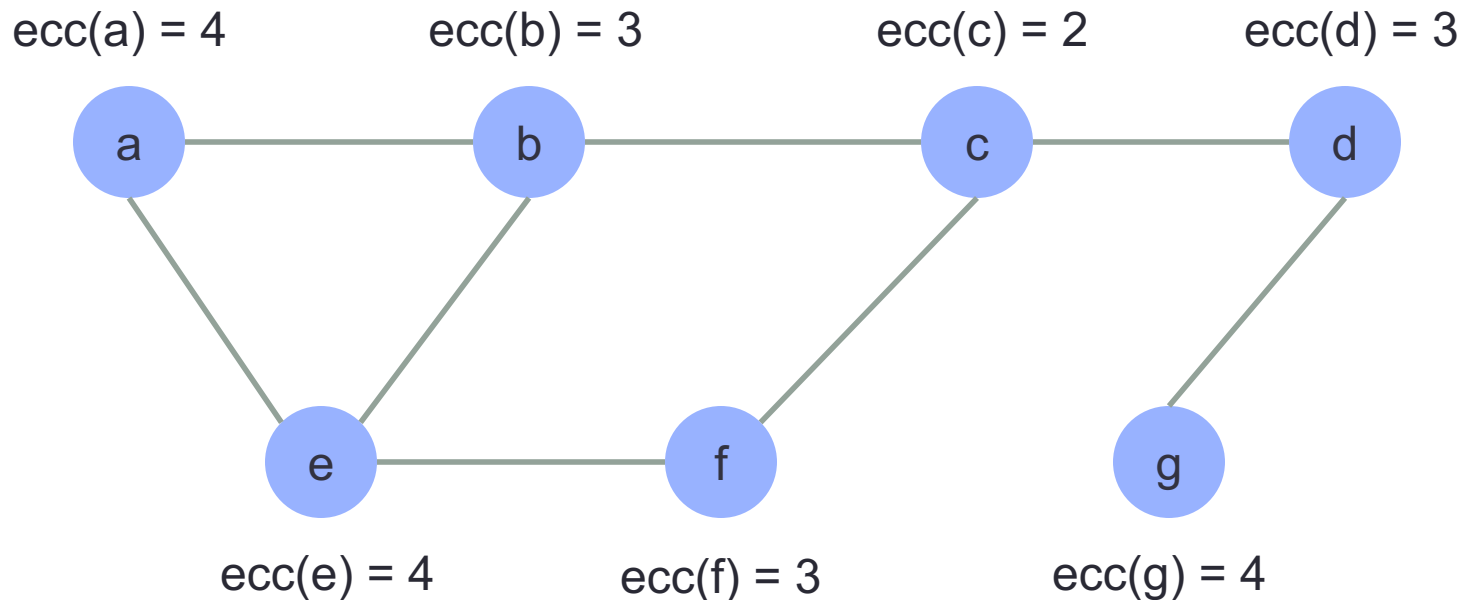


An Evaluation of Parallel Eccentricity Estimation Algorithms on Undirected Real-World Graphs

Julian Shun

Graph Eccentricities

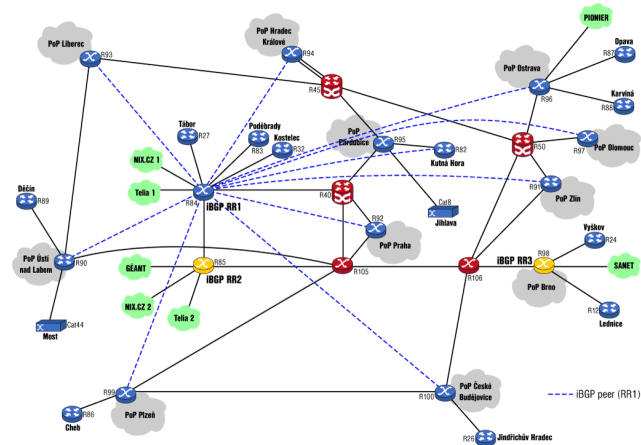
- The **eccentricity** of a vertex v is the distance to furthest reachable vertex from v



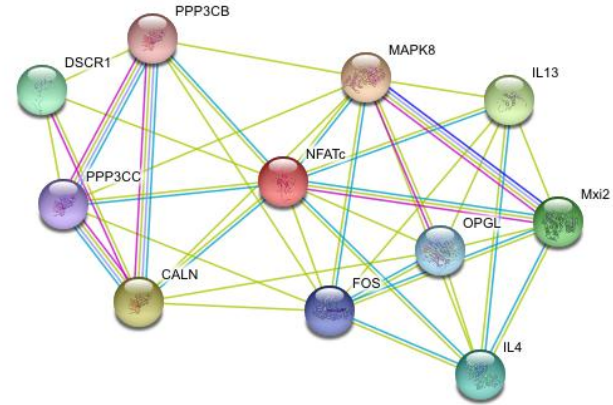
- The **diameter** of a graph is the maximum eccentricity value
- Extends to directed and/or weighted graphs

Applications

Routing networks



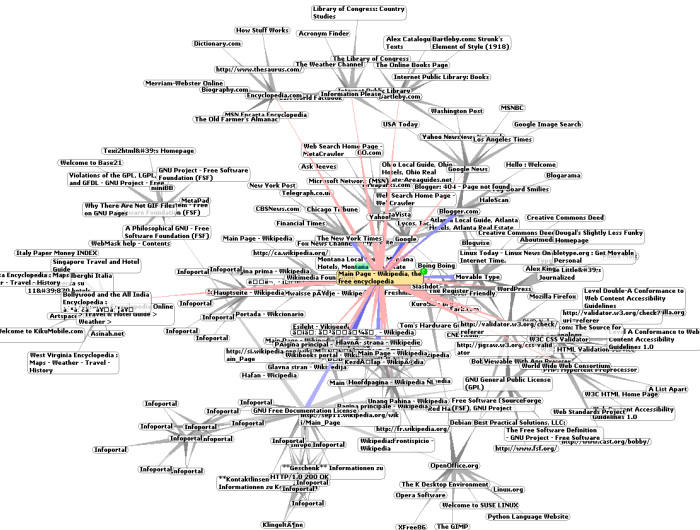
Protein networks



Social networks



Web graphs

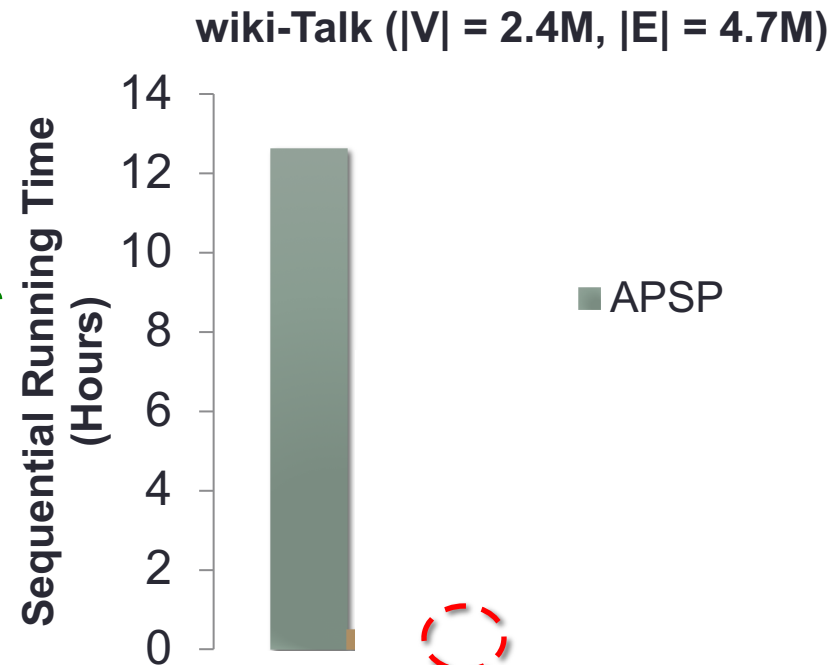


- *Core-periphery structure*

How to compute eccentricities?

