

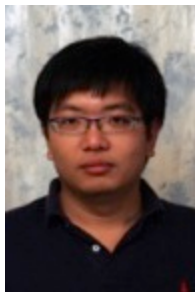


## Large-Scale Distributed Machine Learning

Carlos Guestrin



Joseph  
Gonzalez



Yucheng  
Low



Aapo  
Kyrola



Haijie  
Gu



Joseph  
Bradley



Danny  
Bickson

# Needless to Say, We Need Machine Learning for Big Data



6 Billion  
Flickr Photos



28 Million  
Wikipedia Pages



1 Billion  
Facebook Users



72 Hours a Minute  
YouTube

The New York Times  
**SundayReview**

WORLD U.S. N.Y. / REGION BUSINESS TEC

NEWS ANALYSIS

## The Age of Big Data

By STEVE LOHR

Published: February 11, 2012

“... data a new class of economic asset,  
like currency or gold.”

# Big Learning

How will we  
**design and implement**  
*parallel* learning systems?

Part 1

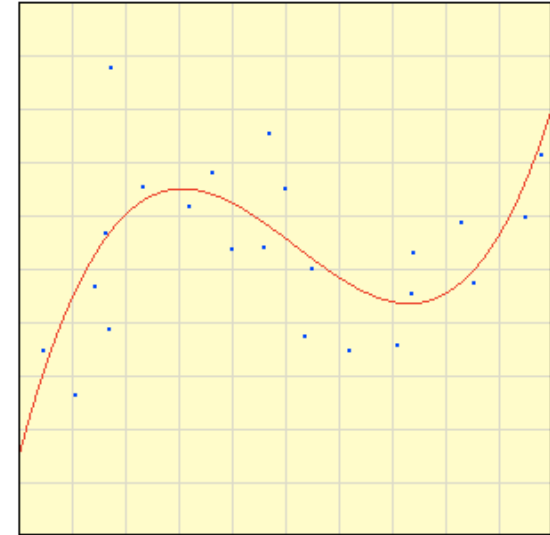
# **ASYNCHRONOUS DATA-PARALLEL ALGORITHMS**



# Sparse Regression

[Tibshirani, 1996]

$$\begin{array}{lcl} y & \approx & w_0 + aw_1 + a^2w_2 + a^3w_3 + \dots \\ \uparrow & & \uparrow \quad \leftarrow \text{weights} \\ \text{target} & = & Aw \\ & & \uparrow \\ & & \text{basis functions} \end{array}$$



**LASSO: find sparse weight vector  $w^*$**

$$\min_w F(w)$$

$$F(w) = ||y - Aw||_2^2 + \lambda ||w||_1$$

least  
squares

sparsity  
inducing  
regularizer

- Fundamental machine learning task
- Huge number of applications (many thousands of papers)
  - Computational biology, computer vision, compressed sensing...

# Shooting: Stochastic Coordinate Descent (SCD)

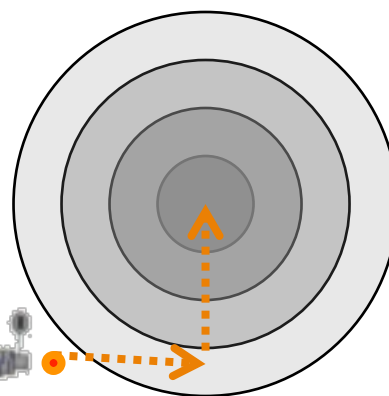
[e.g., Shalev-Shwartz & Tewari '09]

While not converged

- Choose random coordinate  $j$
- Optimize  $w_j$   
(closed-form minimization)



$F(w)$  contour



Lasso:  $\min_w F(w)$  where  $F(w) = \frac{1}{2} \|\mathbf{A}w - \mathbf{y}\|_2^2 + \lambda \|w\|_1$

# Coordinate Descent for LASSO (aka Shooting Algorithm)

- Repeat until convergence
  - Pick a coordinate  $j$  at (random or sequentially)

- Set:

$$\hat{w}_\ell = \begin{cases} (c_\ell + \lambda)/a_\ell & c_\ell < -\lambda \\ 0 & c_\ell \in [-\lambda, \lambda] \\ (c_\ell - \lambda)/a_\ell & c_\ell > \lambda \end{cases}$$

- Where:

$$a_\ell = 2 \sum_{j=1}^N (h_\ell(\mathbf{x}_j))^2$$

$$c_\ell = 2 \sum_{j=1}^N h_\ell(\mathbf{x}_j) \left( t(\mathbf{x}_j) - (w_0 + \sum_{i \neq \ell} w_i h_i(\mathbf{x}_j)) \right)$$

# Analysis of SCD [Shalev-Shwartz, Tewari '09/'11]

- Theorem: With iterations  $T$ , expected error decreases as:

$$\mathcal{O} \left( \frac{d \gamma ||w^*||^2}{T} \right)$$

- For  $d$  dimensions, and optimum  $w^*$
- For (coordinate-wise) strongly convex functions ( $\Delta w = \delta_{wj} e_j$ ):  
$$F(w + \Delta w) \leq F(w) + |\Delta w|(\nabla F(w))_j + \frac{\gamma |\Delta w|^2}{2}$$
- For LASSO  $\gamma=1$ , for Logistic Regression  $\gamma=1/4$

*Great rate...*

*but gets expensive in high dimensions*

# Shotgun: Data-Parallel SCD

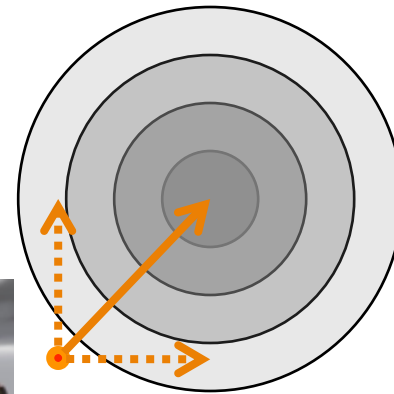
[Bradley, Kyrola, Bickson, G. '11]

While not converged

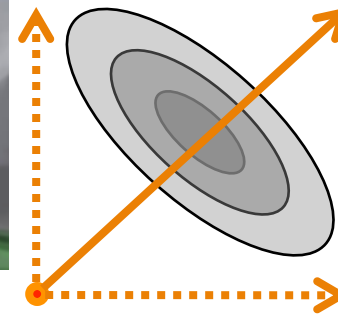
- On each of  $P$  processors
  - Choose random coordinate  $j$
  - Optimize  $w_j$  (as in Shooting)



*Is coordinate descent inherently sequential?*



Nice case:  
Uncorrelated  
features



Bad case:  
Correlated  
features

Lasso:  $\min_w F(w)$  where  $F(w) = \frac{1}{2} \| \mathbf{A}w - \mathbf{y} \|^2_2 + \lambda \| w \|_1$

# Is SCD inherently sequential?

Lasso:  $\min_w F(w)$  where  $F(w) = \|Xw - \mathbf{y}\|_2^2 + \lambda \|w\|_1$

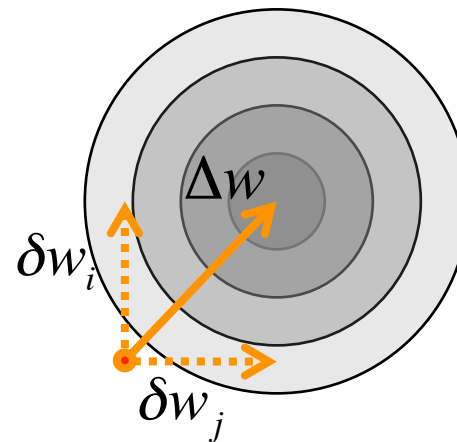
Coordinate update:

$$w_j \leftarrow w_j + \delta w_j$$

(closed-form minimization)

Collective update:

$$\Delta w = \begin{pmatrix} \delta w_i \\ 0 \\ 0 \\ \delta w_j \\ 0 \end{pmatrix}$$



# Is SCD inherently sequential?

$$\text{Lasso: } \min_w F(w) \quad \text{where} \quad F(w) = \|Xw - \mathbf{y}\|_2^2 + \lambda \|w\|_1$$

Lemma: If  $X$  is normalized s.t.  $\text{diag}(X^T X) = 1$ ,

$$F(w + \Delta w) - F(w) \leq \underbrace{-\sum_{i_j \in \mathcal{P}} (\delta w_{i_j})^2}_{\text{"Positive" progress}} + \sum_{\substack{i_j, i_k \in \mathcal{P}, \\ j \neq k}} \underbrace{(X^T X)_{i_j, i_k}}_{\text{Can be positive or negative } \odot} \delta w_{i_j} \delta w_{i_k}$$

**Key term!**

(Measures "correlation" between features...)

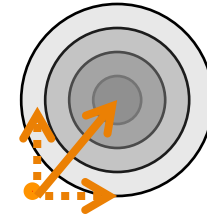
**"interference" between updates**

# Theorem: Shotgun Convergence

Assume  $P < d/\rho + 1$   
 where  $\rho$  = largest eigenvalue of  $\mathbf{A}^T\mathbf{A}$

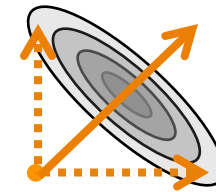
*Then:* can achieve  
 linear speed ups with  
 up to  $P$  processors

Nice case:  
 Uncorrelated  
 features



$$\rho = 1 \Rightarrow P_{\max} = d$$

Bad case:  
 Correlated  
 features

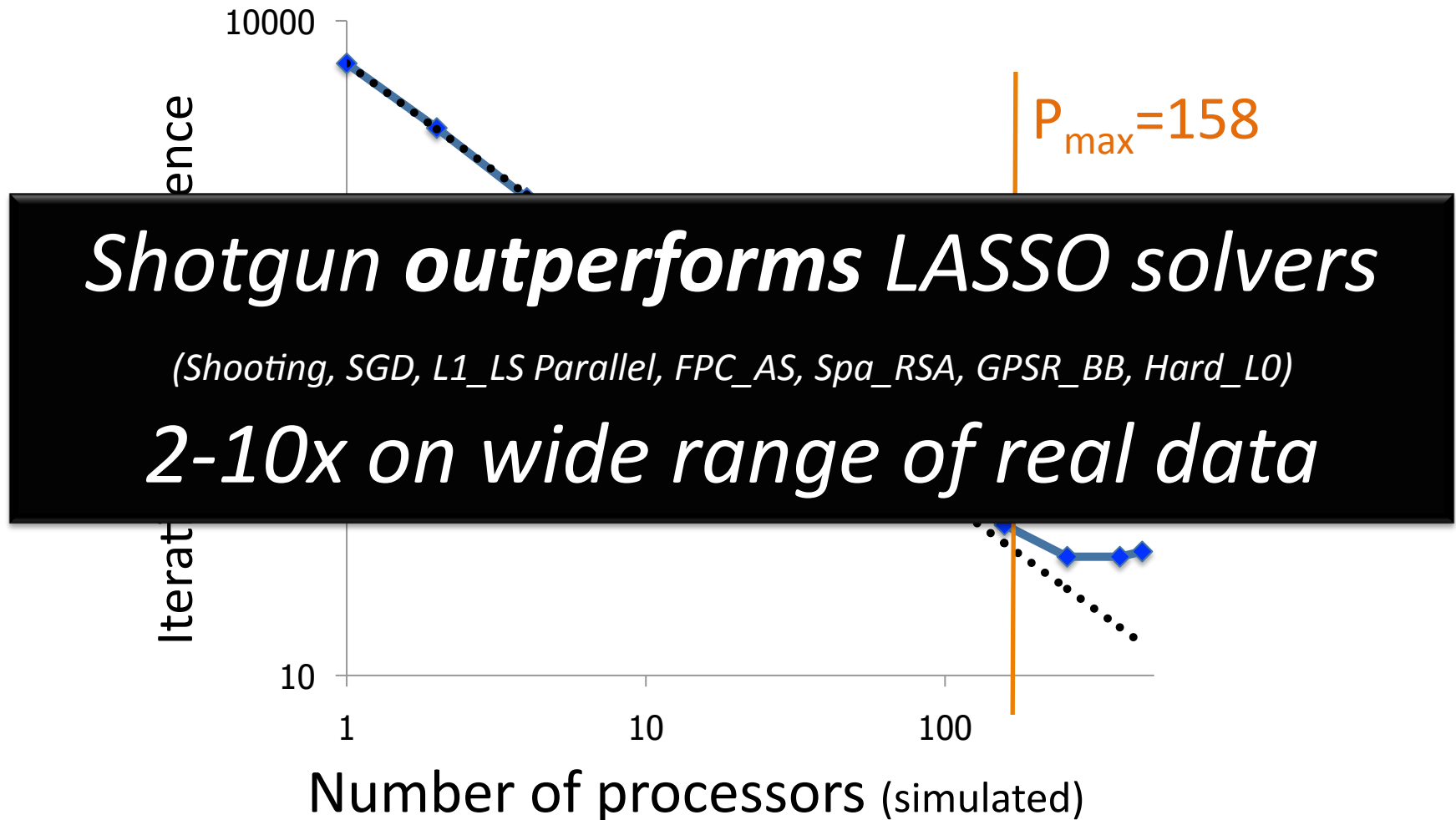


$$\rho = d \Rightarrow P_{\max} = 1 \text{ (at worst)}$$

$$E \left[ \overbrace{F(w^{(T)})}^{\text{final - opt objective}} \right] - F(w^*) \leq \frac{d \left( \frac{1}{2} \|w^*\|_2^2 + F(w^{(0)}) \right)}{\overbrace{TP}^{\substack{\text{iterations} \quad \# \text{ parallel updates}}}}$$



# Experiments Match Theory!



Mug32, single pixel camera dataset

# Key Proof Technique

Parallel optimization problem

```
graph TD; A[Parallel optimization problem] --> B[Potential interference between updates]; B --> C[Guarantee based on bounding magnitude of interference];
```

Potential interference  
between updates

Guarantee based on bounding  
magnitude of interference

# Stepping Back...

- Stochastic coordinate ascent (SCD)
  - Optimization: Pick a coordinate  $j$ , find  $\operatorname{argmin}_{w_j} F(w)$
  - Parallel SCD: Pick  $p$  coordinates and update at once
  - Issue: Updates may interfere on  $p$  coordinates
  - Solution: Bound possible interference using spectral norm
- Natural counterpart: Stochastic gradient descent (SGD)
  - Optimization: Pick a data point and take a small gradient step on all coordinates
  - Parallel: Pick  $p$  data points and update at once
  - Issue: Updates may interfere on *all* coordinates
  - Solution: Bound interference using sparsity of data points

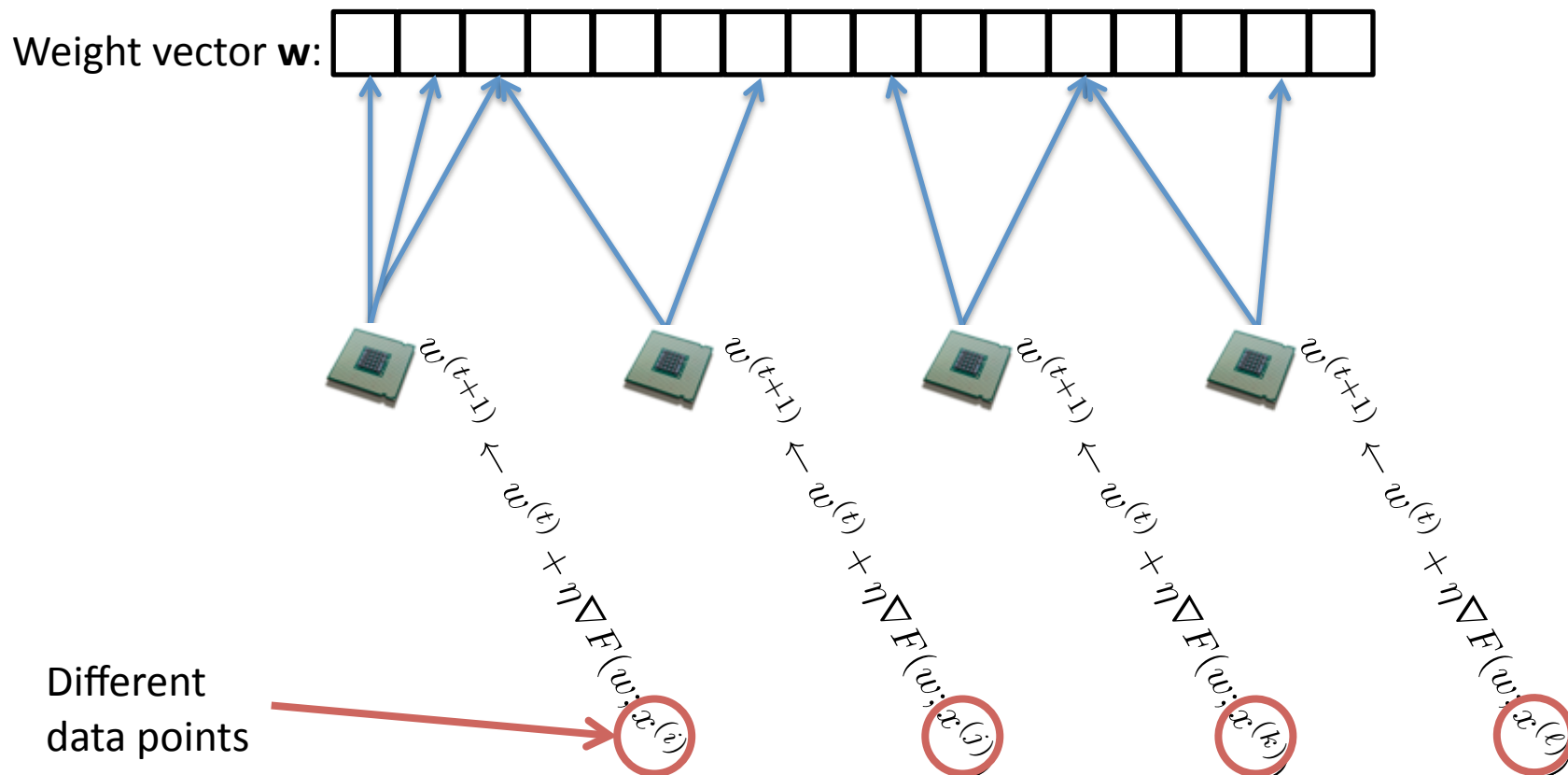
# Stochastic Gradient Descent

- Coordinate descent updates one coordinate  $w_j$ , using all data points
- Stochastic gradient descent updates all coordinates, using one data point  $x^{(i)}$ :

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \nabla F(w; x^{(i)})$$

# Parallel Stochastic Gradient Descent

- Each processor does update using a different data point



*Risk versus coordinate descent:  
SGD could interfere on all coordinates simultaneously*

# Parallel SGD with No Locks

[e.g., Hogwild!, Niu et al. '11]

- Each processor in parallel:
  - Pick data point  $i$  at random
  - For  $j = 1 \dots d$ :

$$w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta \left( \nabla F(w; x^{(i)}) \right)_j$$

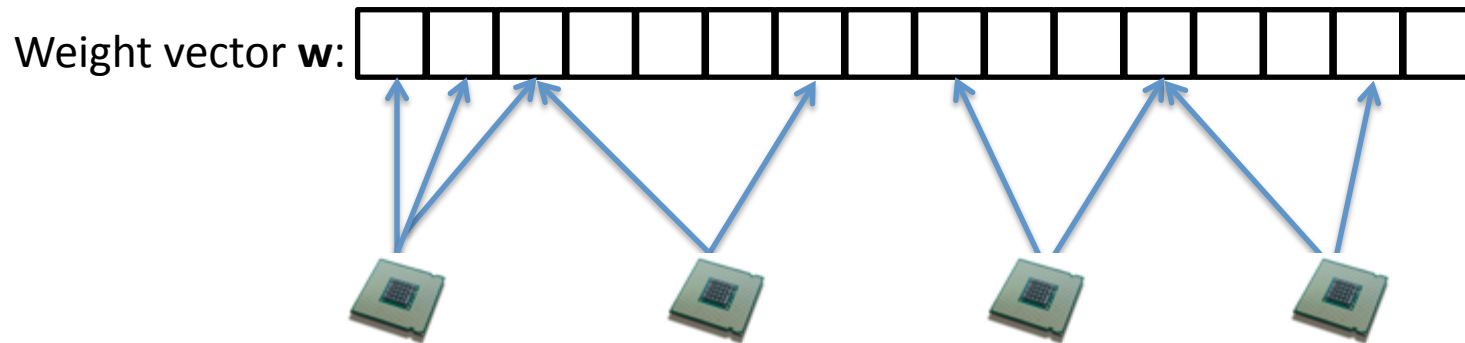
- Assume atomicity of sum operation for a coordinate:

$$w_j \leftarrow w_j + a$$

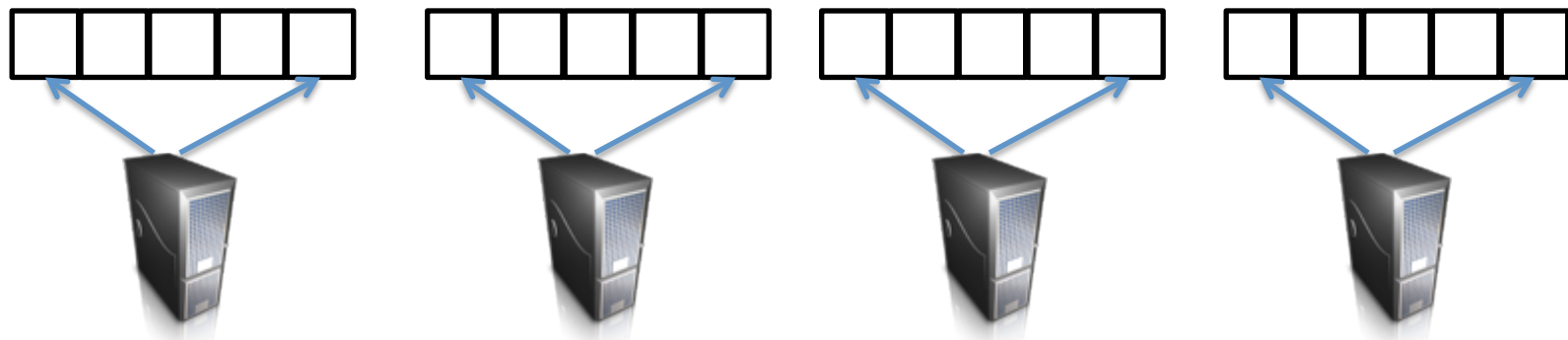
*Key to proof of bounded interference:  
Assume data points are sparse  $\rightarrow$   
update interferes at most on a few coordinates*

# Shared Memory versus Distributed Memory

- **Shared memory:** all machines can access same memory space

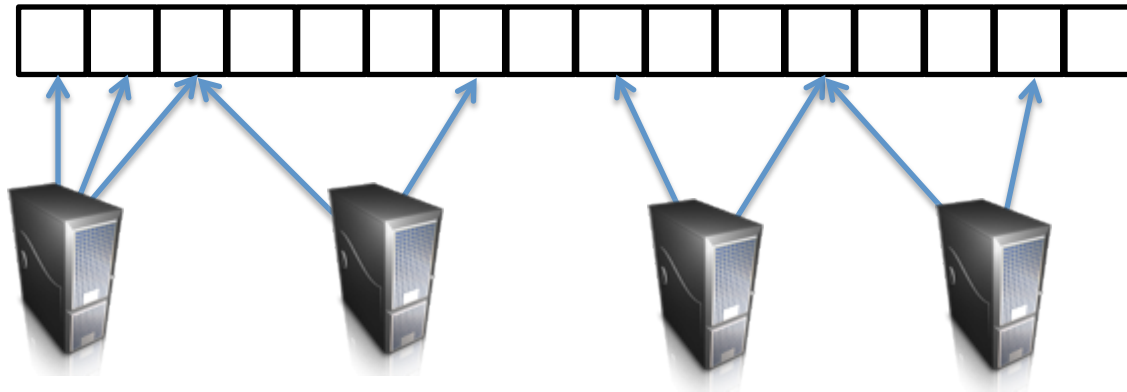


- **Distributed memory:** machines can only access local memory



- Much harder to implement Shotgun or Hogwild!, because of need to synchronize parameters across machines
  - Synchronization can be extremely slow

# Distributed Hash Tables (DHTs)



- **DHT:** Distributed memory that looks like shared memory from the programmer's perspective



Easy to program

Guarantees consistency of values read/written



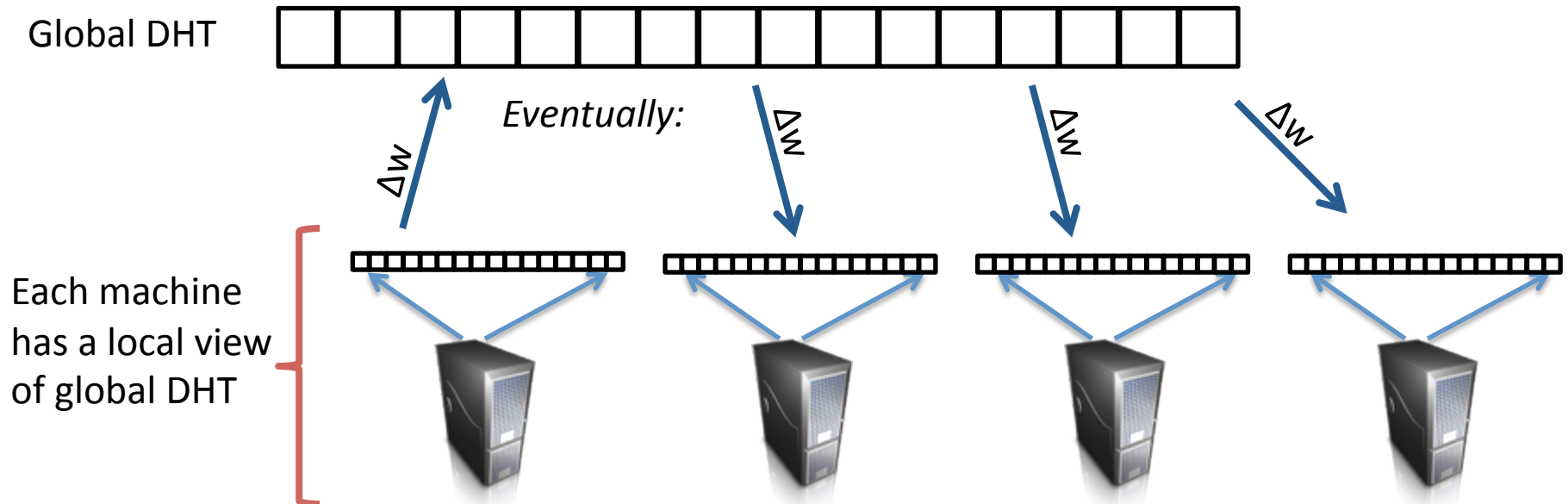
Only really efficient when “large” objects are written/read

