# End-of-semester Review

cs378h
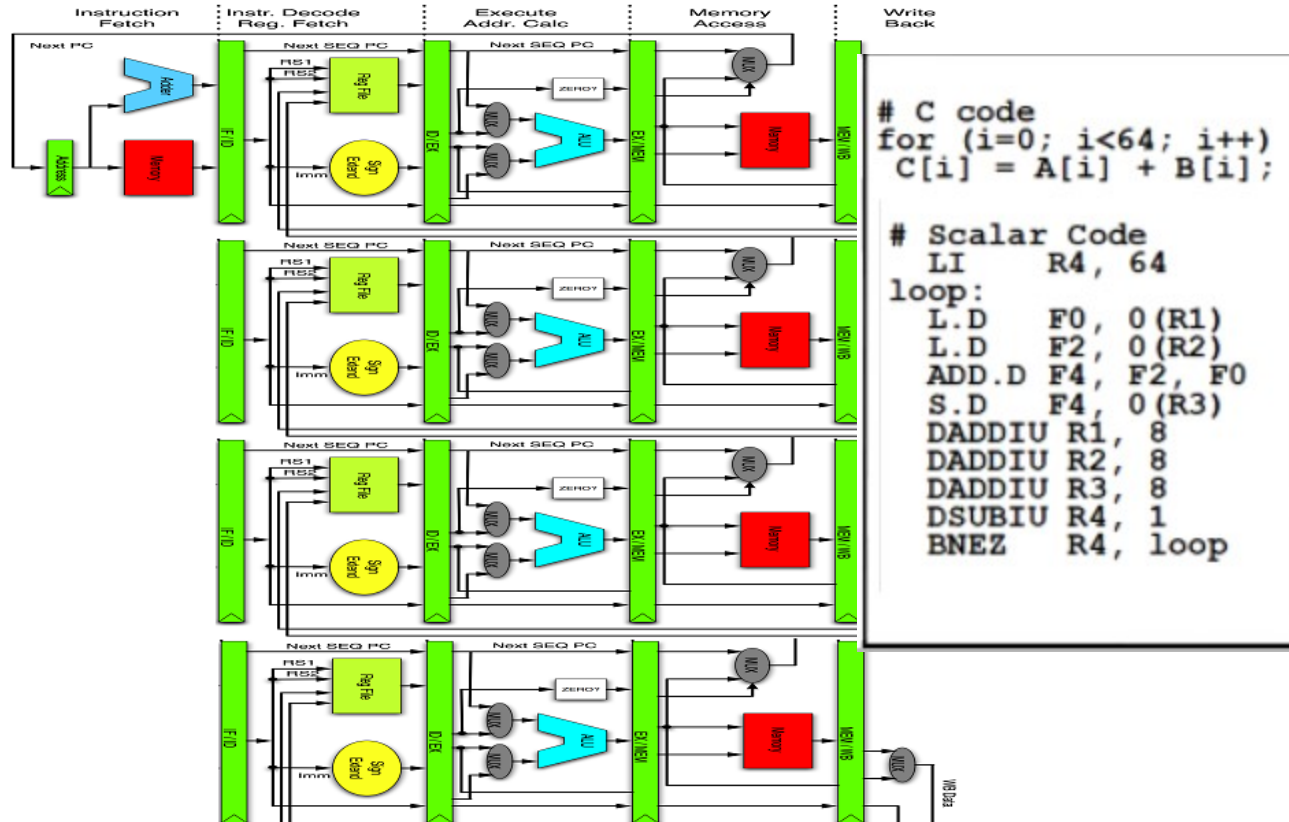
# Outline/Administrivia

- Questions?
- Review
  - Can someone please act as scribe?
  - Requested review content:
    - GPUs: SIMT vs SIMD, schedulers, limitations on threads/blocks and num blocks, divergence, sharing global memory
    - FPGAs/Verilog: CLB, BRAM, and LUT
    - MPI, distributed systems, shared nothing architectures, PGAS
    - Distributed systems (like CAP and NoSQL)
    - Consistency guarantees?
    - Linearizability vs. Serializability

# 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
```
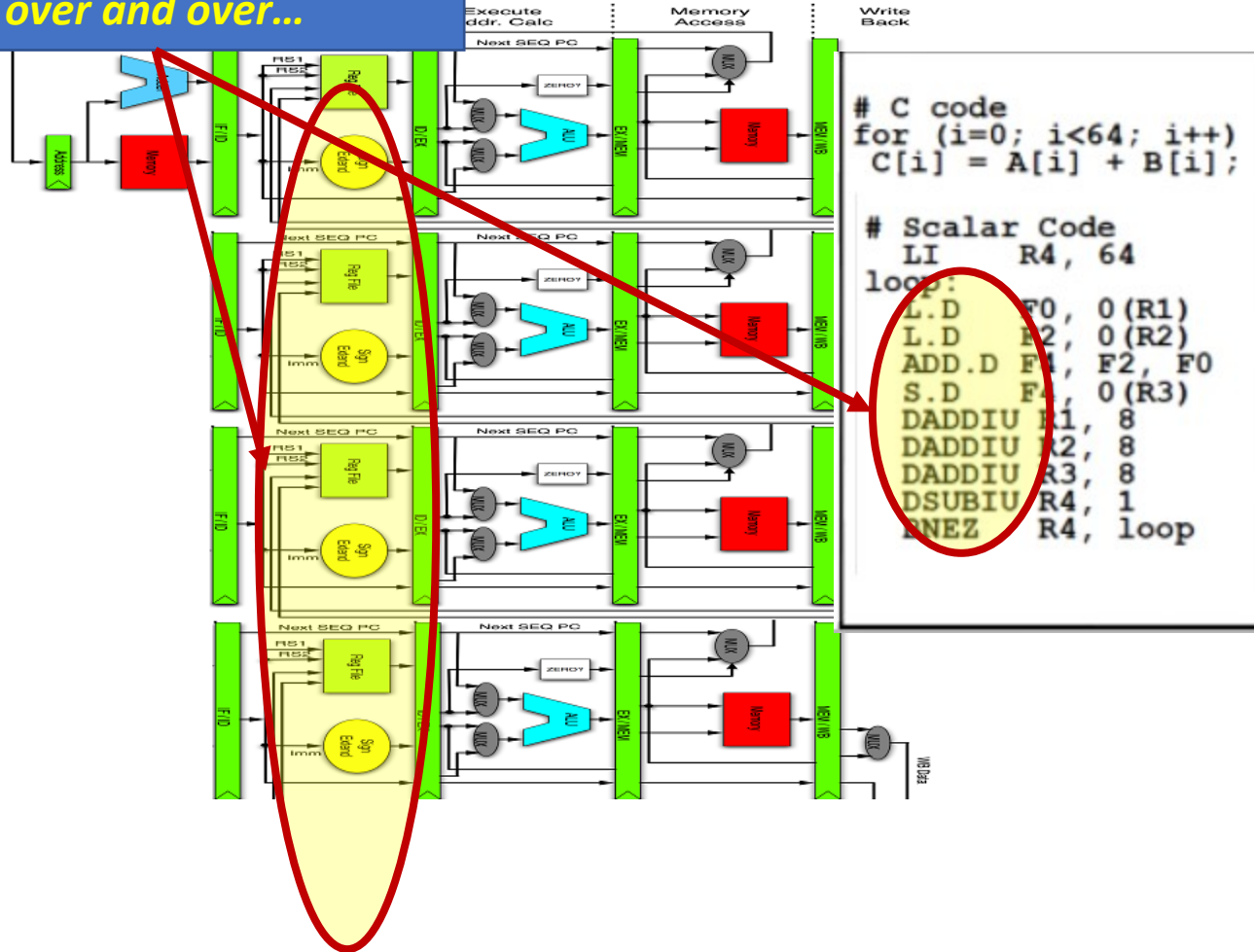
# Review: what is a vector processor?



```
# 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
```
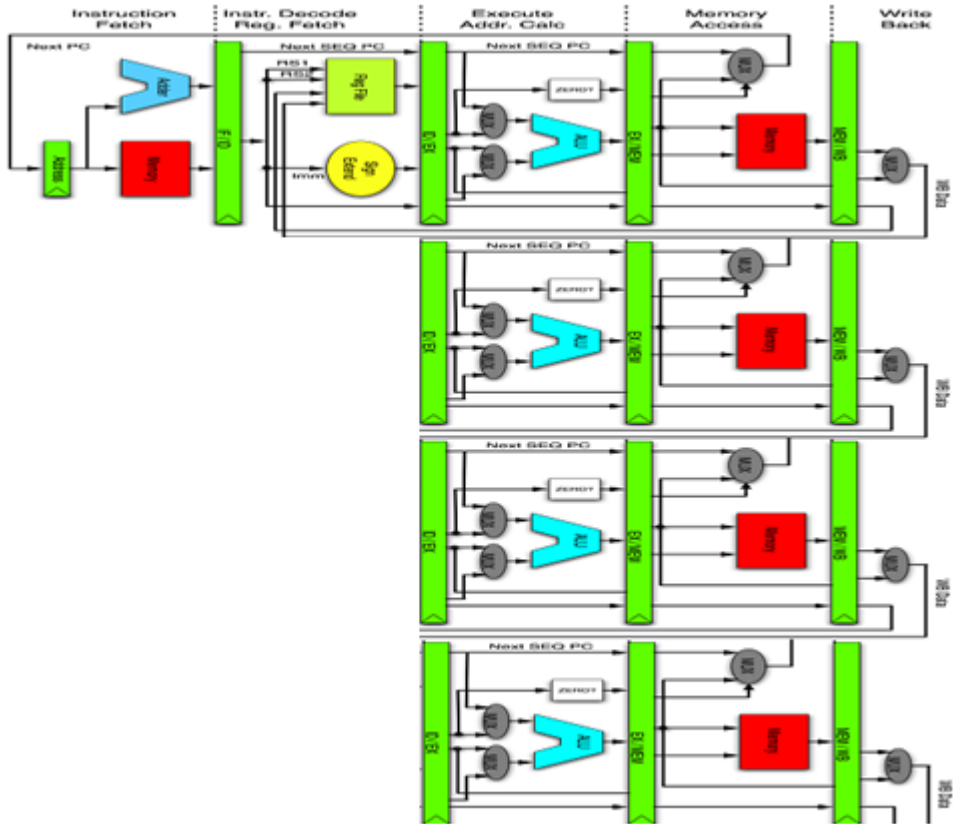
# Review: what is a vector processor?



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
```

3

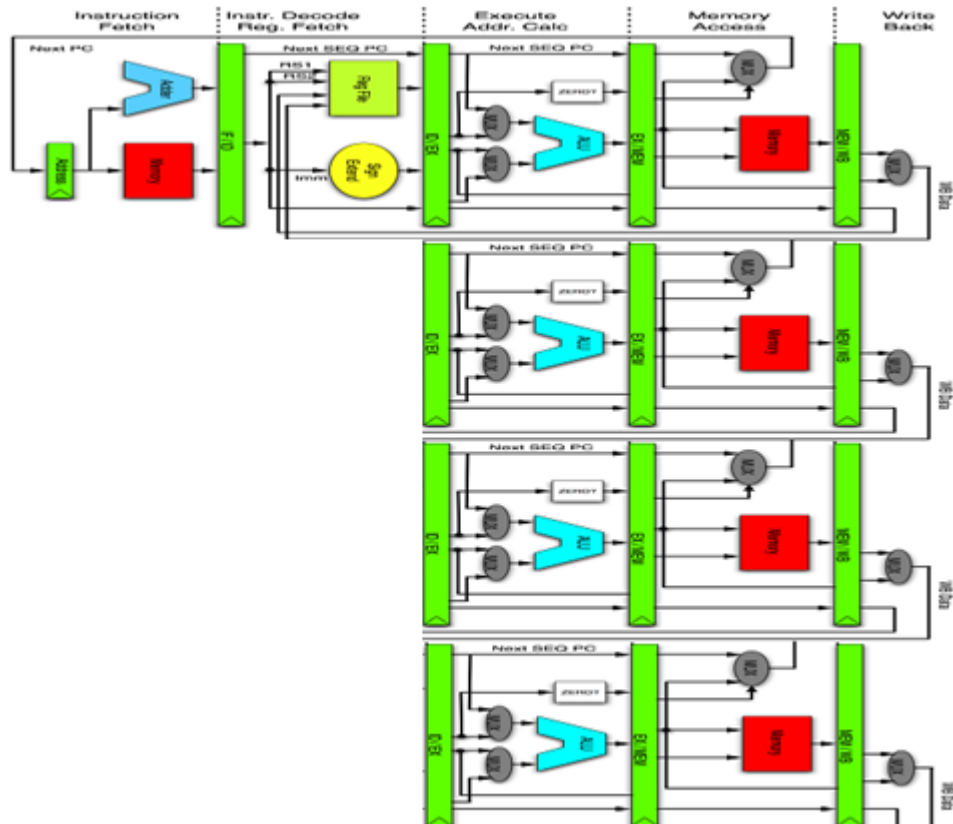# Review: what is a vector processor?

# Review: what is a vector processor?
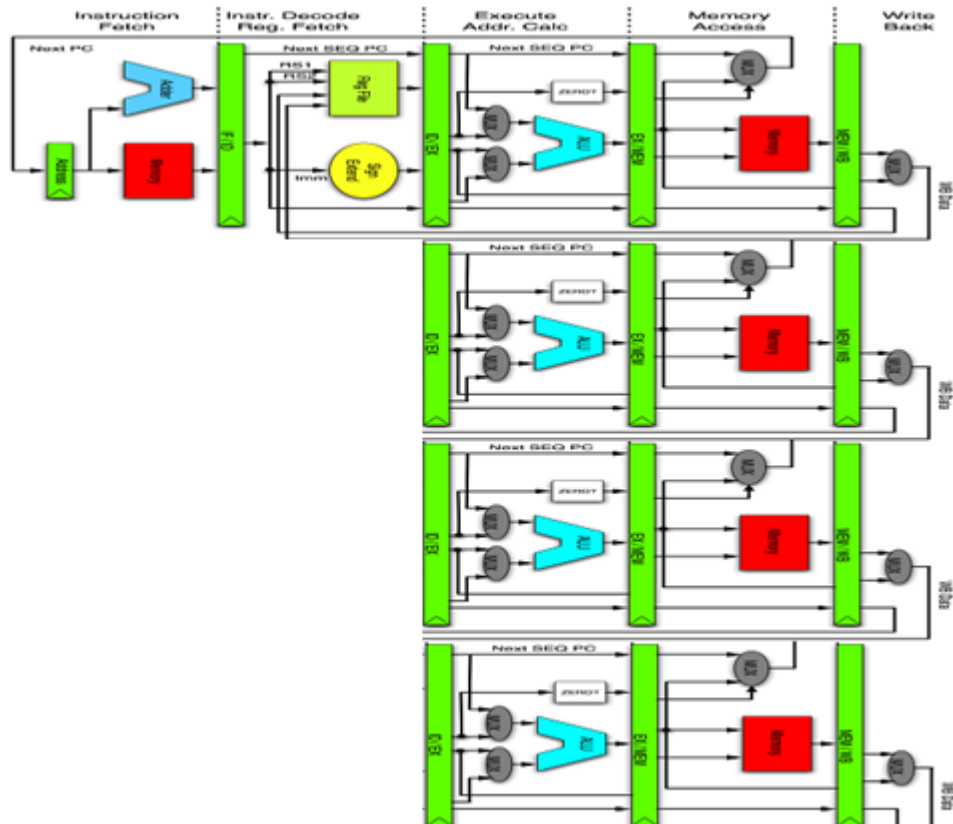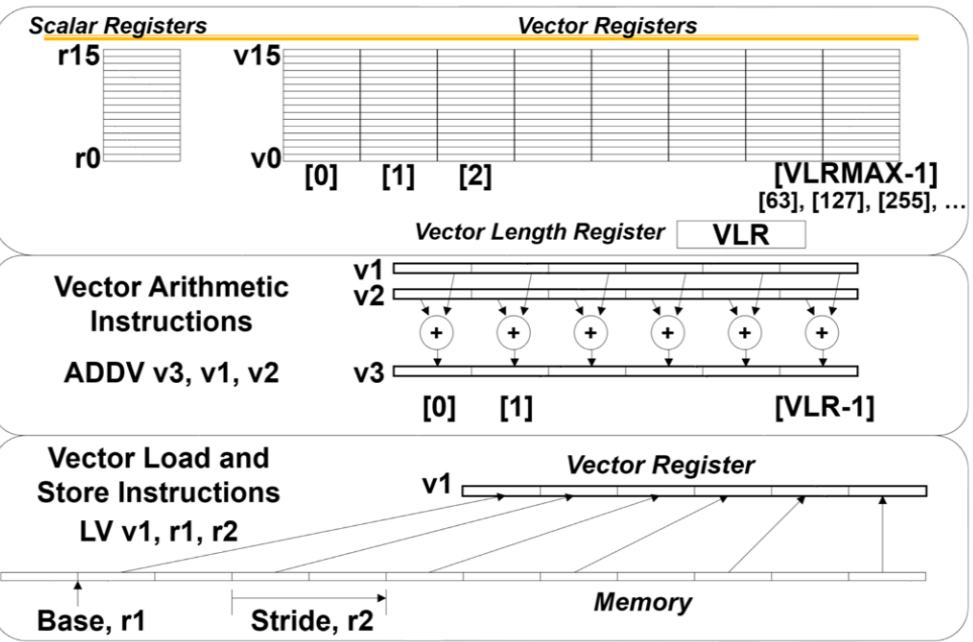


```
# C code
for (i=0; i<64; i++)
  C[i] = A[i] + B[i];
```

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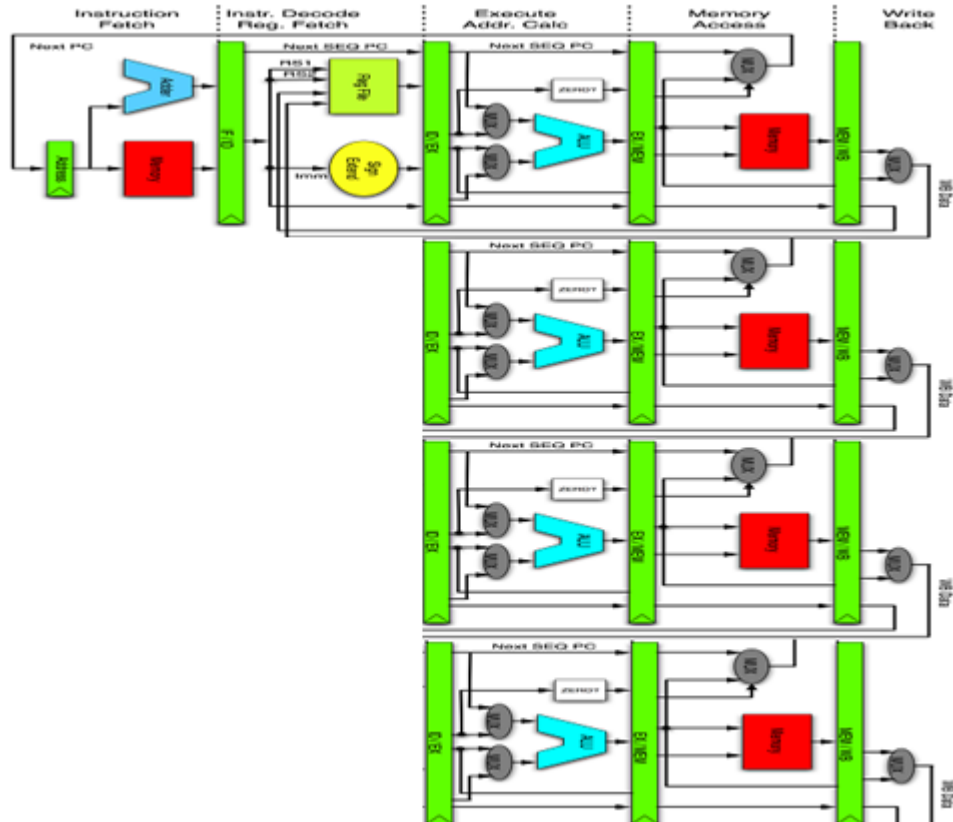
```
# Scalar Code
  LI      R4, 64
loop:
  L.D     F0, 0(R1)
  L.D     F2, 0(R2)
  ADD.D   F4, F2, F0
  S.D     F4, 0(R3)
  DADDIU  R1, 8
  DADDIU  R2, 8
  DADDIU  R3, 8
  DSUBIU  R4, 1
  BNEZ    R4, loop
```

```
# 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...



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

Imp...

- In...
- Sa...
- M...
- 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
```

