


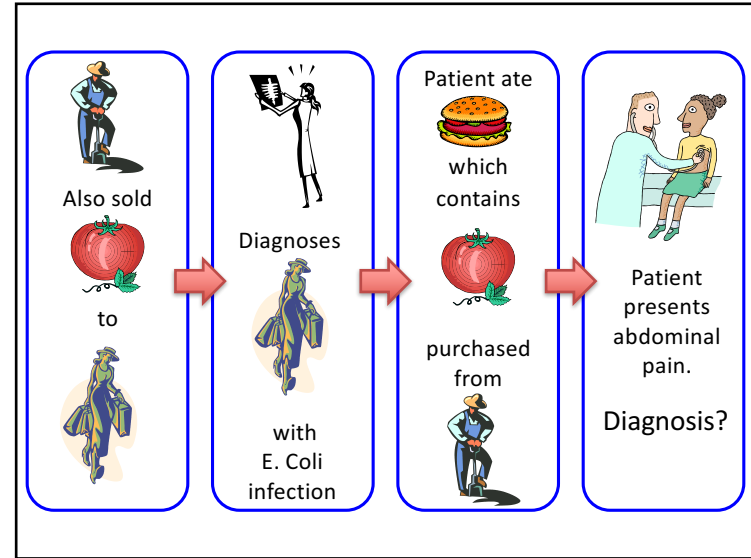
Big Data: Graph Processing



COS 418: *Distributed Systems*
Lecture 21

Kyle Jamieson

[Content adapted from J. Gonzalez]



Big Data is Everywhere



6 Billion
Flickr Photos



28 Million
Wikipedia Pages



900 Million
Facebook Users



72 Hours a Minute
YouTube

- **Machine learning** is a reality
- How will we design and implement “**Big Learning**” systems?

3

We could use ...

Threads, Locks, & Messages

“Low-level parallel primitives”

Shift Towards Use Of Parallelism in ML



- Programmers **repeatedly** solve the same **parallel design challenges**:
 - Race conditions, distributed state, communication...
- Resulting code is very **specialized**:
 - **Difficult** to maintain, extend, debug...

Idea: Avoid these problems by using **high-level abstractions**

5

... a *better* answer:

MapReduce / Hadoop

Build learning algorithms on top of high-level parallel abstractions

MapReduce – Map Phase



Embarrassingly Parallel independent computation
No Communication needed

7

MapReduce – Map Phase

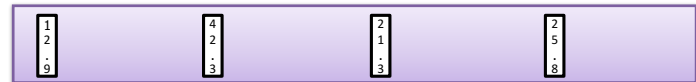
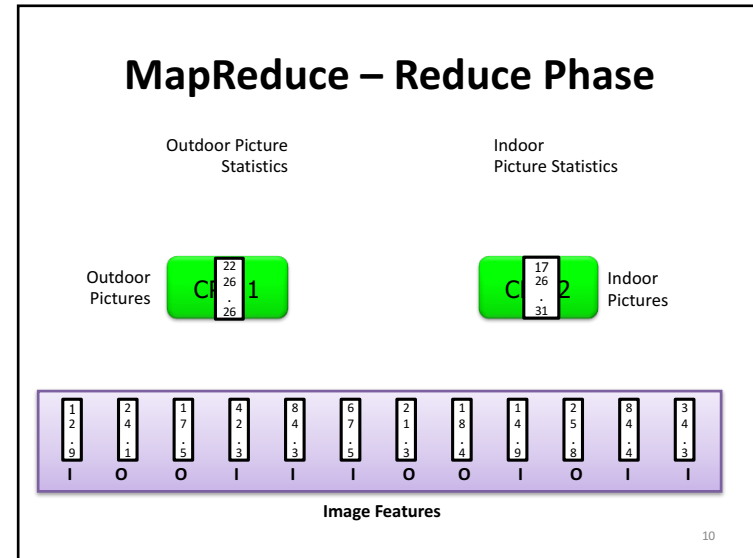
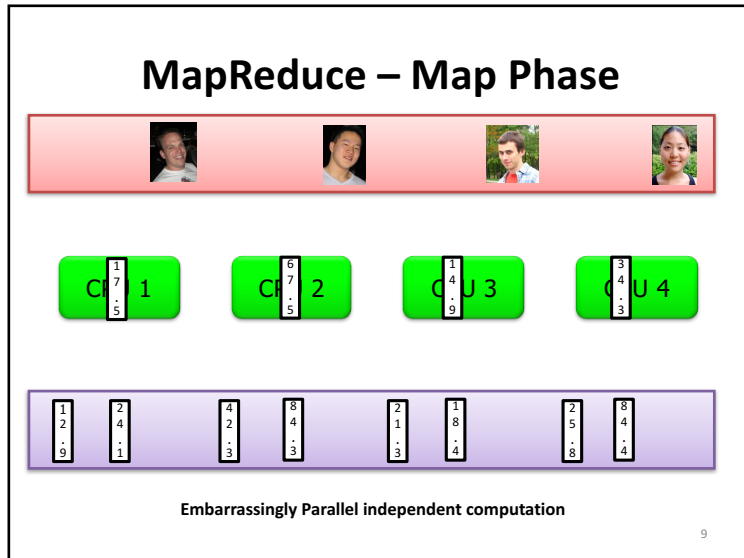


