

Counting Triangles and the **CURSE OF THE LAST REDUCER**

A paper by: Siddharth Suri and Sergei
Vassilvitskii

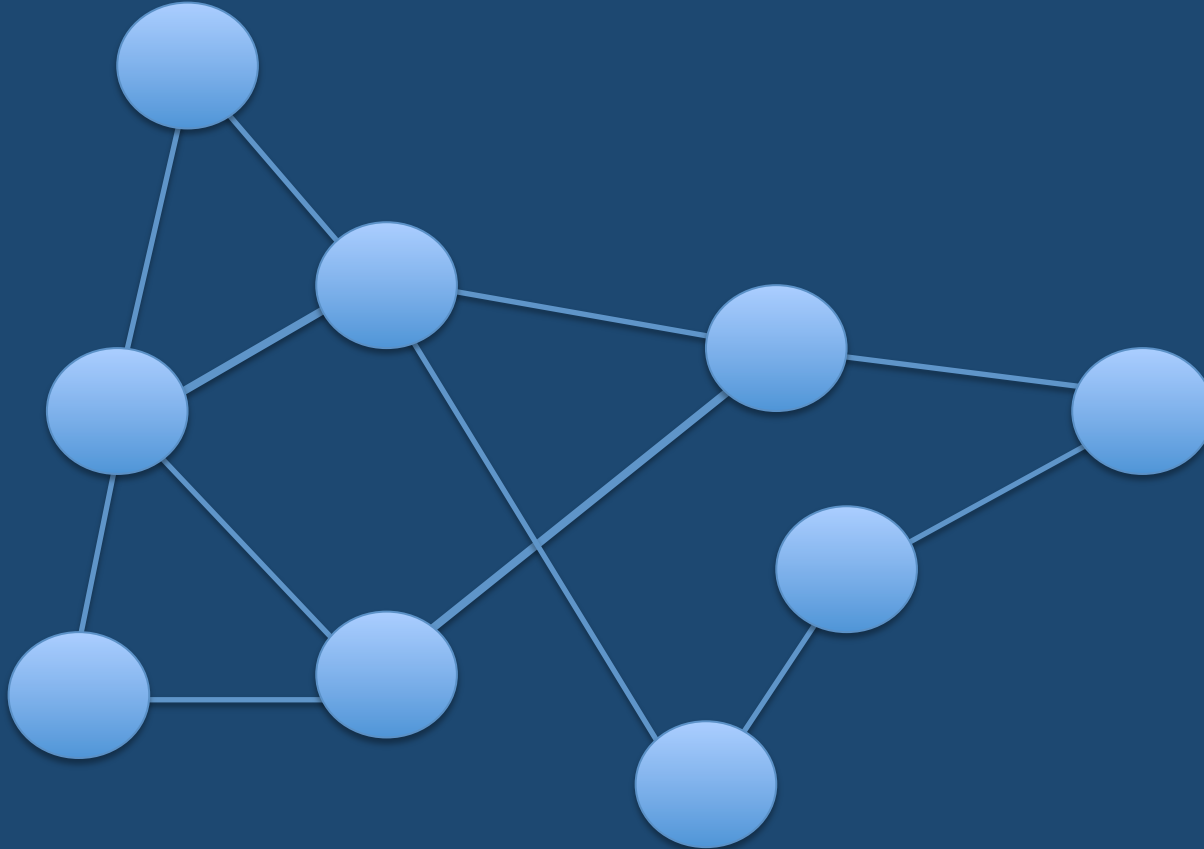
Presented by: Ryan Rogers (with some
slides from Sergei's Presentation)

Introduction

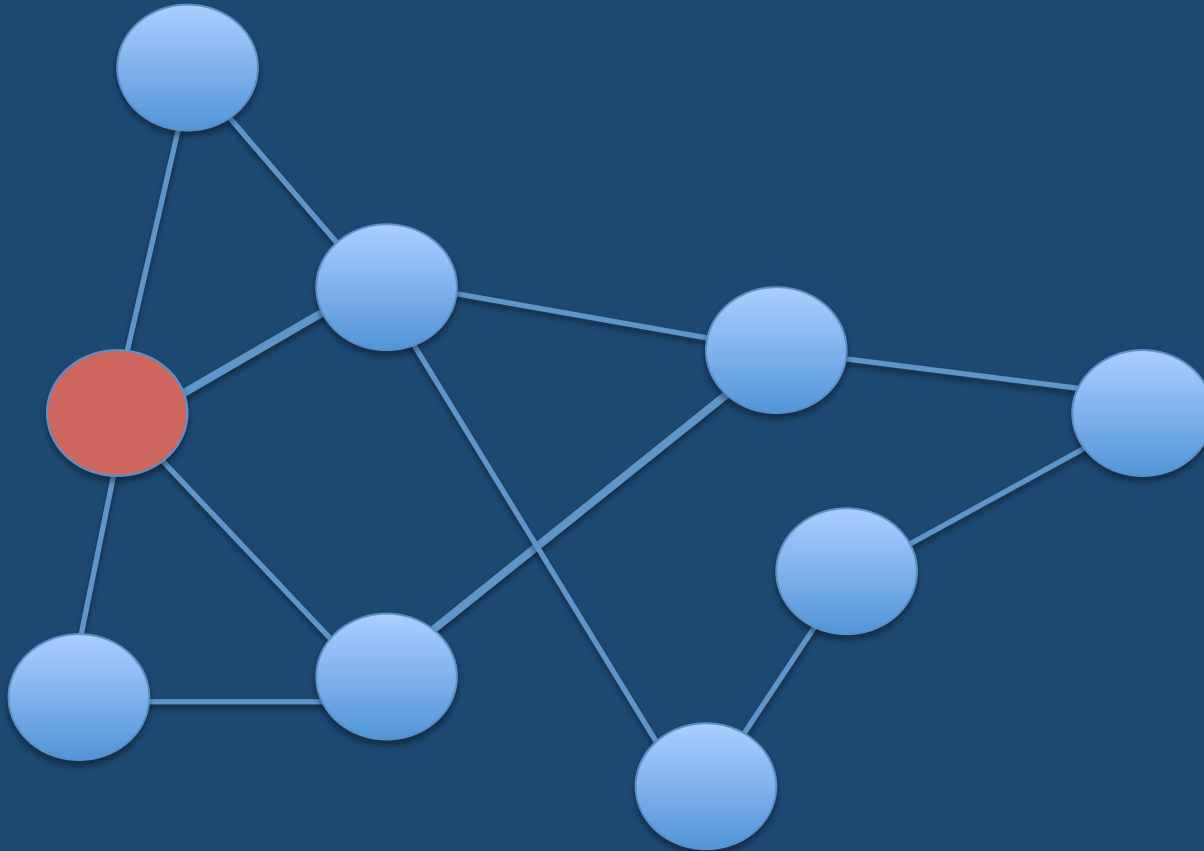


- Study Social Networks
- Main metric for analyzing Social Networks: Clustering Coefficient of each node
- Problem of finding the Clustering Coefficient of a node is basically the same as counting the number of ▲'s incident to that node.

