# **D**ata **I**ntensive **S**calable **C**omputing

## *Finding the Right Programming Models*

# Randal E. Bryant
# Carnegie Mellon University

`http://www.cs.cmu.edu/~bryant`

# Outline

**Data-Intensive Scalable Computing**

- **Focus on petabytes, not petaflops**
- **Rethinking machine design**

**Map/Reduce Programming Model**

- **Suitability for wide class of computing tasks**

**Strengths & Limitations**

- **Scalability**
- **Performance**

**Beyond Map/Reduce**

- **Small variations**
- **Other programming models**

# Our Data-Driven World

## Science

- **Data bases from astronomy, genomics, natural languages, seismic modeling, …**

## Humanities

- **Scanned books, historic documents, …**

## Commerce

- **Corporate sales, stock market transactions, census, airline traffic, …**

## Entertainment

- **Internet images, Hollywood movies, MP3 files, …**

## Medicine

- **MRI & CT scans, patient records, …**

# Why So Much Data?

## We Can Get It

- **Automation + Internet**

## We Can Keep It

- **Seagate Barracuda**
- **1.5 TB @ $76**
  5¢ / GB
  (vs. 40¢ in 2007)

## We Can Use It

- **Scientific computation**
- **Business applications**
- **Harvesting Internet data**

**Deal of the Day**

**Seagate Barracuda 1.5 TB Desktop Hard Drive**

Seagate Barracuda drives deliver reliable service and high performance. The Barracuda family has been developed and refined over 12 generations and embodies the highest levels of design for reliability and performance. While Barracuda LP is the low power leader, the entire family of Barracuda drives lead with... read more

| | |
|---|---|
| List Price: | ~~$199.99~~ |
| Yesterday's Price: | $94.99 |
| Today's Discount: | - $19.00 |
| **Gold Box Price:** | **$75.99 (62% off)** |

11 Comments | ★★★☆☆ ▽ (337)

✓Prime    Add to cart

# Oceans of Data, Skinny Pipes

No more blaming connection speeds for your losses.

Verizon FiOS – the fastest Internet available.

Plans as low $39.99/month (up to 5 Mbps).
Plus, order online & get your first month FREE!

Enter your home phone number below to check availability.

GO!

Don't have a Verizon phone number? Qualify your address.

Seagate

## 1 Terabyte

- Easy to store
- Hard to move

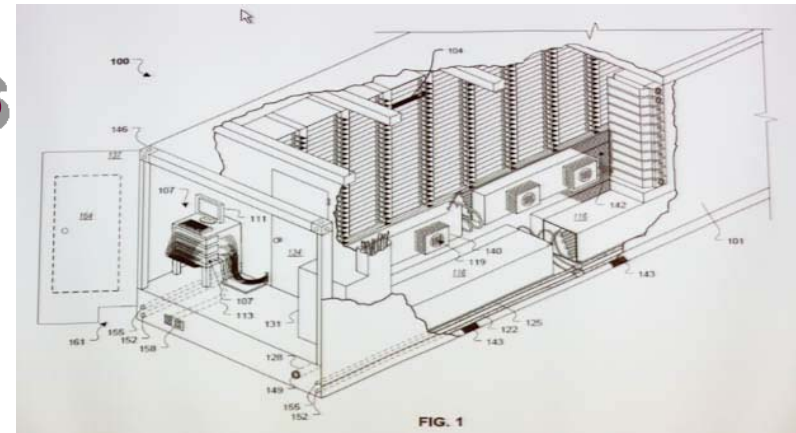| Disks | MB / s | Time |
|---|---|---|
| Seagate Barracuda | 115 | 2.3 hours |
| Seagate Cheetah | 125 | 2.2 hours |
| Networks | MB / s | Time |
| Home Internet | < 0.625 | > 18.5 days |
| Gigabit Ethernet | < 125 | > 2.2 hours |
| PSC Teragrid Connection | < 3,750 | > 4.4 minutes |

# Data-Intensive System Challenge

## For Computation That Accesses 1 TB in 5 minutes

- **Data distributed over 100+ disks**
  - Assuming uniform data partitioning
- **Compute using 100+ processors**
- **Connected by gigabit Ethernet (or equivalent)**

## System Requirements

- **Lots of disks**
- **Lots of processors**
- **Located in close proximity**
  - Within reach of fast, local-area network

# Google Data Centers





FIG. 1



## Dalles, Oregon

- **Hydroelectric power @ 2¢ / KW Hr**
- **50 Megawatts**
  - Enough to power 6,000 homes

- **Engineered for maximum modularity & power efficiency**
- **Container: 1160 servers, 250KW**
- **Server: 2 disks, 2 processors**

# High-Performance Distributed Computing: Two Versions

## Grid Computing

**Connect small number of big machines**

- Allow resource sharing among supercomputers

**Issues**

- Programs must usually be specialized for machine
- Hard to get cooperation from multiple organizations

## DISC System

**Build big system from many small machines**

- Use distributed systems principles to construct large-scale machine
- File system provides distribution, reliability, recovery
- Dynamically scheduled task as basic processing unit

# Programming Model Comparison

**Bulk Synchronous**

- **Commonly used for compute-intensive applications**

**Map/Reduce**

- **Commonly used for data-intensive applications**

**Issues**

- **Raw performance**
- **Scalability**
  - Cost required to increase machine size by K
    - » $\geq$ K
  - Performance achieved by increasing machine size by K
    - » $\leq$ K
  - Frequency and impact of failures

