

Introduction to Database Systems

CSE 444

Lecture 21: Parallel and Distributed Databases

Where We Are

- How to use a DBMS as a:
 - Data analyst: SQL, SQL, SQL,...
 - Application programmer: JDBC, XML,...
 - Database administrator: tuning, triggers, security
 - Massive-scale data analyst:
 - Parallel DBMSs and Pig/MapReduce
- How DBMSs work
 - Transactions
 - Data storage and indexing
 - Query execution
 - Databases as a service

Outline

- **Parallel data management motivation**
 - All database vendors offer parallel DBMSs
 - Full support for entire SQL and ACID transactions
- Basic concepts behind parallel DBMSs
- MapReduce data processing approach
- Readings:
 - Jeffrey Dean and Sanjay Ghemawat. [MapReduce: Simplified Data Processing on Large Clusters](#). OSDI 2004. Sections 1 through 3 only.

Parallel v.s. Distributed Databases

- Parallel database system
 - Improve performance through parallel implementation
- Distributed database system
 - Data is stored across several sites, each site managed by a DBMS capable of running independently

Parallel DBMSs

- Goal
 - Improve performance by executing multiple operations in parallel
 - Want to scale transactions per second
 - Want to scale large analytics queries
- Key benefit
 - Cheaper to scale than relying on a single increasingly more powerful processor
- Key challenge
 - Ensure overhead and contention do not kill performance

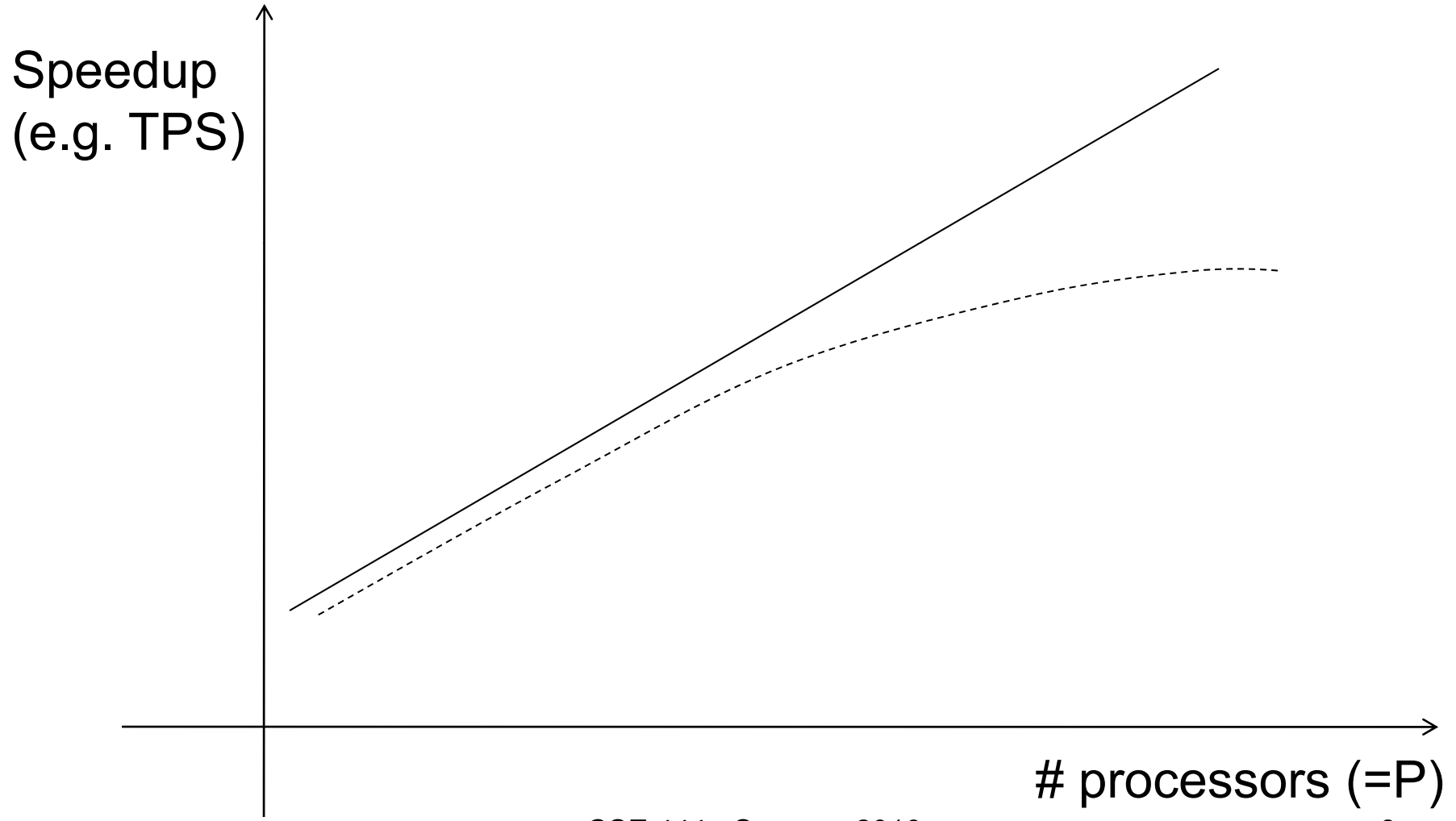
Outline

- Parallel data management motivation
 - All database vendors offer parallel DBMSs
 - Full support for entire SQL and ACID transactions
- **Basic concepts behind parallel DBMSs**
- MapReduce data processing approach
- Readings:
 - Jeffrey Dean and Sanjay Ghemawat. [MapReduce: Simplified Data Processing on Large Clusters](#). OSDI 2004. Sections 1 through 3 only.