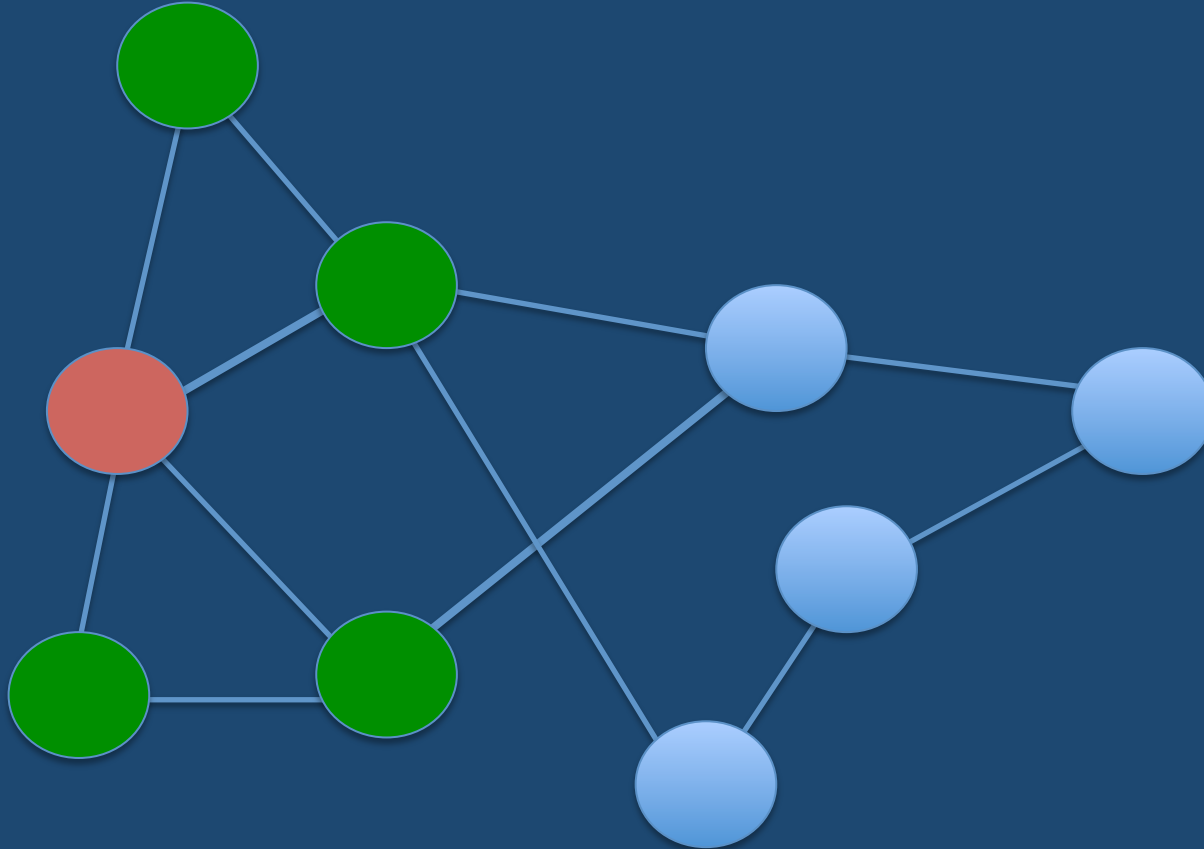
Clustering and Triangles



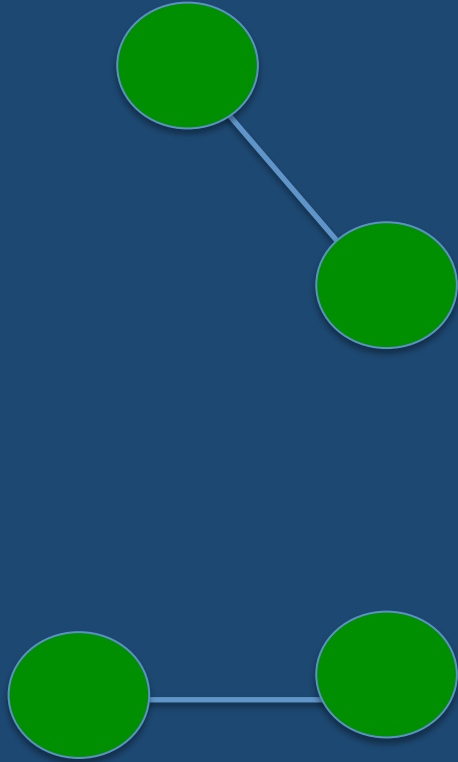
Clustering and Triangles



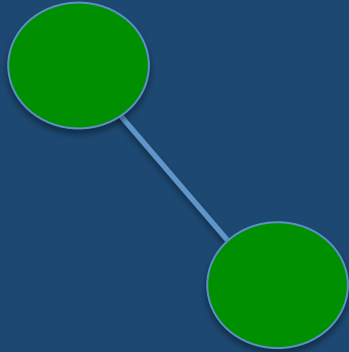
Clustering and Triangles



Clustering and Triangles



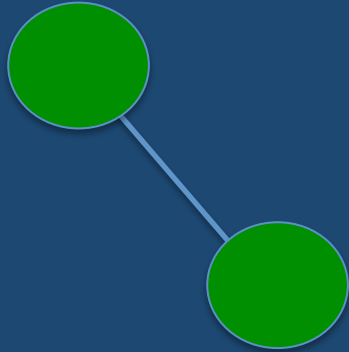
Clustering and Triangles



$$CC(\text{red circle}) = 2 / \binom{4}{2} \\ = 1/3$$



Clustering and Triangles

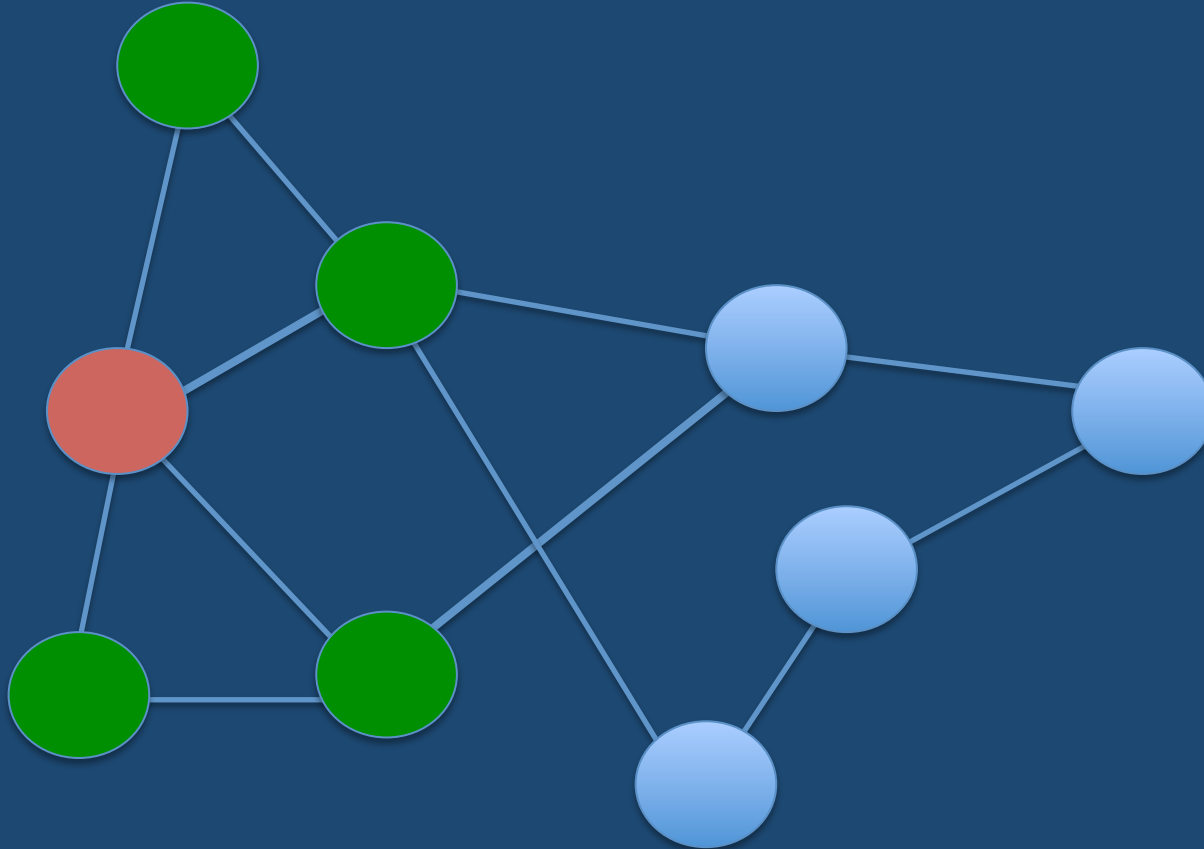


$$CC(\text{red circle}) = 2 / \binom{4}{2} \\ = 1/3$$

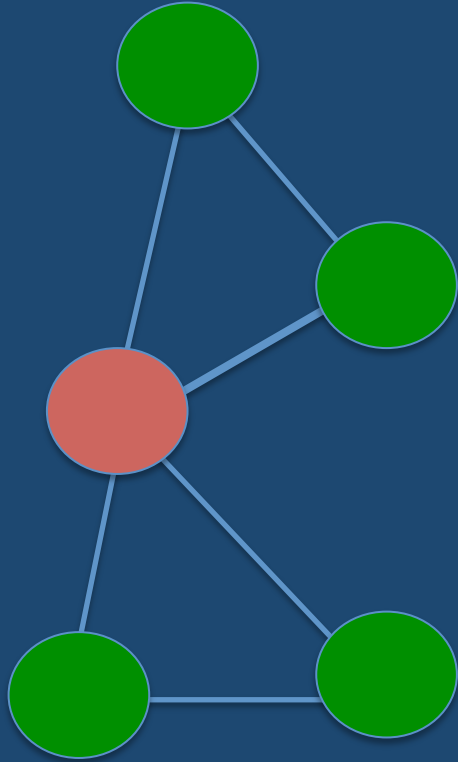


OR...

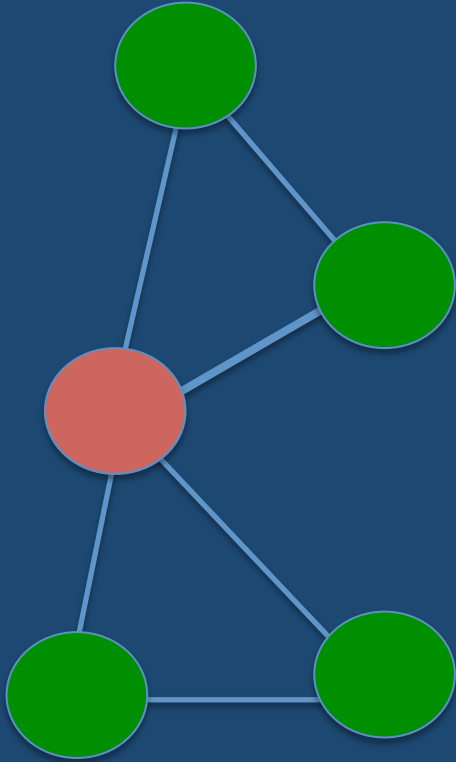
Clustering and Triangles



Clustering and Triangles



Clustering and Triangles



Number of  's

$$= \binom{d}{2} \times \text{CC}(\text{red node})$$

Past Work

- Coppersmith and Kumar ('04) and Buriol et al. ('04): Streaming algorithms to find **total** number of triangles with **high accuracy**
- Becchetti et al. ('08): **Estimate** the number of triangles incident on **each node**.
- Tsourakakis et al. ('09): Randomized MapReduce procedure that gives the **total** number of triangles accurately in **expectation**.

Contributions

- Count the **exact** number of triangles
- Count the number of triangles incident on **each node**, exactly.
- Comparable speedup as the randomized MapReduce procedure.