Experiments done on one thread of a 40-core (2.4 GHz) Intel Nehalem machine with 256 GB memory

- All pairs shortest paths (APSP)
 - At least quadratic work
 - Would take over 2670 years for Yahoo! Web graph (1.4B vertices, 6.4B edges)
- Takes and Kosters [*Algorithms 2013*] algorithm
 - Quadratic work in the worst case
- Approximation algorithms
 - Orders of magnitude faster
 - Can process Yahoo! Web graph in minutes
- Parallelism



4 seconds, < 0.0001 average relative error

This work

- First comprehensive comparison of **parallel** implementations of eccentricity algorithms on **large** (undirected/unweighted) real-world graphs

Algorithm	Guarantee	Complexity
APSP	Exact	$O(E V)$
Takes-Kosters [<i>Algorithms 2013</i>]	Exact	$O(E V)$
Single BFS	2-approx	$O(E)$
Chechik et al. [<i>SODA 2014</i>]	1.66-approx	$O(E (V \log V)^{1/2})$
Roditty-Vassilevska Williams [<i>STOC 2013</i>]	1.5-approx	$O(E (V \log V)^{1/2})$
ANF/HADI [Palmer et al. <i>KDD '02</i> , Kang et al. <i>TKDD '10</i>]		$O(k E D)$
HyperANF [Boldi et al. <i>WWW 2011</i>]		$O(k (\log \log V / \log V) E D)$
k-BFS [this work]		$O(k E \min(1, D / \log V))$

$|V|$ = # vertices

$|E|$ = # edges

k = # probabilistic counters/BFS's

D = diameter

- Simple shared-memory implementations in the Ligra graph processing framework [*Shun and Blelloch PPOPP 2013*]
- k-BFS as a parallel primitive for fast, scalable, and accurate eccentricity estimation

Simple 2-approximation

- **Step 1:** Perform breadth-first search (BFS) from arbitrary vertex v to compute eccentricity
- **Step 2:** Assign all vertices $\text{ecc}(v)$
- **Guarantee:** $\text{ecc}(w)/2 \leq \hat{\text{ecc}}(w) \leq 2 \text{ecc}(w)$ for all w
- BFS complexity is $O(|E|)$ and easily parallelizable with parallel time proportional to diameter

This algorithm is fast but its estimates are useless!

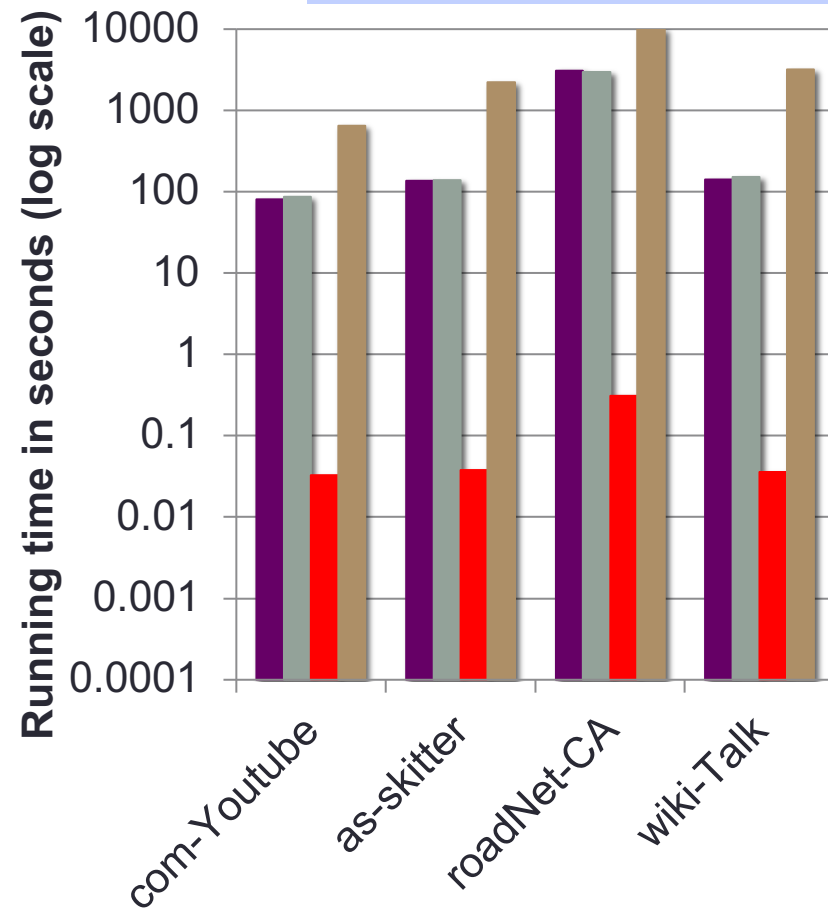
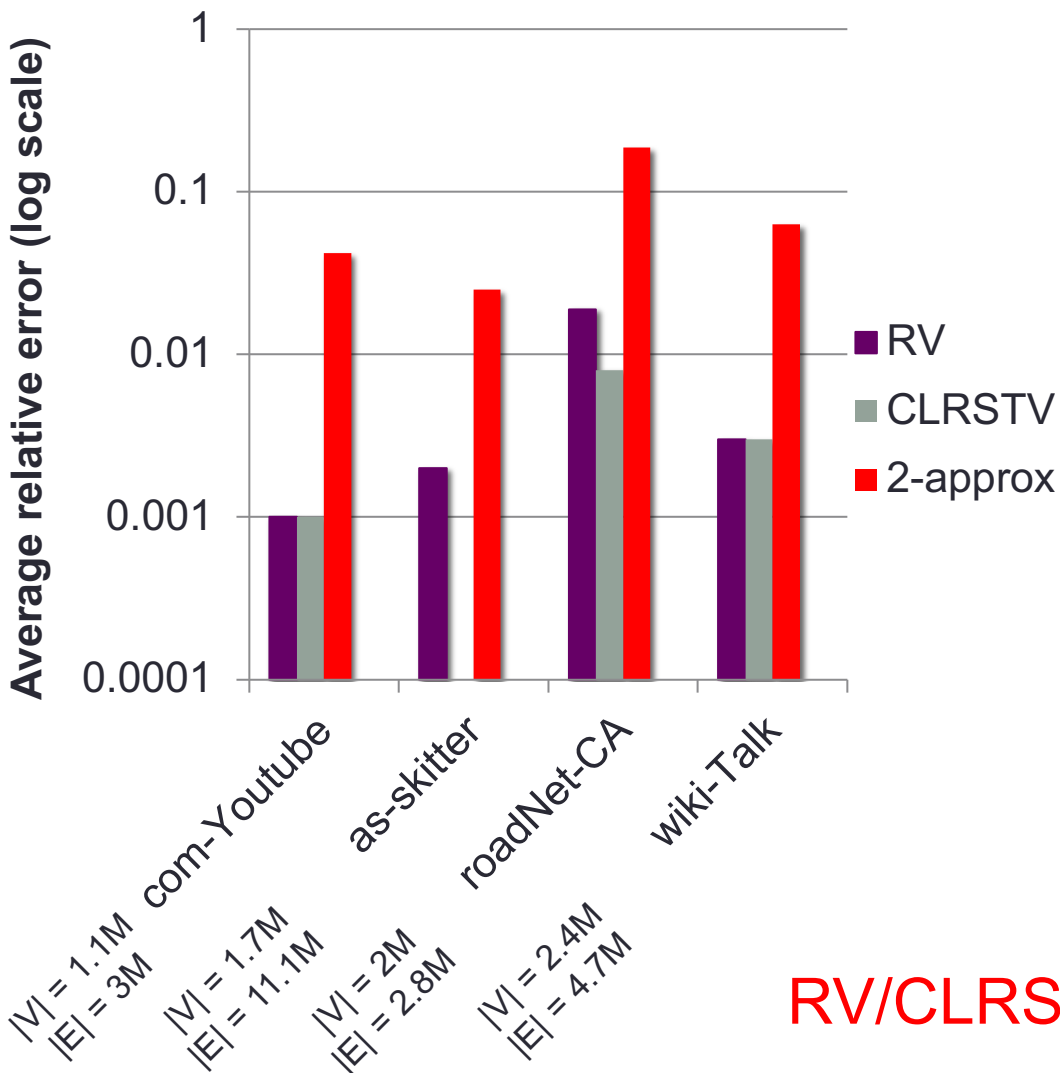
Stronger provable approx algorithms