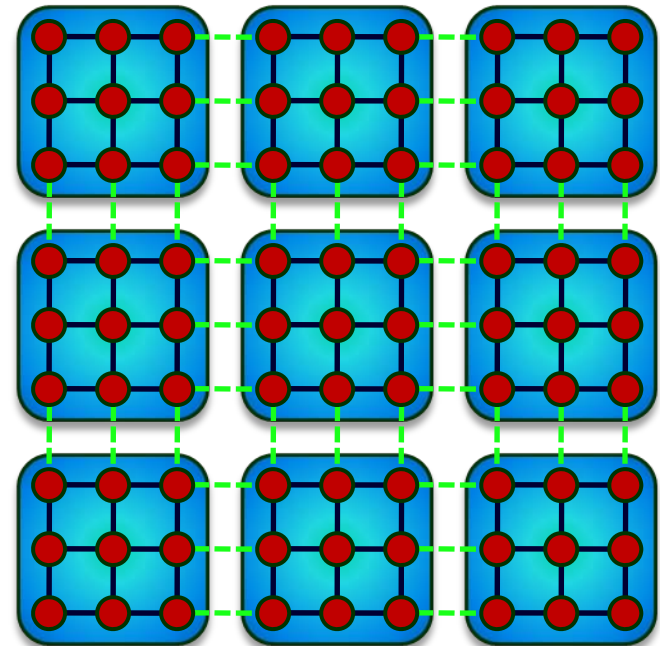# Bulk Synchronous Programming

## Solving Problem Over Grid

- **E.g., finite-element computation**

## Partition into Regions
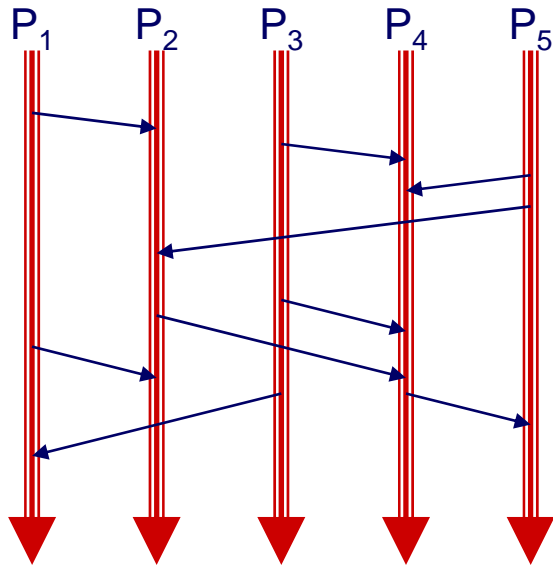
- **p regions for p processors**

## Map Region per Processor

- **Local computation sequential**
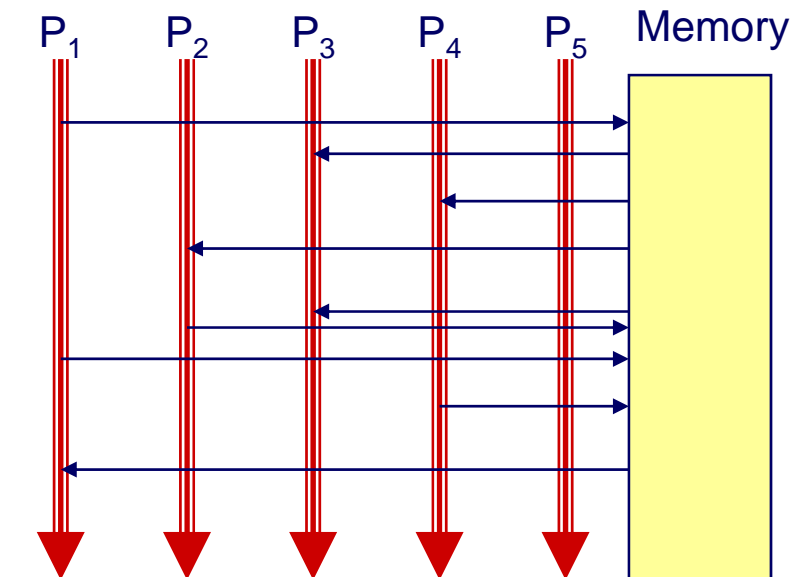- **Periodically communicate boundary values with neighbors**

# Typical HPC Operation

**Message Passing**

P$_1$   P$_2$   P$_3$   P$_4$   P$_5$

**Shared Memory**

P$_1$   P$_2$   P$_3$   P$_4$   P$_5$   Memory

## Characteristics

- **Long-lived processes**
- **Make use of spatial locality**
- **Hold all program data in memory (no disk access)**
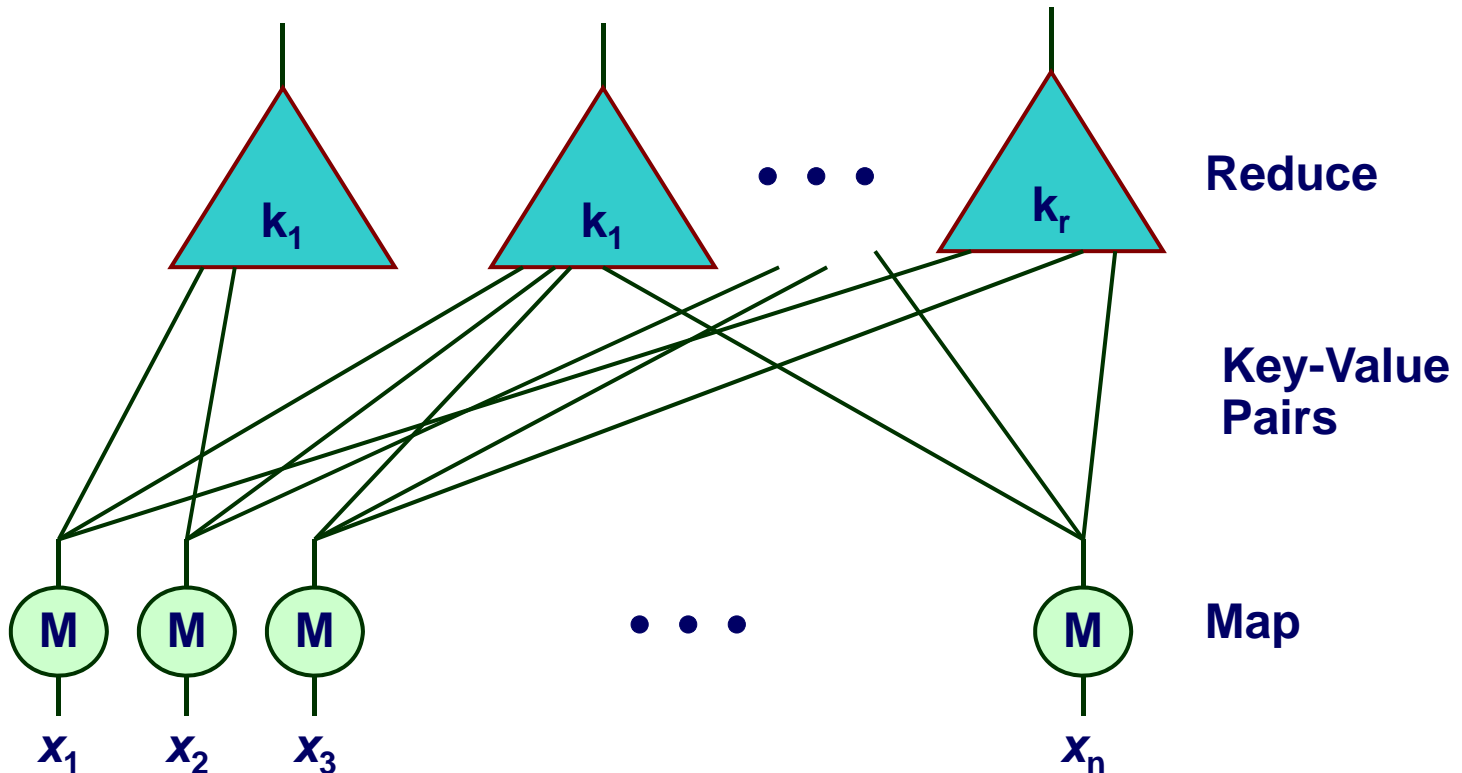- **High bandwidth communication**