Counting Triangles (Naïve)

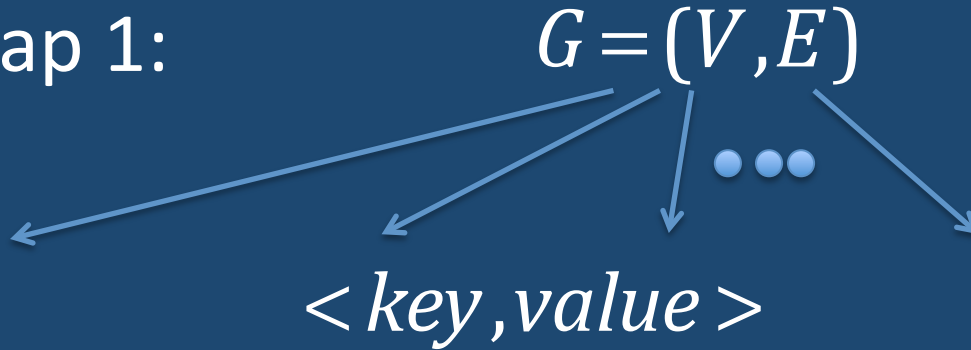
- Let $T \leftarrow 0$
 - for $v \in V$
 - for each $u \in \Gamma(v)$
 - for each $w \in \Gamma(v)$
 - » if $(u, w) \in E$
 $T \leftarrow T + 1/2$
- Output $T \leftarrow T / 3$

RUN TIME

$$O\left(\sum_{u \in V} d_u^2\right)$$

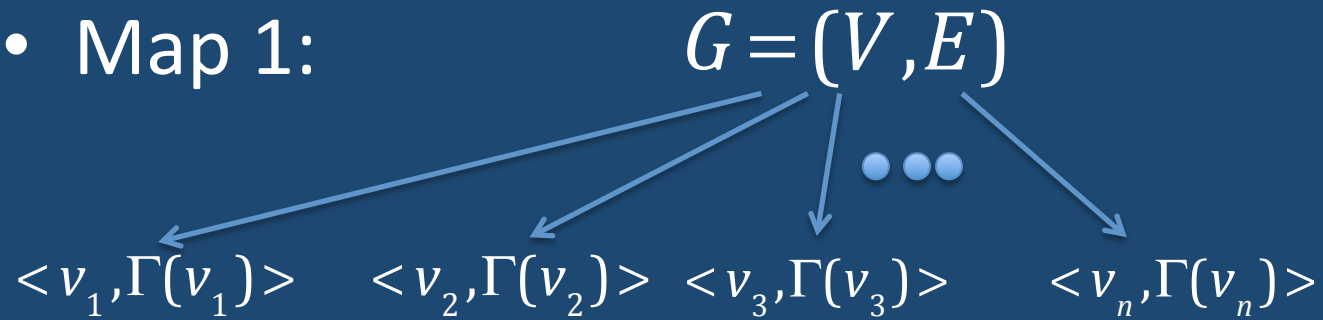
MapReduce (Naïve)

- Map 1:



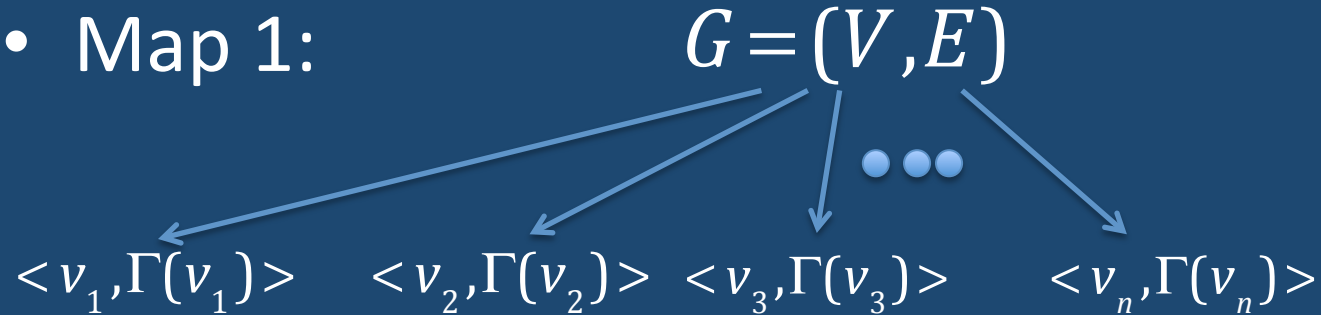
MapReduce (Naïve)

- Map 1:



MapReduce (Naïve)

- Map 1:

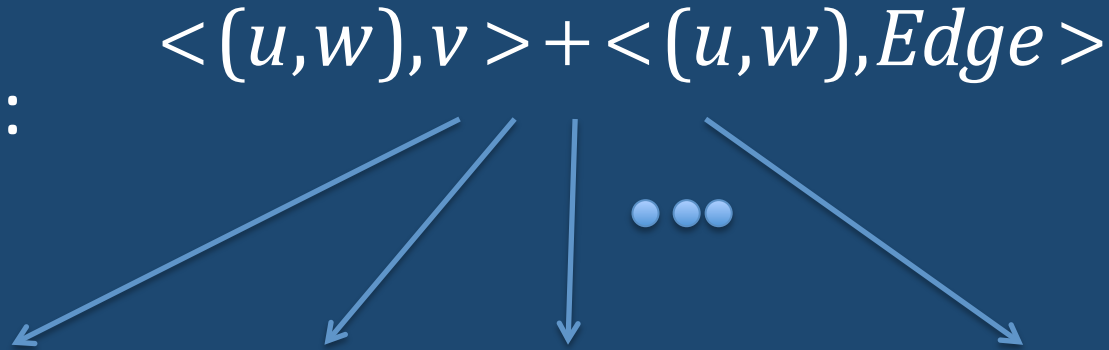


- Reduce 1:

$$\langle v, \Gamma(v) \rangle \longrightarrow \left\{ \langle (u_1, u_2), v \rangle : u_1, u_2 \in \Gamma(v) \right\}$$

MapReduce (Naïve)

- Map 2:



$$\langle (u, w), \{v_1, \dots, v_k, Edge?\} \rangle$$

- Reduce 2:
 - If *Edge* then.
 - For $v \in \{v_1, \dots, v_k\}$ emit $\langle v, 1 \rangle$

What's Wrong with this?

- Does this improve our running time?
- There still may be a very high degree vertex in the network
- Thus, one machine may be stuck with a lot of data!