Image Features

8



Map-Reduce for Data-Parallel ML

- Excellent for **large data-parallel** tasks!

← Data-Parallel

Map Reduce

Feature Extraction	Algorithm Tuning
--------------------	------------------

Basic Data Processing

Is there more to Machine Learning ?

11

Exploiting Dependencies

Graphs are Everywhere

Social Network

Collaborative Filtering

Probabilistic Analysis

Text Analysis

Concrete Example

Label Propagation

Label Propagation Algorithm

- Social Arithmetic:**
 - 50% What I list on my profile
 - 40% Sue Ann Likes
 - + 10% Carlos Like

I Like: 60% Cameras, 40% Biking
- Recurrence Algorithm:**

$$Likes[i] = \sum_{j \in Friends[i]} W_{ij} \times Likes[j]$$
 - iterate until convergence
- Parallelism:**
 - Compute all $Likes[j]$ in parallel

Node	Weight	Cameras	Biking
Sue Ann	40%	80%	20%
Profile	50%	50%	50%
Carlos	10%	30%	70%

Properties of Graph Parallel Algorithms

Dependency Graph

Factored Computation

Iterative Computation

Map-Reduce for Data-Parallel ML

- Excellent for **large data-parallel** tasks!

<h4>MapReduce</h4> <ul style="list-style-type: none"> Feature Extraction Algorithm Tuning Basic Data Processing 	<div style="border: 1px solid black; padding: 5px; display: inline-block; margin-bottom: 10px;">MapReduce?</div> <ul style="list-style-type: none"> Lasso Tensor Factorization Deep Belief Networks Label Propagation Kernel Methods Belief Propagation PageRank Neural Networks
--------------------------------------------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

17

Problem: Data Dependencies

- MapReduce **doesn't** efficiently express data dependencies
 - User **must code** substantial data transformations
 - Costly **data replication**

