

Jaccard Coefficients as a Potential Graph Benchmark

Peter M. Kogge
McCourtney Prof. of CSE
Univ. of Notre Dame
IBM Fellow (retired)



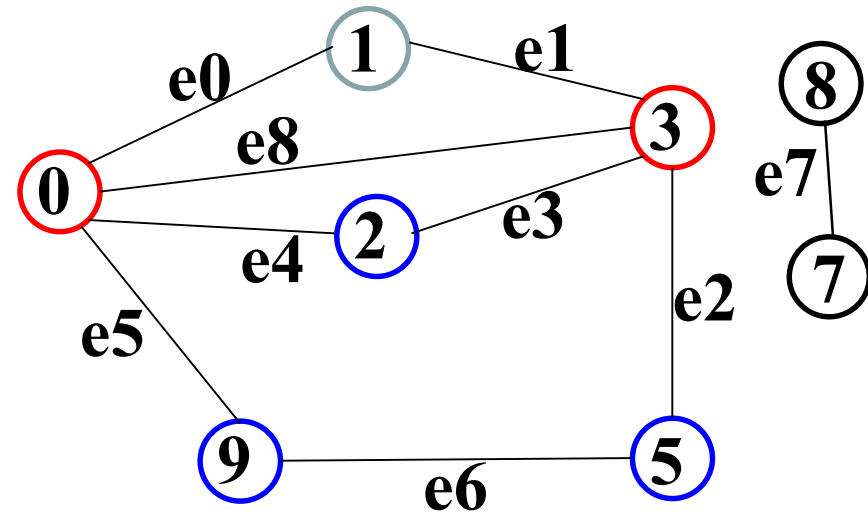
Outline

- Motivation
- Jaccard Coefficients
- A MapReduce Baseline
- Variations
- A Possible Heuristic
- Suggested Benchmarks
- Key Questions



Graph 500

- Start with a root, find reachable vertices
- Simplifications: only 1 kind of edge, no weights
- Performance metric: TEPS: Traversed Edges/sec



Level	Scale	Size	Vertices (Billion)	TB	Bytes /Vertex
10	26	Toy	0.1	0.02	281.8048
11	29	Mini	0.5	0.14	281.3952
12	32	Small	4.3	1.1	281.472
13	36	Medium	68.7	17.6	281.4752
14	39	Large	549.8	141	281.475
15	42	Huge	4398.0	1,126	281.475
				Average	281.5162

Scale = $\log_2(\# \text{ vertices})$

Starting at 1: 1, 0, 3, 2, 9, 5



Limitations of BFS as a Benchmark

Has provided rich set of new algorithms & implementations, BUT:

- Complexity only $O(E)$
 - $E = \#$ edges
- Only a batch algorithm
 - Must investigate virtually entire static graph
- No natural incremental variant
- Little non-academic applicability
 - Commercial apps much more focused on limited neighborhoods



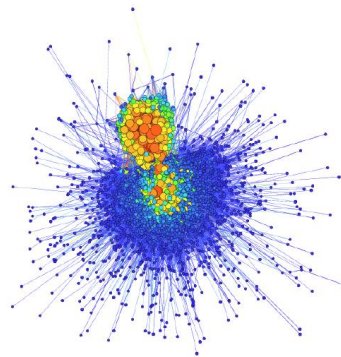
Big Graph Relationship Problems

Social scale. . .