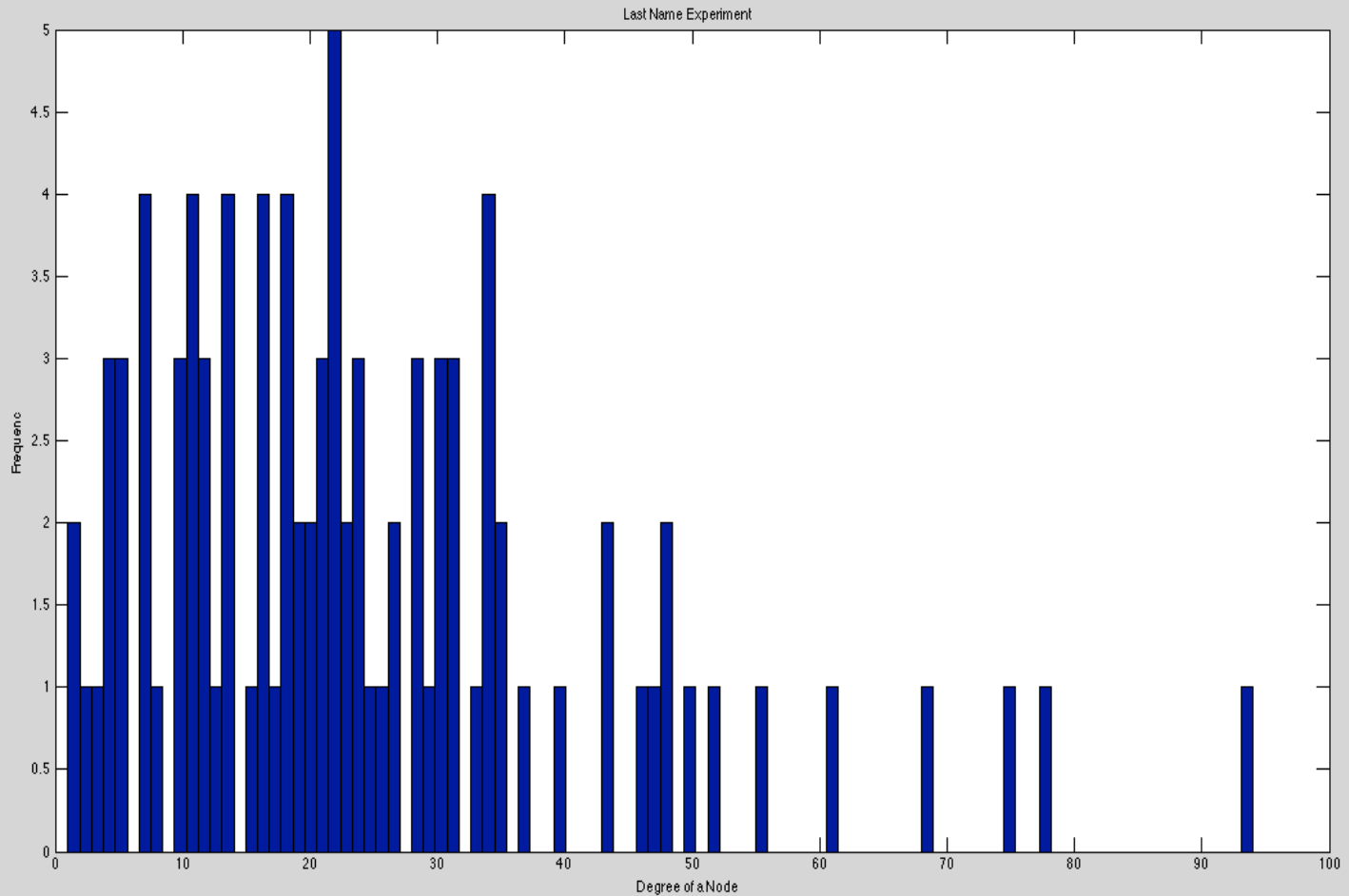
STILL HAS RUN TIME

$$O(d_{\max}^2)$$

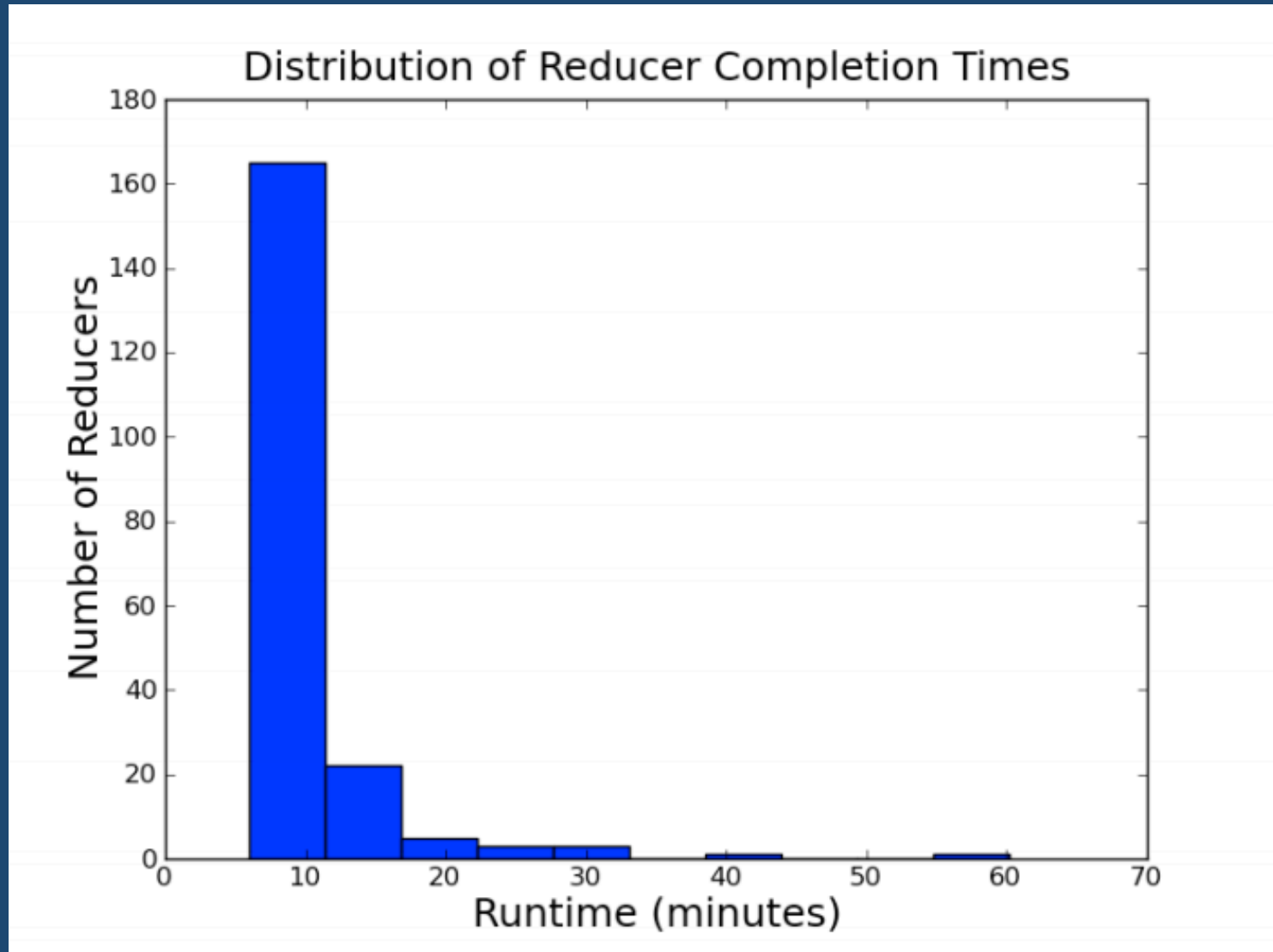
Reality

- Social Networks are typically sparse
- However, there may be few nodes with very high degree.