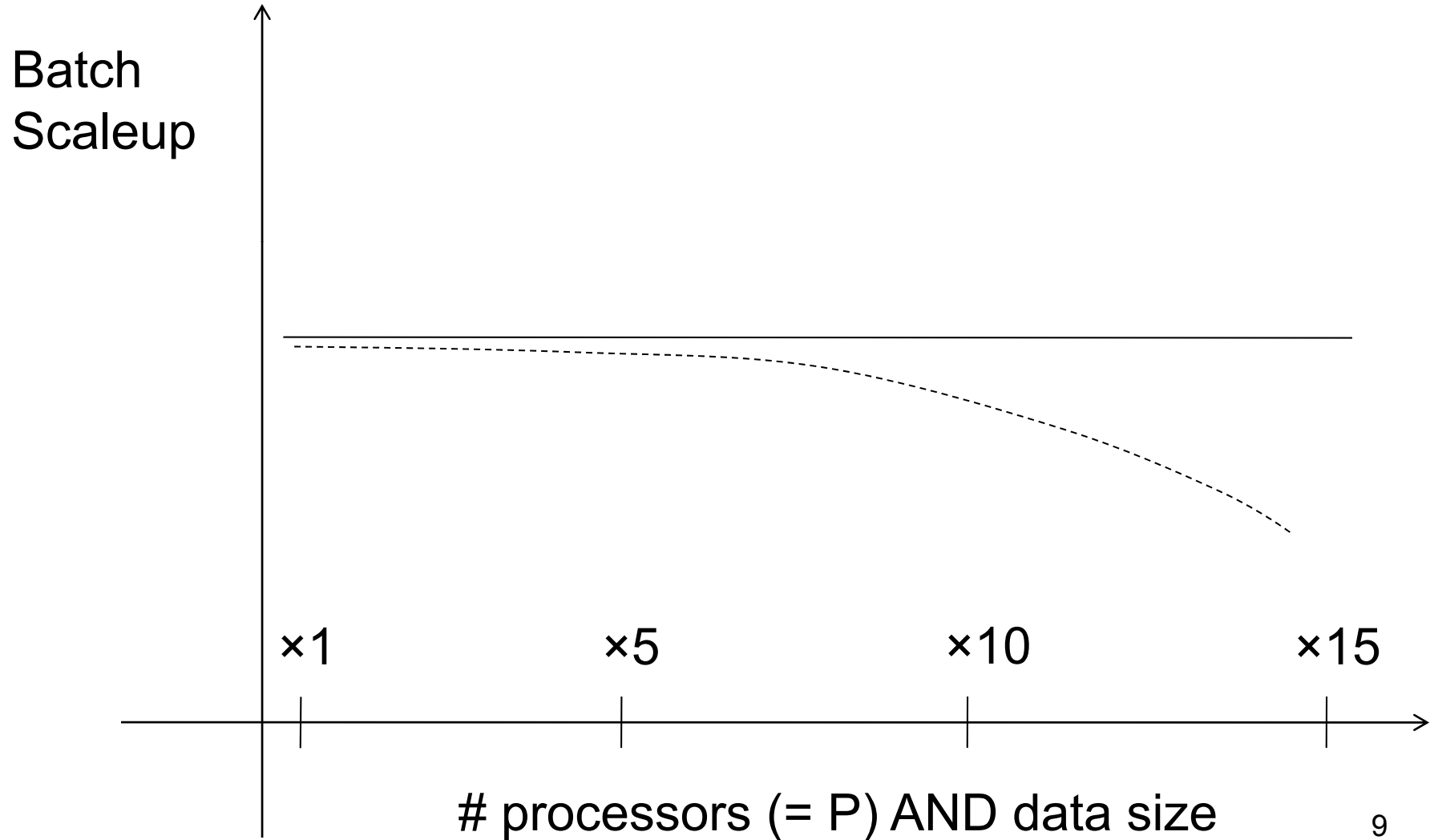
Performance Metrics for Parallel DBMSs

- **Speedup**
 - More processors → higher speed
 - Individual queries should run faster
 - Should have more Transactions per Second (TPS)
- **Scaleup**
 - More processors → can process more data
 - Batch scaleup: same time for same query on larger data
 - Transaction scaleup: N-times TPS on N-times larger database