3

# Hardware multi-threading
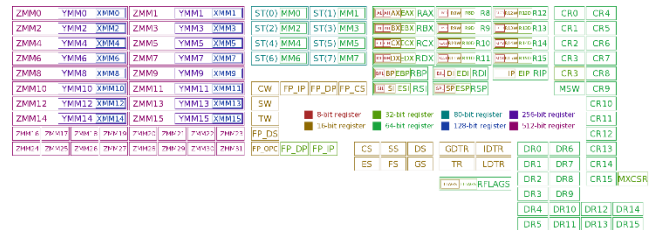
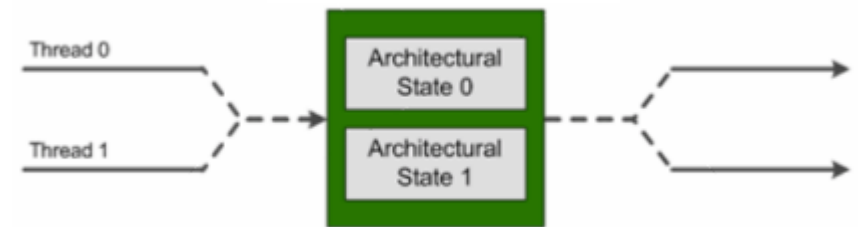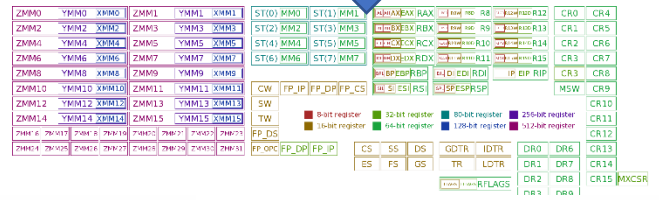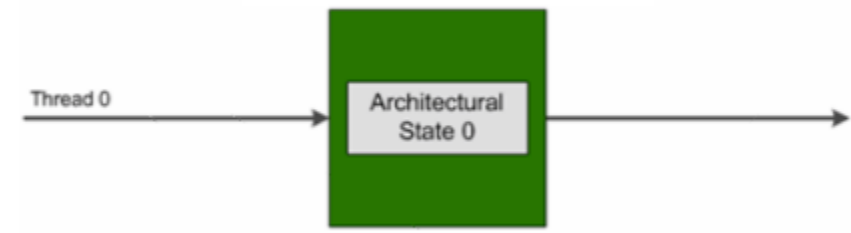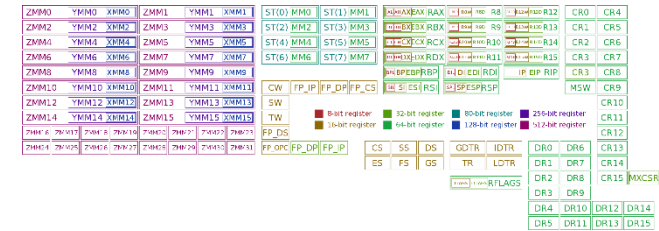# Hardware multi-threading

- 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

# Hardware multi-threading

- Address memory bottleneck

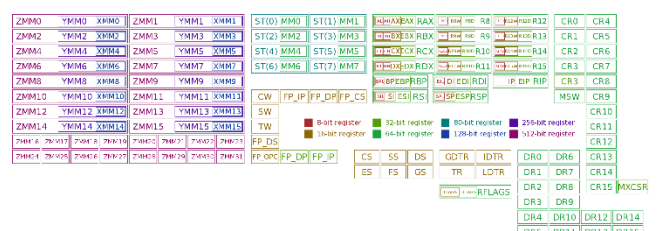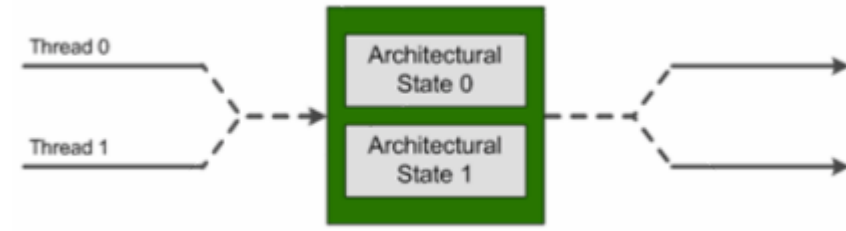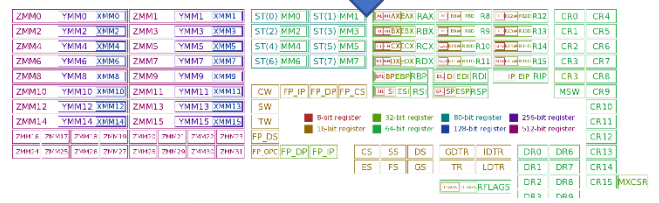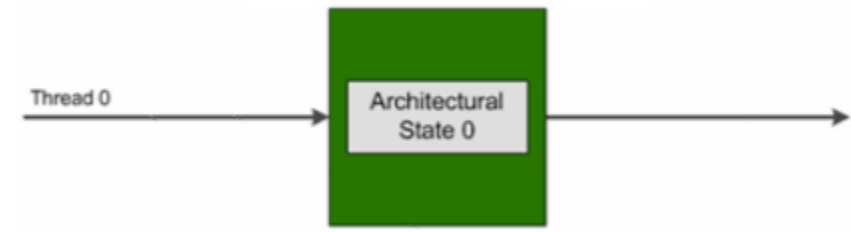- Share exec unit across
  - Instruction streams
  - Switch on stalls

# Hardware multi-threading

- Address memory bottleneck

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

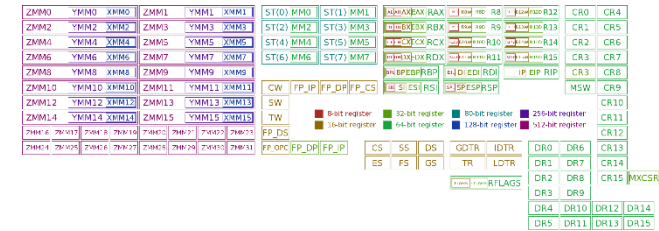- Looks like multiple cores to the OS

# Hardware multi-threading

- Address memory bottleneck

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

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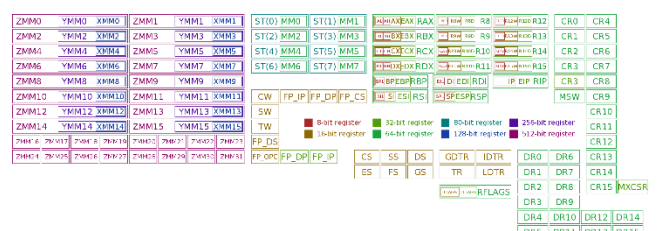
- Three variants:
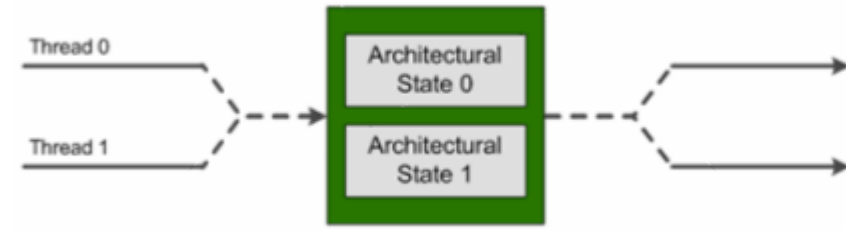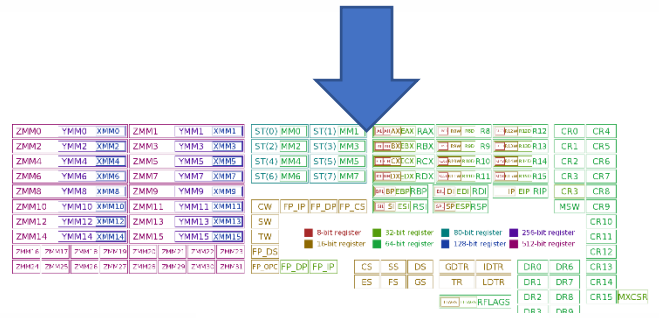  - Coarse
  - Fine-grain
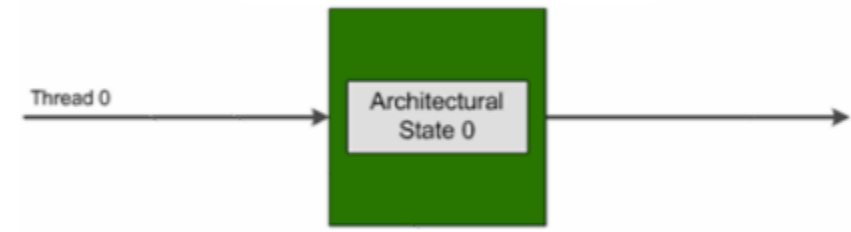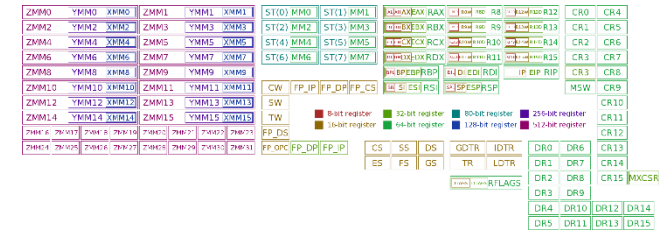  - Simultaneous

# 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

**SIMT = SIMD + Hw MT**

# SIMD vs. SIMT

# Review

# Review



- Each SM has multiple vector units (4)
  - 32 lanes wide → warp size

# Review



- Each SM has multiple vector units (4)
  - 32 lanes wide → warp size
- Vector units use *hardware multi-threading*

# Review



- Each SM has multiple vector units (4)
  - 32 lanes wide → warp size
- Vector units use ***hardware multi-threading***
- Execution → a grid of thread blocks (TBs)
  - Each TB has some number of threads

# Review



- Each SM has multiple vector units (4)
  - 32 lanes wide → warp size
- Vector units use *hardware multi-threading*
- Execution → a grid of thread blocks (TBs)
  - Each TB has some number of threads

6

# Review



Thread block scheduler

- Each SM has multiple vector units (4)
  - 32 lanes wide → warp size
- Vector units use *hardware multi-threading*
- Execution → a grid of thread blocks (TBs)
  - Each TB has some number of threads

# Review



Thread block scheduler

warp (thread) scheduler

- Each SM has multiple vector units (4)
  - 32 lanes wide → warp size
- Vector units use *hardware multi-threading*
- Execution → a grid of thread blocks (TBs)
  - Each TB has some number of threads

6

# GPU Performance Metric: *Occupancy*

# GPU Performance Metric: *Occupancy*

- Occupancy = (#Active Warps) /(#MaximumActive Warps)
    - Measures how well concurrency/parallelism is utilized

# GPU Performance Metric: *Occupancy*

- Occupancy = (#Active Warps) /(#MaximumActive Warps)
  - Measures how well concurrency/parallelism is utilized
- Occupancy captures
  - *which resources* can be dynamically shared
  - how to reason about resource demands of a CUDA kernel
  - Enables device-specific online tuning of kernel parameters

# GPU Performance Metric: *Occupancy*

- Occupancy = (#Active Warps) /(#MaximumActive Warps)
  - Measures how well concurrency/parallelism is utilized
- Occupancy captures
  - *which resources* can be dynamically shared
  - how to reason about resource demands of a CUDA kernel
  - Enables device-specific online tuning of kernel parameters

# GPU Performance Metric: *Occupancy*

- Occupancy = (#Active Warps) /(#MaximumActive Warps)
  - Measures how well concurrency/parallelism is utilized
- Occupancy captures
  - *which resources* can be dynamically shared
  - how to reason about resource demands of a CUDA kernel
  - Enables device-specific online tuning of kernel parameters

# GPU Performance Metric: *Occupancy*

- Occupancy = (#Active Warps) /(#MaximumActive Warps)
  - Measures how well concurrency/parallelism is utilized
- Occupancy captures
  - *which resources* can be dynamically shared
  - how to reason about resource demands of a CUDA kernel
  - Enables device-specific online tuning of kernel parameters

# GPU Performance Metric: *Occupancy*

- Occupancy = (#Active Warps) /(#MaximumActive Warps)
  - Measures how well concurrency/parallelism is utilized
- Occupancy captures
  - *which resources* can be dynamically shared
  - how to reason about resource demands of a CUDA kernel
  - Enables device-specific online tuning of kernel parameters

# GPU Performance Metric: *Occupancy*

- Occupancy = (#Active Warps) /(#MaximumActive Warps)
  - Measures how well concurrency/parallelism is utilized
- Occupancy captures
  - *which resources* can be dynamically shared
  - how to reason about resource demands of a CUDA kernel
  - Enables device-specific online tuning of kernel parameters

Shouldn't we just create as many threads as possible?

# Hardware Resources Are Finite



SM – Stream Multiprocessor

SP – Stream Processor

# Hardware Resources Are Finite



SM – Stream Multiprocessor

SP – Stream Processor

# Hardware Resources Are Finite



Kernel Distributor

SM Scheduler

SM  SM  SM  SM

DRAM

SM – Stream Multiprocessor

SP – Stream Processor

Thread Block Control

TB 0

Limits the #thread blocks

Warp Schedulers

Warp Context

Limits the #threads

SP SP SP SP
SP SP SP SP
SP SP SP SP
SP SP SP SP

Register File

L1/Shared Memory

# Hardware Resources Are Finite



Kernel Distributor

SM Scheduler

SM  SM  SM  SM

DRAM

SM – Stream Multiprocessor

SP – Stream Processor

Thread Block Control

TB 0

Limits the #thread blocks

Warp Schedulers

Warp Context

Limits the #threads

SP SP SP SP
SP SP SP SP
SP SP SP SP
SP SP SP SP

Register File

Limits the #threads

L1/Shared Memory

# Hardware Resources Are Finite



SM – Stream Multiprocessor

SP – Stream Processor

# Hardware Resources Are Finite



**Kernel Distributor** → **SM Scheduler**

**Thread Block Control**
- TB 0

Limits the **#thread blocks**

**Warp Schedulers**

**Warp Context**

Limits the **#threads**

SP SP SP SP
SP SP SP SP
SP SP SP SP
SP SP SP SP

**Register File**

Limits the **#threads**

**L1/Shared Memory**

Limits the **#thread blocks**

**Occupancy:**

- (#Active Warps) /(#MaximumActive Warps)

- Limits on the numerator:
  - Registers/thread
  - Shared memory/thread block
  - Number of scheduling slots: blocks, warps

- Limits on the denominator:
  - Memory bandwidth
  - Scheduler slots

# Hardware Resources Are Finite



**Kernel Distributor**

**SM Scheduler**

**Thread Block Control**

TB 0

Limits the #thread blocks

**Warp Schedulers**

**Warp Context**

Limits the #threads

SP SP SP SP
SP SP SP SP
SP SP SP SP
SP SP SP SP

**Register File**

Limits the #threads

**L1/Shared Memory**

Limits the #thread blocks
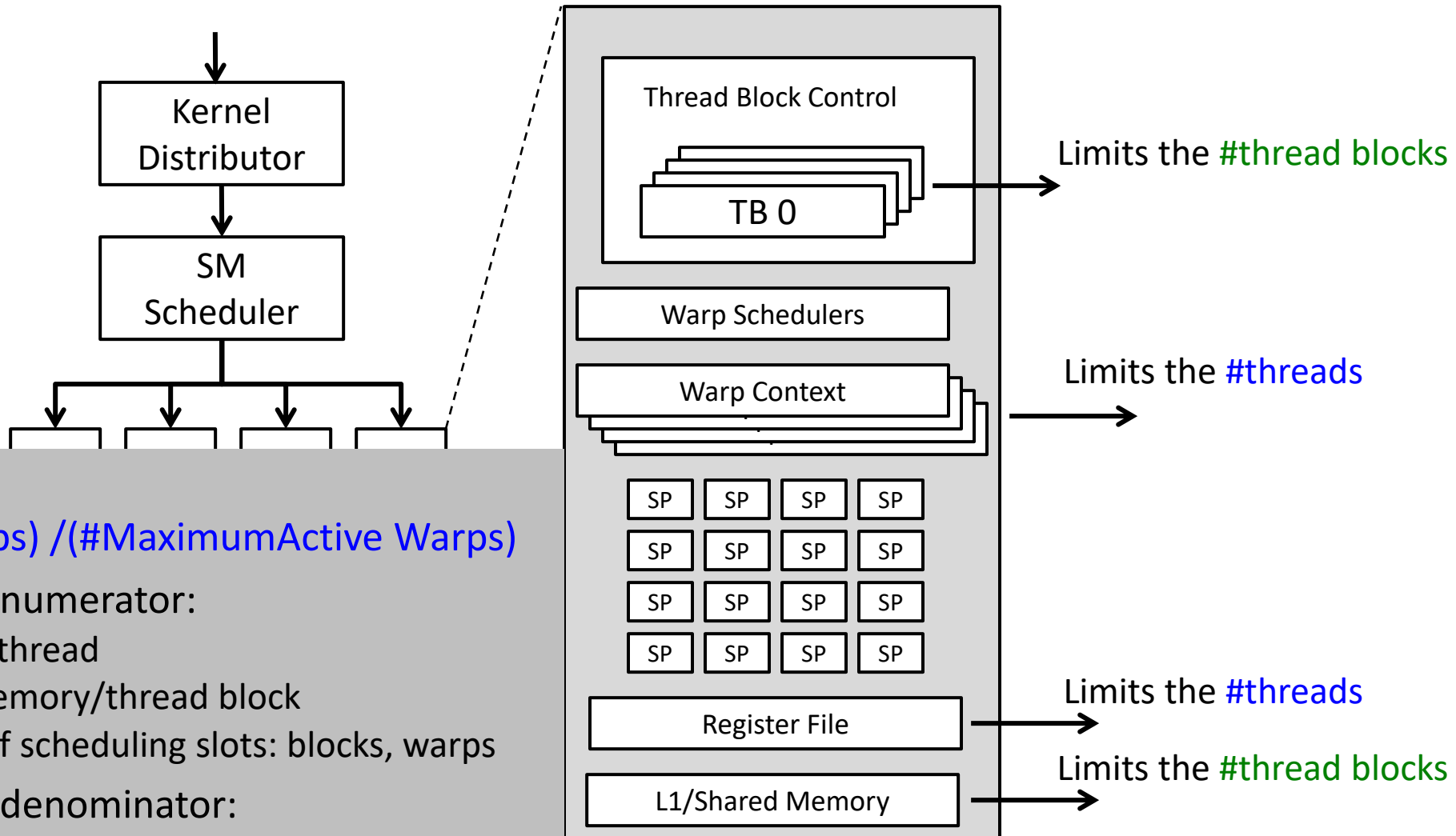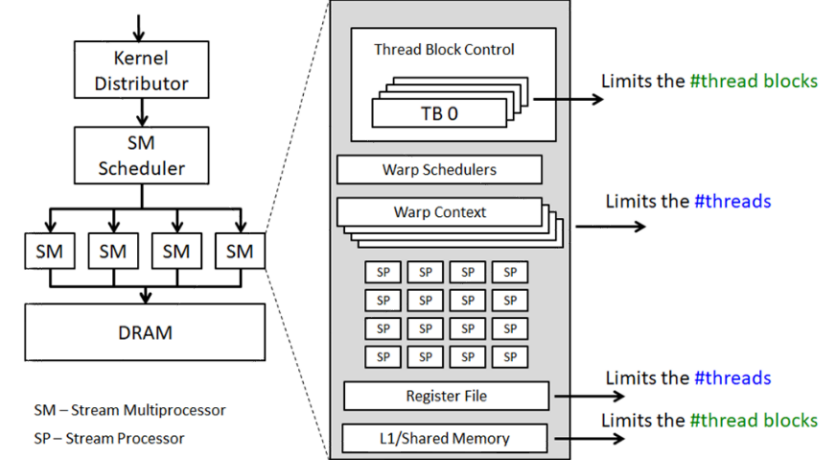
Occupancy:

- (#Active Warps) /(#MaximumActive Warps)

- Limits on the numerator:
  - Registers/thread
  - Shared memory/thread block
  - Number of scheduling slots: blocks, warps
- Limits on the denominator:
  - Memory bandwidth
  - Scheduler slots

What is the performance impact of varying kernel resource demands?