Reality



Live Journal Data



THE CURSE OF THE LAST REDUCER

- The idea that 99% of the computation finishes quickly, but the last 1% takes a HUGE amount of time.



Possible Fixes

- Generating 2-paths around **high**-degree nodes is expensive – concentrate on **low** degree.



- Divide the graph into **overlapping** subgraphs and somehow account for the overlap.

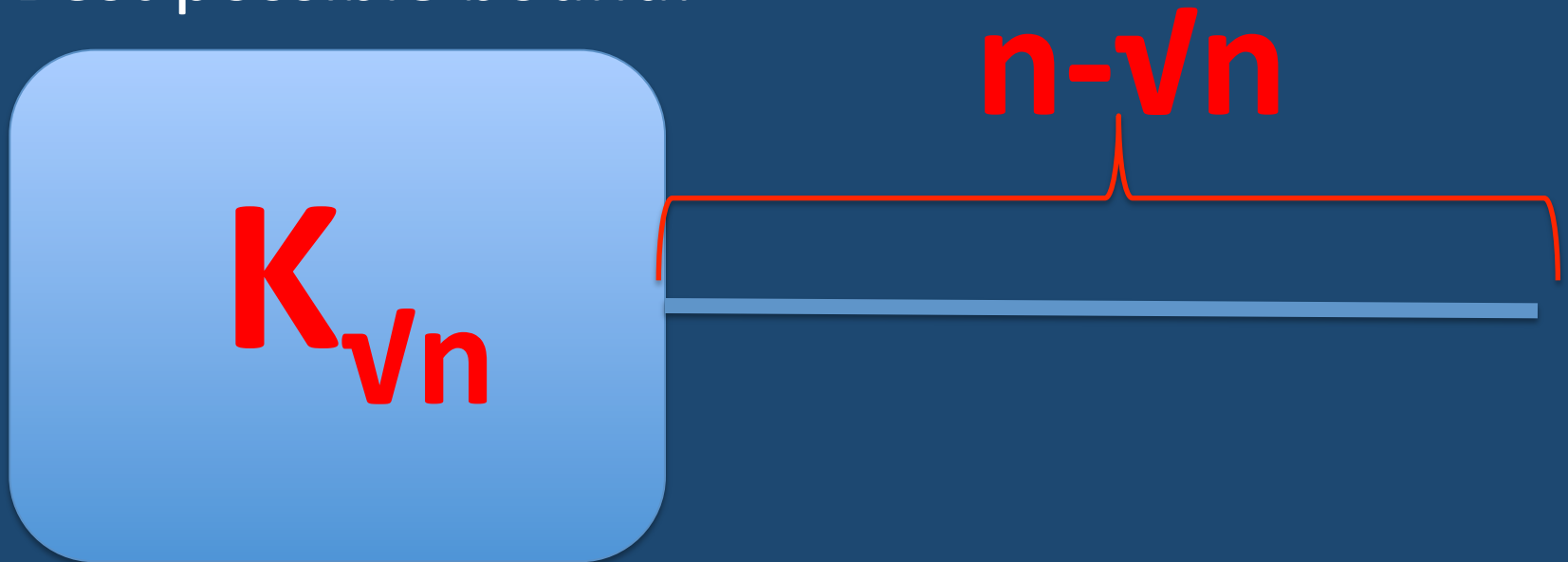
Counting Triangles (Optimal)

- Nodelterator++(V, E)
 - $T \longleftarrow 0$
 - For $v \in V$
 - For $u \in \Gamma(v)$ and $u \succ v$
 - for $w \in \Gamma(v)$ and $w \succ u$
 - » if $(u, w) \in E$
 - $T \longleftarrow T + 1$
- Return T


$$d_u > d_v$$

Properties of Nodelterator++

- Has running time $O(m^{3/2})$ and gives the exact number of triangles incident to each node [Schank '07]
- Best possible bound:



MR-NodeIterator++

- Map 1':
 - If $v \succ u$
 - Emit $\langle u, v \rangle$
- Reduce 1':
$$\langle u, S \subseteq \Gamma(u) \rangle \longrightarrow \left\{ \langle u, (v, w) \rangle : v, w \in S \right\}$$
- Map 2, Reduce 2.

Memory Required per Machine

- Lemma: The input to any reduce instance in first round has $O(\sqrt{m})$ edges (Sublinear space)
- Proof:

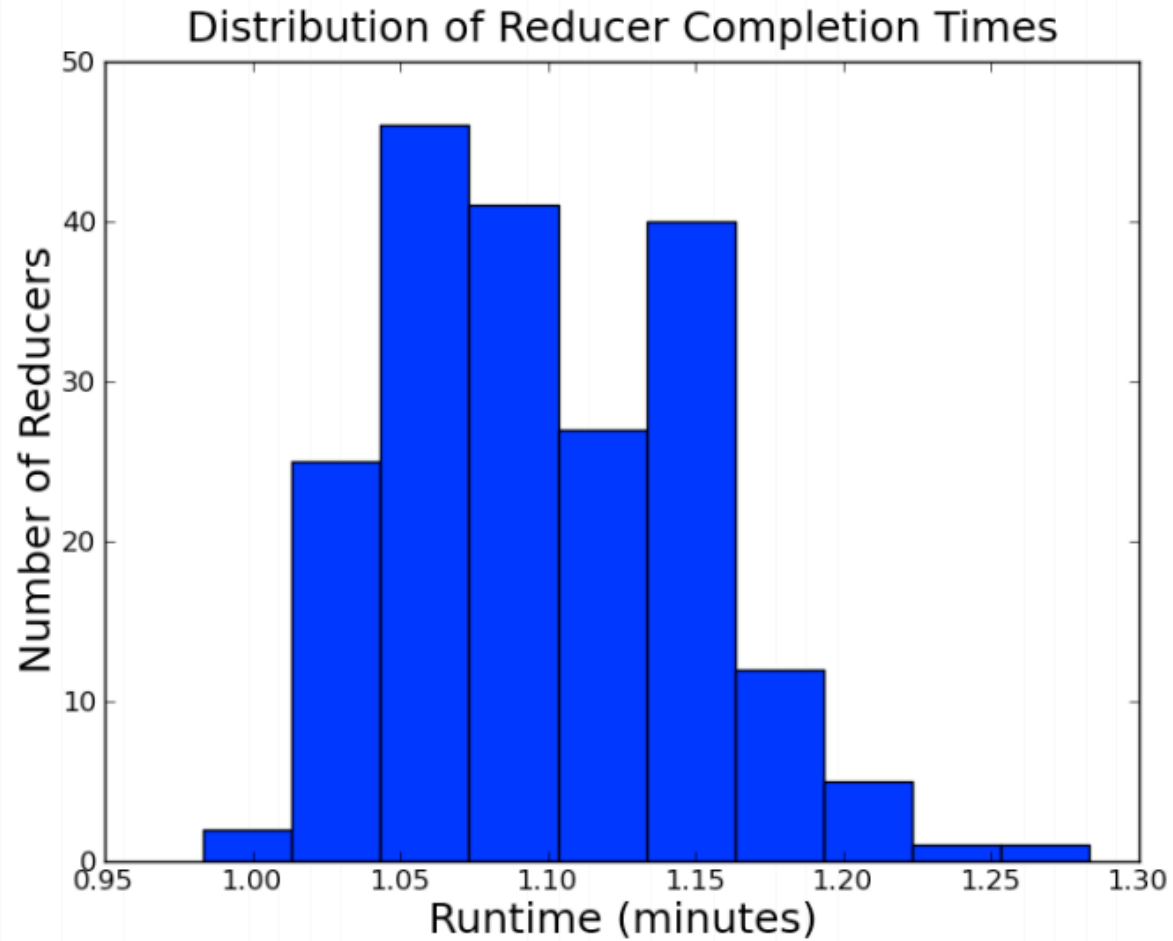
$$\mathcal{L} = \left\{ v \in V : d_v < \sqrt{m} \right\}$$