- RV: [Roditty-Vassilevska Williams *STOC 2013*]
 - $(2/3) \text{ecc}(v) \leq \hat{e}c(v) \leq (3/2) \text{ecc}(v)$
 - Complexity: $O(|E| (|V| \log |V|)^{1/2})$
- CLRSTV: [Chechik, Larkin, Roditty, Schoenebeck, Tarjan, and Vassilevska Williams *SODA 2014*]
 - $(3/5) \text{ecc}(v) \leq \hat{e}c(v) \leq \text{ecc}(v)$
 - Complexity: $O(|E| (|V| \log |V|)^{1/2})$
- We provide the first empirical evaluation of these algorithms
 - Shared-memory parallel implementations in Ligra

How do they perform in practice?

• Average relative error = $\frac{1}{|V|} \sum_{v \in V} \left| \frac{e\hat{c}c(v) - ecc(v)}{ecc(v)} \right|$

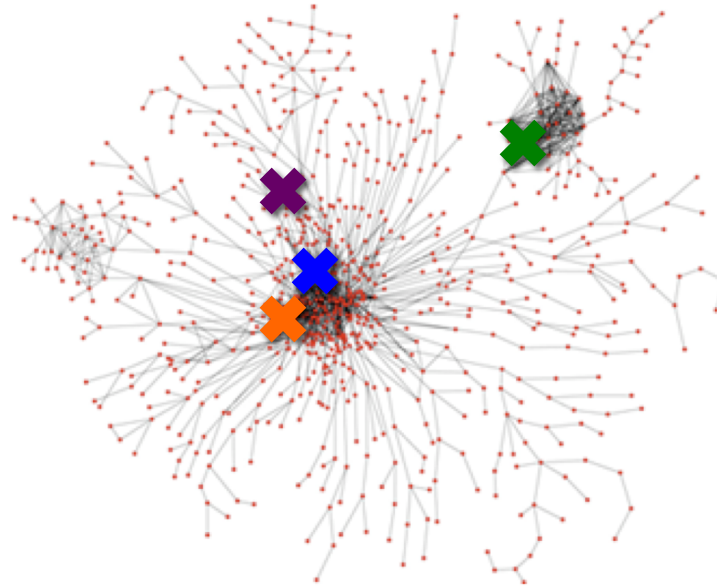
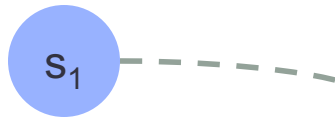
Experiments done on a 40-core Intel Nehalem machine with 256 GB memory



RV/CLRSTV accurate but not scalable

Multiple BFS's

- Run multiple BFS's from a sample of vertices and compute the maximum eccentricity



Source: C.A. Hidalgo, B. Klinger, A.-L. Barabási, R. Hausmann. *Science* 317 (2007)

$$\hat{ecc}(v) = \max(d(v, s_1), d(v, s_2), d(v, s_3), d(v, s_4)) \quad \text{for all } v$$

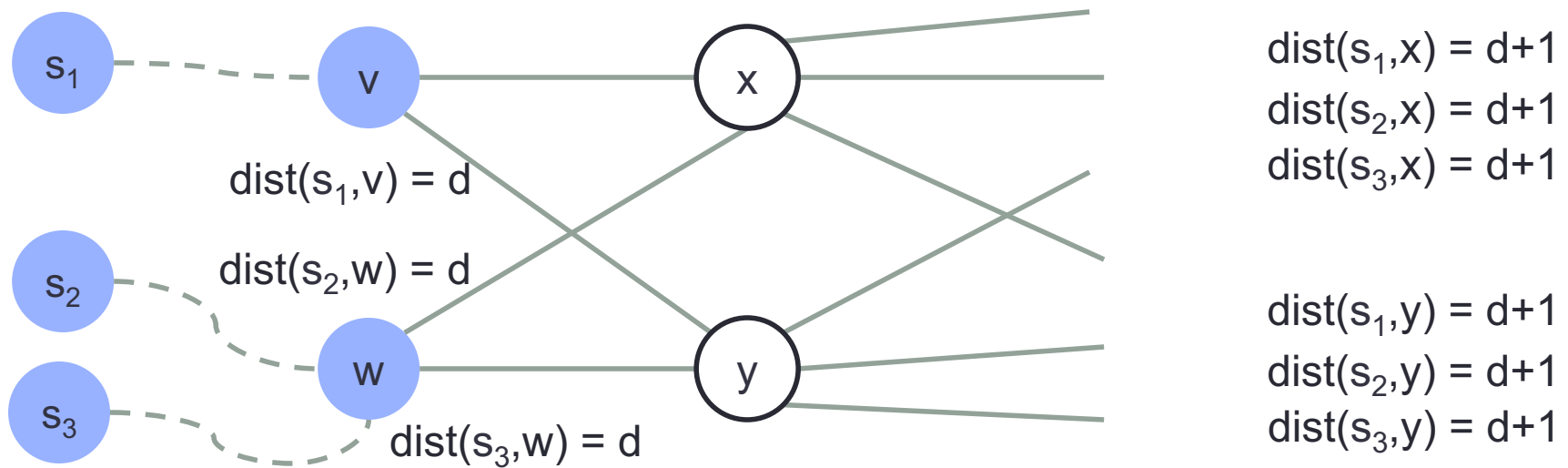
- Does not work that well in practice
- Double-sweep: find furthest vertices from sample of random vertices, then run BFS's from the set of furthest vertices (used before for diameter estimation [Corneil et al. 2001, Magnien et al. 2009])

k-BFS Implementation (second sweep)