In ML, an “object” is a parameter, just a double  
➔ standard DHTs are too slow



# Parameter Servers (e.g., Smola et al.)

- A parameter server is a **Lazy DHT** with **commutative-associative operations**, e.g.,  $w_j \leftarrow w_j + a$



- Parameter servers only guarantee eventual consistency
- But, often good enough for many distributed learning procedures

# Summary of Part 1

- Shotgun/Hogwild! solve distributed optimization by ignoring dependencies in problem
- Key proof method: bound interference in updates
- Implement in distributed settings using parameter servers

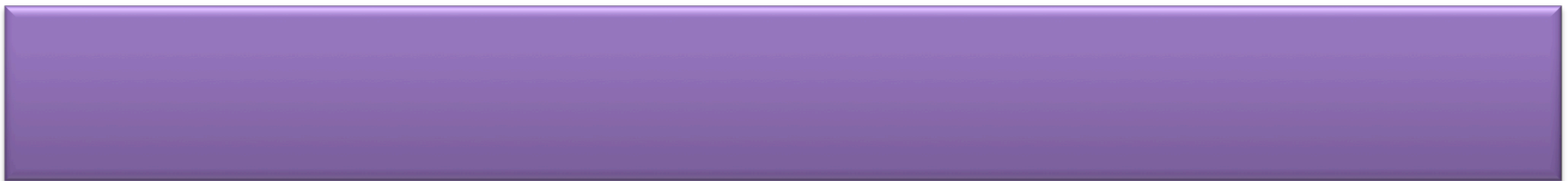
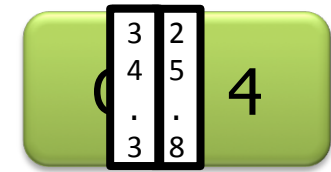
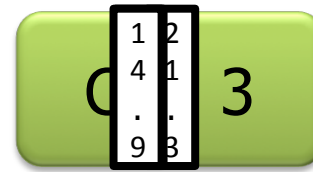
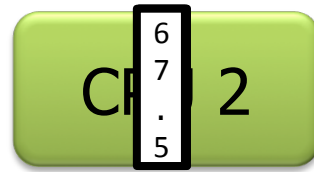
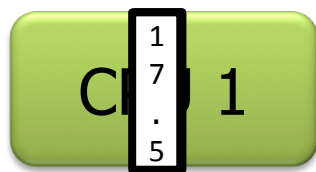
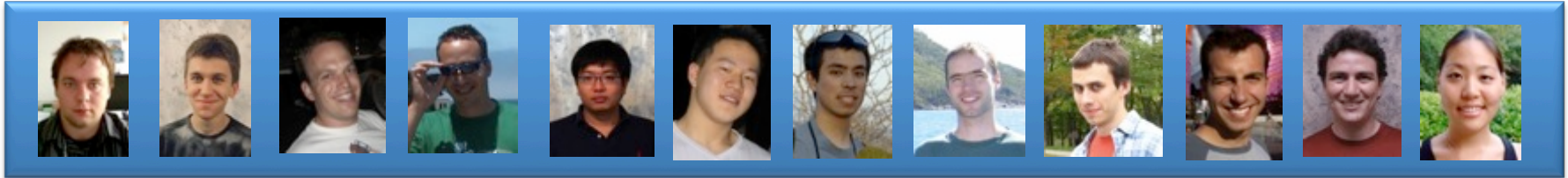
Part 2

# **ASYNCHRONOUS GRAPH-PARALLEL ALGORITHMS**

DATA PARALLEL  
versus  
GRAPH PARALLEL

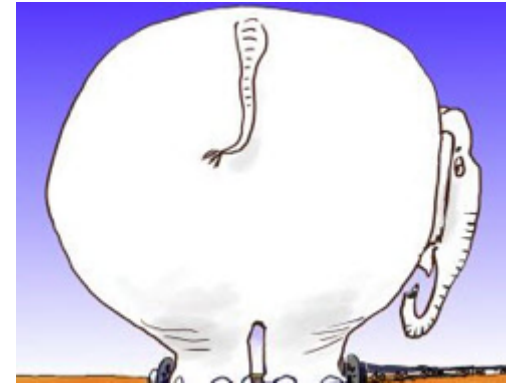
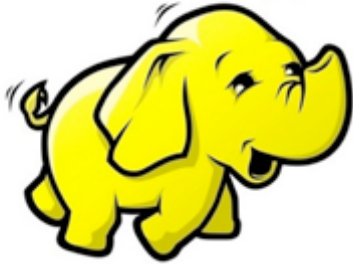
abstractions

# Data Parallelism (MapReduce)



*Solve a huge number of independent subproblems*

**hadoop**

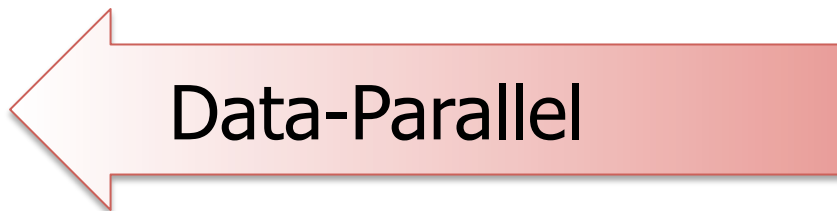


“A white elephant is a valuable but burdensome possession of which its owner cannot dispose and **whose cost (particularly cost of upkeep) is out of proportion to its usefulness or worth.**” *Wikipedia*

Everyone knows has limitations,  
nobody happy,  
but **what to do next???**

# MapReduce for Data-Parallel ML

Excellent for large data-parallel tasks!



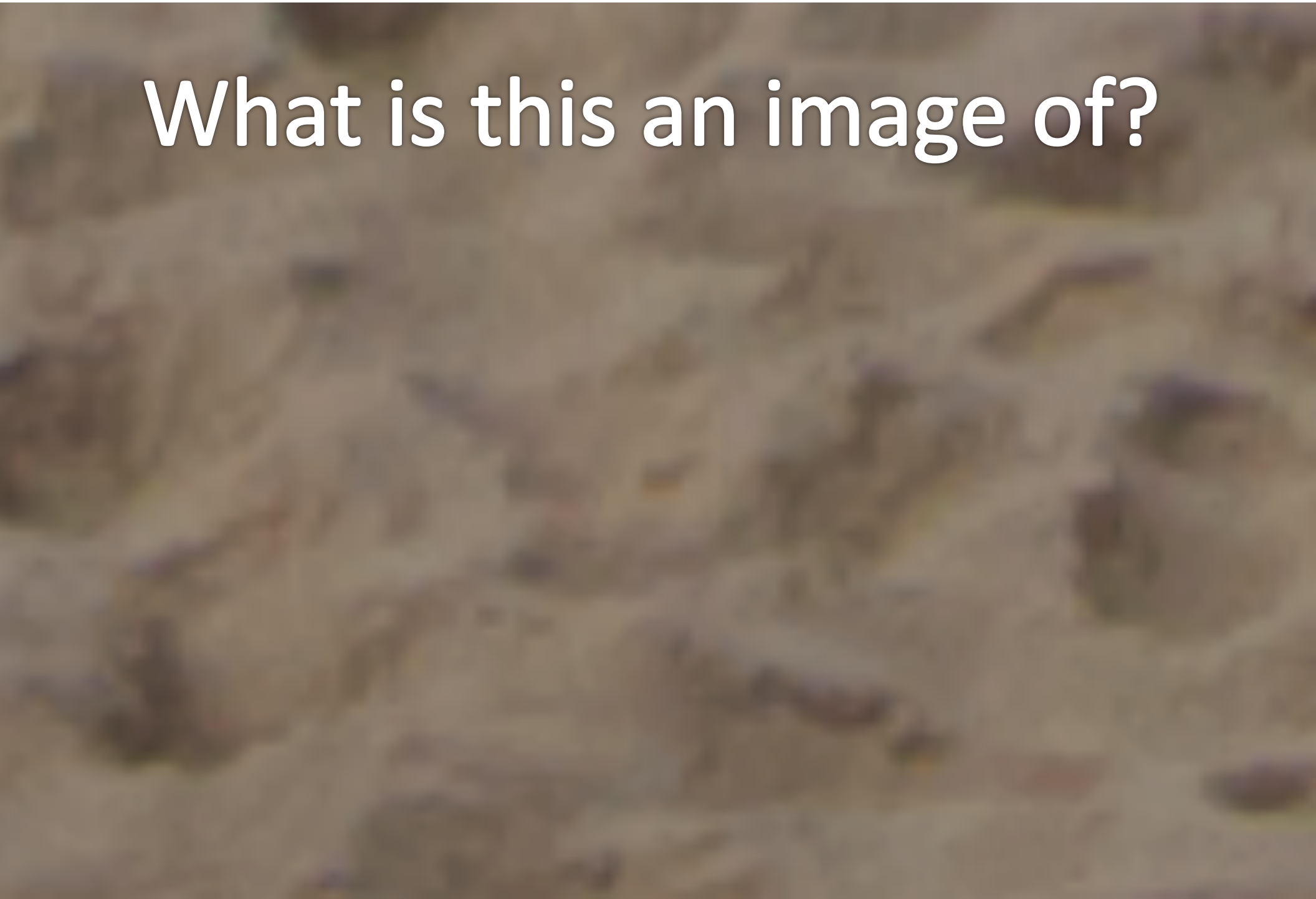
## MapReduce

Feature Extraction      Cross Validation  
Computing Sufficient Statistics

Is there more to  
Machine Learning

?

What is this an image of?





It's next to this...





# The Power of Dependencies

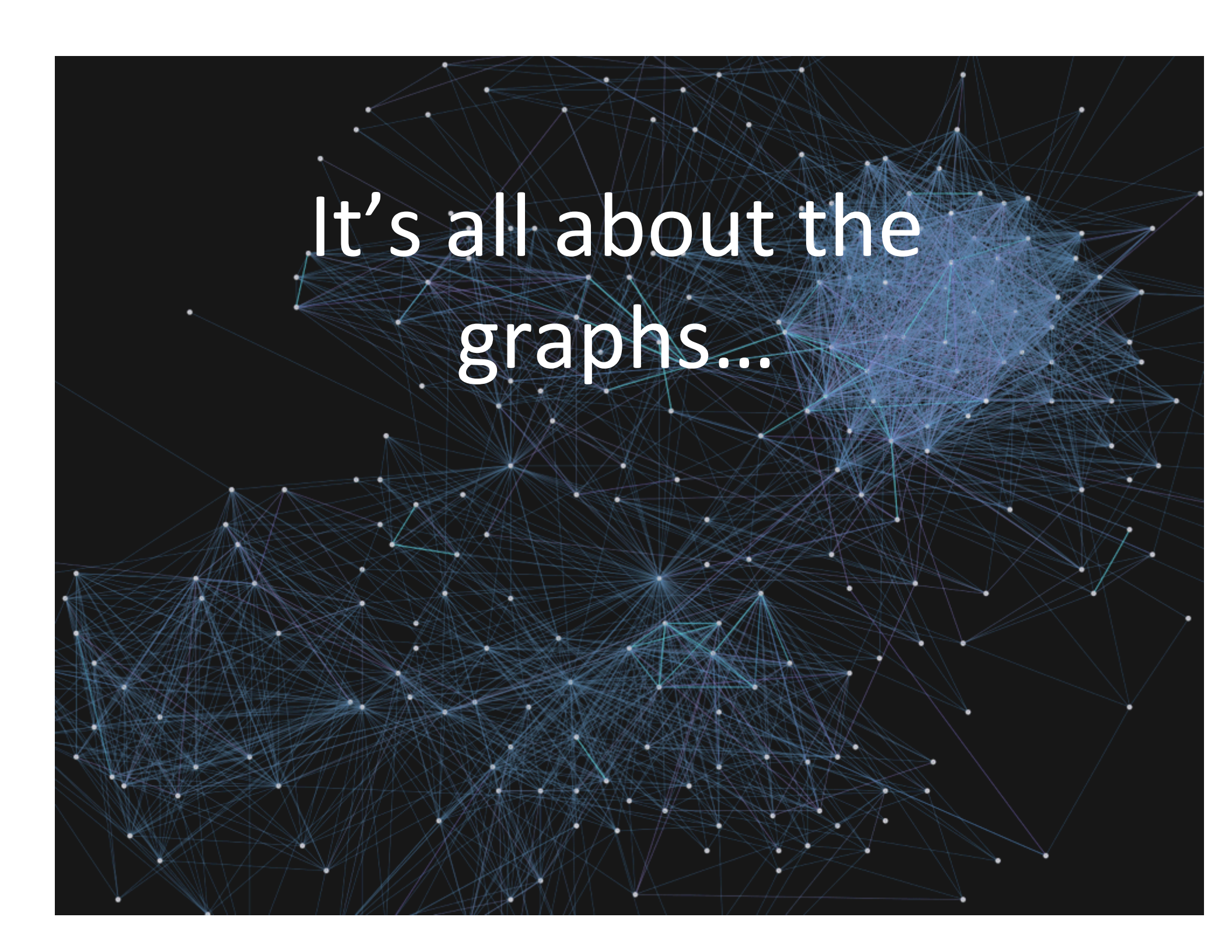
***where the value is!***

# Flashback to 1998



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



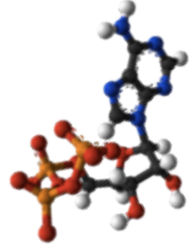


It's all about the  
graphs...

## Social Media



## Science



## Advertising



## Web



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

People

Products

Ideas

Facts

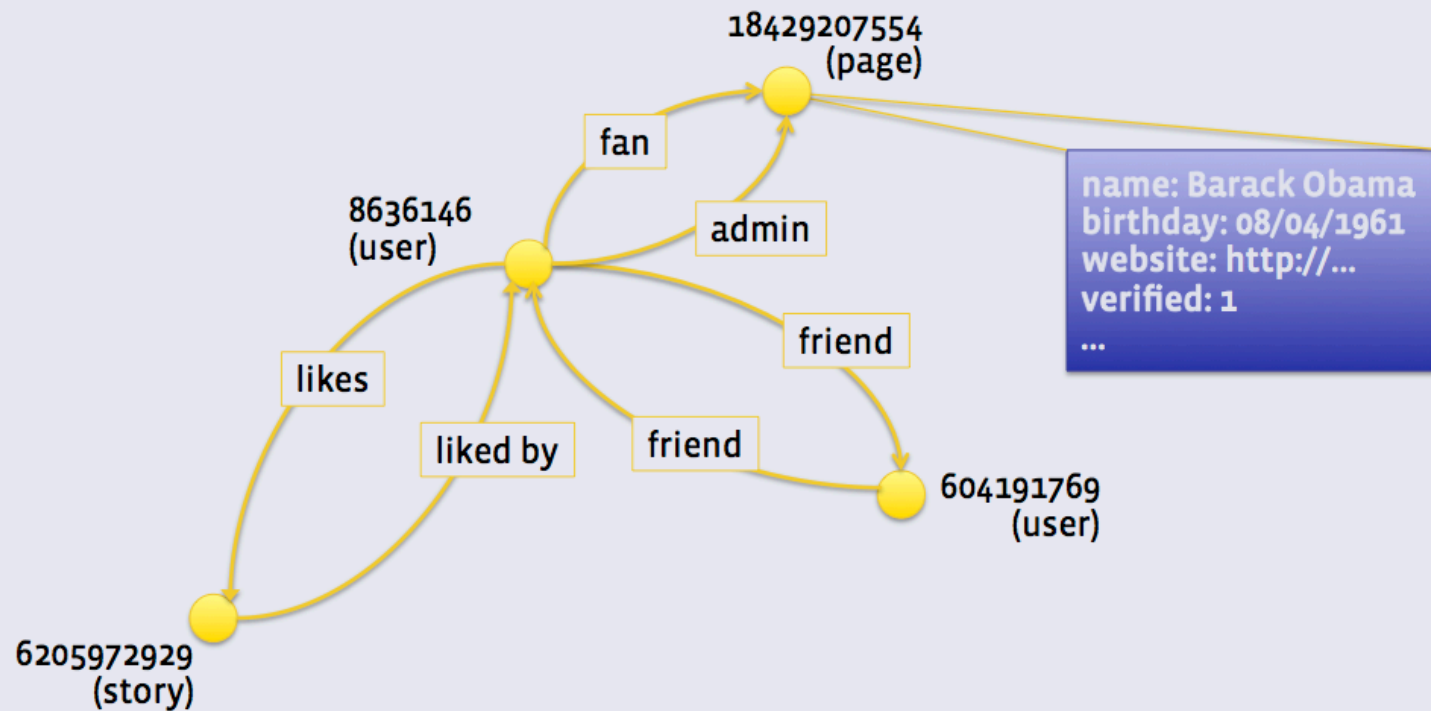
Interests

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

# Facebook Graph

## Data model

### Objects & Associations



# Examples of Graphs in Machine Learning



# Label a Face and Propagate



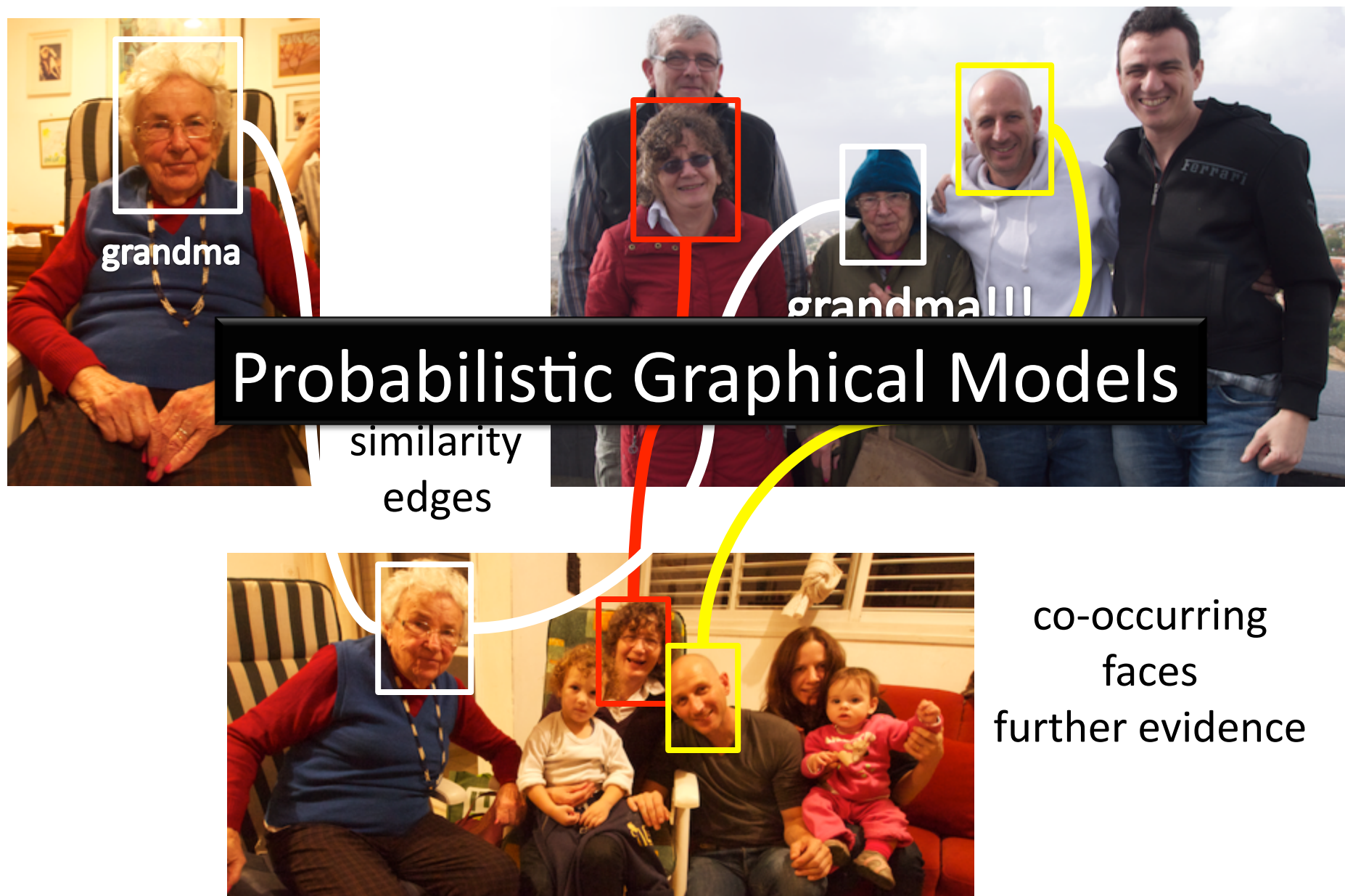
# Pairwise similarity not enough...



Not similar enough  
to be sure



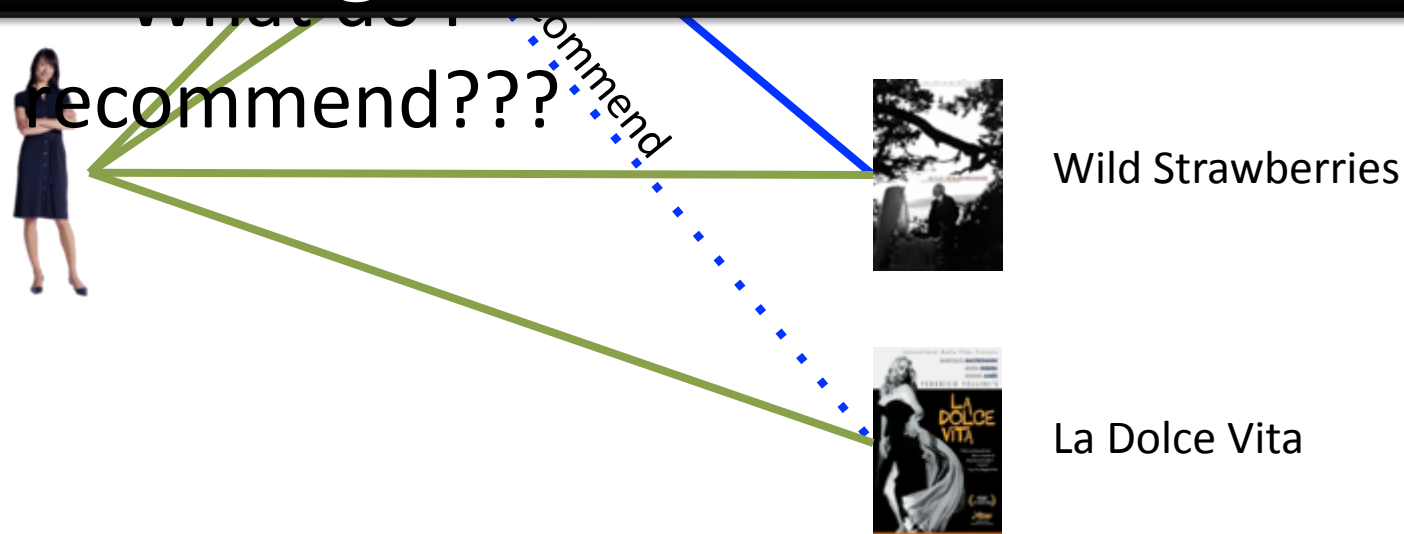
# Propagate Similarities & Co-occurrences for Accurate Predictions