$$\mathcal{H} = \left\{ v \in V : d_v \geq \sqrt{m} \right\}$$

Size of Output after Round 1

- Lemma: The total number of records output at the end of the first reduce is $O(m^{3/2})$
- Proof:
 - There are at most $n = O(m^{1/2})$ machines with low degree nodes, and each machine produces an output of size $O(m)$
 - There are at most $O(m^{1/2})$ machines with high degree nodes and each machine must output pairs with other high degree nodes $\Rightarrow O(m)$ output size

Did it Help?



Possible Fixes

- Generating 2-paths around **high**-degree nodes is expensive – concentrate on **low** degree.



- Divide the graph into **overlapping** subgraphs and somehow account for the overlap.

MR-GraphPartition

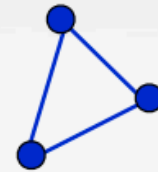
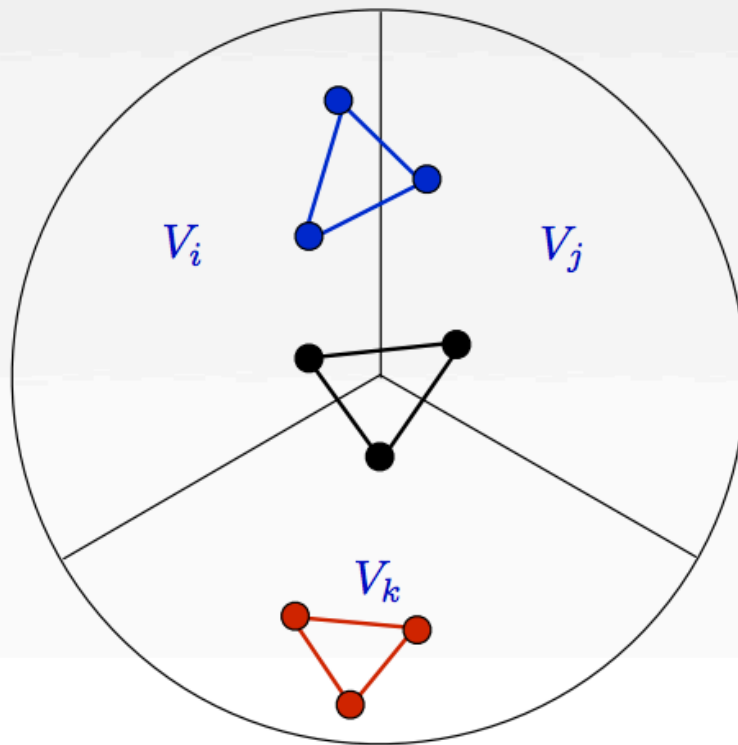
- Input: (V, E, ρ)
- Partition vertices into ρ equal sized $V_0, \dots, V_{\rho-1}$
- Consider all triples (V_i, V_j, V_k) and the induced graph $G_{ijk} = G[V_i, V_j, V_k]$ for $i < j < k$
- Compute Triangles on each graph separately
 - You can use your favorite triangle counting algorithm on each!
- Map nodes to index i by using a universal hash

MR-GraphPartition

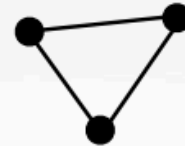
- Map 1'': Input $\langle (u,v), 1 \rangle$
 - for $a < b < c \leq \rho - 1$
 - if $\{h(u), h(v)\} \subseteq \{a, b, c\}$
 - emit $\langle (a, b, c), (u, v) \rangle$
- Reduce 1'': Input: $\langle (i, j, k), E_{ijk} \rangle$
 - Count Triangles and weight accordingly.



May Over Count # of 's



in $p-2$ subgraphs



in 1 subgraph



in $\sim p^2$ subgraphs

Can count exactly how many subgraphs each triangle will be in

Analysis

- The expected size of the input to any machine instance is $O(m / \rho^2)$
- The expected total space used at the end of map phase is $O(m\rho)$
- Proof: SEE BOARD