## Strengths

- **High utilization of resources**
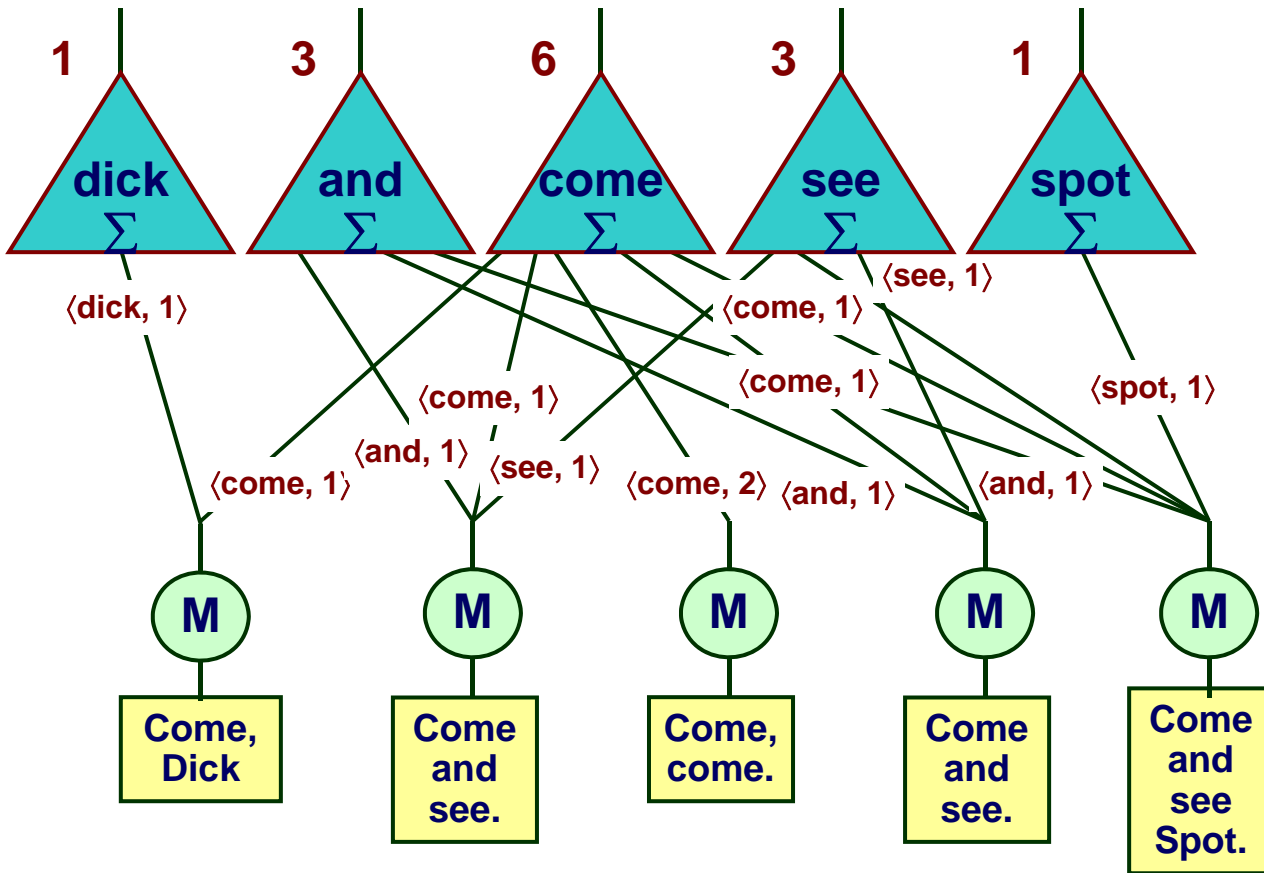- **Effective for many scientific applications**

## Weaknesses

- **Requires careful tuning of application to resources**
- **Intolerant of any variability**

# Map/Reduce Programming Model



Reduce

Key-Value Pairs

Map

$x_1$ $x_2$ $x_3$ $x_n$

- **Map computation across many objects**
  - E.g., $10^{10}$ Internet web pages

- **Aggregate results in many different ways**

- **System deals with issues of resource allocation & reliability**

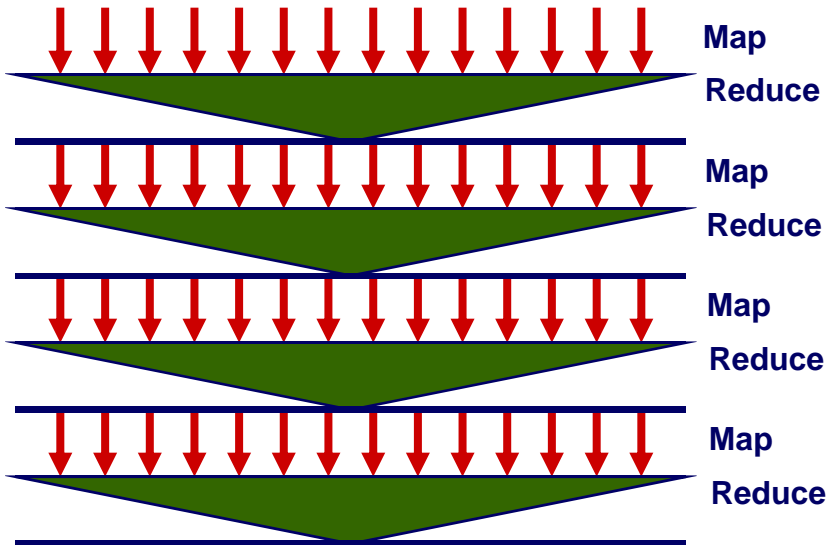Dean & Ghemawat: "MapReduce: Simplified Data Processing on Large Clusters", OSDI 2004

# Map/Reduce Example



**1** dick $\Sigma$    **3** and $\Sigma$    **6** come $\Sigma$    **3** see $\Sigma$    **1** spot $\Sigma$

⟨dick, 1⟩
⟨see, 1⟩
⟨come, 1⟩
⟨come, 1⟩
⟨come, 1⟩
⟨spot, 1⟩
⟨come, 1⟩
⟨and, 1⟩
⟨see, 1⟩
⟨come, 2⟩
⟨and, 1⟩
⟨and, 1⟩

**Sum**

**Word-Count Pairs**

Come and see.

Come and see Spot.

**Extract**

M — Come, Dick

M — Come and see.

M — Come, come.

M — Come and see.

M — Come and see Spot.