# Collaborative Filtering: Exploiting Dependencies

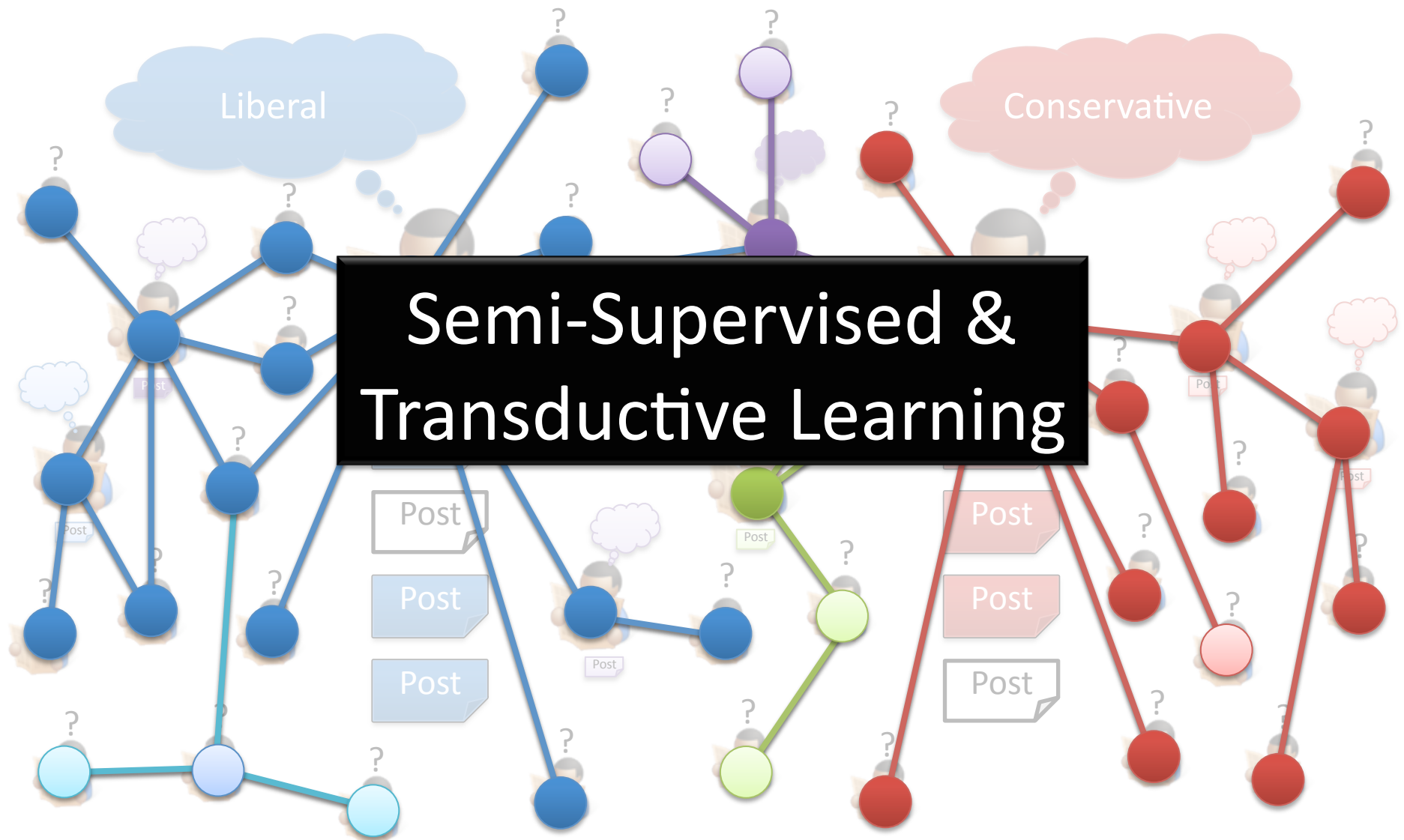


## Latent Factor Models Non-negative Matrix Factorization

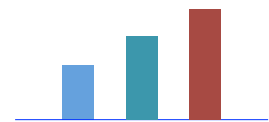
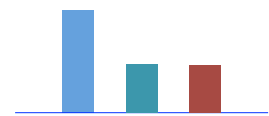
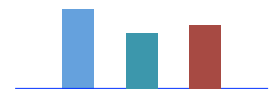
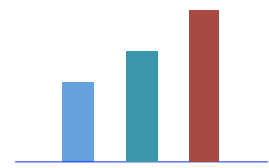
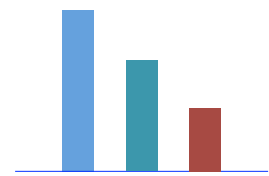




# *Estimate Political Bias*



# Topic Modeling



Cat

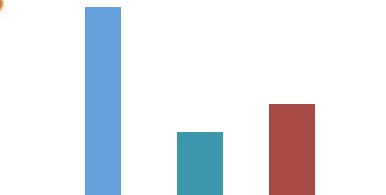
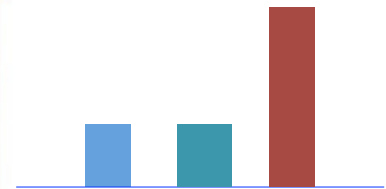
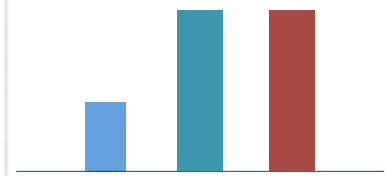
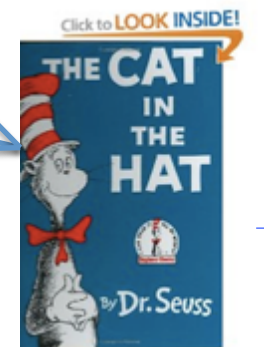
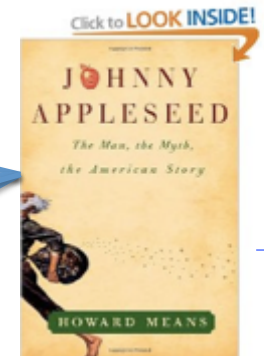
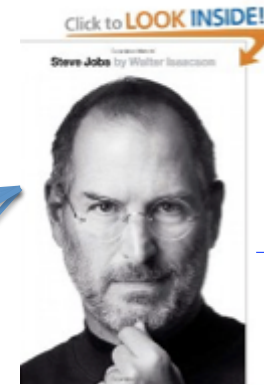
Apple

Growth

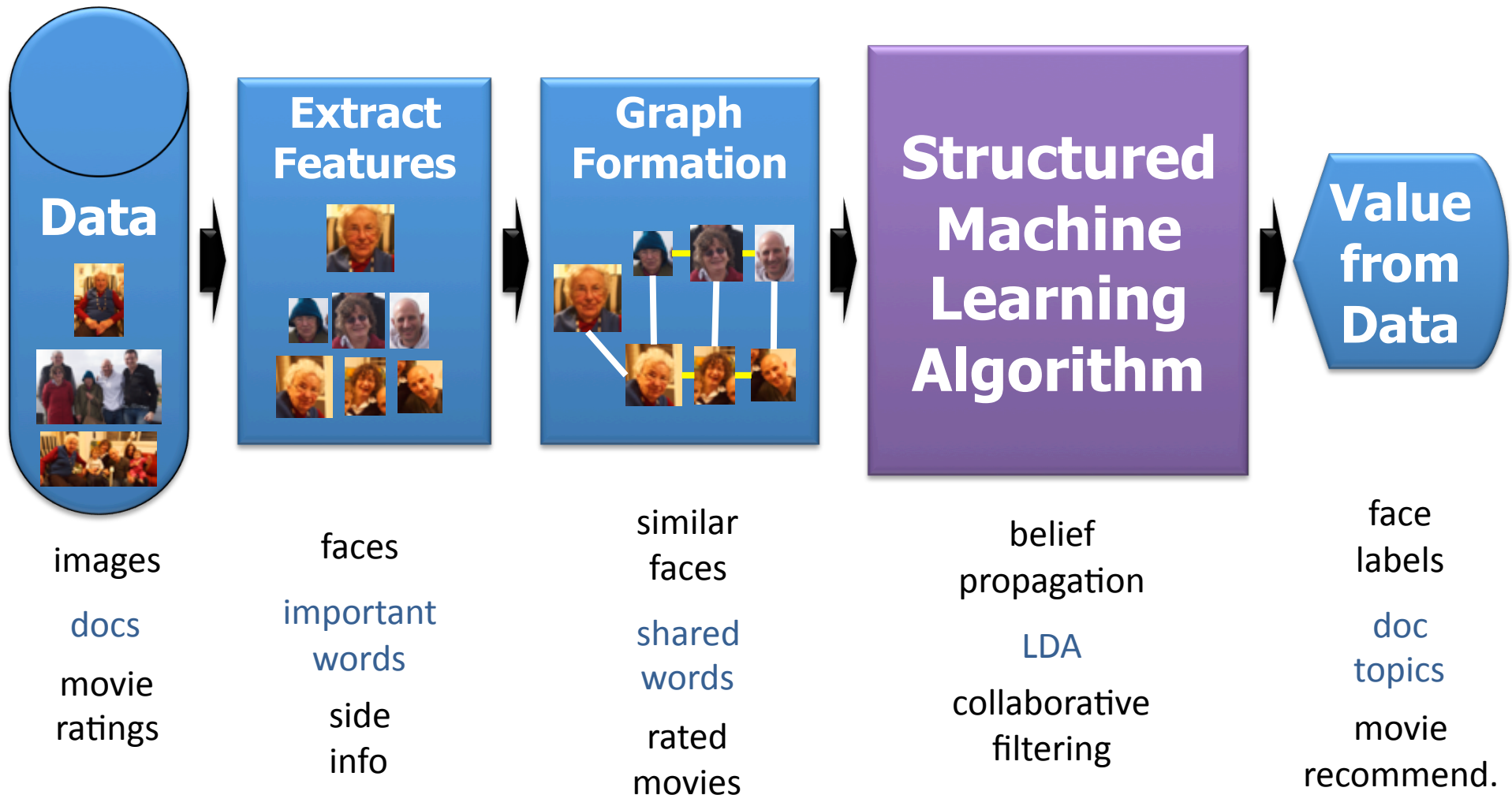
Hat

Plant

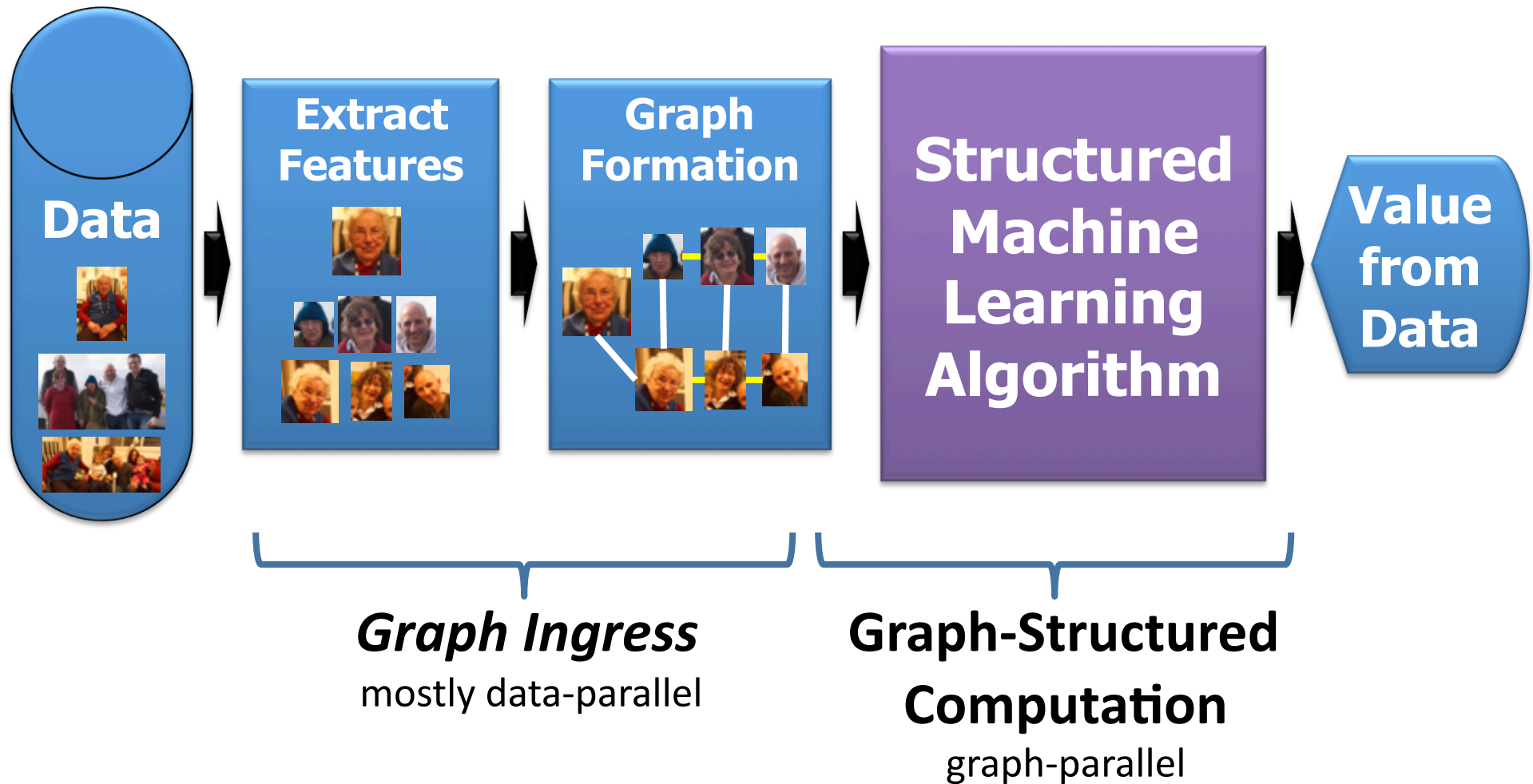
LDA and co.



# Machine Learning Pipeline



# Parallelizing Machine Learning





# 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

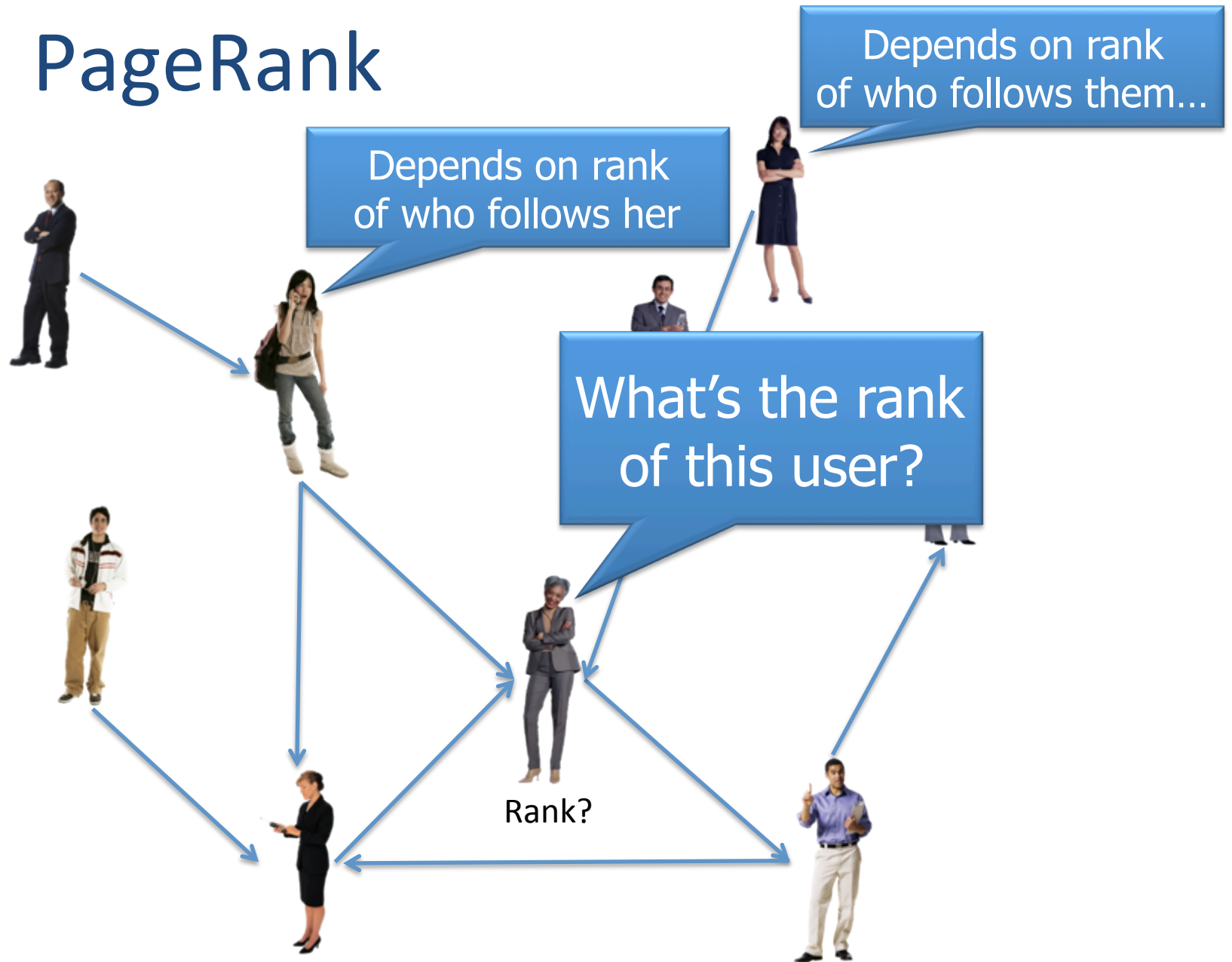
Label Propagation  
CoEM

### Graph Analysis

PageRank  
Triangle Counting

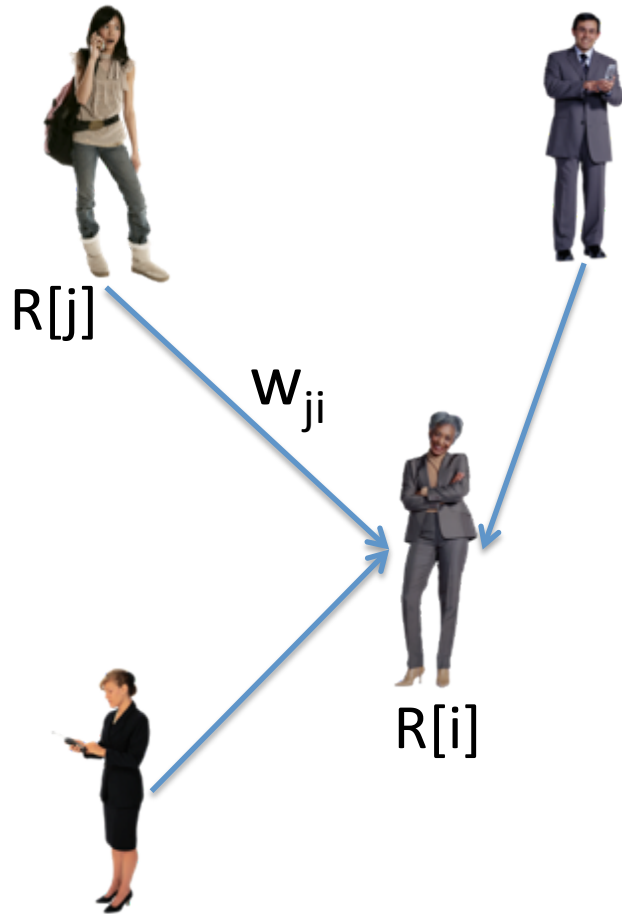
# Example of a Graph-Parallel Algorithm

# PageRank



**Loops in graph → Must iterate!**

# PageRank Iteration



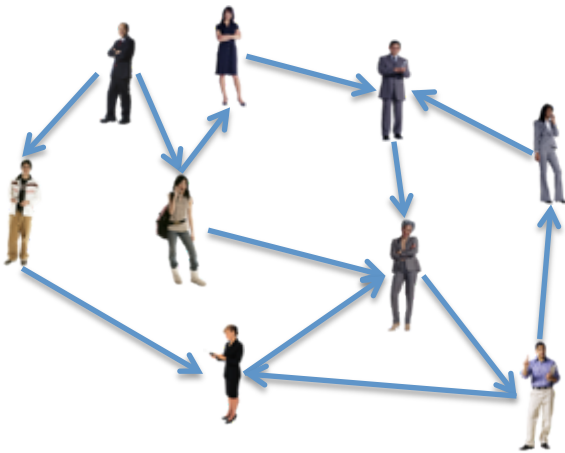
Iterate until convergence:  
“My rank is weighted  
average of my friends’ ranks”

$$R[i] = \alpha + (1 - \alpha) \sum_{(j,i) \in E} w_{ji} R[j]$$

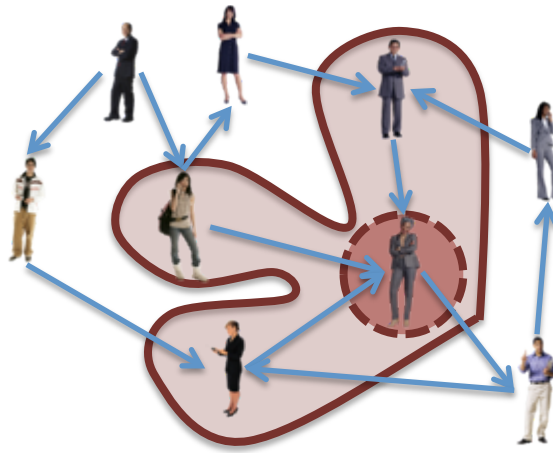
- $\alpha$  is the random reset probability
- $w_{ji}$  is the prob. transitioning (similarity) from  $j$  to  $i$

# Properties of Graph Parallel Algorithms

Dependency Graph



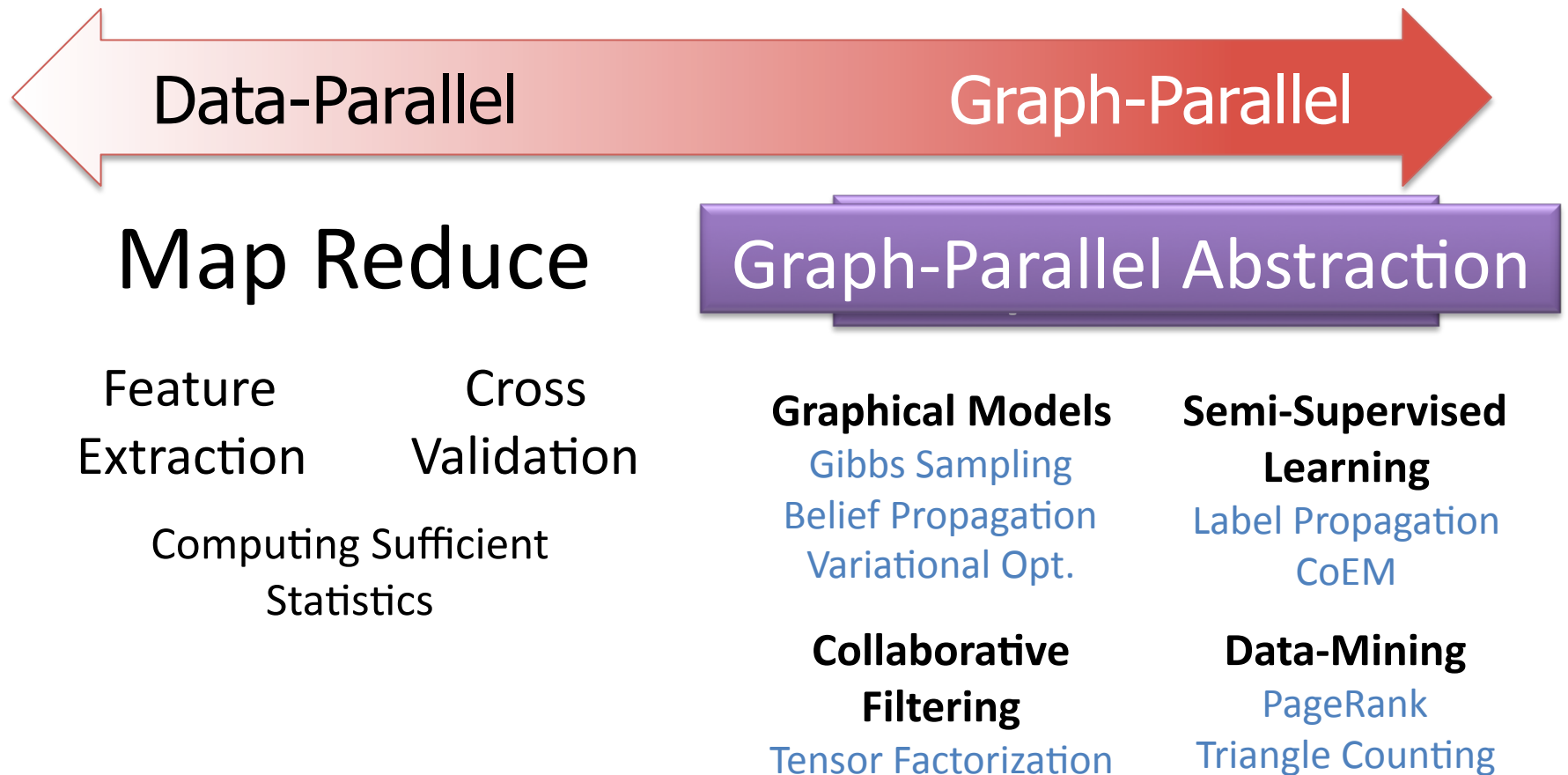
Local Updates



Iterative Computation



# Addressing Graph-Parallel ML



Graph Computation:

