

Case Study 4: Collaborative Filtering

Graph-Parallel Problems

Synchronous v. Asynchronous Computation

Machine Learning for Big Data
 CSE547/STAT548, University of Washington
 Carlos Guestrin, guest lecturer
 May 14th, 2015

©Emily Fox 2015

1

ML in the Context of Parallel Architectures



- But scalable ML in these systems is hard, especially in terms of:
 1. Programmability
 2. Data distribution
 3. Failures

we'll go through
these ideas...
explore with MapReduce

©Emily Fox 2015

2

Move Towards Higher-Level Abstraction

- Distributed computing challenges are hard and annoying!
 - Programmability
 - Data distribution
 - Failures
 - High-level abstractions try to simplify distributed programming by hiding challenges:
 - Provide different levels of robustness to failures, optimizing data movement and communication, protect against race conditions...
 - ~~Generally, you are still on your own WRT designing parallel algorithms~~
 - Some common parallel abstractions:
 - Lower-level:
 - Pthreads: abstraction for distributed threads on single machine
 - MPI: abstraction for distributed communication in a cluster of computers
 - Higher-level:
 - Map-Reduce (Hadoop: open-source version): mostly data-parallel problems
 - GraphLab: for graph-structured distributed problems
- this quarter*

©Emily Fox 2015

3

Simplest Type of Parallelism: Data Parallel Problems

- You have already learned a classifier
 - What's the test error?
 - You have 10B labeled documents and 1000 machines
-
- $\text{err}_{\text{test}} = \frac{\sum_{i=1}^{1000} |y_i - \text{sign}(w^* \cdot x_i)|}{N_{\text{test}}}$
- $\text{error on local data}$
- $\text{err}_{\text{test}} \leftarrow \text{avg of } e_1 \dots e_{1000}$

- Problems that can be broken into independent subproblems are called data-parallel (or embarrassingly parallel)
- Map-Reduce is a great tool for this...
 - Focus of today's lecture
 - but first a simple example

©Emily Fox 2015

4

Data Parallelism (MapReduce)



*Solve a huge number of **independent** subproblems,
e.g., extract features in images*

Map-Reduce Abstraction

- Map: *Transform a data element*
 - Data-parallel over elements, e.g., documents
 - Generate (key,value) pairs
 - "value" can be any data type

In our example: $(\text{UW}, 17)$
 $(\text{UW}, 1)$
 $(\text{Mary}, 1)$
 $(\text{UW}, 1)$

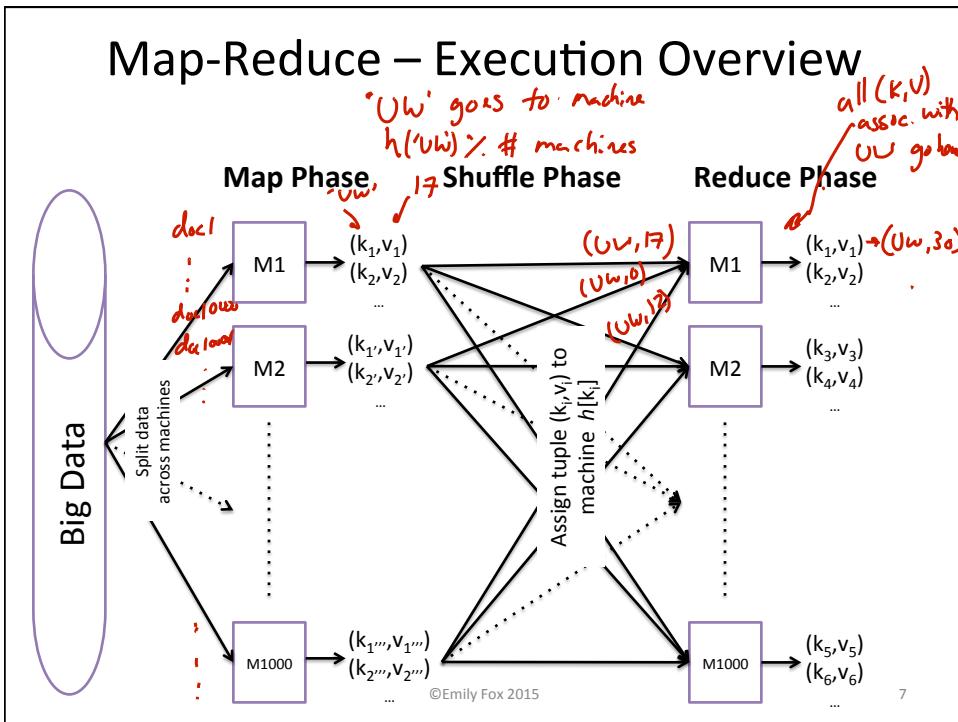
- Reduce: *Take all values associated with a key and aggregate*
 - Aggregate values for each key
 - Must be commutative-associate operation
 - Data-parallel over keys
 - Generate (key,value) pairs

reduce ('UW', [1, 17, 0, 0, 1])
emit ('UW', 30)
Key *total count*

- Map-Reduce has long history in functional programming
 - But popularized by Google, and subsequently by open-source Hadoop implementation from Yahoo!

Example:
word count
map (document)
for word in doc
emit (word, 1)
key *value*

Reduce (word, count: list/int)
C = 0 *key* *list of counts*
for i in count
c += count[i]
emit (word, c)
total



Issues with Map-Reduce Abstraction

- Often all data gets moved around cluster
 - Very bad for iterative settings

- Definition of Map & Reduce functions can be unintuitive in many apps
 - Graphs are challenging

- Computation is synchronous
 - all mappers are done before can complete reduce Phase*

SGD for Matrix Factorization in Map-Reduce?

$$\begin{matrix} \text{user } u \\ \text{item } v \end{matrix} \rightarrow \begin{bmatrix} L_u^{(t+1)} \\ R_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u) L_u^{(t)} - \eta_t \epsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v) R_v^{(t)} - \eta_t \epsilon_t L_u^{(t)} \end{bmatrix}$$

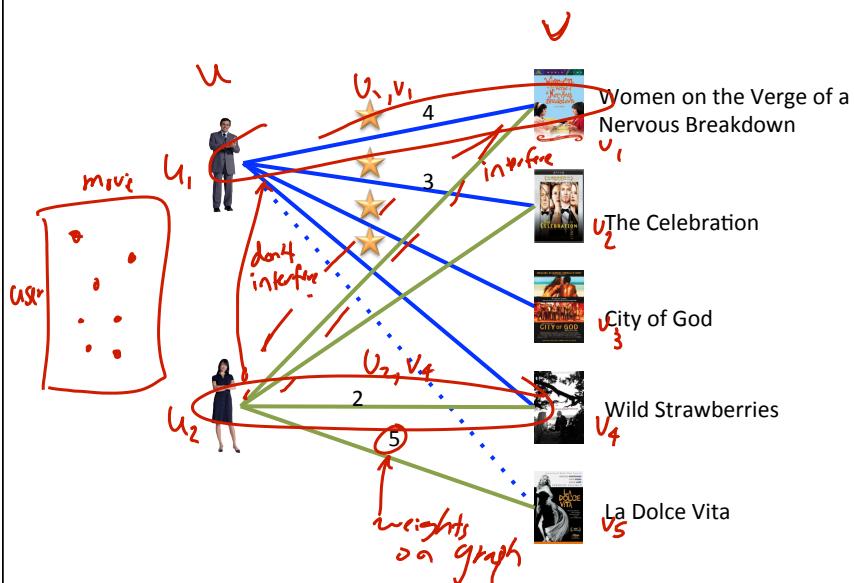
$\epsilon_t = L_u^{(t)} \cdot R_v^{(t)} - r_{uv}$

- Map and Reduce functions???
 - Map-Reduce:
 - Data-parallel over all mappers
 - Data-parallel over reducers with same key
 - Here, one update at a time!
- don't fit into
data parallel framework

©Emily Fox 2015

9

Matrix Factorization as a Graph



©Emily Fox 2015

10

Flashback to 1998

Why?