- **Create an word index of set of documents**
- **Map: generate ⟨word, count⟩ pairs for all words in document**
- **Reduce: sum word counts across documents**

# Map/Reduce Operation

## Map/Reduce



## Characteristics

- **Computation broken into many, short-lived tasks**
  - Mapping, reducing
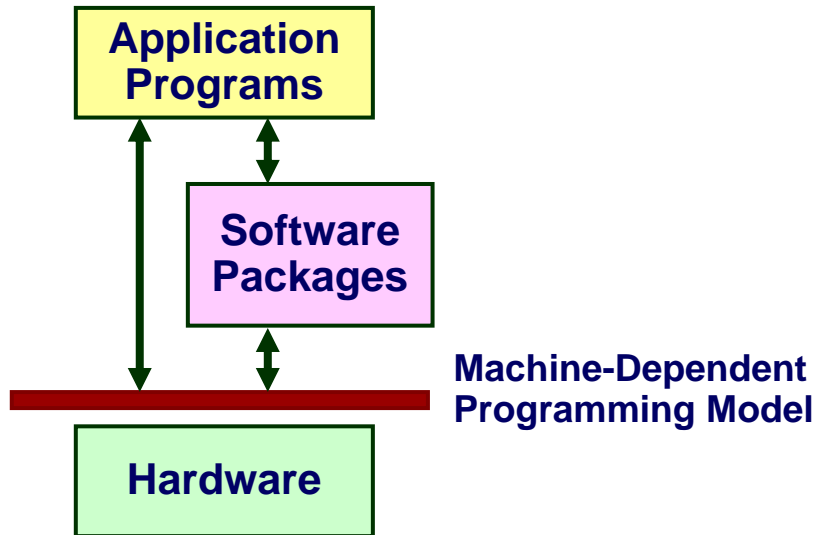- **Use disk storage to hold intermediate results**

## Strengths

- **Great flexibility in placement, scheduling, and load balancing**
- **Can access large data sets**
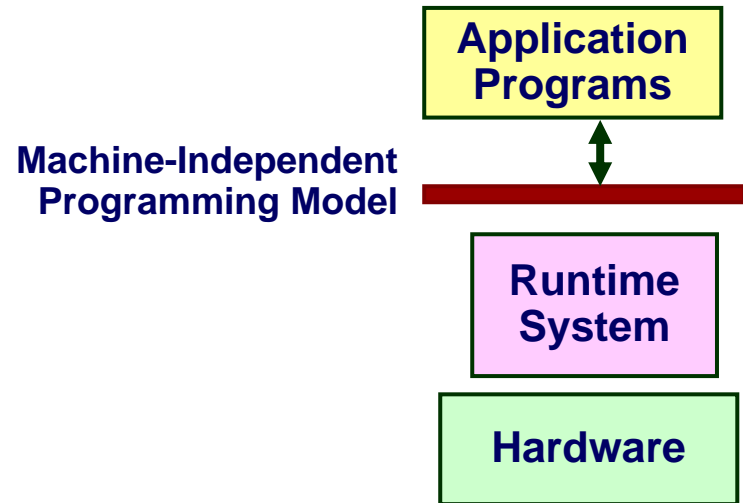
## Weaknesses

- **Higher overhead**
- **Lower raw performance**

# System Comparison: Programming Models

## Conventional HPC

```
┌─────────────────┐
│   Application   │
│    Programs     │
└─────────────────┘
        ↕
┌─────────────────┐
│    Software     │
│    Packages     │
└─────────────────┘
        ↕
━━━━━━━━━━━━━━━━━━━━   Machine-Dependent
                      Programming Model
┌─────────────────┐
│    Hardware     │
└─────────────────┘
```
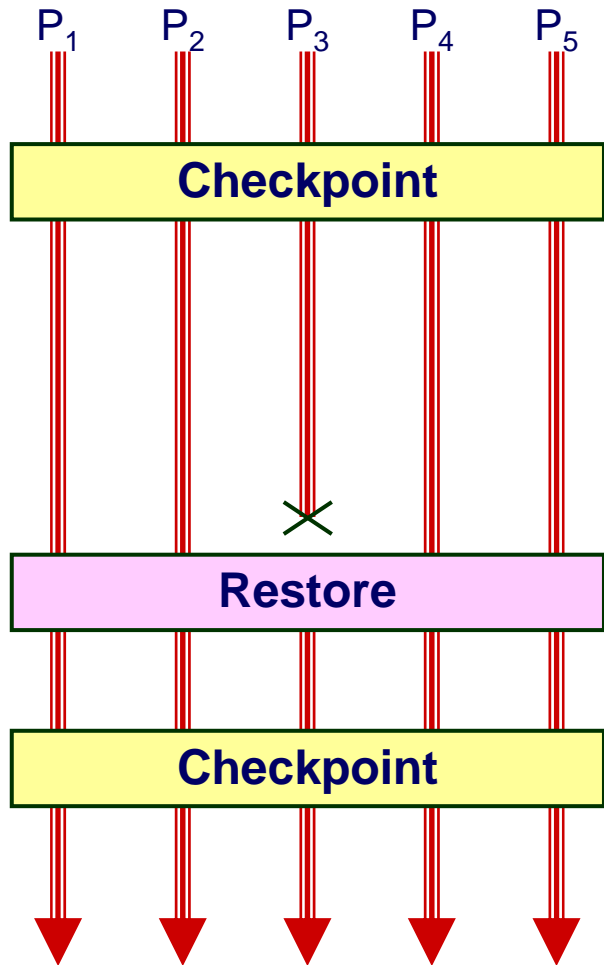
- **Programs described at very low level**
  - Specify detailed control of processing & communications
- **Rely on small number of software packages**
  - Written by specialists
  - Limits classes of problems & solution methods

## DISC

```
┌─────────────────┐
│   Application   │
│    Programs     │
└─────────────────┘
        ↕
━━━━━━━━━━━━━━━━━━━━
Machine-Independent
Programming Model
┌─────────────────┐
│     Runtime     │
│     System      │
└─────────────────┘

┌─────────────────┐
│    Hardware     │
└─────────────────┘
```

- **Application programs written in terms of high-level operations on data**
- **Runtime system controls scheduling, load balancing, …**

