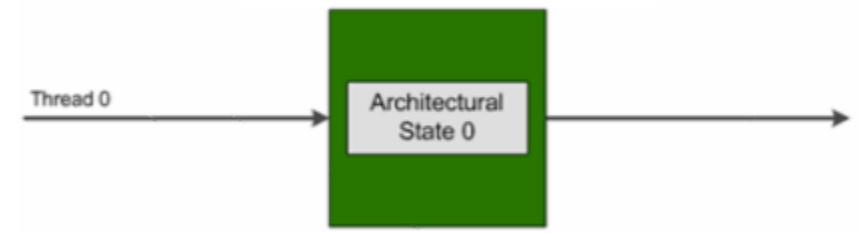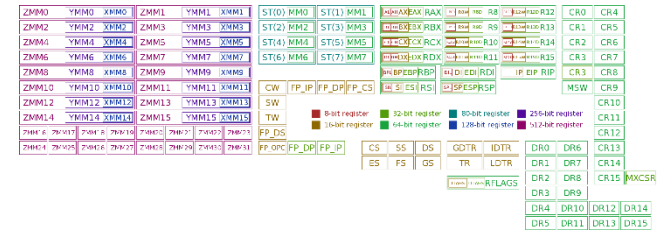
- Three variants:
  - Coarse
  - Fine-grain
  - Simultaneous

# Running example

**Thread A**    **Thread B**    **Thread C**    **Thread D**



- Colors → pipeline full
- White → stall

# Coarse- grained multithreading

# Coarse- grained multithreading

- Single thread runs until a costly stall
  - E.g. 2nd level cache miss

# Coarse- grained multithreading

- Single thread runs until a costly stall
  - E.g. 2nd level cache miss
- Another thread starts during stall
  - Pipeline fill time requires several cycles!

# Coarse- grained multithreading

- Single thread runs until a costly stall
  - E.g. 2nd level cache miss
- Another thread starts during stall
  - Pipeline fill time requires several cycles!

# Coarse- grained multithreading

- Single thread runs until a costly stall
  - E.g. 2nd level cache miss
- Another thread starts during stall
  - Pipeline fill time requires several cycles!
- Hardware support required
  - PC and register file for each thread
  - Looks like another physical CPU to OS/software

# Coarse- grained multithreading

- Single thread runs until a costly stall
  - E.g. 2nd level cache miss
- Another thread starts during stall
  - Pipeline fill time requires several cycles!
- Hardware support required
  - PC and register file for each thread
  - Looks like another physical CPU to OS/software

# Coarse- grained multithreading

- Single thread runs until a costly stall
  - E.g. 2nd level cache miss
- Another thread starts during stall
  - Pipeline fill time requires several cycles!
- Hardware support required
  - PC and register file for each thread
  - Looks like another physical CPU to OS/software

*Pros? Cons?*

# Fine-grained multithreading

# Fine-grained multithreading

- Threads interleave instructions
  - Round-robin
  - Skip stalled threads

# Fine-grained multithreading

- Threads interleave instructions
  - Round-robin
  - Skip stalled threads

# Fine-grained multithreading

- Threads interleave instructions
  - Round-robin
  - Skip stalled threads

- Hardware support required
  - Separate PC and register file per thread
  - Hardware to control alternating pattern

# Fine-grained multithreading

- Threads interleave instructions
  - Round-robin
  - Skip stalled threads

- Hardware support required
  - Separate PC and register file per thread
  - Hardware to control alternating pattern

- Naturally hides delays
  - Data hazards, Cache misses
  - Pipeline runs with rare stalls

# Fine-grained multithreading

- Threads interleave instructions
  - Round-robin
  - Skip stalled threads
- Hardware support required
  - Separate PC and register file per thread
  - Hardware to control alternating pattern
- Naturally hides delays
  - Data hazards, Cache misses
  - Pipeline runs with rare stalls

# Fine-grained multithreading

- Threads interleave instructions
  - Round-robin
  - Skip stalled threads

- Hardware support required
  - Separate PC and register file per thread
  - Hardware to control alternating pattern

- Naturally hides delays
  - Data hazards, Cache misses
  - Pipeline runs with rare stalls

***Pros? Cons?***

# Simultaneous Multithreading (SMT)

# Simultaneous Multithreading (SMT)

- Instructions from multiple threads issued on same cycle
  - Uses register renaming
  - dynamic scheduling facility of multi-issue architecture

# Simultaneous Multithreading (SMT)

- Instructions from multiple threads issued on same cycle
    - Uses register renaming
    - dynamic scheduling facility of multi-issue architecture

Skip C

Skip A

# Simultaneous Multithreading (SMT)

- Instructions from multiple threads issued on same cycle
    - Uses register renaming
    - dynamic scheduling facility of multi-issue architecture
- Hardware support:
    - Register files, PCs per thread
    - Temporary result registers pre commit
    - Support to sort out which threads get results from which instructions

Skip C

Skip A

# Simultaneous Multithreading (SMT)

- Instructions from multiple threads issued on same cycle
  - Uses register renaming
  - dynamic scheduling facility of multi-issue architecture

- Hardware support:
  - Register files, PCs per thread
  - Temporary result registers pre commit
  - Support to sort out which threads get results from which instructions

Skip C

Skip A

# Simultaneous Multithreading (SMT)

- Instructions from multiple threads issued on same cycle
  - Uses register renaming
  - dynamic scheduling facility of multi-issue architecture
- Hardware support:
  - Register files, PCs per thread
  - Temporary result registers pre commit
  - Support to sort out which threads get results from which instructions

**Pros? Cons?**

Skip C

Skip A

# Why Vector and Multithreading Background?

# Why Vector and Multithreading Background?

GPU:

- A very wide vector machine

- Massively multi-threaded to hide memory latency

- *Originally designed for graphics pipelines…*

# Graphics ~= Rendering

# Graphics ~= Rendering

Inputs

# Graphics ~= Rendering

Inputs

- 3D world model(objects, materials)
  - Geometry modeled w triangle meshes, surface normals
  - GPUs subdivide triangles into "fragments" (rasterization)
  - Materials modeled with "textures"
  - Texture coordinates, sampling "map" textures $\rightarrow$ geometry

# Graphics ~= Rendering

Inputs

- 3D world model(objects, materials)
  - Geometry modeled w triangle meshes, surface normals
  - GPUs subdivide triangles into "fragments" (rasterization)
  - Materials modeled with "textures"
  - Texture coordinates, sampling "map" textures → geometry
- Light locations and properties
  - Attempt to model surtface/light interactions with modeled objects/materials

# Graphics ~= Rendering

Inputs

- 3D world model(objects, materials)
  - Geometry modeled w triangle meshes, surface normals
  - GPUs subdivide triangles into "fragments" (rasterization)
  - Materials modeled with "textures"
  - Texture coordinates, sampling "map" textures $\rightarrow$ geometry
- Light locations and properties
  - Attempt to model surtface/light interactions with modeled objects/materials
- View point