First Google advantage:
a Graph Algorithm & a System to Support it!

Page Rank *at Scale*

©Emily Fox 2015 11

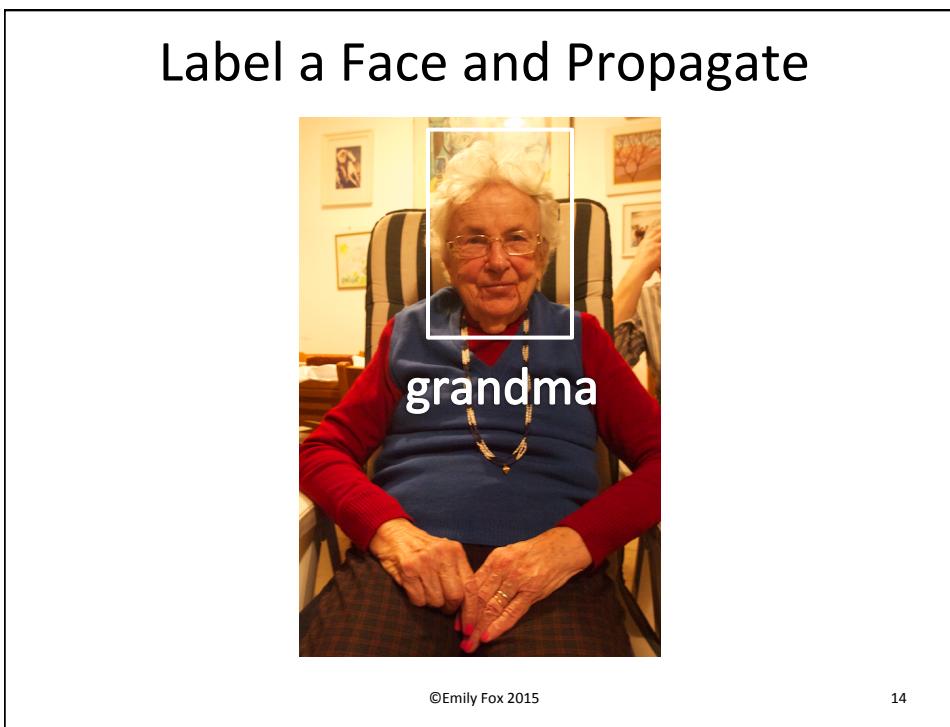
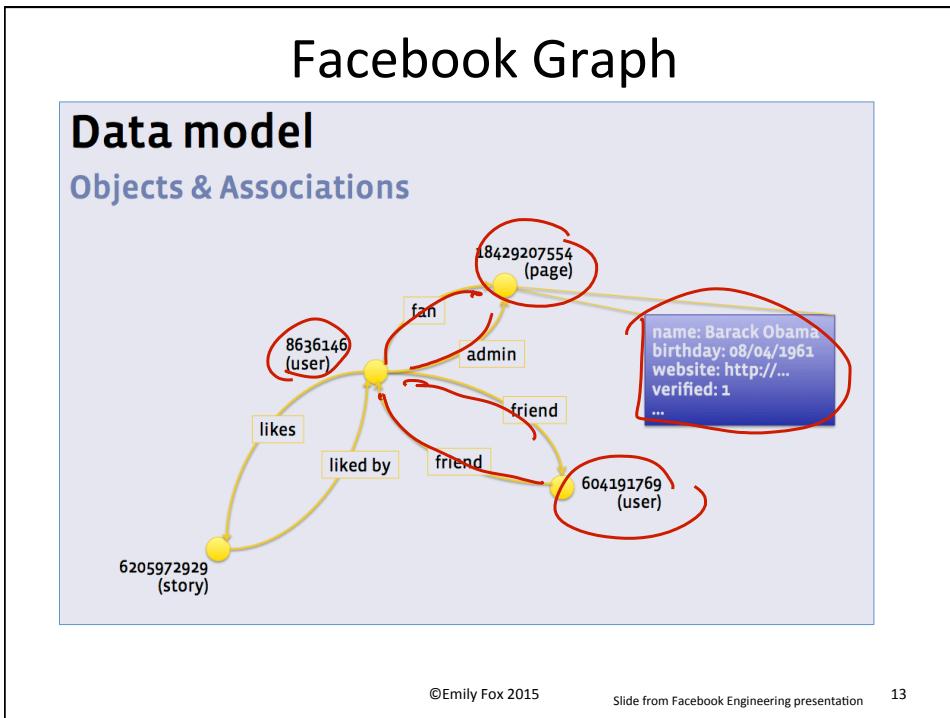
Social Media	Science	Advertising	Web

- **Graphs** encode the **relationships** between:

People	Products	Ideas
Facts	Interests	

- **Big: 100 billions of vertices and edges and rich metadata**
 - Facebook (10/2012): 1B users, 144B friendships
 - Twitter (2011): 15B follower edges

©Emily Fox 2015 12



Pairwise similarity not enough...