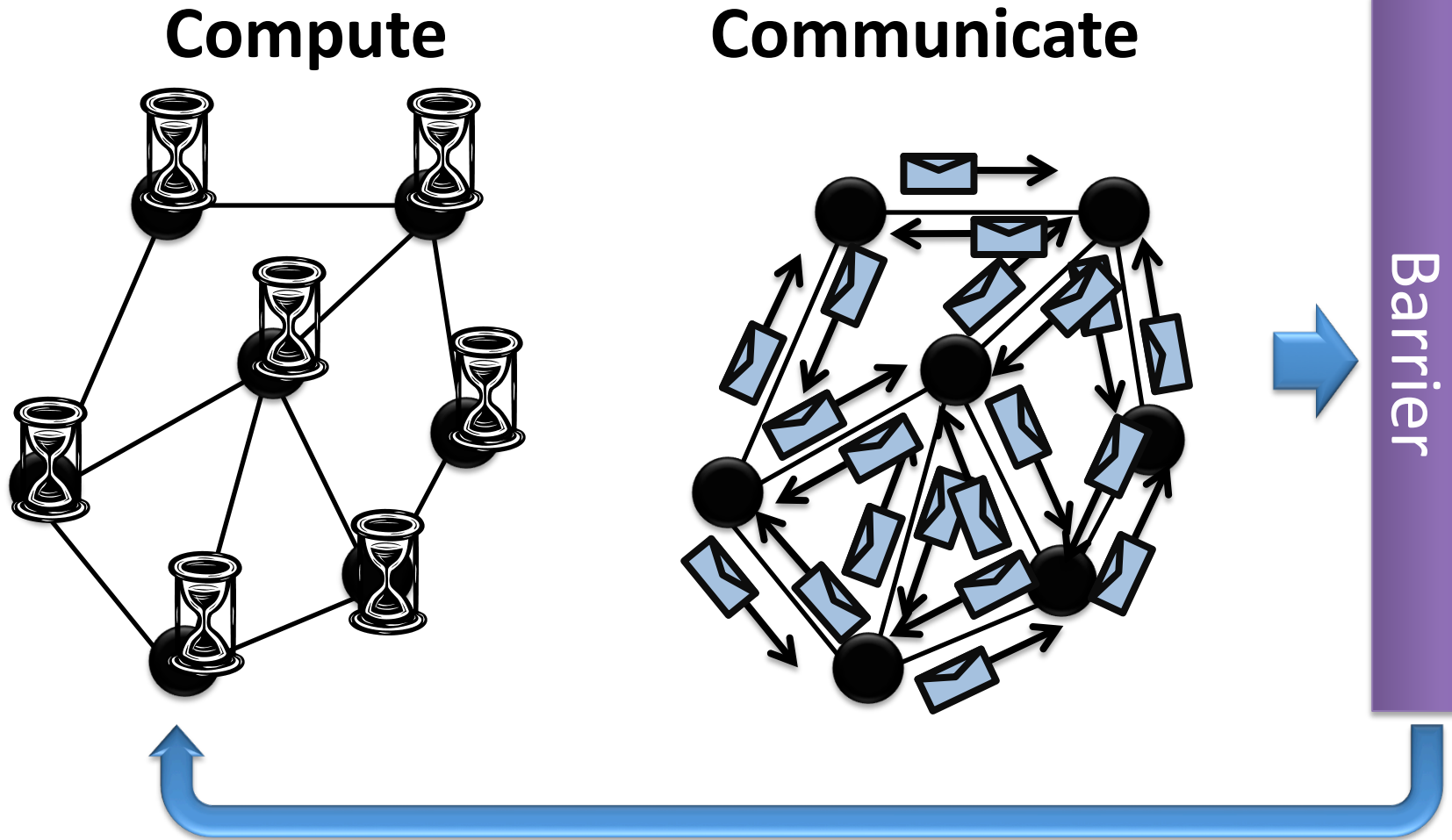
*Synchronous*

*v.*

*Asynchronous*

# Bulk Synchronous Parallel Model: Pregel (Giraph)

[Valiant '90]

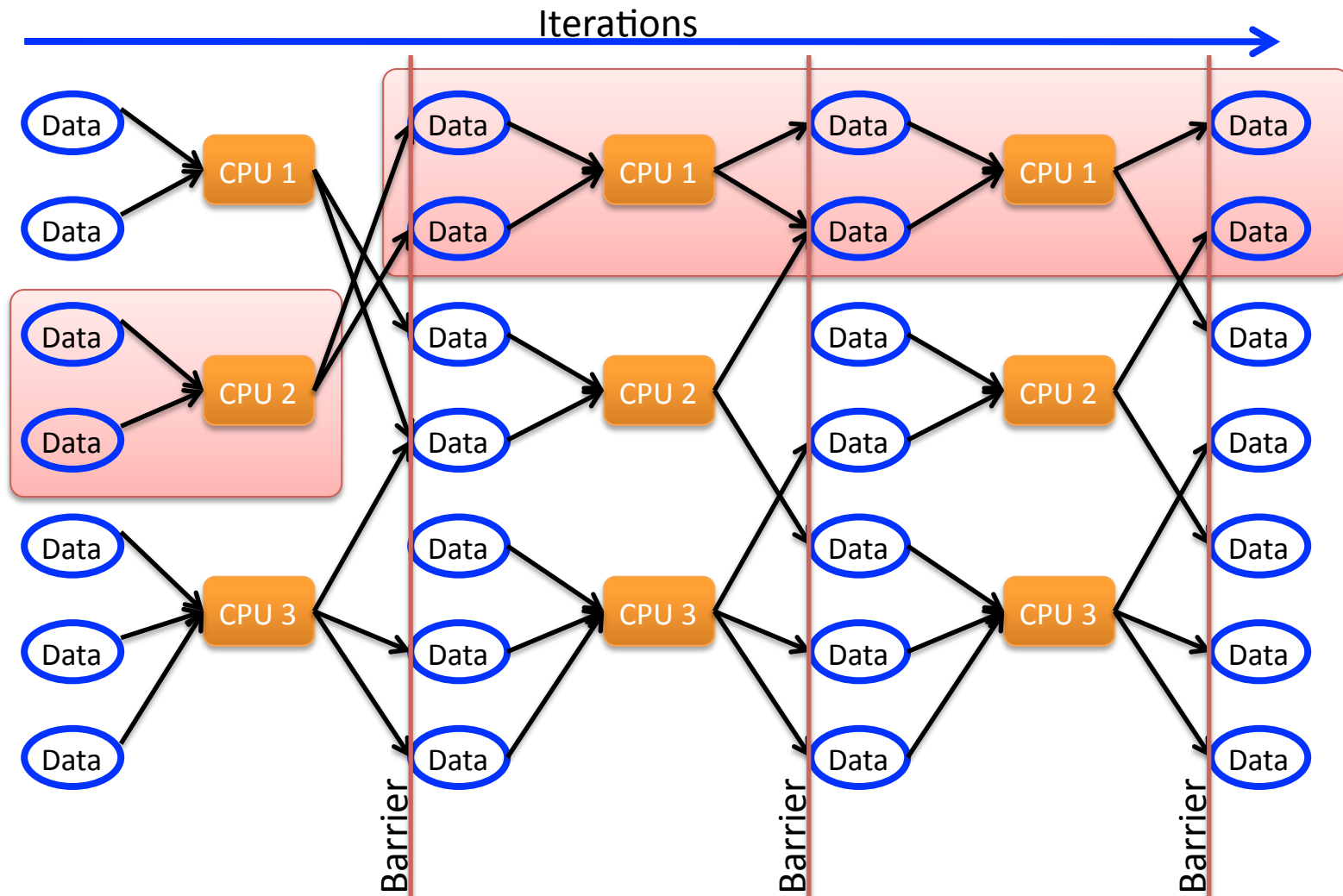




Problem:

*Bulk synchronous  
parallel systems can  
be highly inefficient*

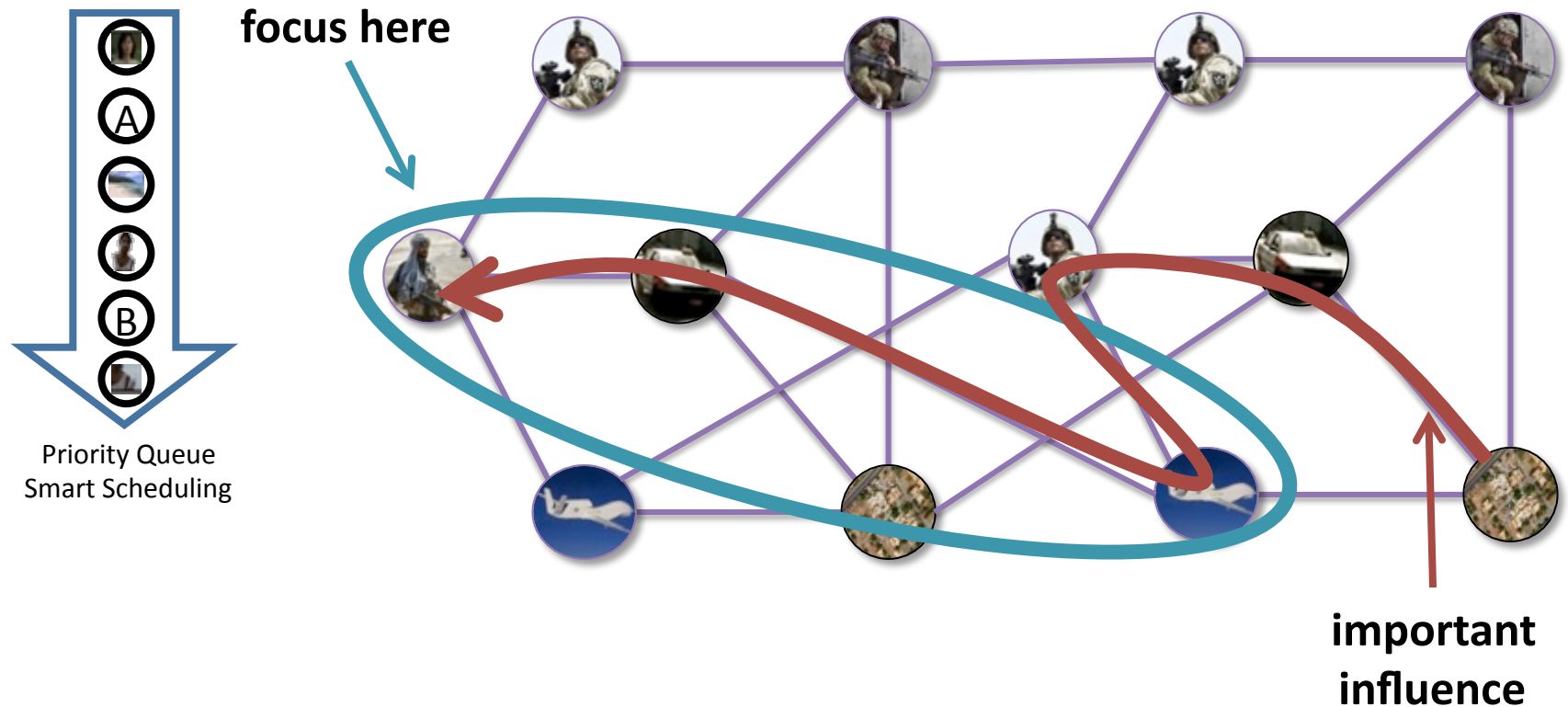
# BSP Systems Problem: Curse of the Slow Job



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

# Analyzing Belief Propagation

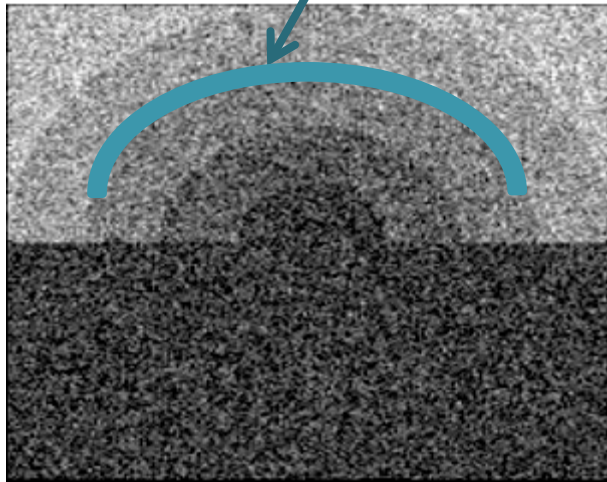
[Gonzalez, Low, G. '09]



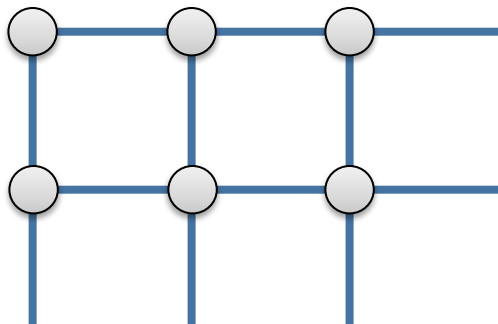
Asynchronous Parallel Model (rather than BSP)  
fundamental for efficiency

# Asynchronous Belief Propagation

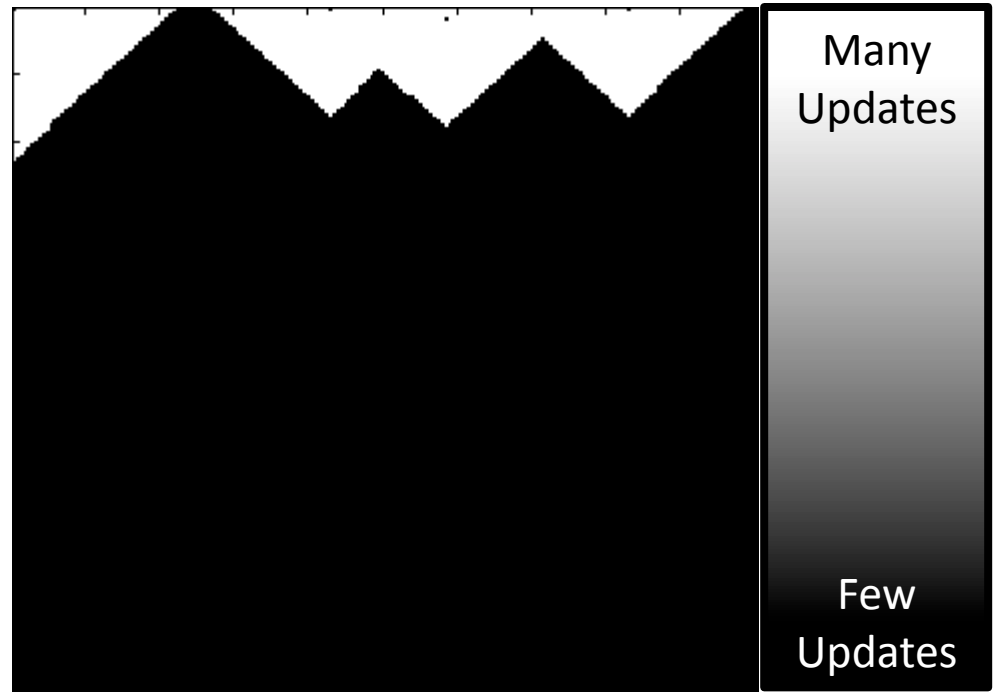
Challenge = Boundaries



Synthetic Noisy Image



Graphical Model



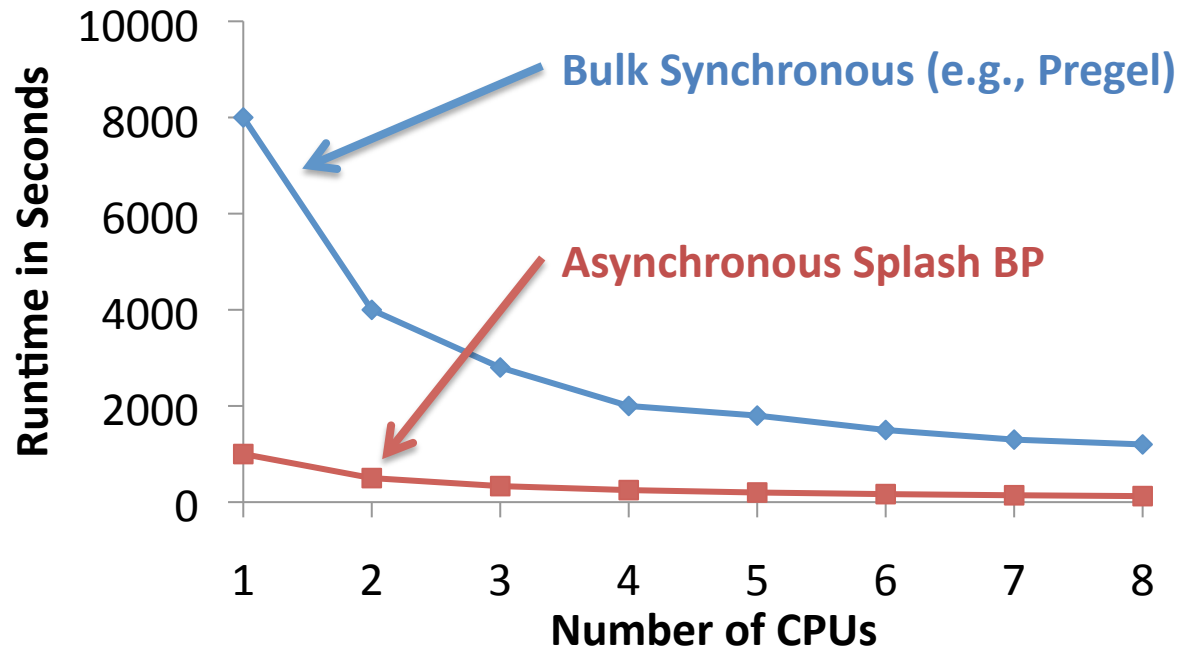
Cumulative Vertex Updates

Algorithm identifies and focuses on hidden sequential structure

# BSP ML Problem:

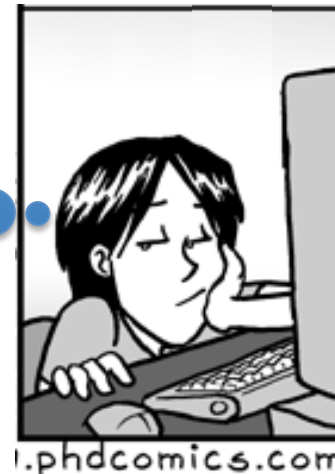
## Synchronous Algorithms can be **Inefficient**

[Gonzalez, Low, G. '09]



**Theorem:**  
Bulk Synchronous BP  
 $O(\#vertices)$  slower  
than Asynchronous BP

Efficient parallel  
implementation was  
painful, painful, painful...



# The Need for a New Abstraction

- Need: Asynchronous, Dynamic Parallel Computations



## 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

### Data-Mining

PageRank  
Triangle Counting

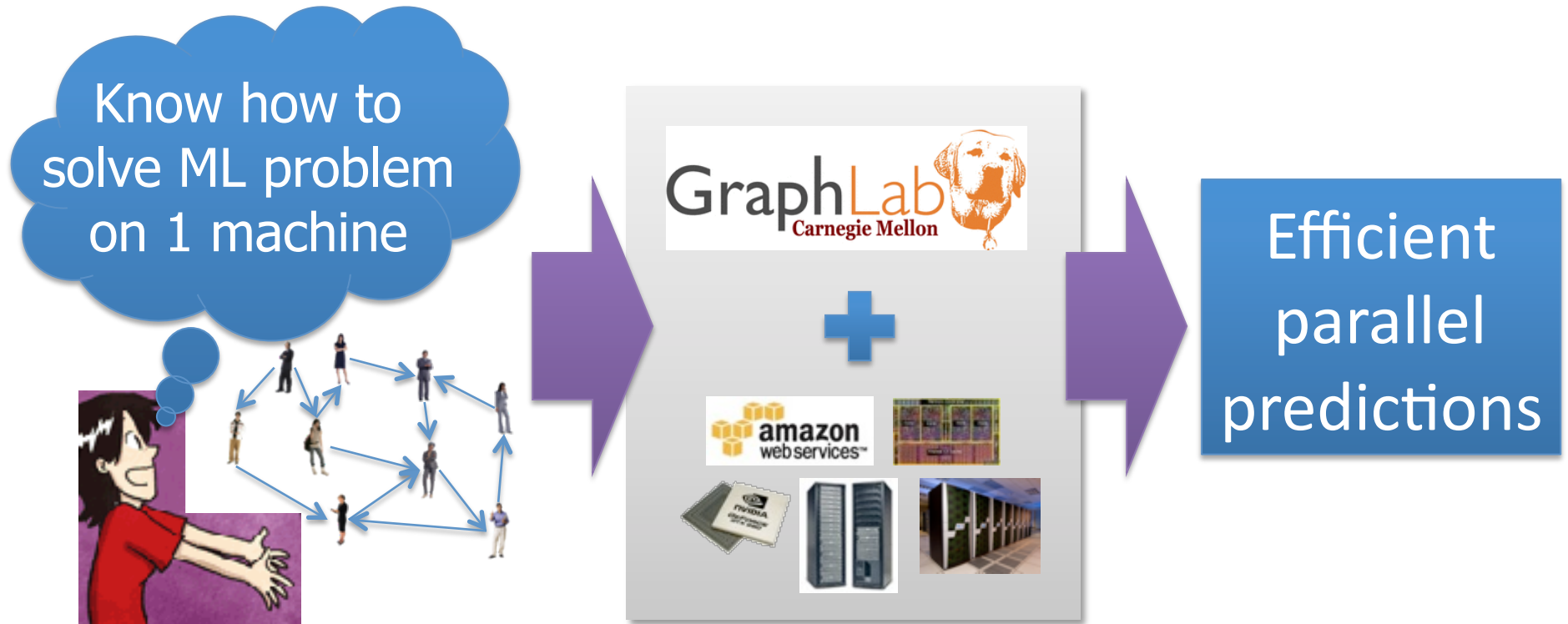
# The GraphLab Goals

- Designed specifically for ML
  - Graph dependencies
  - Iterative
  - Asynchronous
  - Dynamic
- Simplifies design of parallel programs:
  - Abstract away hardware issues
  - Automatic data synchronization
  - Addresses multiple hardware architectures





# The GraphLab Goals



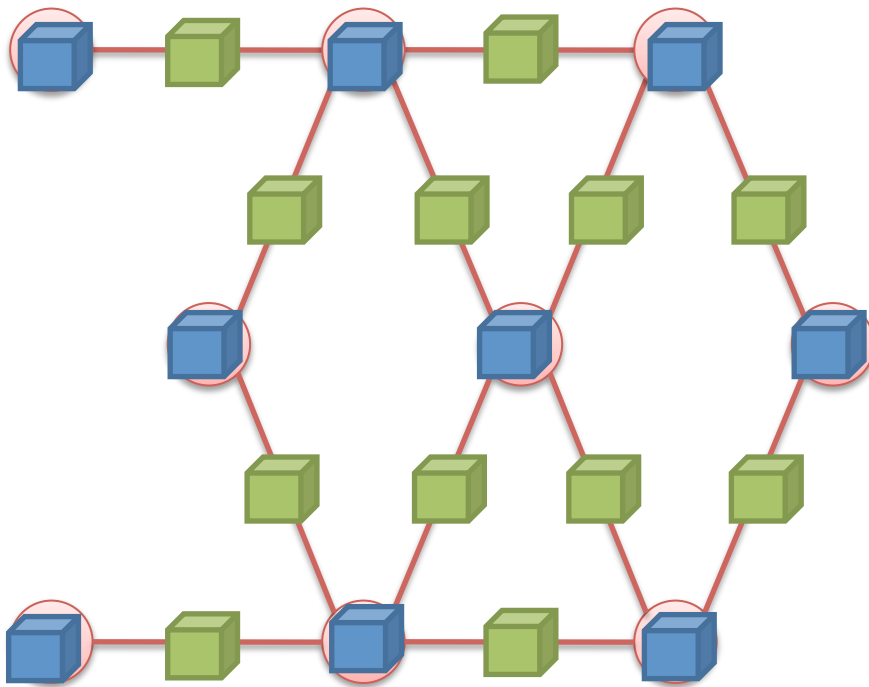


POSSIBILITY



# Data Graph

Data associated with vertices and edges



Graph: 

- Social Network

Vertex Data: 

- User profile text
- Current interests estimates

Edge Data: 

- Similarity weights

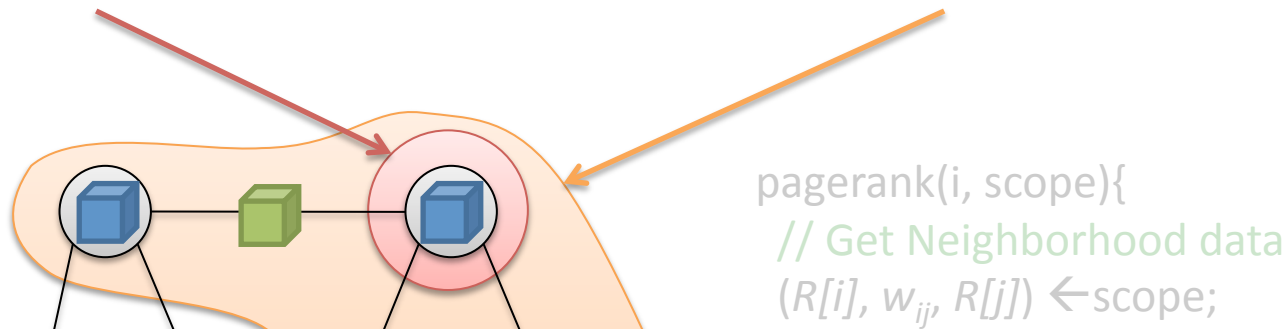
How do we *program*  
graph computation?

“Think like a Vertex.”