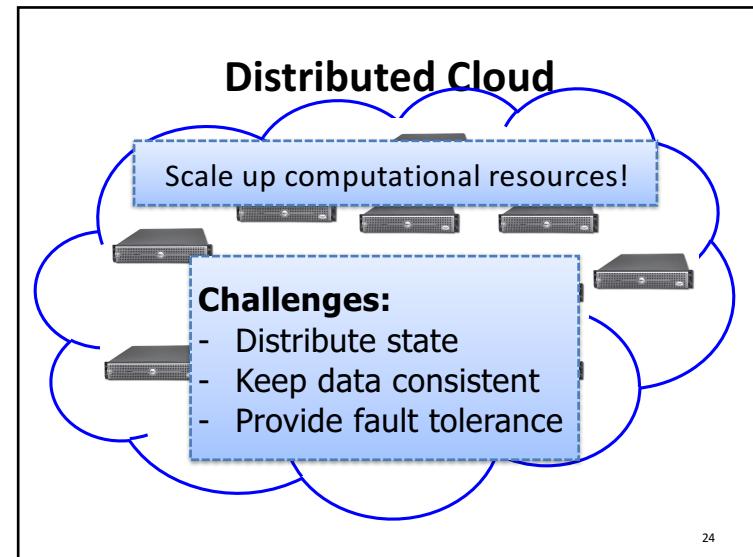
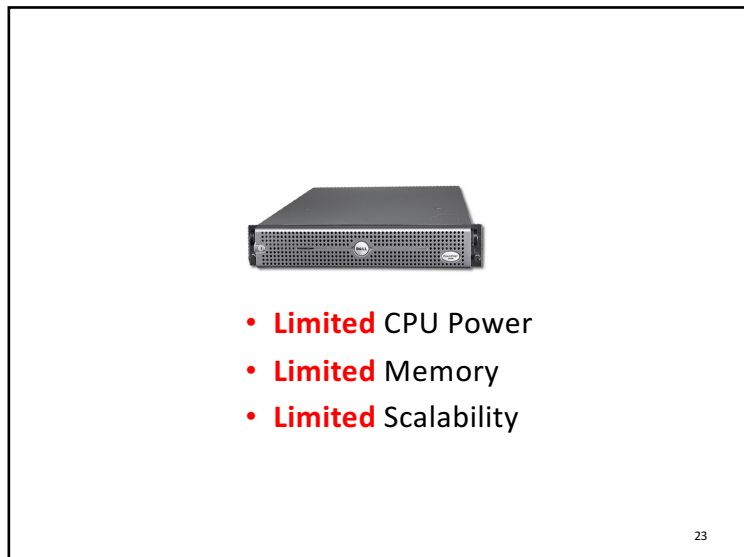
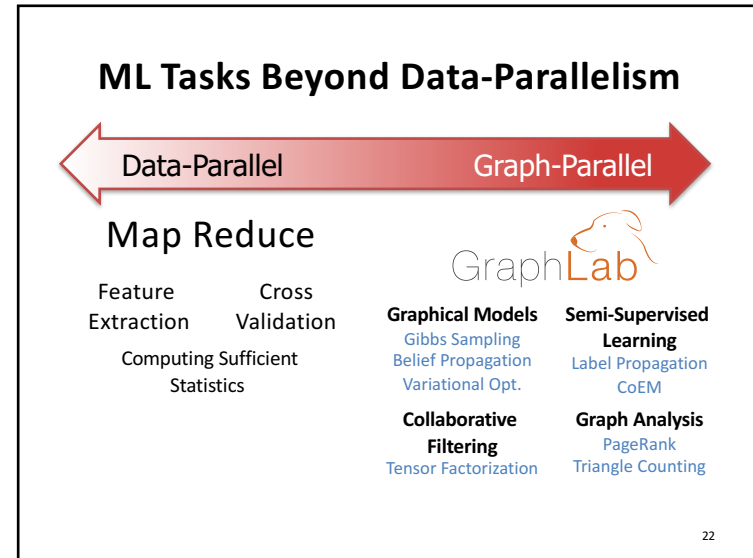
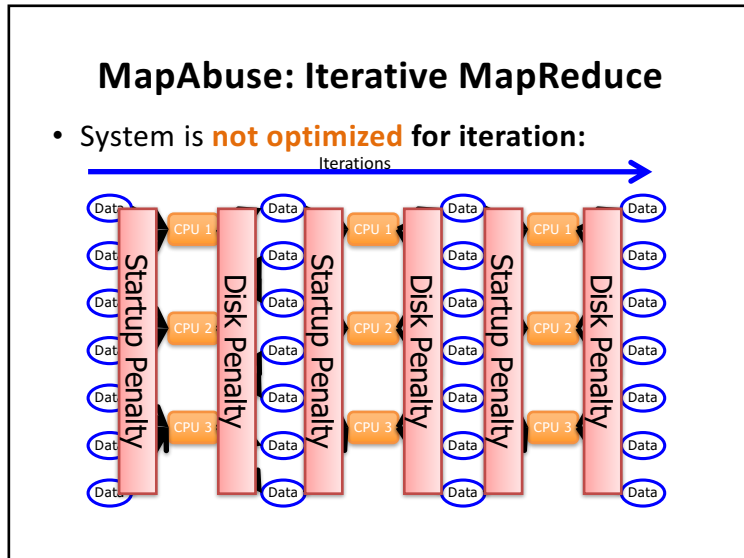
Independent Data Rows