©Emily Fox 2015

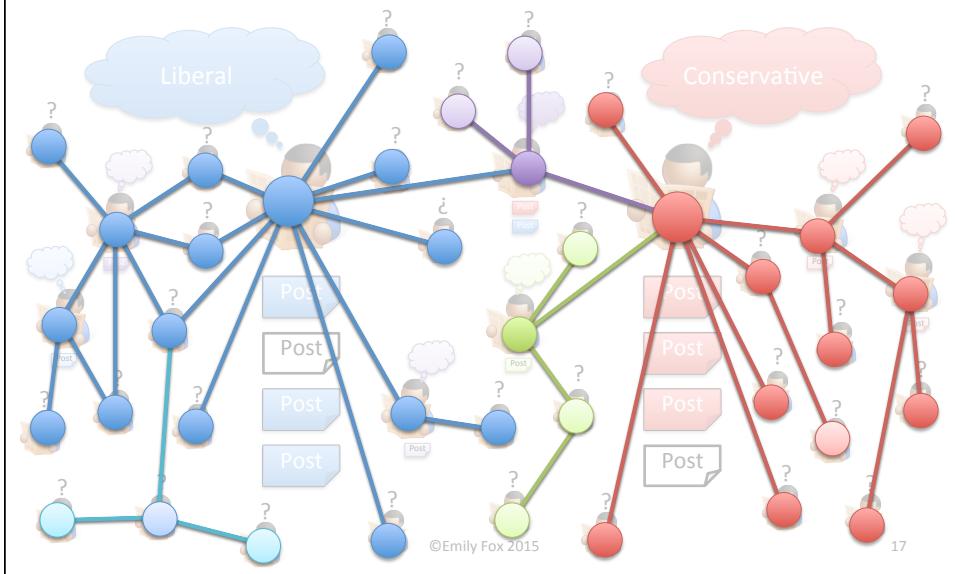
15

Propagate Similarities & Co-occurrences for Accurate Predictions

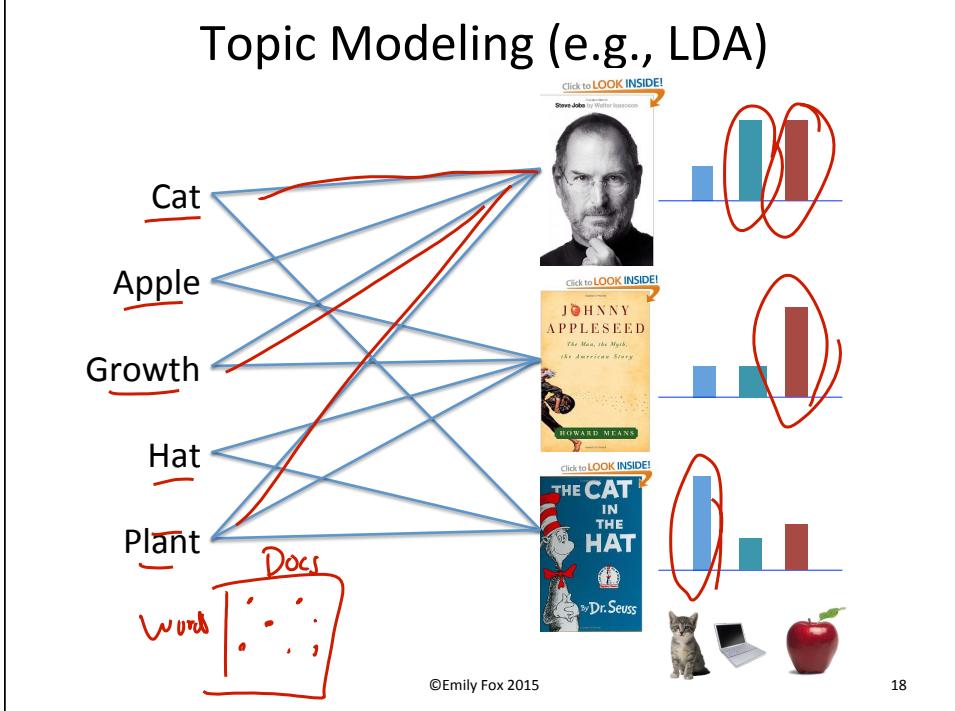


16

Example: Estimate Political Bias



Topic Modeling (e.g., LDA)



ML Tasks Beyond Data-Parallelism



Map Reduce

Feature Extraction
Cross Validation
Computing Sufficient Statistics

Graphical Models
Gibbs Sampling
Belief Propagation
Variational Opt.
Collaborative Filtering
Tensor Factorization

Semi-Supervised Learning
Label Propagation
CoEM
Graph Analysis
PageRank
Triangle Counting

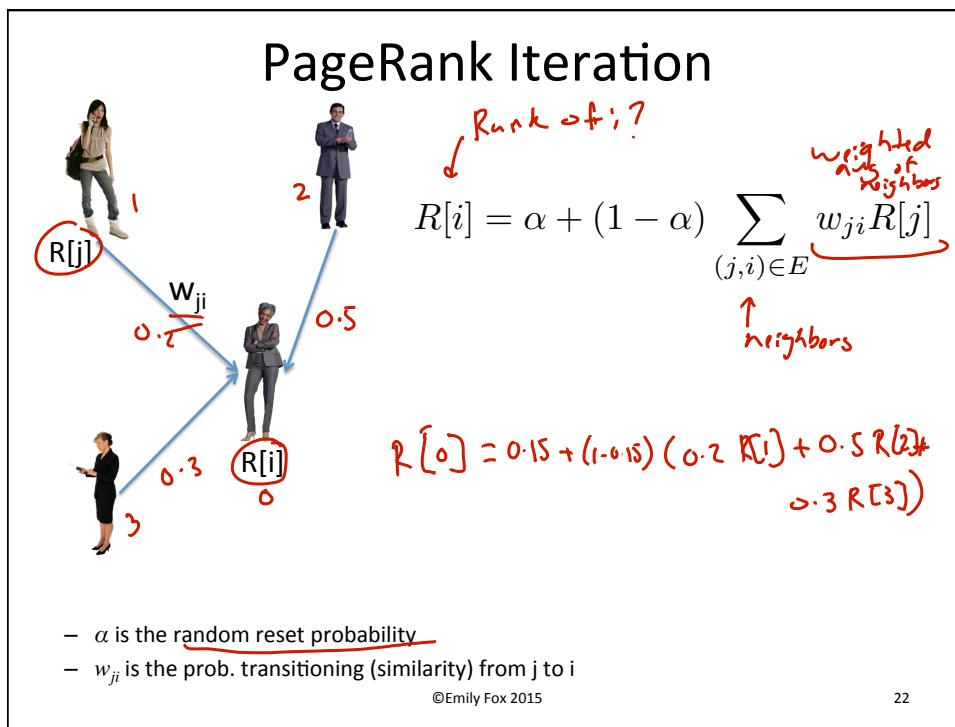
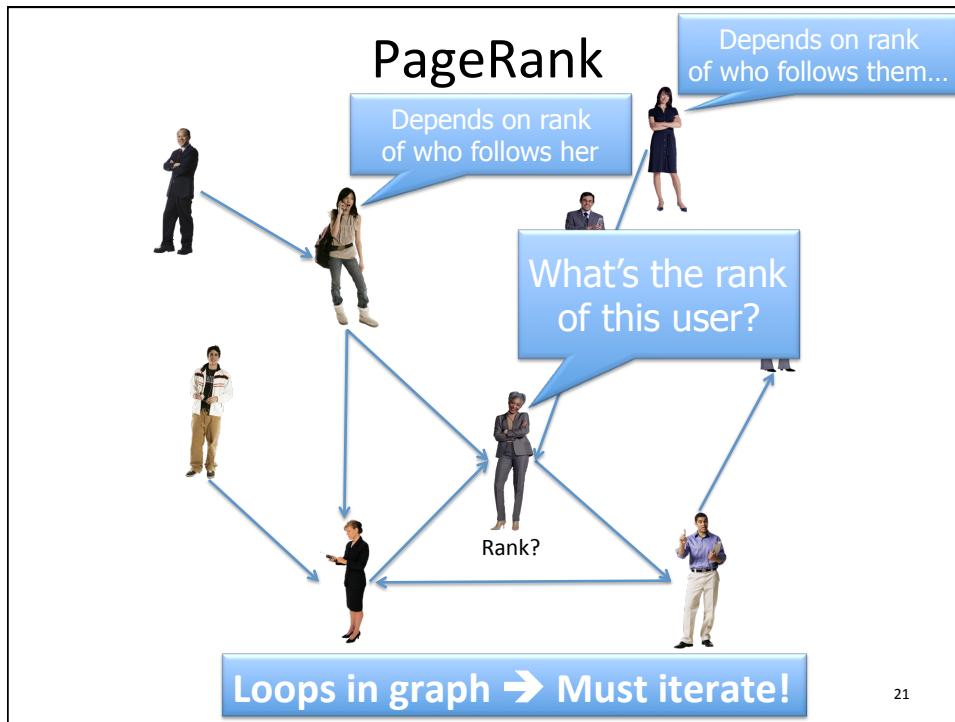
©Emily Fox 2015

19

Example of a Graph-Parallel Algorithm

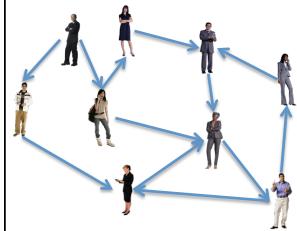
©Emily Fox 2015

20

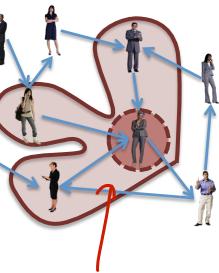


Properties of Graph Parallel Algorithms

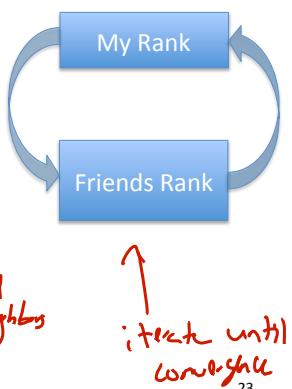
Dependency Graph



Local Updates



Iterative Computation



©Emily Fox 2015

23

Addressing Graph-Parallel ML



Map Reduce

Feature Extraction
Cross Validation
Computing Sufficient Statistics

Graph-Parallel

Graph-Parallel Abstraction

Graphical Models
Gibbs Sampling
Belief Propagation
Variational Opt.
Collaborative Filtering
Tensor Factorization

Semi-Supervised Learning
Label Propagation
CoEM
Data-Mining
PageRank
Triangle Counting

©Emily Fox 2015

24

Key Question.