- k BFS sources
- Each vertex \mathbf{v} stores estimate $e\hat{c}(\mathbf{v}) = 0$
- Run parallel BFS from each source, updating each encountered vertex \mathbf{v} at distance \mathbf{d} to $\max(e\hat{c}(\mathbf{v}), \mathbf{d})$
- **Observation:** There is shared work among the different BFS's
 - Vertices updated multiple times
 - Vertices placed on frontiers multiple times \rightarrow edges traversed multiple times

k-BFS Implementation

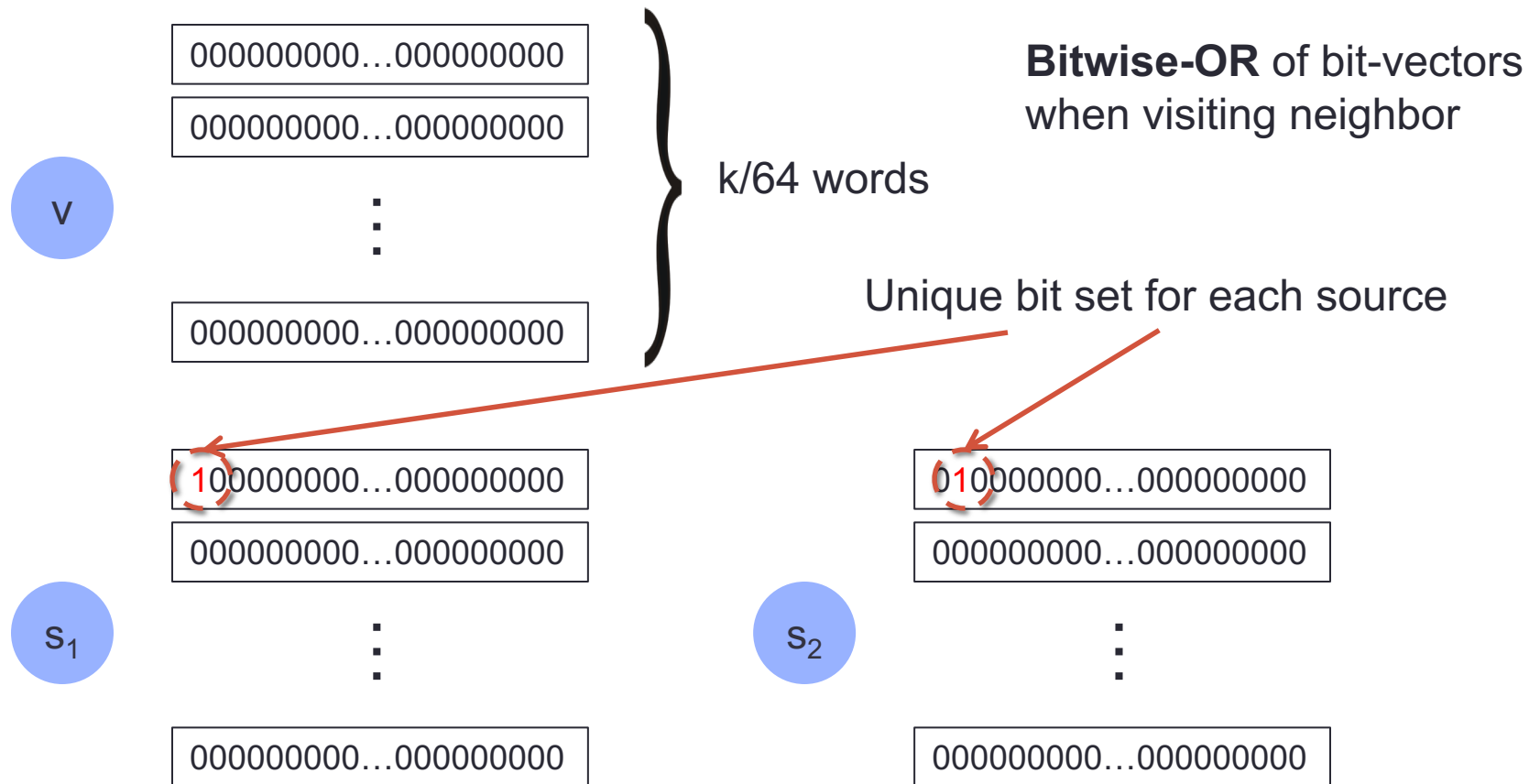
- **Observation:** There is shared work among the different BFS's



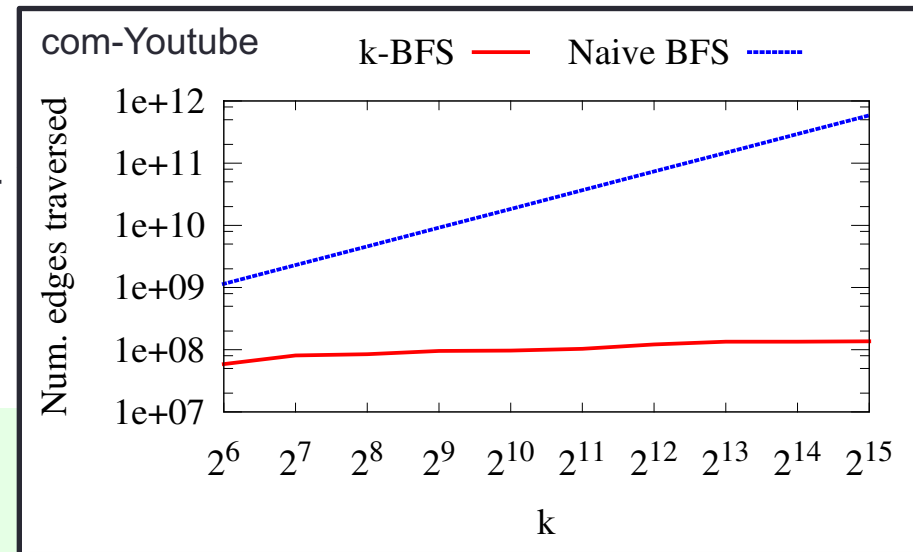
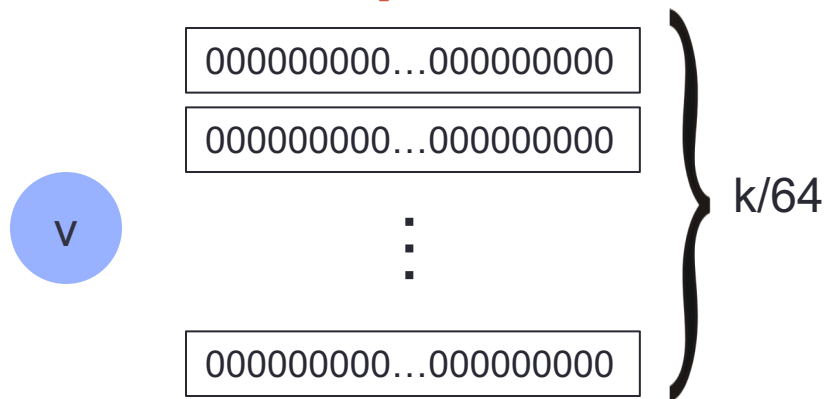
- **Goal:** Reduce redundant computation on vertices
- **Goal:** Reduce the number of times each visited vertex is placed onto the frontier

k-BFS Implementation

- Run all k BFS's simultaneously
- Take advantage of bit-level parallelism to store "visited" information



k-BFS Implementation



- Initial **frontier** = $\{s_1, s_2, \dots, s_k\}$

- $d = 0$

- While **frontier** not empty:

- nextFrontier** = $\{\}$

- $d = d+1$

- For each vertex v in **frontier**:

- For each neighbor ng :

- Do bitwise-OR of v 's words with ng 's words and store in ng

- If any of ng 's words changed:

- $e\hat{c}(ng) = \max(e\hat{c}(ng), d)$ and place ng on **nextFrontier** if not there

- frontier** = **nextFrontier**

//Advance all BFS's by 1 level

//pass "visited" information

Parallel k-BFS

- Initial **frontier** = $\{s_1, s_2, \dots, s_k\}$
- $d = 0$
- While **frontier** not empty:
 - **nextFrontier** = $\{\}$
 - $d = d+1$
 - For each vertex v in **frontier**:
 - For each neighbor ng :
 - Do bitwise-OR of v 's words with ng 's words and store in ng
 - If any of ng 's words changed:
 - $e\hat{c}(ng) = \max(e\hat{c}(ng), d)$ and place ng on **nextFrontier** if not there
 - **frontier** = **nextFrontier**

atomic bitwise-OR
using compare-
and-swap

//Advance all BFS's by 1 level

parallel for-loops

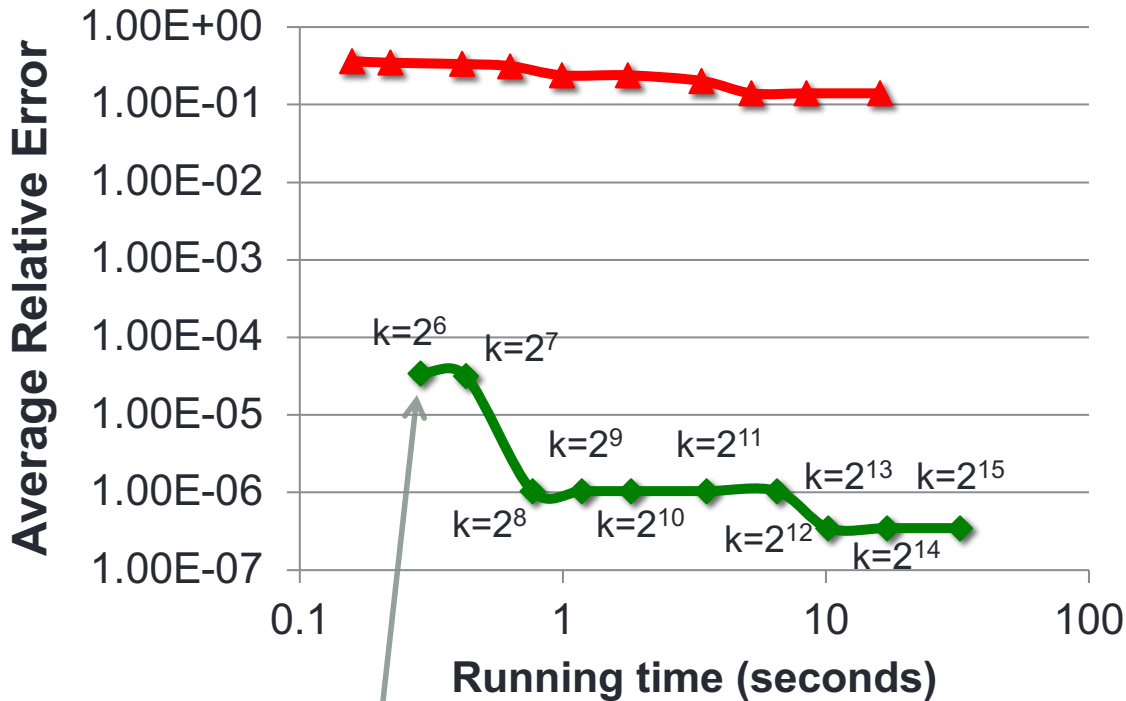
//pass "visited" information

remove
duplicates

- Ligra framework that we use takes care of most details
 - User only specifies function to apply on each edge traversed
 - Performs "direction-optimizing" BFS [Beamer '12] automatically

k-BFS Performance (varying k)

wiki-Talk ($|V| = 2.4M$, $|E| = 4.7M$)



Experiments done on a 40-core Intel Nehalem machine with 256 GB memory

$k = \{2^6, 2^7, 2^8, \dots, 2^{15}\}$

k-BFS (k = 64)