Iterative Algorithms

- MR **doesn't** efficiently **express** iterative algorithms:

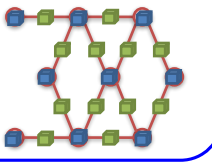
MapAbuse: Iterative MapReduce

- Only a **subset** of data needs computation:

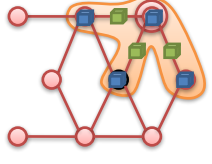


The GraphLab Framework


Graph Based
Data Representation



Update Functions
User Computation



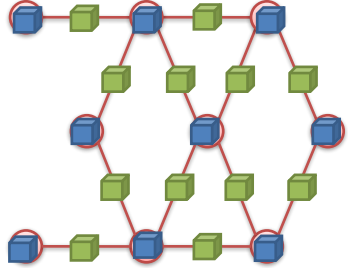
Consistency Model





25


Data Graph

Data is associated with both vertices and edges




Graph:  

- Social Network

Vertex Data: 

- User profile
- Current interests estimates

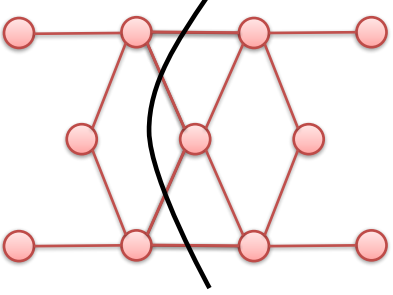
Edge Data: 

- Relationship (friend, classmate, relative)

26

Distributed Data Graph

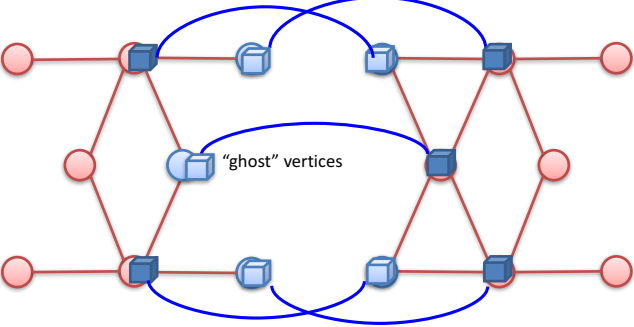
Partition the graph across multiple machines:



27

Distributed Data Graph

- **Ghost vertices** maintain adjacency structure and replicate remote data.



28

Distributed Data Graph

- Cut efficiently using **HPC Graph partitioning** tools (ParMetis / Scotch / ...)

The diagram illustrates graph partitioning. On the left, a graph is divided into two parts. The right part contains 'ghost' vertices (blue circles) which are connected to nodes in the left part. On the right, the graph is shown again, but the ghost vertices are now part of the right partition, and the original nodes in the left part are now ghost vertices.

29

The GraphLab Framework

Graph Based
Data Representation

The diagram shows the GraphLab framework components. It includes 'Update Functions' and 'User Computation' (highlighted in a blue box), 'Graph Based Data Representation' (a grid of nodes), and a 'Consistency Model' (a central node with a circular neighborhood).

30

Update Function

A user-defined **program**, applied to a **vertex**; transforms data in **scope** of vertex

The diagram shows a vertex being updated. A green box highlights the text: 'Update function applied (asynchronously) in parallel until convergence. Many schedulers available to prioritize computation'. A blue box at the bottom says 'Selectively triggers computation at neighbors'. A code snippet shows: `Pagerank(scope){` and `reschedule_neighbors() if needed if vertex PageRank changes then reschedule_all_neighbors;`

31

Distributed Scheduling

Each machine maintains a **schedule** over the vertices it owns

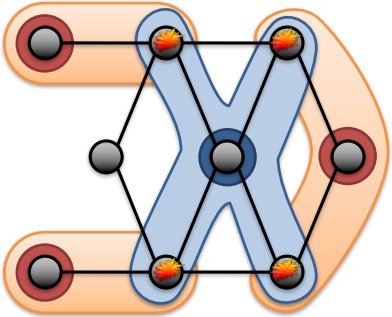
The diagram shows two server machines. The left machine has vertices a, b, e, h, i. The right machine has vertices c, d, f, g, j, k. A blue circle highlights a dependency between vertices b and f. Blue arrows on the sides indicate the flow of information.