Graph Computation:

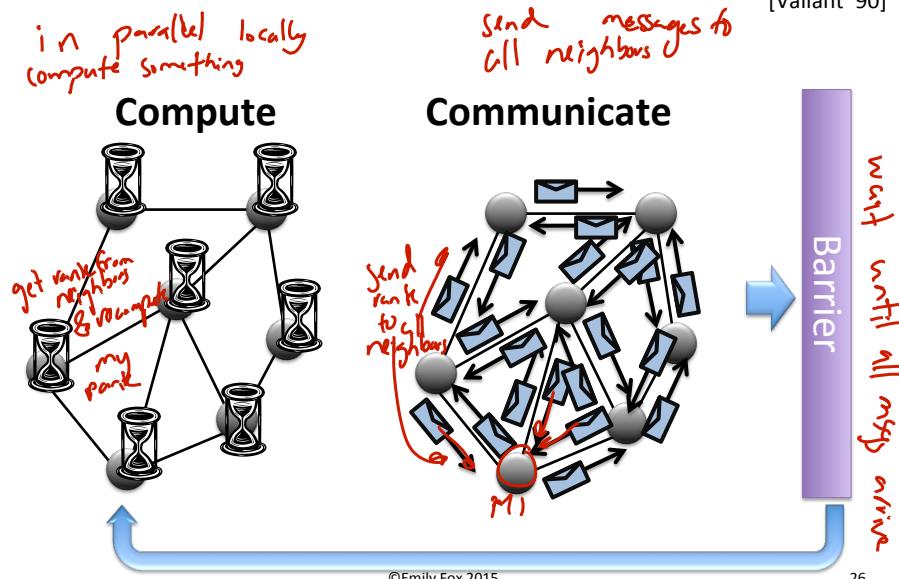
Synchronous

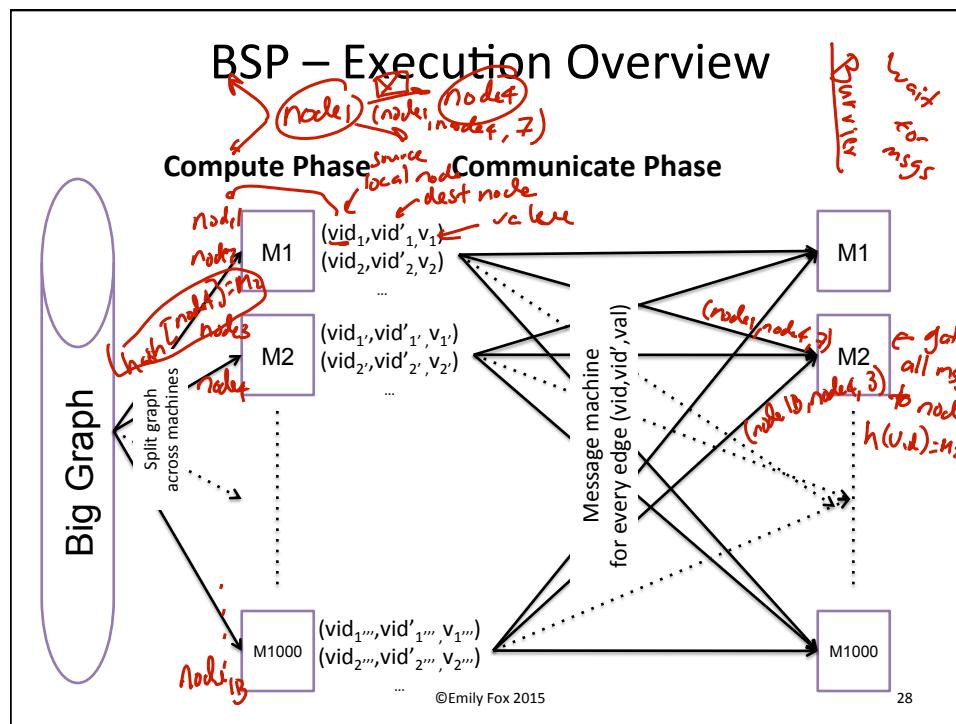
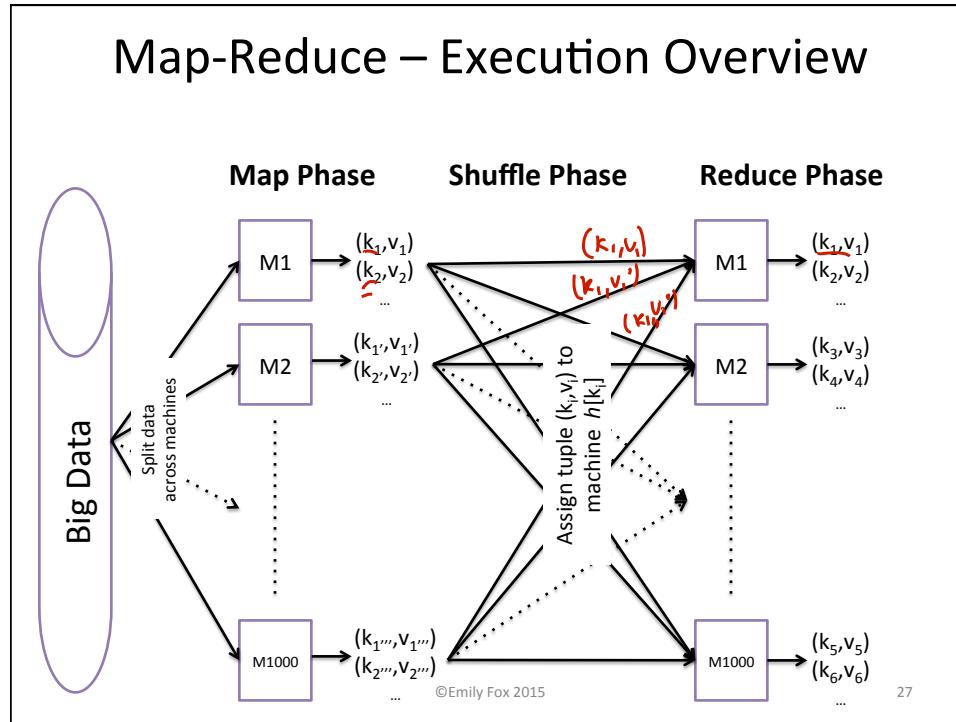
v.