Linear v.s. Non-linear Speedup



Linear v.s. Non-linear Scaleup



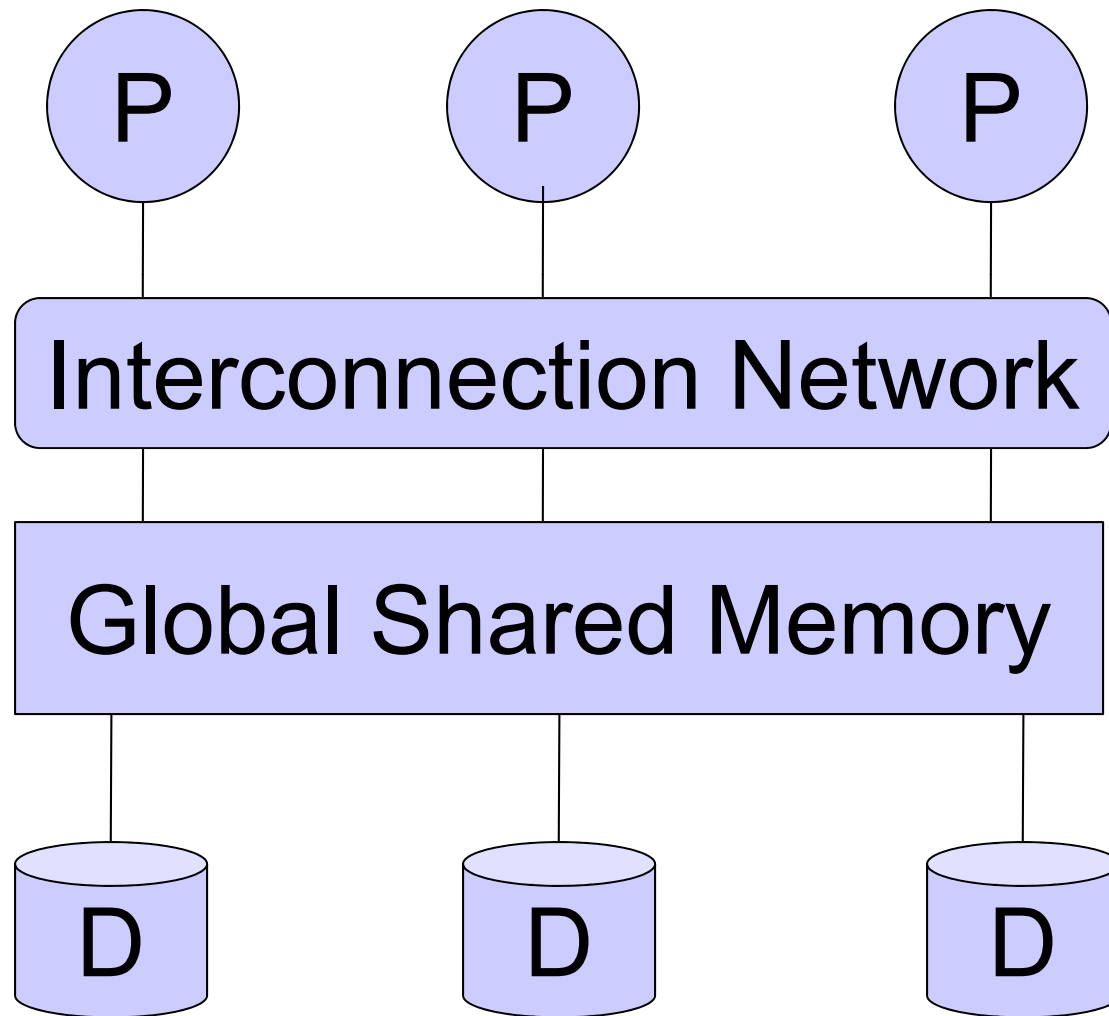
Challenges to Speedup & Scaleup

- **Startup cost**
 - Cost of starting an operation on many processors
- **Interference**
 - Contention for resources between processors
- **Skew**
 - Slowest step becomes the bottleneck