Distributed Consensus used to identify completion

32

Ensuring Race-Free Code

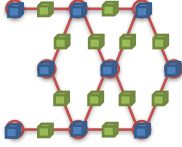
- How much can computation overlap?



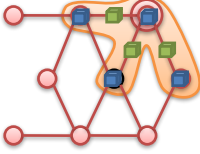
33

The GraphLab Framework

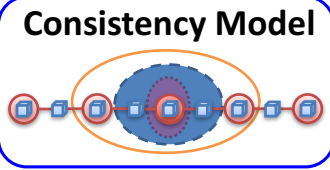
Graph Based
Data Representation



Update Functions
User Computation



Consistency Model



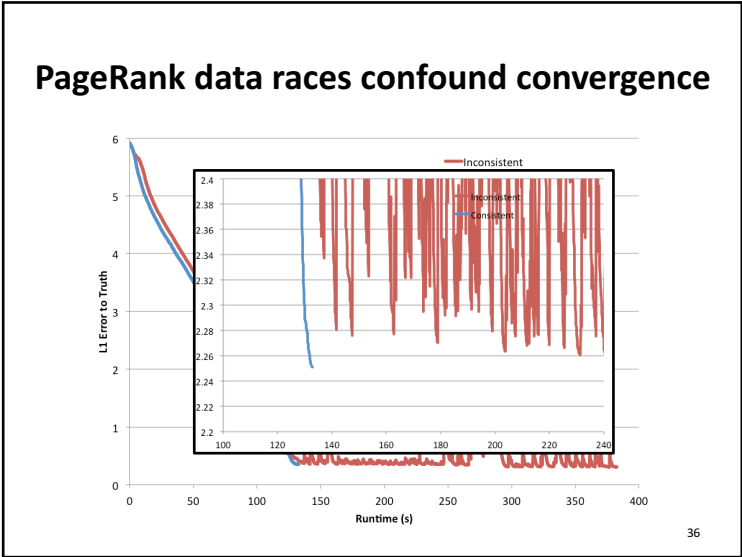
34

PageRank Revisited

```

Pagerank(scope) {
  vertex.PageRank =  $\alpha$ 
  ForEach inPage:
    vertex.PageRank +=  $(1 - \alpha) \times \text{inPage.PageRank}$ 
  ...
}
    
```

35



Racing PageRank: Bug

```

Pagerank(scope) {
  vertex.PageRank =  $\alpha$ 
  ForEach inPage:
    vertex.PageRank += (1 -  $\alpha$ ) × inPage.PageRank
  ...
}
    
```

37

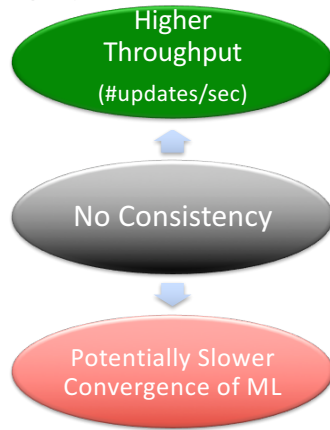
Racing PageRank: Bug Fix

```

Pagerank(scope) {
  tmp =  $\alpha$ 
  ForEach inPage:
    tmp += (1 -  $\alpha$ ) × inPage.PageRank
  vertex.PageRank = tmp
  ...
}
    
```

38

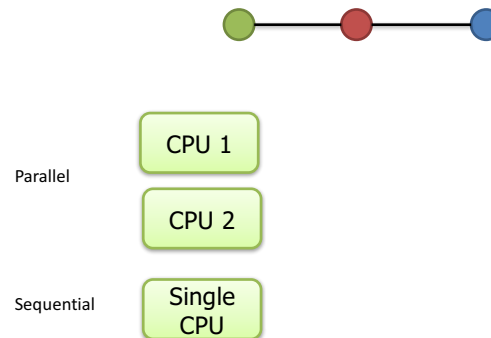
Throughput != Performance



39

Serializability

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



40

Serializability Example

Stronger / Weaker consistency levels available

User-tunable consistency levels trades off parallelism & consistency

Overlapping regions are only read.