Asynchronous

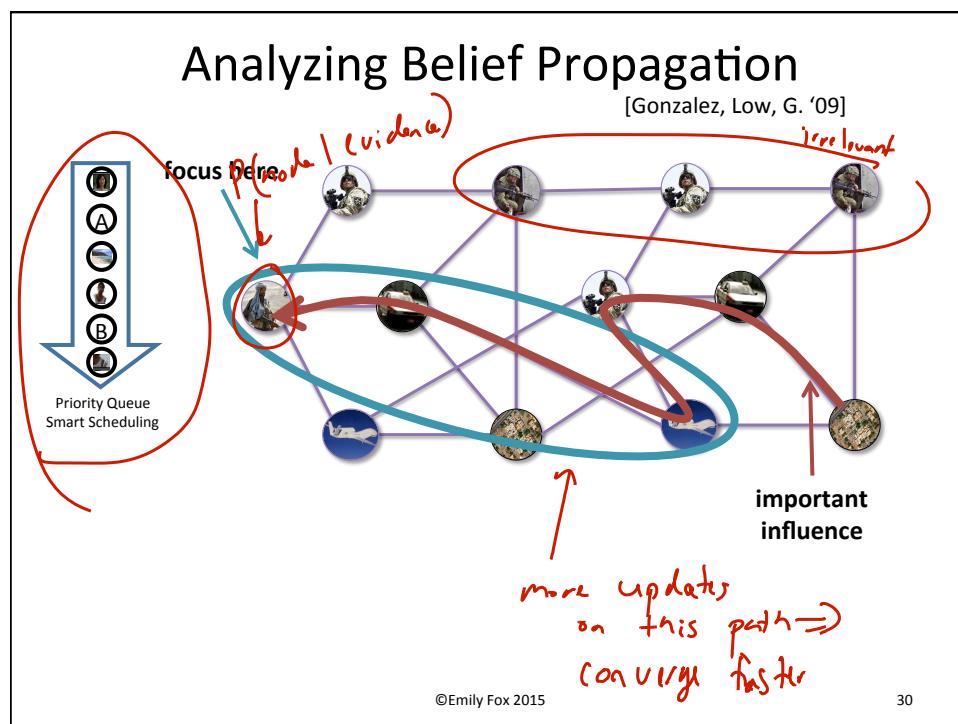
Bulk Synchronous Parallel Model: Pregel (Giraph)

[Valiant '90]





*Bulk synchronous
parallel model
provably inefficient
for some ML tasks*



Asynchronous Belief Propagation

Challenge = Boundaries

Synthetic Noisy Image

Cumulative Vertex Updates

Many Updates

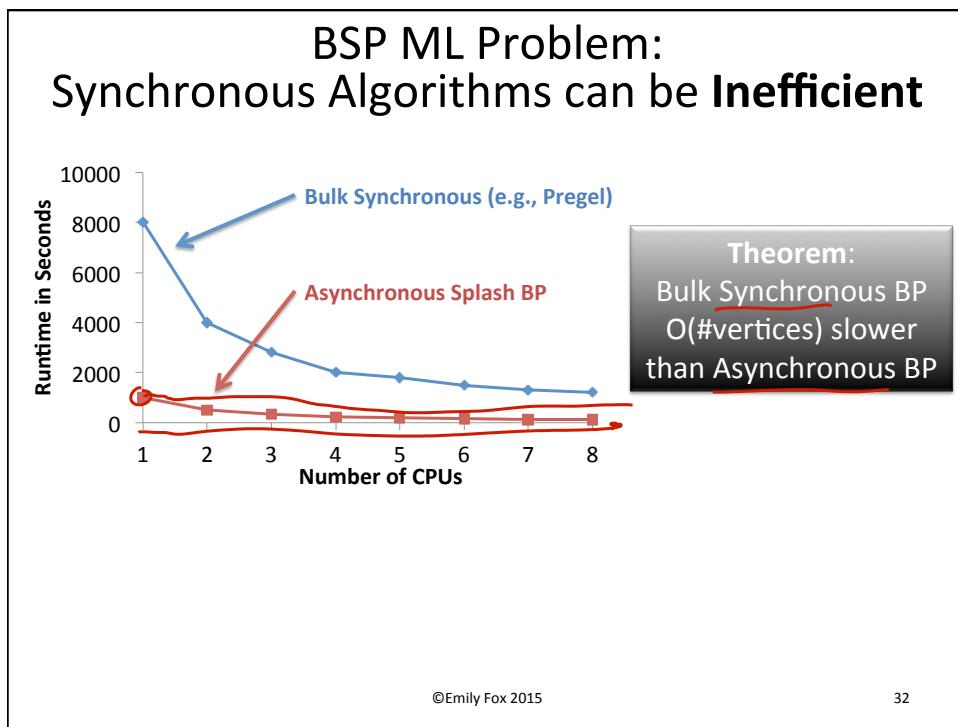
Few Updates

Graphical Model

Algorithm identifies and focuses on hidden sequential structure

©Emily Fox 2015

31



Synchronous v. Asynchronous

- Bulk synchronous processing:
 - Computation in phases
 - All vertices participate in a phase
 - Though OK to say no-op
 - All messages are sent
 - Simpler to build, like Map-Reduce
 - No worries about race conditions, barrier guarantees data consistency
 - Simpler to make fault-tolerant, save data on barrier
 - Slower convergence for many ML problems
 - In matrix-land, called Jacobi Iteration
 - Implemented by Google Pregel 2010
- Asynchronous processing:
 - Vertices see latest information from neighbors
 - Most closely related to sequential execution
 - Harder to build:
 - Race conditions can happen all the time
 - Must protect against this issue
 - More complex fault tolerance
 - When are you done?
 - Must implement scheduler over vertices
 - Faster convergence for many ML problems
 - In matrix-land, called Gauss-Seidel Iteration
 - Implemented by GraphLab 2010, 2012

©Emily Fox 2015

33

Case Study 4: Collaborative Filtering

GraphLab

Machine Learning for Big Data
 CSE547/STAT548, University of Washington

Carlos Guestrin, guest lecturer

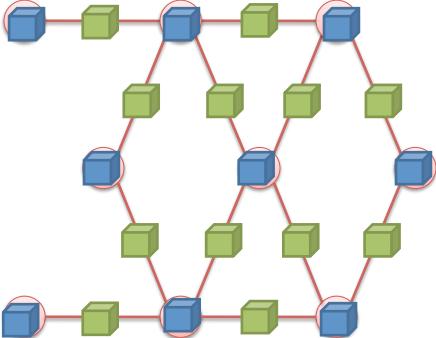
May 14th, 2015

©Emily Fox 2015

34

Data Graph

Data associated with vertices and edges



Graph:
• Social Network

Vertex Data:
• User profile text
• Current interests estimates

Edge Data:
• Similarity weights

©Emily Fox 2015

35

How do we *program*
graph computation?

“Think like a Vertex.”

-Malewicz et al. [SIGMOD’10]

Update Functions

User-defined program: applied to **vertex** transforms data in scope of vertex

`pagerank(i, scope){`

$$R[i] = \alpha + (1-\alpha) \sum_{j \in \text{scope.neighbors}} \frac{\text{Rank}[j]}{|\text{scope.neighbors}|}$$

©Emily Fox 2015 37

Update Function Example: Connected Components

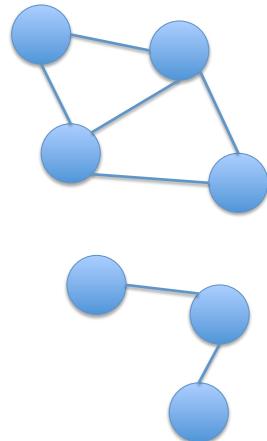
1. initialize: `vertex_component = vertex_id`

2. iterate
pick a node i :
 $\text{component}(i) = \min(\text{self}, \text{neighbors})$
 if change, add neighbors to queue

©Emily Fox 2015 38

Update Function Example: Connected Components

init: $\text{Component}[i] = i$



update (i , scope):

