# MapReduce Algorithm Design

## WWW 2013 Tutorial, Rio de Janeiro

Jimmy Lin
University of Maryland
Monday, May 13, 2013

**From the Ivory Tower…**

... to building sh*t that works

... and back.

# More about me…

- Past MapReduce teaching experience:
  - Numerous tutorials
  - Several semester-long MapReduce courses
    `http://lintool.github.io/MapReduce-course-2013s/`

- Lin & Dyer MapReduce textbook
  `http://mapreduce.cc/`



Follow me at @lintool

# What we'll cover

- Big data

- MapReduce overview

- Importance of local aggregation

- Sequencing computations

- Iterative graph algorithms

- MapReduce and abstract algebra

**Focus on design patterns and general principles**

# What we won't cover

- MapReduce for machine learning (supervised and unsupervised)

- MapReduce for similar item detection

- MapReduce for information retrieval

- Hadoop for data warehousing

- Extensions and alternatives to MapReduce

**Big Data**

**Google** processes 20 PB a day (2008)
crawls 20B web pages a day (2012)

**JPMorganChase** 150 PB on 50k+ servers
running 15k apps (6/2011)

**ebay** >10 PB data, 75B DB
calls per day (6/2012)

**INTERNET ARCHIVE** Wayback Machine: 240B web
pages archived, 5 PB (1/2013)

>100 PB of user data +
500 TB/day (8/2012) **facebook**

LHC: ~15 PB a year **CERN**

**amazon webservices** S3: 449B objects, peak 290k
request/second (7/2011)
1T objects (6/2012)

LSST: 6-10 PB a year
(~2015)

640K ought to be
enough for anybody.

SKA: 0.3 – 1.5 EB
per year (~2020)

# How much data?

# Why big data?

Science
Engineering
Commerce

# Science

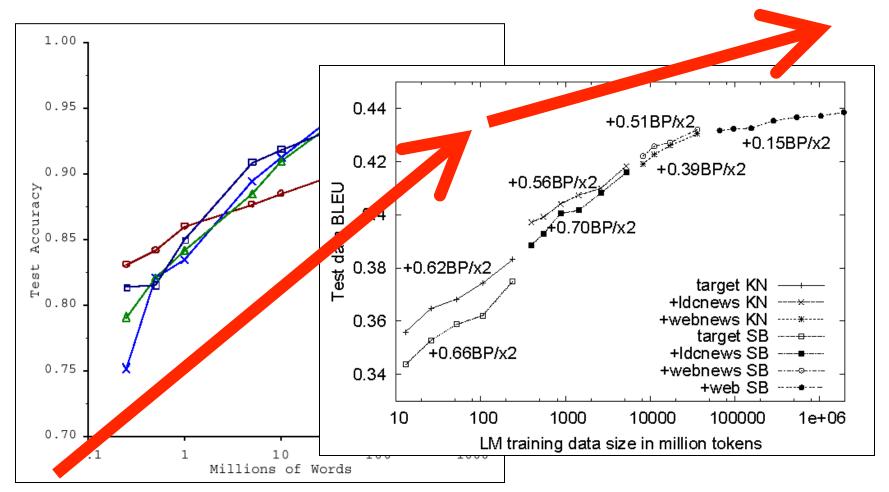Emergence of the 4<sup>th</sup> Paradigm

Data-intensive e-Science

# Engineering

The unreasonable effectiveness of data

Count and normalize!

# No data like more data!

s/knowledge/data/g;



+0.51BP/x2

+0.15BP/x2

+0.56BP/x2

+0.39BP/x2

+0.70BP/x2

+0.62BP/x2

+0.66BP/x2

Test data BLEU

| | |
|---|---|
| target KN | |
| +ldcnews KN | |
| +webnews KN | |
| target SB | |
| +ldcnews SB | |
| +webnews SB | |
| +web SB | |

LM training data size in million tokens

Test Accuracy

Millions of Words

(Banko and Brill, ACL 2001)
(Brants et al., EMNLP 2007)

Know thy customers

Data → Insights → Competitive advantages

# Commerce

**Why big data?**
**How big data?**

Source: Wikipedia (Noctilucent cloud)

MapReduce

# Typical Big Data Problem

- Iterate over a large number of records

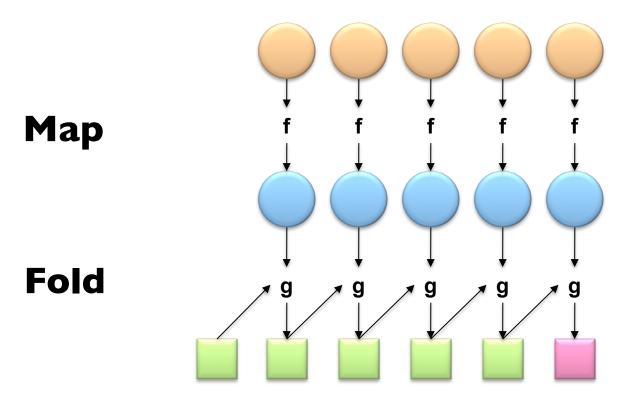**Map** - Extract something of interest from each

- Shuffle and sort intermediate results

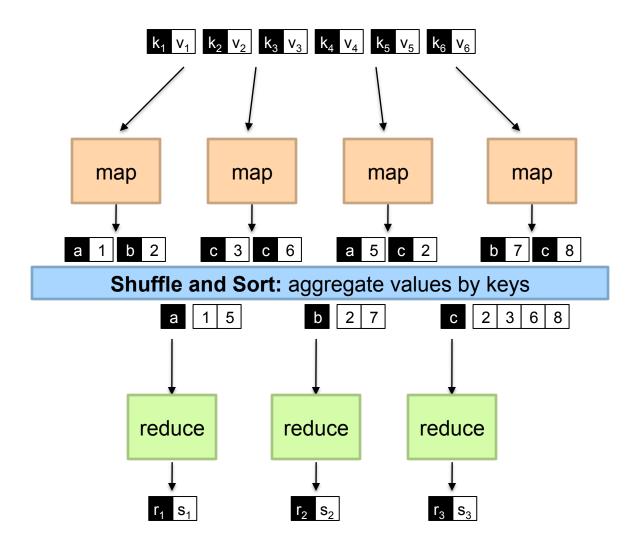- Aggregate intermediate results **Reduce**

- Generate final output

**Key idea: provide a functional abstraction for these two operations**

# Roots in Functional Programming

**Map**

**Fold**

# MapReduce

- Programmers specify two functions:

  **map** $(k_1, v_1) \rightarrow [<k_2, v_2>]$
  **reduce** $(k_2, [v_2]) \rightarrow [<k_3, v_3>]$
  - All values with the same key are sent to the same reducer

- The execution framework handles everything else…

# MapReduce

- Programmers specify two functions:

  **map** (k, v) → <k', v'>*
  **reduce** (k', v') → <k', v'>*
  - All values with the same key are sent to the same reducer

- The execution framework handles everything else…

**What's "everything else"?**
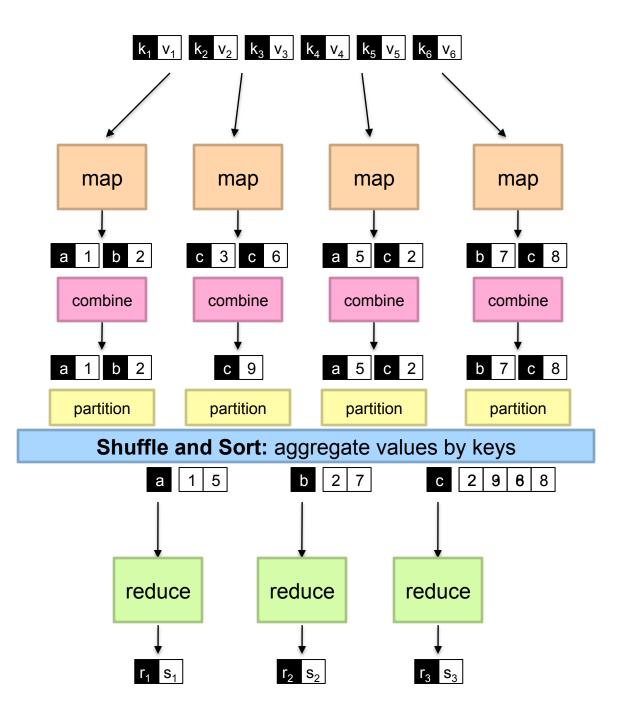
# MapReduce "Runtime"

- Handles scheduling
  - Assigns workers to map and reduce tasks
- Handles "data distribution"
  - Moves processes to data
- Handles synchronization
  - Gathers, sorts, and shuffles intermediate data
- Handles errors and faults
  - Detects worker failures and restarts
- Everything happens on top of a distributed filesystem