# Impact of Thread Block Size

# Impact of Thread Block Size

Example: v100:



Kernel Distributor

SM Scheduler

SM SM SM SM

DRAM

SM – Stream Multiprocessor
SP – Stream Processor

Thread Block Control

TB 0

Limits the #thread blocks

Warp Schedulers

Warp Context

Limits the #threads

SP SP SP SP
SP SP SP SP
SP SP SP SP
SP SP SP SP

Register File

Limits the #threads

L1/Shared Memory

Limits the #thread blocks
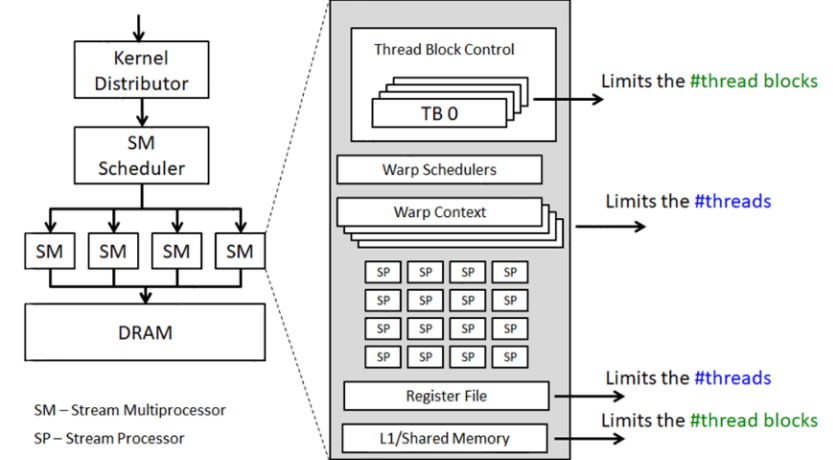
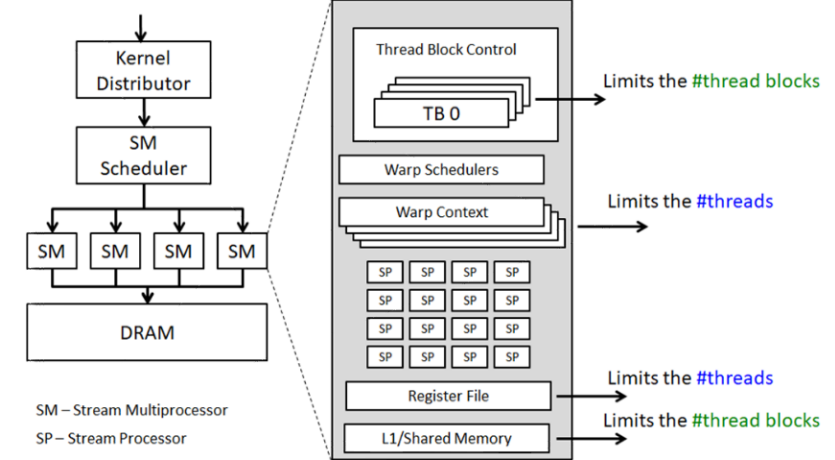# Impact of Thread Block Size



Example: v100:
- max active warps/SM == 64 (limit: warp context)

# Impact of Thread Block Size



Example: v100:
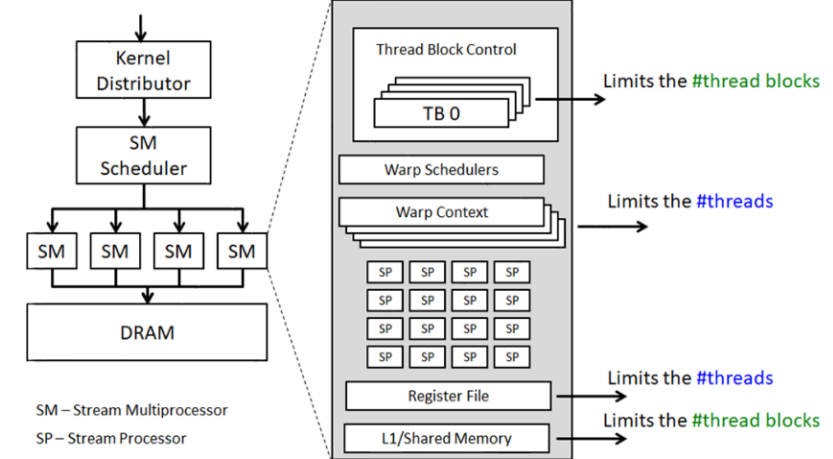- max active warps/SM == 64 (limit: warp context)
- max active blocks/SM == 32 (limit: block control)

# Impact of Thread Block Size



Example: v100:

- max active warps/SM == 64 (limit: warp context)
- max active blocks/SM == 32 (limit: block control)
  - With 512 threads/block how many blocks can execute (per SM) concurrently?
  - Max active warps * threads/warp = 64*32 = 2048 threads →

# Impact of Thread Block Size



SM – Stream Multiprocessor
SP – Stream Processor

Example: v100:

- max active warps/SM == 64 (limit: warp context)
- max active blocks/SM == 32 (limit: block control)
  - With 512 threads/block how many blocks can execute (per SM) concurrently?
  - Max active warps * threads/warp = 64*32 = 2048 threads → 4

# Impact of Thread Block Size



SM – Stream Multiprocessor
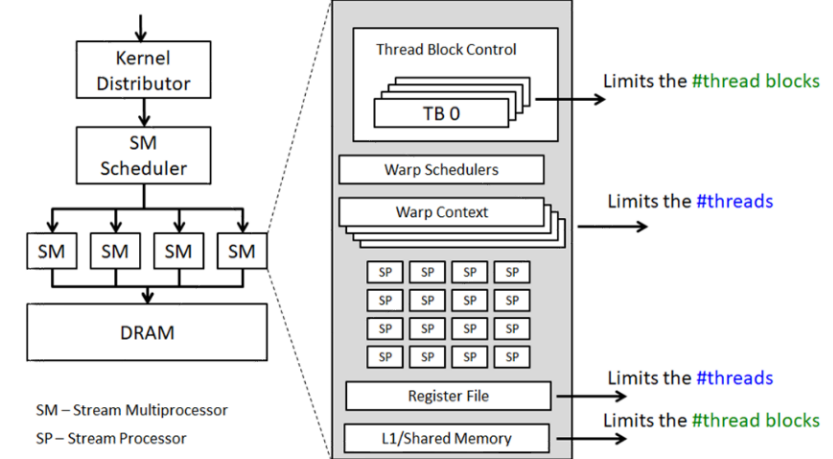SP – Stream Processor

Example: v100:

- max active warps/SM == 64 (limit: warp context)
- max active blocks/SM == 32 (limit: block control)
  - With 512 threads/block how many blocks can execute (per SM) concurrently?
  - Max active warps * threads/warp = 64*32 = 2048 threads → 4
  - With 128 threads/block? →

# Impact of Thread Block Size



SM – Stream Multiprocessor
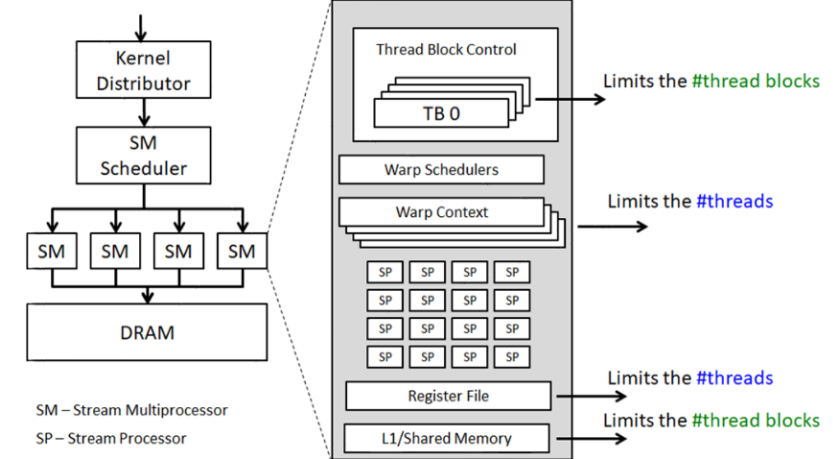SP – Stream Processor

Example: v100:

- max active warps/SM == 64 (limit: warp context)
- max active blocks/SM == 32 (limit: block control)
  - With 512 threads/block how many blocks can execute (per SM) concurrently?
  - Max active warps * threads/warp = 64*32 = 2048 threads → 4
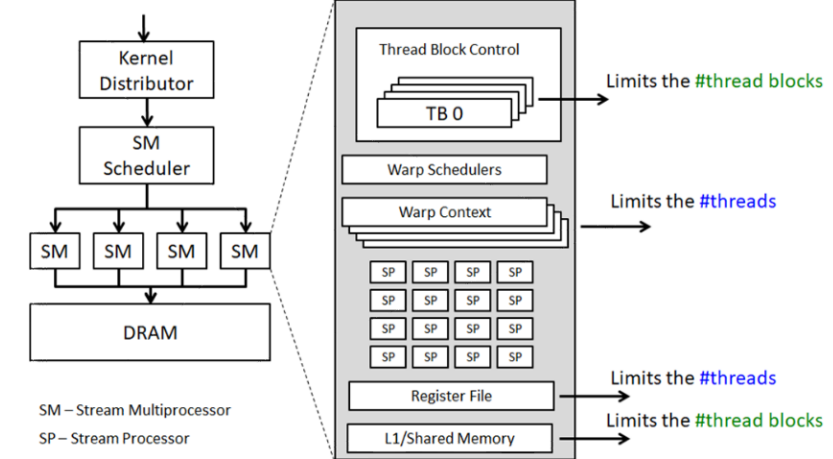  - With 128 threads/block? → 16

# Impact of Thread Block Size



Example: v100:

- max active warps/SM == 64 (limit: warp context)
- max active blocks/SM == 32 (limit: block control)
  - With 512 threads/block how many blocks can execute (per SM) concurrently?
  - Max active warps * threads/warp = 64*32 = 2048 threads → 4
  - With 128 threads/block? → 16
- Consider HW limit of 32 thread blocks/SM @ 32 threads/block:
  - Blocks are maxed out, but max active threads = 32*32 = 1024
  - Occupancy = .5 (1024/2048)

# Impact of Thread Block Size
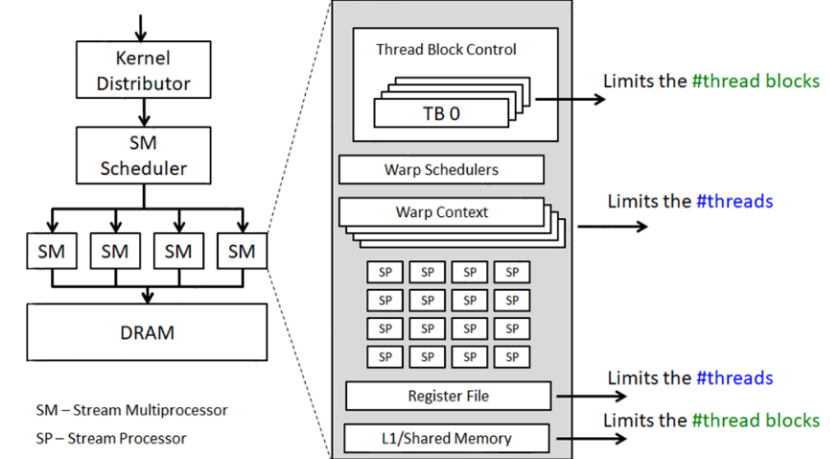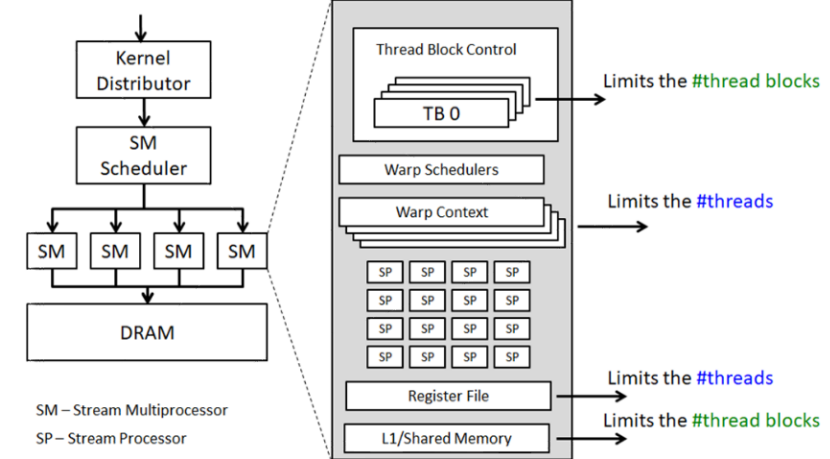


Example: v100:

- max active warps/SM == 64 (limit: warp context)
- max active blocks/SM == 32 (limit: block control)
  - With 512 threads/block how many blocks can execute (per SM) concurrently?
  - Max active warps * threads/warp = 64*32 = 2048 threads → 4
  - With 128 threads/block? → 16
- Consider HW limit of 32 thread blocks/SM @ 32 threads/block:
  - Blocks are maxed out, but max active threads = 32*32 = 1024
  - Occupancy = .5 (1024/2048)
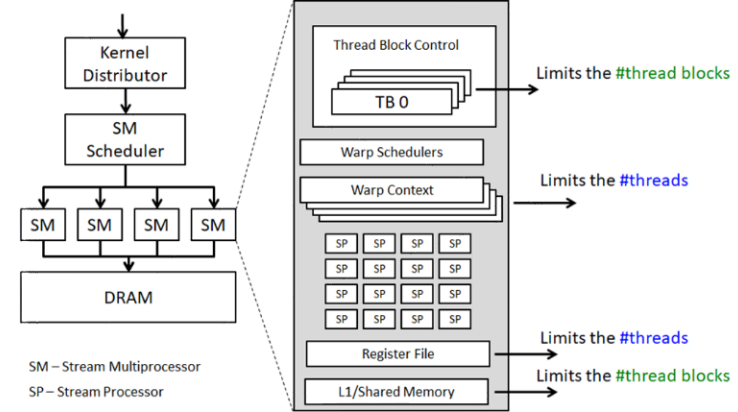- To maximize utilization, thread block size should balance
  - Limits on active thread blocks vs.
  - Limits on active warps

9

# Impact of #Registers Per Thread

# Impact of #Registers Per Thread



Registers/thread can limit number of active threads!

# Impact of #Registers Per Thread



Registers/thread can limit number of active threads!

V100:

# Impact of #Registers Per Thread



Registers/thread can limit number of active threads!

V100:
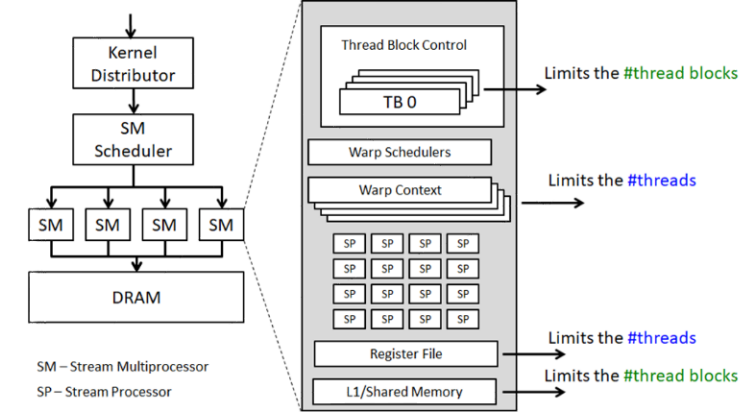
- Registers per thread max: 255

# Impact of #Registers Per Thread



Registers/thread can limit number of active threads!

V100:

- Registers per thread max: 255
- 64K registers per SM

# Impact of #Registers Per Thread



SM – Stream Multiprocessor
SP – Stream Processor

Registers/thread can limit number of active threads!

V100:

- Registers per thread max: 255
- 64K registers per SM

Assume a kernel uses 32 registers/thread, thread block size of 256

# Impact of #Registers Per Thread



Registers/thread can limit number of active threads!

V100:

- Registers per thread max: 255
- 64K registers per SM

Assume a kernel uses 32 registers/thread, thread block size of 256

- Thus, A TB requires 8192 registers for a maximum of 8 thread blocks per SM
  - Uses all 2048 thread slots (8 blocks * 256 threads/block)
  - 8192 *regs/block * 8 block/SM = 64k registers*
  - *FULLY Occupied!*
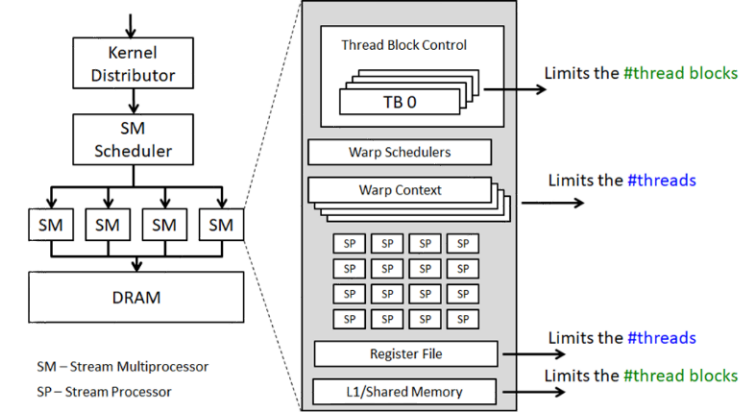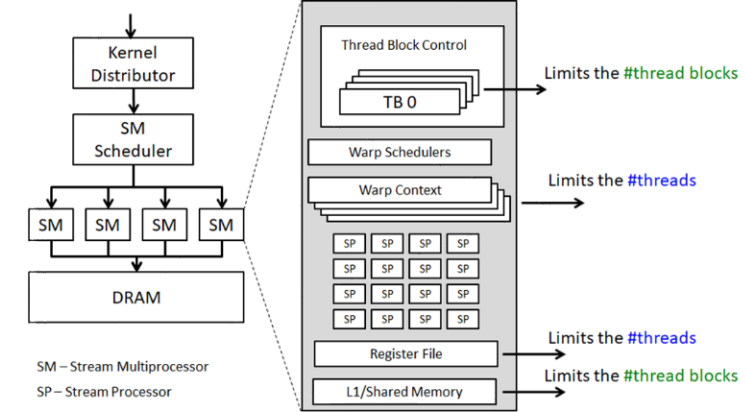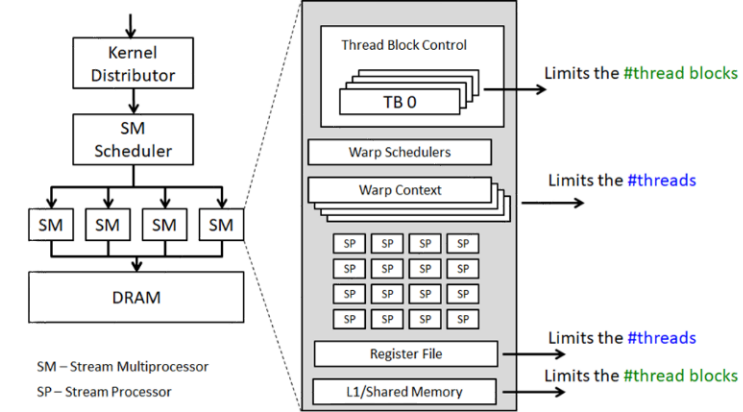
# Impact of #Registers Per Thread



Registers/thread can limit number of active threads!

V100:

- Registers per thread max: 255
- 64K registers per SM

Assume a kernel uses 32 registers/thread, thread block size of 256

- Thus, A TB requires 8192 registers for a maximum of 8 thread blocks per SM
  - Uses all 2048 thread slots (8 blocks * 256 threads/block)
  - 8192 *regs/block * 8 block/SM = 64k registers*
  - *FULLY Occupied!*
- What is the impact of increasing number of registers by 2?

# Impact of #Registers Per Thread



SM – Stream Multiprocessor
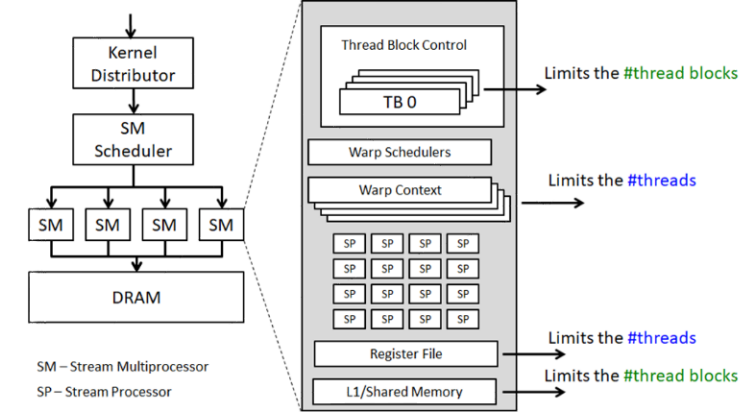SP – Stream Processor

Registers/thread can limit number of active threads!

V100:

- Registers per thread max: 255
- 64K registers per SM

Assume a kernel uses 32 registers/thread, thread block size of 256

- Thus, A TB requires 8192 registers for a maximum of 8 thread blocks per SM
  - Uses all 2048 thread slots (8 blocks * 256 threads/block)
  - 8192 *regs/block * 8 block/SM = 64k registers*
  - *FULLY Occupied!*
- What is the impact of increasing number of registers by 2?
  - Recall: granularity of management is a thread block!
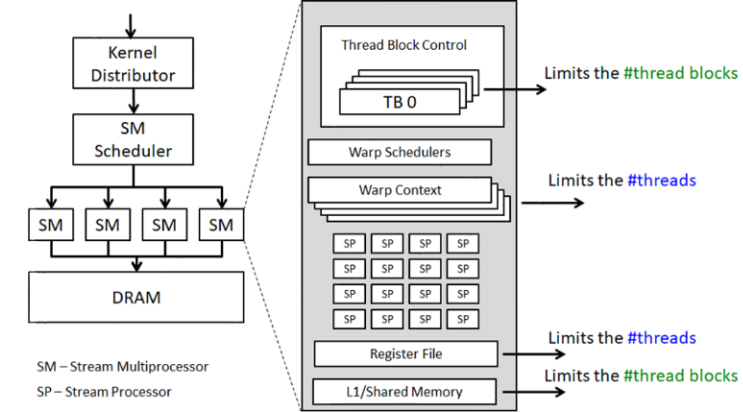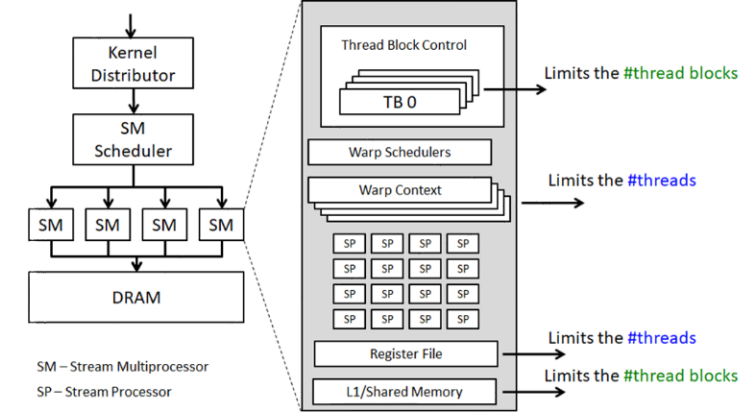
# Impact of #Registers Per Thread



Registers/thread can limit number of active threads!