-Malewicz et al. [SIGMOD'10]

# Update Functions

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



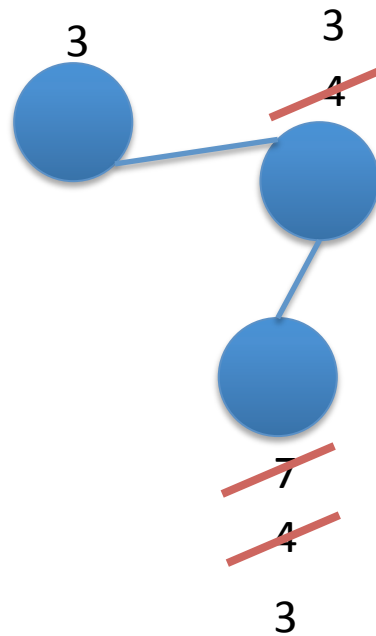
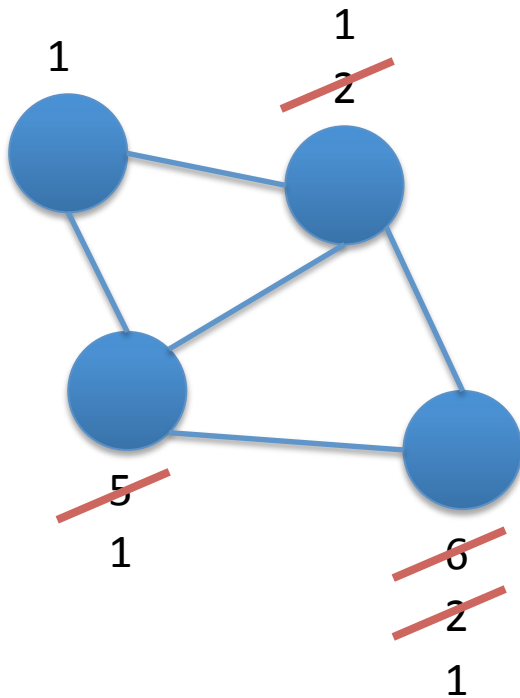
Update function applied (asynchronously)  
in parallel until convergence

Many schedulers available to prioritize computation



Dynamic  
computation

# Update Function Example: Connected Components



## Initialize:

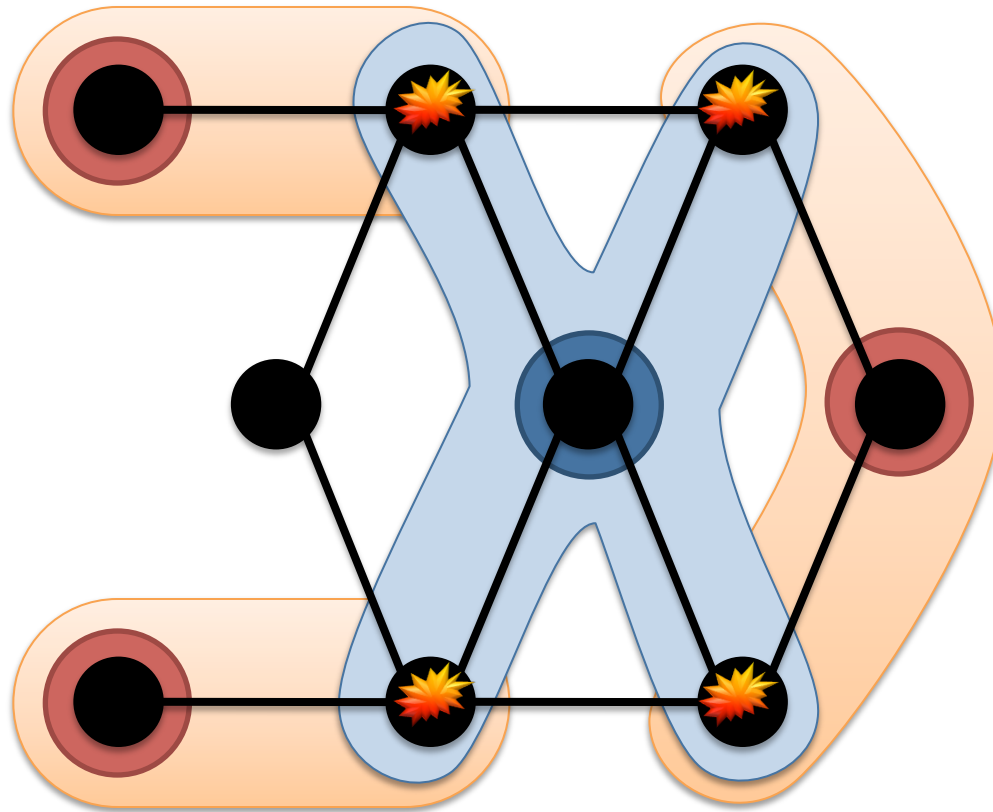
Assign component id  
to vertex id

## Update(v):

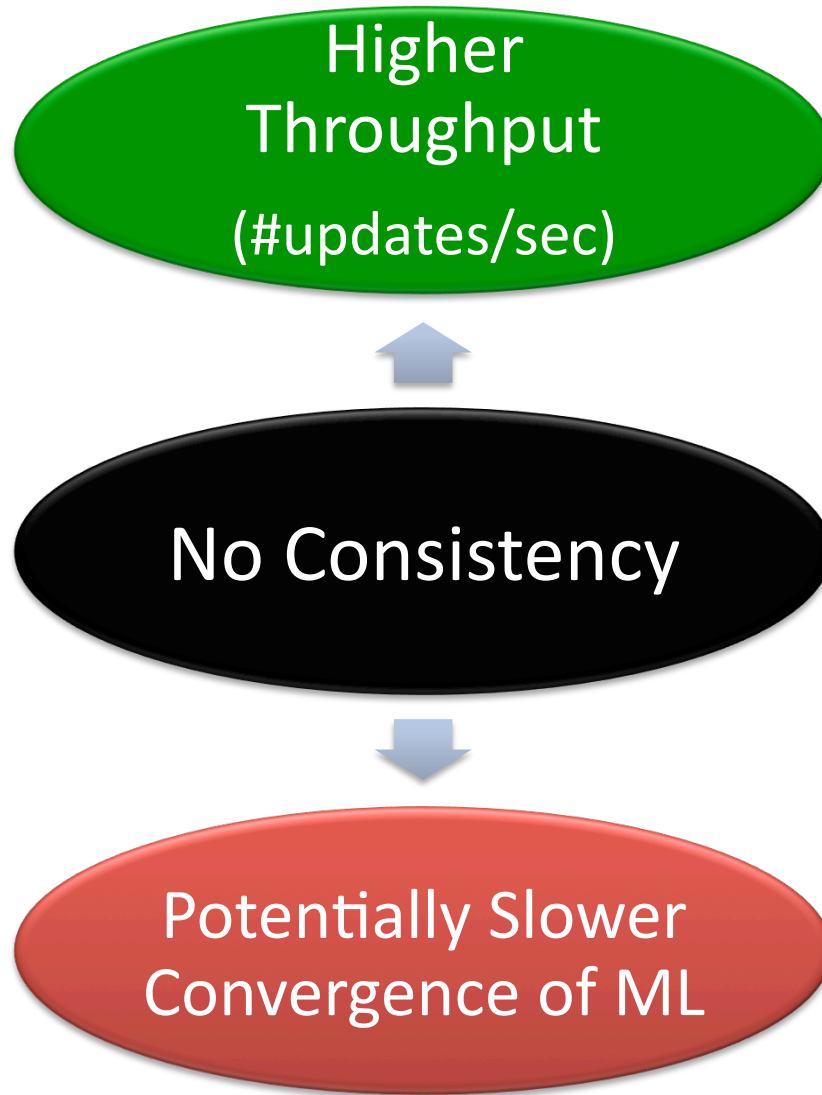
$v.component =$   
 $\min(\text{self \&}$   
 $\text{neighbor components})$

# Ensuring Race-Free Code

How much can computation **overlap**?

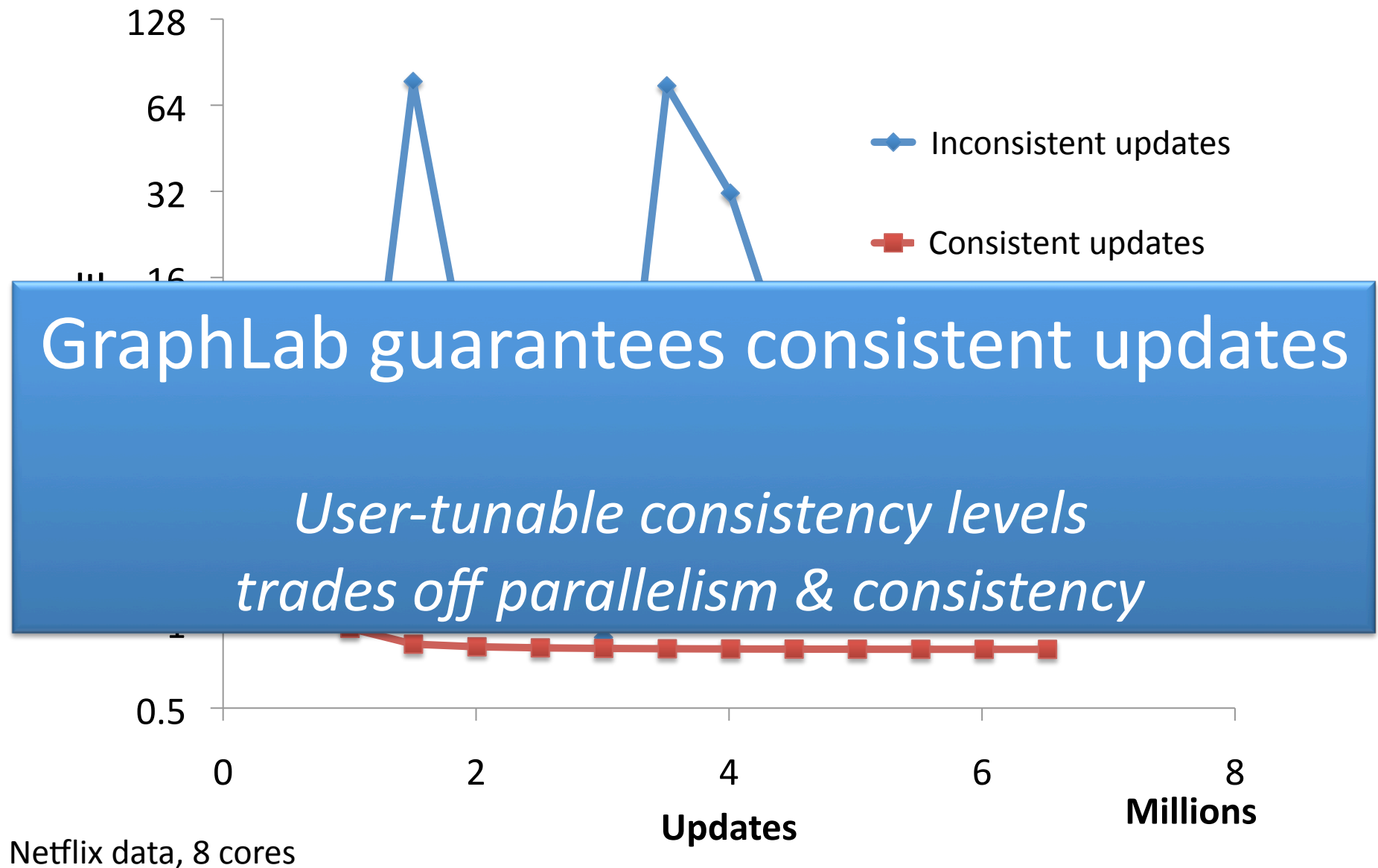


# Need for Consistency?





# Consistency in Collaborative Filtering

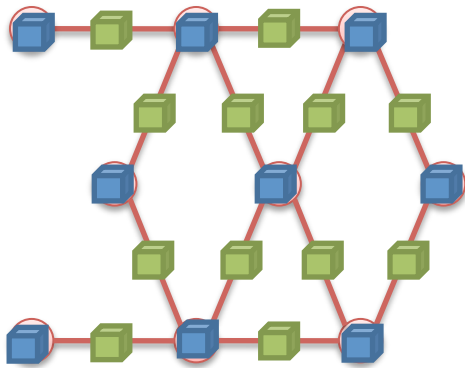


Fix text

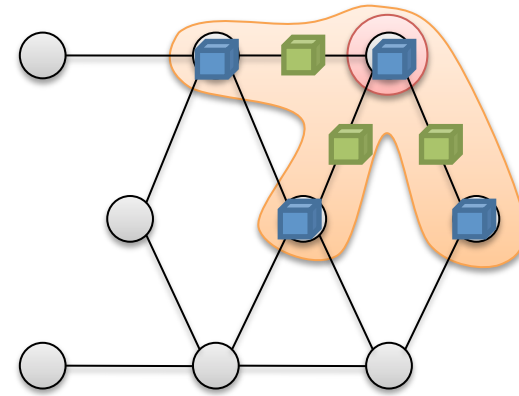
**MORE SLIDES ABOUT  
CONSISTENCY???**

# The GraphLab Framework

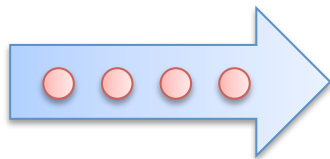
Graph Based  
*Data Representation*



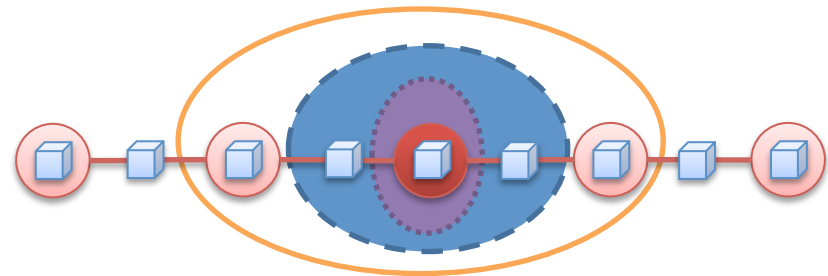
Update Functions  
*User Computation*



Scheduler



Consistency Model



Alternating Least  
Squares

SVD

Splash Sampler

CoEM

Bayesian Tensor  
Factorization

Lasso

Belief Propagation

PageRank

LDA



SVM

Gibbs Sampling

Dynamic Block Gibbs Sampling

K-Means

**...Many others...**

Matrix  
Factorization

Linear Solvers

# Never Ending Learner Project (CoEM)

Hadoop	95 Cores	7.5 hrs
<b>Distributed GraphLab</b>	<b>32 EC2 machines</b>	<b>80 secs</b>

**0.3% of Hadoop time**

2 orders of mag faster →  
2 orders of mag cheaper



- ML algorithms as vertex programs
- Asynchronous execution and consistency models

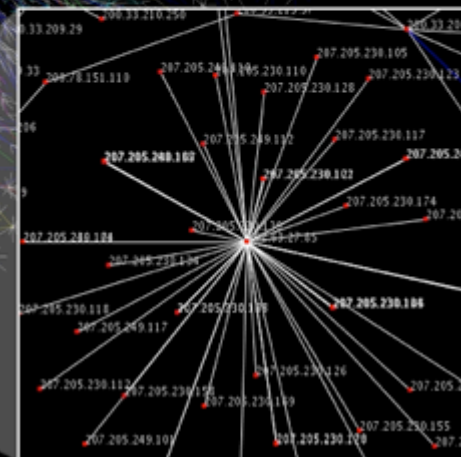
Thus far...

GraphLab 1 provided exciting  
scaling performance

But...

**We couldn't scale up to  
Altavista Webgraph 2002  
1.4B vertices, 6.7B edges**

# Natural Graphs



[Image from WikiCommons]

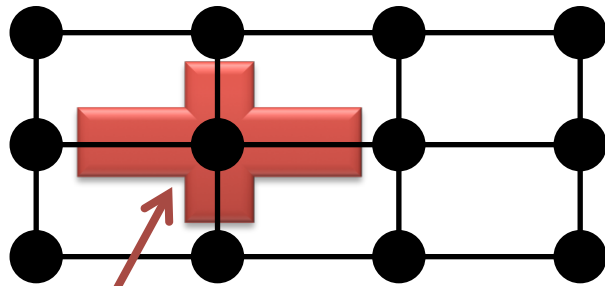


# Problem:

Existing *distributed* graph  
computation systems perform  
poorly on **Natural Graphs**

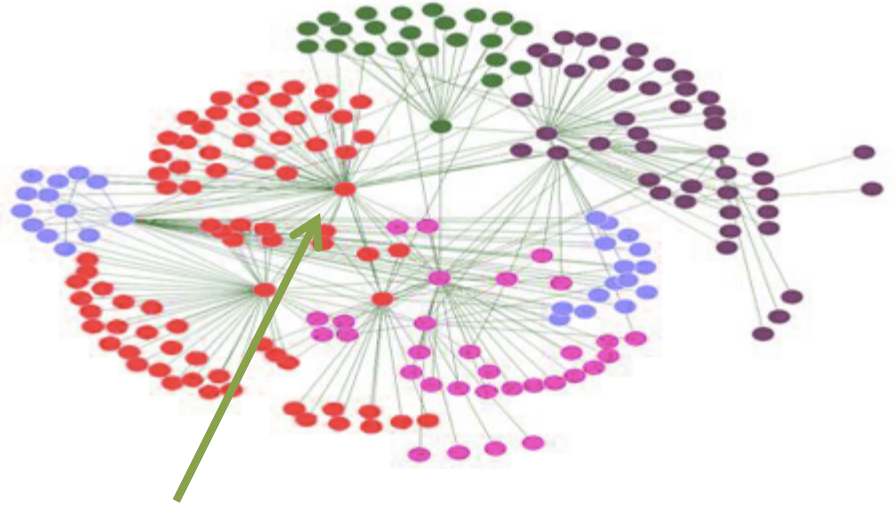
# *Achilles Heel:* Idealized Graph Assumption

Assumed...



Small degree →  
Easy to partition

But, Natural Graphs...

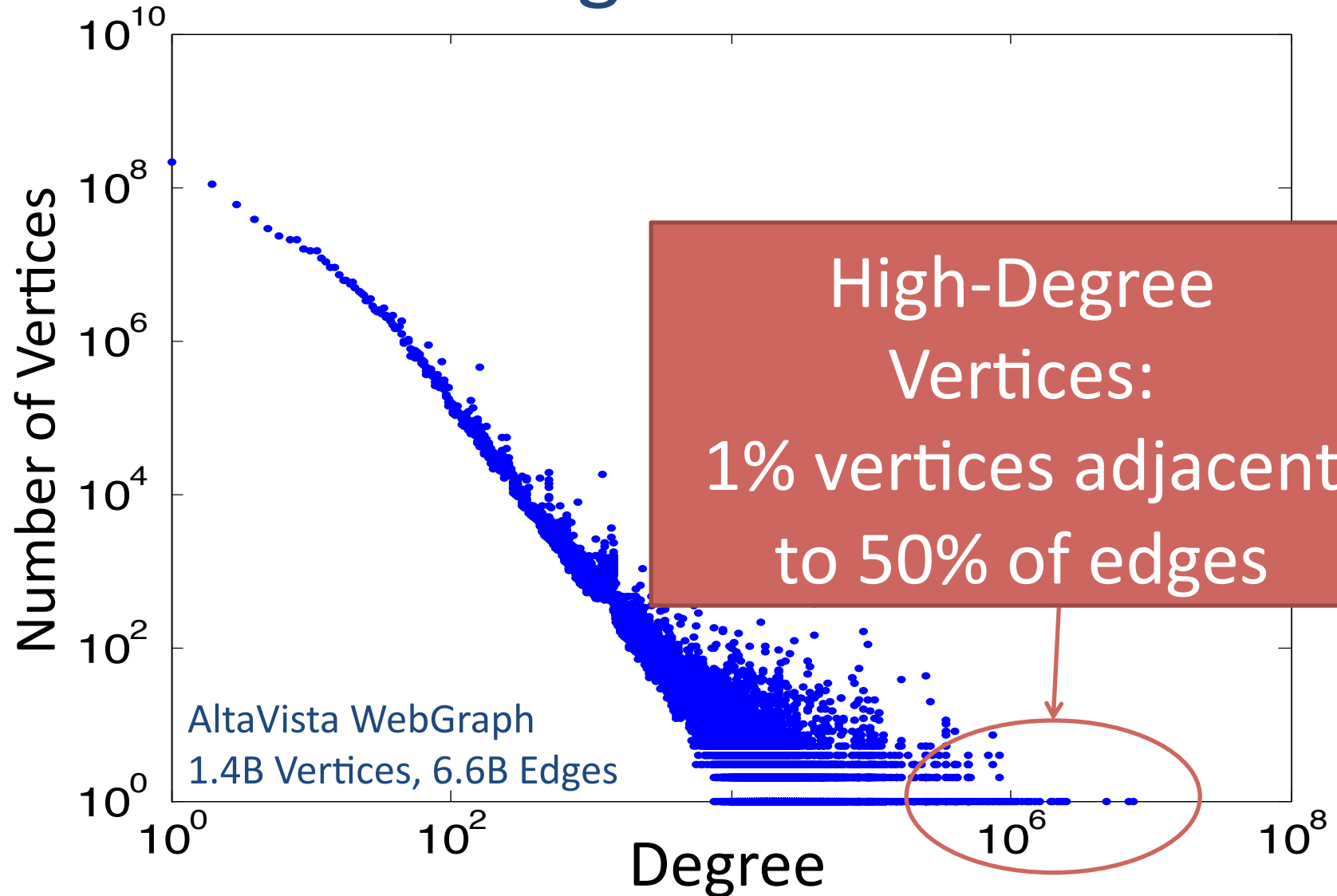


Many high degree vertices  
(power-law degree distribution)



Very hard to partition

# Power-Law Degree Distribution

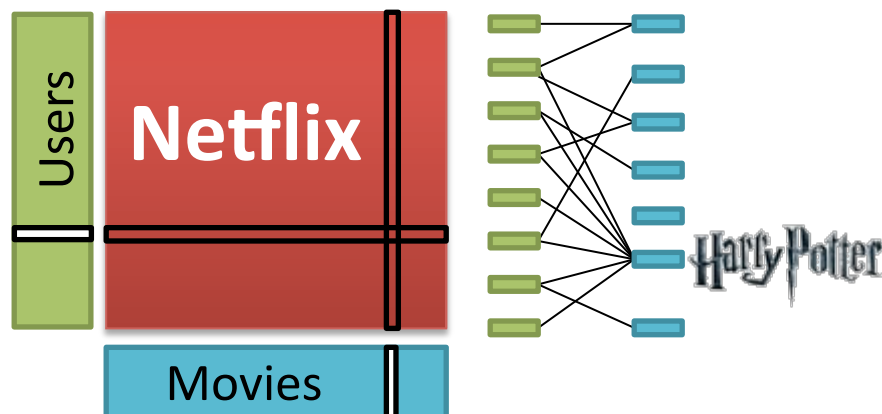


# High Degree Vertices are Common

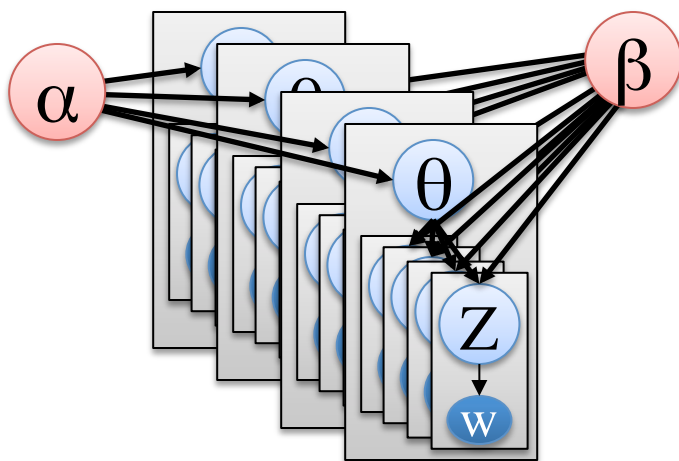
“Social” People



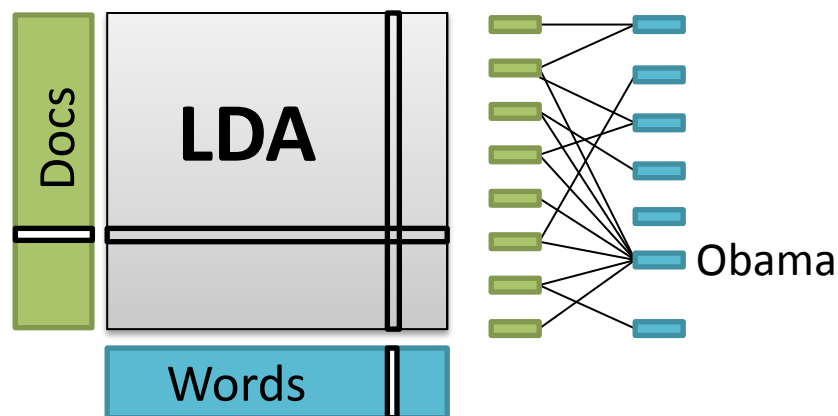
Popular Movies



Hyper Parameters

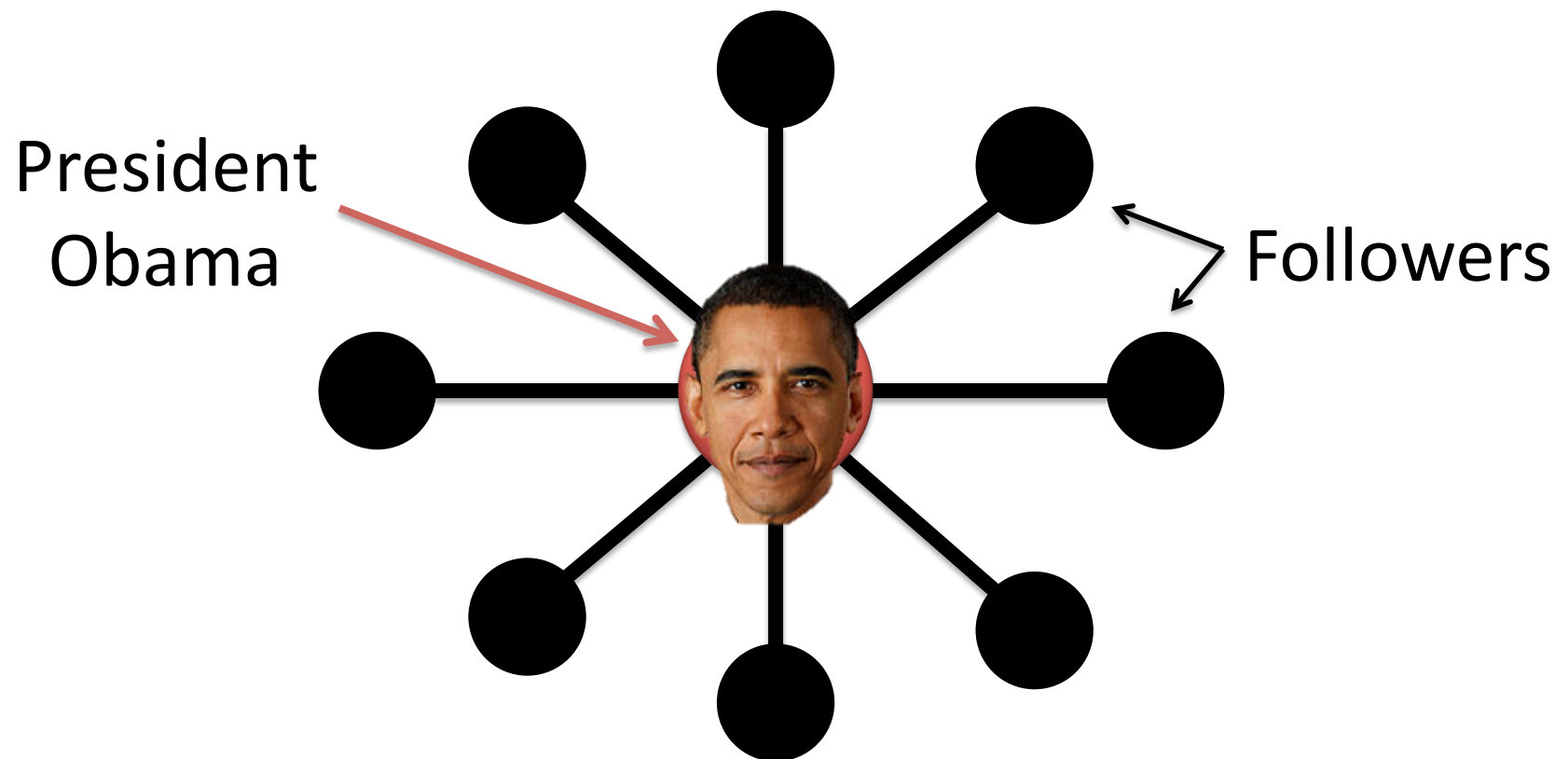


Common Words



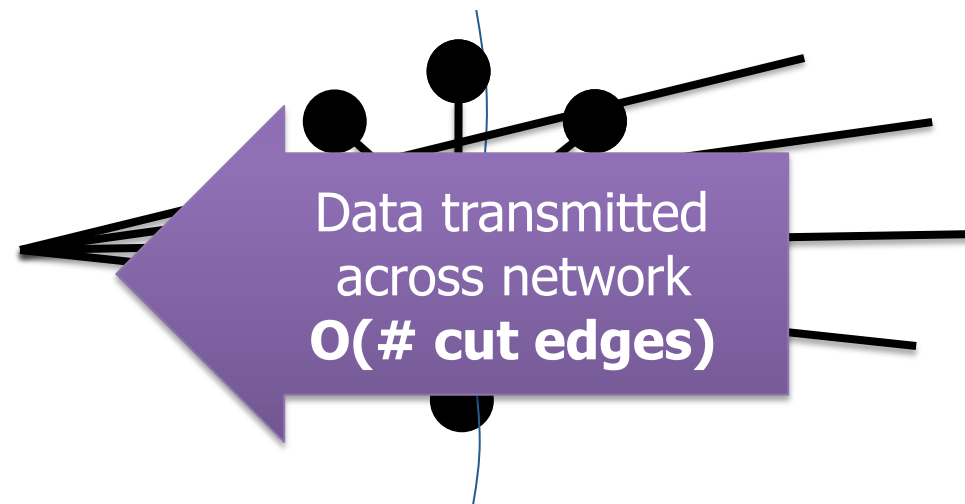
# Power-Law Degree Distribution

## “Star Like” Motif



Problem:

**High Degree Vertices → High  
Communication for Distributed Updates**



Natural graphs do not have low-cost balanced cuts

*[Leskovec et al. 08, Lang 04]*

Popular partitioning tools (Metis, Chaco,...) perform poorly

*[Abou-Rjeili et al. 06]*

*Extremely slow and require substantial memory*

# Random Partitioning

- Both GraphLab 1, Pregel, Twitter, Facebook,... rely on Random (hashed) partitioning for Natural Graphs

For  $p$  Machines:

$$\mathbb{E} \left[ \frac{|Edges\ Cut|}{|E|} \right] = 1 - \frac{1}{p}$$

...