# MapReduce

○ Programmers specify two functions:

**map** (k, v) → <k', v'>*
**reduce** (k', v') → <k', v'>*
- All values with the same key are reduced together

○ The execution framework handles everything else…

○ Not quite…usually, programmers also specify:

**partition** (k', number of partitions) → partition for k'
- Often a simple hash of the key, e.g., hash(k') mod n
- Divides up key space for parallel reduce operations

**combine** (k', v') → <k', v'>*
- Mini-reducers that run in memory after the map phase
- Used as an optimization to reduce network traffic

# Two more details...

- Barrier between map and reduce phases

  - But intermediate data can be copied over as soon as mappers finish

- Keys arrive at each reducer in sorted order

  - No enforced ordering *across* reducers

# What's the big deal?

○ Developers need the right level of abstraction

- Moving beyond the von Neumann architecture
- We need better programming models

○ Abstractions hide low-level details from the developers

- No more race conditions, lock contention, etc.

○ MapReduce separating the *what* from *how*

- Developer specifies the computation that needs to be performed
- Execution framework ("runtime") handles actual execution

# The datacenter *is* the computer!

Source: Google

# MapReduce can refer to...

- The programming model

- The execution framework (aka "runtime")

- The specific implementation

**Usage is usually clear from context!**

# MapReduce Implementations

- Google has a proprietary implementation in C++
  - Bindings in Java, Python

- Hadoop is an open-source implementation in Java
  - Development led by Yahoo, now an Apache project
  - Used in production at Yahoo, Facebook, Twitter, LinkedIn, Netflix, …
  - The *de facto* big data processing platform
  - Rapidly expanding software ecosystem

- Lots of custom research implementations
  - For GPUs, cell processors, etc.

# MapReduce algorithm design

○ The execution framework handles "everything else"…

- Scheduling: assigns workers to map and reduce tasks
- "Data distribution": moves processes to data
- Synchronization: gathers, sorts, and shuffles intermediate data
- Errors and faults: detects worker failures and restarts

○ Limited control over data and execution flow

- All algorithms must expressed in m, r, c, p

○ You don't know:

- Where mappers and reducers run
- When a mapper or reducer begins or finishes
- Which input a particular mapper is processing
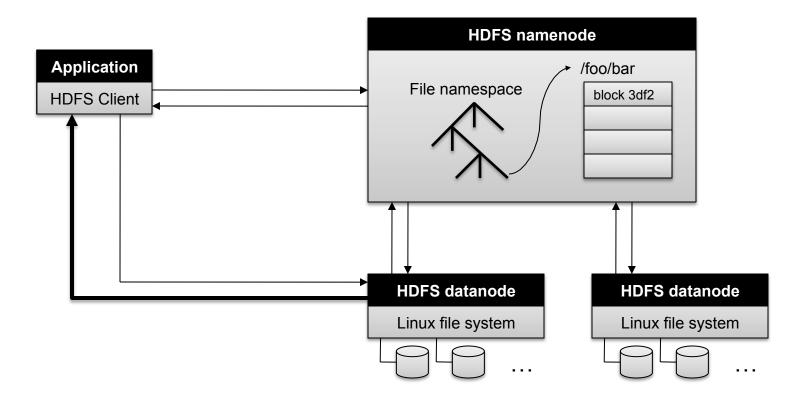- Which intermediate key a particular reducer is processing

# Implementation Details

# HDFS Architecture
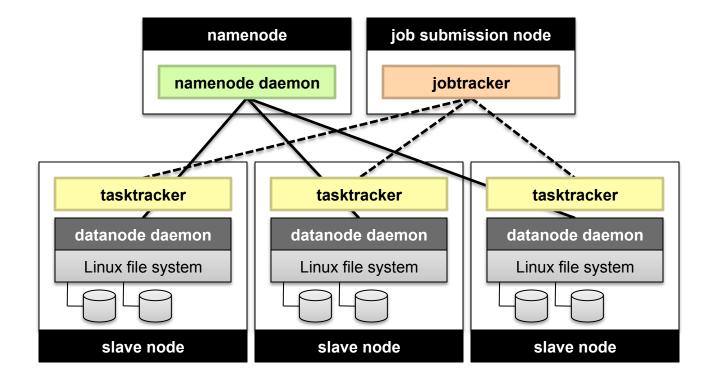


**Application**
HDFS Client

**HDFS namenode**
File namespace
/foo/bar
block 3df2

**HDFS datanode**
Linux file system
…

**HDFS datanode**
Linux file system
…

# Putting everything together...

# Shuffle and Sort

Mapper

circular buffer
(in memory)

spills (on disk)

Combiner

merged spills
(on disk)

intermediate files
(on disk)

Combiner

Reducer

other reducers

other mappers

# Preserving State

# Implementation **Don'ts**

- Don't unnecessarily create objects

  - Object creation is costly

  - Garbage collection is costly

- Don't buffer objects

  - Processes have limited heap size (remember, commodity machines)

  - May work for small datasets, but won't scale!

# Secondary Sorting

- MapReduce sorts input to reducers by key

  - Values may be arbitrarily ordered

- What if want to sort value also?

  - E.g., $k \rightarrow (v_1, r), (v_3, r), (v_4, r), (v_8, r)\ldots$

# Secondary Sorting: Solutions

- Solution 1:
    - Buffer values in memory, then sort
    - Why is this a bad idea?

- Solution 2:
    - "Value-to-key conversion" design pattern: form composite intermediate key, $(k, v_1)$
    - Let execution framework do the sorting
    - Preserve state across multiple key-value pairs to handle processing
    - Anything else we need to do?

**Local Aggregation**

# Importance of Local Aggregation

- Ideal scaling characteristics:
    - Twice the data, twice the running time
    - Twice the resources, half the running time

- Why can't we achieve this?
    - Synchronization requires communication
    - Communication kills performance (network is slow!)

- Thus… avoid communication!
    - Reduce intermediate data via local aggregation
    - Combiners can help

# Word Count: Baseline

```
1: class MAPPER
2:     method MAP(docid a, doc d)
3:         for all term t ∈ doc d do
4:             EMIT(term t, count 1)

1: class REDUCER
2:     method REDUCE(term t, counts [c_1, c_2, . . .])
3:         sum ← 0
4:         for all count c ∈ counts [c_1, c_2, . . .] do
5:             sum ← sum + c
6:         EMIT(term t, count s)
```

**What's the impact of combiners?**

# Word Count: Version 1

```
1:  class MAPPER
2:      method MAP(docid a, doc d)
3:          H ← new ASSOCIATIVEARRAY
4:          for all term t ∈ doc d do
5:              H{t} ← H{t} + 1                    ▷ Tally counts for entire document
6:          for all term t ∈ H do
7:              EMIT(term t, count H{t})
```

**Are combiners still needed?**

# Word Count: Version 2

```
1: class MAPPER
2:     method INITIALIZE
3:         H ← new ASSOCIATIVEARRAY
4:     method MAP(docid a, doc d)
5:         for all term t ∈ doc d do
6:             H{t} ← H{t} + 1                    ▷ Tally counts across documents
7:     method CLOSE
8:         for all term t ∈ H do
9:             EMIT(term t, count H{t})
```

Key idea: preserve state across input key-value pairs!

**Are combiners still needed?**

# Design Pattern for Local Aggregation

- "In-mapper combining"

  - Fold the functionality of the combiner into the mapper by preserving state across multiple map calls

- Advantages

  - Speed

  - Why is this faster than actual combiners?