# Graphics ~= Rendering

Inputs

- 3D world model(objects, materials)
  - Geometry modeled w triangle meshes, surface normals
  - GPUs subdivide triangles into "fragments" (rasterization)
  - Materials modeled with "textures"
  - Texture coordinates, sampling "map" textures → geometry
- Light locations and properties
  - Attempt to model surtface/light interactions with modeled objects/materials
- View point

Output

# Graphics ~= Rendering

## Inputs

- 3D world model(objects, materials)
  - Geometry modeled w triangle meshes, surface normals
  - GPUs subdivide triangles into "fragments" (rasterization)
  - Materials modeled with "textures"
  - Texture coordinates, sampling "map" textures $\rightarrow$ geometry
- Light locations and properties
  - Attempt to model surtface/light interactions with modeled objects/materials
- View point

## Output

- 2D projection seen from the view-point

# Graphics ~= Rendering

## Inputs

- 3D world model(objects, materials)
  - Geometry modeled w triangle meshes, surface normals
  - GPUs subdivide triangles into "fragments" (rasterization)
  - Materials modeled with "textures"
  - Texture coordinates, sampling "map" textures → geometry
- Light locations and properties
  - Attempt to model surtface/light interactions with modeled objects/materials
- View point

## Output

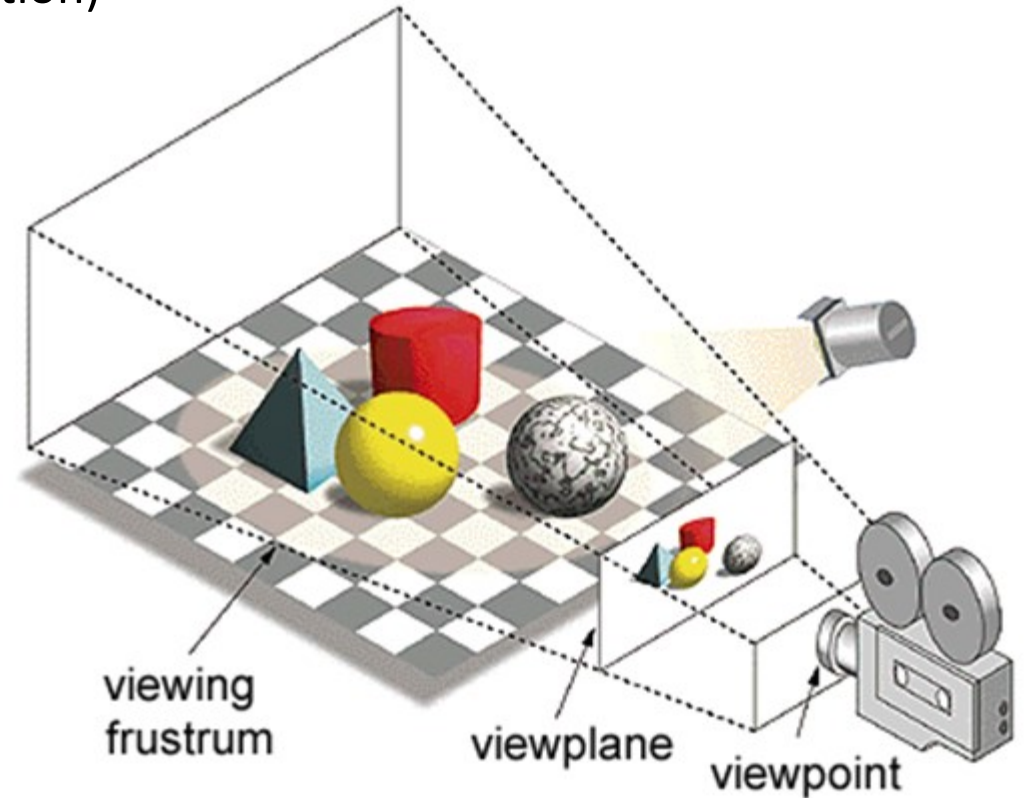- 2D projection seen from the view-point

viewing
frustrum

viewplane

viewpoint

# Grossly over-simplified rendering algorithm

# Grossly over-simplified rendering algorithm

foreach(vertex v in model)

# Grossly over-simplified rendering algorithm

foreach(vertex v in model)

$\qquad$ map $v_{model} \rightarrow v_{view}$

# Grossly over-simplified rendering algorithm

foreach(vertex v in model)

map $v_{model} \rightarrow v_{view}$

# Grossly over-simplified rendering algorithm



foreach(vertex v in model)

      map $v_{model} \rightarrow v_{view}$

fragment[] frags = {};

# Grossly over-simplified rendering algorithm



foreach(vertex v in model)

       map $v_{model}$ → $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0$, $v_1$, $v_2$)

# Grossly over-simplified rendering algorithm



foreach(vertex v in model)

      map $v_{model}$ $\rightarrow$ $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0$, $v_1$, $v_2$)

      frags.add(rasterize(t));

# Grossly over-simplified rendering algorithm



foreach(vertex v in model)

      map $v_{model}$ $\rightarrow$ $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0$, $v_1$, $v_2$)

      frags.add(rasterize(t));

foreach fragment f in frags

# Grossly over-simplified rendering algorithm



foreach(vertex v in model)

      map $v_{model}$ $\rightarrow$ $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0,$ $v_1,$ $v_2$)

      frags.add(rasterize(t));

foreach fragment f in frags

      choose_color(f);

# Grossly over-simplified rendering algorithm

foreach(vertex v in model)

      map $v_{model}$ → $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0$, $v_1$, $v_2$)

      frags.add(rasterize(t));

foreach fragment f in frags

      choose_color(f);

display(visible_fragments(frags));

# Grossly over-simplified rendering algorithm

foreach(vertex v in model)

    map $v_{model}$ → $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0$, $v_1$, $v_2$)

    frags.add(rasterize(t));
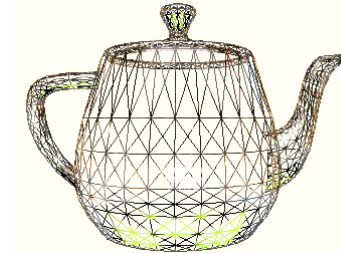
foreach fragment f in frags

    choose_color(f);

display(visible_fragments(frags));

# Algorithm ➔ Graphics Pipeline

foreach(vertex v in model)

map $v_{model}$ ➔ $v_{view}$

fragment[] frags = {};

foreach triangle t $(v_0, v_1, v_2)$

frags.add(rasterize(t));

foreach fragment f in frags

choose_color(f);

display(visible_fragments(frags));