...

**10 Machines → 90% of edges cut**

**100 Machines → 99% of edges cut!**

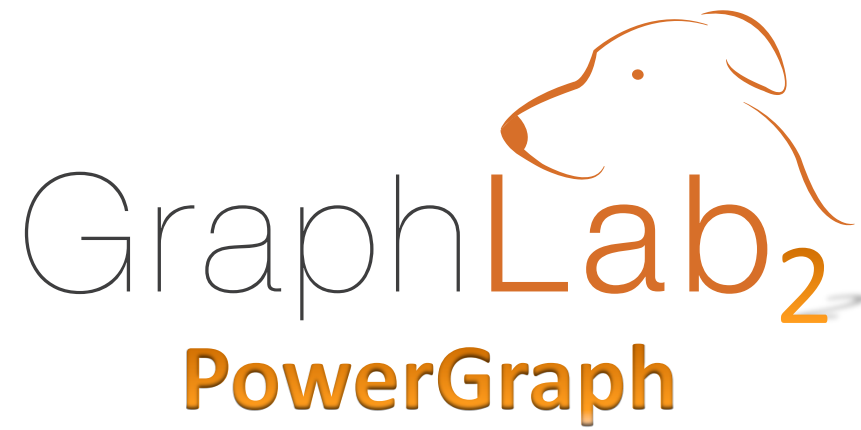
All data is communicated... Little advantage over MapReduce

# In Summary

**GraphLab 1 and Pregel are not well suited for natural graphs**

- Poor performance on high-degree vertices
- Low Quality Partitioning

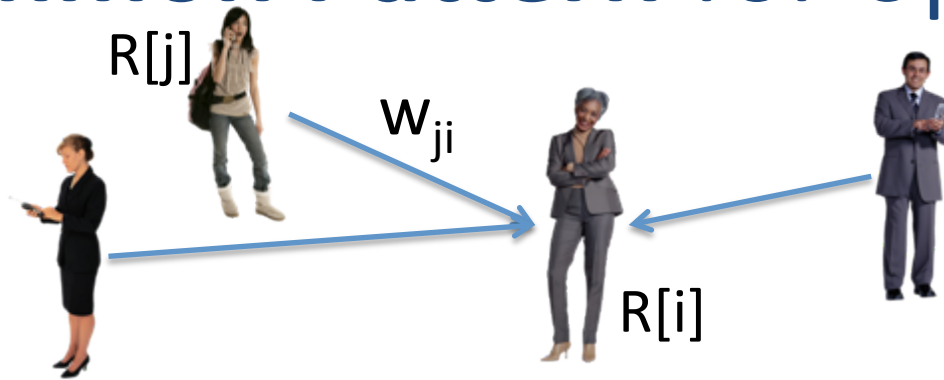




## SCALABILITY



# Common Pattern for Update Fncs.



GraphLab\_PageRank(i)

```
// Compute sum over neighbors
total = 0
foreach( j in in_neighbors(i)):
    total = total + R[j] *  $w_{ji}$ 
```

***Gather* Information  
About Neighborhood**

```
// Update the PageRank
R[i] = 0.1 + total
```

***Apply* Update to Vertex**

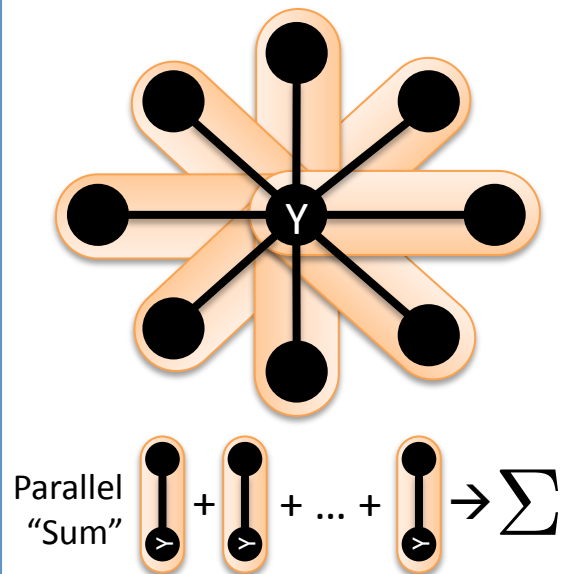
```
// Trigger neighbors to run again
if R[i] not converged then
    foreach( j in out_neighbors(i))
        signal vertex-program on j
```

***Scatter* Signal to Neighbors  
& Modify Edge Data**

# GAS Decomposition

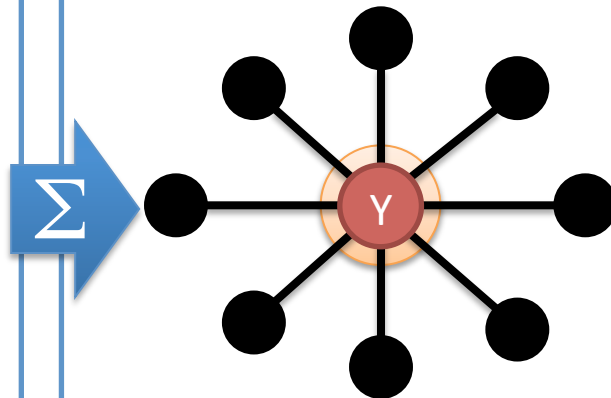
## **G**ather (Reduce)

Accumulate information  
about neighborhood



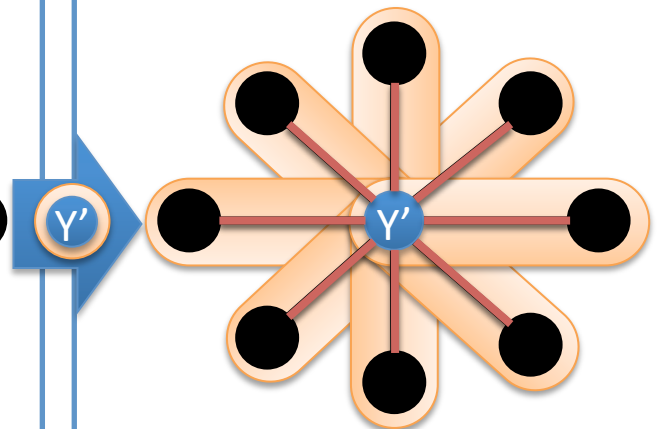
## **A**pply

Apply the accumulated  
value to center vertex



## **S**catter

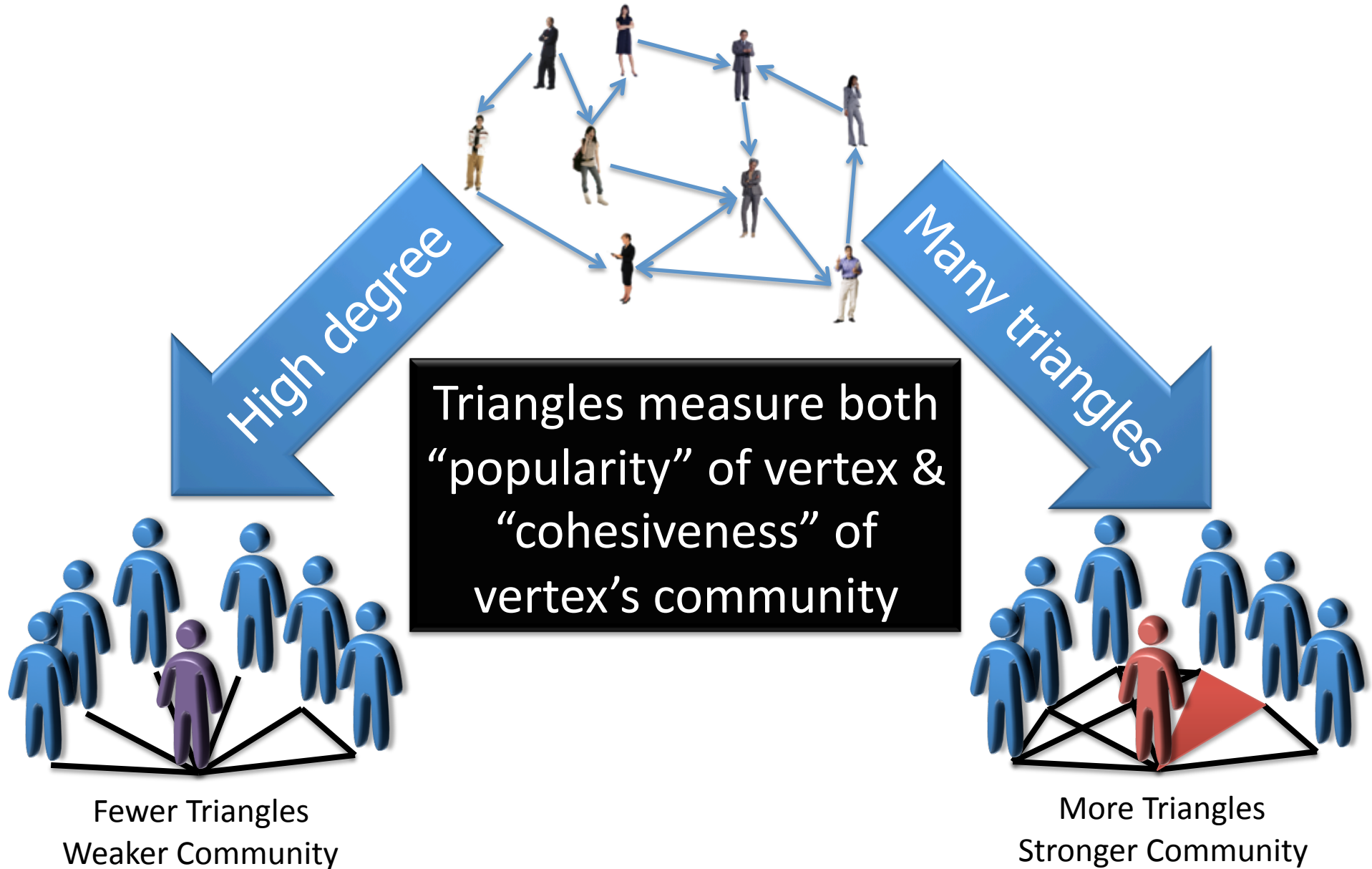
Update adjacent edges  
and vertices.



# Many ML Algorithms fit into GAS Model

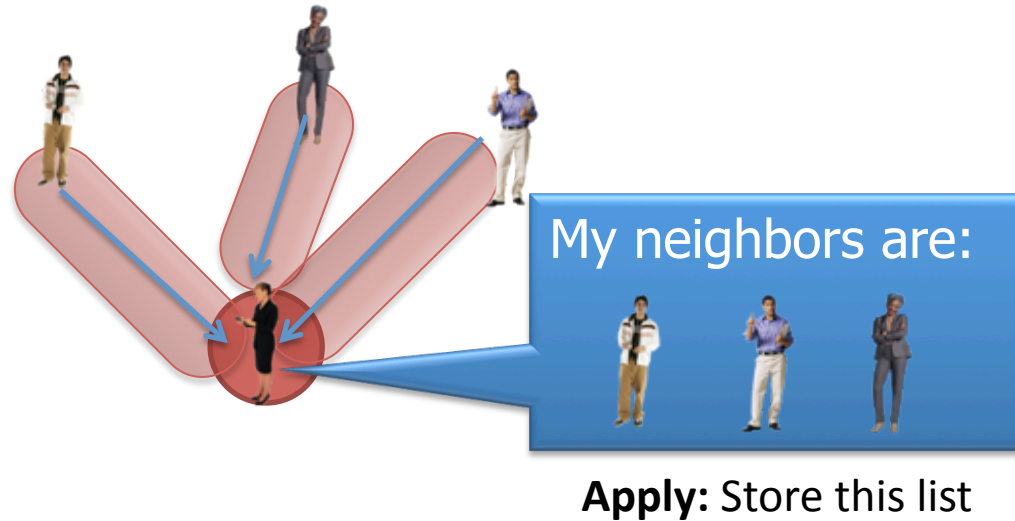
graph analytics, inference in graphical  
models, matrix factorization,  
collaborative filtering, clustering, LDA, ...

# Discovering *Influencers* in Social Networks

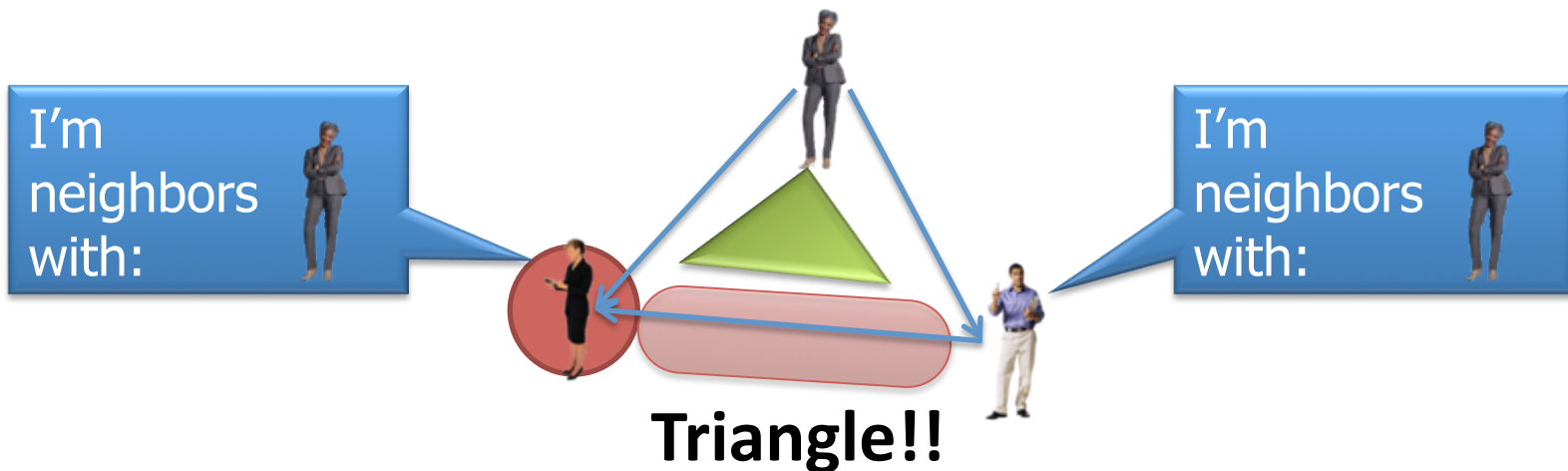


# Gather/Apply/Scatter Triangle Counting

**Gather:**






















**Scatter:**



# Triangle Counting on Twitter (2010)

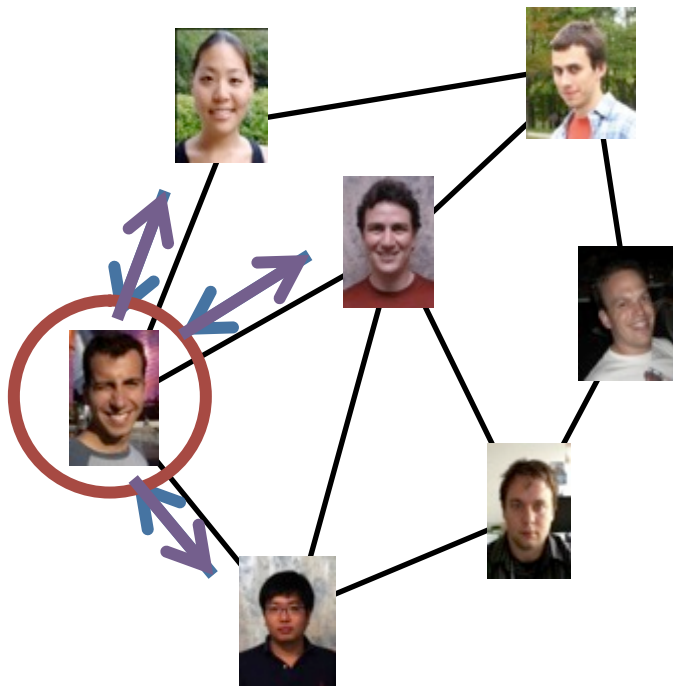
Popular People

Degree	
	<a href="#">Britney Spears</a> 3081108
	<a href="#">ashton kutcher</a> 2997653
	<a href="#">Ellen DeGeneres</a> 2679666
	<a href="#">Barack Obama</a> 2653045
	<a href="#">CNN Breaking News</a> 2450768
	<a href="#">Oprah Winfrey</a> 1994945
	<a href="#">Twitter</a> 1959765
	<a href="#">Ryan Seacrest</a> 1885917
	<a href="#">SHAQ</a> 1844123

Triangles / Following	
	<a href="#">Women's Wear Daily</a> 6449985
	<a href="#">wefollow</a> 5.962666e+06
	<a href="#">Stephen Colbert</a> 5607368
	<a href="#">Jonas Brothers</a> 5272692
	<a href="#">Rev Run</a> 4483823
	<a href="#">Defamer.com</a> 3564740
	<a href="#">You Look Great</a> 3207562
	<a href="#">Oprah Winfrey</a> 2.936561e+06
	<a href="#">Al Gore</a> 2.488950e+06
	<a href="#">CNN Breaking News</a> 2.474015e+06

Popular People  
With  
Strong  
Communities

# Factorized Belief Propagation



- **Gather:** Accumulates product of *in messages*
- **Apply:** Updates central belief
- **Scatter:** Computes out messages & schedules neighbors as needed



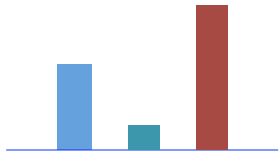
# Collaborative Filtering (via Alternating Least Squares)



# Factorized Collaborative Filtering Updates

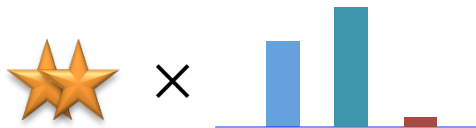
**Apply:**

Compute user's  
new factor weights

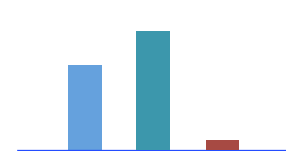
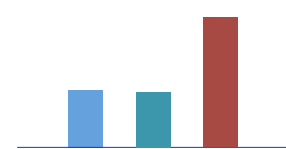
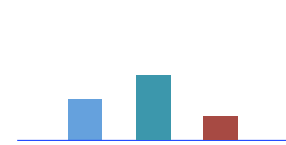
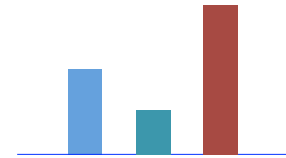
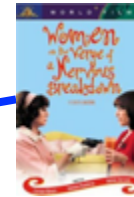


**Gather:**

sum over movies,  
product of ratings  
& factor weights  
(and a little more info)

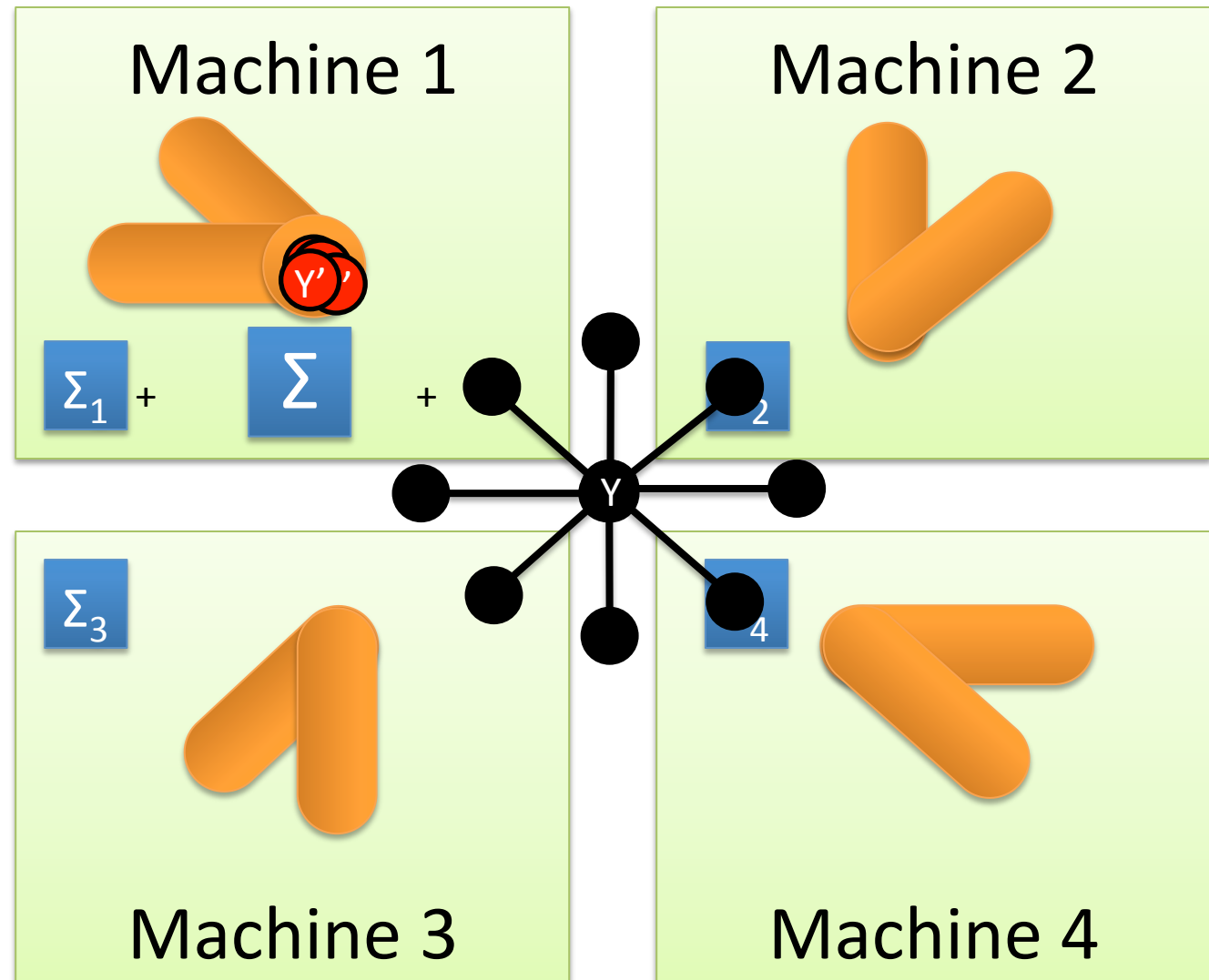


Iterate over  
users & movies

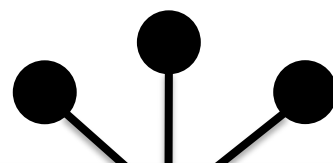


# Distributed Execution of a GL2 PowerGraph Vertex-Program

**Gather**  
**Apply**  
**Scatter**



# Minimizing Communication in GL2 PowerGraph: Vertex Cuts



Communication linear  
in # connected machines

GL2 PowerGraph includes novel vertex cut algorithms



Provides order of magnitude gains in performance  
# machines per vertex

*Percolation theory suggests Power Law graphs can be split  
by removing only a small set of vertices [Albert et al. 2000]*