- Disadvantages

  - Explicit memory management required

  - Potential for order-dependent bugs

# Combiner Design

- Combiners and reducers share same method signature

  - Sometimes, reducers can serve as combiners
  - Often, not…

- Remember: combiner are optional optimizations

  - Should not affect algorithm correctness
  - May be run 0, 1, or multiple times

- Example: find average of integers associated with the same key

# Computing the Mean: Version 1

```
1: class MAPPER
2:     method MAP(string t, integer r)
3:         EMIT(string t, integer r)

1: class REDUCER
2:     method REDUCE(string t, integers [r_1, r_2, . . .])
3:         sum ← 0
4:         cnt ← 0
5:         for all integer r ∈ integers [r_1, r_2, . . .] do
6:             sum ← sum + r
7:             cnt ← cnt + 1
8:         r_avg ← sum/cnt
9:         EMIT(string t, integer r_avg)
```

**Why can't we use reducer as combiner?**

# Computing the Mean: Version 2

```
1: class MAPPER
2:     method MAP(string t, integer r)
3:         EMIT(string t, integer r)

1: class COMBINER
2:     method COMBINE(string t, integers [r_1, r_2, ...])
3:         sum ← 0
4:         cnt ← 0
5:         for all integer r ∈ integers [r_1, r_2, ...] do
6:             sum ← sum + r
7:             cnt ← cnt + 1
8:         EMIT(string t, pair (sum, cnt))                    ▷ Separate sum and count

1: class REDUCER
2:     method REDUCE(string t, pairs [(s_1, c_1), (s_2, c_2) ...])
3:         sum ← 0
4:         cnt ← 0
5:         for all pair (s, c) ∈ pairs [(s_1, c_1), (s_2, c_2) ...] do
6:             sum ← sum + s
7:             cnt ← cnt + c
8:         r_avg ← sum/cnt
9:         EMIT(string t, integer r_avg)
```

**Why doesn't this work?**

# Computing the Mean: Version 3

```
1: class MAPPER
2:     method MAP(string t, integer r)
3:         EMIT(string t, pair (r, 1))

1: class COMBINER
2:     method COMBINE(string t, pairs [(s₁, c₁), (s₂, c₂)...])
3:         sum ← 0
4:         cnt ← 0
5:         for all pair (s, c) ∈ pairs [(s₁, c₁), (s₂, c₂)...] do
6:             sum ← sum + s
7:             cnt ← cnt + c
8:         EMIT(string t, pair (sum, cnt))

1: class REDUCER
2:     method REDUCE(string t, pairs [(s₁, c₁), (s₂, c₂)...])
3:         sum ← 0
4:         cnt ← 0
5:         for all pair (s, c) ∈ pairs [(s₁, c₁), (s₂, c₂)...] do
6:             sum ← sum + s
7:             cnt ← cnt + c
8:         r_avg ← sum/cnt
9:         EMIT(string t, pair (r_avg, cnt))
```

**Fixed?**

# Computing the Mean: Version 4

```
1: class MAPPER
2:     method INITIALIZE
3:         S ← new ASSOCIATIVEARRAY
4:         C ← new ASSOCIATIVEARRAY
5:     method MAP(string t, integer r)
6:         S{t} ← S{t} + r
7:         C{t} ← C{t} + 1
8:     method CLOSE
9:         for all term t ∈ S do
10:            EMIT(term t, pair (S{t}, C{t}))
```

**Are combiners still needed?**

# Sequencing Computations

# Sequencing Computations

1. Turn synchronization into a sorting problem

   - Leverage the fact that keys arrive at reducers in sorted order
   - Manipulate the sort order and partitioning scheme to deliver partial results at appropriate junctures

2. Create appropriate algebraic structures to capture computation

   - Build custom data structures to accumulate partial results

# Algorithm Design: Running Example

- Term co-occurrence matrix for a text collection

  - M = N x N matrix (N = vocabulary size)
  - $M_{ij}$: number of times $i$ and $j$ co-occur in some context
    (for concreteness, let's say context = sentence)

- Why?

  - Distributional profiles as a way of measuring semantic distance
  - Semantic distance useful for many language processing tasks
  - Basis for large classes of more sophisticated algorithms

# MapReduce: Large Counting Problems

○ Term co-occurrence matrix for a text collection
  = specific instance of a large counting problem

  - A large event space (number of terms)

  - A large number of observations (the collection itself)

  - Goal: keep track of interesting statistics about the events

○ Basic approach

  - Mappers generate partial counts

  - Reducers aggregate partial counts

**How do we aggregate partial counts efficiently?**

# First Try: "Pairs"

- Each mapper takes a sentence:
  - Generate all co-occurring term pairs
  - For all pairs, emit (a, b) → count

- Reducers sum up counts associated with these pairs

- Use combiners!

# Pairs: Pseudo-Code

```
1: class MAPPER
2:     method MAP(docid a, doc d)
3:         for all term w ∈ doc d do
4:             for all term u ∈ NEIGHBORS(w) do
5:                 EMIT(pair (w, u), count 1)        ▷ Emit count for each co-occurrence

1: class REDUCER
2:     method REDUCE(pair p, counts [c₁, c₂, . . .])
3:         s ← 0
4:         for all count c ∈ counts [c₁, c₂, . . .] do
5:             s ← s + c                              ▷ Sum co-occurrence counts
6:         EMIT(pair p, count s)
```

# "Pairs" Analysis

- Advantages

  - Easy to implement, easy to understand

- Disadvantages

  - Lots of pairs to sort and shuffle around (upper bound?)
  - Not many opportunities for combiners to work

# Another Try: "Stripes"

○ Idea: group together pairs into an associative array

$(a, b) \rightarrow 1$
$(a, c) \rightarrow 2$
$(a, d) \rightarrow 5$          $a \rightarrow \{ b: 1, c: 2, d: 5, e: 3, f: 2 \}$
$(a, e) \rightarrow 3$
$(a, f) \rightarrow 2$

○ Each mapper takes a sentence:

  ● Generate all co-occurring term pairs

  ● For each term, emit $a \rightarrow \{ b: count_b, c: count_c, d: count_d \dots \}$

○ Reducers perform element-wise sum of associative arrays

$a \rightarrow \{ b: 1, \qquad d: 5, e: 3 \}$
$+ \quad a \rightarrow \{ b: 1, c: 2, d: 2, \qquad f: 2 \}$
_____
$a \rightarrow \{ b: 2, c: 2, d: 7, e: 3, f: 2 \}$

Key idea: cleverly-constructed data structure
for aggregating partial results

# Stripes: Pseudo-Code

```
1: class MAPPER
2:     method MAP(docid a, doc d)
3:         for all term w ∈ doc d do
4:             H ← new ASSOCIATIVEARRAY
5:             for all term u ∈ NEIGHBORS(w) do
6:                 H{u} ← H{u} + 1              ▷ Tally words co-occurring with w
7:             EMIT(Term w, Stripe H)

1: class REDUCER
2:     method REDUCE(term w, stripes [H₁, H₂, H₃, . . .])
3:         H_f ← new ASSOCIATIVEARRAY
4:         for all stripe H ∈ stripes [H₁, H₂, H₃, . . .] do
5:             SUM(H_f, H)                              ▷ Element-wise sum
6:         EMIT(term w, stripe H_f)
```

# "Stripes" Analysis

- Advantages

  - Far less sorting and shuffling of key-value pairs
  - Can make better use of combiners

- Disadvantages

  - More difficult to implement
  - Underlying object more heavyweight
  - Fundamental limitation in terms of size of event space

# Comparison of "pairs" vs. "stripes" for computing word co-occurrence matrices



$R^2 = 0.999$

$R^2 = 0.992$

**Cluster size:** 38 cores
**Data Source:** Associated Press Worldstream (APW) of the English Gigaword Corpus (v3),
which contains 2.27 million documents (1.8 GB compressed, 5.7 GB uncompressed)

# Effect of cluster size on "stripes" algorithm

relative size of EC2 cluster



$R^2 = 0.997$

running time (seconds)

relative speedup

size of EC2 cluster (number of slave instances)

# Relative Frequencies

- How do we estimate relative frequencies from counts?

$$f(B|A) = \frac{N(A,B)}{N(A)} = \frac{N(A,B)}{\sum_{B'} N(A,B')}$$

- Why do we want to do this?