NSA-RD-2013-056001v1

1 billion vertices, 100 billion edges

- 111 PB adjacency matrix
- 2.92 TB adjacency list
- 2.92 TB edge list



Twitter graph from Gephi dataset (<http://www.gephi.org>)

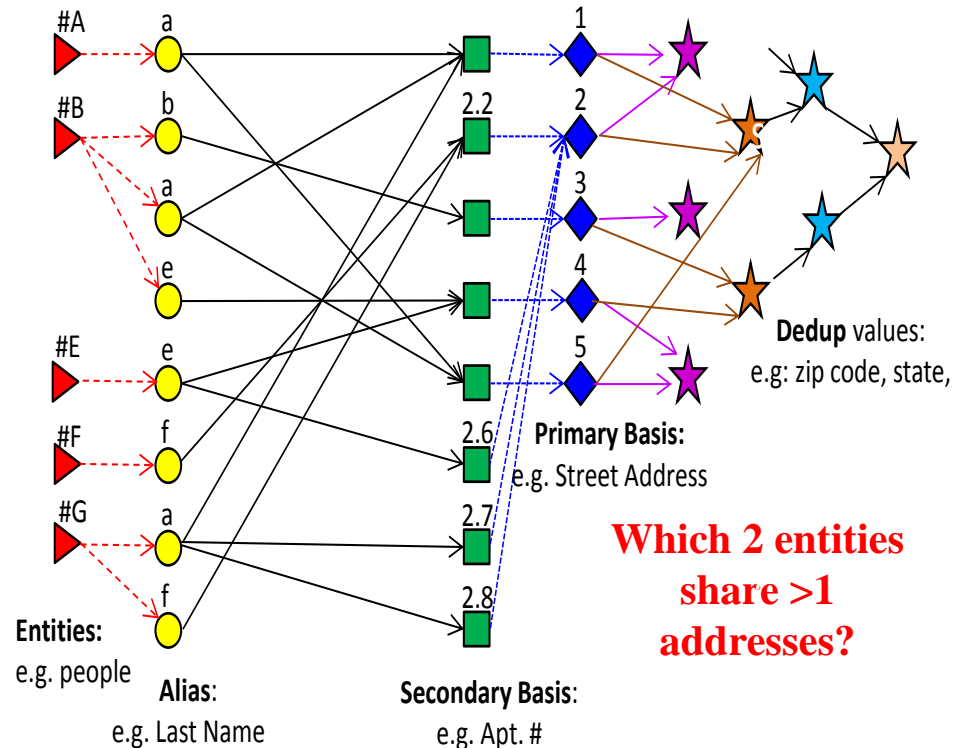


What are the pairs of people that follow the same set of twitter feeds

Paul Burkhardt, Chris Waring

An NSA Big Graph experiment

from "Burkhardt & Waring, An NSA Big Graph Experiment"
http://www.pdl.cmu.edu/SDI/2013/slides/big_graph_nsa_rd_2013_56002v1.pdf



- Tough Problem: Find vertex **pairs** that "share some common property"
- Related graph problem: computing "Jaccard coefficients"
- Commercial version: "Non-obvious Relationship Problems" (NORA)



Sample Real World Problem

Auto Insurance Co: “Tell me about giving auto policy to Jane Doe” in < 0.1sec

- 40+ TB of Raw Data
 - Periodically clean up & combine to 4-7 TB
 - Weekly “**Boil the Ocean**” to precompute answers to all standard queries
 - Does W have financial difficulties?
 - Does X have legal problems?
 - Has Y had significant driving problems? **Relationships**
- Who has shared addresses with Z?**
- ...