V100:

- Registers per thread max: 255
- 64K registers per SM

Assume a kernel uses 32 registers/thread, thread block size of 256

- Thus, A TB requires 8192 registers for a maximum of 8 thread blocks per SM
  - Uses all 2048 thread slots (8 blocks * 256 threads/block)
  - 8192 *regs/block * 8 block/SM = 64k registers*
  - *FULLY Occupied!*
- What is the impact of increasing number of registers by 2?
  - Recall: granularity of management is a thread block!
  - Loss of concurrency of 256 threads!
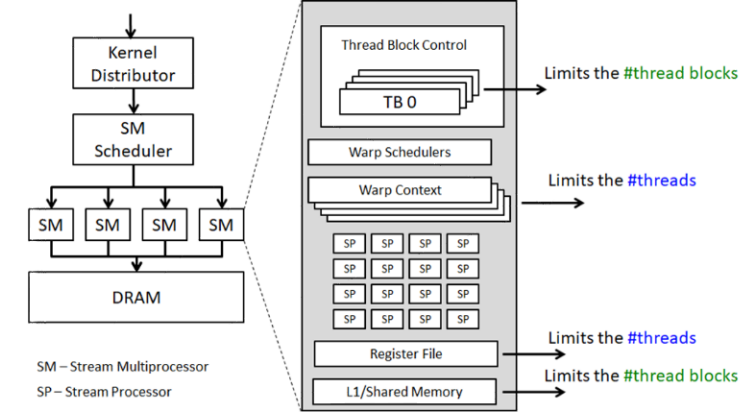  - *34 regs/thread * 256 threads/block * 7 blocks/SM = 60k registers,*
  - *8 blocks would over-subscribe register file*
  - *Occupancy drops to .875!*

# Control Flow Divergence

- Performance concern with branching: divergence
  - Threads within a single warp take different paths
  - Different execution paths are serialized
    - The control paths taken by the threads in a warp are traversed one at a time until there is no more.
- Common case: branch condition is a function of thread ID
  - Example with divergence:
    - `If (threadIdx.x > 2) { }`
    - This creates two different control paths for threads in a block
    - Branch granularity < warp size; threads 0, 1 and 2 follow different path than the rest of the threads in the first warp
  - Example without divergence:
    - `If (threadIdx.x / WARP_SIZE > 2) { }`
    - Also creates two different control paths for threads in a block
    - Branch granularity is a whole multiple of warp size; all threads in any given warp follow the same path

# FPGAs/Verilog

# FPGAs/Verilog

- CLB, BRAM, and LUT?

# FPGAs/Verilog

- CLB, BRAM, and LUT?
- CLB: combinational logic block

# FPGAs/Verilog

- CLB, BRAM, and LUT?
- CLB: combinational logic block
- BRAM: block random access memory

# FPGAs/Verilog

- CLB, BRAM, and LUT?
- CLB: combinational logic block
- BRAM: block random access memory
- LUT: lookup table

# FPGAs/Verilog

- CLB, BRAM, and LUT?
- CLB: combinational logic block
- BRAM: block random access memory
- LUT: lookup table
- Other questions?

# Blocking vs Non-blocking Behavior

- A sequence of nonblocking assignments don't communicate

a = 1;
b = a;
c = b;

Blocking assignment:
a = b = c = 1

a <= 1;
b <= a;
c <= b;
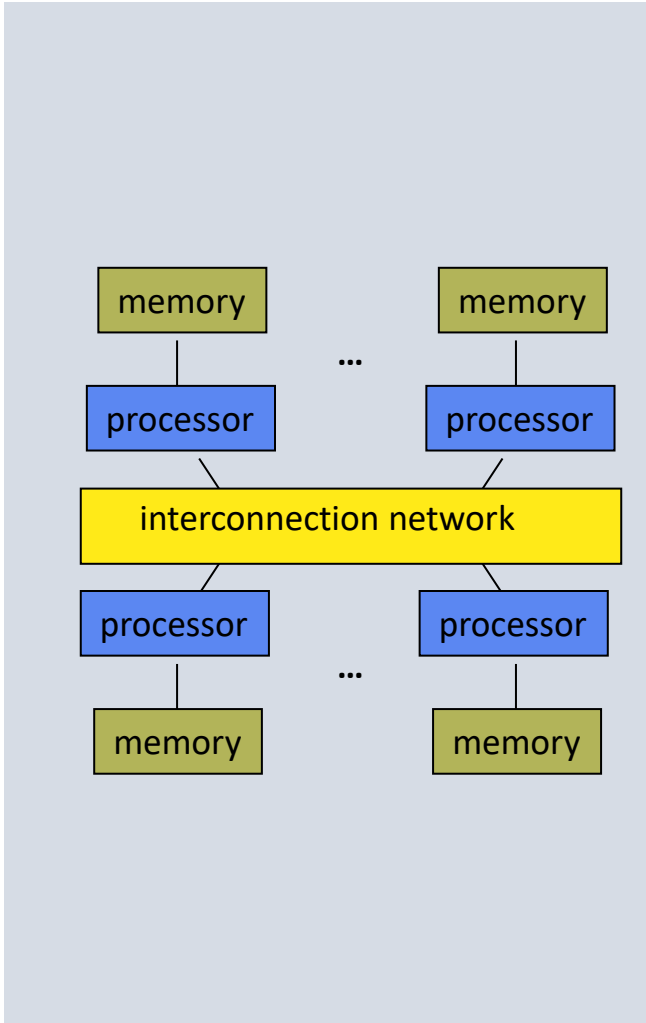
Nonblocking assignment:
a = 1
b = old value of a
c = old value of b

# MPI

# MPI

# MPI

Distributed Memory
Multiprocessor

Messaging between nodes

# MPI

Distributed Memory
Multiprocessor

Messaging between nodes

memory    ...    memory

processor    processor

interconnection network

processor    processor

...

memory    memory

Massively Parallel Processor (MPP)

Many, many processors

# MPI

Distributed Memory
Multiprocessor

Messaging between nodes

memory    ...    memory

processor    processor

interconnection network

processor    processor

...

memory    memory

Massively Parallel Processor (MPP)

Many, many processors

# MPI



Distributed Memory Multiprocessor

Messaging between nodes

interconnection network

Massively Parallel Processor (MPP)

Many, many processors

Cluster of SMPs
- Shared memory in SMP node
- Messaging ←→ SMP nodes

network interface

interconnection network

- also regarded as MPP if processor # is large
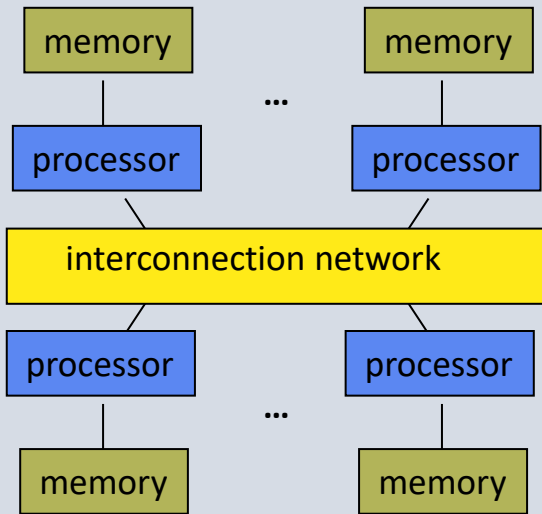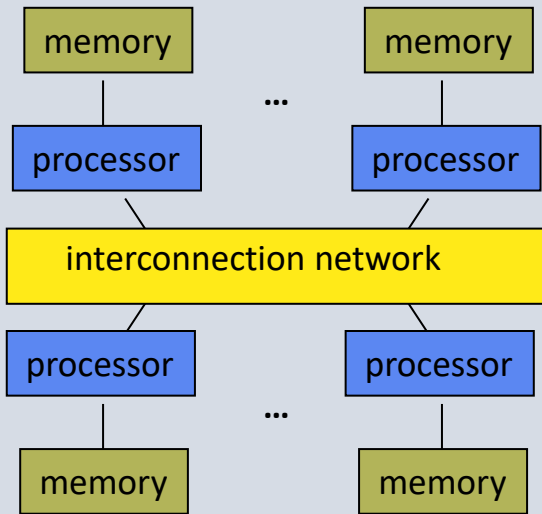
# MPI



**Distributed Memory Multiprocessor**

Messaging between nodes

Massively Parallel Processor (MPP)

Many, many processors

**Cluster of SMPs**
- Shared memory in SMP node
- Messaging ←→ SMP nodes

network interface

- also regarded as MPP if processor # is large

**Multicore SMP+GPU Cluster**
- Shared mem in SMP node
- Messaging between nodes

- GPU accelerators attached

# MPI



Distributed Memory Multiprocessor
Messaging between nodes
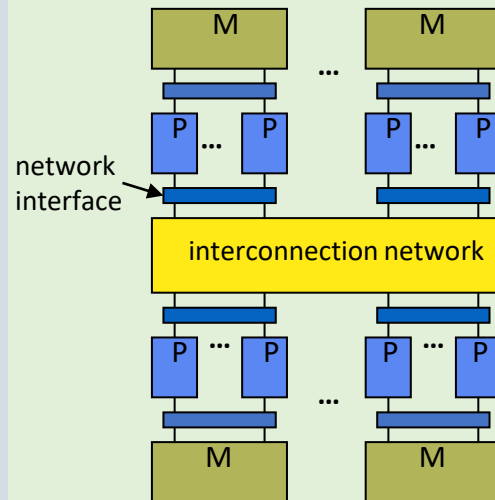
memory — processor — interconnection network — processor — memory

Massively Parallel Processor (MPP)
Many, many processors

Cluster of SMPs
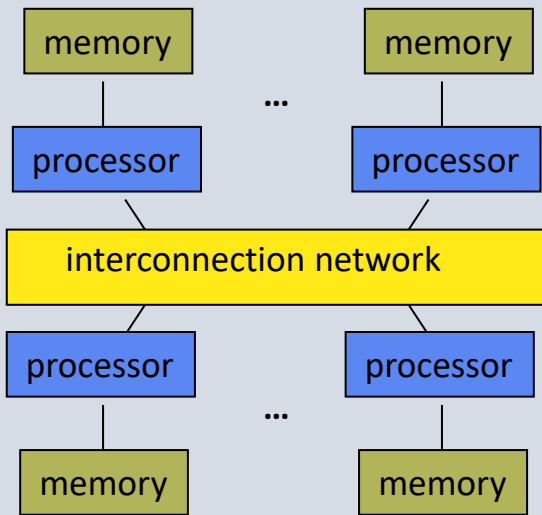- Shared memory in SMP node
- Messaging ⟵⟶ SMP nodes

network interface

interconnection network

- also regarded as MPP if processor # is large

Multicore SMP+GPU Cluster
- Shared mem in SMP n...
- Messaging between n...

interconnection network

- GPU accelerators attached

PGAS = partitioned global address space
How is that different from shared nothing?

# What is NoSQL?

- Next Generation Compute/Storage engines (databases)
    - **non-relational**
    - **distributed**
    - **open-source**
    - **horizontally scalable**
- One view: "no" → elide SQL/database functionality to achieve scale
- Another view: "NoSQL" is actually misleading.
    - more appropriate term is actually "Not Only SQL"

# What is NoSQL?

- Next Generation Compute/Storage engines (databases)
  - **non-relational**
  - **distributed**
  - **open-source**
  - **horizontally scalable**
- One view: "no" → elide SQL/da
- Another view: "NoSQL" is actua
  - more appropriate term is actually "Not Only SQL"

What NoSQL gives up in exchange for scale:
- Relationships between entities are non-existent
- Limited or no ACID transactions
- No standard language for queries (SQL)
- Less structured

# What is NoSQL?

- Next Generation Compute/Storage engines (databases)
  - **non-relational**
  - **distributed**
  - **open-source**
  - **horizontally scalable**
- One view: "no" → elide SQL/dat
- Another view: "NoSQL" is actua
  - more appropriate term is actually "No

What NoSQL gives up in exchange for scale:

- *Why talk about NoSQL in concurrency class?*
  - Principle
    - Most tradeoffs are a ***direct result*** of concurrency
  - Practice
    - NoSQL systems are ubiquitous
  - Relevant aspects
    - scale/performance tradeoff space
    - Correctness/programmability tradeoff space

# Review: noSQL Taxonomy

# Review: noSQL Taxonomy

*Consistency*

# Review: noSQL Taxonomy

*Consistency*

*Data Model*

# Review: noSQL Taxonomy

*Data Model*

*Consistency*

*Implementation Techniques*

# Review: noSQL Taxonomy



Data Model

Strong: ACID     Consistency     Eventual: BASE

Implementation Techniques

16

# Review: noSQL Taxonomy



Data Model

- Atomicity
- Consistency
- Isolation
- Durability

Strong: ACID

Consistency

Eventual: BASE

Implementation Techniques

# Review: noSQL Taxonomy



**Strong: ACID**
- Atomicity
- Consistency
- Isolation
- Durability

**Eventual: BASE**
- Basically Available
- Soft State
- Eventually Consistent

*Consistency*

*Data Model*

*Implementation Techniques*

16

# Review: noSQL Taxonomy



Data Model

Strong: ACID   *Consistency*   Eventual: BASE

Implementation Techniques

# Review: noSQL Taxonomy



Data Model

Strong: ACID    Consistency    Eventual: BASE

Implementation Techniques

16

# Review: noSQL Taxonomy

*Key Value Stores*

*Data Model*

**Strong: ACID** *Consistency* **Eventual: BASE**

*Implementation Techniques*

16

# Review: noSQL Taxonomy



*Key Value Stores*

*Document Stores*

*Data Model*

**Strong: ACID**  *Consistency*  **Eventual: BASE**
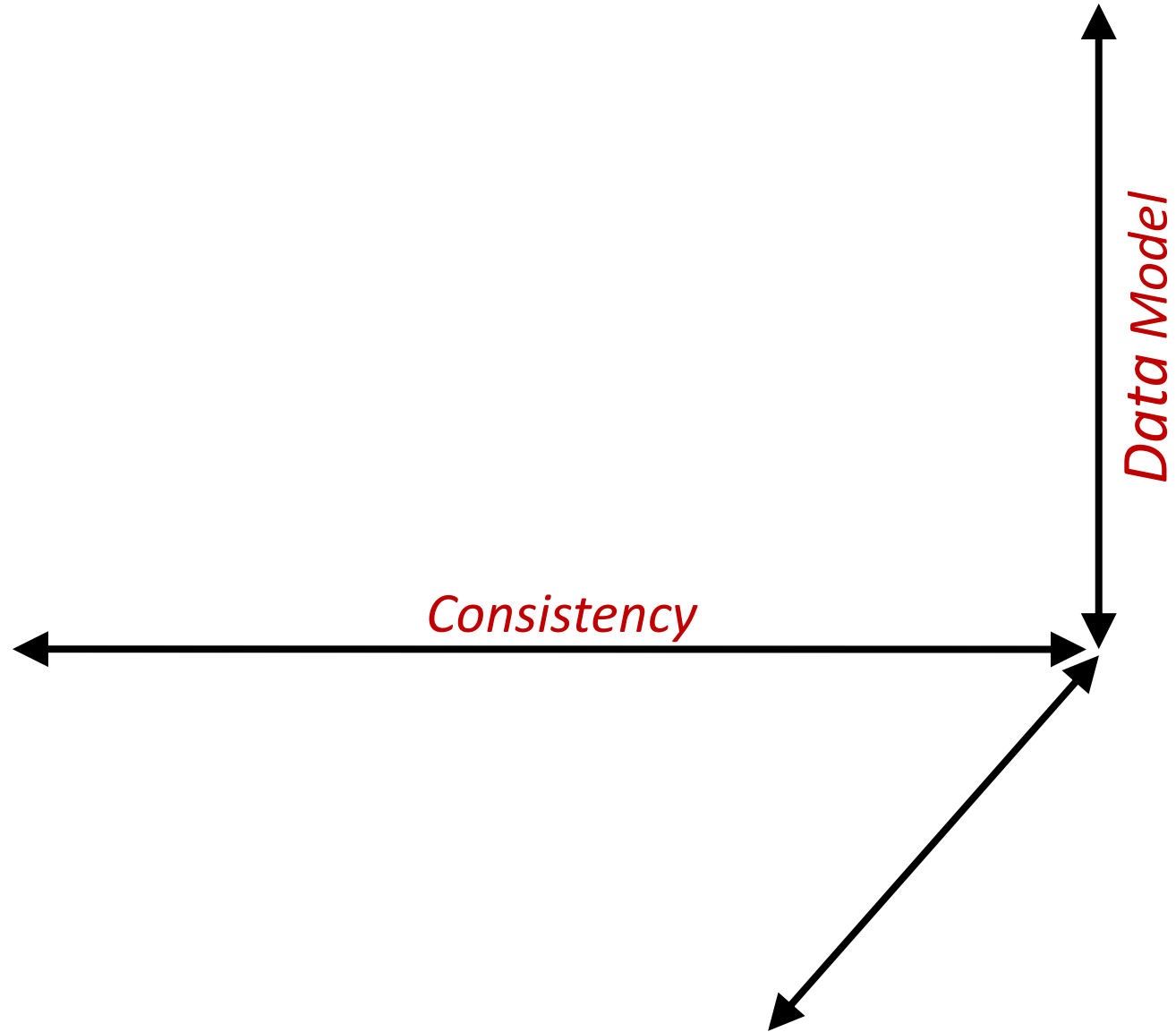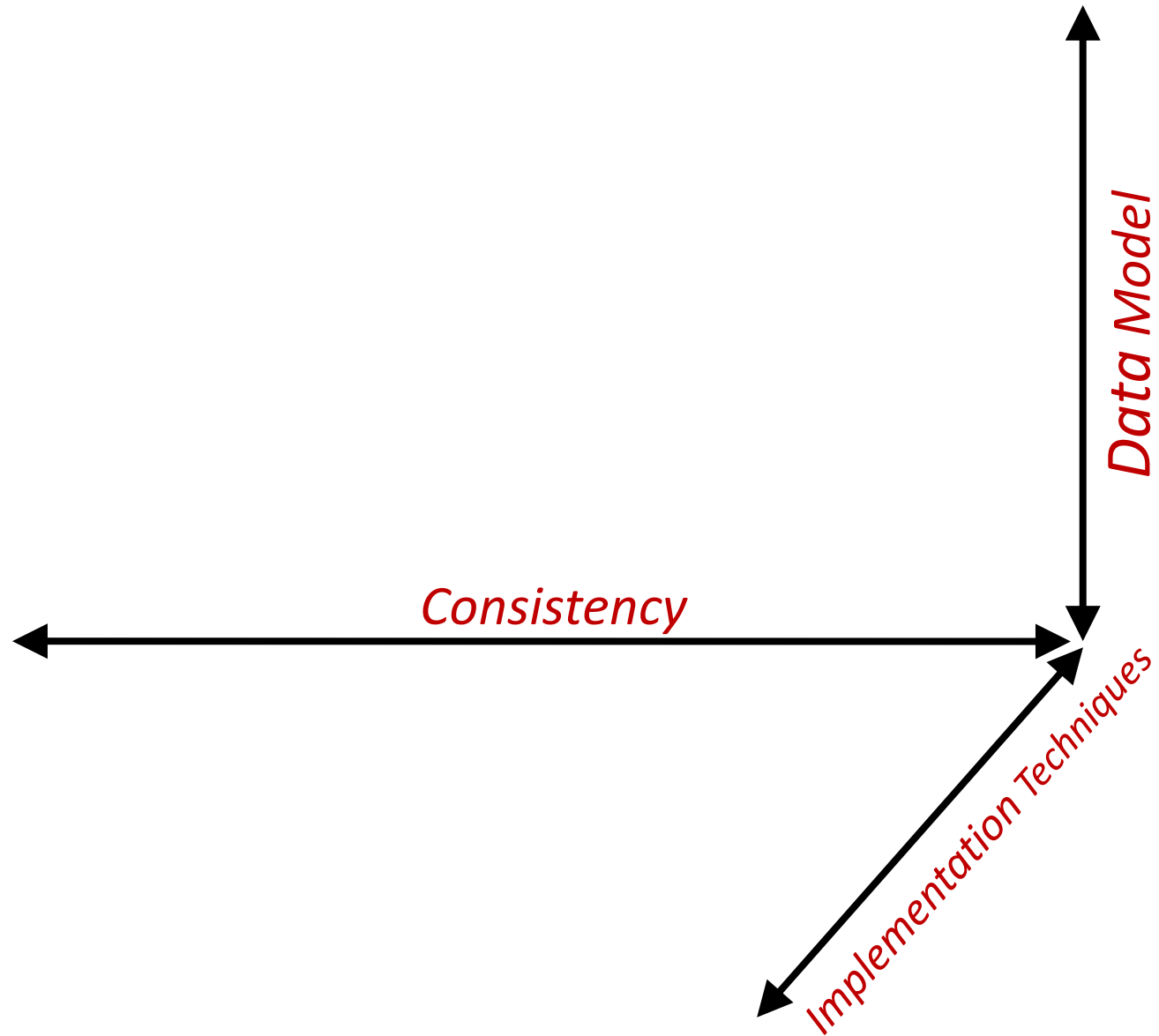
*Implementation Techniques*