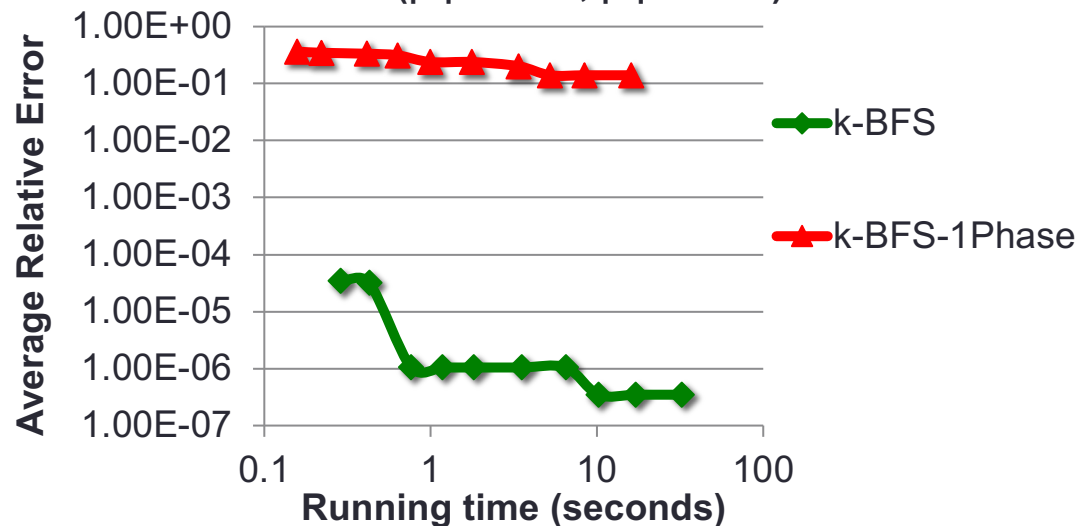
0.288 sec

$< 10^{-4}$

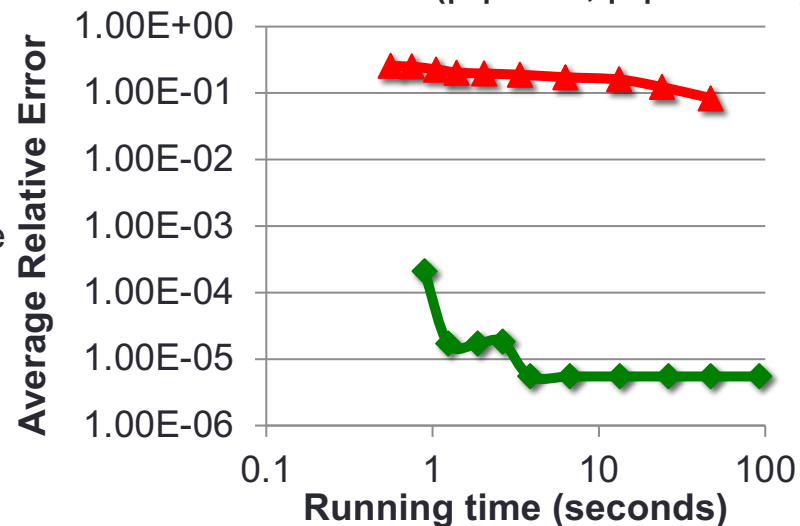
k-BFS Performance

Experiments done on a 40-core Intel Nehalem machine with 256 GB memory

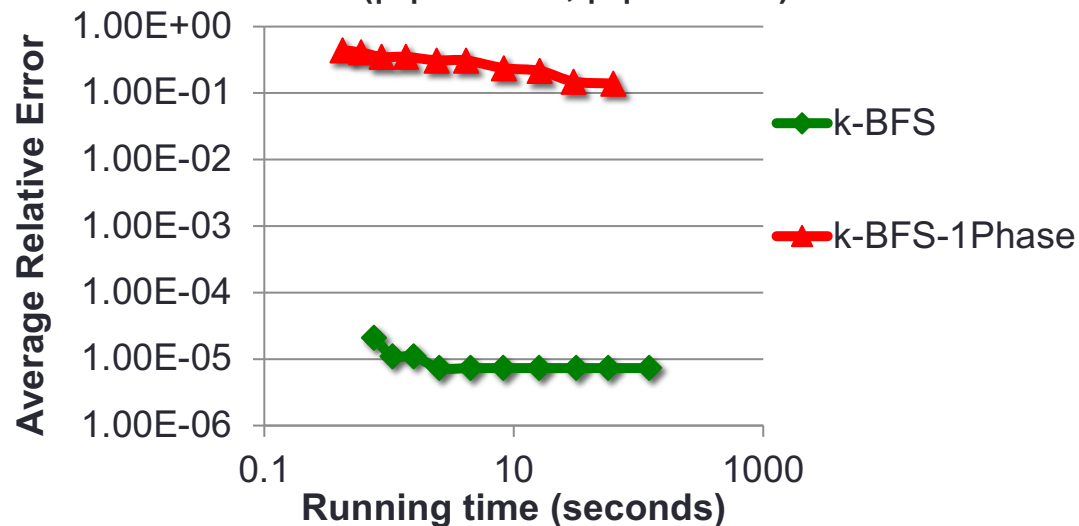
wiki-Talk ($|V| = 2.4M$, $|E| = 4.7M$)



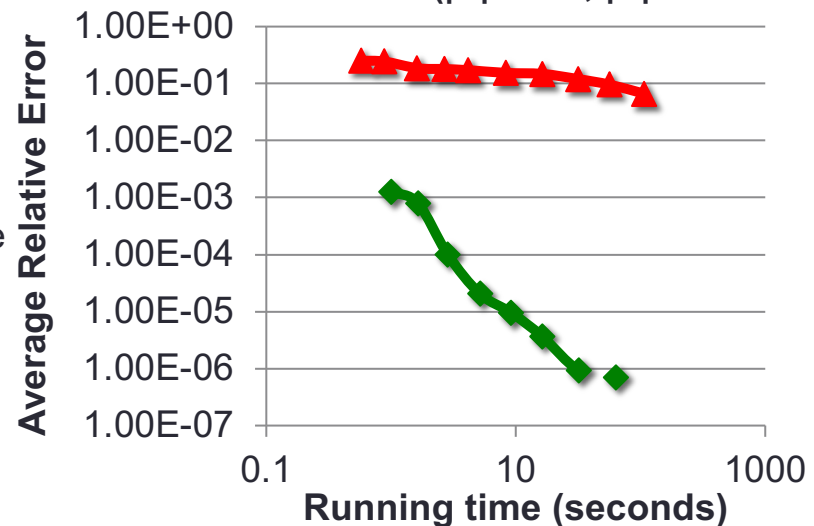
cit-Patents ($|V| = 6M$, $|E| = 16.5M$)



soc-LJ ($|V| = 4.85M$, $|E| = 42.9M$)



com-Orkut ($|V| = 3M$, $|E| = 117.2M$)

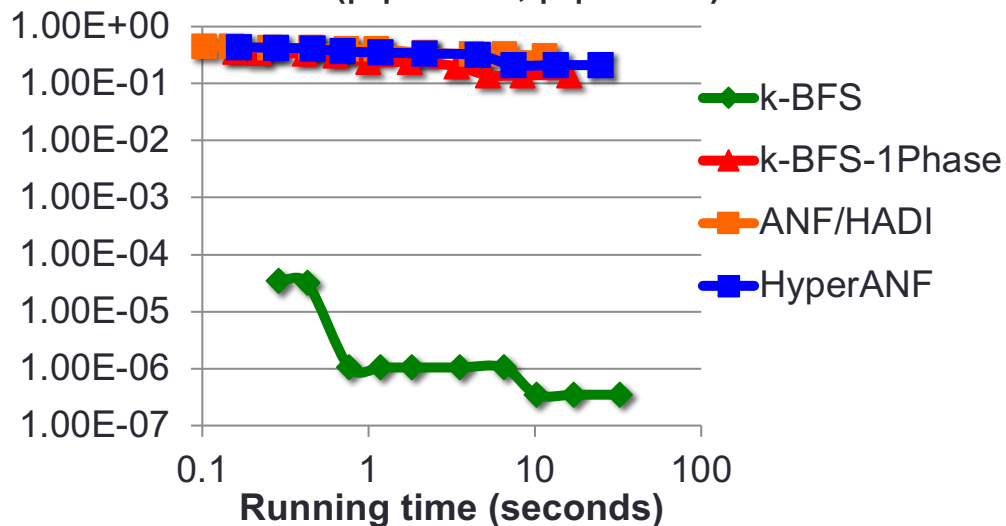


Probabilistic Counter-based Algorithms

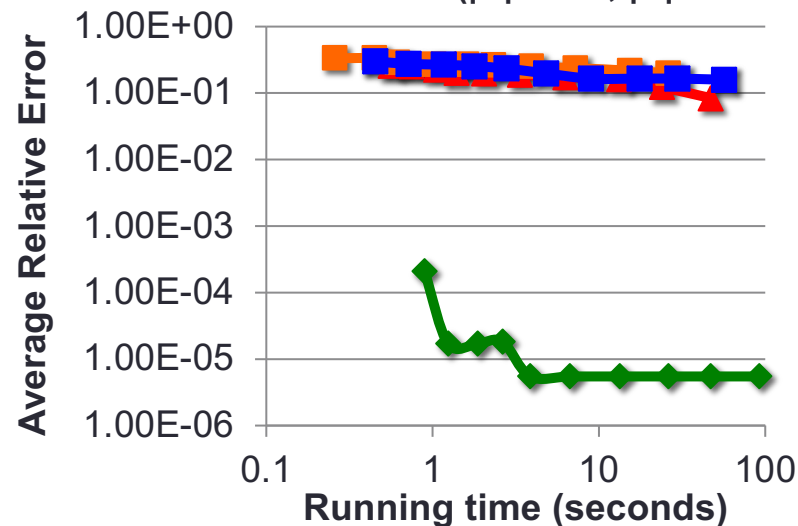
- Assign k probabilistic counters per vertex
- Bitwise-OR counters with all of neighbors' counters until all counters stabilize
- Eccentricity estimate for a vertex is the last round in which any of its counters changed
- ANF/HADI algorithm [Palmer et al. '03, Kang et al. '10] used Flajolet-Martin counters
- HyperANF algorithm [Boldi et al. '11] use more space-efficient HyperLogLog counters [Flajolet et al. '08]
- Shared-memory implementations of variants of ANF/HADI and HyperANF in Ligra

