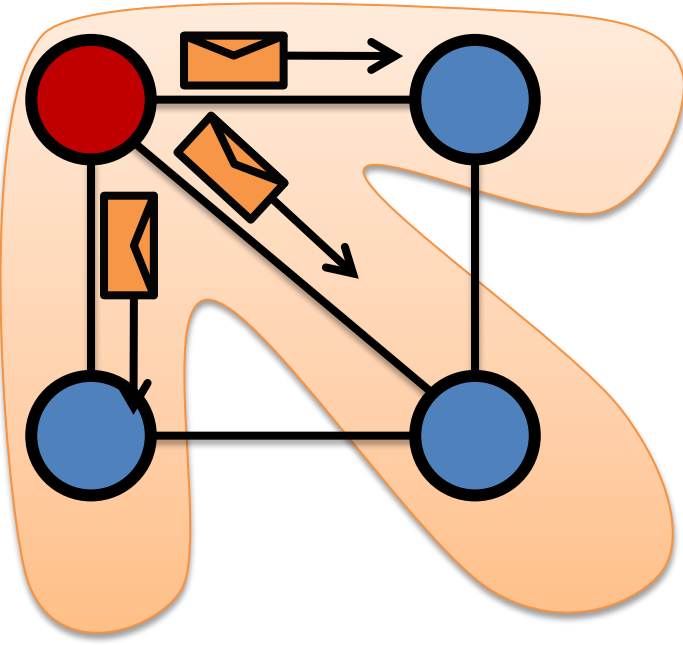


# GraphLab2: Distributed Graph-Parallel Computation on Natural Graphs

Joseph Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, Carlos Guestrin

## Graph Parallel Abstraction

- **Vertex-Program** associated with each vertex
- **Graph** constrains the interaction along edges
  - **Pregel**: Programs interact through messages
  - **GraphLab 1**: Programs can read each-others state



### Example: PageRank

Iterate:  $R[i] = 0.15 + 0.85 \sum_{(j,i) \in E} \frac{R[j]}{L[j]}$

Rank of Page  $i$     Random Reset Probability    Sum over in-links    Number of out-links form  $j$

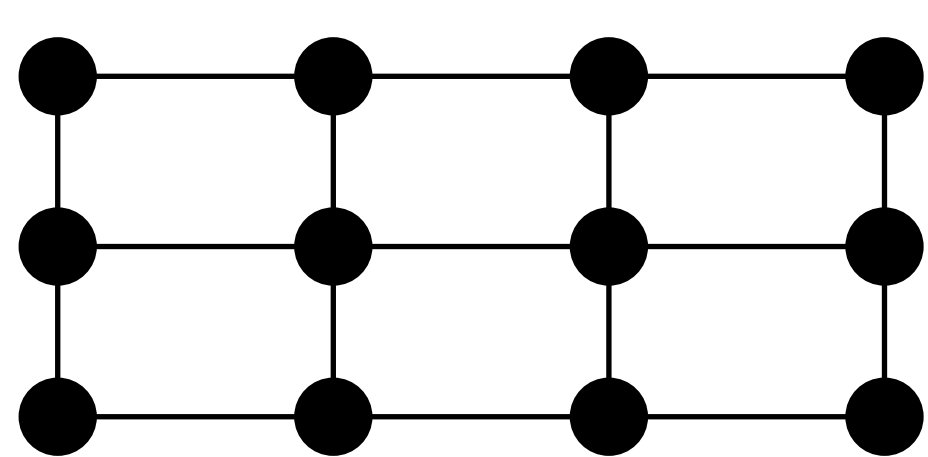
$R[5] = 0.15 + 0.85 \left( \frac{1}{3} R[1] + \frac{1}{1} R[4] \right)$