Look up answers to precomputed queries for “Jane Doe” and combine

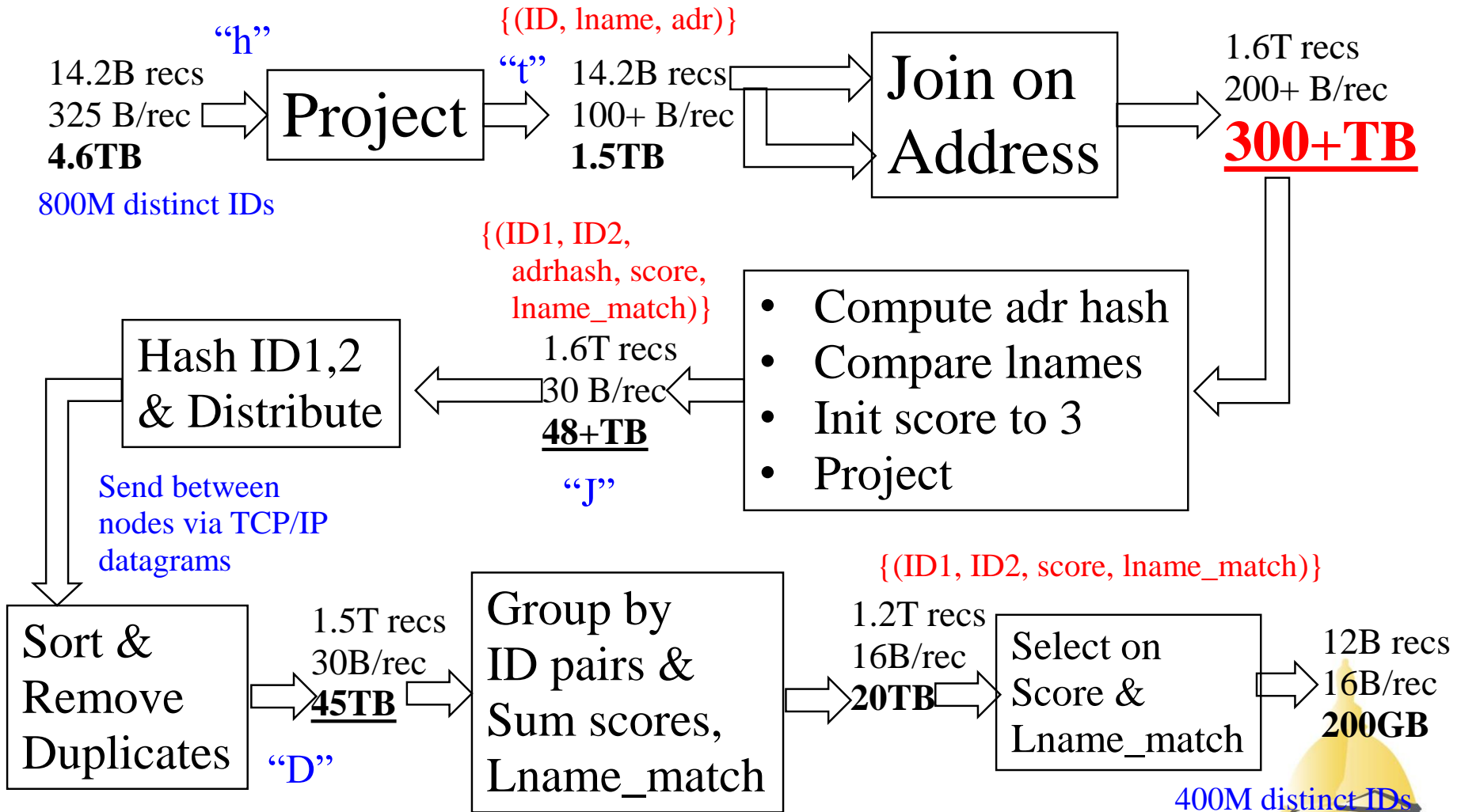
“Jane Doe has no indicators
But
she has shared multiple addresses with Joe Scofflaw
Who has the following negative indicators”