16

# Review: noSQL Taxonomy

# Review: noSQL Taxonomy

# Review: noSQL Taxonomy



*Key Value Stores*

*Document Stores*

*Wide-Column Stores*

**Strong: ACID**   *Consistency*   **Eventual: BASE**

*Sharding/Partitioning*

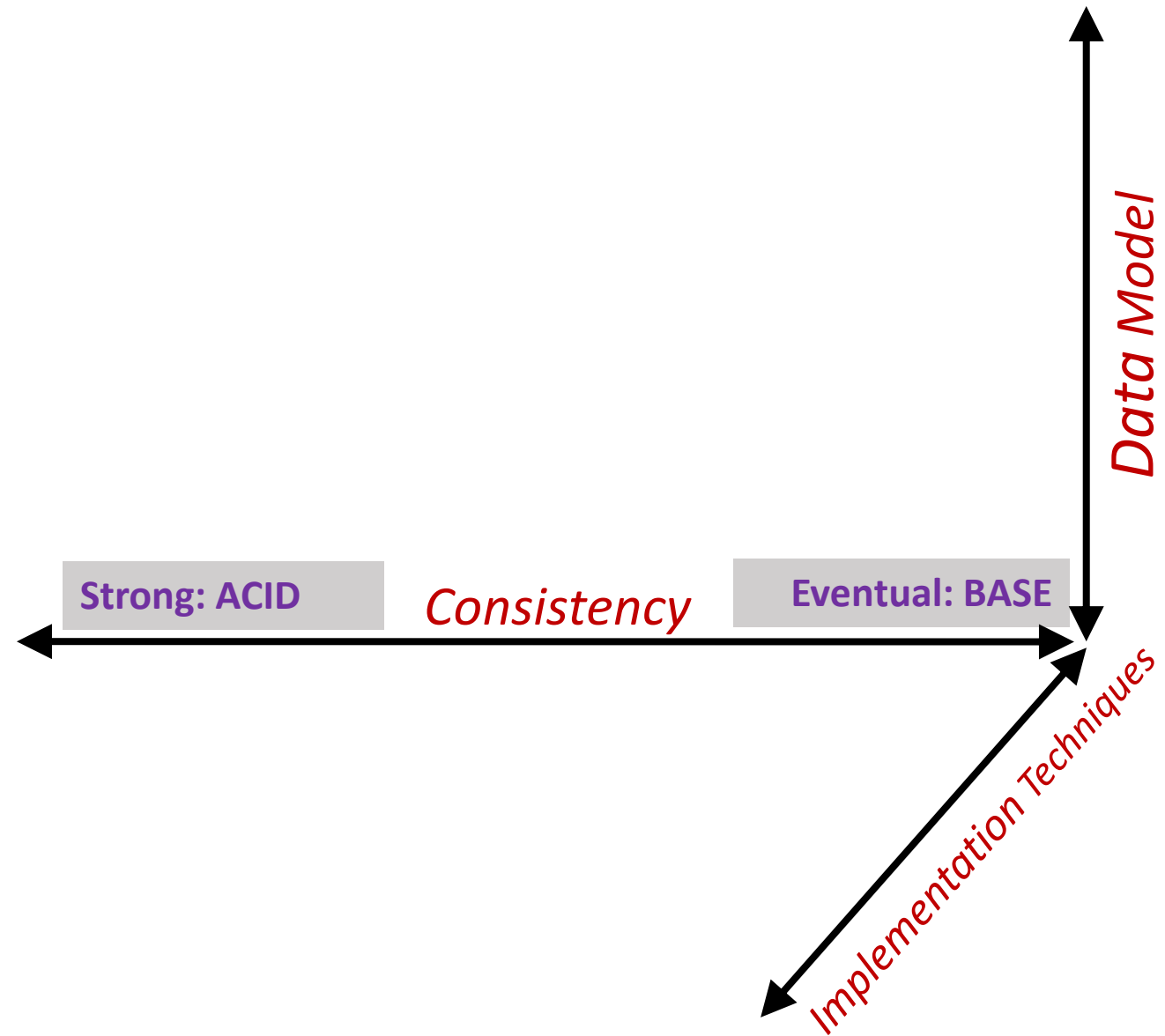*Replication*

*Data Model*

*Implementation Techniques*

# Review: noSQL Taxonomy
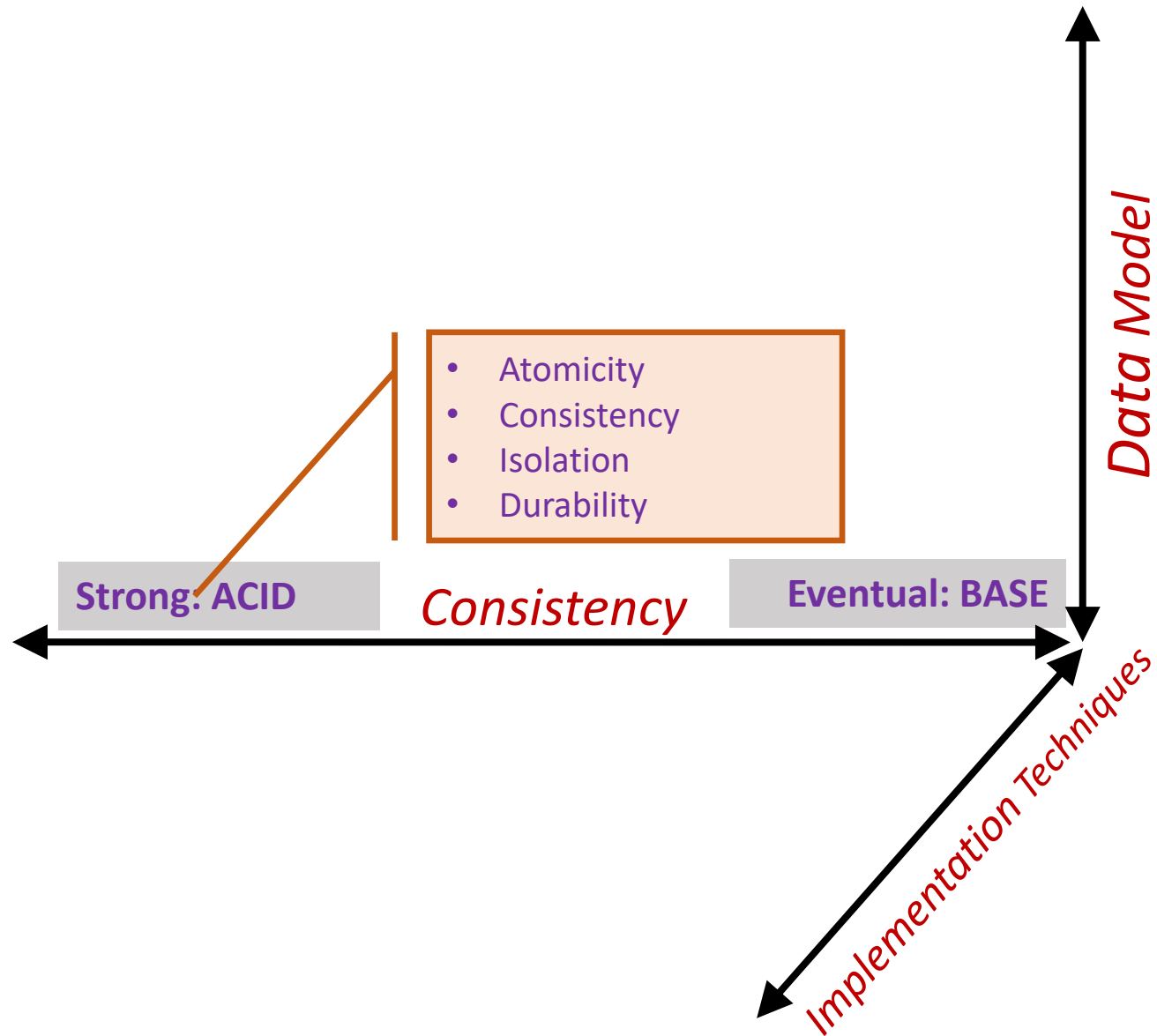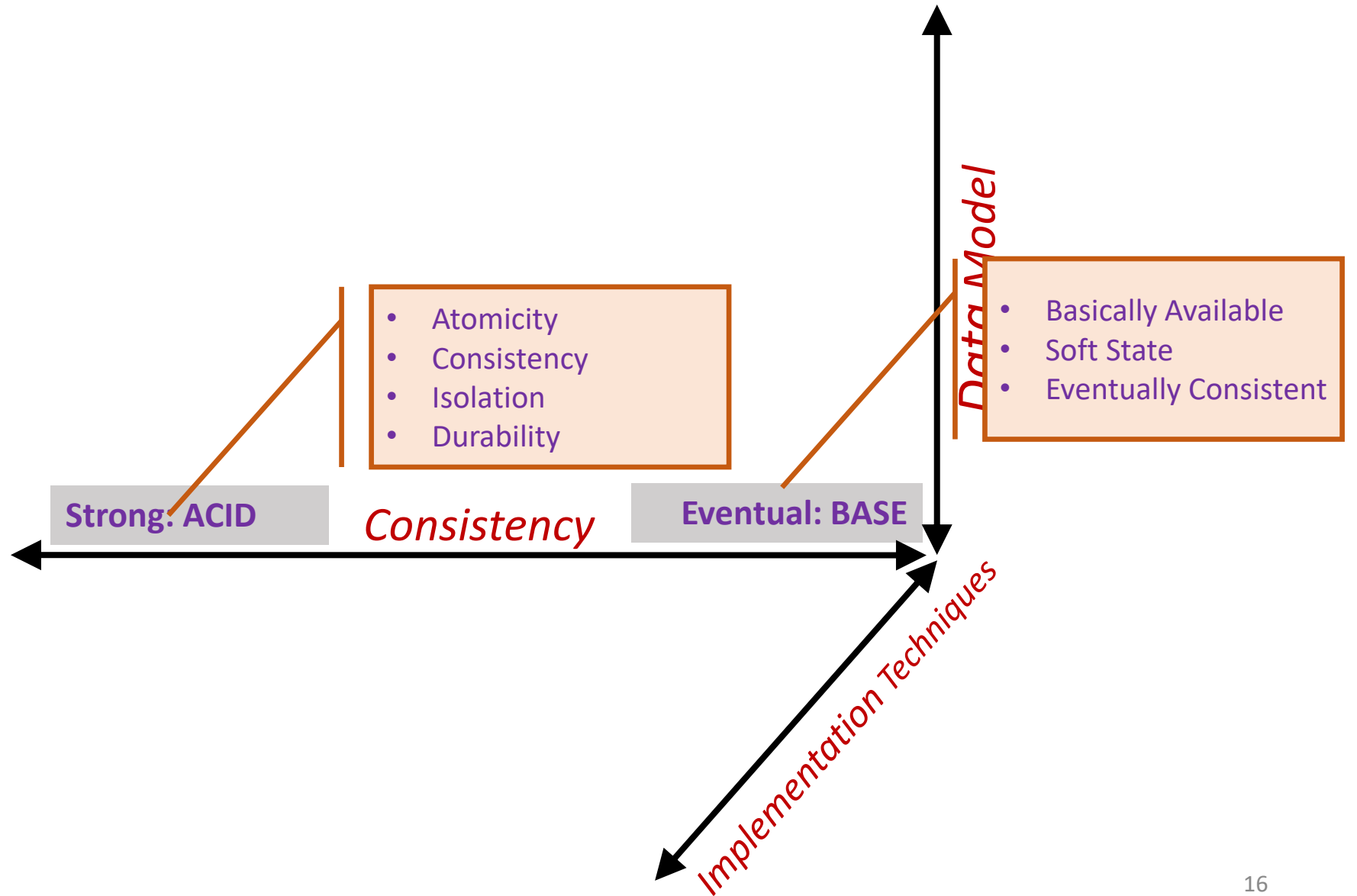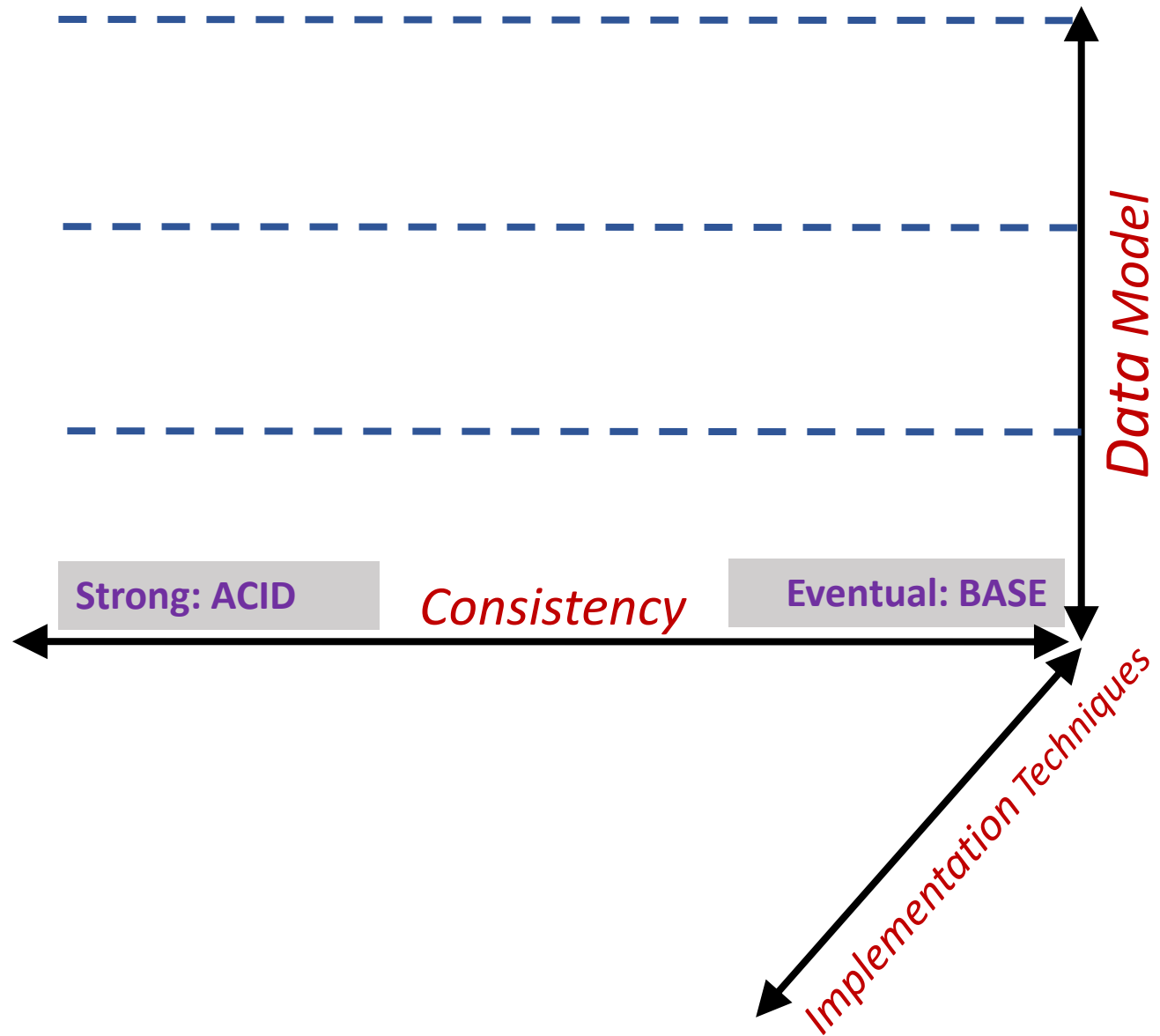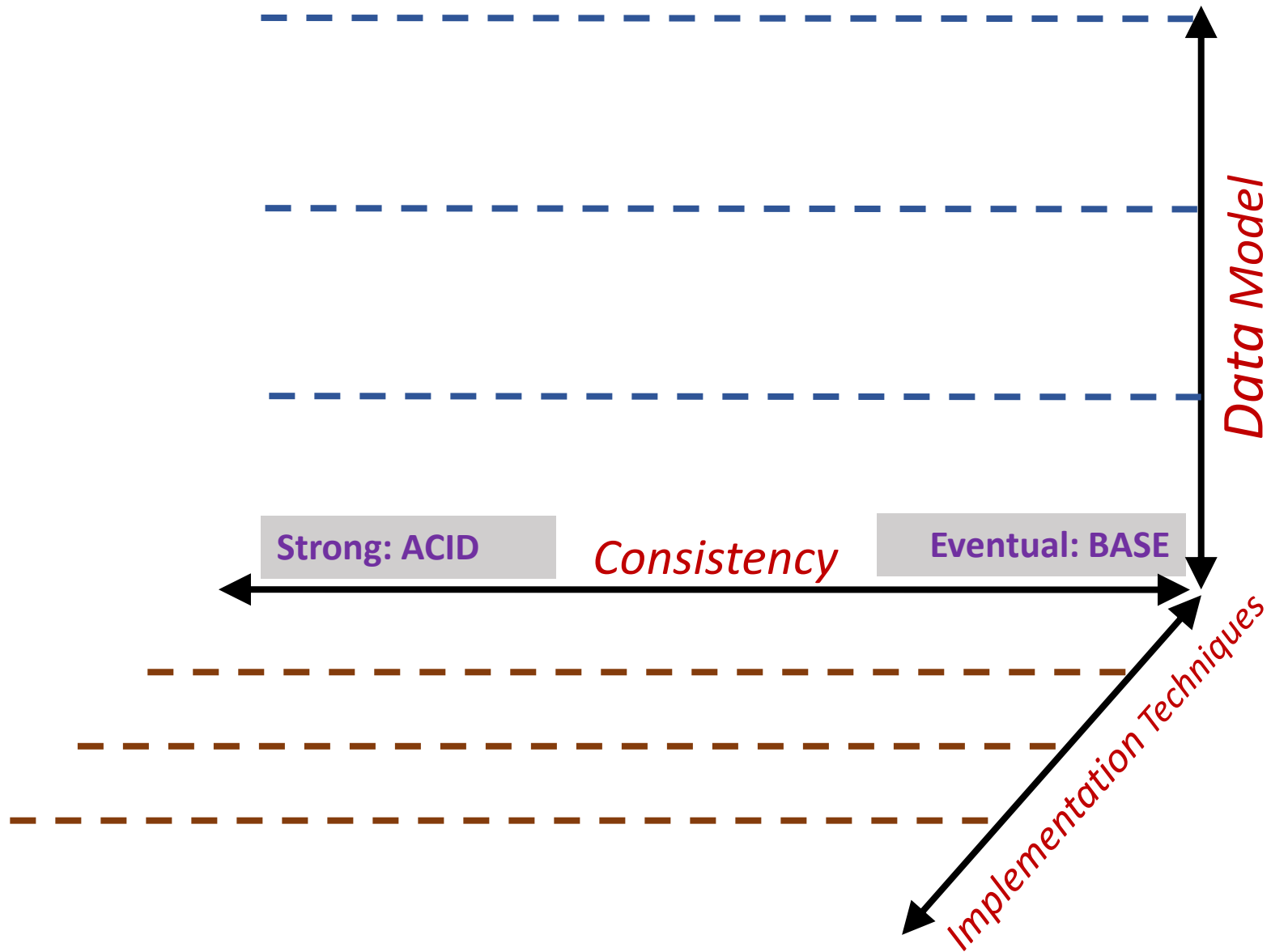
# Review: noSQL Taxonomy