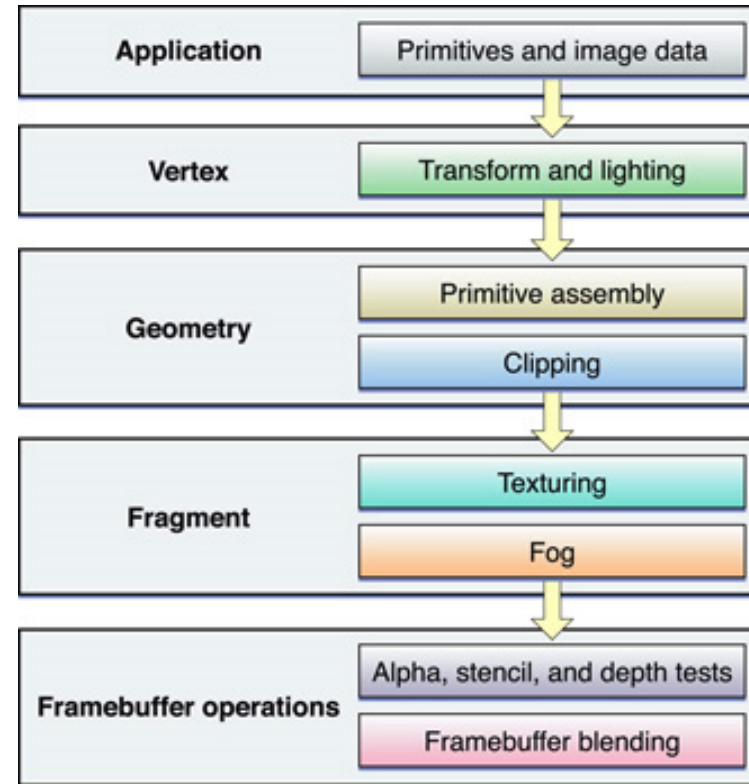

OpenGL pipeline

To first order, DirectX looks the same!

# Algorithm → Graphics Pipeline

foreach(vertex v in model)

  map $v_{model}$ → $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0$, $v_1$, $v_2$)

  frags.add(rasterize(t));

foreach fragment f in frags

  choose_color(f);

display(visible_fragments(frags));



OpenGL pipeline

To first order, DirectX looks the same!

# Algorithm → Graphics Pipeline

foreach(vertex v in model)

    map $v_{model}$ → $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0$, $v_1$, $v_2$)

    frags.add(rasterize(t));

foreach fragment f in frags
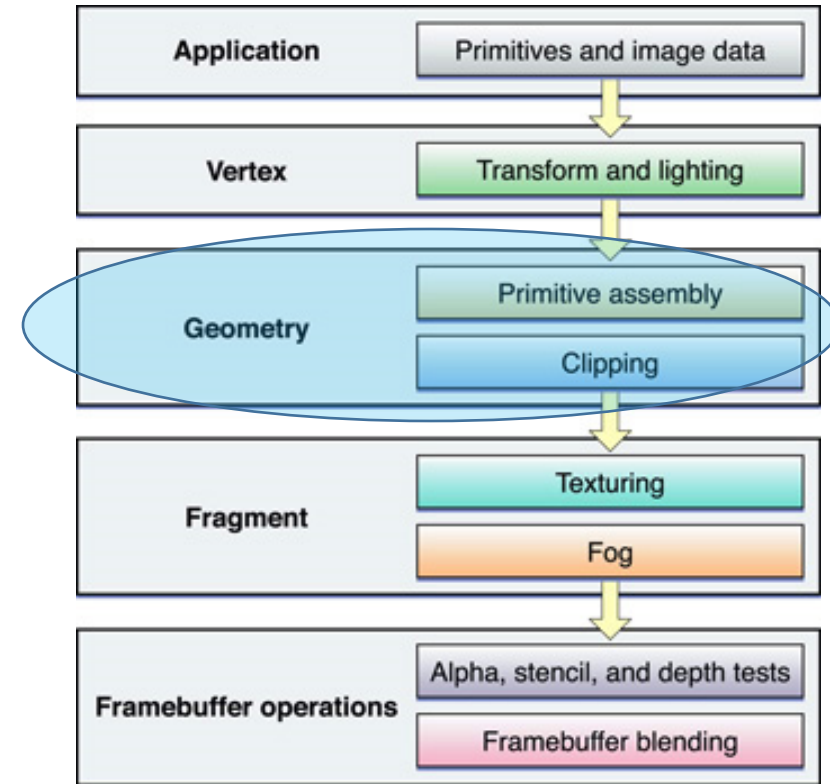
    choose_color(f);

display(visible_fragments(frags));



OpenGL pipeline

To first order, DirectX looks the same!

# Algorithm → Graphics Pipeline

foreach(vertex v in model)

   map $v_{model}$ → $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0$, $v_1$, $v_2$)

   frags.add(rasterize(t));

foreach fragment f in frags

   choose_color(f);

display(visible_fragments(frags));



OpenGL pipeline

To first order, DirectX looks the same!

# Algorithm → Graphics Pipeline

foreach(vertex v in model)

    map $v_{model}$ → $v_{view}$

fragment[] frags = {};

foreach triangle t ($v_0$, $v_1$, $v_2$)

    frags.add(rasterize(t));

foreach fragment f in frags

    choose_color(f);
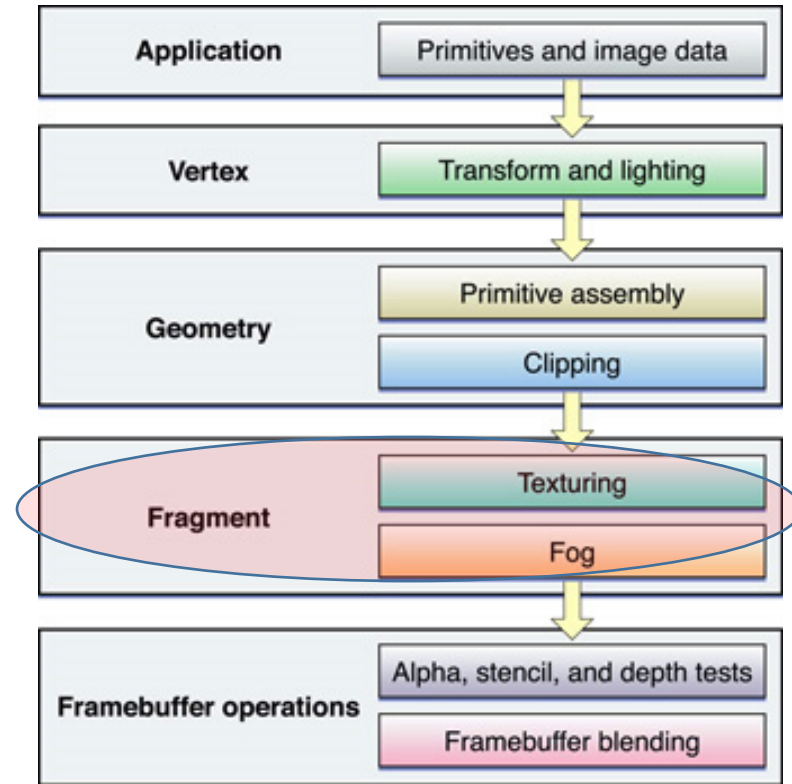
display(visible_fragments(frags));
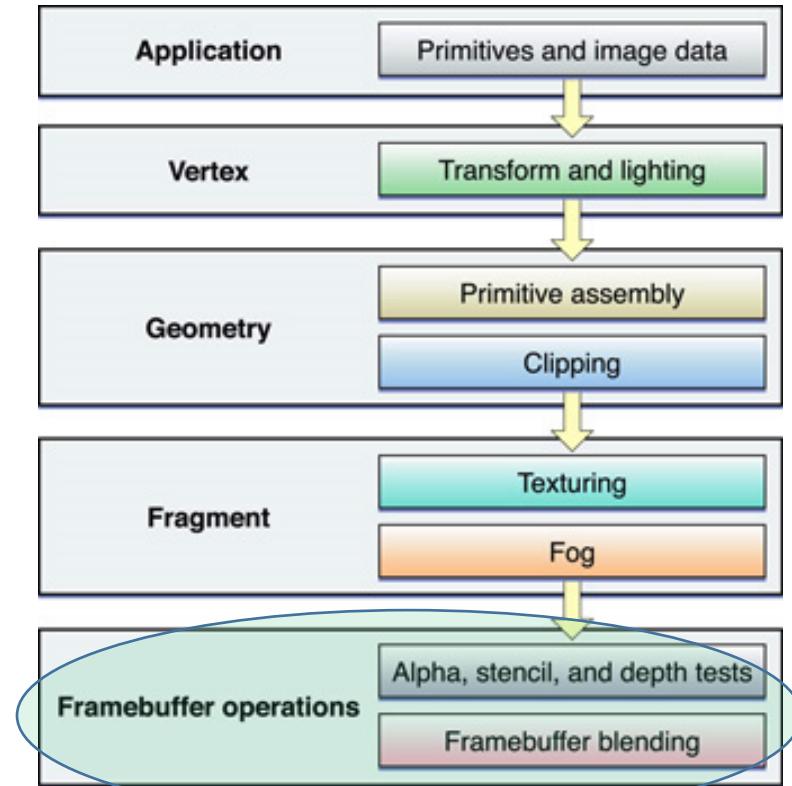