A 2012 Relationship

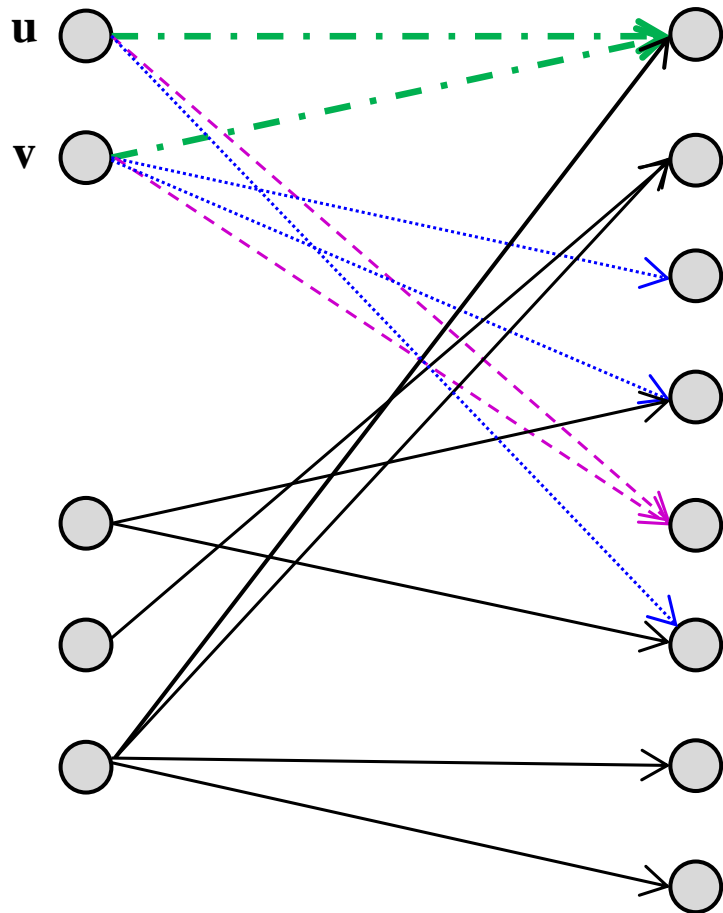
- Given: 14.2 billion records from
 - 800 million **entities** (North American people, businesses)
 - 100 million **addresses**
- Goal: given entity, find all other entities that
 - Share at least 2 addresses in common
 - Or have one address in common and last name that is “close”
- Matching last names requires processing to check for typos (“Levenshtein distance”)
- Above one of dozens or relationships