# Review: noSQL Taxonomy



Key Value Stores

Document Stores

Wide-Column Stores

Strong: ACID

Eventual: BASE

*Consistency*

*Data Model*

Sharding/Partitioning

Replication

Storage

Query Support

*Implementation Techniques*

- Shared-Disk
- Range-Sharding
- Hash-Sharding
- Consistent Hashing

# Review: noSQL Taxonomy

# Review: noSQL Taxonomy



*Key Value Stores*

*Document Stores*

*Wide-Column Stores*

**Data Model**

**Strong: ACID**     *Consistency*     **Eventual: BASE**

*Sharding/Partitioning*

*Replication*

*Storage*

*Query Support*

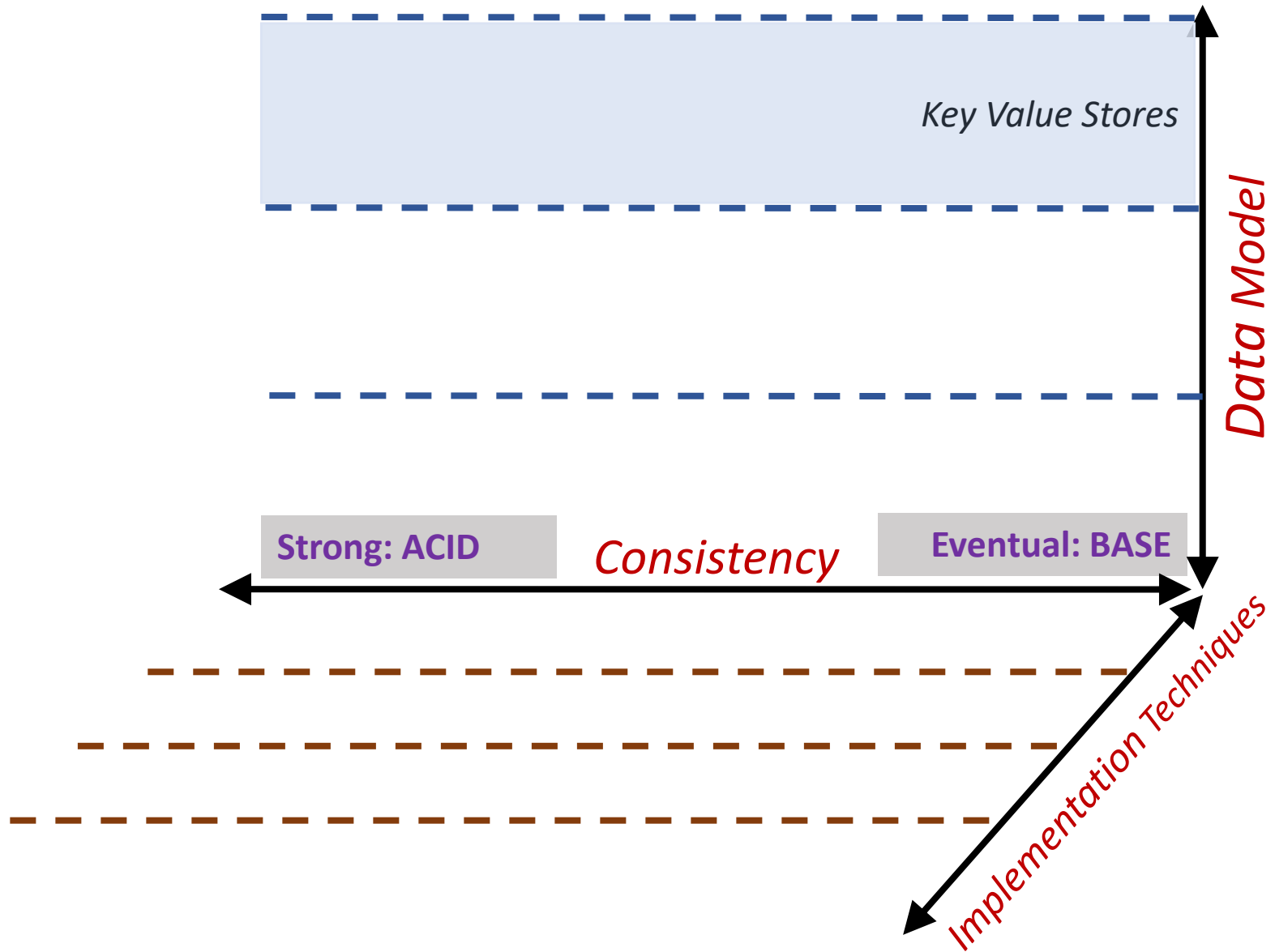**Implementation Techniques**

- Primary-Backup
- Commit-Consensus Protocol
- Sync/Async

16

# Review: noSQL Taxonomy

# Review: noSQL Taxonomy



Key Value Stores

Document Stores

Wide-Column Stores

*Data Model*

**Strong: ACID**      *Consistency*      **Eventual: BASE**

*Sharding/Partitioning*

*Replication*

*Storage*

*Query Support*

*Implementation Techniques*

- Logging
- Update In Place
- Caching
- In-Memory Storage

16

# Review: noSQL Taxonomy



Key Value Stores

Document Stores

Wide-Column Stores

**Strong: ACID**   *Consistency*   **Eventual: BASE**

*Data Model*

Sharding/Partitioning

Replication

Storage

Query Support

*Implementation Techniques*

16

# Review: noSQL Taxonomy

# Review: noSQL Taxonomy

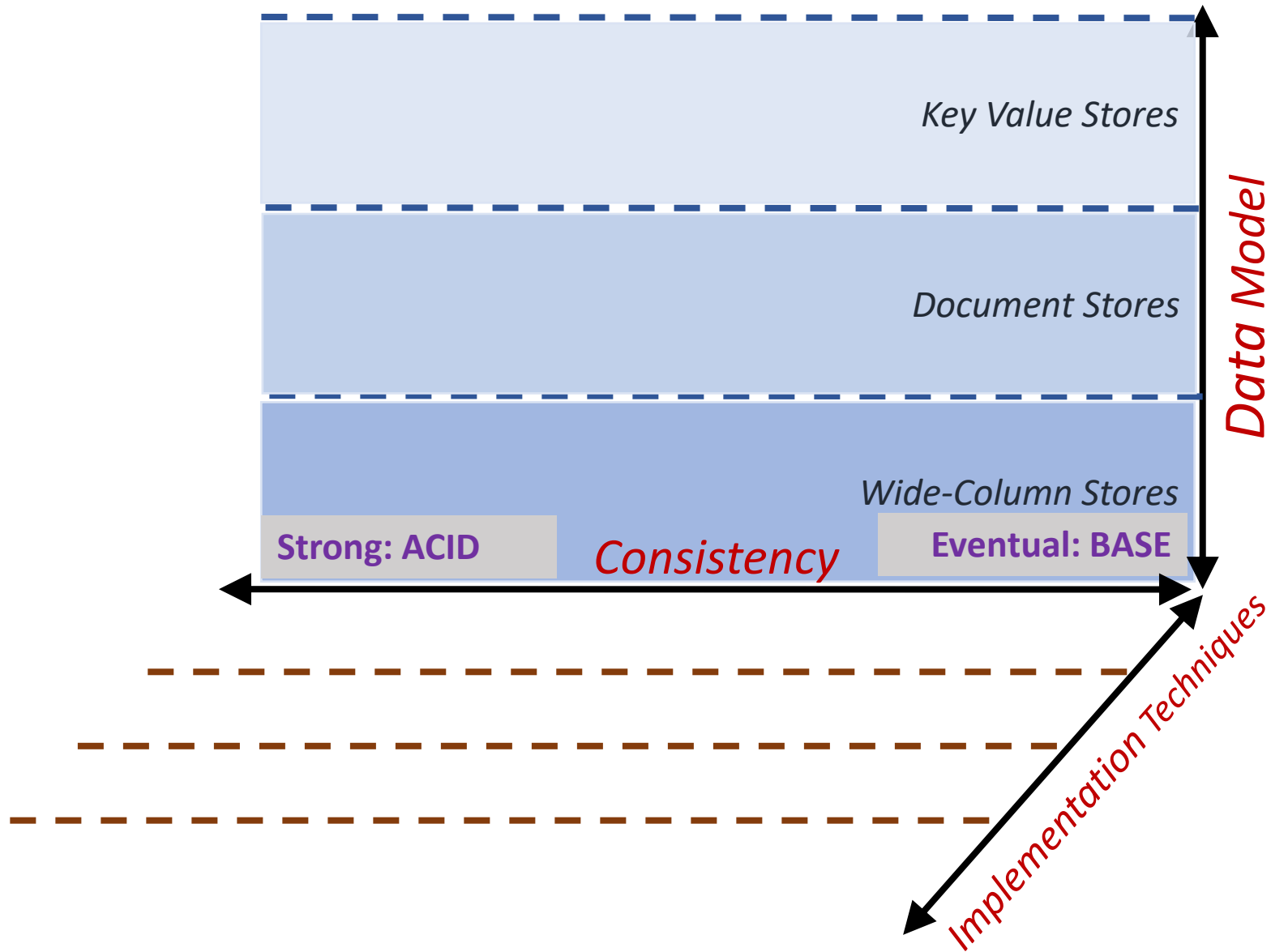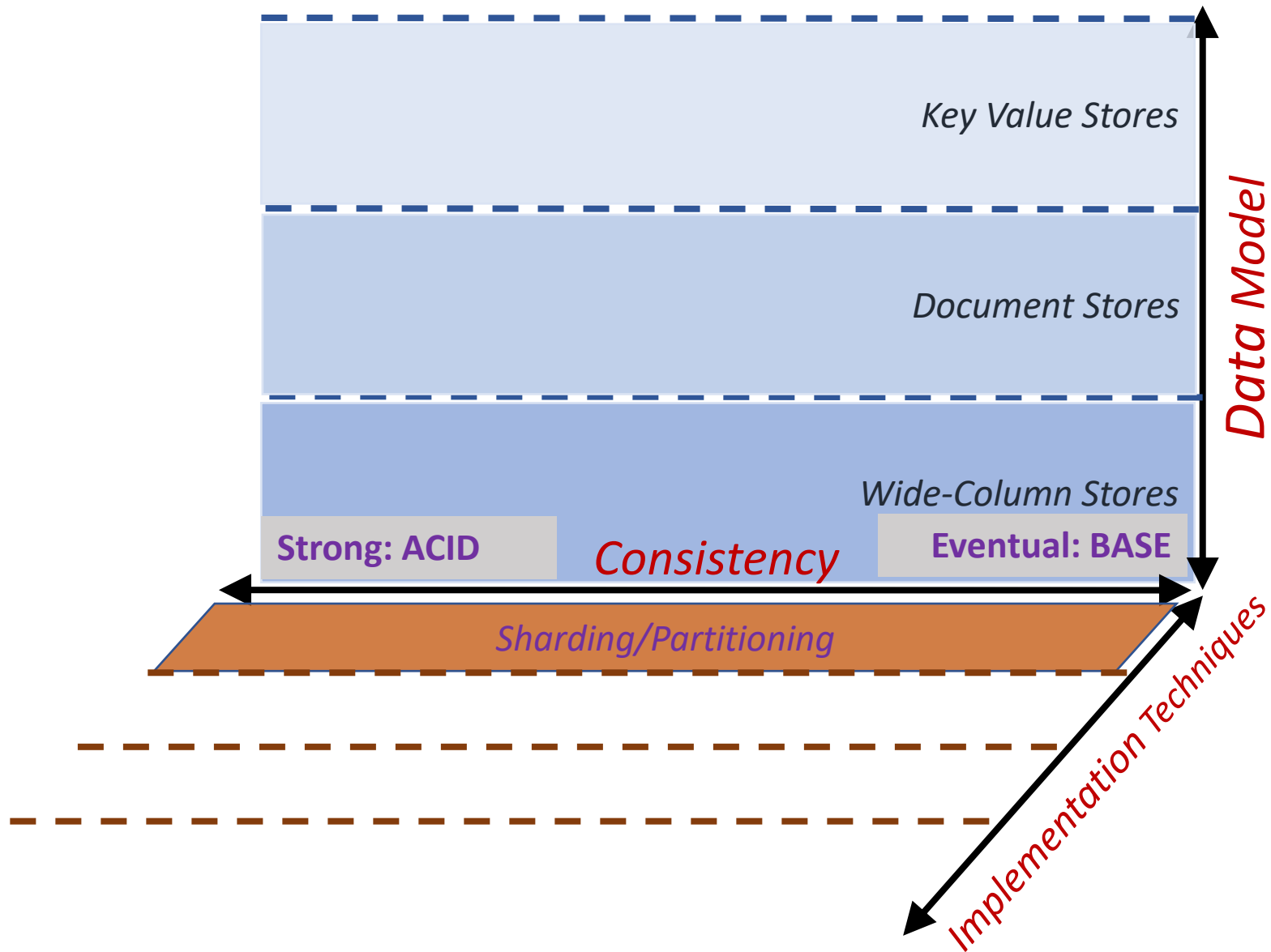# Consistency



| col | col | col$_2$ | . . . | col$_c$ |
|-----|-----|---------|-------|---------|
| 0 | 1 | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

# Consistency



| col | col | col$_2$ | . . . | col$_c$ |
|-----|-----|---------|-------|---------|
| 0 | 1 | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

# Consistency



| col | col | col$_2$ | . . . | col$_c$ |
|-----|-----|---------|-------|---------|
| 0   | 1   |         |       |         |
|     |     |         |       |         |
|     |     |         |       |         |
|     |     |         |       |         |
|     |     |         |       |         |

Partitions

# Consistency



Partitions

# Consistency



Partitions

# Consistency



| col | col | col$_2$ | . . . | col$_c$ |
|-----|-----|---------|-------|---------|
| 0 | 1 | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

- Clients perform reads and writes

Partitions

# Consistency

| col | col | col$_2$ | . . . | col$_c$ |
|-----|-----|---------|-------|---------|
| 0   | 1   |         |       |         |
|     |     |         |       |         |
|     |     |         |       |         |
|     |     |         |       |         |
|     |     |         |       |         |

Partitions

- Clients perform reads and writes
- Data is replicated among a set of servers

# Consistency



- Clients perform reads and writes
- Data is replicated among a set of servers
- Writes must be performed at all servers

Partitions

# Consistency



Partitions

- Clients perform reads and writes
- Data is replicated among a set of servers
- Writes must be performed at all servers
- Reads return the result of one or more past writes

# Consistency



Partitions

- Clients perform reads and writes
- Data is replicated among a set of servers
- Writes must be performed at all servers
- Reads return the result of one or more past writes
- How to keep data in sync?

# Consistency



| col | col | col$_2$ | . . . | col$_c$ |
|---|---|---|---|---|
| 0 | 1 | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

Partitions

- Clients perform reads and writes
- Data is replicated among a set of servers
- Writes must be performed at all servers
- Reads return the result of one or more past writes
- How to keep data in sync?

Consistency != Correctess
- consistency: no internal contradictions
- Correct: higher-level property
- Inconsistency → code does wrong things

# Consistency: CAP Theorem

# Consistency: CAP Theorem



- A distributed system can satisfy at most 2/3 guarantees of:

# Consistency: CAP Theorem



- A distributed system can satisfy at most 2/3 guarantees of:
  1. **Consistency**:

# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:

  1. **Consistency**:
     - all nodes see same data at any time
     - or reads return latest written value by any client

# Consistency: CAP Theorem



- A distributed system can satisfy at most 2/3 guarantees of:
  1. **Consistency**:
     - all nodes see same data at any time
     - or reads return latest written value by any client
  2. **Availability**:

# Consistency: CAP Theorem



- A distributed system can satisfy at most 2/3 guarantees of:

  1. **Consistency**:
     - all nodes see same data at any time
     - or reads return latest written value by any client

  2. **Availability**:
     - system allows operations all the time,
     - and operations return quickly

# Consistency: CAP Theorem



- A distributed system can satisfy at most 2/3 guarantees of:
  1. **Consistency**:
     - all nodes see same data at any time
     - or reads return latest written value by any client
  2. **Availability**:
     - system allows operations all the time,
     - and operations return quickly
  3. **Partition-tolerance**:

# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:

    1. **Consistency**:
        - all nodes see same data at any time
        - or reads return latest written value by any client

    2. **Availability**:
        - system allows operations all the time,
        - and operations return quickly

    3. **Partition-tolerance**:
        - system continues to work in spite of network partitions

# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:

   1. **Consistency**:
      - all nodes see same data at any time
      - or reads return latest written value by any client

   2. **Availability**:
      - system allows operations all the time,
      - and operations return quickly

   3. **Partition-tolerance**:
      - system continues to work in spite of netwo

# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:

  1. **Consistency**:
     - all nodes see same data at any time
     - or reads return latest written value by any client

  2. **Availability**:
     - system allows operations all the time,
     - and operations return quickly

  3. **Partition-tolerance**:
     - system continues to work in spite of netwo

**Why care about CAP Properties?**
**Availability**
  - Reads/writes complete reliably and quickly.
  - E.g. Amazon, each ms latency → $6M yearly loss.

**Partitions**
  - Internet router outages
  - Under-sea cables cut
  - rack switch outage
  - *system should continue functioning normally!*

**Consistency**
  - all nodes see same data at any time, or reads return latest written value by any client.
  - ***This basically means correctness!***

# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:
    1. **Consistency**:
        - all nodes see same data at any time
        - or reads return latest written value by any client
    2. **Availability**:
        - system allows operations all the time,
        - and operations return quickly
    3. **Partition-tolerance**:
        - system continues to work in spite of netwo
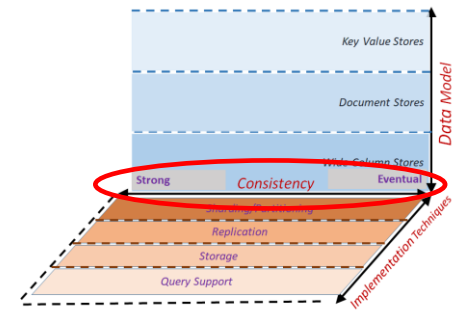
# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:
  1. **Consistency**:
     - all nodes see same data at any time
     - or reads return latest written value by any client
  2. **Availability**:
     - system allows operations all the time,
     - and operations return quickly
  3. **Partition-tolerance**:
     - system continues to work in spite of netwo

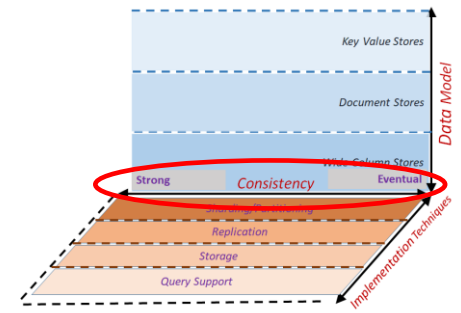**Why is this "theorem" true?**

# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:
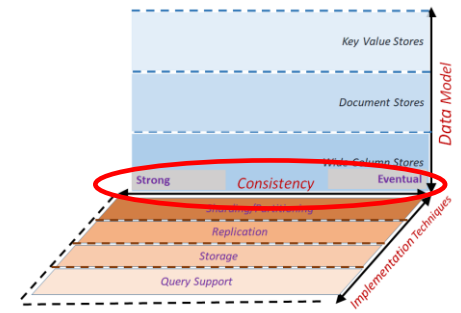
    1. **Consistency**:
        - all nodes see same data at any time
        - or reads return latest written value by any client

    2. **Availability**:
        - system allows operations all the time,
        - and operations return quickly

    3. **Partition-tolerance**:
        - system continues to work in spite of netwo

**Why is this "theorem" true?**

# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:
    1. **Consistency**:
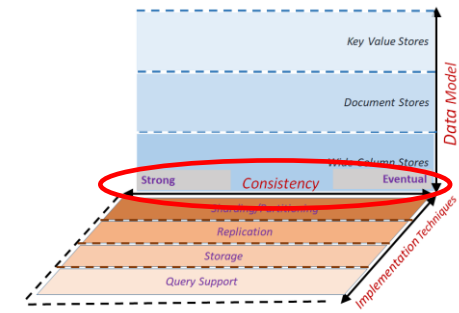        - all nodes see same data at any time
        - or reads return latest written value by any client
    2. **Availability**:
        - system allows operations all the time,
        - and operations return quickly
    3. **Partition-tolerance**:
        - system continues to work in spite of netwo

**Why is this "theorem" true?**

# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:

    1. **Consistency**:
        - all nodes see same data at any time
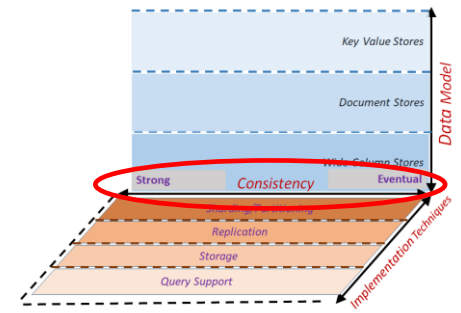        - or reads return latest written value by any client

    2. **Availability**:
        - system allows operations all the time,
        - and operations return quickly

    3. **Partition-tolerance**:
        - system continues to work in spite of netwo

**Why is this "theorem" true?**

Write(k,v) → writer → R1 | R2 → Read(k,v) → reader

if(partition) { keep going } → !consistent && available

# Consistency: CAP Theorem

- A distributed system can satisfy at most 2/3 guarantees of:

  1. **Consistency**:
     - all nodes see same data at any time
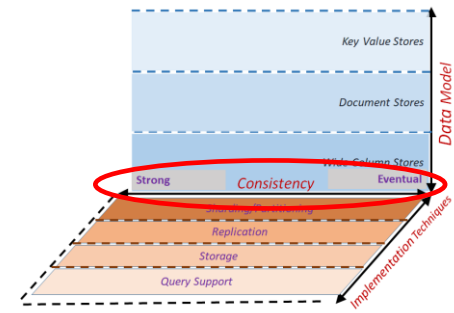     - or reads return latest written value by any client

  2. **Availability**:
     - system allows operations all the time,
     - and operations return quickly

  3. **Partition-tolerance**:
     - system continues to work in spite of netwo

**Why is this "theorem" true?**

if(partition) { keep going } → !consistent && available

if(partition) { stop } → consistent && !available

# CAP Implications

- A distributed storage system can achieve at most two of C, A, and P.

- When partition-tolerance is important, you have to choose between consistency and availability

**Consistency**

*HBase, HyperTable, BigTable, Spanner*

*RDBMSs (non-replicated)*



**Partition-tolerance** **Availability**

*Cassandra, RIAK, Dynamo, Voldemort*

# CAP Implications

- A distributed storage system can achieve at most two of C, A, and P.

- When partition-tolerance is important, you have to choose between consistency and availability



**Consistency**

*HBase, HyperTable, BigTable, Spanner*

*RDBMSs (non-replicated)*

**Partition-tolerance** **Availability**

*Cassandra, RIAK, Dynamo, Voldemort*

CAP is flawed

# CAP Implications

- A distributed storage system can achieve at most two of C, A, and P.

- When partition-tolerance is important, you have to choose between consistency and availability

**Consistency**

*HBase, HyperTable, BigTable, Spanner*

*RDBMSs (non-replicated)*

**Partition-tolerance  Availability**

*Cassandra, RIAK, Dynamo, Voldemort*

**PACELC:**

if(partition) {
    choose A or C
} else {
    choose latency or consistency
}

CAP is flawed

# Consistency Spectrum



- Eventual Consistency
  - If writes to a key stop, all replicas of key will converge
  - Originally from Amazon's Dynamo and LinkedIn's Voldemort systems

BASE:

- Basically Available

- Soft State

- Eventually Consistent

- **Strict:**
  - Absolute time ordering of all shared accesses, reads always return last write

- **Linearizability**:
  - Each operation is visible (or available) to all other clients in real-time order

- **Sequential Consistency** [Lamport]:
  - *"... the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program.*
  - After the fact, find a "reasonable" ordering of the operations (can re-order operations) that obeys sanity (consistency) at all clients, and across clients.

- **ACID** properties



*Faster reads and writes*

*More consistency*

Eventual → Strong (e.g., Sequential)

# Consistency Spectrum

- Eventual Consistency
  - If writes to a key stop, all replicas of key will converge
  - Originally from Amazon's Dynamo and LinkedIn's Voldemort systems

BASE:

- Basically Available

- Soft State

- Eventually Consistent

- **Strict:**
  - Absolute time ordering of all shared accesses, reads always return last write

- **Linearizability**:
  - Each operation is visible (or available) to all other clients in real-time order

- **Sequential Consistency** [Lamport]:
  - *"... the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program.*
  - After the fact, find a "reasonable" ordering of the operations (can re-order operations) that obeys sanity (consistency) at all clients, and across clients.

- **ACID** properties

**Eventual: BASE** ← 

*Faster reads and writes*

*More consistency* →

Eventual

**Strong: ACID**

Strong
(e.g., Sequential)

# Sequential Consistency

- weaker than strict/strong consistency
    - All operations are executed in *some* sequential order
    - each process issues operations in program order
        - Any valid interleaving is allowed
        - All  agree on the same interleaving
        - Each process preserves its program order

```
P1: W(x)a                               P1:  W(x)a
P2:        W(x)b                        P2:        W(x)b
P3:                R(x)b      R(x)a     P3:                R(x)b        R(x)a
P4:                   R(x)b  R(x)a      P4:                      R(x)a  R(x)b
              (a)                                       (b)
```

# Sequential Consistency

- weaker than strict/strong consistency
  - All operations are executed in *some* sequential order
  - each process issues operations in program order
    - Any valid interleaving is allowed
    - All agree on the same interleaving
    - Each process preserves its program order

```
P1: W(x)a                           P1: W(x)a
P2:       W(x)b                     P2:       W(x)b
P3:             R(x)b     R(x)a     P3:             R(x)b     R(x)a
P4:                  R(x)b  R(x)a   P4:                  R(x)a  R(x)b
                                                (b)
```

- ***Why is this weaker than strict/strong?***

# Sequential Consistency

- weaker than strict/strong consistency
  - All operations are executed in *some* sequential order
  - each process issues operations in program order
    - Any valid interleaving is allowed
    - All agree on the same interleaving
    - Each process preserves its program order

```
P1: W(x)a                          P1: W(x)a
P2:      W(x)b                      P2:      W(x)b
P3:           R(x)b    R(x)a        P3:           R(x)b    R(x)a
P4:           R(x)b  R(x)a          P4:           R(x)a  R(x)b
           (a)                              (b)
```

- *Why is this weaker than strict/strong?*
- *Nothing is said about "most recent write"*

# Causal consistency

# Causal consistency

- Causally related writes seen by all processes in same order.