OpenGL pipeline

To first order, DirectX looks the same!

# Graphics pipeline → GPU architecture
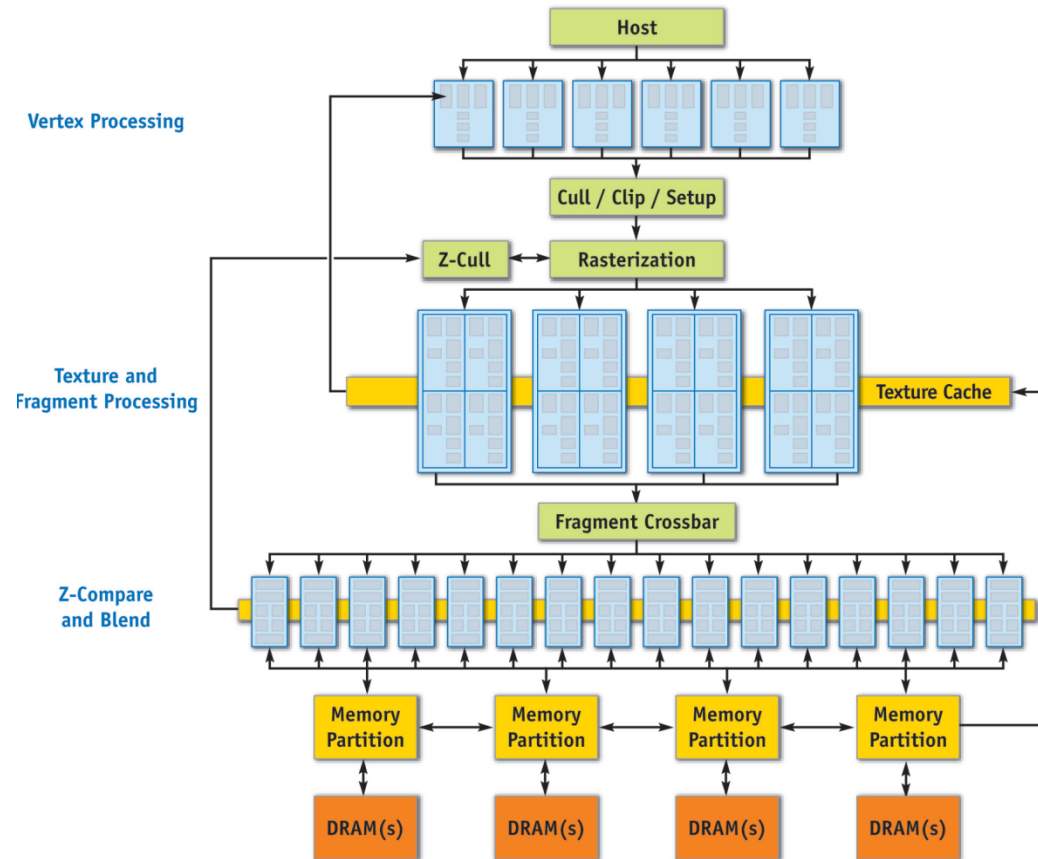




GeForce 6 series

**Limited "programmability" of shaders:**
**Minimal/no control flow**
**Maximum instruction count**

# Graphics pipeline → GPU architecture



**Limited "programmability" of shaders:**
**Minimal/no control flow**
**Maximum instruction count**

GeForce 6 series

# Graphics pipeline → GPU architecture



**Limited "programmability" of shaders:**
**Minimal/no control flow**
**Maximum instruction count**

GeForce 6 series

# Graphics pipeline ➔ GPU architecture



**Limited "programmability" of shaders:**
**Minimal/no control flow**
**Maximum instruction count**
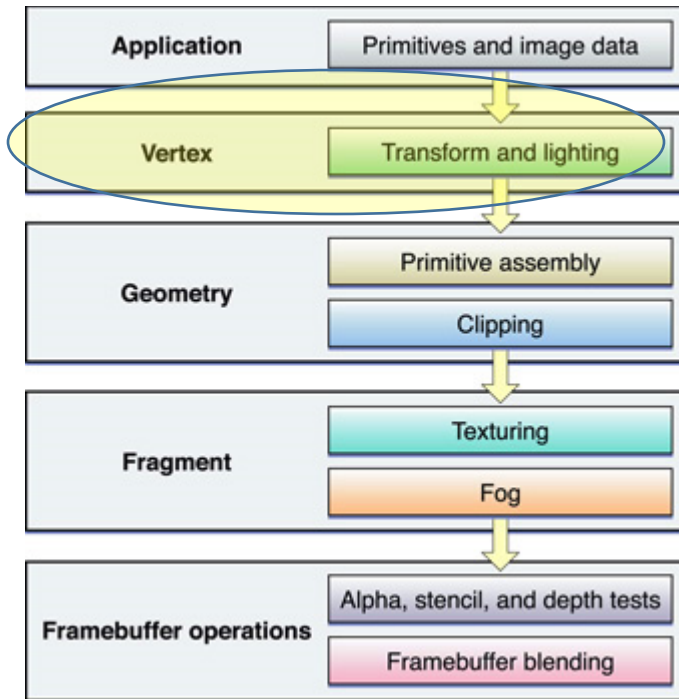
GeForce 6 series