Architectures for Parallel Databases

- Shared memory
- Shared disk
- Shared nothing

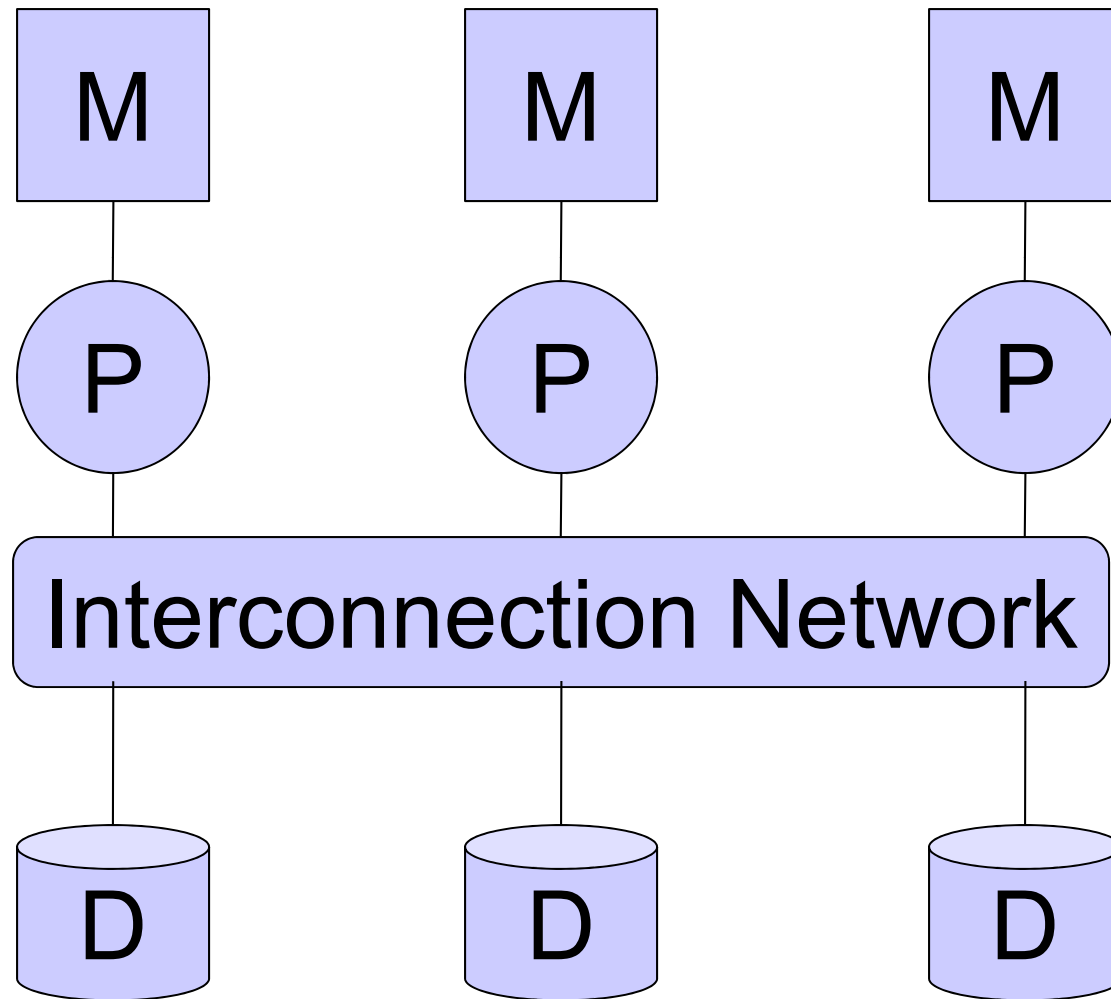
Shared Memory



Shared Memory

- Single address space
- Typically each processor has own memory, but the the combined memory forms unique address space
 - NUMA (nonuniform memory access):
 - Local memory faster than remote memory
 - *Any* memory access is faster than disk

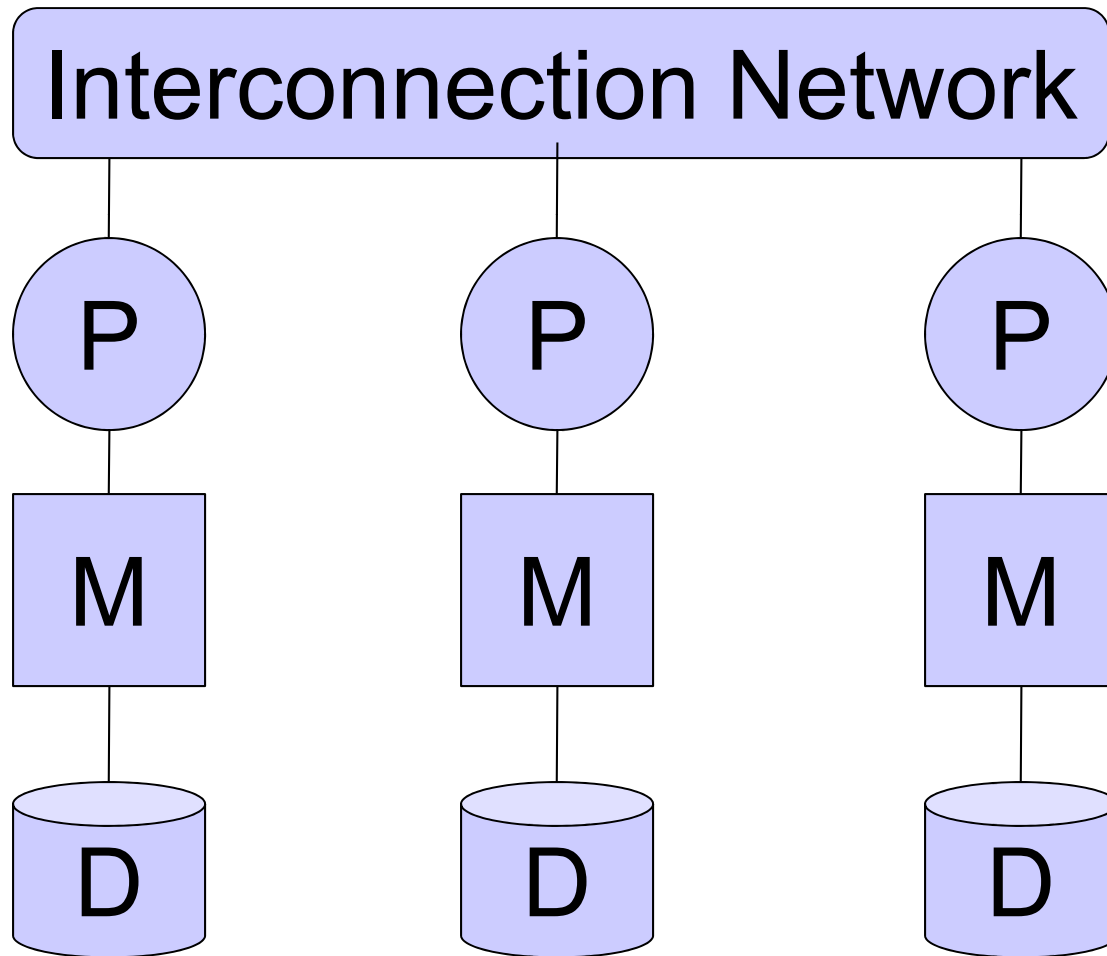
Shared Disk



Shared Disk

- Separate disks (not attached to processors)
- Disk controllers manage competing requests from processors