## Natural Graphs

Graphs Encode Relationships between Entities

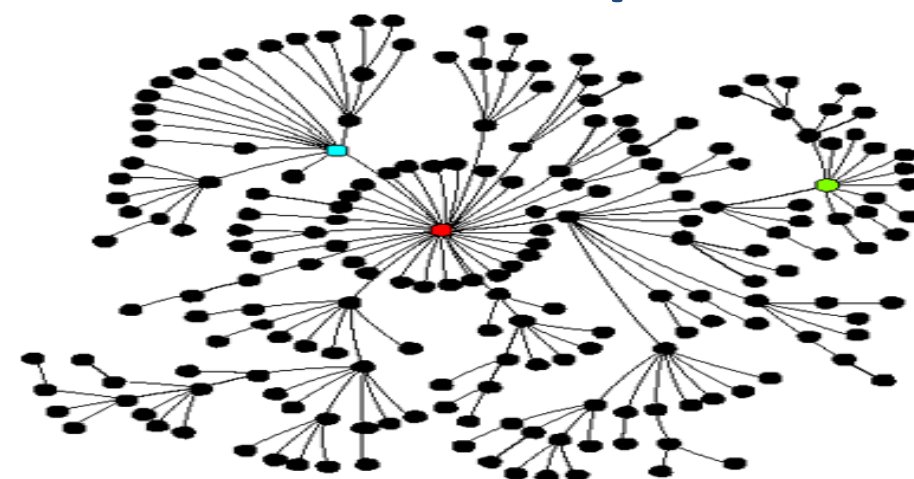


Assumed Structure



- **Small** neighborhoods
- **Similar** degree vertices
- Easy to **partition**

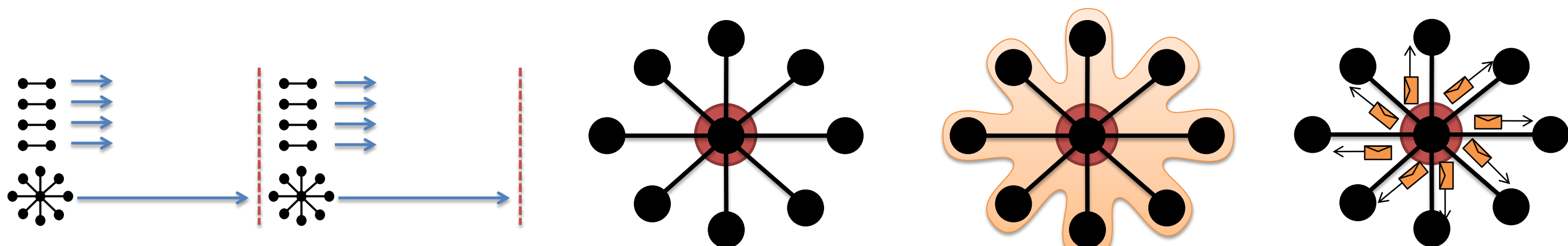
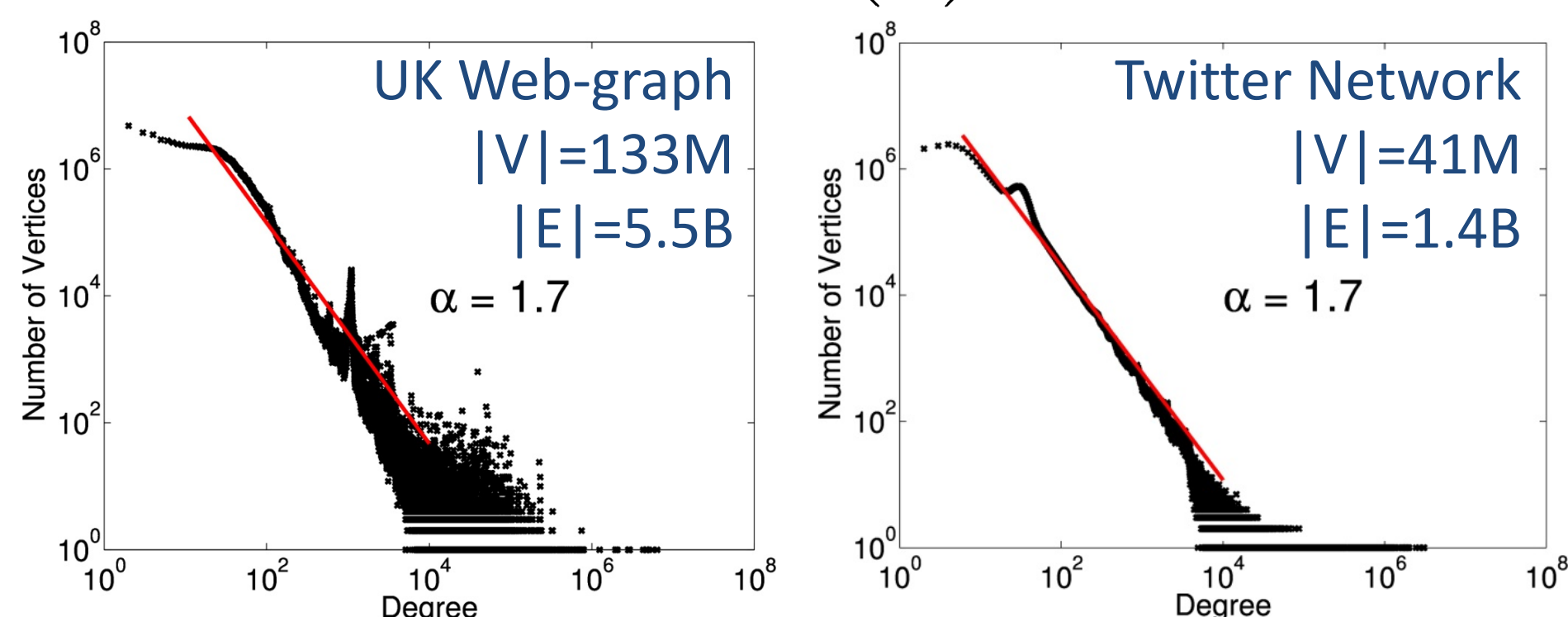
Natural Graph