# HPC Fault Tolerance

P$_1$   P$_2$   P$_3$   P$_4$   P$_5$

| Checkpoint |

Wasted
Computation

| Restore |

| Checkpoint |

## Checkpoint

- **Periodically store state of all processes**
- **Significant I/O traffic**

## Restore

- **When failure occurs**
- **Reset state to that of last checkpoint**
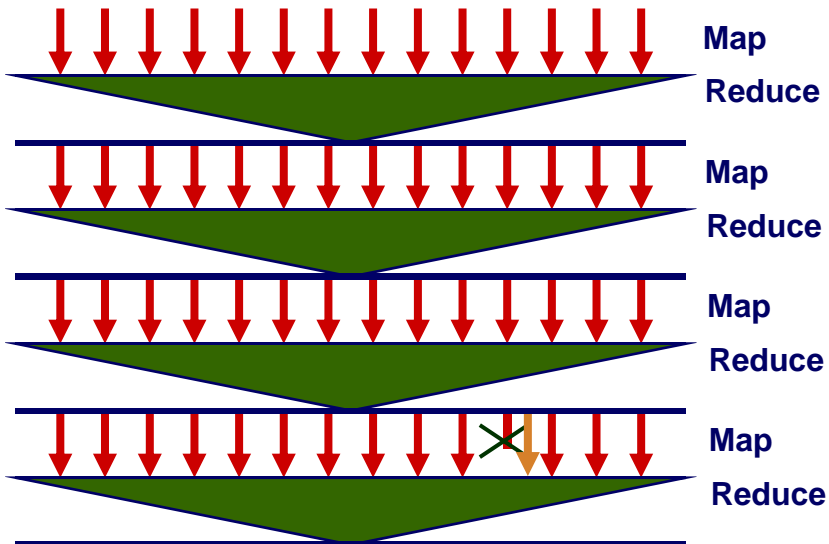- **All intervening computation wasted**

## Performance Scaling

- **Very sensitive to number of failing components**

# Map/Reduce Fault Tolerance

## Map/Reduce



Map
Reduce
Map
Reduce
Map
Reduce
Map
Reduce

## Data Integrity

- **Store multiple copies of each file**
- **Including intermediate results of each Map / Reduce**
  - Continuous checkpointing

## Recovering from Failure

- **Simply recompute lost result**
  - Localized effect
- **Dynamic scheduler keeps all processors busy**

– 17 –

# DISC Scalability Advantages

- **Distributed system design principles lead to scalable design**
- **Dynamically scheduled tasks with state held in replicated files**

## Provisioning Advantages

- **Can use consumer-grade components**
  - maximizes cost-peformance
- **Can have heterogenous nodes**
  - More efficient technology refresh

## Operational Advantages

- **Minimal staffing**
- **No downtime**

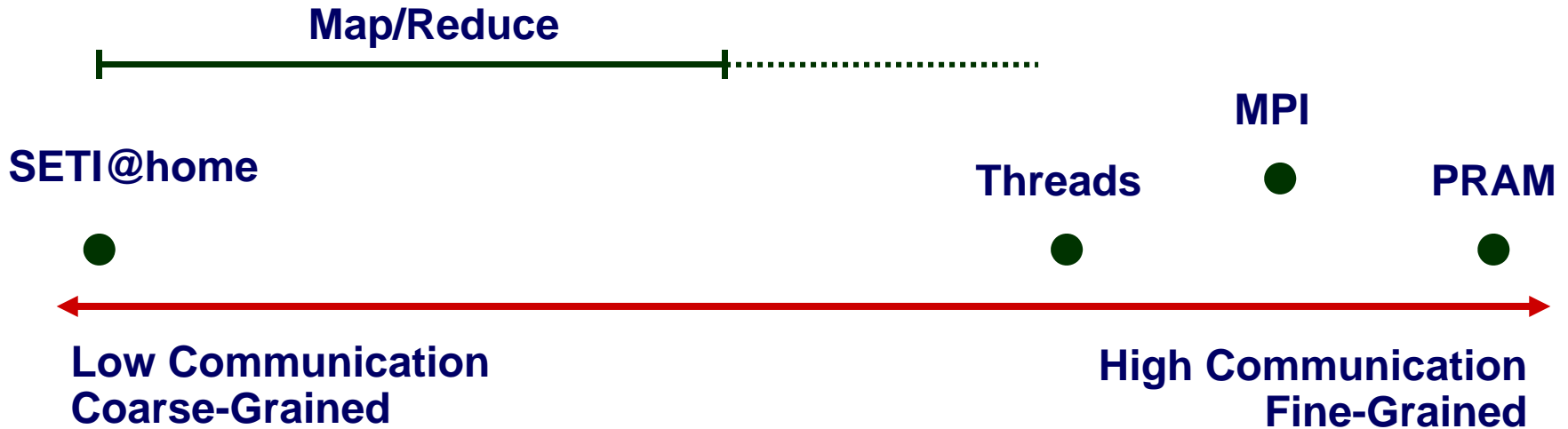# Getting Started

## Goal

- **Get people involved in DISC**

## Software

- **Hadoop Project**
  - Open source project providing file system and Map/Reduce
  - Supported and used by Yahoo
  - Rapidly expanding user/developer base
  - Prototype on single machine, map onto cluster

# Exploring Parallel Computation Models

**Map/Reduce**

**SETI@home**

**MPI**

**Threads**

**PRAM**

**Low Communication
Coarse-Grained**

**High Communication
Fine-Grained**

## DISC + Map/Reduce Provides Coarse-Grained Parallelism

- **Computation done by independent processes**
- **File-based communication**

## Observations

- **Relatively "natural" programming model**
- **Research issue to explore full potential and limits**

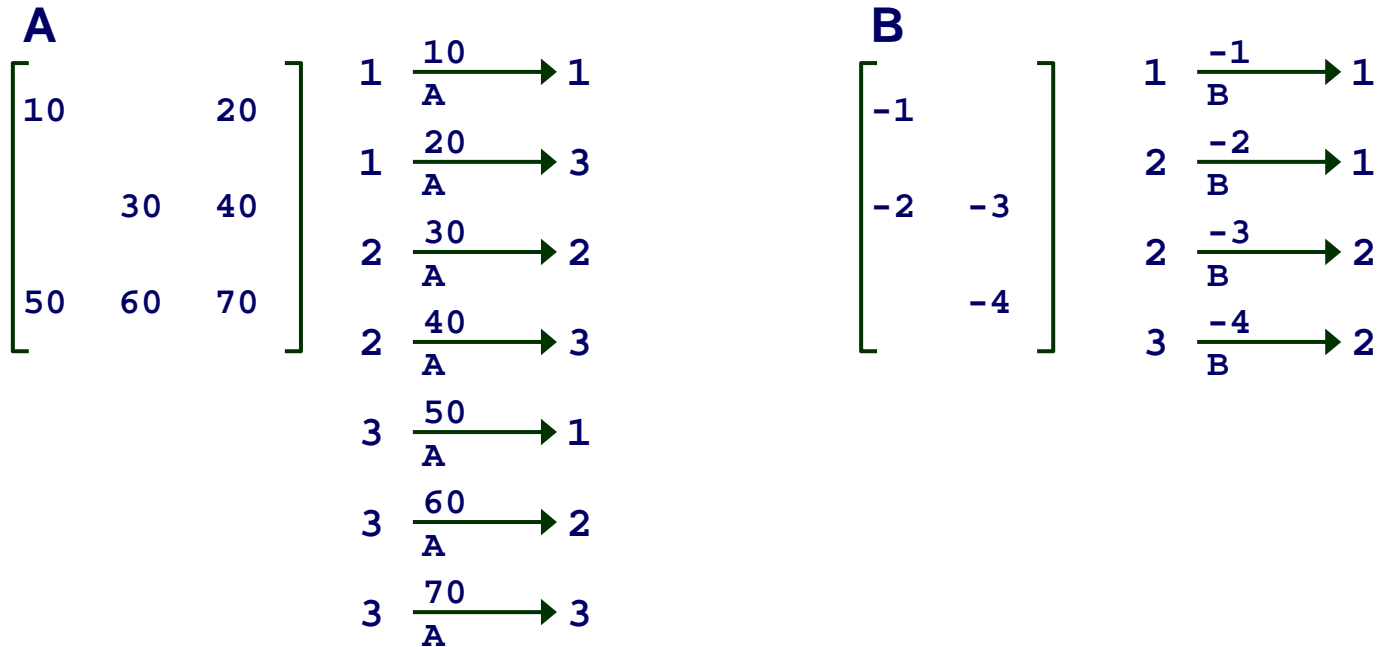# Example: Sparse Matrices with Map/Reduce