Shared Nothing



Shared Nothing

- Most scalable architecture
 - Minimizes interference by minimizing resource sharing
 - Can use commodity hardware
- Also most difficult to program and manage
- Processor = server = node
- P = number of nodes

We will focus on shared nothing

Taxonomy for Parallel Query Evaluation

- Inter-query parallelism
 - Each query runs on one processor
- Inter-operator parallelism
 - A query runs on multiple processors
 - An operator runs on one processor
- Intra-operator parallelism
 - An operator runs on multiple processors

We study only intra-operator parallelism: most scalable

Horizontal Data Partitioning

- Relation R split into P chunks R_0, \dots, R_{P-1} , stored at the P nodes
- **Round robin**: tuple t_i to chunk $(i \bmod P)$
- **Hash based partitioning on attribute A** :
 - Tuple t to chunk $h(t.A) \bmod P$
- **Range based partitioning on attribute A** :
 - Tuple t to chunk i if $v_{i-1} < t.A < v_i$

Parallel Selection

Compute $\sigma_{A=v}(R)$, or $\sigma_{v_1 < A < v_2}(R)$

- On a conventional database: cost = $B(R)$
- Q: What is the cost on a parallel database with P processors ?
 - Round robin
 - Hash partitioned
 - Range partitioned

Parallel Selection

- Q: What is the cost on a parallel database with P processors ?
- A: $B(R) / P$ in all cases
- However, different processors do the work:
 - Round robin: all servers do the work
 - Hash: one server for $\sigma_{A=v}(R)$, all for $\sigma_{v_1 < A < v_2}(R)$
 - Range: one server only

Data Partitioning Revisited

What are the pros and cons ?

- Round robin
 - Good load balance but always needs to read all the data
- Hash based partitioning
 - Good load balance but works only for equality predicates and full scans
- Range based partitioning
 - Works well for range predicates but can suffer from data skew

Parallel Group By

- Compute $\gamma_{A, \text{sum}(B)}(R)$
- Step 1: server i partitions chunk R_i using a hash function $h(t.A) \bmod P$: $R_{i0}, R_{i1}, \dots, R_{i,P-1}$
- Step 2: server i sends partition R_{ij} to server j
- Step 3: server j computes $\gamma_{A, \text{sum}(B)}$ on $R_{0j}, R_{1j}, \dots, R_{P-1,j}$

Cost of Parallel Group By

Recall conventional cost = $3B(R)$

- Cost of Step 1: $B(R)/P$ I/O operations
- Cost of Step 2: $(P-1)/P B(R)$ blocks are sent
 - Network costs assumed to be much lower than I/O
- Cost of Step 3: $2 B(R)/P$
 - Why ?
 - When can we reduce it to 0 ?

Total = $3B(R) / P$ + communication costs

Parallel Join

- Step 1
 - For all servers in $[0,k]$, server i partitions chunk R_i using a hash function $h(t.A) \bmod P$: $R_{i0}, R_{i1}, \dots, R_{i,P-1}$
 - For all servers in $[k+1,P]$, server j partitions chunk S_j using a hash function $h(t.A) \bmod P$: $S_{j0}, S_{j1}, \dots, S_{j,P-1}$
- Step 2:
 - Server i sends partition R_{iu} to server u
 - Server j sends partition S_{ju} to server u
- Step 3: Server u computes the join of R_{iu} with S_{ju}

Cost of Parallel Join

- Step 1: $(B(R) + B(S))/P$
- Step 2: 0
 - $(P-1)/P (B(R) + B(S))$ blocks are sent, but we assume network costs to be \ll disk I/O costs
- Step 3:
 - 0 if smaller table fits in main memory: $B(S)/p \leq M$
 - $4(B(R)+B(S))/P$ otherwise

Parallel Dataflow Implementation

- Use relational operators unchanged
- Add special split and merge operators
 - Handle data routing, buffering, and flow control
- Example: exchange operator
 - Inserted between consecutive operators in the query plan
 - Can act as either a producer or consumer
 - Producer pulls data from operator and sends to n consumers
 - Producer acts as driver for operators below it in query plan
 - Consumer buffers input data from n producers and makes it available to operator through getNext interface

Outline

- Parallel data management motivation
 - All database vendors offer parallel DBMSs
 - Full support for entire SQL and ACID transactions
- Basic concepts behind parallel DBMSs
- MapReduce data processing approach
- Readings:
 - Jeffrey Dean and Sanjay Ghemawat. [MapReduce: Simplified Data Processing on Large Clusters](#). OSDI 2004. Sections 1 through 3 only.

Map Reduce

- Google: paper published 2004
- Open source variant: Hadoop
- Map-reduce = high-level programming model and implementation for large-scale parallel data processing

MapReduce (MR) tools

MR implementation:



One MR query language:

Pig Latin

Query engine:



MR Data Model

- Files !
- A file = a bag of (key, value) pairs
- A map-reduce program:
 - Input: a bag of (input key, value) pairs
 - Output: a bag of (output key, value) pairs

Step 1: the MAP Phase

- User provides the MAP-function:
 - Input: one (input key, value)
 - Output: a bag of (intermediate key, value) pairs
- System applies map function in parallel to all (input key, value) pairs in the input file

Step 2: the REDUCE Phase

- User provides the REDUCE function:
 - Input: intermediate key, and bag of values
 - Output: bag of output values
- System groups all pairs with the same intermediate key, and passes the bag of values to the REDUCE function

Example

- Counting the number of occurrences of each word in a large collection of documents