- How do we do this with MapReduce?
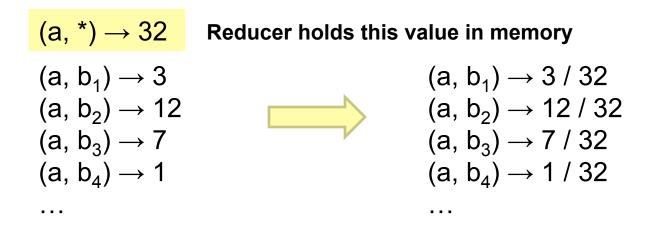
# f(B|A): "Stripes"

$a \rightarrow \{b_1:3, b_2:12, b_3:7, b_4:1, \ldots\}$

- Easy!
  - One pass to compute (a, *)
  - Another pass to directly compute f(B|A)

# f(B|A): "Pairs"

- What's the issue?

  - Computing relative frequencies requires marginal counts
  - But the marginal cannot be computed until you see all counts
  - Buffering is a bad idea!

- Solution:

  - What if we could get the marginal count to arrive at the reducer first?

# f(B|A): "Pairs"

(a, *) → 32     **Reducer holds this value in memory**

(a, $b_1$) → 3          (a, $b_1$) → 3 / 32
(a, $b_2$) → 12         (a, $b_2$) → 12 / 32
(a, $b_3$) → 7          (a, $b_3$) → 7 / 32
(a, $b_4$) → 1          (a, $b_4$) → 1 / 32
…                       …

○ For this to work:

- Must emit extra (a, *) for every $b_n$ in mapper
- Must make sure all a's get sent to same reducer (use partitioner)
- Must make sure (a, *) comes first (define sort order)
- Must hold state in reducer across different key-value pairs

# "Order Inversion"

- Common design pattern:

  - Take advantage of sorted key order at reducer to sequence computations

  - Get the marginal counts to arrive at the reducer before the joint counts

- Optimization:

  - Apply in-memory combining pattern to accumulate marginal counts

# Synchronization: Pairs vs. Stripes

- Approach 1: turn synchronization into an ordering problem
    - Sort keys into correct order of computation
    - Partition key space so that each reducer gets the appropriate set of partial results
    - Hold state in reducer across multiple key-value pairs to perform computation
    - Illustrated by the "pairs" approach

- Approach 2: construct data structures to accumulate partial results
    - Each reducer receives all the data it needs to complete the computation
    - Illustrated by the "stripes" approach

# Issues and Tradeoffs

○ Number of key-value pairs

- Object creation overhead
- Time for sorting and shuffling pairs across the network

○ Size of each key-value pair

- De/serialization overhead

**Lots are algorithms are just fancy conditional counts!**

# Hidden Markov Models

An HMM $\lambda = (A, B, \Pi)$ is characterized by:

- N states: $Q = \{q_1, q_2, \ldots q_N\}$
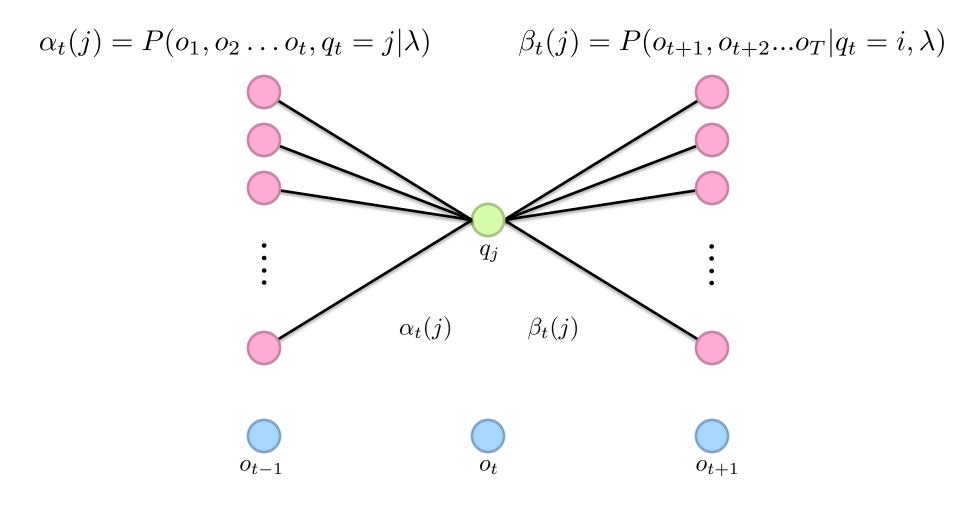- N x N Transition probability matrix $A = [a_{ij}]$

$$a_{ij} = p(q_j | q_i) \qquad \sum_j a_{ij} = 1 \quad \forall i$$

- V observation symbols: $O = \{o_1, o_2, \ldots o_V\}$
- N x |V| Emission probability matrix $B = [b_{iv}]$

$$b_{iv} = b_i(o_v) = p(o_v | q_i)$$

- Prior probabilities vector $\Pi = [\pi_i, \pi_2, \ldots \pi_N]$

$$\sum_{i=1}^{N} \pi_i = 1$$

# Forward-Backward

$$\alpha_t(j) = P(o_1, o_2 \ldots o_t, q_t = j | \lambda) \qquad \beta_t(j) = P(o_{t+1}, o_{t+2} \ldots o_T | q_t = i, \lambda)$$



$q_j$

$\alpha_t(j)$     $\beta_t(j)$

$o_{t-1}$           $o_t$           $o_{t+1}$

# Estimating Emissions Probabilities

- Basic idea:

$$b_j(v_k) = \frac{\text{expected number of times in state } j \text{ and observing symbol } v_k}{\text{expected number of times in state } j}$$

- Let's define:

$$\gamma_t(j) = \frac{P(q_t = j, O|\lambda)}{P(O|\lambda)} = \frac{\alpha_t(j)\beta_t(j)}{P(O|\lambda)}$$

- Thus:

$$\hat{b}_j(v_k) = \frac{\sum_{i=1 \cap O_t = v_k}^{T} \gamma_t(j)}{\sum_{i=1}^{T} \gamma_t(j)}$$

# Forward-Backward

# Estimating Transition Probabilities

- Basic idea:

$$a_{ij} = \frac{\text{expected number of transitions from state } i \text{ to state } j}{\text{expected number of transitions from state } i}$$

- Let's define:

$$\xi_t(i,j) = \frac{\alpha_t(i)a_{ij}b_j(o_{t+1})\beta_{t+1}(j)}{P(O|\lambda)}$$

- Thus:

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i,j)}{\sum_{t=1}^{T-1} \sum_{j=1}^{N} \xi_t(i,j)}$$

# MapReduce Implementation: Mapper

```
 1: class MAPPER
 2:    method INITIALIZE(integer iteration)
 3:        ⟨S, O⟩ ← READMODEL
 4:        θ ← ⟨A, B, π⟩ ← READMODELPARAMS(iteration)
 5:    method MAP(sample id, sequence x)
 6:        α ← FORWARD(x, θ)
 7:        β ← BACKWARD(x, θ)
 8:        I ← new ASSOCIATIVEARRAY
 9:        for all q ∈ S do
10:            I{q} ← α₁(q) · β₁(q)
11:        O ← new ASSOCIATIVEARRAY of ASSOCIATIVEARRAY
12:        for t = 1 to |x| do
13:            for all q ∈ S do
14:                O{q}{xₜ} ← O{q}{xₜ} + αₜ(q) · βₜ(q)
15:            t ← t + 1
16:        T ← new ASSOCIATIVEARRAY of ASSOCIATIVEARRAY
17:        for t = 1 to |x| − 1 do
18:            for all q ∈ S do
19:                for all r ∈ S do
20:                    T{q}{r} ← T{q}{r} + αₜ(q) · A_q(r) · B_r(x_{t+1}) · β_{t+1}(r)
21:            t ← t + 1
22:        EMIT(string 'initial', stripe I)
23:        for all q ∈ S do
24:            EMIT(string 'emit from ' + q, stripe O{q})
25:            EMIT(string 'transit from ' + q, stripe T{q})
```

$$\hat{b}_j(v_k) = \frac{\sum_{i=1 \cap O_t = v_k}^{T} \gamma_t(j)}{\sum_{i=1}^{T} \gamma_t(j)}$$

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i,j)}{\sum_{t=1}^{T-1} \sum_{j=1}^{N} \xi_t(i,j)}$$

$$\gamma_t(j) = \frac{\alpha_t(j)\beta_t(j)}{P(O|\lambda)}$$

$$\xi_t(i,j) = \frac{\alpha_t(i)a_{ij}b_j(o_{t+1})\beta_{t+1}(j)}{P(O|\lambda)}$$