2012 Processing Flow



Jaccard Coefficient $\Gamma(u, v)$



$N(u)$ = set of neighbors of u

$\Gamma(u, v)$ = fraction of neighbors of u and v that are in common

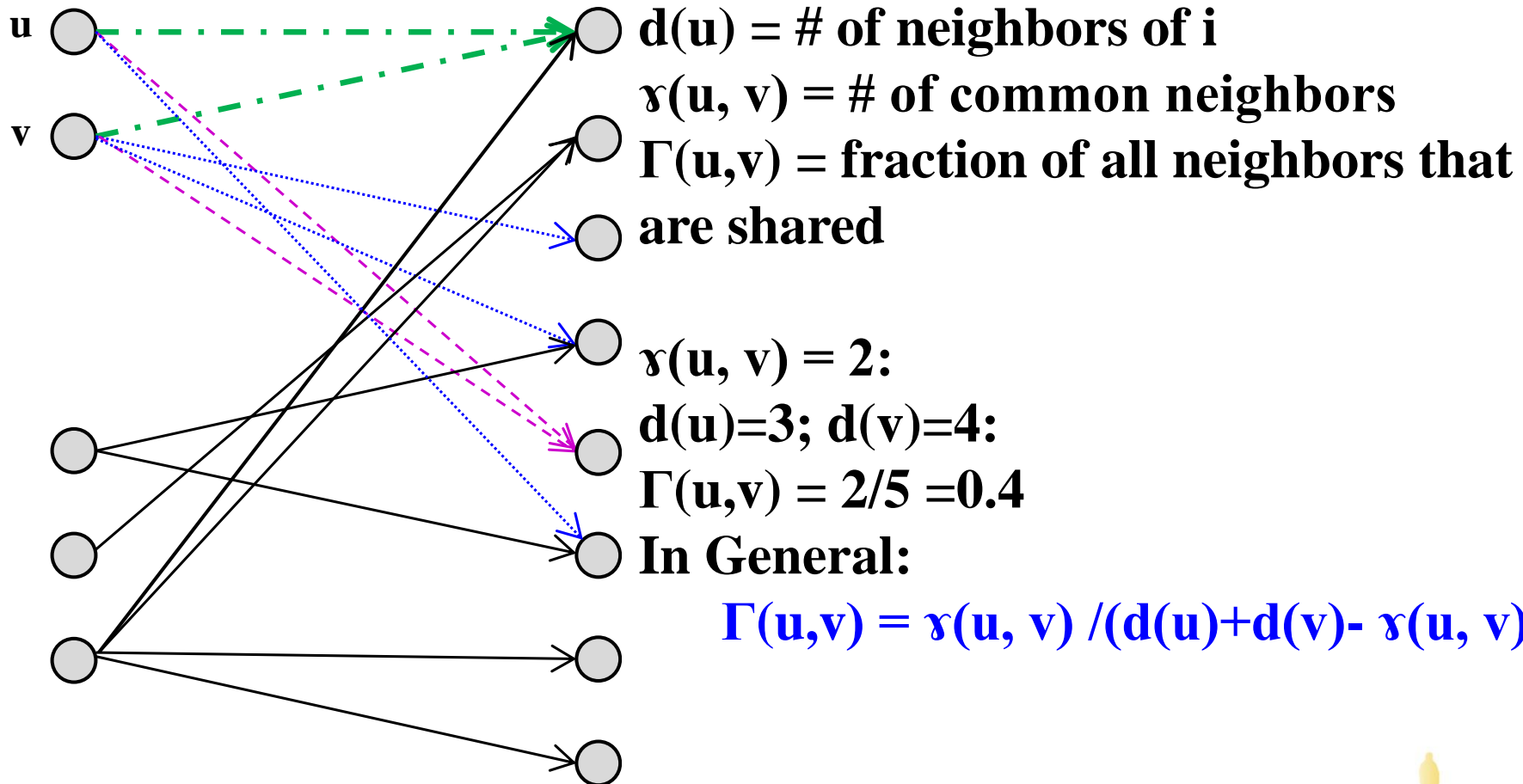
$$\Gamma(u, v) = \frac{|N(u) \cap N(v)|}{|N(u) \cup N(v)|}$$

Green and Purple lead to common neighbors

Blue lead to non-common neighbors



Jaccard Coefficient $\Gamma(u, v)$



Green and Purple lead to common neighbors
Blue lead to non-common neighbors



Computing a Single γ

- $\gamma(u,v) = \#$ of common neighbors
- Do not need to “visit” neighbors, only enumerate them
- Algorithm requires comparing 2 lists of lengths $d(u)$ & $d(v)$
 - Let $d = \max(d(u), d(v))$
- If lists are sorted, then $O(d)$
- If lists not sorted, then $O(d \cdot \log(d))$
 - Clever hash algorithms may reduce to almost $O(d)$



Computing all γ s

- Apply prior single γ to all pairs
 - $O(V^2d)$ to $O(V^2d \cdot \log(d))$
 - Factor of 2 reduction if only do (u,v) where $i < j$
- Avoid computing all clearly 0 terms
 - For each u , explore each w , (u,w) an edge
 - Compute $\gamma(u,v)$ for each v where (v,w) an edge & $u < v$
 - $O(V \cdot d^3)$ to $O(V \cdot d^4)$
- Backwards: compute γ s incrementally
 - Group edges so for each w we have $\{x | (x,w)\}$
 - For each x & y in this set, increment $\gamma(x,y)$
 - $O(Vd^2)$: Avoid $O(V^2)$ initialization by dynamic creation & initialization
- Sparse Matrix Equivalency: $A =$ adjacency matrix
 - $\gamma = AA^T$; $O(\text{nnz}(A))$