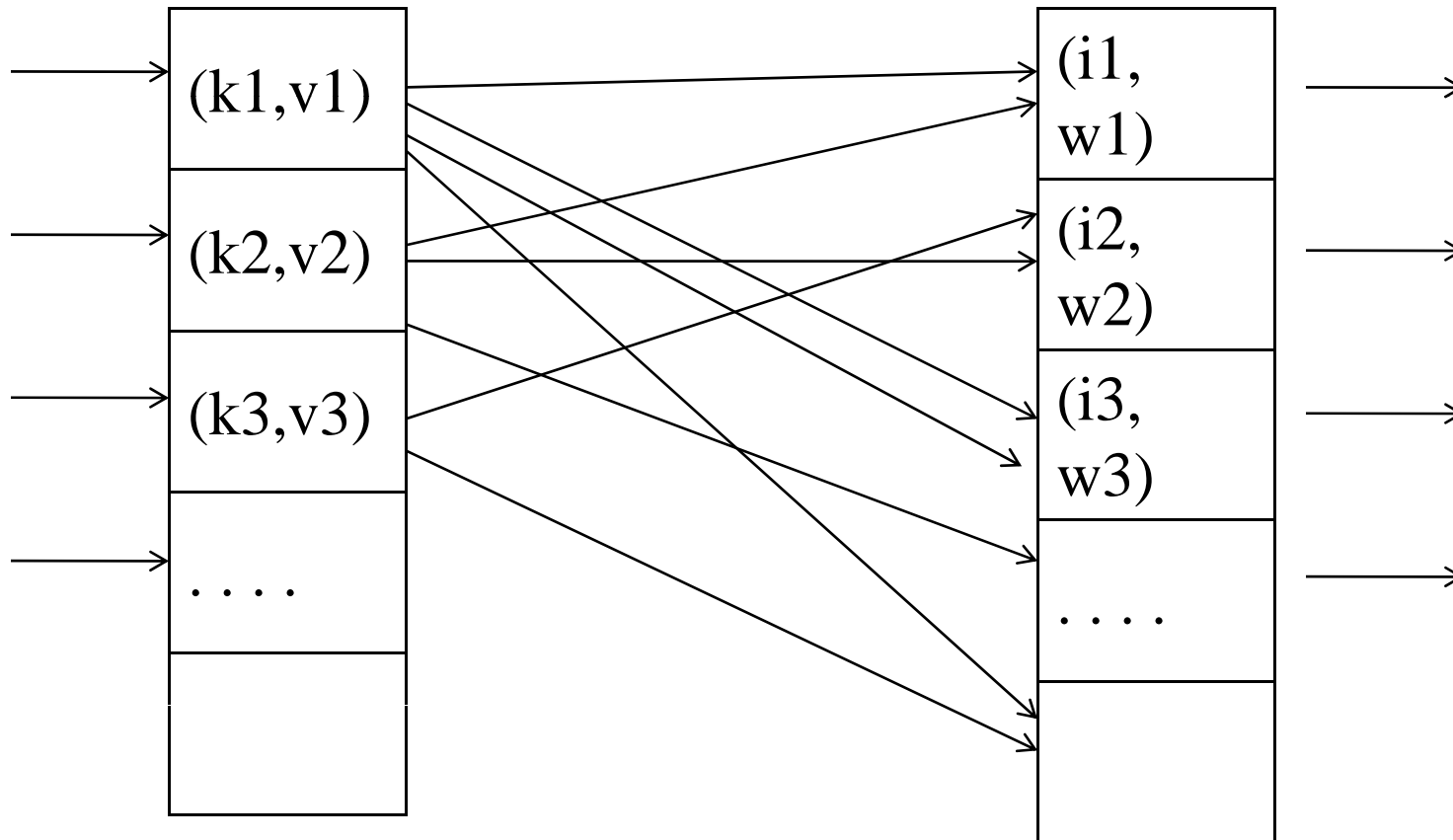
```
map(String key, String value):  
  // key: document name  
  // value: document contents  
  for each word w in value:  
    EmitIntermediate(w, "1");
```

```
reduce(String key, Iterator values):  
  // key: a word  
  // values: a list of counts  
  int result = 0;  
  for each v in values:  
    result += ParseInt(v);  
  Emit(AsString(result));
```

MapReduce Execution

MAP

REDUCE



Map = GROUP BY,
Reduce = Aggregate

R(documentKey, word)

```
SELECT word, sum(1)
FROM R
GROUP BY word
```

Example 2: MR word length count

Abridged Declaration of Independence

A Declaration By the Representatives of the United States of America, in General Congress Assembled. When in the course of human events it becomes necessary for a people to advance from that subordination in which they have hitherto remained, and to assume among powers of the earth the equal and independent station to which the laws of nature and of nature's god entitle them, a decent respect to the opinions of mankind requires that they should declare the causes which impel them to the change.