- **Large** Neighborhoods
- **Power-Law Degree**
- Difficult to partition

### Power-Law Graphs

- Most vertices have relatively few neighbors while a few vertices have many neighbors
- Probability of having degree  $d$ :  $P(d) \propto d^{-\alpha}$



Synchronous execution is prone to stragglers.

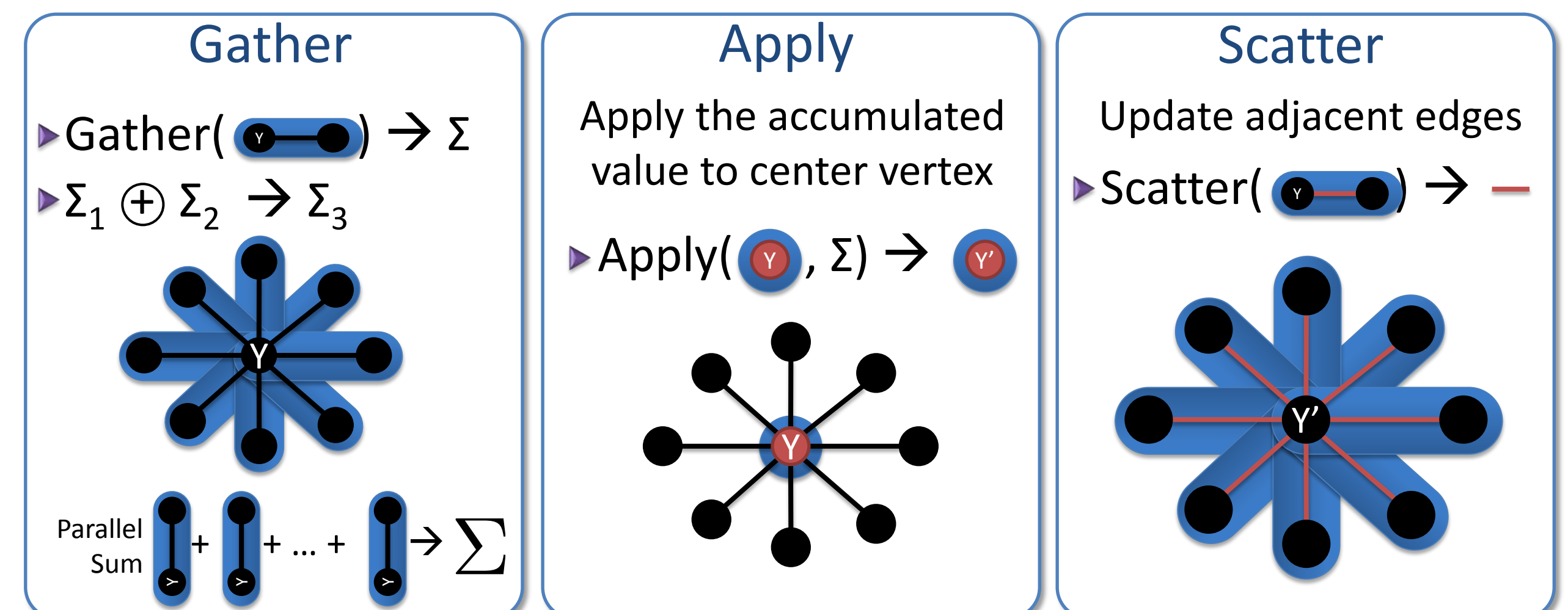
Edge information overwhelms single machine