k-BFS outperforms ANF/HADI and HyperANF

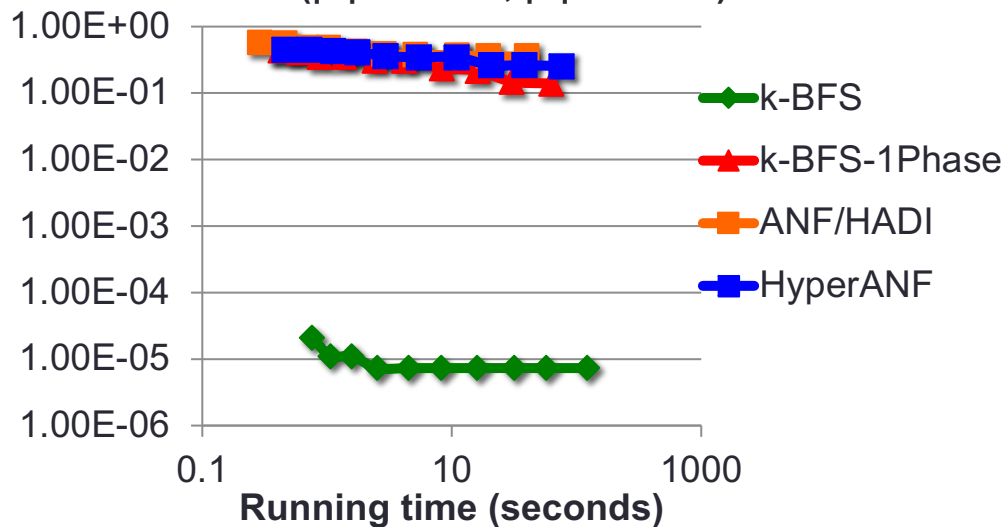
wiki-Talk ($|V| = 2.4M$, $|E| = 4.7M$)



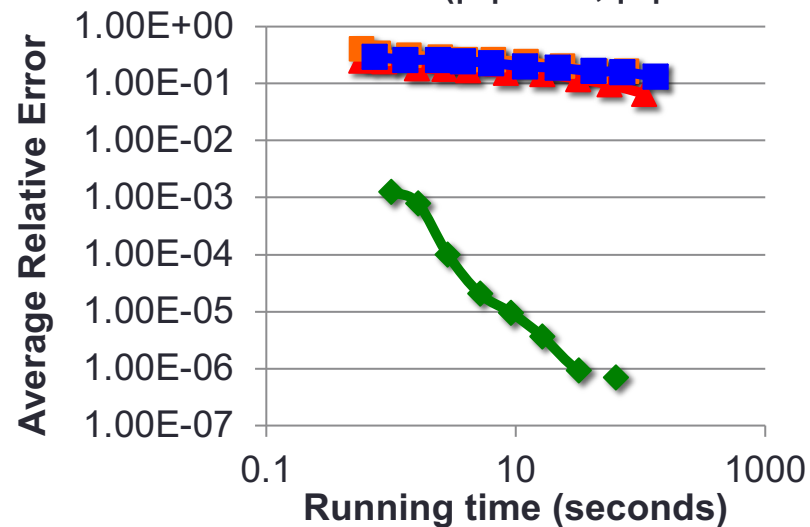
cit-Patents ($|V| = 6M$, $|E| = 16.5M$)



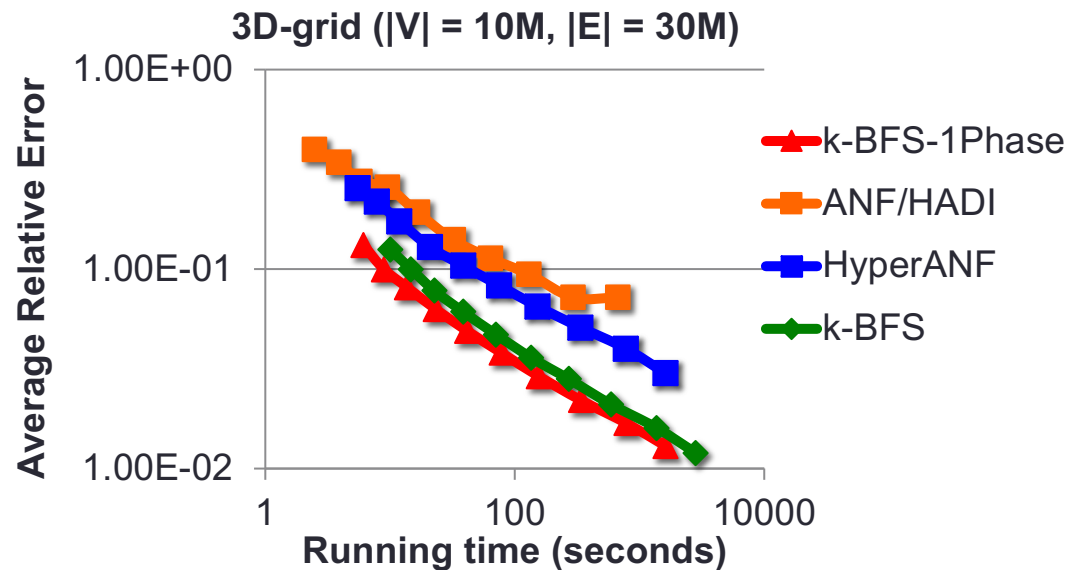
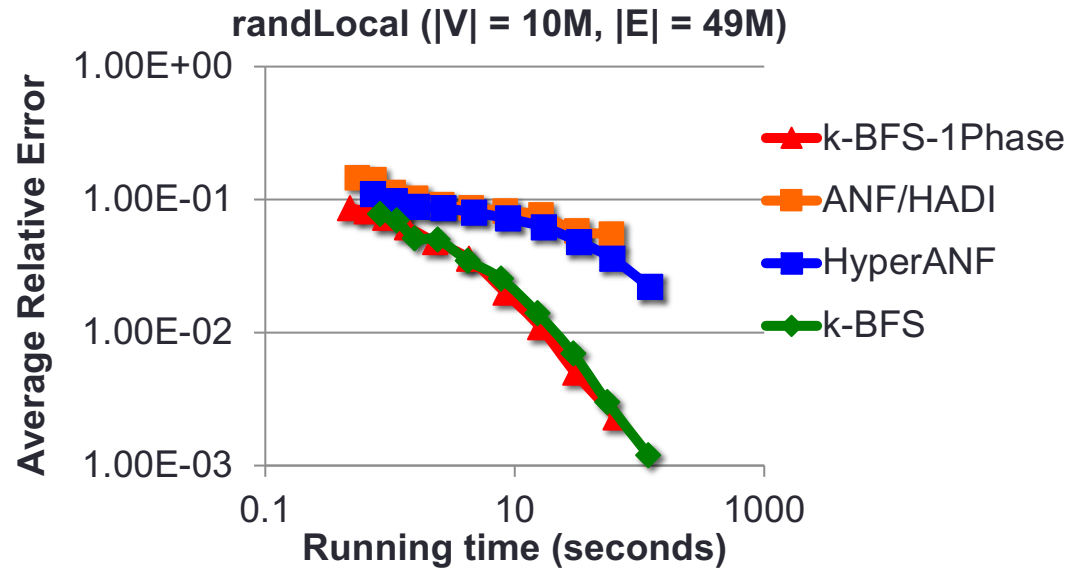
soc-LJ ($|V| = 4.85M$, $|E| = 42.9M$)



com-Orkut ($|V| = 3M$, $|E| = 117.2M$)