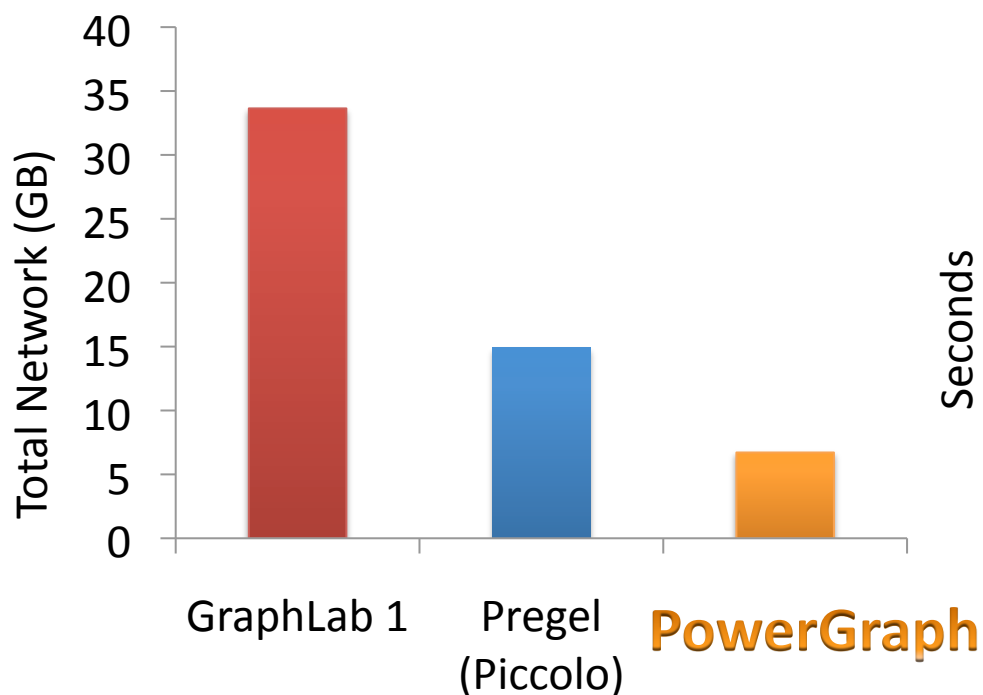


*Small vertex cuts possible!*

# PageRank on the Twitter Follower Graph

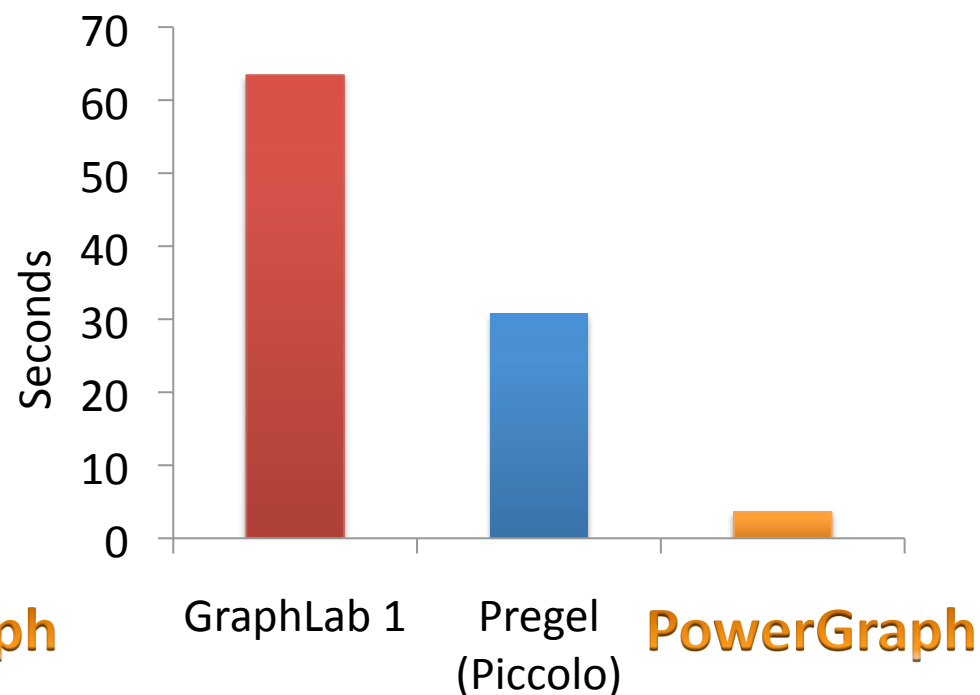
## Natural Graph with 41M Users, 1.4 Billion Links

### Communication



Reduces Communication

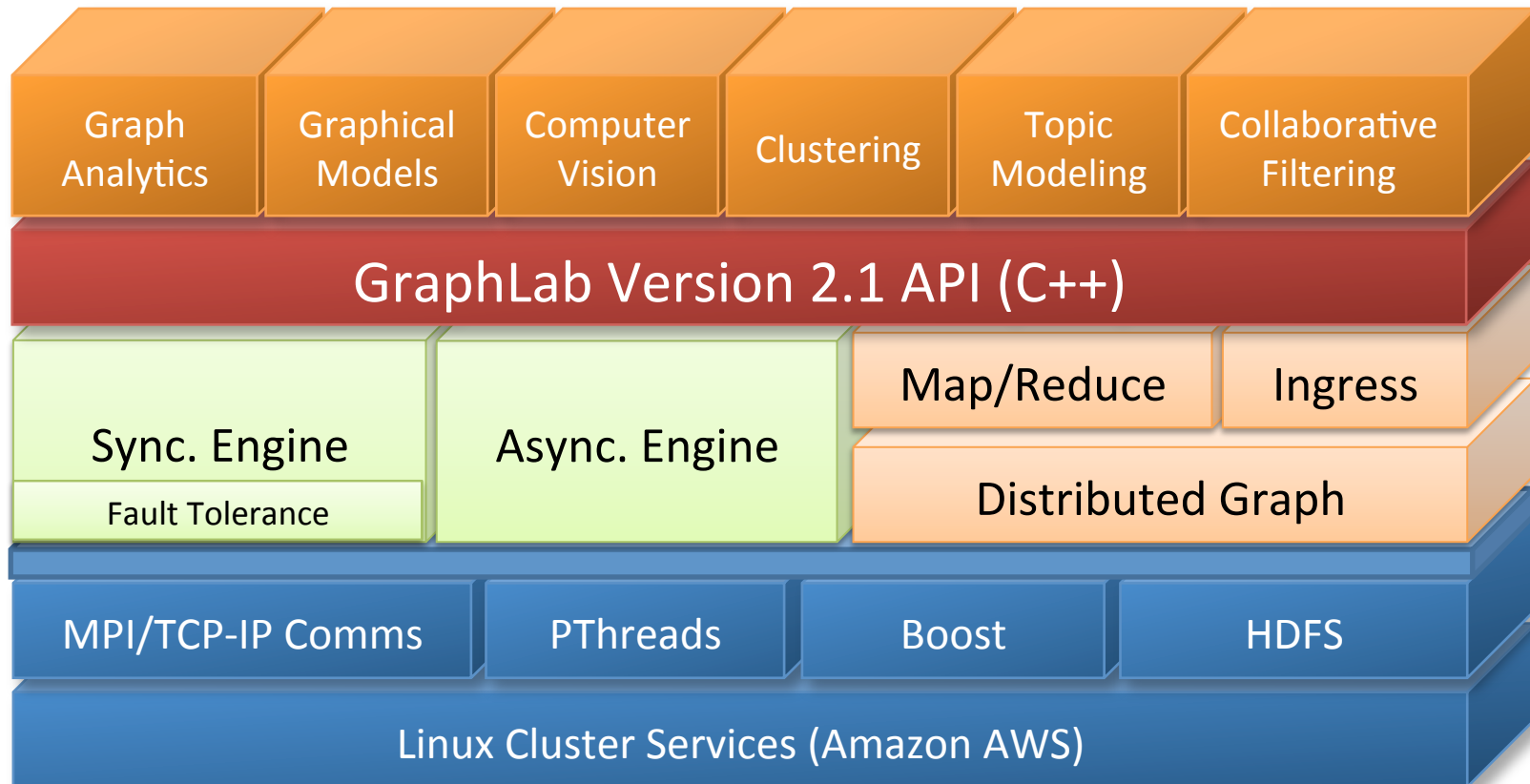
### Running time



Order of Magnitude Faster



# From the Abstraction to a System



# Triangle Counting on Twitter Graph

## 34.8 Billion Triangles

Hadoop [WWW'11]	1636 Machines 423 Minutes
GL2 PowerGraph	64 Machines 15 Seconds

*Why? Wrong Abstraction* →  
Broadcast  $O(\text{degree}^2)$  messages per Vertex



# Topic Modeling (LDA)



- English language Wikipedia
  - 2.6M Documents, 8.3M Words, 500M Tokens
  - Computationally intensive algorithm

## Million Tokens Per Second

0 20 40 60 80 100 120 140 160

Smola et al. 100 Yahoo! Machines  
**Specifically engineered for this task**

GL2 PowerGraph 64 cc2.8xlarge EC2 Nodes  
**200 lines of code & 4 human hours**

# How well does GraphLab scale?

Yahoo Altavista Web Graph (2002):

One of the largest publicly available webgraphs

**1.4B Webpages, 6.7 Billion Links**

**7 seconds per iter.**

**1B links processed per second**

**30 lines of user code**



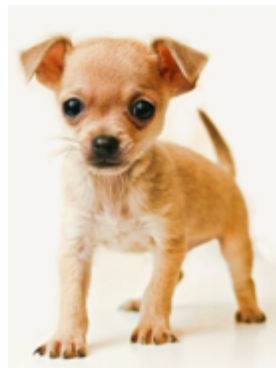
**1024 Cores (2048 HT)**



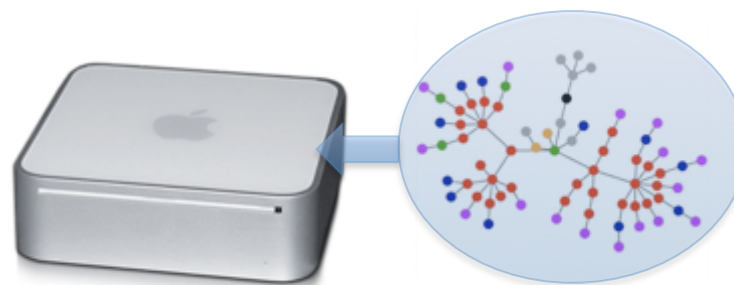
**4.4 TB RAM**

# GraphChi: Going small with GraphLab

GraphLab



Solve huge problems on  
small or embedded  
devices?

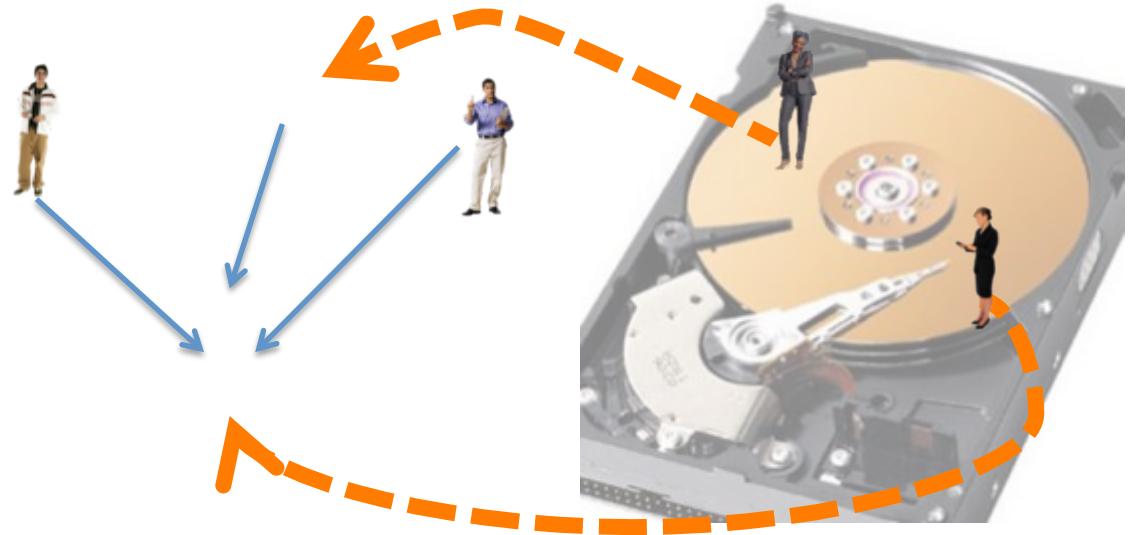


**Key: Exploit non-volatile memory  
(starting with SSDs and HDs)**

# GraphChi – disk-based GraphLab

## Challenge:

*Random Accesses*



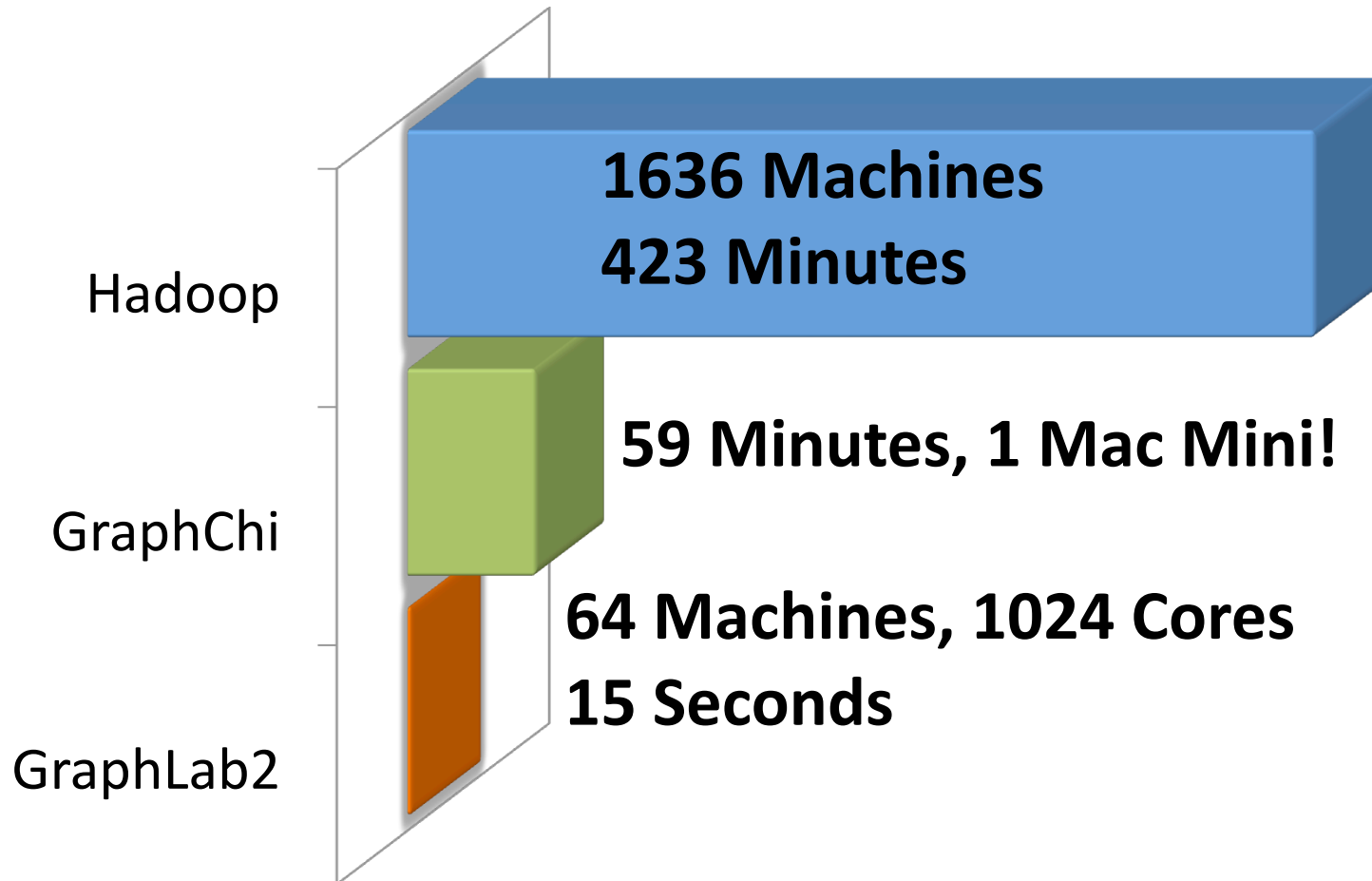
## Novel GraphChi solution:

*Parallel sliding windows method →  
minimizes number of random accesses*

# Triangle Counting on Twitter Graph

**40M Users**  
**1.2B Edges**

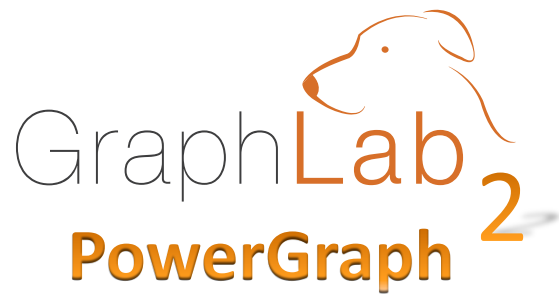
**Total: 34.8 Billion Triangles**



Hadoop results from [Suri & Vassilvitskii '11]



- ML algorithms as vertex programs
- Asynchronous execution and consistency models



- Natural graphs change the nature of computation
- Vertex cuts and gather/apply/scatter model

GL2 PowerGraph  
focused on  
**Scalability**

at the loss of  
**Usability**

# GraphLab 1

```
PageRank(i, scope){  
  acc = 0  
  for (j in InNeighbors) {  
    acc += pr[j] * edge[j].weight  
  }  
  pr[i] = 0.15 + 0.85 * acc  
}
```

**Explicitly described operations**

**Code is intuitive**



# GraphLab 1

```
PageRank(i, scope){  
  acc = 0  
  for (j in InNeighbors) {  
    acc += pr[j] * edge[j].weight  
  }  
  pr[i] = 0.15 + 0.85 * acc  
}
```

Explicitly described operations

**Code is intuitive**

# GL2 PowerGraph

Implicit operation

```
gather(edge) {  
  return edge.source.value *  
    edge.weight  
}
```

```
merge(acc1, acc2) {  
  return accum1 + accum2  
}
```

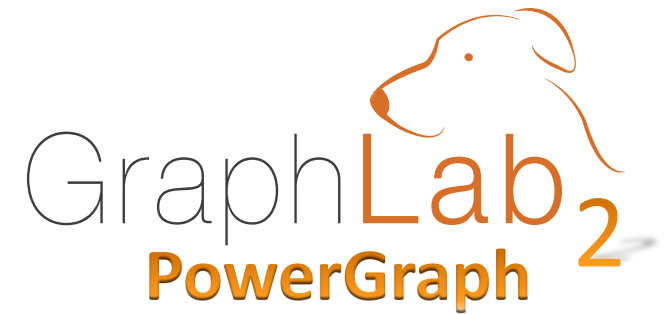
Implicit aggregation

```
apply(v, accum) {  
  v.pr = 0.15 + 0.85 * acc  
}
```

**Need to understand engine  
to understand code**

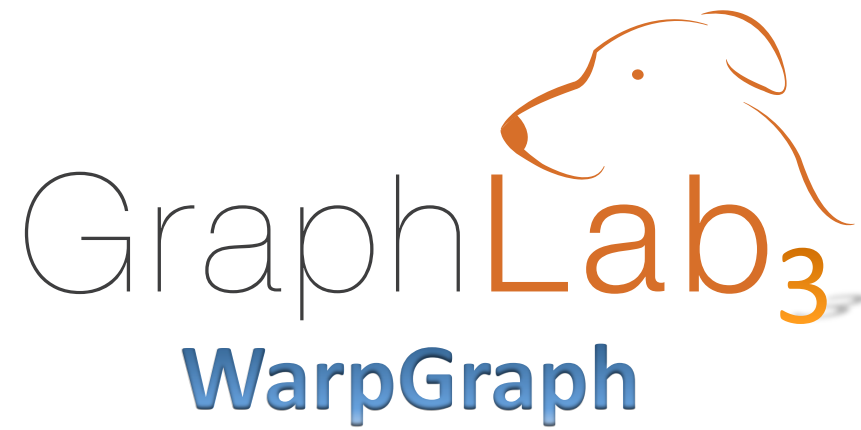


Great flexibility,  
but hit scalability wall



Scalability,  
but very rigid abstraction  
(many contortions needed to implement  
SVD++, Restricted Boltzmann Machines)



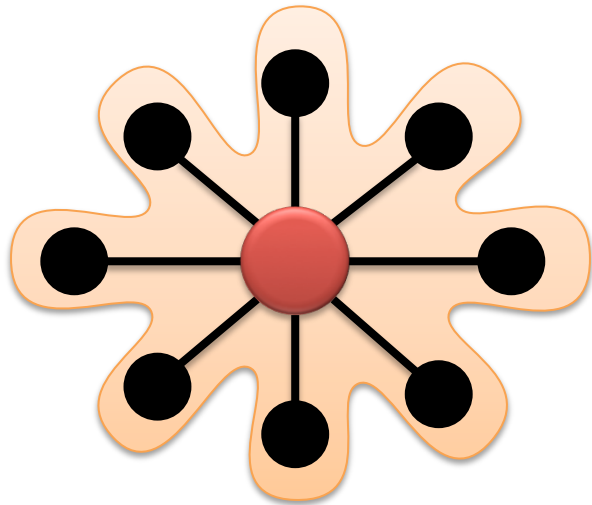


## USABILITY

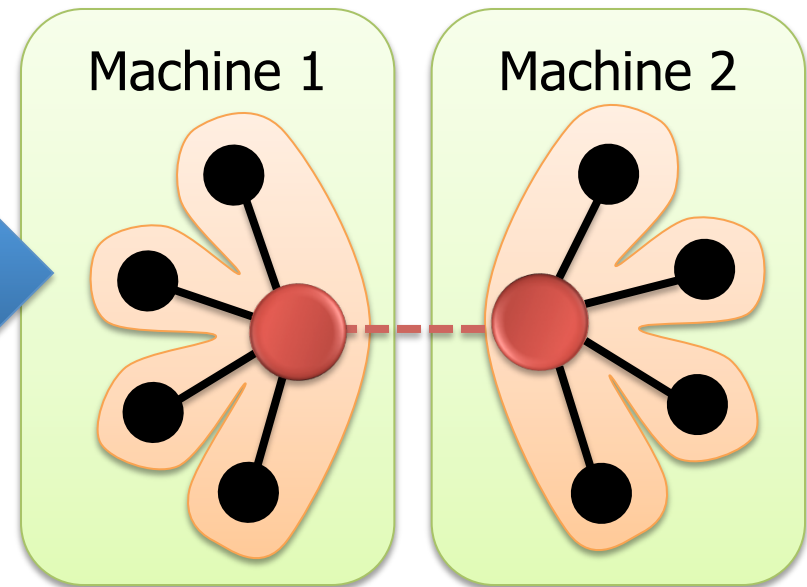


# GL3 WarpGraph Goals

**Program  
Like GraphLab 1**



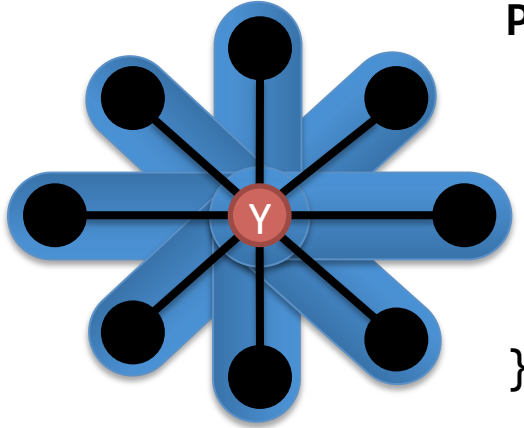
**Run Like  
GraphLab 2**



# Fine-Grained Primitives

Expose Neighborhood Operations through Parallelizable Iterators

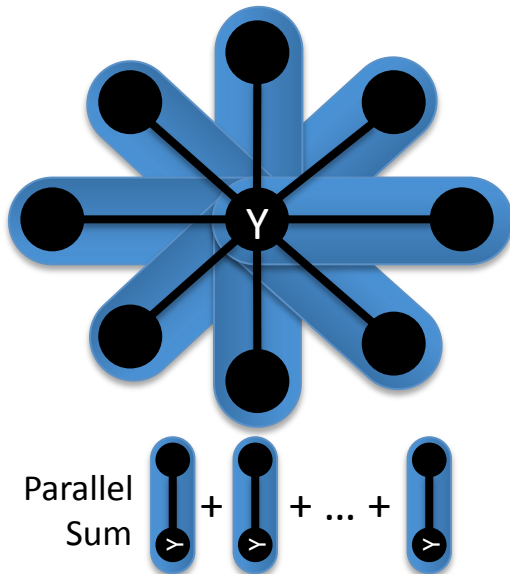
$$R[i] = 0.15 + 0.85 \sum_{(j,i) \in E} w[j,i] * R[j]$$