Dandelion

# Graphics pipeline → GPU architecture



GeForce 6 series

**Limited "programmability" of shaders:**
**Minimal/no control flow**
**Maximum instruction count**

# Late Modernity: unified shaders



Mapping to Graphics pipeline no longer apparent

Processing elements no longer specialized to a particular role

Model supports *real* control flow, larger instr count

Dandelion

# Mostly Modern: Pascal

# Definitely Modern: Turing



Turing SM
16 TFLOPS + 16 TIPS
Concurrent FP & INT Execution
Unified L1 Cache
Variable Rate Shading

Display
Native HDR
8K DisplayPort
VirtualLink

RT Core
10 Giga Rays/sec
Ray Triangle Intersection
BVH Traversal

NVLINK
100 GB/sec
GPU-GPU Memory Access

Tensor Core
125 TFLOPS FP16
250 TOPS INT8
500 TOPS INT4

Video
HEVC 8K Real Time Encode
25% Improved Bitrate

Memory
6MB L2 Cache
384-bit G6 @ 14Gbps
672 GB/sec

# Cross-generational GPU observations

GPUs designed for parallelism in graphics pipeline:

- Data
  - Per-vertex
  - Per-fragment
  - Per-pixel

- Task
  - Vertex processing
  - Fragment processing
  - Rasterization
  - Hidden-surface elimination

- MLP
  - HW multi-threading for hiding memory latency

- Simple cores

- Single instruction stream
  - Vector instructions (SIMD) OR
  - Implicit HW-managed sharing (SIMT)

- Hide memory latency with HW multi-threading

# Cross-generational GPU observations

GPUs designed for parallelism in graphics pipeline:

- Data
    - Per-vertex
    - Per-fragment
    - Per-pixel
- Task
    - Vertex processing
    - Fragment processing
    - Rasterization
    - Hidden-surface elimination
- MLP
    - HW multi-threading for hiding memory latency

- Simple cores
- Single instruction stream

OR

ing

n

Even as GPU architectures become more general, certain assumptions persist:
1. Data parallelism is *trivially* exposed
2. **All** problems look like painting a box with colored dots

# Cross-generational GPU observations

GPUs designed for parallelism in graphics pipeline:

- Data
  - Per-vertex
  - Per-fragment
  - Per-pixel

- Task
  - Vertex processing
  - Fragment processing
  - Rasterization
  - Hidden-surface elimination

- MLP
  - HW multi-threading for hiding memory latency

- Simple cores
- Single instruction stream

OR
...ing

Even as GPU architectures become more general, certain assumptions persist:
1. Data parallelism is *trivially* exposed
2. **All** problems look like painting a box with colored dots

*But what if my problem isn't painting a box?!!?!*

# Programming Model

- ***GPUs are I/O devices, managed by user-code***

- "kernels" == "shader programs"

- 1000s of HW-scheduled threads per kernel

- Threads grouped into independent blocks.
    - Threads in a block can synchronize (barrier)
    - This is the *only* synchronization

- "Grid" == "launch" == "invocation" of a kernel
    - a group of blocks (or warps)

# Programming Model

- ***GPUs are I/O devices, managed by user-code***

- "kernels" == "shader programs"

- 1000s of HW-scheduled threads per kernel

- Threads grouped into independent blocks.
  - Threads in a block can synchronize (barrier)
  - This is the *only* synchronization

- "Grid" == "launch" == "invocation" of a kernel
  - a group of blocks (or warps)

*Need codes that are 1000s-X parallel....*

# Parallel Algorithms

- Sequential algorithms often do not permit easy parallelization
  - Does not mean there work has no parallelism
  - A different approach can yield parallelism
  - but often changes the algorithm
  - Parallelizing != just adding locks to a sequential algorithm
- Parallel Patterns
  - Map
  - Scatter, Gather
  - Reduction
  - Scan
  - Search, Sort

# Parallel Algorithms

- Sequential algorithms often do not permit easy parallelization
  - Does not mean there work has no parallelism
  - A different approach can yield parallelism
  - but often changes the algorithm
  - Parallelizing != just adding locks to a sequential algorithm
- Parallel Patterns
  - Map
  - Scatter, Gather
  - Reduction
  - Scan
  - Search, Sort

If you can express your algorithm using these patterns, an apparently fundamentally sequential algorithm can be made parallel

# Map

- Inputs
  - Array A
  - Function f(x)
- map(A, f) → apply f(x) on all elements in A
- Parallelism trivially exposed
  - f(x) can be applied in parallel to all elements, in principle

# Map

- Inputs
  - Array A
  - Function f(x)
- map(A, f) → apply f(x) on all elements in A
- Parallelism trivially exposed
  - f(x) can be applied in parallel to all elements, in principle

```
for(i=0; i<numPoints; i++) {
    labels[i] = findNearestCenter(points[i]);
}
```

```
map(points, findNearestCenter)
```

# Scatter and Gather

# Scatter and Gather

- Gather:
  - Read multiple items to single /packed location

# Scatter and Gather

- Gather:
  - Read multiple items to single /packed location
- Scatter:
  - Write single/packed data item to multiple locations

# Scatter and Gather

- Gather:
  - Read multiple items to single /packed location
- Scatter:
  - Write single/packed data item to multiple locations
- Inputs: x, y, indeces, N

# Scatter and Gather

- Gather:
  - Read multiple items to single /packed location
- Scatter:
  - Write single/packed data item to multiple locations
- Inputs: x, y, indeces, N

```
for (i=0; i<N; ++i)
  x[i] = y[idx[i]];
```
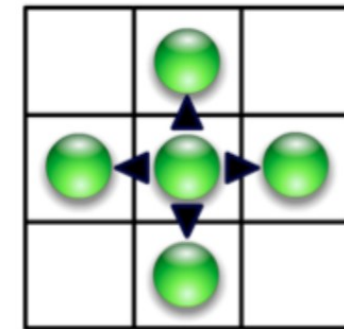→ gather(x, y, idx)

```
for (i=0; i<N; ++i)
  y[idx[i]] = x[i];
```
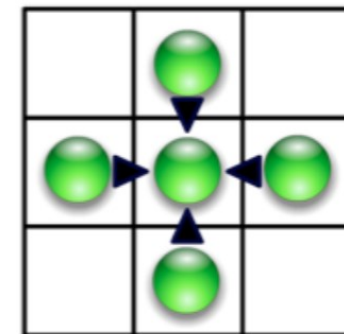→ scatter(x, y, idx)

# Scatter and Gather

- Gather:
  - Read multiple items to single /packed location
- Scatter:
  - Write single/packed data item to multiple locations
- Inputs: x, y, indeces, N

```
for (i=0; i<N; ++i)
  x[i] = y[idx[i]];
```
→ gather(x, y, idx)

```
for (i=0; i<N; ++i)
  y[idx[i]] = x[i];
```
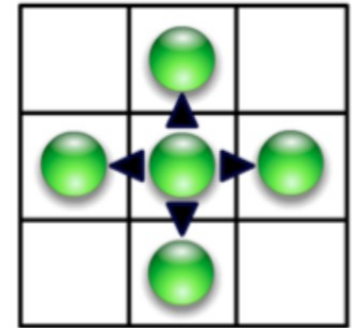→ scatter(x, y, idx)