$$
\mathbf{A} \quad\quad\quad\quad \mathbf{B} \quad\quad\quad\quad \mathbf{C}
$$

$$
\begin{bmatrix} 10 & & 20 \\ & 30 & 40 \\ 50 & 60 & 70 \end{bmatrix}
\;\mathbf{X}\;
\begin{bmatrix} -1 & \\ -2 & -3 \\ & -4 \end{bmatrix}
\;=\;
\begin{bmatrix} -10 & -80 \\ -60 & -250 \\ -170 & -460 \end{bmatrix}
$$

- **Task: Compute product C = A·B**
- **Assume most matrix entries are 0**
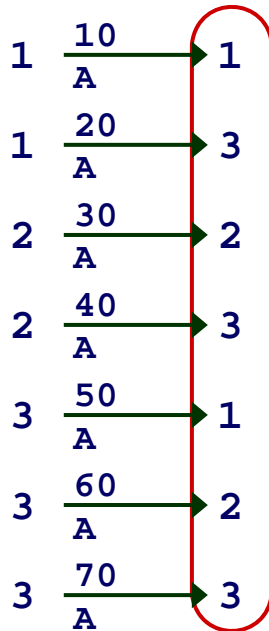
## Motivation

- **Core problem in scientific computing**
- **Challenging for parallel execution**
- **Demonstrate expressiveness of Map/Reduce**
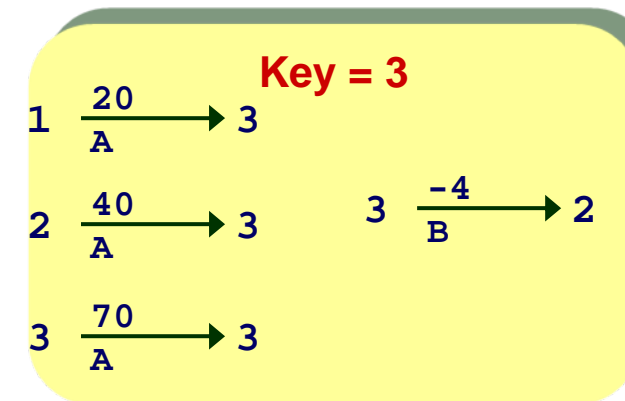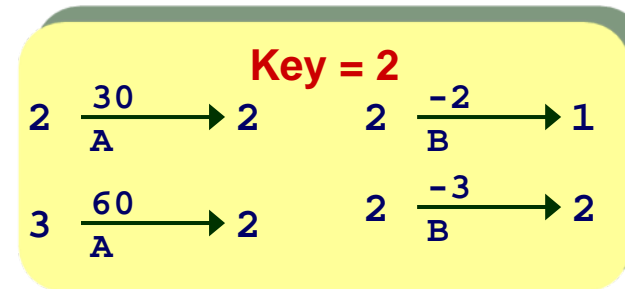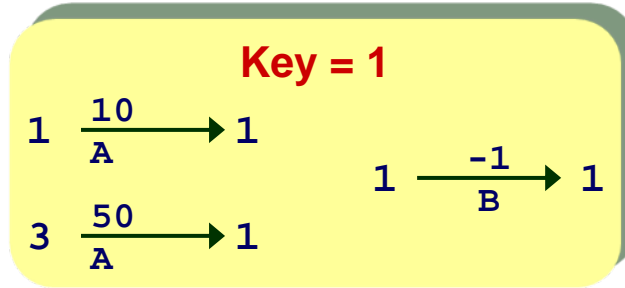
# Computing Sparse Matrix Product

**A**

$$\begin{bmatrix} 10 & & 20 \\ & 30 & 40 \\ 50 & 60 & 70 \end{bmatrix}$$

1 $\xrightarrow[A]{10}$ 1

1 $\xrightarrow[A]{20}$ 3

2 $\xrightarrow[A]{30}$ 2

2 $\xrightarrow[A]{40}$ 3

3 $\xrightarrow[A]{50}$ 1

3 $\xrightarrow[A]{60}$ 2

3 $\xrightarrow[A]{70}$ 3

**B**

$$\begin{bmatrix} -1 & & \\ -2 & -3 & \\ & & -4 \end{bmatrix}$$

1 $\xrightarrow[B]{-1}$ 1

2 $\xrightarrow[B]{-2}$ 1

2 $\xrightarrow[B]{-3}$ 2

3 $\xrightarrow[B]{-4}$ 2

- **Represent matrix as list of nonzero entries**
  - ⟨**row, col, value, matrixID**⟩

- **Strategy**
  - **Phase 1: Compute all products $a_{i,k} \cdot b_{k,j}$**
  - **Phase 2: Sum products for each entry i,j**
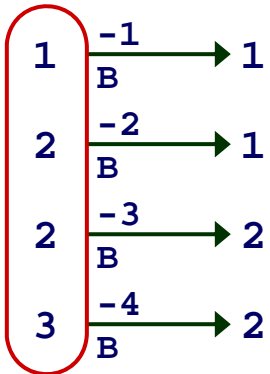  - **Each phase involves a Map/Reduce**