We hold these truths to be self-evident; that all men are created equal and independent; that from that equal creation they derive rights inherent and inalienable, among which are the preservation of life, and liberty, and the pursuit of happiness; that to secure these ends, governments are instituted among men, deriving their just power from the consent of the governed; that whenever any form of government shall become destructive of these ends, it is the right of the people to alter or to abolish it, and to institute new government, laying it's foundation on such principles and organizing it's power in such form, as to them shall seem most likely to effect their safety and happiness. Prudence indeed will dictate that governments long established should not be changed for light and transient causes: and accordingly all experience hath shewn that mankind are more disposed to suffer while evils are sufferable, than to right themselves by abolishing the forms to which they are accustomed. But when a long train of abuses and usurpations, begun at a distinguished period, and pursuing invariably the same object, evinces a design to reduce them to arbitrary power, it is their right, it is their duty, to throw off such government and to provide new guards for future security. Such has been the patient sufferings of the colonies; and such is now the necessity which constrains them to expunge their former systems of government. the history of his present majesty is a history of unremitting injuries and usurpations, among which no one fact stands single or solitary to contradict the uniform tenor of the rest, all of which have in direct object the establishment of an absolute tyranny over these states. To prove this, let facts be submitted to a candid world, for the truth of which we pledge a faith yet unsullied by falsehood.

Example 2: MR word length count

Abridged Declaration of Independence

Map Task 1
(204 words)

Yellow: 10+

Red: 5..9

Blue: 2..4

Pink: = 1

A Declaration By the Representatives of the United States of America, in General Congress Assembled.
 When in the course of human events it becomes necessary for a people to advance from that subordination in which they have hitherto remained, and to assume among powers of the earth the equal and independent station to which the laws of nature and of nature's god entitle them, a decent respect to the opinions of mankind requires that they should declare the causes which impel them to the change.
 We hold these truths to be self-evident; that all men are created equal and independent; that from that equal creation they derive rights inherent and inalienable, among which are the preservation of life, and liberty, and the pursuit of happiness; that to secure these ends, governments are instituted among men, deriving their just power from the consent of the governed; that whenever any form of government shall become destructive of these ends, it is the right of the people to alter or to abolish it, and to institute new government, laying it's foundation on such principles and organizing it's power in such form, as to them shall seem most likely to effect their safety and happiness. Prudence indeed will

(key, value)

(yellow, 17)

(red, 77)

(blue, 107)

(pink, 3)

Map Task 2
(190 words)

dictate that governments long established should not be changed for light and transient causes: and accordingly all experience hath shewn that mankind are more disposed to suffer while evils are sufferable, than to right themselves by abolishing the forms to which they are accustomed. But when a long train of abuses and usurpations, begun at a distinguished period, and pursuing invariably the same object, evinces a design to reduce them to arbitrary power, it is their right, it is their duty, to throw off such government and to provide new guards for future security. Such has been the patient sufferings of the colonies; and such is now the necessity which constrains them to expunge their former systems of government. the history of his present majesty is a history of unremitting injuries and usurpations, among which no one fact stands single or solitary to contradict the uniform tenor of the rest, all of which have in direct object the establishment of an absolute tyranny over these states. To prove this, let facts be submitted to a candid world, for the truth of which we pledge a faith yet unsullied by falsehood.

(yellow, 20)

(red, 71)

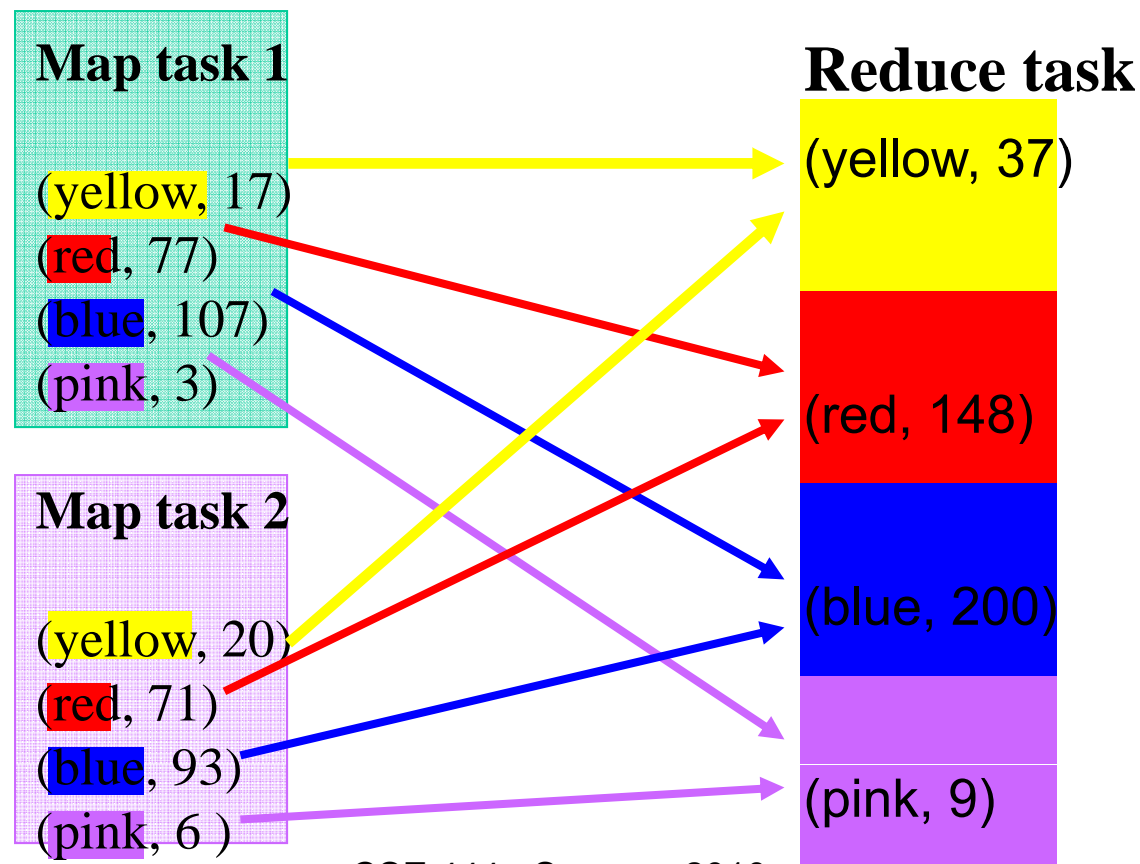
(blue, 93)