# MapReduce Implementation: Reducer

```
1: class COMBINER
2:     method COMBINE(string t, stripes [C_1, C_2, ...])
3:         C_f ← new ASSOCIATIVEARRAY
4:         for all stripe C ∈ stripes [C_1, C_2, ...] do
5:             SUM(C_f, C)
6:         EMIT(string t, stripe C_f)

1: class REDUCER
2:     method REDUCE(string t, stripes [C_1, C_2, ...])
3:         C_f ← new ASSOCIATIVEARRAY
4:         for all stripe C ∈ stripes [C_1, C_2, ...] do
5:             SUM(C_f, C)
6:         z ← 0
7:         for all ⟨k, v⟩ ∈ C_f do
8:             z ← z + v
9:         P_f ← new ASSOCIATIVEARRAY
10:        for all ⟨k, v⟩ ∈ C_f do
11:            P_f{k} ← v/z
12:        EMIT(string t, stripe P_f)
```

$$\hat{b}_j(v_k) = \frac{\sum_{i=1 \cap O_t=v_k}^{T} \gamma_t(j)}{\sum_{i=1}^{T} \gamma_t(j)}$$

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i,j)}{\sum_{t=1}^{T-1} \sum_{j=1}^{N} \xi_t(i,j)}$$

$$\gamma_t(j) = \frac{\alpha_t(j)\beta_t(j)}{P(O|\lambda)}$$

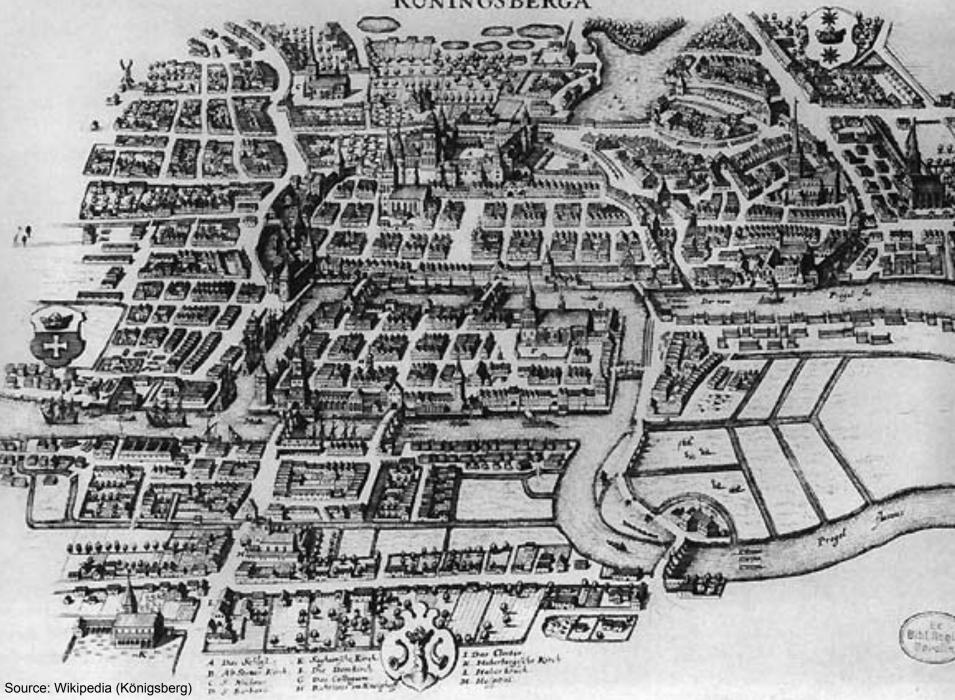$$\xi_t(i,j) = \frac{\alpha_t(i)a_{ij}b_j(o_{t+1})\beta_{t+1}(j)}{P(O|\lambda)}$$

# Iterative Algorithms: Graphs

# What's a graph?

- G = (V,E), where

  - V represents the set of vertices (nodes)
  - E represents the set of edges (links)
  - Both vertices and edges may contain additional information

- Different types of graphs:

  - Directed vs. undirected edges
  - Presence or absence of cycles

- Graphs are everywhere:

  - Hyperlink structure of the web
  - Physical structure of computers on the Internet
  - Interstate highway system
  - Social networks

# KONINGSBERGA

# Some Graph Problems

- Finding shortest paths

  - Routing Internet traffic and UPS trucks

- Finding minimum spanning trees

  - Telco laying down fiber

- Finding Max Flow

  - Airline scheduling

- Identify "special" nodes and communities

  - Breaking up terrorist cells, spread of avian flu

- Bipartite matching

  - Monster.com, Match.com

- And of course... PageRank

# Graphs and MapReduce

- A large class of graph algorithms involve:
  - Performing computations at each node: based on node features, edge features, and local link structure
  - Propagating computations: "traversing" the graph

- Key questions:
  - How do you represent graph data in MapReduce?
  - How do you traverse a graph in MapReduce?

*In reality: graph algorithms in MapReduce suck!*
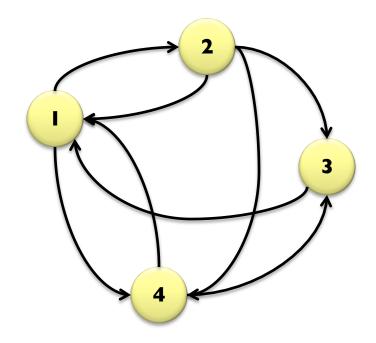
# Representing Graphs

- G = (V, E)

- Two common representations
  - Adjacency matrix
  - Adjacency list

# Adjacency Matrices

Represent a graph as an *n* x *n* square matrix *M*

- $n = |V|$
- $M_{ij} = 1$ means a link from node *i* to *j*

|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| **1** | 0 | 1 | 0 | 1 |
| **2** | 1 | 0 | 1 | 1 |
| **3** | 1 | 0 | 0 | 0 |
| **4** | 1 | 0 | 1 | 0 |

# Adjacency Matrices: Critique

○ Advantages:

- Amenable to mathematical manipulation
- Iteration over rows and columns corresponds to computations on outlinks and inlinks

○ Disadvantages:

- Lots of zeros for sparse matrices
- Lots of wasted space

# Adjacency Lists

Take adjacency matrices… and throw away all the zeros

|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| **1** | 0 | 1 | 0 | 1 |
| **2** | 1 | 0 | 1 | 1 |
| **3** | 1 | 0 | 0 | 0 |
| **4** | 1 | 0 | 1 | 0 |

1: 2, 4
2: 1, 3, 4
3: 1
4: 1, 3

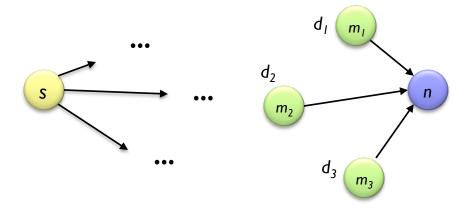# Adjacency Lists: Critique

- Advantages:

  - Much more compact representation

  - Easy to compute over outlinks

- Disadvantages:

  - Much more difficult to compute over inlinks

# Single-Source Shortest Path

- **Problem:** find shortest path from a source node to one or more target nodes

  - Shortest might also mean lowest weight or cost

- Single processor machine: Dijkstra's Algorithm

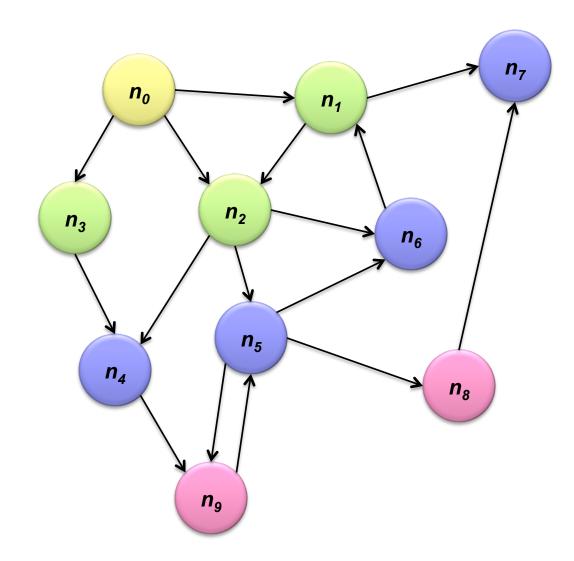- MapReduce: parallel breadth-first search (BFS)