# Phase 1 Map of Matrix Multiply

1 $\xrightarrow[\text{A}]{10}$ 1

1 $\xrightarrow[\text{A}]{20}$ 3

2 $\xrightarrow[\text{A}]{30}$ 2

2 $\xrightarrow[\text{A}]{40}$ 3

3 $\xrightarrow[\text{A}]{50}$ 1

3 $\xrightarrow[\text{A}]{60}$ 2

3 $\xrightarrow[\text{A}]{70}$ 3

**Key = col**

**Key = row**

1 $\xrightarrow[\text{B}]{-1}$ 1

2 $\xrightarrow[\text{B}]{-2}$ 1

2 $\xrightarrow[\text{B}]{-3}$ 2

3 $\xrightarrow[\text{B}]{-4}$ 2

## Key = 1

1 $\xrightarrow[\text{A}]{10}$ 1

3 $\xrightarrow[\text{A}]{50}$ 1

1 $\xrightarrow[\text{B}]{-1}$ 1

## Key = 2

2 $\xrightarrow[\text{A}]{30}$ 2

3 $\xrightarrow[\text{A}]{60}$ 2

2 $\xrightarrow[\text{B}]{-2}$ 1

2 $\xrightarrow[\text{B}]{-3}$ 2

## Key = 3

1 $\xrightarrow[\text{A}]{20}$ 3

2 $\xrightarrow[\text{A}]{40}$ 3

3 $\xrightarrow[\text{A}]{70}$ 3

3 $\xrightarrow[\text{B}]{-4}$ 2

- **Group values $a_{i,k}$ and $b_{k,j}$ according to key k**

# Phase 1 "Reduce" of Matrix Multiply

**Key = 1**

$1 \xrightarrow[A]{10} 1$

$3 \xrightarrow[A]{50} 1$

**X**  $1 \xrightarrow[B]{-1} 1$

$1 \xrightarrow[C]{-10} 1$

$3 \xrightarrow[A]{-50} 1$

**Key = 2**

$2 \xrightarrow[A]{30} 2$

$3 \xrightarrow[A]{60} 2$

**X**  $2 \xrightarrow[B]{-2} 1$

$2 \xrightarrow[B]{-3} 2$

$2 \xrightarrow[C]{-60} 1$

$2 \xrightarrow[C]{-90} 2$

$3 \xrightarrow[C]{-120} 1$

$3 \xrightarrow[C]{-180} 2$

**Key = 3**

$1 \xrightarrow[A]{20} 3$

$2 \xrightarrow[A]{40} 3$

**X**  $3 \xrightarrow[B]{-4} 2$

$3 \xrightarrow[A]{70} 3$

$1 \xrightarrow[C]{-80} 2$
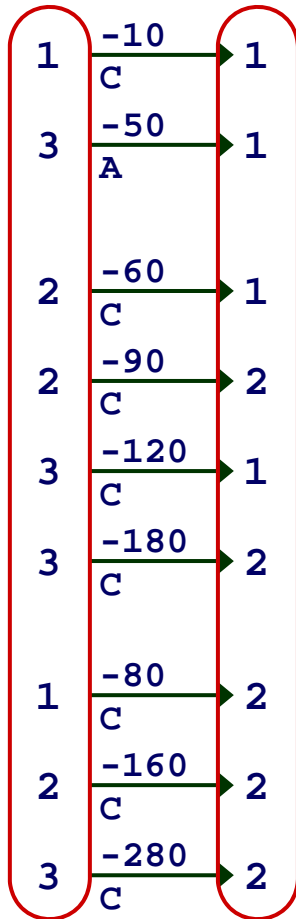
$2 \xrightarrow[C]{-160} 2$