(pink, 6)

Example 2: MR word length count

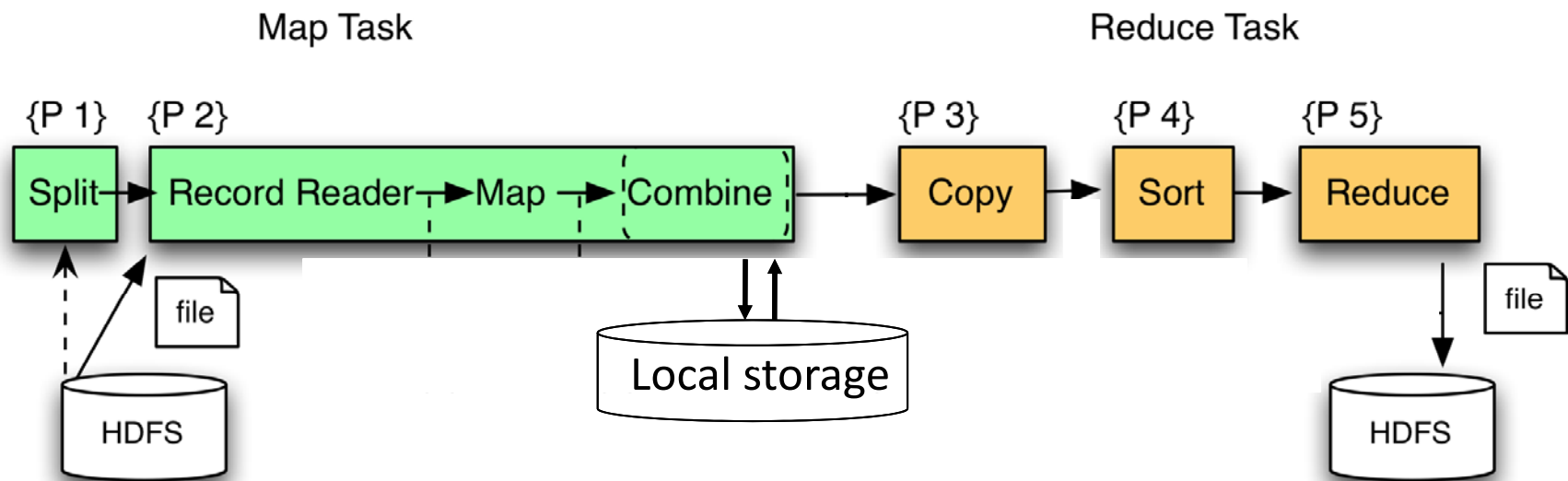
Map is a **GROUP BY** operation

Reduce is an **AGGREGATE** operation



MR Phases

- Each Map and Reduce task has multiple phases:



Implementation

- There is one master node
- Master partitions input file into M *splits*, by key
- Master assigns workers (=servers) to the M map tasks, keeps track of their progress
- Workers write their output to local disk, partition into R regions
- Master assigns workers to the R reduce tasks
- Reduce workers read regions from the map workers' local disks

Interesting Implementation Details

- Worker failure:
 - Master pings workers periodically,
 - If down then reassigns the task to another worker
- Choice of M and R:
 - Larger is better for load balancing
 - Limitation: master needs $O(M \times R)$ memory

Interesting Implementation Details

- Backup tasks:
 - “Straggler” = a machine that takes unusually long time to complete one of the last tasks. Eg:
 - Bad disk forces frequent correctable errors (30MB/s → 1MB/s)
 - The cluster scheduler has scheduled other tasks on that machine
 - Stragglers are a main reason for slowdown
 - Solution: pre-emptive backup execution of the last few remaining in-progress tasks

Map-Reduce Summary

- Hides scheduling and parallelization details
- However, very limited queries
 - Difficult to write more complex tasks
 - Need multiple map-reduce operations
- Solution: PIG Latin (Y!) – open source, etc.
- Others
 - Scope (MS): SQL ! But proprietary...
 - DryadLINQ (MS): LINQ ! But also proprietary...
 - Clustera (other UW) : SQL ! Not publicly available

MapReduce vs Parallel DBMS

- Both provide massive-scale parallel data processing
- **Parallel DBMS benefits**
 - Indexing and partition awareness!
 - Lower per-query start-up overhead
 - Offers much better query processing performance overall
 - Can often run on much fewer nodes than MR
- **MapReduce benefits**
 - Easier to setup and configure
 - No data loading
 - Great approach for parallelizing user-defined functions
 - Intra-query fault-tolerance.

Preview of Final Project

- Process large data using Pig/Hadoop
 - Develop locally on desktop/laptop
 - Large runs on Amazon Web Services
- Read *Pig Latin: A not-so-foreign language for data processing*, by Olston, Reed, Srivastava, Kumar, and Tomkins, 2008
 - Next couple of lectures based on this