$\text{Component}[i] = \min ($

$\text{Component}[i],$

$\min_{j \in \text{neighbors}} (\text{Component}[j]))$

; if $\text{Component}[i]$ changes
{ schedule neighbors of i }

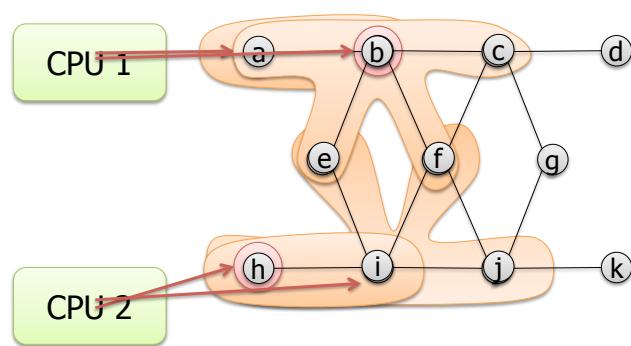
©Emily Fox 2015

39

The Scheduler

The **scheduler** determines order vertices are updated

Scheduler



©Emily Fox 2015

40

Example Schedulers

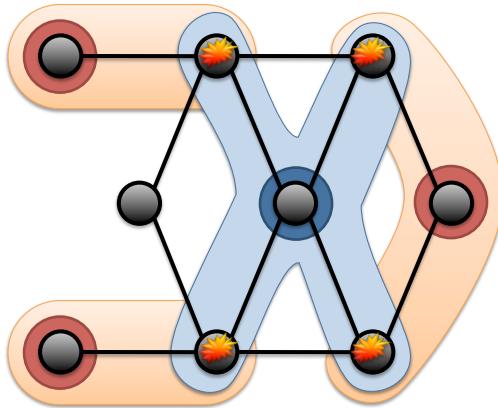
- Round-robin
- Selective scheduling (skipping):
 - round robin but jump over un-scheduled vertex
- FIFO
- Prioritize scheduling
 - Hard to implement in a distributed fashion
 - Approximations used (each machine has its own priority queue)

©Emily Fox 2015

41

Ensuring Race-Free Code

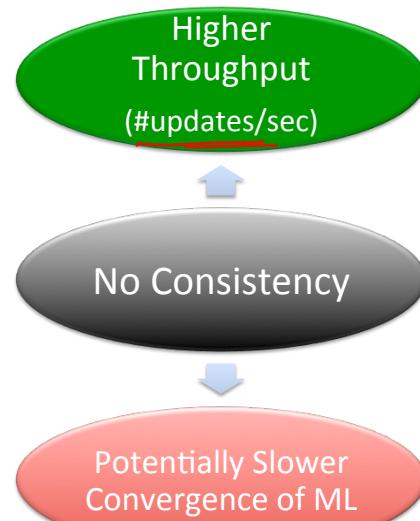
How much can computation **overlap**?



©Emily Fox 2015

42

Need for Consistency?

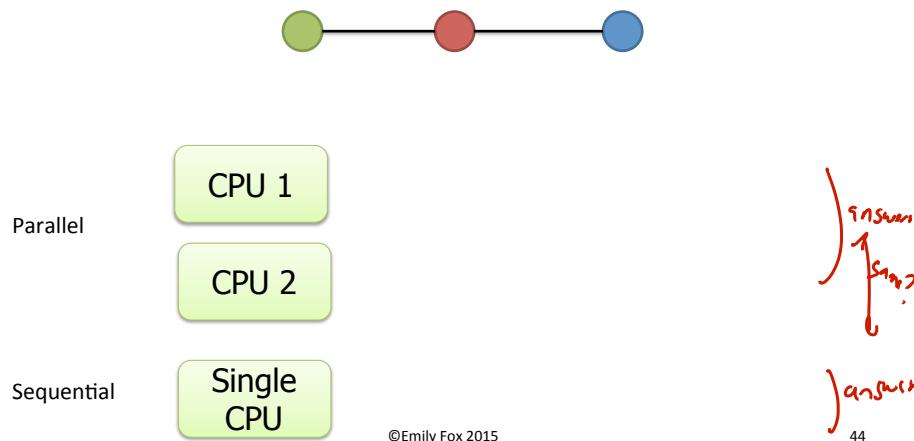


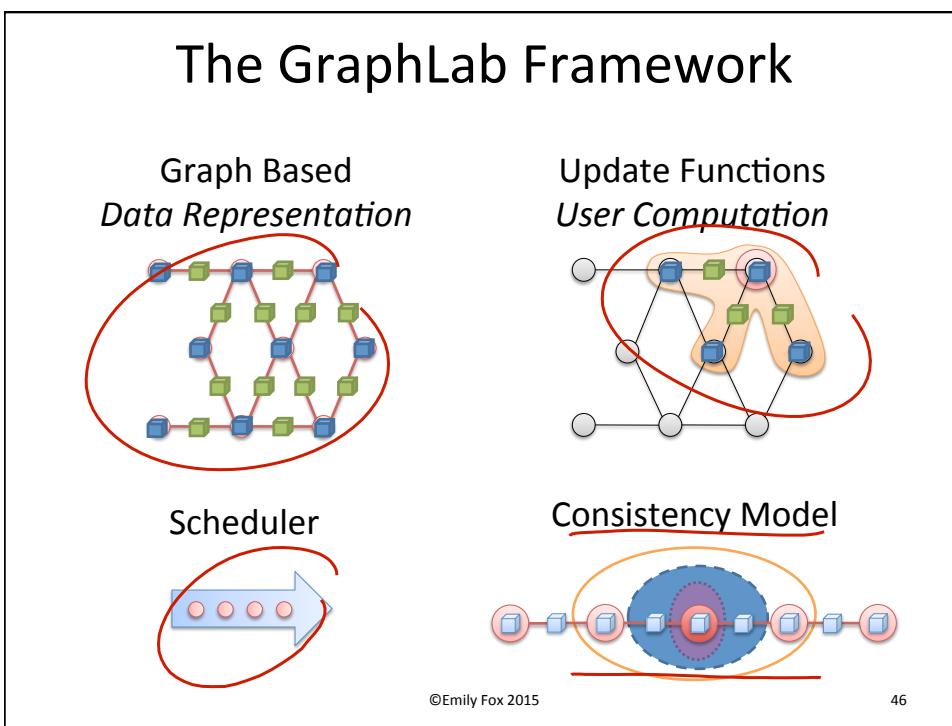
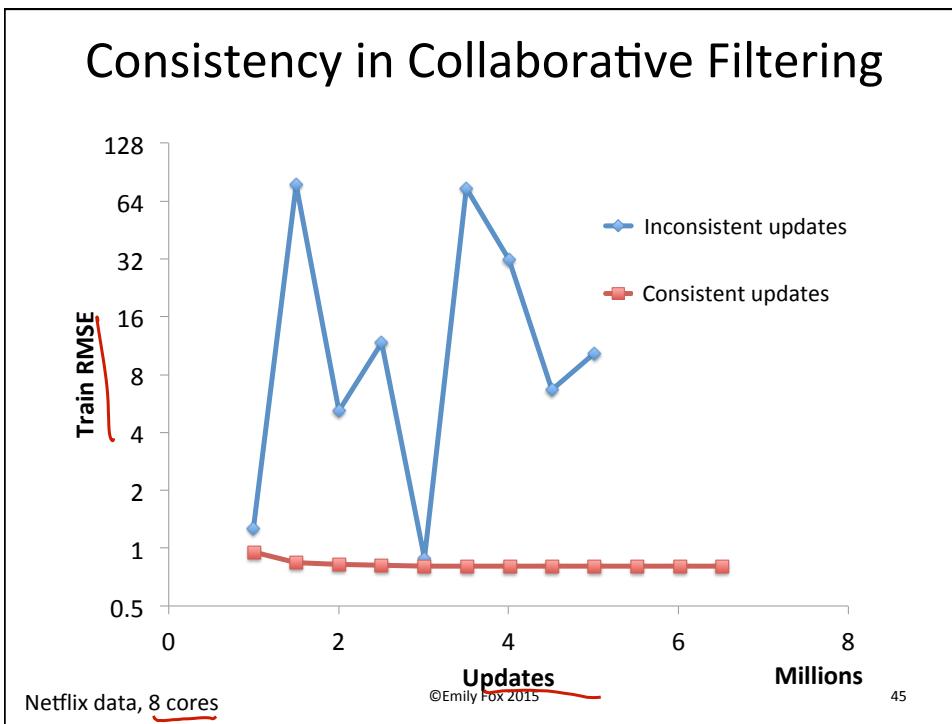
©Emily Fox 2015