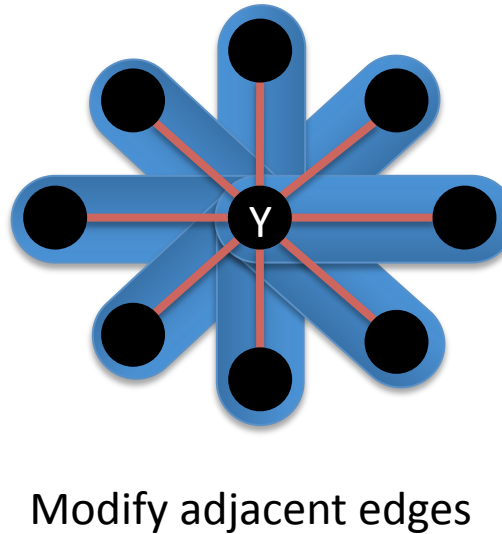
```
PageRankUpdateFunction(Y) {  
    Y.pagerank = 0.15 + 0.85 *  
}
```

# Expressive, Extensible Neighborhood API

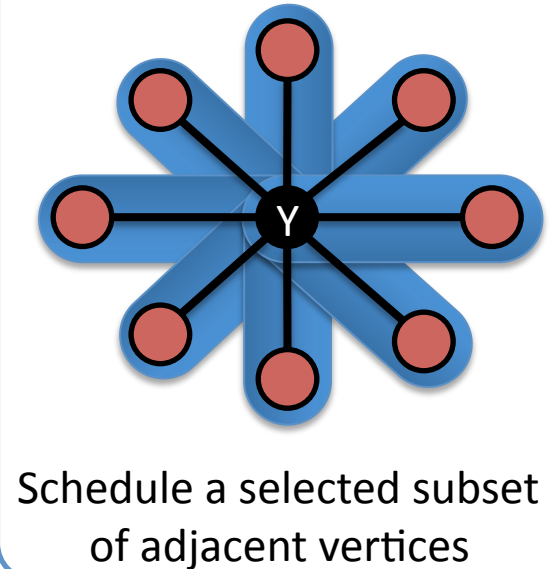
## MapReduce over Neighbors



## Parallel Transform Adjacent Edges

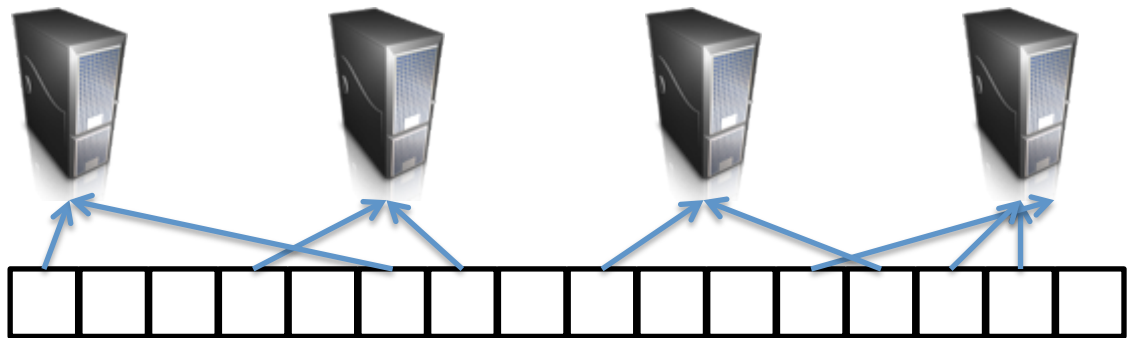


## Broadcast



DHT Get Keys

DHT Update Keys



# Can express every GL2 PowerGraph program (more easily) in GL3 WarpGraph

But GL3 is more  
expressive

```
UpdateFunction(v) {  
  if (v.data == 1)  
    accum = MapReduceNeighs(g,m)  
  else ...  
}
```

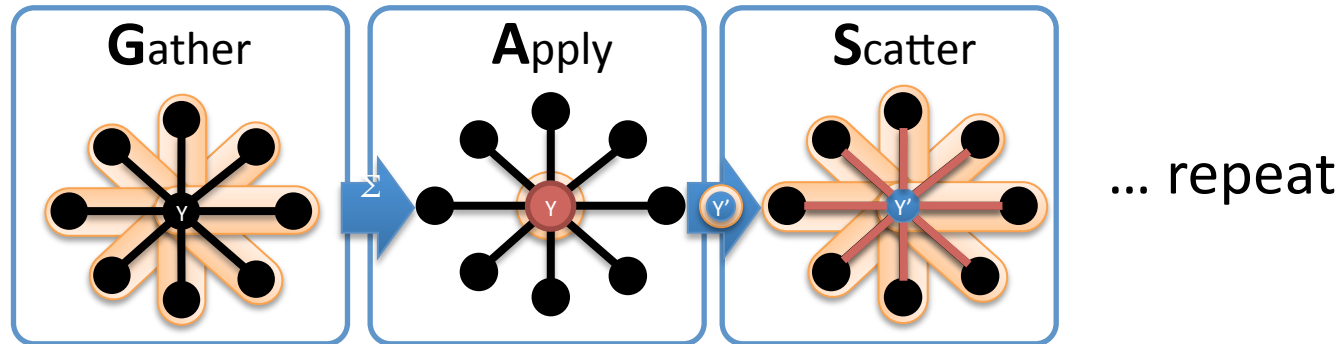
Multiple  
gathers

Scatter before  
gather

Conditional  
execution

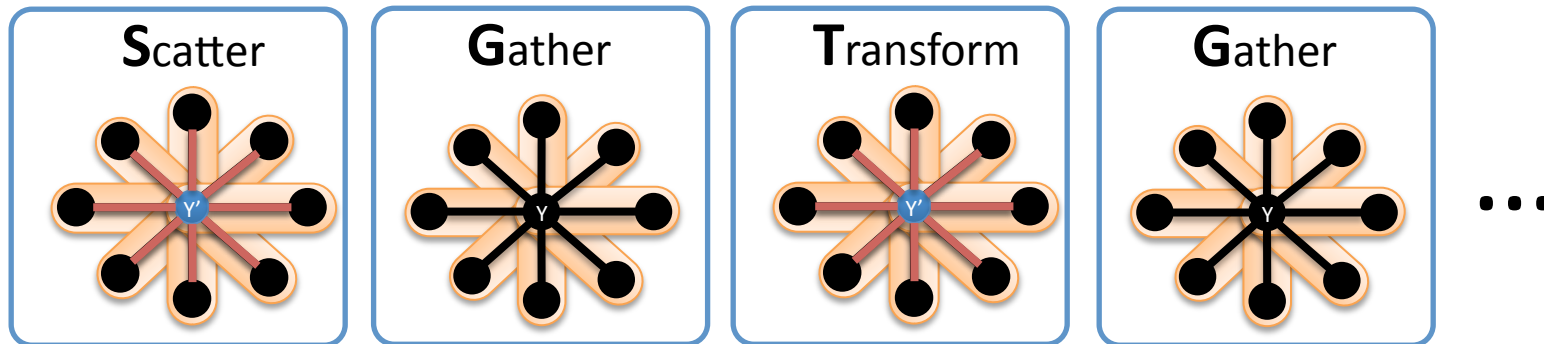
## GL2 PowerGraph:

Fast because **communication** phases are very **predictable**



## GL3 WarpGraph:

**Communication highly unpredictable**



**Risk: High Latency**

(spend all our time waiting for a reply...)

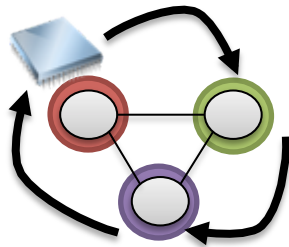


# Hide Latency

## Do Something Else while Waiting

Create 1000s of threads, each running an update function on a different vertex

### Performance Bottleneck: Context Switching



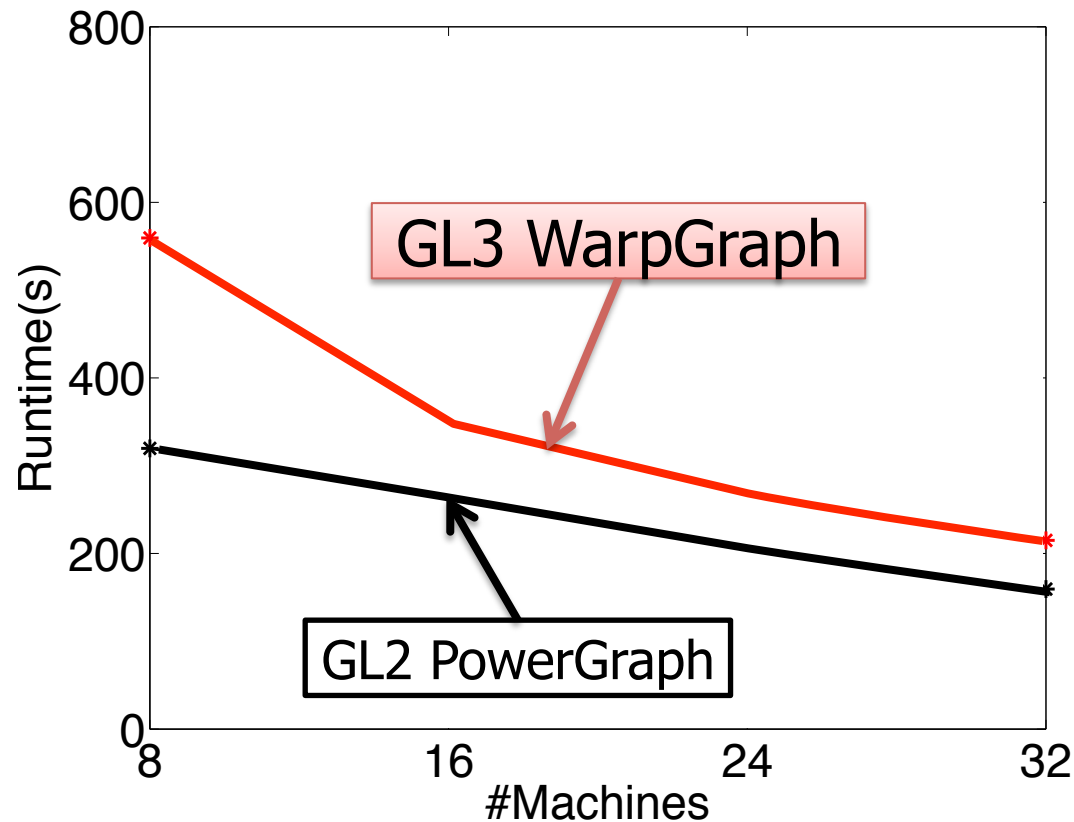
Every cycle used in  
context switching is wasted  
(OS context switch is slow requiring 10K-100k cycles)

### GL3 WarpGraph: Novel user-mode threading

8M context switches per second

**100x faster than OS**

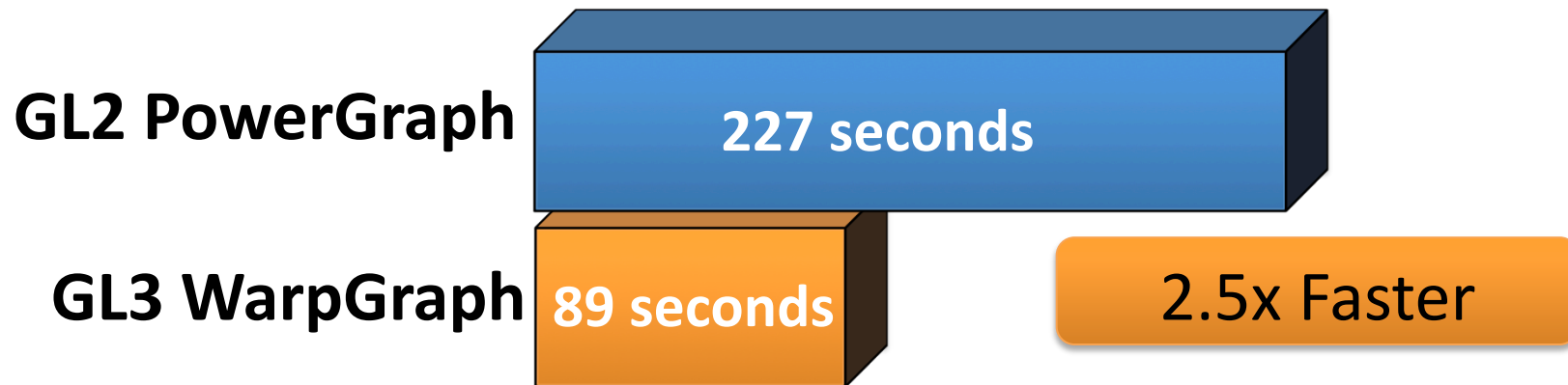
# PageRank Twitter Graph: 41M Vertices 1.4B Edges



WarpGraph only 25% slower, with much improved programmability  
**But, here, asynchrony not fundamental for performance**

# Graph Coloring

Twitter Graph: 41M Vertices 1.4B Edges



**Asynchrony fundamental here →  
WarpGraph outperforms PowerGraph with simpler code**

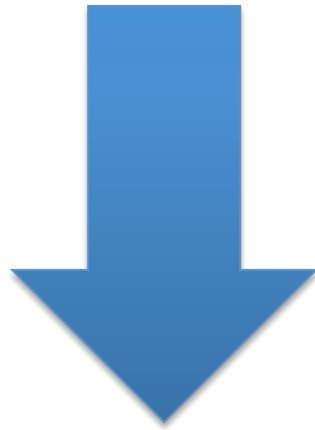
# Usability

Consensus that WarpGraph is much easier to use than PowerGraph

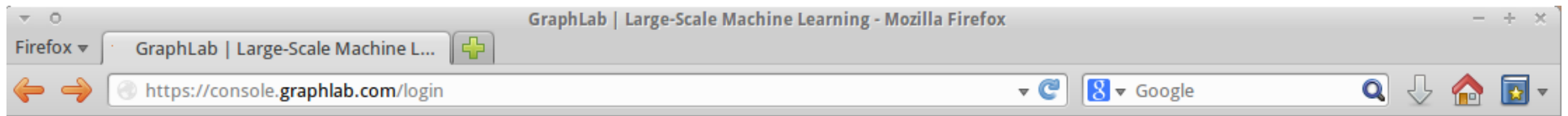
User study size = 2 :-)

**Bigger + Real User Study in Progress,  
as we release new open-source version of GraphLab**

**New abstraction simplifies  
writing programs in GraphLab**



**But you still need to get a cluster,  
install GraphLab, configure system...**



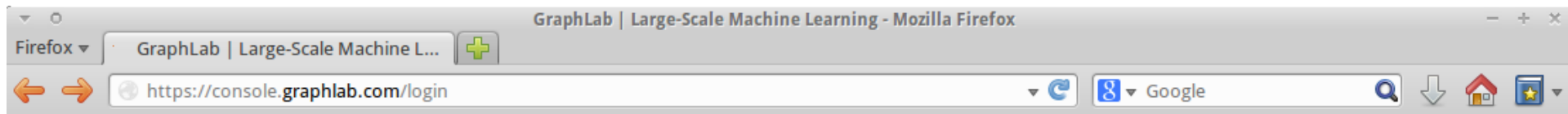
Your login :

Password :

Log in

Cancel

You don't have an account yet? [Register here](#)



Demo User | [Help](#)

### Launch GraphLab Cluster

Name:

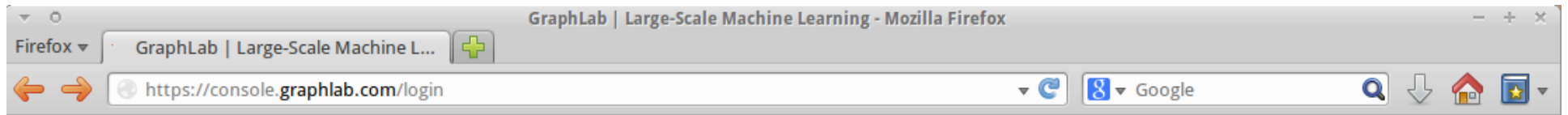
Node Configuration | Number:

Type:

m2.4xlarge

[Options](#) | [SSL Certificate](#)

Launch

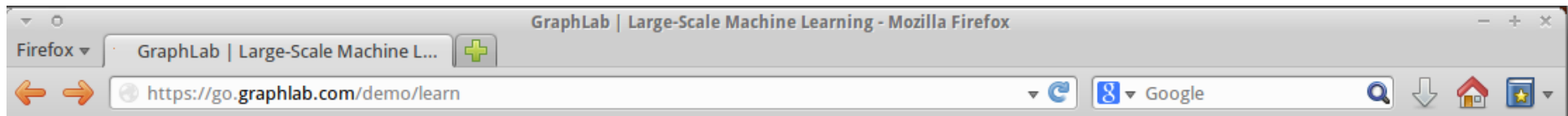


Demo User | [Help](#)

### Clusters

Name	Connect URL	Configuration	Status	Last
demo	<a href="https://go.graphlab.com/user/demo/learn">https://go.graphlab.com/user/demo/learn</a> <a href="#">Connect</a>	1 Node (m2.4xlarge)	Running	today





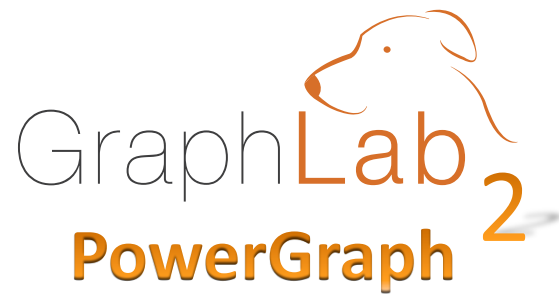
[Getting Started](#) | [Help](#)

```
In [2]: import GraphLab
gl = GraphLab.UndirectedTriangleCount()
file = '/home/graphlab/python-demo/1M.tsv'
(seconds, triangles) = gl.execute(input_file = file)
print "Finding %d undirected triangles in '%s' took %f seconds." % (triangles, file, seconds)
```

Finding 329024 undirected triangles in '/home/graphlab/python-demo/1M.tsv' took 2.439280 seconds.



- ML algorithms as vertex programs
- Asynchronous execution and consistency models

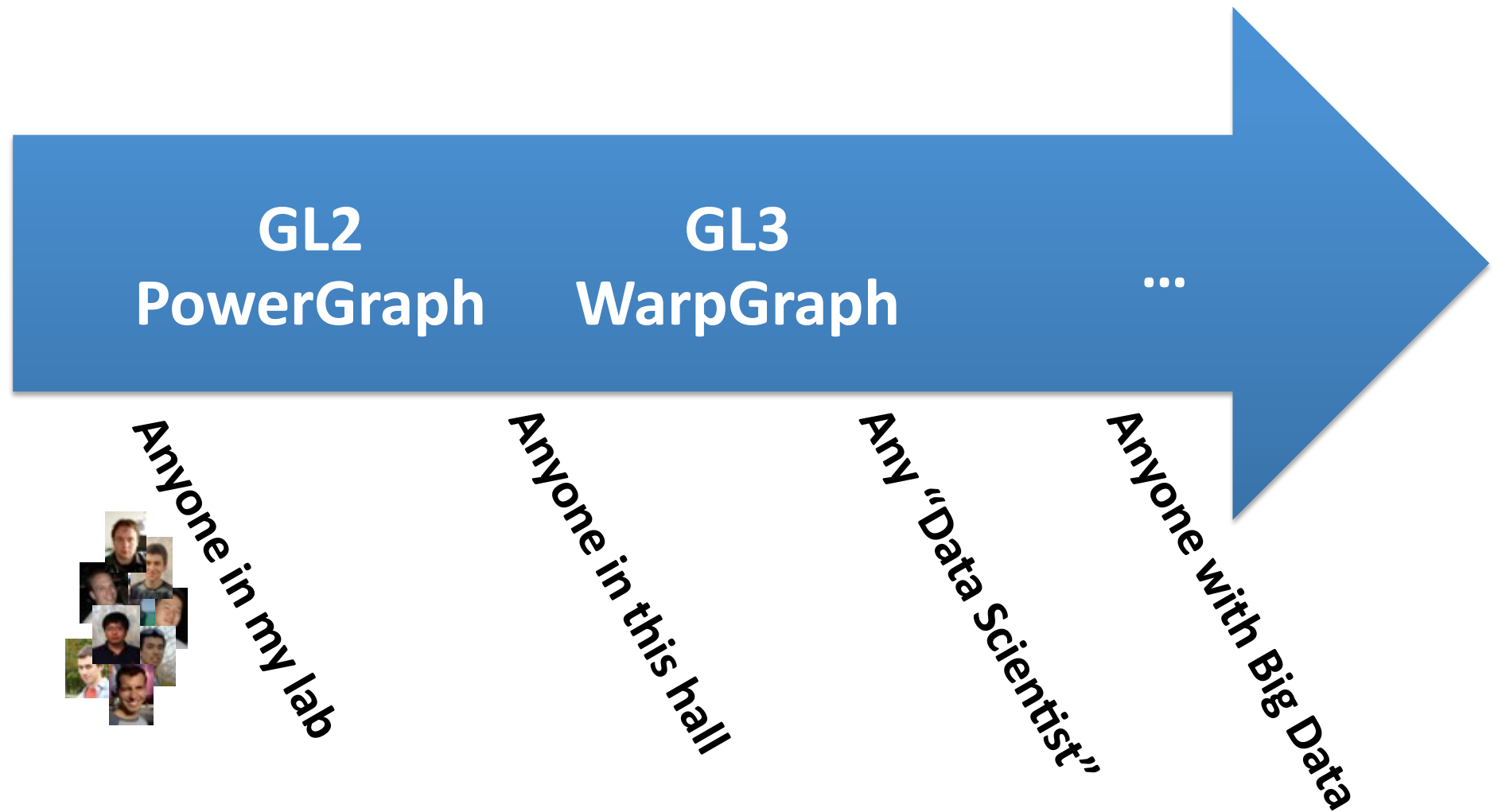


- Natural graphs change the nature of computation
- Vertex cuts and gather/apply/scatter model



- Usability is key
- Access neighborhood through parallelizable iterators and latency hiding

# Usability for Whom???



# Machine Learning

## PHASE 3

USABILITY



# Exciting Time to Work in ML



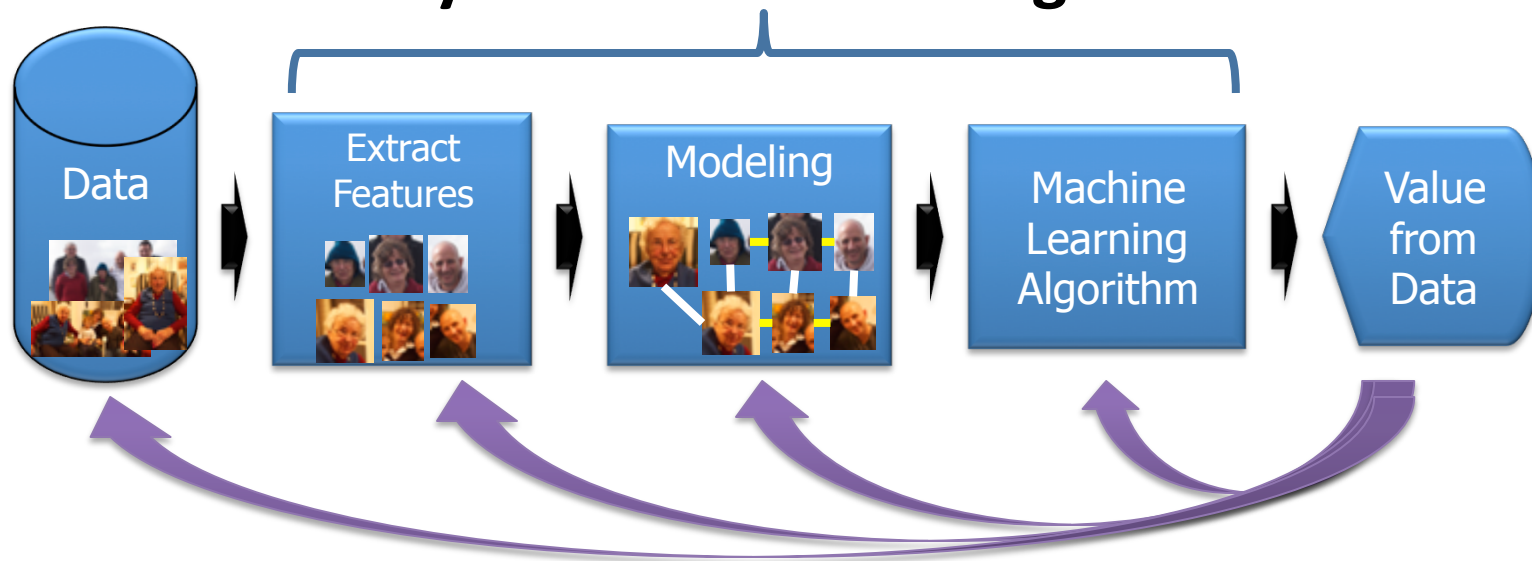
Unique opportunities to change the world!! 😊

But, **every deployed system** is an one-off solution, and **requires PhDs** to make work... 😞

And, Usability for ML is not just “Engineering” –  
Must Be Easy to Iterate through Models to Solve Task

But, when ML doesn't work, need a PhD to understand why...

GraphLab One kind of ML usability:  
Fast and easy iterations over huge datasets



Interpretable feature engineering?  
No parameters to tune, please...



Why was this prediction made?  
How can I give *valuable* feedback?



**v1** Possibility

**v2** Scalability

**v3** Usability

GraphLab 2.2 available now: **graphlab.com**