Update functions **one** vertex apart can be run in parallel.

Edge Consistency

41

Distributed Consistency

- **Solution 1: Chromatic Engine**
 - Edge Consistency via **Graph Coloring**

- **Solution 2: Distributed Locking**

Chromatic Distributed Engine

Execute tasks on all vertices of color 0

Execute tasks on all vertices of color 0

Ghost Synchronization Completion + Barrier

Execute tasks on all vertices of color 1

Execute tasks on all vertices of color 1

Ghost Synchronization Completion + Barrier

Time

43

Matrix Factorization

- **Netflix Collaborative Filtering**
 - Alternating Least Squares Matrix Factorization

Model: 0.5 million nodes, 99 million edges

Users

Netflix

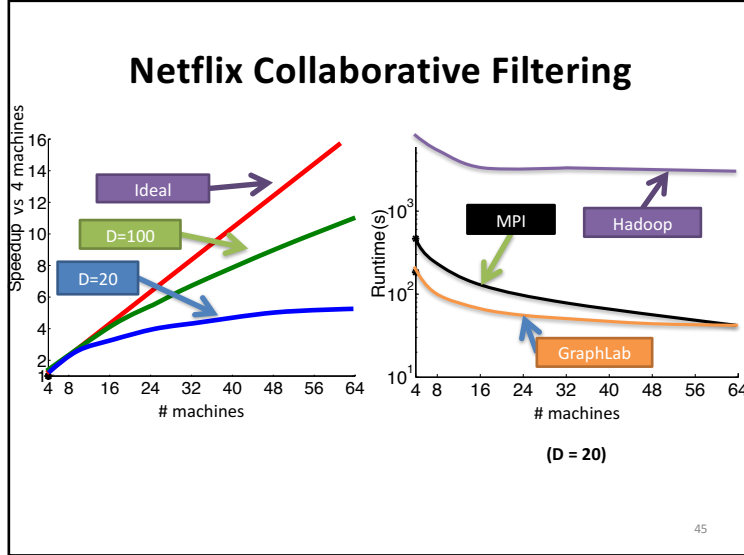
Movies

Users

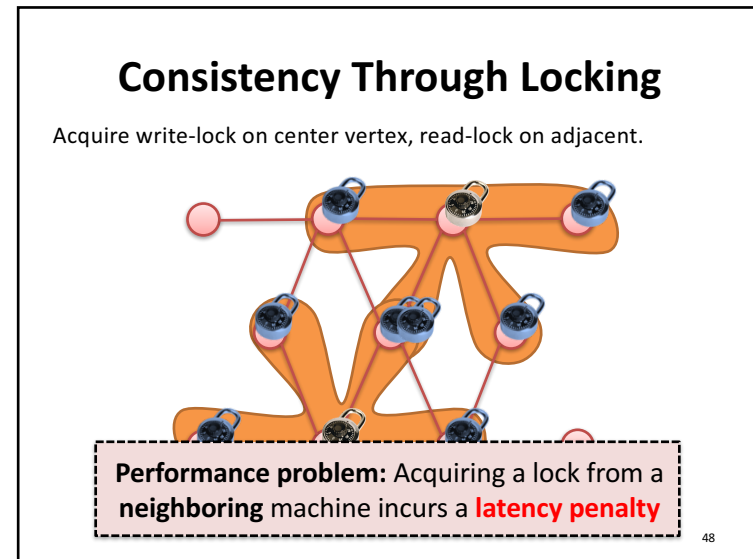
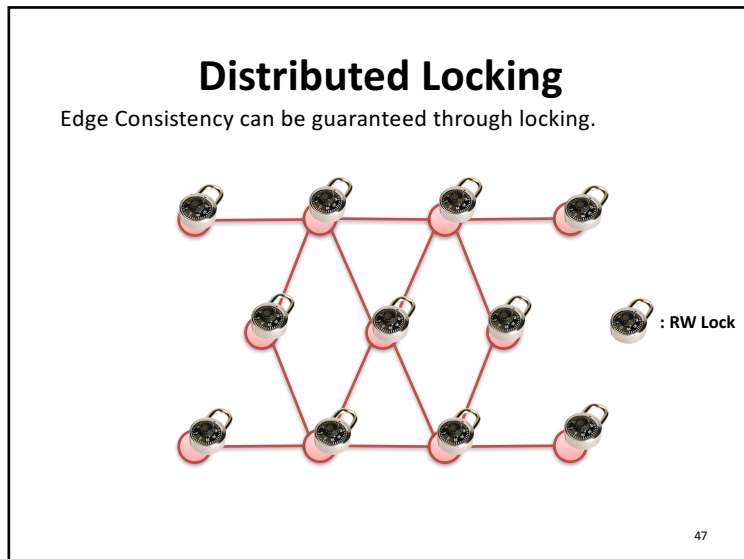
Movies