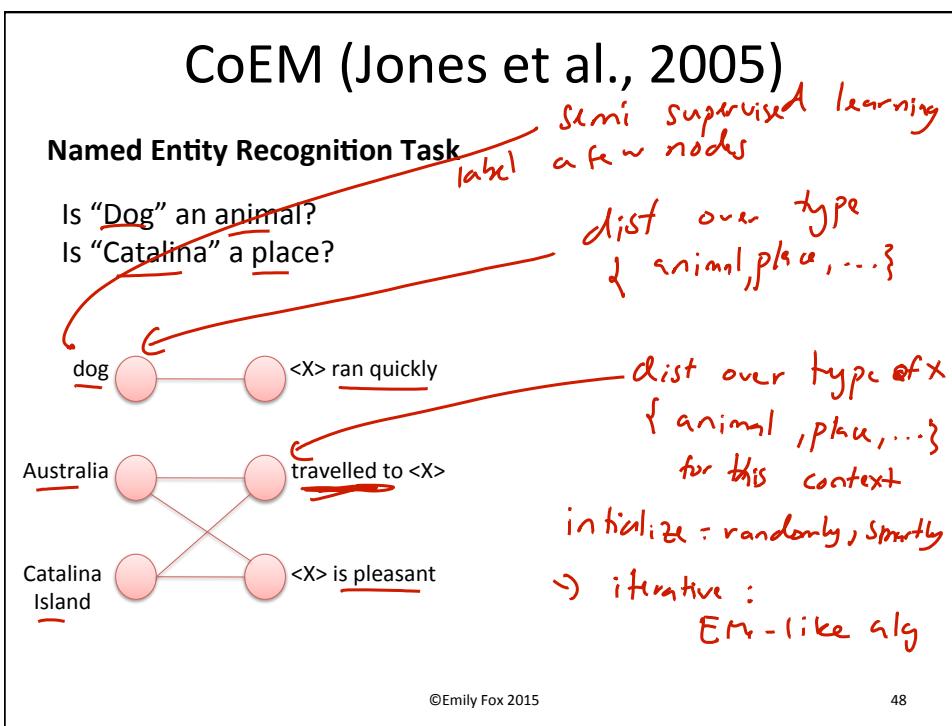
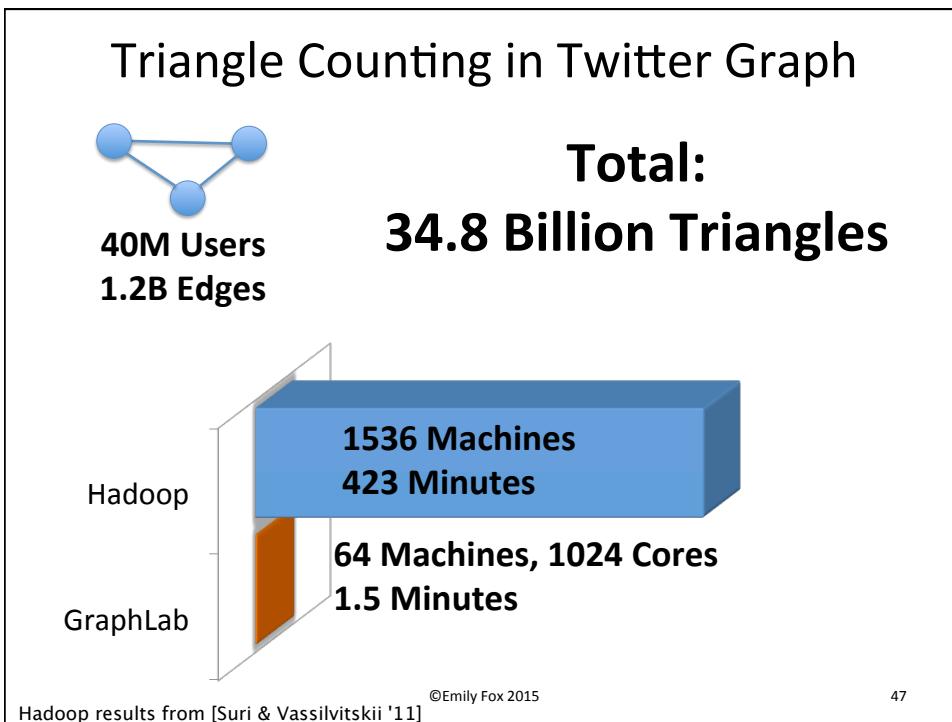
43

GraphLab Ensures Sequential Consistency

For each parallel execution, there exists a sequential execution of update functions which produces the same result







Never Ending Learner Project (CoEM)