# Finding the Shortest Path

- Consider simple case of equal edge weights

- Solution to the problem can be defined inductively

- Here's the intuition:

  - Define: $b$ is reachable from $a$ if $b$ is on adjacency list of $a$

    $\textsc{DistanceTo}(s) = 0$

  - For all nodes $p$ reachable from $s$,

    $\textsc{DistanceTo}(p) = 1$

  - For all nodes $n$ reachable from some other set of nodes $M$,

    $\textsc{DistanceTo}(n) = 1 + \min(\textsc{DistanceTo}(m), m \in M)$

# Visualizing Parallel BFS

# From Intuition to Algorithm

○ Data representation:

- Key: node $n$
- Value: $d$ (distance from start), adjacency list (nodes reachable from $n$)
- Initialization: for all nodes except for start node, $d = \infty$

○ Mapper:

- $\forall m \in$ adjacency list: emit $(m, d + 1)$

○ Sort/Shuffle

- Groups distances by reachable nodes

○ Reducer:

- Selects minimum distance path for each reachable node
- Additional bookkeeping needed to keep track of actual path

# Multiple Iterations Needed

- Each MapReduce iteration advances the "frontier" by one hop

  - Subsequent iterations include more and more reachable nodes as frontier expands
  - Multiple iterations are needed to explore entire graph

- Preserving graph structure:

  - Problem: Where did the adjacency list go?
  - Solution: mapper emits ($n$, adjacency list) as well

# BFS Pseudo-Code

```
1: class MAPPER
2:    method MAP(nid n, node N)
3:       d ← N.DISTANCE
4:       EMIT(nid n, N)                              ▷ Pass along graph structure
5:       for all nodeid m ∈ N.ADJACENCYLIST do
6:          EMIT(nid m, d + 1)                       ▷ Emit distances to reachable nodes

1: class REDUCER
2:    method REDUCE(nid m, [d₁, d₂, . . .])
3:       d_min ← ∞
4:       M ← ∅
5:       for all d ∈ counts [d₁, d₂, . . .] do
6:          if ISNODE(d) then
7:             M ← d                                 ▷ Recover graph structure
8:          else if d < d_min then                   ▷ Look for shorter distance
9:             d_min ← d
10:      M.DISTANCE ← d_min                          ▷ Update shortest distance
11:      EMIT(nid m, node M)
```

# Stopping Criterion

- When a node is first discovered, we've found the shortest path

  - Maximum number of iterations is equal to the diameter of the graph

- Practicalities of implementation in MapReduce

# Comparison to Dijkstra

○ Dijkstra's algorithm is more efficient

  ● At each step, only pursues edges from minimum-cost path inside frontier

○ MapReduce explores all paths in parallel

  ● Lots of "waste"

  ● Useful work is only done at the "frontier"

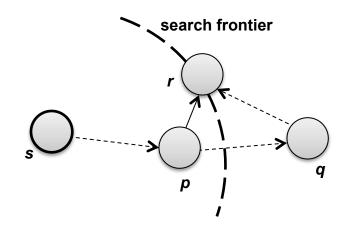○ Why can't we do better using MapReduce?
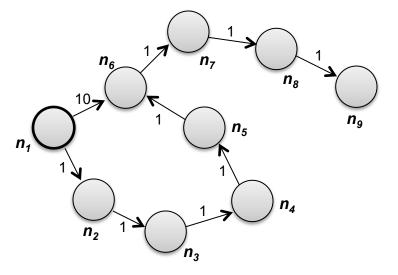
# Single Source: Weighted Edges

- Now add positive weights to the edges

  - Why can't edge weights be negative?

- Simple change: add weight $w$ for each edge in adjacency list

  - In mapper, emit $(m, d + w_p)$ instead of $(m, d + 1)$ for each node $m$

- That's it?

# Stopping Criterion

○ How many iterations are needed in parallel BFS (positive edge weight case)?

○ When a node is first discovered, we've found the shortest path

*Not true!*

# Additional Complexities

# Stopping Criterion

- How many iterations are needed in parallel BFS (positive edge weight case)?

- Practicalities of implementation in MapReduce

# All-Pairs?

- Floyd-Warshall Algorithm: difficult to MapReduce-ify…

- Multiple-source shortest paths in MapReduce: run multiple parallel BFS *simultaneously*

  - Assume source nodes $\{s_0, s_1, \dots s_n\}$
  - Instead of emitting a single distance, emit an array of distances, with respect to each source
  - Reducer selects minimum for each element in array

- Does this scale?

# Application: Social Search

# Social Search

- When searching, how to rank friends named "John"?

  - Assume undirected graphs
  - Rank matches by distance to user

- Naïve implementations:

  - Precompute all-pairs distances
  - Compute distances at query time

- Can we do better?

# Landmark Approach (aka sketches)

○ Select $n$ seeds $\{s_0, s_1, \ldots s_n\}$

○ Compute distances from seeds to every node:

$$A = [2, 1, 1]$$
$$B = [1, 1, 2]$$
$$C = [4, 3, 1]$$
$$D = [1, 2, 4]$$

- What can we conclude about distances?
- Insight: landmarks bound the maximum path length

○ Lots of details:

- How to more tightly bound distances
- How to select landmarks (random isn't the best…)

○ Use multi-source parallel BFS implementation in MapReduce!
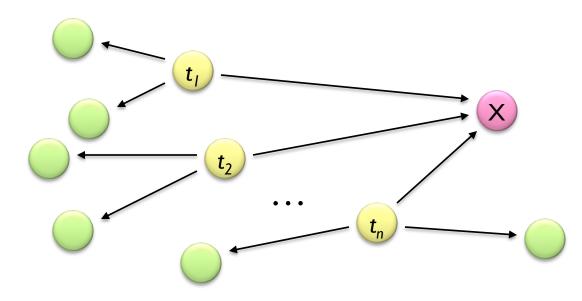
**&lt;pause/&gt;**

# Graphs and MapReduce

○ A large class of graph algorithms involve:

- Performing computations at each node: based on node features, edge features, and local link structure
- Propagating computations: "traversing" the graph

○ Generic recipe:

- Represent graphs as adjacency lists
- Perform local computations in mapper
- Pass along partial results via outlinks, keyed by destination node
- Perform aggregation in reducer on inlinks to a node
- Iterate until convergence: controlled by external "driver"
- Don't forget to pass the graph structure between iterations

# PageRank

Given page *x* with inlinks $t_1 \ldots t_n$, where

- *C(t)* is the out-degree of *t*
- $\alpha$ is probability of random jump
- *N* is the total number of nodes in the graph

$$PR(x) = \alpha \left( \frac{1}{N} \right) + (1 - \alpha) \sum_{i=1}^{n} \frac{PR(t_i)}{C(t_i)}$$
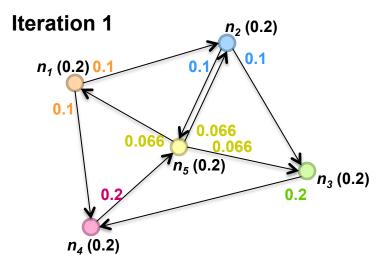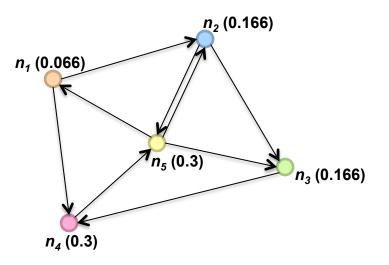
# Computing PageRank

- Properties of PageRank

  - Can be computed iteratively

  - Effects at each iteration are local

- Sketch of algorithm:

  - Start with seed $PR_i$ values

  - Each page distributes $PR_i$ "credit" to all pages it links to

  - Each target page adds up "credit" from multiple in-bound links to compute $PR_{i+1}$

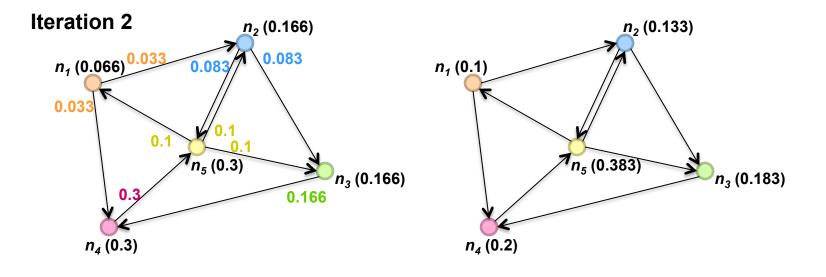  - Iterate until values converge

# Simplified PageRank

- First, tackle the simple case:
  - No random jump factor
  - No dangling nodes

- Then, factor in these complexities…
  - Why do we need the random jump?
  - Where do dangling nodes come from?