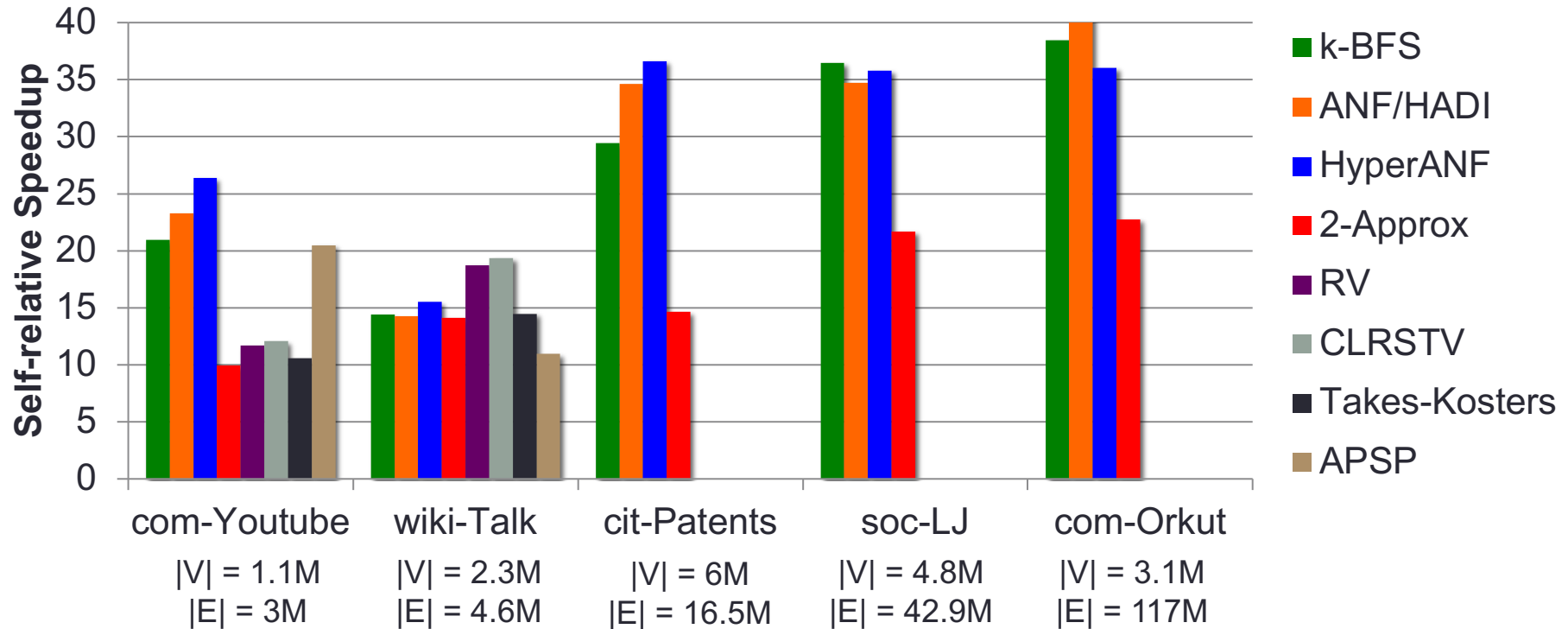


Performance on Synthetic Graphs



Parallel Scalability

Experiments done on a 40-core Intel Nehalem machine with 256 GB memory



- k-BFS achieves 14—38x speedup on 40 cores (higher for larger graphs)

Scaling to Large Graphs

Experiments done on a 40-core Intel Nehalem machine with 256 GB memory

Size

k-BFS (k=64)

Time

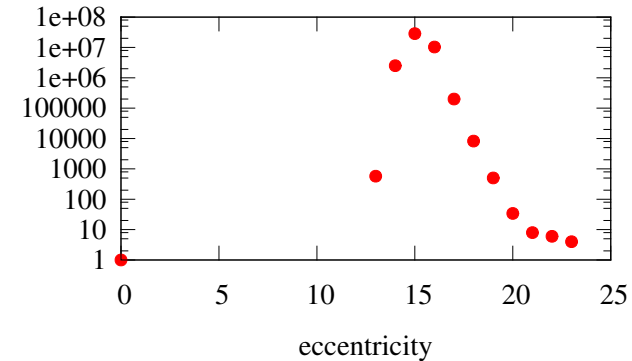
$$|V| = 4.2 \times 10^7$$

$$|E| = 1.2 \times 10^9$$

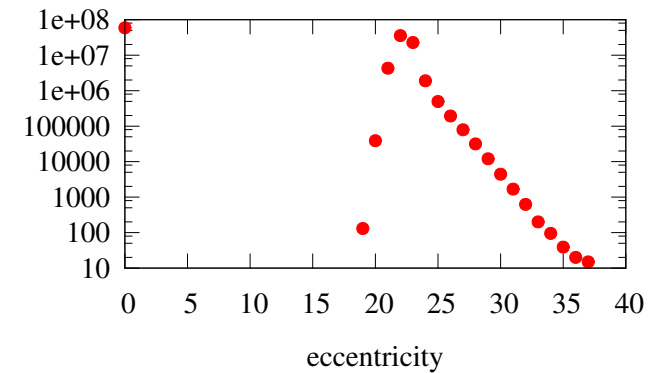
3.5 minutes



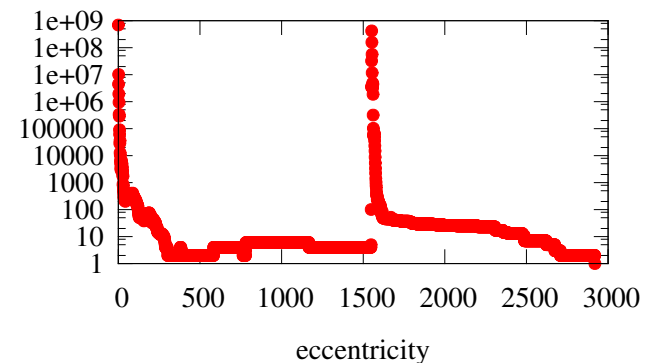
num. vertices with eccentricity



num. vertices with eccentricity



num. vertices with eccentricity



friendster

$$|V| = 1.2 \times 10^8$$

$$|E| = 1.8 \times 10^9$$

1.5 minutes



$$|V| = 1.4 \times 10^9$$

$$|E| = 6.4 \times 10^9$$

11 minutes

Conclusion

- **Comprehensive evaluation** of shared-memory parallel eccentricity estimation algorithms
- k-BFS is orders of magnitude more **accurate** for a fixed running time than other estimation algorithms
- k-BFS is **scalable** to largest publicly-available real-world graphs studied in the literature
- Future work
 - Extensions to directed, weighted graphs
 - Theoretical bounds for (variants of) k-BFS



Code: <http://github.com/jshun/ligra>