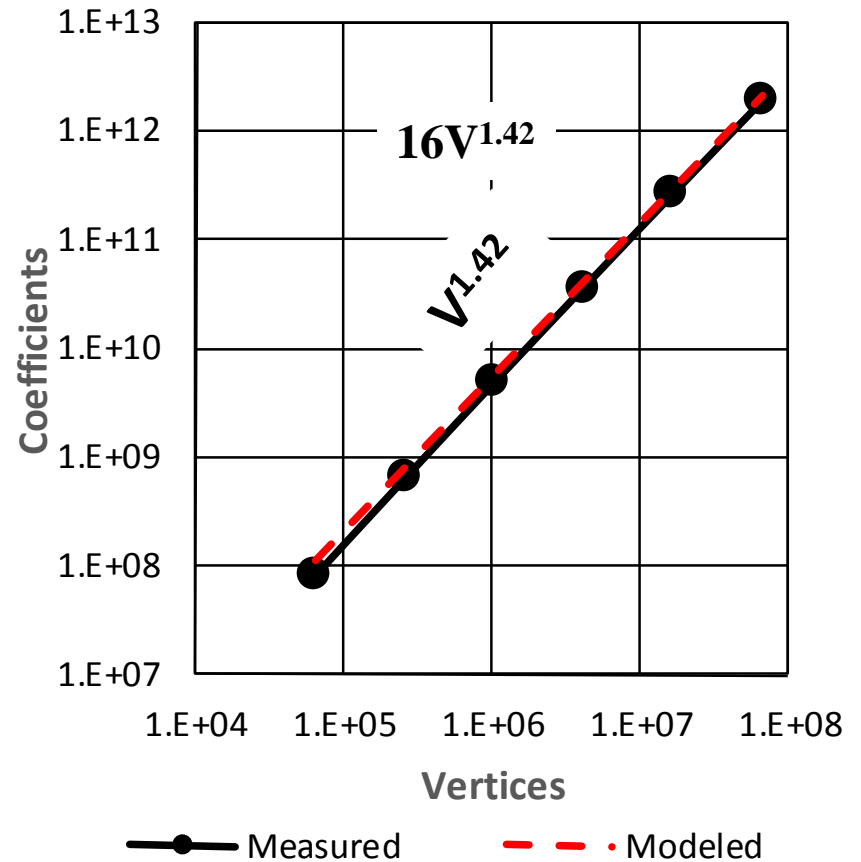
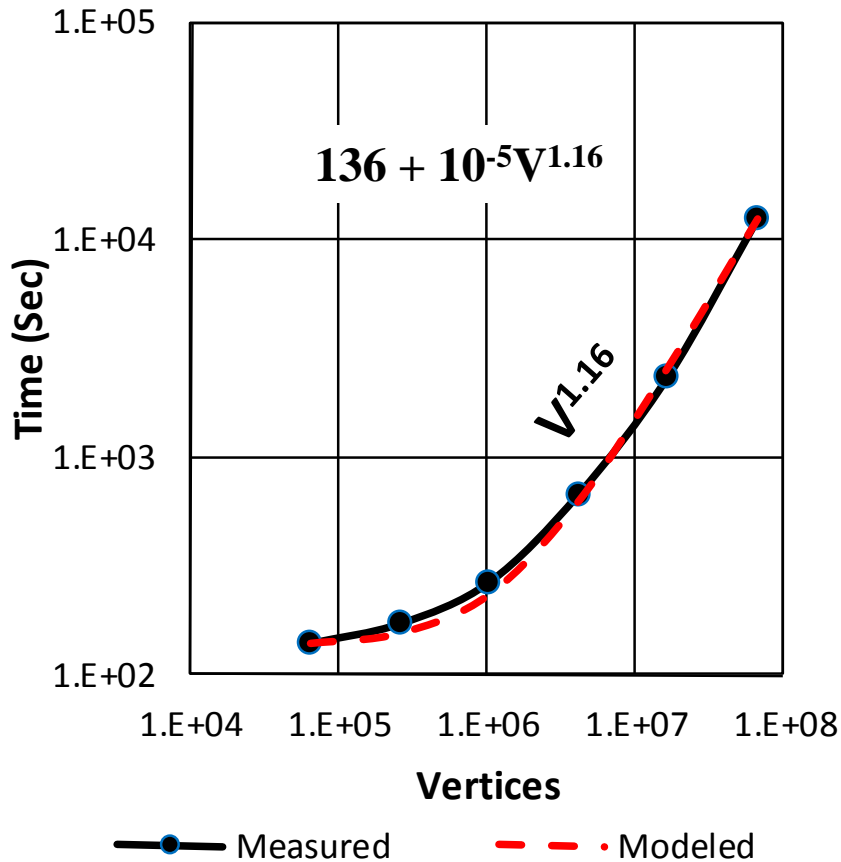