# Causal consistency

- Causally related writes seen by all processes in same order.
  - *Causally?*

# Causal consistency

- Causally related writes seen
  - *Causally?*

**Causal:**
If a write produces a value that
causes another write, they are causally related

```
X = 1
if(X > 0) {
    Y = 1
}
```
Causal consistency → all see X=1, Y=1 in same order

# Causal consistency

- Causally related writes seen by all processes in same order.
  - *Causally?*

# Causal consistency

- Causally related writes seen by all processes in same order.
  - *Causally?*
  - *Concurrent* writes may be seen in different orders on different machines

# Causal consistency

- Causally related writes seen by all processes in same order.
  - *Causally?*
  - *Concurrent* writes may be seen in different orders on different machines

```
P1: W(x)a
P2:          R(x)a       W(x)b
P3:                              R(x)b    R(x)a
P4:                              R(x)a    R(x)b
                (a)
```

# Causal consistency

- Causally related writes seen by all processes in same order.
  - *Causally?*
  - *Concurrent* writes may be seen in different orders on different machines

```
P1: W(x)a
P2:          R(x)a      W(x)b
P3:                              R(x)b    R(x)a
P4:                              R(x)a    R(x)b
                  (a)
```

Not permitted

# Causal consistency

- Causally related writes seen by all processes in same order.
  - *Causally?*
  - *Concurrent* writes may be seen in different orders on different machines

| P1: W(x)a | | | | |
|-----------|--------|--------|--------|--------|
| P2:       | R(x)a  | W(x)b  |        |        |
| P3:       |        |        | R(x)b  | R(x)a  |
| P4:       |        |        | R(x)a  | R(x)b  |

(a)

| P1: W(x)a | | | |
|-----------|--------|--------|--------|
| P2:       | W(x)b  |        |        |
| P3:       |        | R(x)b  | R(x)a  |
| P4:       |        | R(x)a  | R(x)b  |

(b)

Not permitted

# Causal consistency

- Causally related writes seen by all processes in same order.
    - *Causally?*
    - *Concurrent* writes may be seen in different orders on different machines

| P1: W(x)a | | | | |
|---|---|---|---|---|
| P2: | R(x)a | W(x)b | | |
| P3: | | | R(x)b | R(x)a |
| P4: | | | R(x)a | R(x)b |

(a)

Not permitted

| P1: W(x)a | | | |
|---|---|---|---|
| P2: | W(x)b | | |
| P3: | | R(x)b | R(x)a |
| P4: | | R(x)a | R(x)b |

(b)

Permitted

# Linearizability vs. Serializability

http://www.bailis.org/blog/linearizability-versus-serializability/

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency

http://www.bailis.org/blog/linearizability-versus-serializability/

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
    - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
    - Stronger than sequential consistency

- Difference between linearizability and serializability?
    - Granularity: reads/writes versus transactions

http://www.bailis.org/blog/linearizability-versus-serializability/

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If $TS(x) < TS(y)$ then $OP(x)$ should precede $OP(y)$ in the sequence
  - Stronger than sequential consistency
- Difference between linearizability and serializability?
  - Granularity: reads/writes versus transactions

Linearizability:

http://www.bailis.org/blog/linearizability-versus-serializability/

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency
- Difference between linearizability and serializability?
  - Granularity: reads/writes versus transactions

Linearizability:
- Single-operation, single-object, real-time order

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency
- Difference between linearizability and serializability?
  - Granularity: reads/writes versus transactions

Linearizability:
- Single-operation, single-object, real-time order
- Talks about order of ops on single object (e.g. atomic register)

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency
- Difference between linearizability and serializability?
  - Granularity: reads/writes versus transactions

Linearizability:

•Single-operation, single-object, real-time order

•Talks about order of ops on single object (e.g. atomic register)

•Ops should appear instantaneous, reflect real time order

http://www.bailis.org/blog/linearizability-versus-serializability/

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency
- Difference between linearizability and serializability?
  - Granularity: reads/writes versus transactions

Linearizability:
- Single-operation, single-object, real-time order
- Talks about order of ops on single object (e.g. atomic register)
- Ops should appear instantaneous, reflect real time order

Serializability:

http://www.bailis.org/blog/linearizability-versus-serializability/

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency
- Difference between linearizability and serializability?
  - Granularity: reads/writes versus transactions

Linearizability:
- Single-operation, single-object, real-time order
- Talks about order of ops on single object (e.g. atomic register)
- Ops should appear instantaneous, reflect real time order

Serializability:
- Talks about groups of 1 or more ops on one or more objects

http://www.bailis.org/blog/linearizability-versus-serializability/

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency
- Difference between linearizability and serializability?
  - Granularity: reads/writes versus transactions

Linearizability:

•Single-operation, single-object, real-time order

•Talks about order of ops on single object (e.g. atomic register)

•Ops should appear instantaneous, reflect real time order

Serializability:

•Talks about groups of 1 or more ops on one or more objects

•Txns over multiple items equivalent to serial order of txns

http://www.bailis.org/blog/linearizability-versus-serializability/

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency
- Difference between linearizability and serializability?
  - Granularity: reads/writes versus transactions

Linearizability:
- Single-operation, single-object, real-time order
- Talks about order of ops on single object (e.g. atomic register)
- Ops should appear instantaneous, reflect real time order

Serializability:
- Talks about groups of 1 or more ops on one or more objects
- Txns over multiple items equivalent to serial order of txns
- Only requires *some* equivalent serial order

http://www.bailis.org/blog/linearizability-versus-serializability/

# Linearizability vs. Serializability

- Linearizability assumes sequential consistency *and*
  - If TS(x) < TS(y) then OP(x) should precede OP(y) in the sequence
  - Stronger than sequential consistency
- Difference between linearizability and serializability?
  - Granularity: reads/writes versus transactions

Linearizability:
- Single-operation, single-object, real-time order
- Talks about order of ops on single object (e.g. atomic register)
- Ops should appear instantaneous, reflect real time order

Serializability:
- Talks about groups of 1 or more ops on one or more objects
- Txns over multiple items equivalent to serial order of txns
- Only requires *some* equivalent serial order

Serializability + Linearizability == "Strict Serializability"
- Txn order equivalent to some serial order *that respects real time order*
- Linearizability: degenerate case of Strict Ser: txns are single op single object

http://www.bailis.org/blog/linearizability-versus-serializability/

# Some Consistency Guarantees

| | |
|---|---|
| Strong Consistency | See all previous writes. |
| Eventual Consistency | See subset of previous writes. |
| Consistent Prefix | See initial sequence of writes. |
| Bounded Staleness | See all "old" writes. |
| Monotonic Reads | See increasing subset of writes. |
| Read My Writes | See all writes performed by reader. |

# Some Consistency Guarantees

|  |  | consistency | performance | availability |
|---|---|:---:|:---:|:---:|
| Strong Consistency | See all previous writes. | A | D | F |
| Eventual Consistency | See subset of previous writes. | D | A | A |
| Consistent Prefix | See initial sequence of writes. | C | B | A |
| Bounded Staleness | See all "old" writes. | B | C | D |
| Monotonic Reads | See increasing subset of writes. | C | B | B |
| Read My Writes | See all writes performed by reader. | C | C | C |

# NoSQL faux quiz:

- What is the CAP theorem? What does "PACELC" stand for and how does it relate to CAP?
- What is the difference between ACID and BASE?
- Why do NoSQL systems claim to be more horizontally scalable than RDMBSes? List some features NoSQL systems give up toward this goal?
- What is eventual consistency? Give a concrete example of how of why it causes a complex programming model (relative to a strongly consistent model).
- Compare and contrast Key-Value, Document, and Wide-column Stores
- Define and contrast the following consistency properties:
  - strong consistency, eventual consistency, consistent prefix, monotonic reads, read-my-writes, bounded staleness

# NoSQL faux quiz:

- What is the CAP theorem? What does "PACELC" stand for and how does it relate to CAP?
- What is the difference between ACID and BASE?
- Why do NoSQL systems claim to be more horizontally scalable than RDMBSes? ~~List some features NoSQL systems give up toward this goal?~~
- What is eventual consistency? Give a concrete example of how of why it causes a complex programming model (relative to a strongly consistent model).
- ~~Compare and contrast Key-Value, Document, and Wide-column Stores~~
- Define and contrast the following consistency properties:
  - strong consistency, eventual consistency, consistent prefix, monotonic reads, read-my-writes, bounded staleness

# Dataflow

# Dataflow

- MR is a *dataflow* engine

# Dataflow

- MR is a ***dataflow*** engine

# Dataflow



- MR is a *dataflow* engine

- So are Lots of others
  - Dryad
  - DryadLINQ
  - Dandelion
  - CIEL
  - GraphChi/PowerGraph/Pregel
  - Spark

# Dataflow

- MR is a *dataflow* engine
- So are Lots of others
  - Dryad
  - DryadLINQ
  - Dandelion
  - CIEL
  - GraphChi/PowerGraph/Pregel
  - Spark

# Spark faux quiz (5 min, any 2):

- What is the difference between *transformations* and *actions* in Spark?

- Spark supports a persist API. When should a programmer want to use it? When should she [not] use use the "*RELIABLE*" flag?

- Compare and contrast fault tolerance guarantees of Spark to those of MapReduce. How are[n't] the mechanisms different?

- Is Spark a good system for indexing the web? For computing page rank over a web index? Why [not]?

- List aspects of Spark's design that help/hinder multi-core parallelism relative to MapReduce. If the issue is orthogonal, explain why.

# Collections and Iterators

class Collection<T> : IEnumerable<T>;

# Collections and Iterators

class Collection<T> : IEnumerable<T>;



```
public interface IEnumerable<T>  {
        IEnumerator<T> GetEnumerator();
}
```
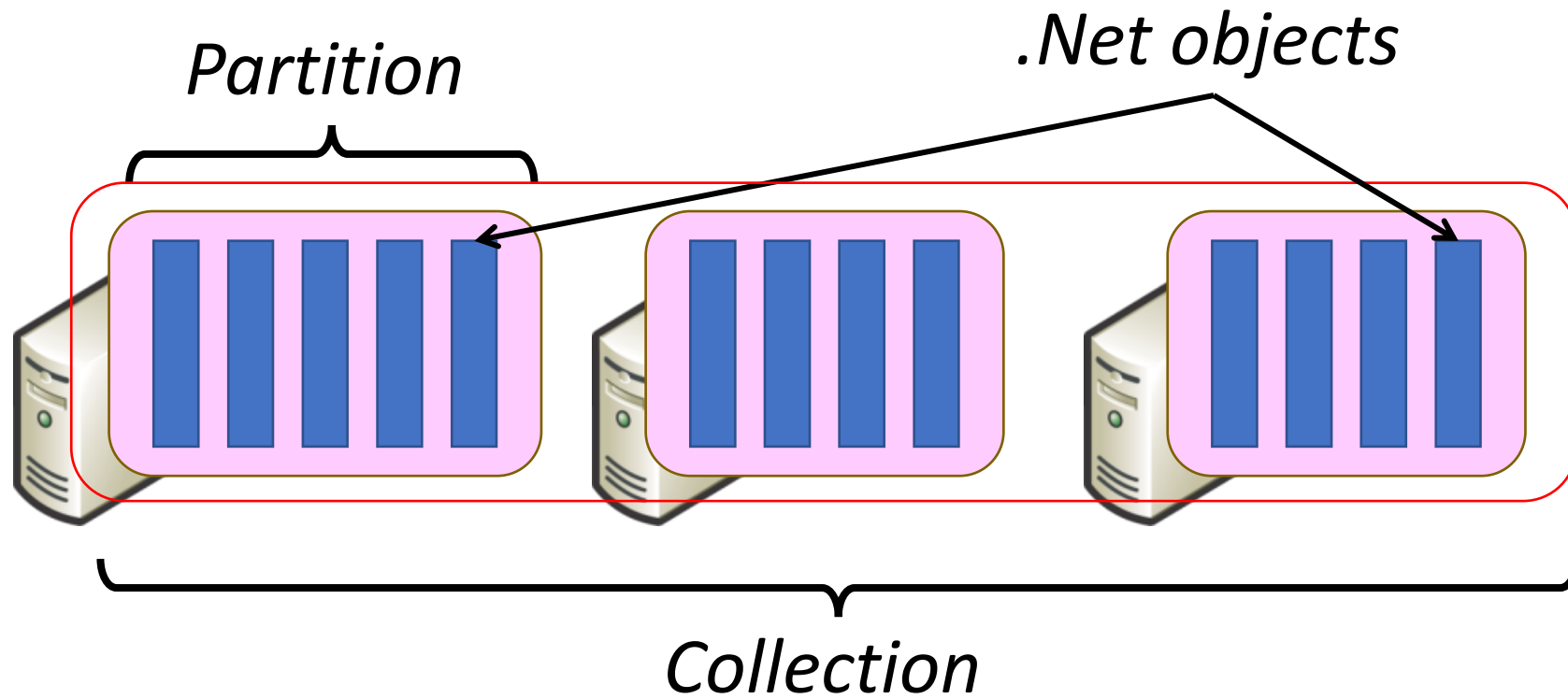
# Collections and Iterators

class Collection<T> : IEnumerable<T>;



```
public interface IEnumerable<T>  {
        IEnumerator<T> GetEnumerator();
}
```

```
public interface IEnumerator <T> {
        T Current { get; }
        bool MoveNext();
        void Reset();
}
```

# Collections and Iterators

class Collection<T> : IEnumerable<T>;



```
public interface IEnumerable<T>  {
        IEnumerator<T> GetEnumerator();
}
```

```
public interface IEnumerator <T> {
        T Current { get; }
        bool MoveNext();
        void Reset();
}
```

# DryadLINQ Data Model



Partition

.Net objects

Collection

# DryadLINQ = LINQ + Dryad



```
Collection<T> collection;

bool IsLegal(Key k);
string Hash(Key);

var results = from c in collection
              where IsLegal(c.key)
              select new { Hash(c.key), c.value};
```
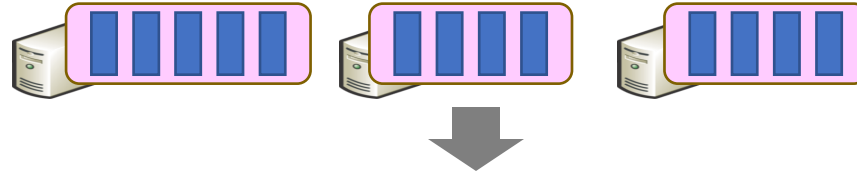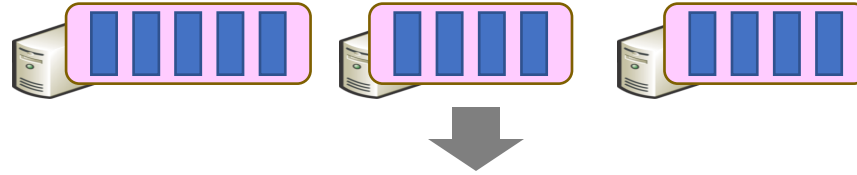
# DryadLINQ = LINQ + Dryad



```
Collection<T> collection;
bool IsLegal(Key k);
string Hash(Key);

var results = from c in collection
              where IsLegal(c.key)
              select new { Hash(c.key), c.value};
```

Data

collection

results

31

# DryadLINQ = LINQ + Dryad



```
Collection<T> collection;
bool IsLegal(Key k);
string Hash(Key);

var results = from c in collection
              where IsLegal(c.key)
              select new { Hash(c.key), c.value};
```

*Data*

*Query plan (Dryad job)*

*collection*

*results*

# DryadLINQ = LINQ + Dryad



Collection<T> collection;

bool IsLegal(Key k);
string Hash(Key);

var results = from c in collection
             where IsLegal(c.key)
             select new { Hash(c.key), c.value};

Vertex code

Data

Query plan (Dryad job)