Touches a large fraction of graph (GraphLab 1)

Produces many messages (Pregel)

## PowerGraph (OSDI '12)

- **Factorized Vertex-Programs**:



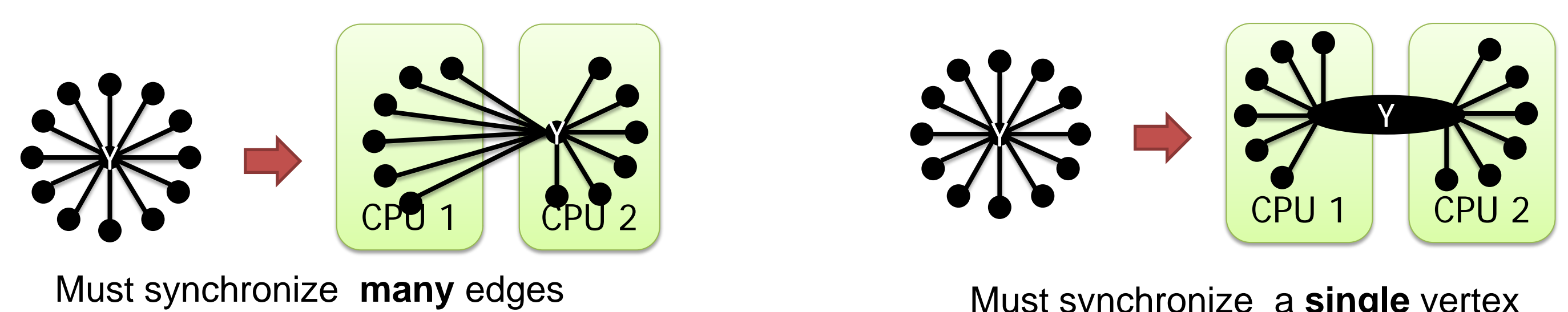
- Distribute vertex-program over several machines and move computation to data

### PowerGraph\_PageRank(v)

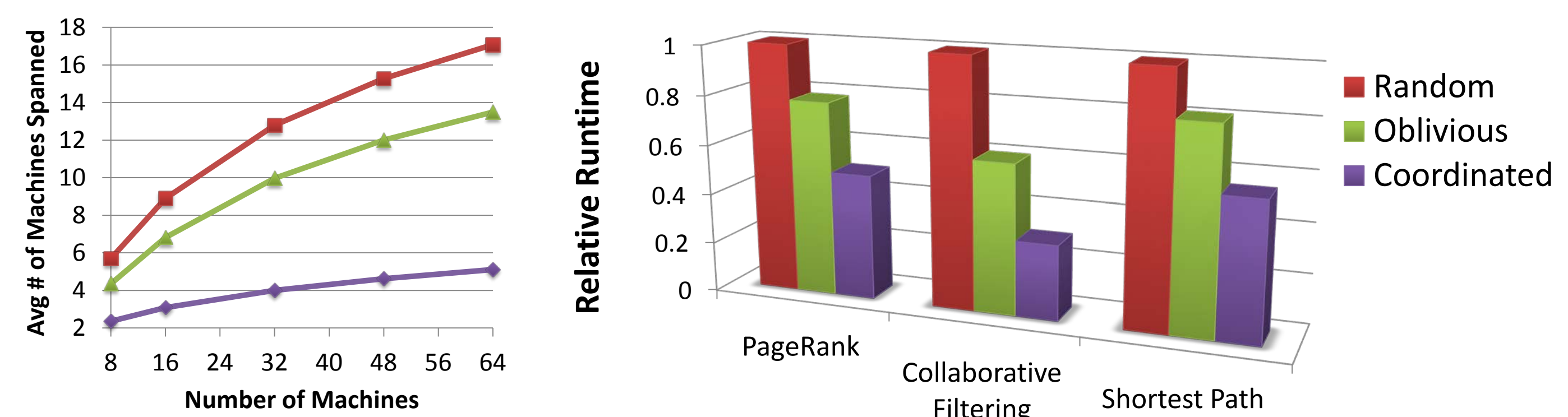
**Gather**(  $u \rightarrow v$  ) : return  $R[u] / \text{\#outlinks}[u]$   
**sum**(  $a, b$  ) : return  $a + b$ ;  
**Apply**(  $v, \Sigma_v$  ) :  $R[v] = 0.15 + 0.85 \Sigma_v$   
**Scatter**(  $v \rightarrow w$  ) :  
 if  $(R[v]$  changes by  $\epsilon$ ) then *activate*( $w$ )