Scatter



Gather

# Scatter and Gather

- Gather:
  - Read multiple items to single /packed location
- Scatter:
  - Write single/packed data item to multiple locations
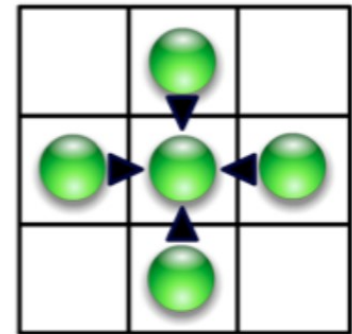- Inputs: x, y, indeces, N

```
for (i=0; i<N; ++i)
  x[i] = y[idx[i]];
```
→ gather(x, y, idx)

```
for (i=0; i<N; ++i)
  y[idx[i]] = x[i];
```
→ scatter(x, y, idx)


Scatter


Gather

24

# Reduce

- Input
  - Associative operator op
  - Ordered set s = [a, b, c, … z]
- Reduce(op, s) returns a op b op c … op z

# Reduce

- Input
  - Associative operator op
  - Ordered set s = [a, b, c, … z]
- Reduce(op, s) returns a op b op c … op z

```
for(i=0; i<N; ++i) {
    accum += (point[i]*point[i])
}
```

accum = reduce(*, point)

# Reduce

- Input
  - Associative operator op
  - Ordered set s = [a, b, c, ... z]
- Reduce(op, s) returns a op b op c ... op z

```
for(i=0; i<N; ++i) {
    accum += (point[i]*point[i])
}
```

accum = reduce(*, point)

Why must op be associative?

# Reduce


$N/2$
$N/4...$
$1$
$N$
O($\log_2 N$) steps, O($N$) work


$MxN/2$
$MxN/4...$
$Mx1$
$MxN$
O($\log_2 N$) steps, O($MN$) work

- Input
  - Associative operator op
  - Ordered set s = [a, b, c, … z]
- Reduce(op, s) returns a op b op c … op z

```
for(i=0; i<N; ++i) {
    accum += (point[i]*point[i])
}
```
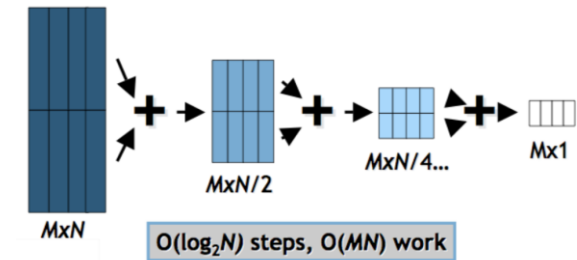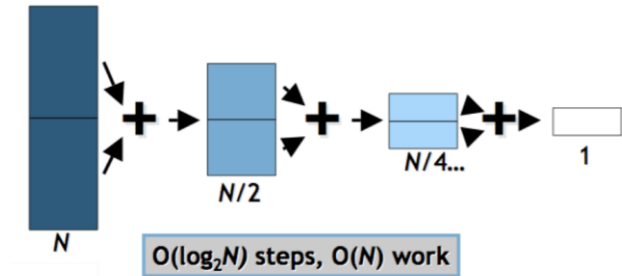
accum = reduce(*, point)
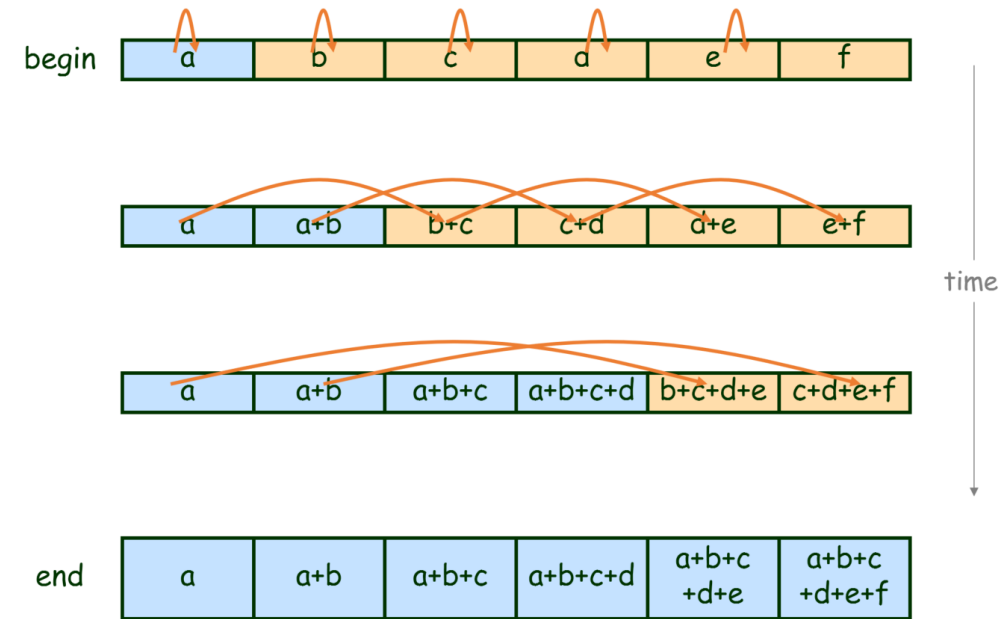
Why must op be associative?

# Scan (prefix sum)

- Input
  - Associative operator op
  - Ordered set s = [a, b, c, … z]
  - Identity I

- scan(op, s) = [I, a, (a op b), (a op b op c) …]

- Scan is the workhorse of parallel algorithms:
  - Sort, histograms, sparse matrix, string compare, …

# GroupBy

- Group a collection by key
- Lambda function maps elements → key

# GroupBy

- Group a collection by key
- Lambda function maps elements → key

```
var res = ints.GroupBy(x => x);
```

# GroupBy

- Group a collection by key
- Lambda function maps elements → key

```
var res = ints.GroupBy(x => x);
```

# GroupBy

- Group a collection by key
- Lambda function maps elements → key
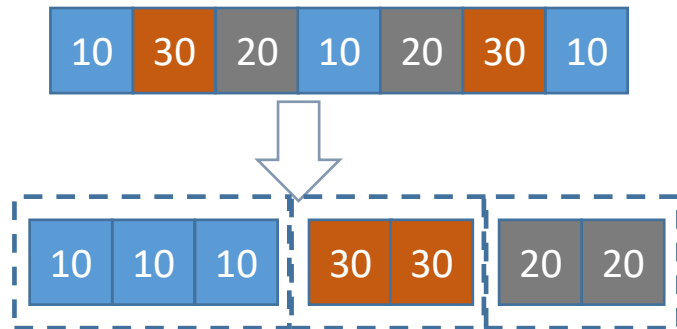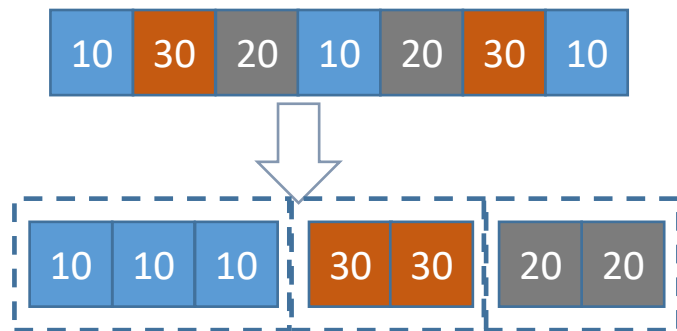
```
var res = ints.GroupBy(x => x);
```