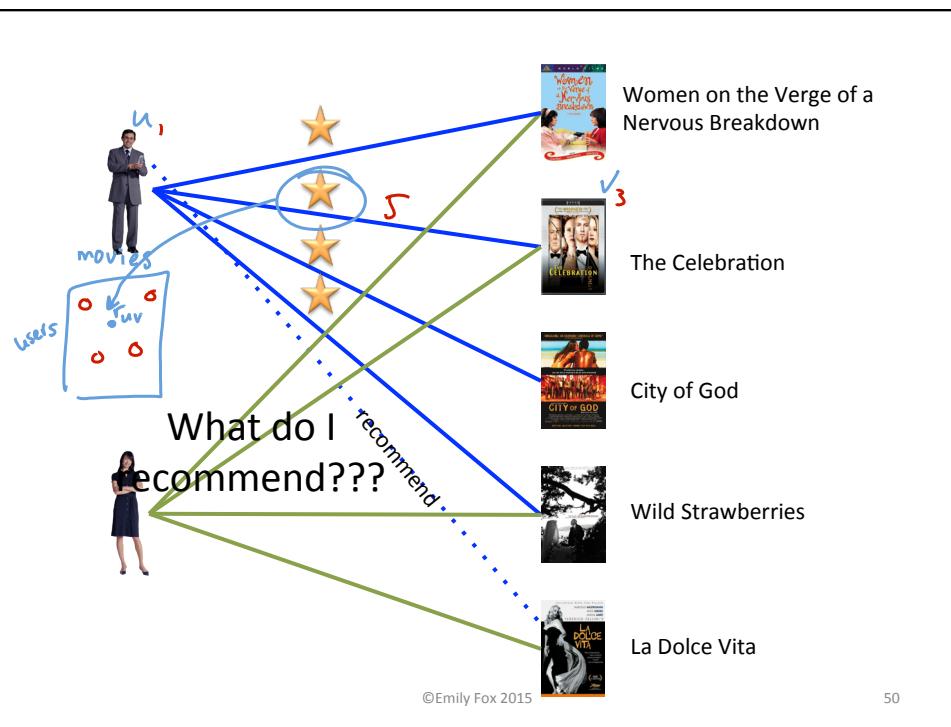
Vertices: 2 Million

Edges: 200 Million

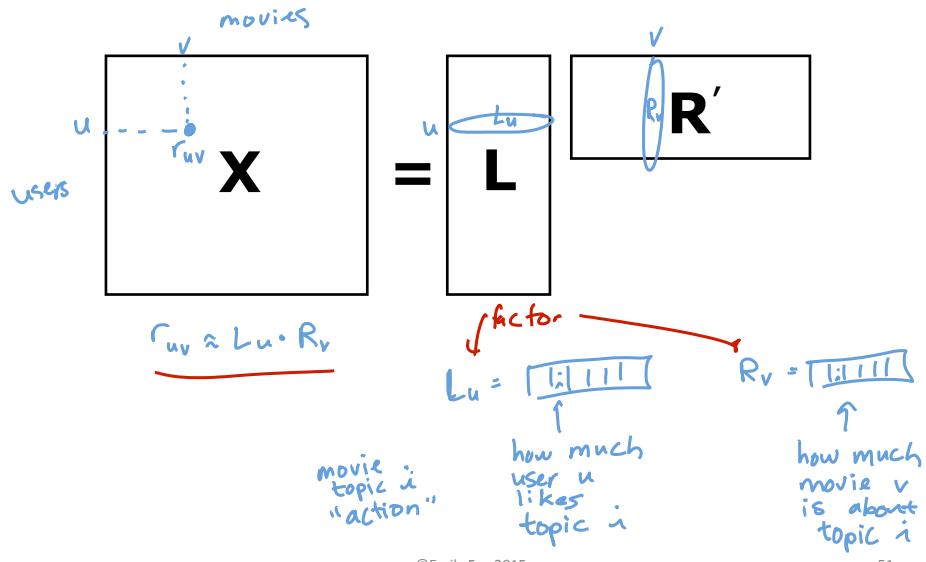
Hadoop	95 Cores _____	7.5 hrs _____
Distributed GraphLab 2012	32 EC2 machines _____	80 secs _____

©Emily Fox 2015

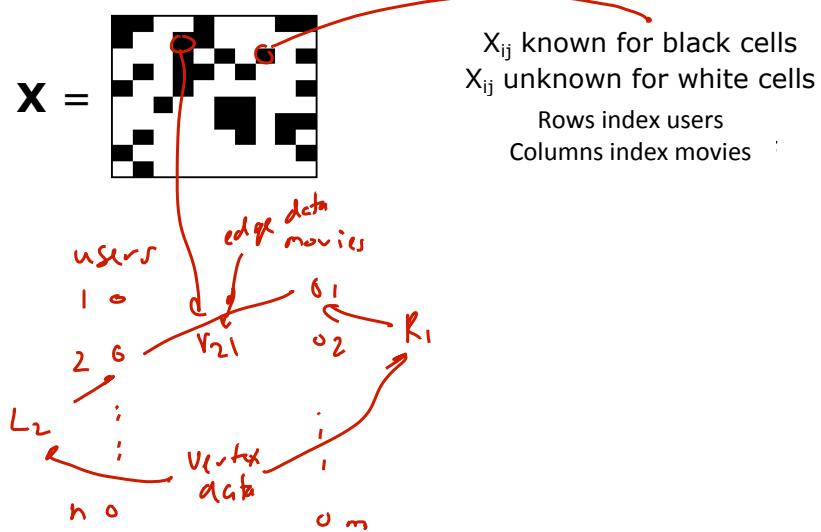
49



Interpreting Low-Rank Matrix Completion (aka Matrix Factorization)



Matrix Completion as a Graph



Coordinate Descent for Matrix Factorization: Alternating Least-Squares

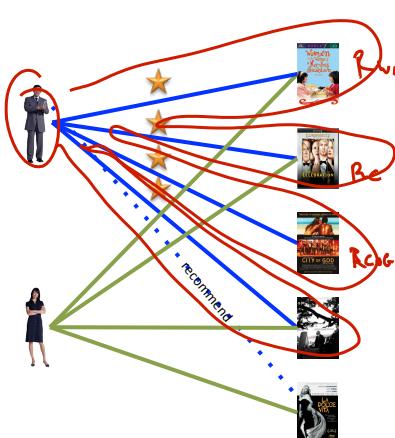
$$\min_{L,R} \sum_{(u,v):r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\| + \lambda_v \|R\|$$

- Fix movie factors, optimize for user factors
 - Independent least-squares over users
$$\min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\|$$
- Fix user factors, optimize for movie factors
 - Independent least-squares over movies
$$\min_{R_v} \sum_{u \in U_v} (L_u \cdot R_v - r_{uv})^2 + \lambda_v \|R\|$$
- System may be underdetermined: use regularization
- Converges to local optima

©Emily Fox 2015

53

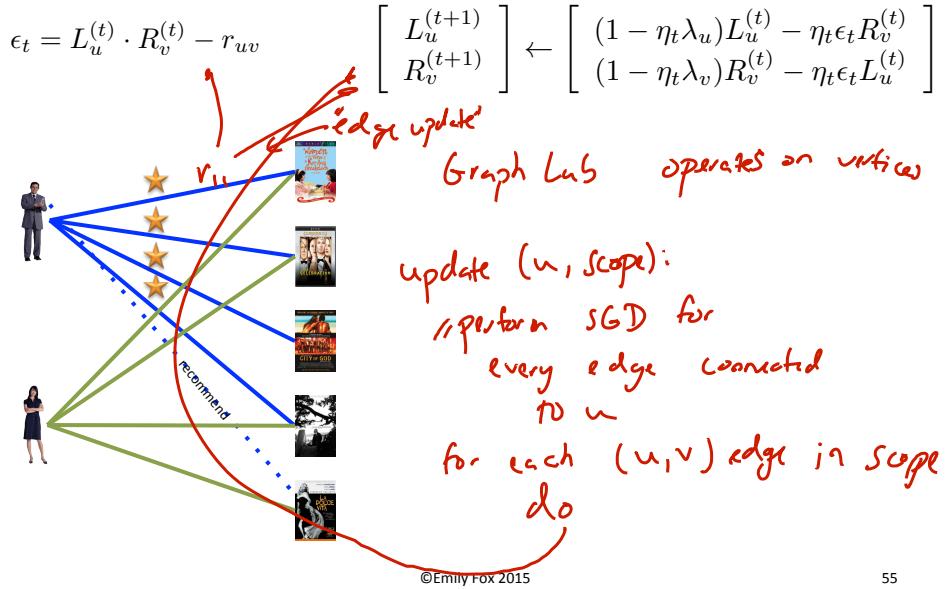
Alternating Least Squares Update Function



©Emily Fox 2015

$$Lu = (X^T X + \lambda_u I)^{-1} X^T Y$$

SGD for Matrix Factorization in GraphLab



55

Out-of-core computation

- Often data doesn't fit in memory
 - Must use disk/SSDs
- Although random accesses to disk are very slow, high performance is possible with optimized memory layout.
Implemented in:
 - GraphChi ↗ 'Mac mini', summarizing time Hadoop Cluster
 - GraphLab Create
- Example performance of GraphLab Create:
 - Common Crawl Graph (3.5 billion Nodes and 128 billion Edges)
 - PageRank iteration in 9 mins on single, commodity machine

©Emily Fox 2015

56

What you need to know...

- Data-parallel versus graph-parallel computation
- Bulk synchronous processing versus asynchronous processing
- GraphLab system for graph-parallel computation
 - Data representation
 - Update functions
 - Scheduling
 - Consistency model
- ALS, SGD for matrix factorization in GraphLab

©Emily Fox 2015

57

Reading

- Papers under “Case Study IV: **Parallel Learning with GraphLab**”
- Optional:
 - Parallel Splash BP
<http://www.ml.cmu.edu/research/dap-papers/dap-gonzalez.pdf>

©Emily Fox 2015

58