## Vertex Separators

Cut vertices instead of edges!

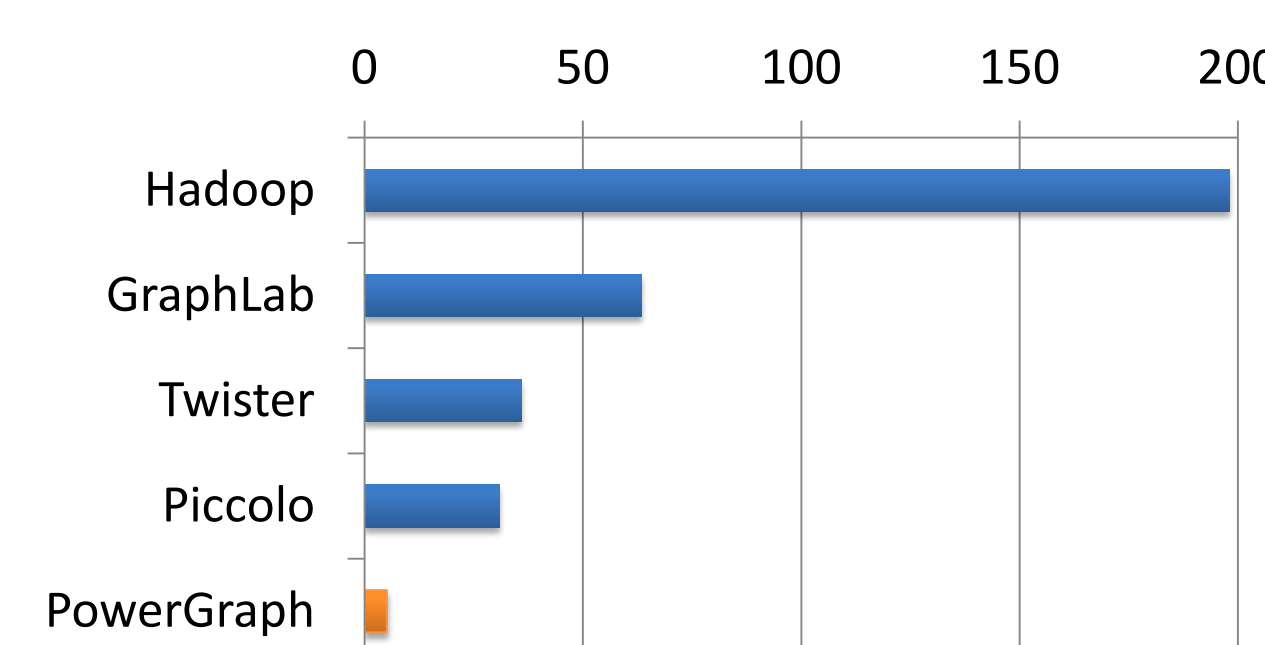


### Three Streaming Partitioning Strategies



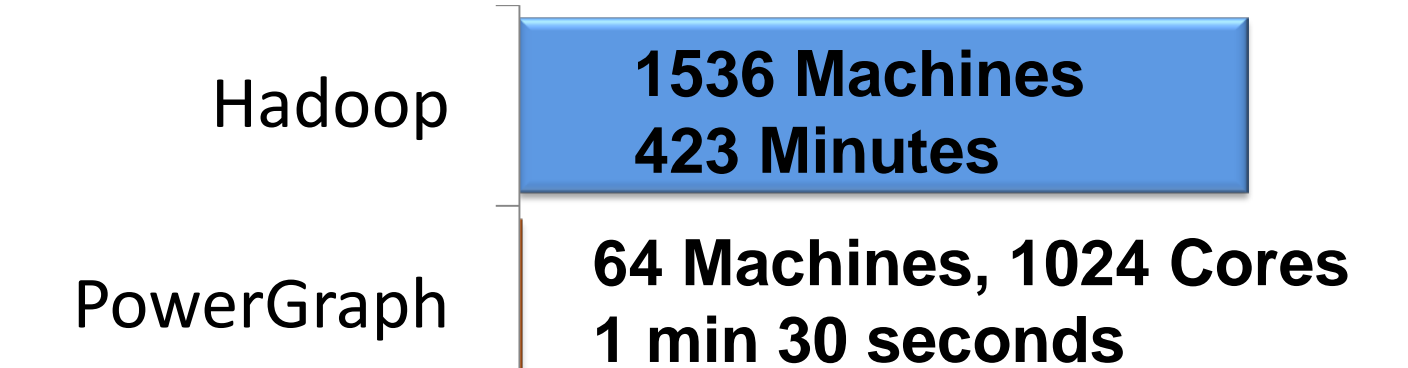
## Results

**PageRank (per Iteration):**  
 40M Webpages, 1.4B Links