# GroupBy

- Group a collection by key
- Lambda function maps elements → key

```
var res = ints.GroupBy(x => x);
```



```
foreach(T elem in ints)
{
    key    = KeyLambda(elem);

    group = GetGroup(key);

    group.Add(elem);
}
```
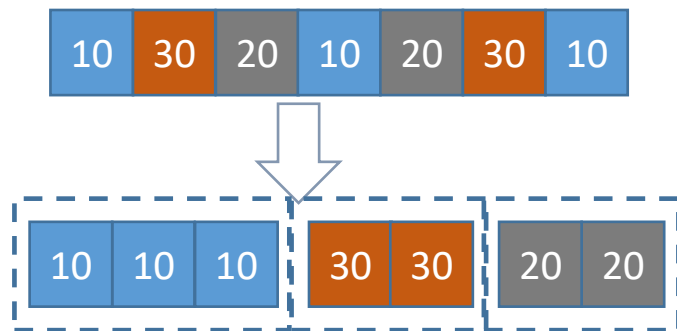
# GroupBy

- Group a collection by key
- Lambda function maps elements → key
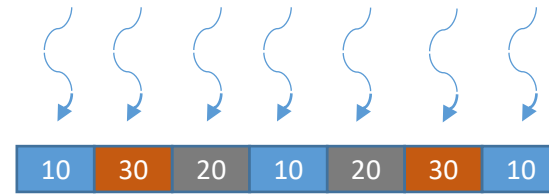
```
var res = ints.GroupBy(x => x);
```



```
foreach(T elem in PF(ints))
{
  key    = KeyLambda(elem);

  group = GetGroup(key);🔒

  group.Add(elem);🔒
}
```

# GroupBy using parallel primitives

| 10 | 30 | 20 | 10 | 20 | 30 | 10 |
|----|----|----|----|----|----|----|

# GroupBy using parallel primitives

# GroupBy using parallel primitives



Assign group IDs

Group ID :

# GroupBy using parallel primitives



| | 10 | 30 | 20 | 10 | 20 | 30 | 10 |

**Assign group IDs**

| | 10 | 20 | 30 |
| --- | --- | --- | --- |
| Group ID : | 0 | 1 | 2 |

**Compute group sizes**

| | 10 | 20 | 30 |
| --- | --- | --- | --- |
| Group ID : | 0 | 1 | 2 |
| Group Size : | 3 | 2 | 2 |

# GroupBy using parallel primitives
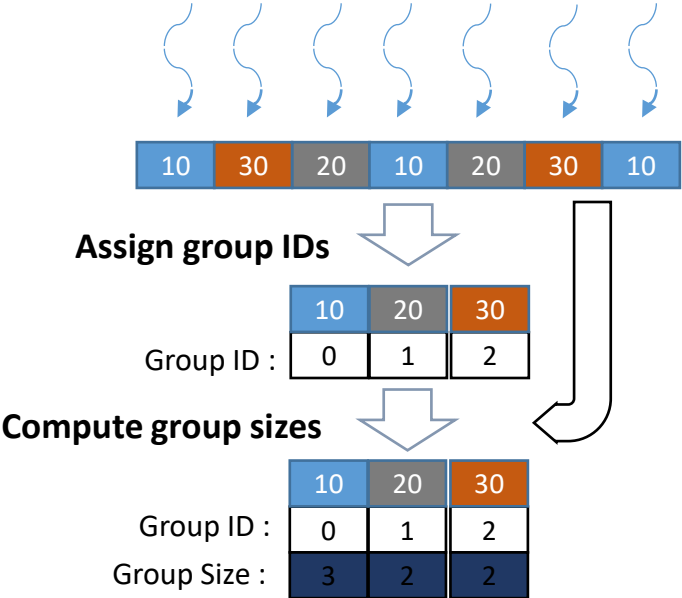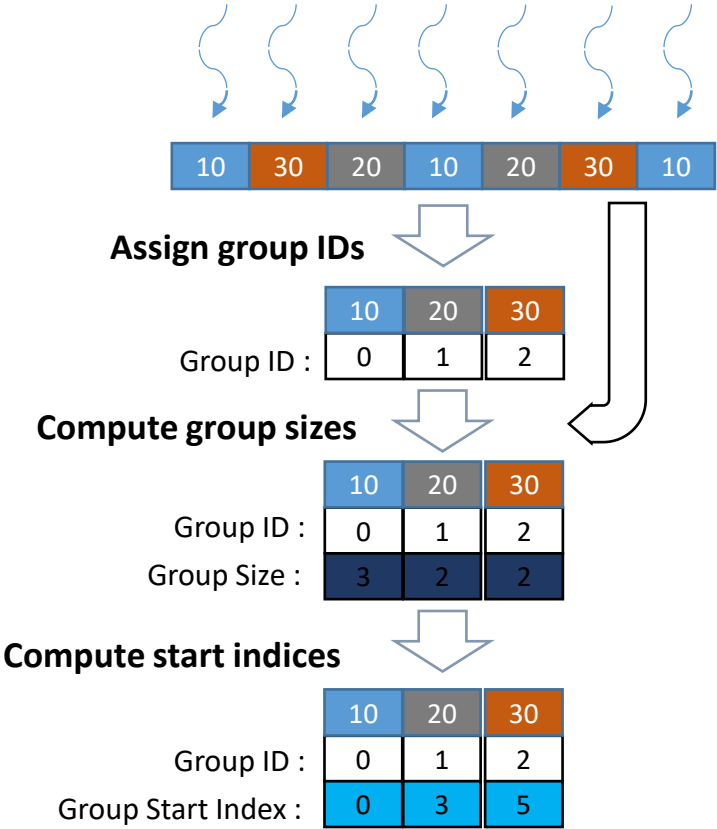
# GroupBy using parallel primitives



| | 10 | 30 | 20 | 10 | 20 | 30 | 10 |

**Assign group IDs**

| | 10 | 20 | 30 |
|---|---|---|---|
| Group ID : | 0 | 1 | 2 |

**Compute group sizes**

| | 10 | 20 | 30 |
|---|---|---|---|
| Group ID : | 0 | 1 | 2 |
| Group Size : | 3 | 2 | 2 |

**Compute start indices**

| | 10 | 20 | 30 |
|---|---|---|---|
| Group ID : | 0 | 1 | 2 |
| Group Start Index : | 0 | 3 | 5 |