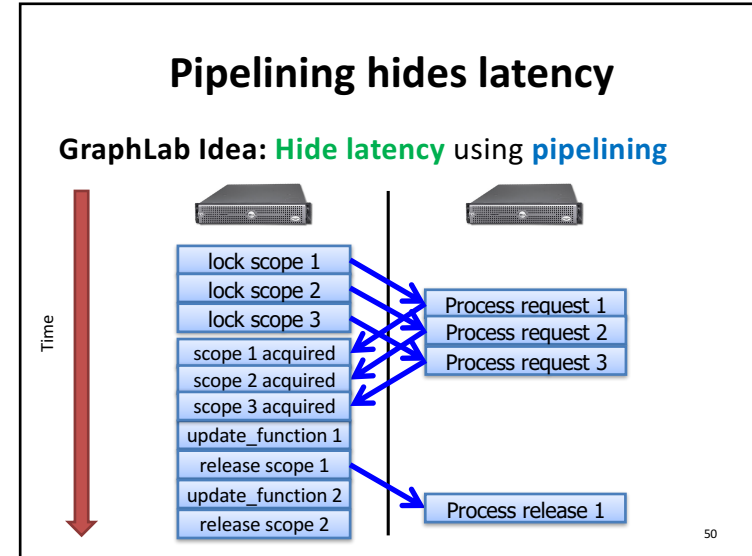
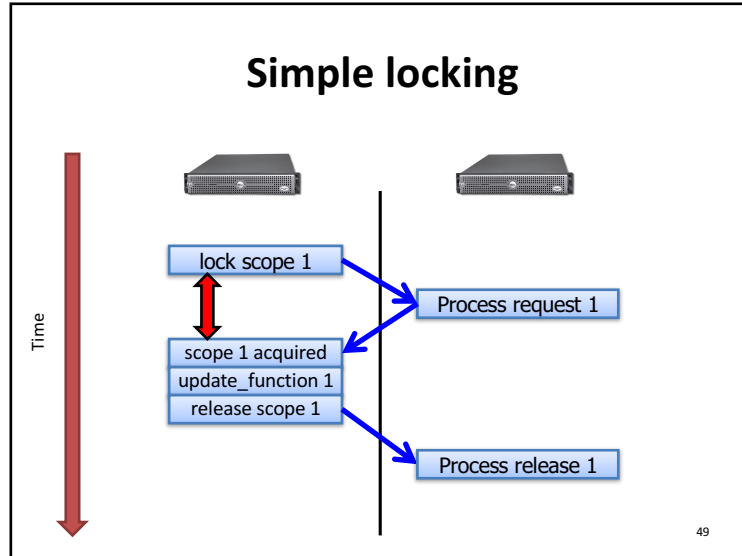
D

44



- ### Distributed Consistency
- **Solution 1: Chromatic Engine**
 - Edge Consistency via **Graph Coloring**
 - **Requires a graph coloring** to be available
 - Frequent barriers → **inefficient** when only **some** vertices active
 - **Solution 2: Distributed Locking**