# Sample PageRank Iteration (1)
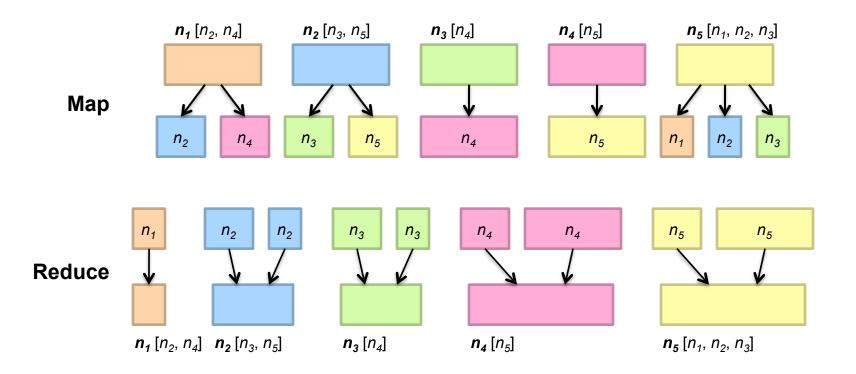


Iteration 1

$n_2$ (0.2)

$n_1$ (0.2) 0.1    0.1    0.1

0.1

0.066
0.066   0.066
$n_5$ (0.2)

$n_3$ (0.2)

0.2

0.2

$n_4$ (0.2)

$n_2$ (0.166)

$n_1$ (0.066)

$n_5$ (0.3)

$n_3$ (0.166)

$n_4$ (0.3)

# Sample PageRank Iteration (2)

# PageRank in MapReduce

**Map**

$n_1$ [$n_2$, $n_4$]  $n_2$ [$n_3$, $n_5$]  $n_3$ [$n_4$]  $n_4$ [$n_5$]  $n_5$ [$n_1$, $n_2$, $n_3$]

$n_2$  $n_4$  $n_3$  $n_5$  $n_4$  $n_5$  $n_1$  $n_2$  $n_3$

**Reduce**

$n_1$  $n_2$  $n_2$  $n_3$  $n_3$  $n_4$  $n_4$  $n_5$  $n_5$

$n_1$ [$n_2$, $n_4$]  $n_2$ [$n_3$, $n_5$]  $n_3$ [$n_4$]  $n_4$ [$n_5$]  $n_5$ [$n_1$, $n_2$, $n_3$]

# PageRank Pseudo-Code

```
1: class MAPPER
2:     method MAP(nid n, node N)
3:         p ← N.PAGERANK/|N.ADJACENCYLIST|
4:         EMIT(nid n, N)                                    ▷ Pass along graph structure
5:         for all nodeid m ∈ N.ADJACENCYLIST do
6:             EMIT(nid m, p)                                ▷ Pass PageRank mass to neighbors

1: class REDUCER
2:     method REDUCE(nid m, [p₁, p₂, ...])
3:         M ← ∅
4:         for all p ∈ counts [p₁, p₂, ...] do
5:             if ISNODE(p) then
6:                 M ← p                                      ▷ Recover graph structure
7:             else
8:                 s ← s + p                                  ▷ Sums incoming PageRank contributions
9:         M.PAGERANK ← s
10:        EMIT(nid m, node M)
```

# Complete PageRank

- Two additional complexities

  - What is the proper treatment of dangling nodes?
  - How do we factor in the random jump factor?

- Solution:

  - Second pass to redistribute "missing PageRank mass" and account for random jumps

$$p' = \alpha \left( \frac{1}{N} \right) + (1 - \alpha) \left( \frac{m}{N} + p \right)$$

  - $p$ is PageRank value from before, $p'$ is updated PageRank value
  - $N$ is the number of nodes in the graph
  - $m$ is the missing PageRank mass

- Additional optimization: make it a single pass!

# PageRank Convergence

- Alternative convergence criteria

  - Iterate until PageRank values don't change

  - Iterate until PageRank rankings don't change

  - Fixed number of iterations

- Convergence for web graphs?

  - Not a straightforward question

- Watch out for link spam:

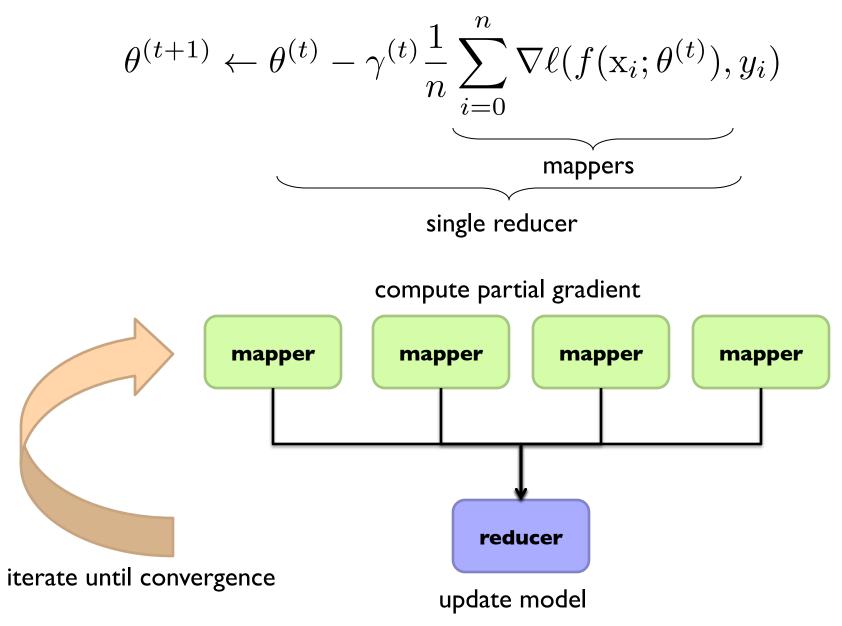  - Link farms

  - Spider traps

  - …

# Beyond PageRank

- Variations of PageRank

  - Weighted edges
  - Personalized PageRank

- Variants on graph random walks

  - Hubs and authorities (HITS)
  - SALSA

# Other Classes of Graph Algorithms

○ Subgraph pattern matching

○ Computing simple graph statistics

  ● Degree vertex distributions

○ Computing more complex graph statistics

  ● Clustering coefficients
  ● Counting triangles
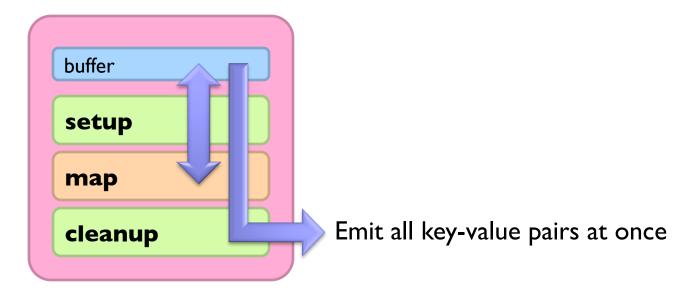
# Batch Gradient Descent in MapReduce

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \frac{1}{n} \sum_{i=0}^{n} \nabla \ell(f(\mathrm{x}_i; \theta^{(t)}), y_i)$$

mappers

single reducer

compute partial gradient



iterate until convergence

update model

# MapReduce sucks at iterative algorithms

- Hadoop task startup time

- Stragglers

- Needless graph shuffling

- Checkpointing at each iteration

# In-Mapper Combining

- Use combiners

  - Perform local aggregation on map output
  - Downside: intermediate data is still materialized