Jaccard via MapReduce

- 1 MapReduce step: 3 phases
 - **Map** some function over all data to (key,value) stream
 - **Group** pairs by key
 - **Reduce** each group
- Two reported Jaccard implementations
 - 3-steps with $V = 1$ million edges on 50 node cluster
 - 5-steps with $V = 64$ million vertices & 1 billion edges
 - Burkhardt “Asking Hard Graph Questions,” Beyond Watson Workshop, Feb. 2014.
 - RMat matrices, average $d(i) = 16$
 - 1000 node system, each with 12 cores & 64GB



Measurements & Trends



JACS (Jaccard Coefficients / Sec) = $1.6E6 * V^{0.26}$



Relevant Batch Variations

- Consider graphs with 2 vertex classes
 - A = set of “people”, B = set of “addresses”
 - Edges from A to B is person a “resided at” address b
 - $\gamma(a,b)$ = # of common addresses between a and b
- Real world: for $|A| = 8E8$, $|B| = 1E8$,
 - $|\gamma| = 1.2E12 = |A|^{1.35}$
- Variations:
 - >2 classes of vertices
 - Exclude paths thru high in-degree neighbors
 - Weight paths on basis of properties (e.g. same last name)
 - Add threshold



Dynamic Variations

- Many applications with dynamic graphs
 - Vertices, edges added/deleted dynamically
- Question: which γ s change with edge addition or deletion
 - Perhaps just those that cross threshold, in either direction



Other Variations

- Don't compute/store all possibly $O(V^2)$ γ s,
But store just statistics of set of γ s for each vertex
 - Number of non-zero, largest, average, etc.
 - Reduces storage to $O(V)$
 - Serve as starting point for more complex queries
- Expand “neighbors” beyond “1 hop”
 - Real commercial applications at “1.5” hops



A Possible Heuristic

- Goal: constant time elimination of γ s that are zero or less than some threshold
- Build bit vector for each vertex that hashes vertex ids into bits
 - Bit $i = 1$ if one or more neighbors fall into set I
- Estimate $\gamma(u,v)$ by ANDing u & v 's bit vectors
 - If results are 0, then $\gamma = 0$
- Can also estimate upper bound on $\gamma(u,v)$
 - Simple function of # of 1s in bit vectors
- If benchmark uses a threshold, this estimate can prevent computation when $\text{bound} < \text{threshold}$
- Based on ideas from SpGEMM package



Suggested Benchmarks

- Simple batch: compute all non-zeros
 - Perhaps use same RMAT as for BFS
 - Performance Metric JACS
- Multi-class benchmark
 - One metric: time to compute all γ s as above
 - Optionally append Jaccard statistics as properties to each vertex of a class
- Real-time event detection
 - Start with precomputed γ statistics
 - Input stream of additions/deletions
 - Check for changes in γ s
 - Performance metric: change throughput



Key Questions for More Precise Benchmark Definitions

- Formal definition of batch and incremental
- Where do non-zeros go? File, properties, ...
- Tie data sets to real applications
- Explore real-world distribution of γ s and Γ s to understand growth rates better
- Build in options such as thresholds
- Develop complexity models, esp. parallel
- Verify correct solutions
- Develop reference implementations



Acknowledgements

- Funded in part by the Univ. of Notre Dame, Notre Dame, IN, USA