**Write Outputs**

| | 10 | 10 | 10 | 20 | 20 | 30 | 30 |

# GroupBy using parallel primitives



| 10 | 30 | 20 | 10 | 20 | 30 | 10 |

**Assign group IDs**

**Sorting** or hashing

| 10 | 20 | 30 |
| Group ID : | 0 | 1 | 2 |

**Compute group sizes**

| 10 | 20 | 30 |
| Group ID : | 0 | 1 | 2 |
| Group Size : | 3 | 2 | 2 |

**Compute start indices**

| 10 | 20 | 30 |
| Group ID : | 0 | 1 | 2 |
| Group Start Index : | 0 | 3 | 5 |

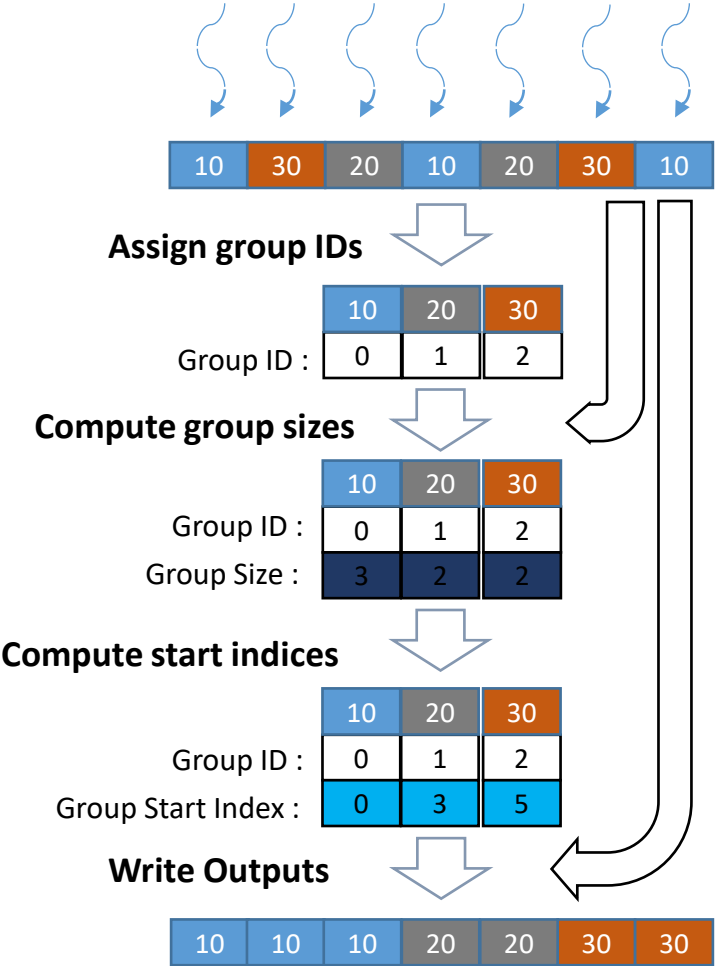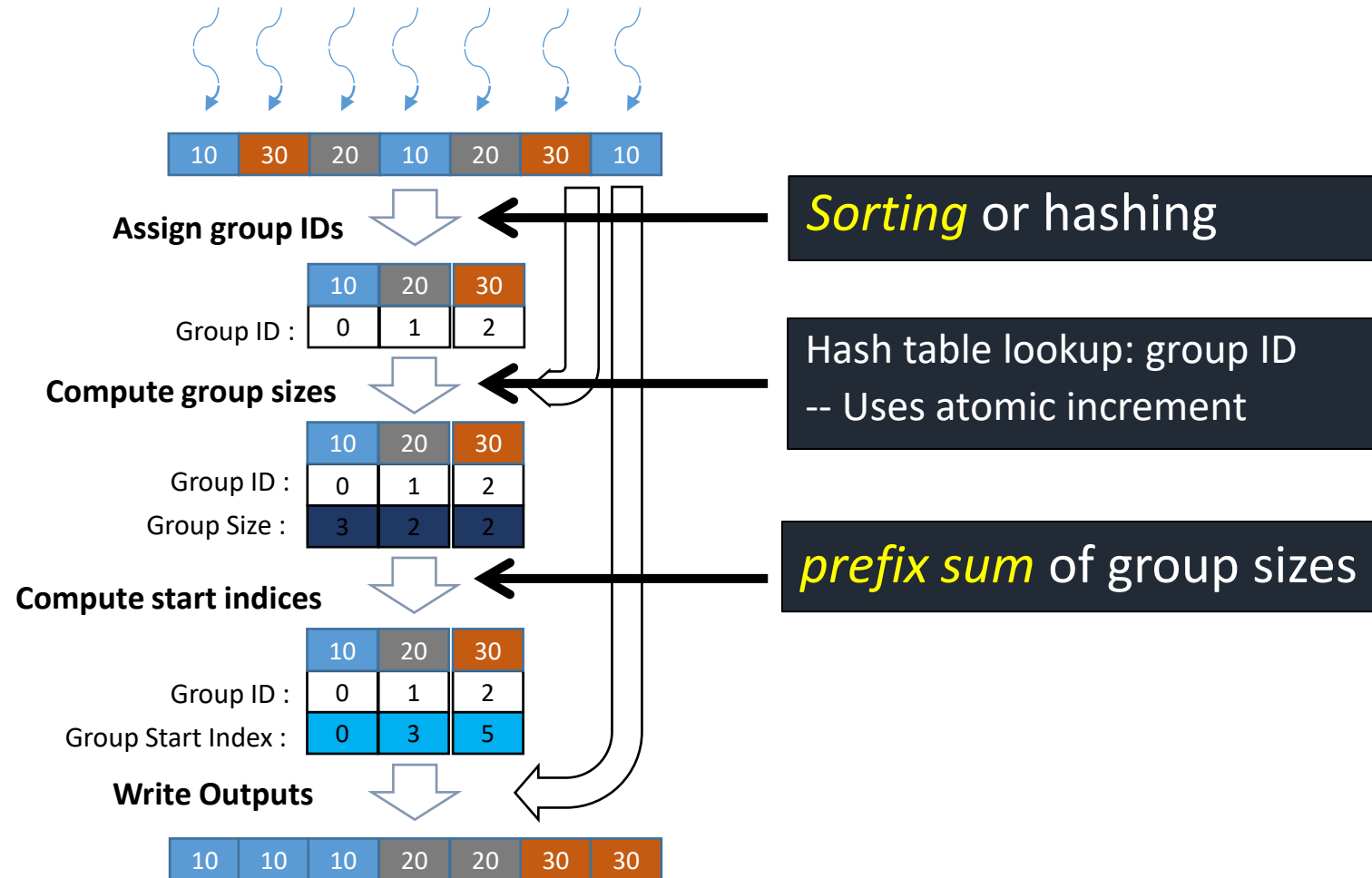**Write Outputs**

| 10 | 10 | 10 | 20 | 20 | 30 | 30 |

# GroupBy using parallel primitives

# GroupBy using parallel primitives

# GroupBy using parallel primitives



*Sorting* or hashing

Hash table lookup: group ID
-- Uses atomic increment

*prefix sum* of group sizes

Write to output location
— Uses atomic increment
— *Scatter gather*