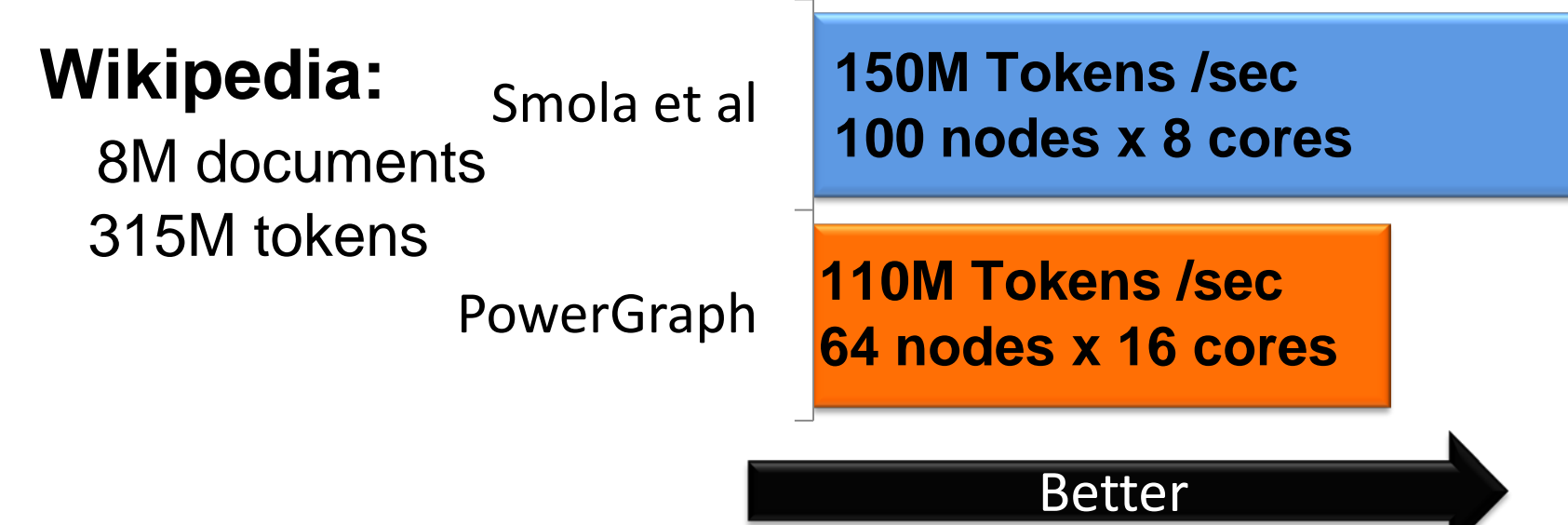


**Triangle Counting:**

**Twitter Graph:**  
 40 Million Vertices  
 1.4 Billion Edges  
 34.8 Billion Triangles



### Latent Dirichlet Allocation



Available Now!

GraphLab

<http://GraphLab.org>