$3 \xrightarrow[C]{-280} 2$

- **Generate all products $a_{i,k} \cdot b_{k,j}$**

# Phase 2 Map of Matrix Multiply

1 $\xrightarrow[\text{C}]{-10}$ 1

3 $\xrightarrow[\text{A}]{-50}$ 1

2 $\xrightarrow[\text{C}]{-60}$ 1

2 $\xrightarrow[\text{C}]{-90}$ 2

3 $\xrightarrow[\text{C}]{-120}$ 1

3 $\xrightarrow[\text{C}]{-180}$ 2

1 $\xrightarrow[\text{C}]{-80}$ 2

2 $\xrightarrow[\text{C}]{-160}$ 2

3 $\xrightarrow[\text{C}]{-280}$ 2

**Key = row,col**

**Key = 1,1**   1 $\xrightarrow[\text{C}]{-10}$ 1

**Key = 1,2**   1 $\xrightarrow[\text{C}]{-80}$ 2

**Key = 2,1**   2 $\xrightarrow[\text{C}]{-60}$ 1

**Key = 2,2**   2 $\xrightarrow[\text{C}]{-90}$ 2

2 $\xrightarrow[\text{C}]{-160}$ 2

**Key = 3,1**   3 $\xrightarrow[\text{C}]{-120}$ 1

3 $\xrightarrow[\text{A}]{-50}$ 1

**Key = 3,2**   3 $\xrightarrow[\text{C}]{-280}$ 2

3 $\xrightarrow[\text{C}]{-180}$ 2

- **Group products $a_{i,k} \cdot b_{k,j}$ with matching values of i and j**

# Phase 2 Reduce of Matrix Multiply

**Key = 1,1**    $1 \xrightarrow[C]{-10} 1$       $1 \xrightarrow[C]{-10} 1$

**Key = 1,2**    $1 \xrightarrow[C]{-80} 2$       $1 \xrightarrow[C]{-80} 2$

**Key = 2,1**    $2 \xrightarrow[C]{-60} 1$       $2 \xrightarrow[C]{-60} 1$

**Key = 2,2**    $2 \xrightarrow[C]{-90} 2$

              $2 \xrightarrow[C]{-160} 2$       $2 \xrightarrow[C]{-250} 2$

**Key = 3,1**    $3 \xrightarrow[C]{-120} 1$

              $3 \xrightarrow[A]{-50} 1$       $3 \xrightarrow[C]{-170} 1$

**Key = 3,2**    $3 \xrightarrow[C]{-280} 2$

              $3 \xrightarrow[C]{-180} 2$       $3 \xrightarrow[C]{-460} 2$

**C**

$$C = \begin{bmatrix} -10 & -80 \\ -60 & -250 \\ -170 & -460 \end{bmatrix}$$

- **Sum products to get final entries**

# Lessons from Sparse Matrix Example

**Associative Matching is Powerful Communication Primitive**

- **Intermediate step in Map/Reduce**

**Similar Strategy Applies to Other Problems**

- **Shortest path in graph**
- **Database join**

**Many Performance Considerations**

- **Kiefer, Volk, Lehner, TU Dresden**
- **Should do systematic comparison to other sparse matrix implementations**
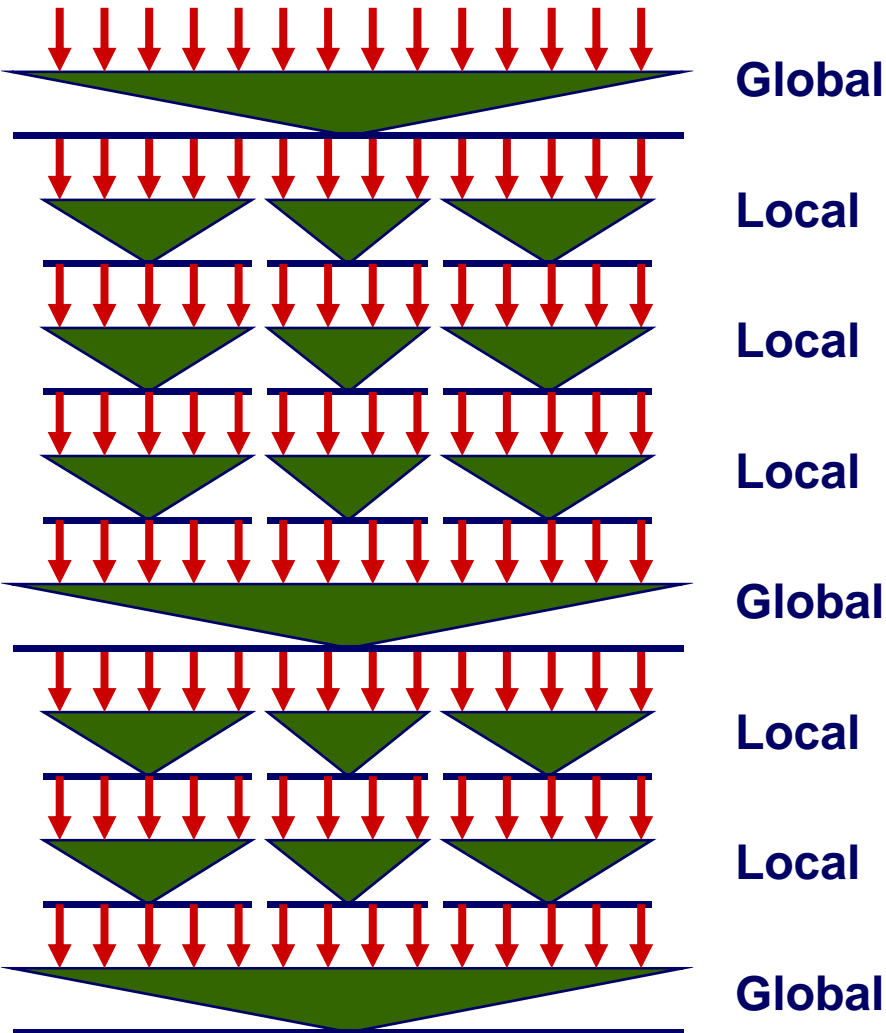
# Tweaking MapReduce

## Map/Reduce/Merge

- **Yang, et al.  Yahoo & UCLA**
- **Add merge step to do tagged merge of data sets**
  - E.g., getting elements from matrices A & B
- **Demonstrate expressive power of relational algebra**

## Local Iterations

- **Kambatla, et al., Purdue**
- **Support local and global iterations**
- **Step toward bulk synchronous model**

# Local / Global Map/Reduce

Global

Local

Local

Local

Global

Local

Local

Global

## Iterative Computations
- **e.g., PageRank**

## Graph Partitioning
- **Partition into clusters**
- **Colocate data for each cluster**

## Computation
- **Compute solution for each cluster, holding inter-cluster values constant**
- **Update inter-cluster values**

# Pig Project

- **Chris Olston, Yahoo!**
- **Part of Apache/Hadoop**

## Merge Database & Programming

- **SQL-like query language**
- **Set-oriented function application**

## Implementation

- **Map onto Hadoop**
- **Automatic selection / optimization of algorithm**
- **Captures low-level tricks programmers have devised for Map/Reduce**
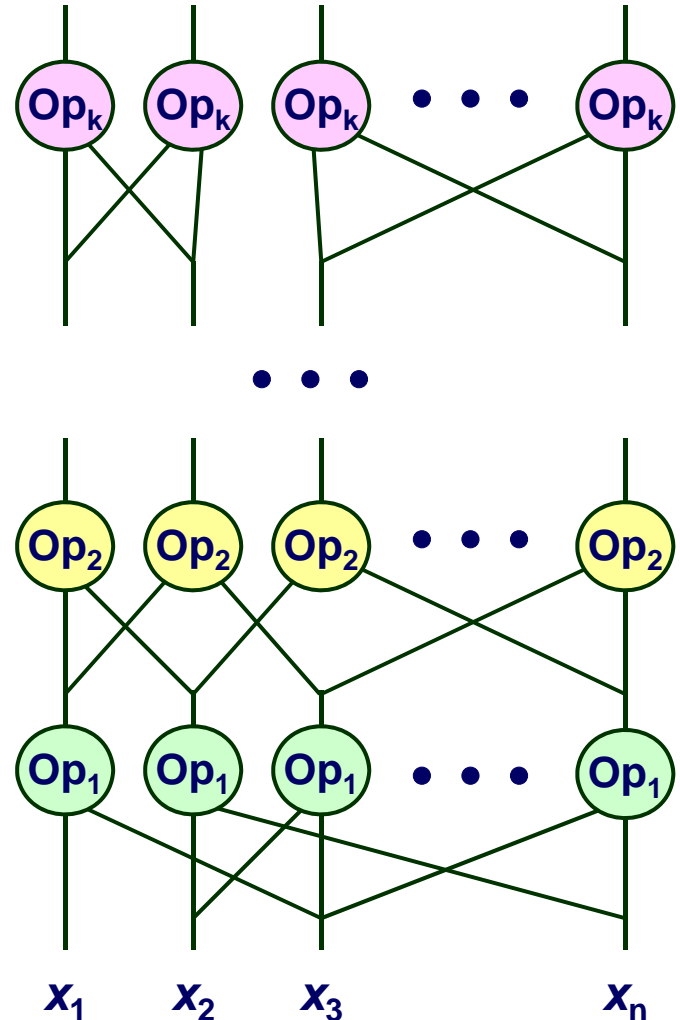
# Generalizing Map/Reduce

- **E.g., Microsoft Dryad Project**

## Computational Model

- **Acyclic graph of operators**
  - But expressed as textual program
- **Each takes collection of objects and produces objects**
  - Purely functional model

## Implementation Concepts

- **Objects stored in files or memory**
- **Any object may be lost; any operator may fail**
- **Replicate & recompute for fault tolerance**
- **Dynamic scheduling**
  - # Operators >> # Processors

# Implementation Challenges

## Hadoop Platform is Blessing

- Large community adding features, improving performance
- Easy to deploy
- Growing body of documentation and materials
- Works well enough for range of applications

## … And Curse

- Map & reduce functionality hardwired
- Not designed as extensible platform

## Dryad Not Widely Adopted

- 305 citations on Google Scholar vs. 1759 Map/Reduce
- Built on .NET

# Conclusions

## Distributed Systems Concepts Lead to Scalable Machines

- Loosely coupled execution model
- Lowers cost of procurement & operation

## Map/Reduce Gaining Widespread Use

- Hadoop makes it widely available
- Great for some applications, good enough for many others

## Lots of Work to be Done

- Richer set of programming models and implementations
- Expanding range of applicability
  - Problems that are data *and* compute intensive
  - The future of supercomputing?