C#    C#    C#    C#

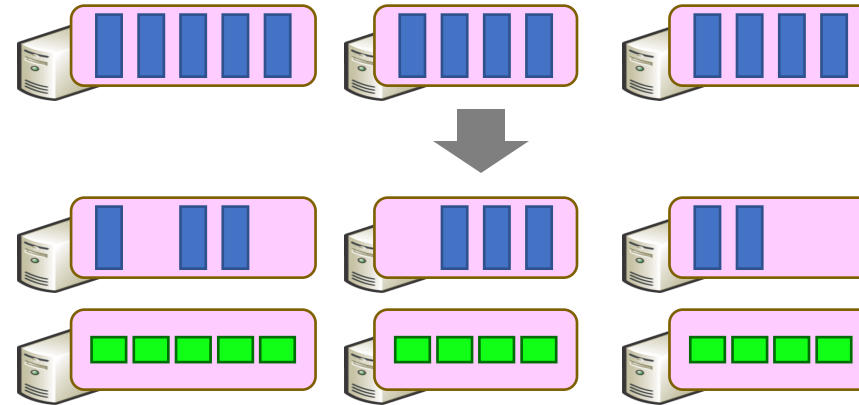collection

results

31

# Language Summary

# Language Summary



Where

# Language Summary

Where

# Language Summary

Where
Select

# Language Summary

Where
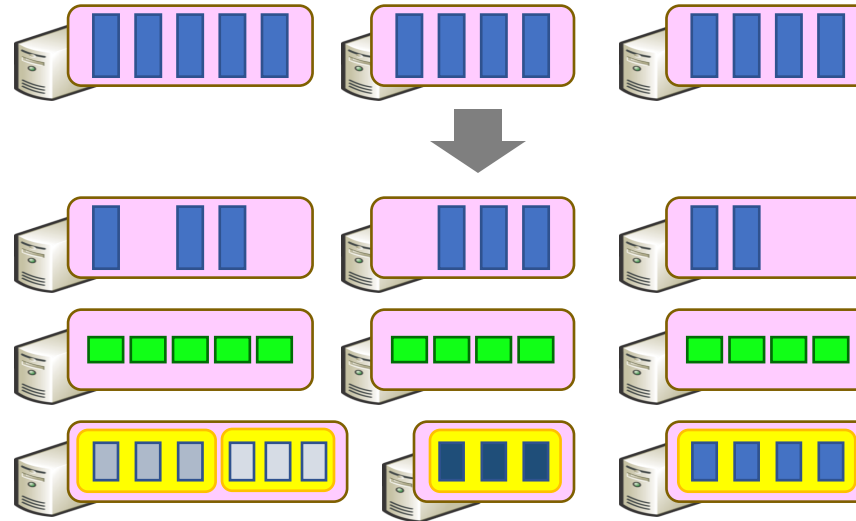Select

# Language Summary

Where
Select
GroupBy

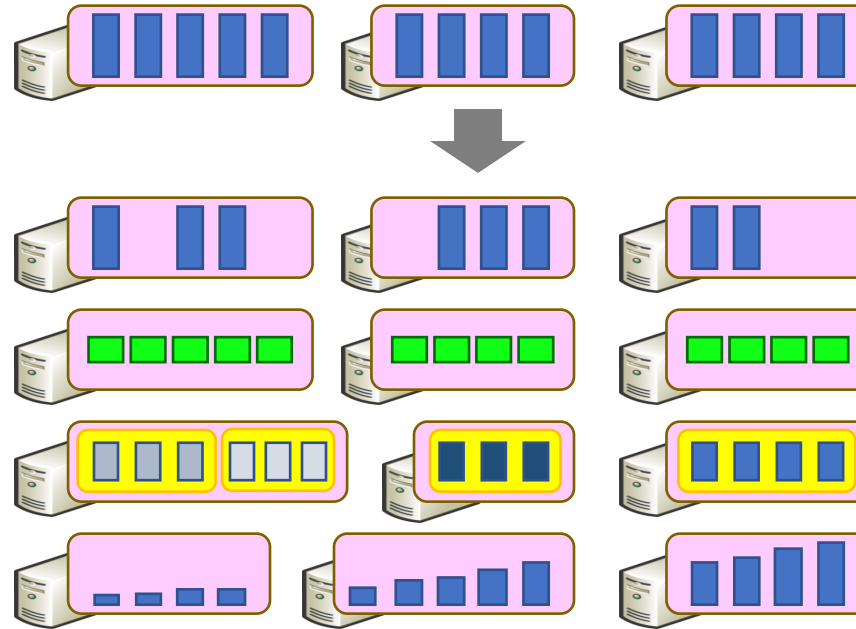# Language Summary

Where
Select
GroupBy

# Language Summary

Where
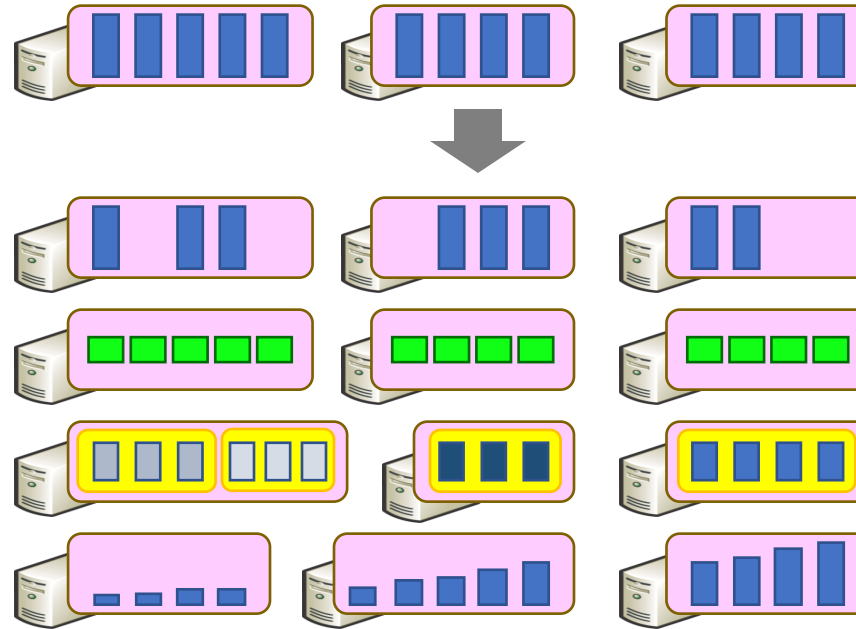Select
GroupBy
OrderBy

# Language Summary

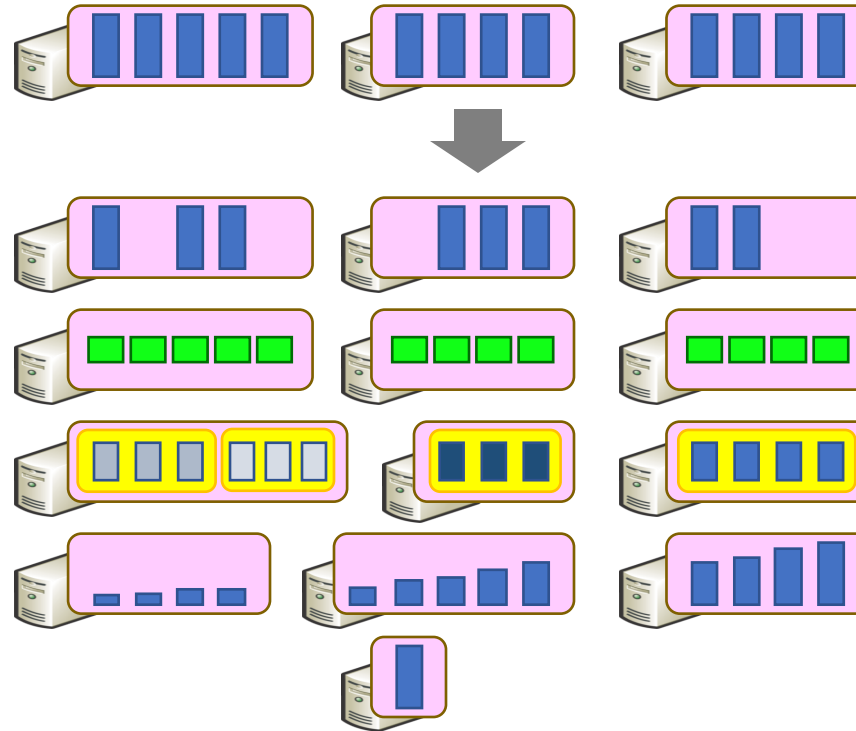Where
Select
GroupBy
OrderBy

# Language Summary

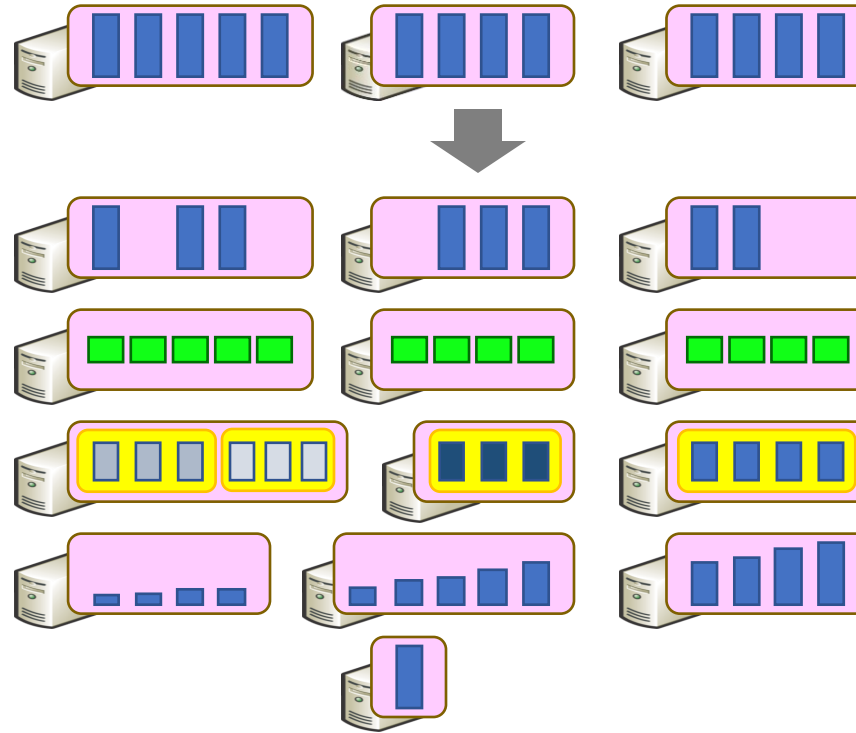Where
Select
GroupBy
OrderBy
Aggregate

# Language Summary

Where
Select
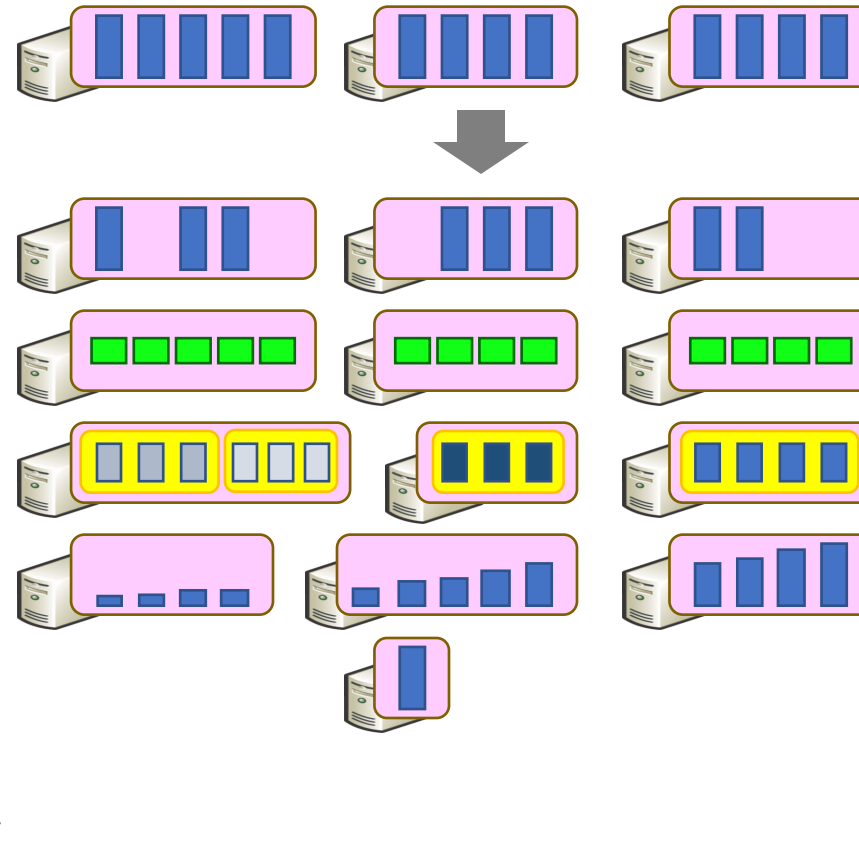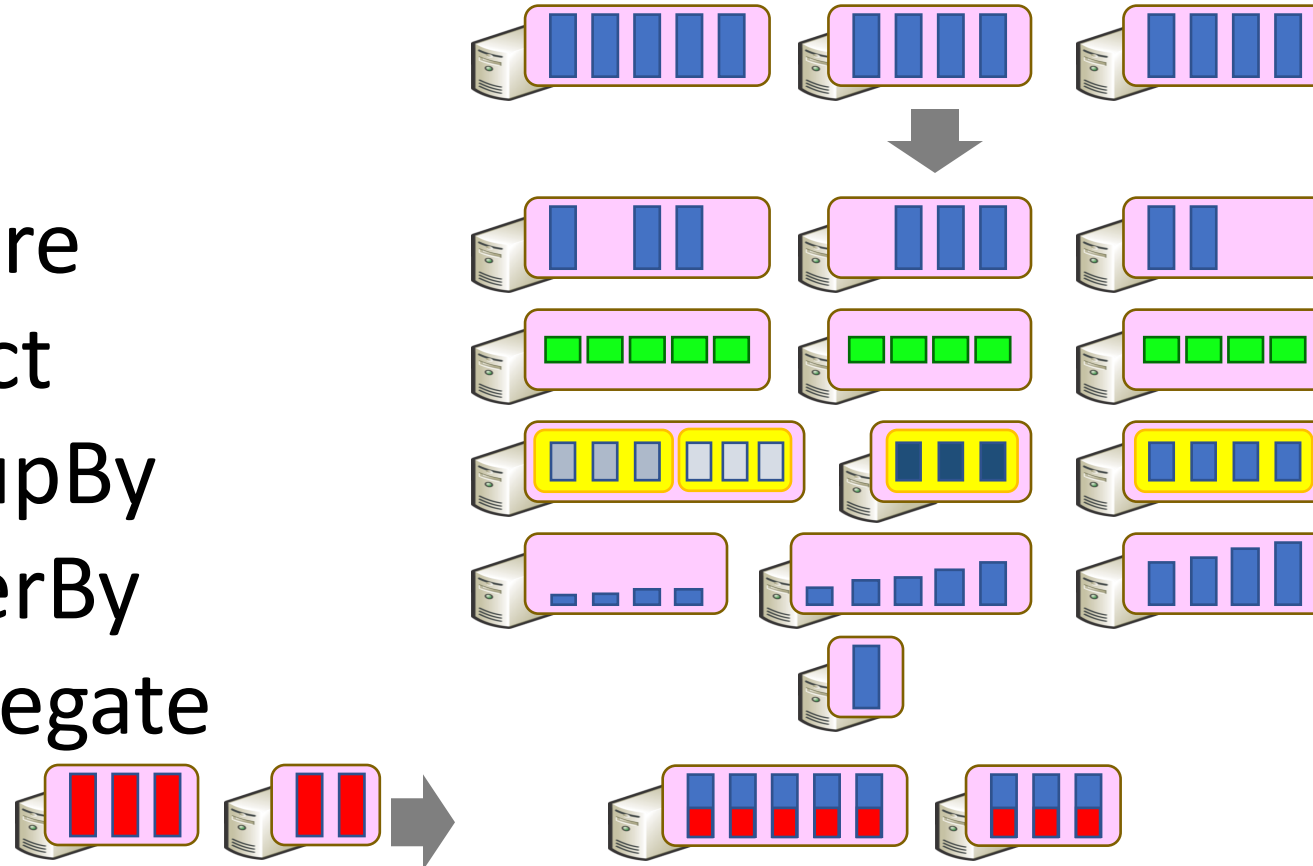GroupBy
OrderBy
Aggregate

# Language Summary

Where
Select
GroupBy
OrderBy
Aggregate
Join

# Language Summary

Where
Select
GroupBy
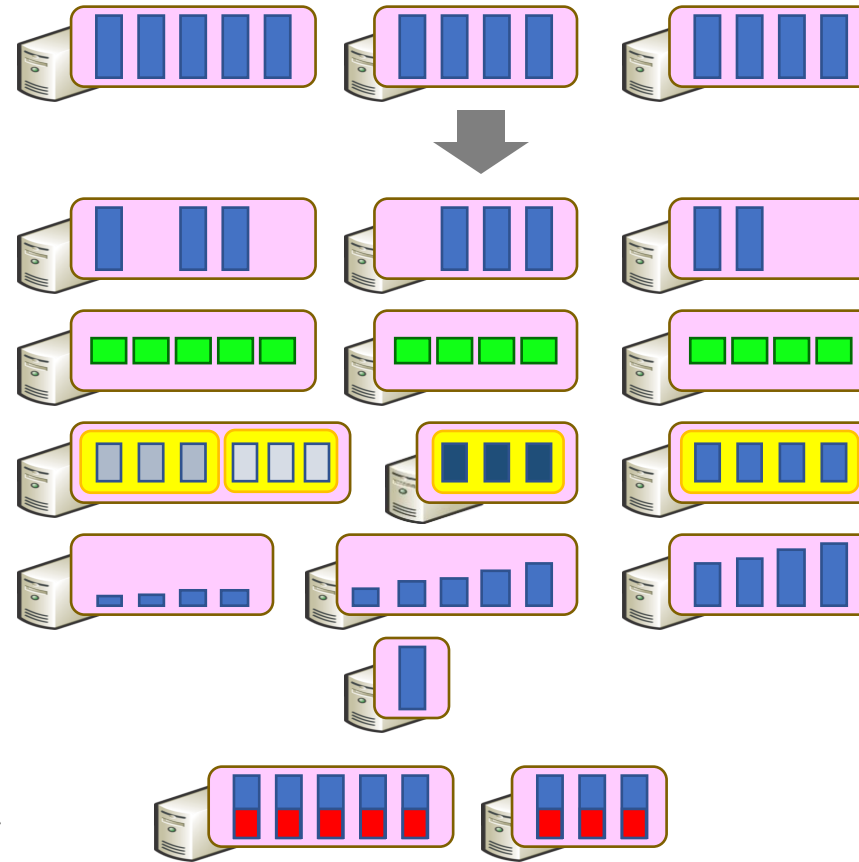OrderBy
Aggregate
Join

# Language Summary

Where
Select
GroupBy
OrderBy
Aggregate
Join

# Language Summary

Where
Select
GroupBy
OrderBy
Aggregate
Join
Apply

# Language Summary

Where
Select
GroupBy
OrderBy
Aggregate
Join
Apply

# Language Summary

Where
Select
GroupBy
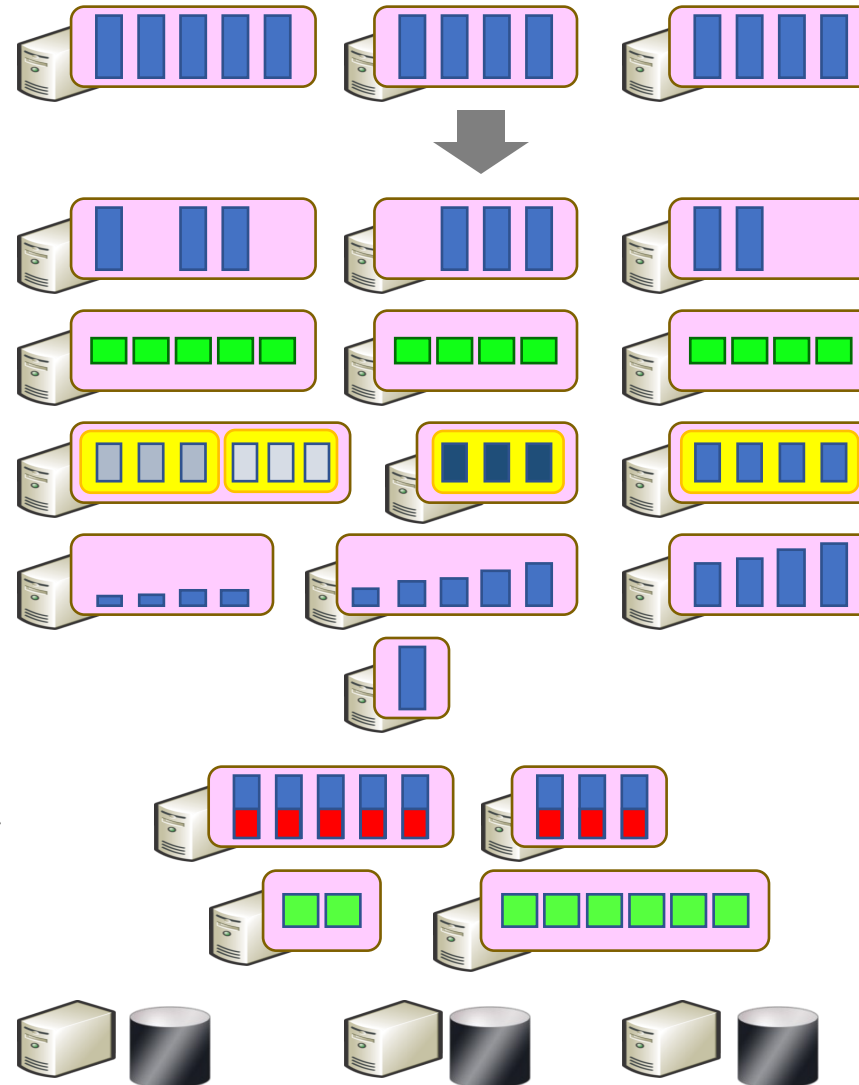OrderBy
Aggregate
Join
Apply
Materialize

# Language Summary

Where
Select
GroupBy
OrderBy
Aggregate
Join
Apply
Materialize

# Example: Histogram

```
public static IQueryable<Pair> Histogram(
    IQueryable<LineRecord> input, int k)
{
    var words = input.SelectMany(x => x.line.Split(' '));
    var groups = words.GroupBy(x => x);
    var counts = groups.Select(x => new Pair(x.Key, x.Count()));
    var ordered = counts.OrderByDescending(x => x.count);
    var top = ordered.Take(k);
    return top;
}
```

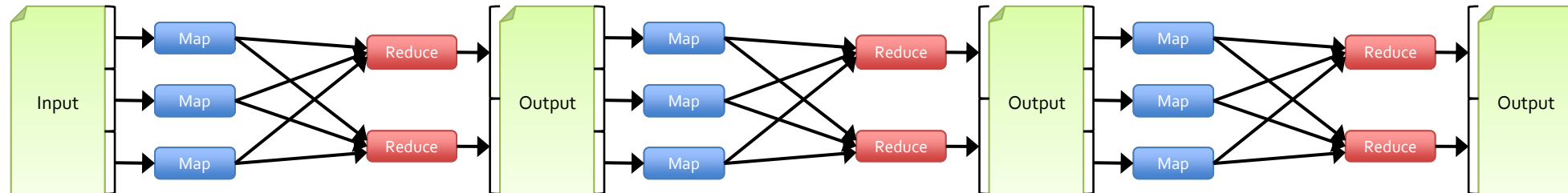| "A line of words of wisdom" |
| --- |
| ["A", "line", "of", "words", "of", "wisdom"] |
| [["A"], ["line"], ["of", "of"], ["words"], ["wisdom"]] |
| [ {"A", 1}, {"line", 1}, {"of", 2}, {"words", 1}, {"wisdom", 1}] |
| [{"of", 2}, {"A", 1}, {"line", 1}, {"words", 1}, {"wisdom", 1}] |
| [{"of", 2}, {"A", 1}, {"line", 1}] |

# Iterative Computations: PageRank

1. Start each page with a rank of 1
2. On each iteration, update each page's rank to

$$\sum_{i \in neighbors} rank_i / |neighbors_i|$$

```
links = // RDD of (url, neighbors) pairs
ranks = // RDD of (url, rank) pairs

for (i <- 1 to ITERATIONS) {
  ranks = links.join(ranks).flatMap {
    (url, (links, rank)) =>
      links.map(dest => (dest, rank/links.size))
  }.reduceByKey(_ + _)
}
```

# RDD Operations

| Transformations (define a new RDD) | Parallel operations (return a result to driver) |
|---|---|
| map<br>filter<br>sample<br>union<br>groupByKey<br>reduceByKey<br>join<br>persist/*cache*<br>… | reduce<br>collect<br>count<br>save<br>lookupKey<br>… |

# RDD Operations

| Transformations (define a new RDD) | Parallel operations (return a result to driver) |
|---|---|
| map<br>filter<br>sample<br>union<br>groupByKey<br>reduceByKey<br>join<br>persist/*cache*<br>… | reduce<br>collect<br>count<br>save<br>lookupKey<br>… |

Where
Select
GroupBy
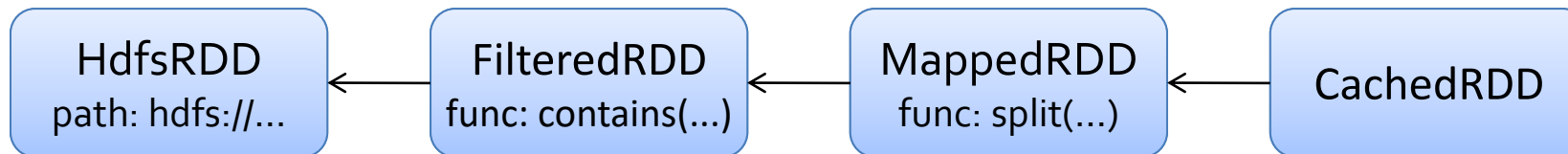OrderBy
Aggregate
Join
Apply
Materialize

# RDD Fault Tolerance

- RDDs maintain *lineage* information that can be used to reconstruct lost partitions

- Ex:

```
cachedMsgs = textFile(...).filter(_.contains("error"))
                          .map(_.split('\t')(2)
                          .persist()
```



| HdfsRDD<br>path: hdfs://... | ← | FilteredRDD<br>func: contains(...) | ← | MappedRDD<br>func: split(...) | ← | CachedRDD |

# RDDs vs Distributed Shared Memory

| Concern | RDDs | Distr. Shared Mem. |
|---------|------|--------------------|
| Reads | Fine-grained | Fine-grained |
| Writes | Bulk transformations | Fine-grained |
| Consistency | Trivial (immutable) | Up to app / runtime |
| Fault recovery | Fine-grained and low-overhead using lineage | Requires checkpoints and program rollback |
| Straggler mitigation | Possible using speculative execution | Difficult |
| Work placement | Automatic based on data locality | Up to app (but runtime aims for transparency) |