- Better: in-mapper combining

  - Preserve state across multiple map calls, aggregate messages in buffer, emit buffer contents at end
  - Downside: requires memory management

buffer

**setup**

**map**

**cleanup**
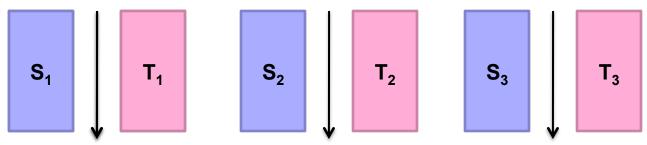
Emit all key-value pairs at once

# Better Partitioning

- Default: hash partitioning
  - Randomly assign nodes to partitions

- Observation: many graphs exhibit local structure
  - E.g., communities in social networks
  - Better partitioning creates more opportunities for local aggregation

- Unfortunately, partitioning is **hard**!
  - Sometimes, chick-and-egg…
  - But cheap heuristics sometimes available
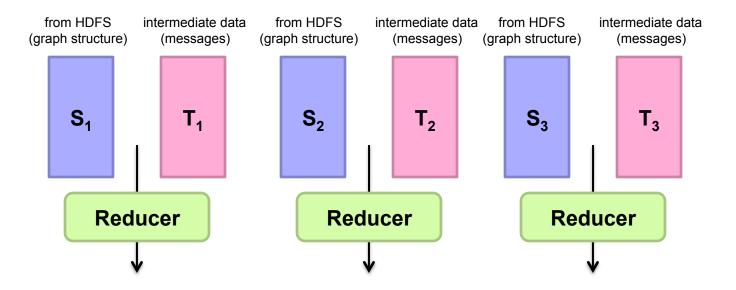  - For webgraphs: range partition on domain-sorted URLs

# Schimmy Design Pattern

○ Basic implementation contains two dataflows:

- Messages (actual computations)
- Graph structure ("bookkeeping")

○ Schimmy: separate the two dataflows, shuffle only the messages

- Basic idea: merge join between graph structure and messages

both relations consistently partitioned and sorted by join key

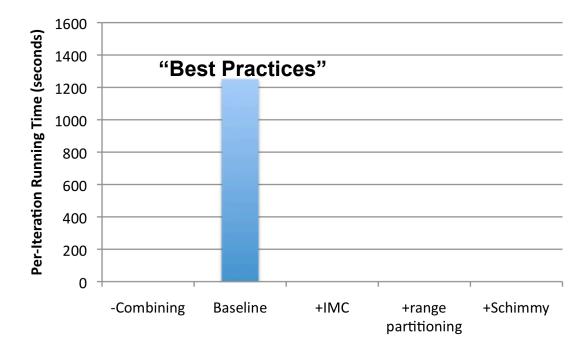| $S_1$ | | $T_1$ | | $S_2$ | | $T_2$ | | $S_3$ | | $T_3$ |

# Do the Schimmy!

- Schimmy = reduce side parallel merge join between graph structure and messages

  - Consistent partitioning between input and intermediate data
  - Mappers emit only messages (actual computation)
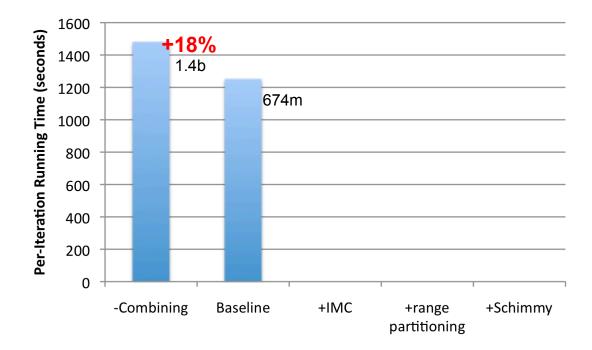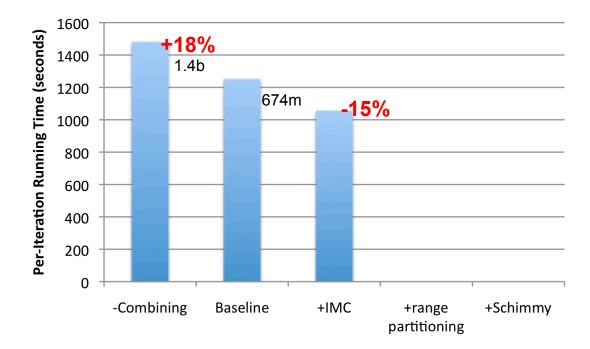  - Reducers read graph structure directly from HDFS

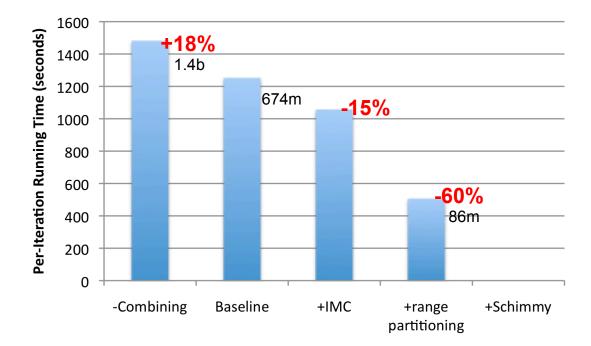| from HDFS (graph structure) | intermediate data (messages) | from HDFS (graph structure) | intermediate data (messages) | from HDFS (graph structure) | intermediate data (messages) |
|---|---|---|---|---|---|
| $S_1$ | $T_1$ | $S_2$ | $T_2$ | $S_3$ | $T_3$ |

Reducer   Reducer   Reducer
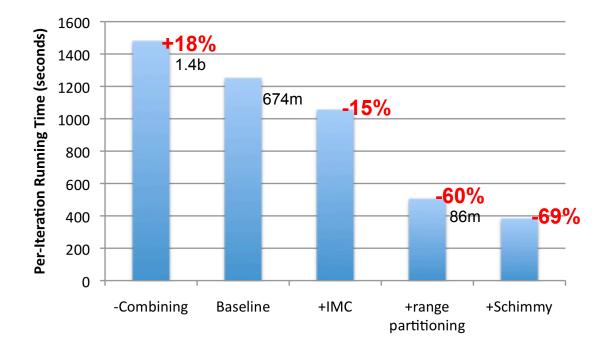
# Experiments

- Cluster setup:

  - 10 workers, each 2 cores (3.2 GHz Xeon), 4GB RAM, 367 GB disk
  - Hadoop 0.20.0 on RHELS 5.3

- Dataset:

  - First English segment of ClueWeb09 collection
  - 50.2m web pages (1.53 TB uncompressed, 247 GB compressed)
  - Extracted webgraph: 1.4 billion edges, 7.0 GB
  - Dataset arranged in crawl order

- Setup:

  - Measured per-iteration running time (5 iterations)
  - 100 partitions

# Results

# Results

# Results

# Results

# Results

# Sequencing Computations

# Sequencing Computations

1. Turn synchronization into a sorting problem
   - Leverage the fact that keys arrive at reducers in sorted order
   - Manipulate the sort order and partitioning scheme to deliver partial results at appropriate junctures

2. Create appropriate algebraic structures to capture computation
   - Build custom data structures to accumulate partial results

Monoids!

# Monoids!

- What's a monoid?

- An algebraic structure with

  - A single associative binary operation
  - An identity

- Examples:

  - Natural numbers form a commutative monoid under + with identity 0
  - Natural numbers form a commutative monoid under × with identity 1
  - Finite strings form a monoid under concatenation with identity ""
  - …

# Monoids and MapReduce

○ Recall averaging example: why does it work?

- AVG is non-associative
- Tuple of (sum, count) forms a monoid under element-wise addition
- Destroy the monoid at end to compute average
- Also explains the various failed algorithms

○ "Stripes" pattern works in the same way!

- Associate arrays form a monoid under element-wise addition

*Go forth and monoidify!*

# Abstract Algebra and MapReduce

- Create appropriate algebraic structures to capture computation

- Algebraic properties

  - Associative: order doesn't matter!
  - Commutative: grouping doesn't matter!
  - Idempotent: duplicates don't matter!
  - Identity: this value doesn't matter!
  - Zero: other values don't matter!
  - …

- Different combinations lead to monoids, groups, rings, lattices, etc.

Recent thoughts, see: Jimmy Lin. Monoidify! Monoids as a Design Principle for Efficient MapReduce Algorithms. arXiv:1304.7544, April 2013.

Questions?