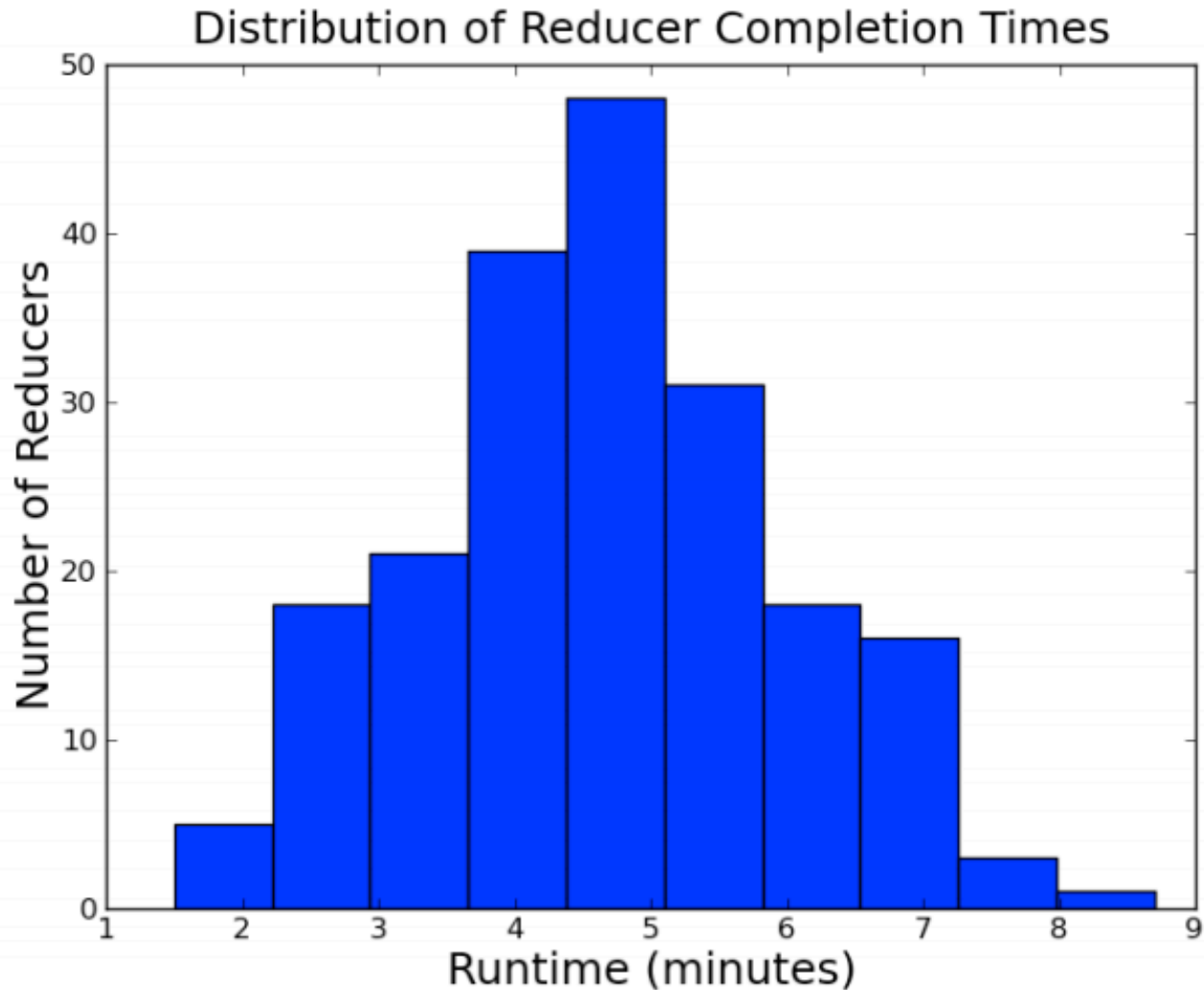
Analysis (continued)

- Theorem: For $\rho \leq \sqrt{m}$, the amount of work done by all the machines is $O(m^{3/2})$
- Proof:
 - $O(1)$ time per edge $\Rightarrow O(m\rho) = O(m^{3/2})$ time for Map 2" phase.

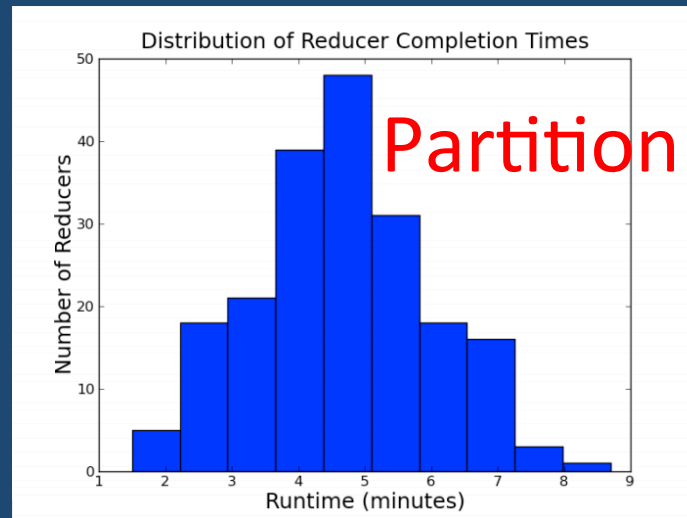
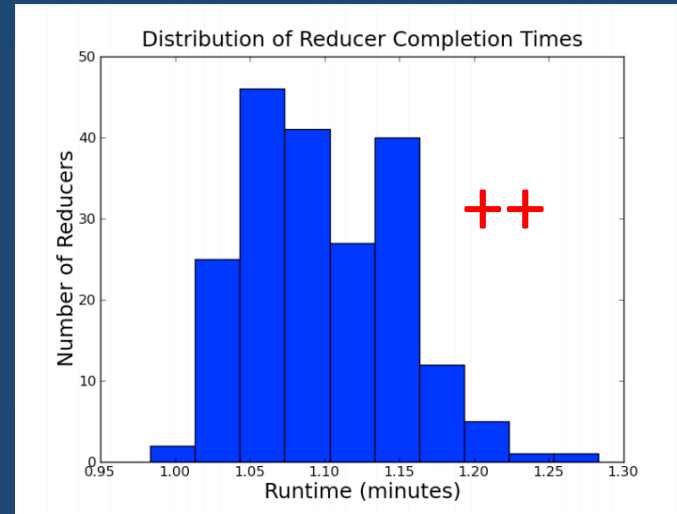
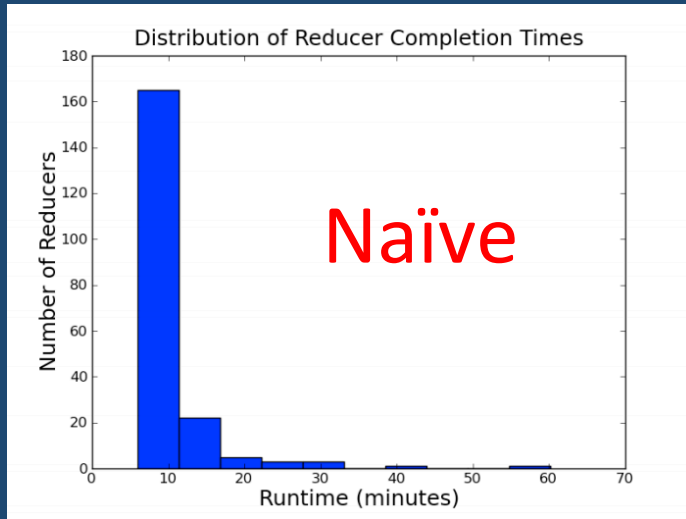
Partition input amongst $O(\rho^3)$ reducers.

Running Time per Reducer:
$$= O(\#Edges^{3/2}) = O\left(\left(\frac{m}{\rho^2}\right)^{3/2}\right)$$

Results for Partition



Comparison of Results





**THE CURSE OF THE
LAST REDUCER**

Questions???