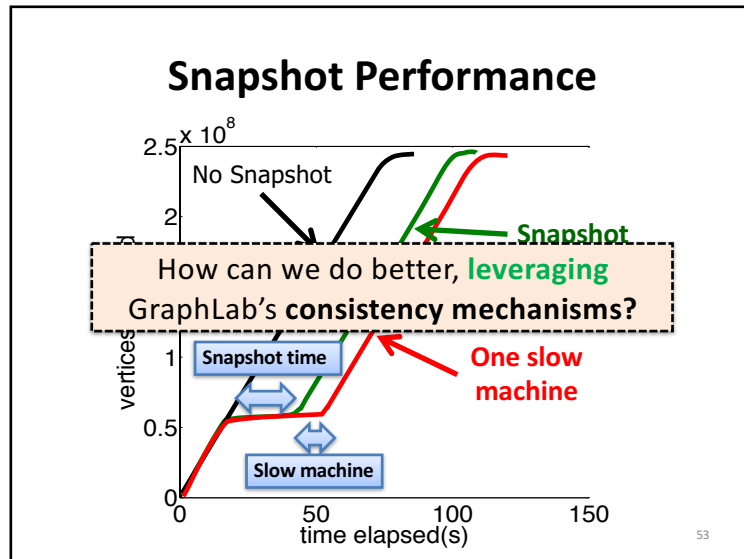


Distributed Consistency

- **Solution 1: Chromatic Engine**
 - Edge Consistency via **Graph Coloring**
 - **Requires a graph coloring** to be available
 - Frequent barriers → **inefficient** when only **some** vertices active
- **Solution 2: Distributed Locking**
 - Residual BP on 190K-vertex/560K-edge graph, 4 machines
 - No pipelining: 472 sec; **with pipelining: 10 sec**

How to handle machine failure?

- *What when machines fail?* **How** do we **provide fault tolerance?**
- Strawman scheme: **Synchronous snapshot** checkpointing
 1. Stop the world
 2. Write each machines' state to disk



Chandy-Lamport checkpointing

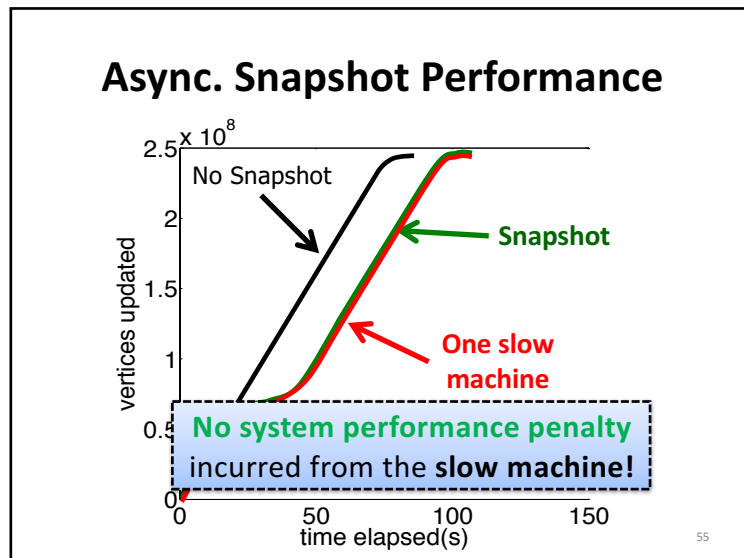
Step 1. Atomically one *initiator*

(a) Turns red, (b) Records its own state
(c) sends *marker* to neighbors

Step 2. On receiving marker **non-red** node atomically: (a) Turns red,
(b) Records its own state, (c) sends markers along all outgoing channels

First-in, first-out channels

Implemented within GraphLab as an **Update Function**



Summary

- Two different methods of achieving **consistency**
 - Graph Coloring
 - Distributed Locking with pipelining
- **Efficient implementations**
- **Asynchronous FT** w/fine-grained Chandy-Lamport

Performance

Efficiency

Scalability

Useability

56

Friday Precept:
Roofnet performance
More Graph Processing

Monday topic:
Streaming Data Processing

57