

# Three-Way Joins on MapReduce: An Experimental Study

Ben Kimmitt

University of Victoria, BC, Canada  
blk@uvic.ca

Alex Thomo

University of Victoria, BC, Canada  
thomo@cs.uvic.ca

S. Venkatesh

University of Victoria, BC, Canada  
venkat@cs.uvic.ca

**Abstract**—We study three-way joins on MapReduce. Joins are very useful in a multitude of applications from data integration and traversing social networks, to mining graphs and automata-based constructions. However, joins are expensive, even for moderate data sets; we need efficient algorithms to perform distributed computation of joins using clusters of many machines. MapReduce has become an increasingly popular distributed computing system and programming paradigm. We consider a state-of-the-art MapReduce multi-way join algorithm by Afrati and Ullman and show when it is appropriate for use on very large data sets. By providing a detailed experimental study, we demonstrate that this algorithm scales much better than what is suggested by the original paper. However, if the join result needs to be summarized or aggregated, as opposed to being only enumerated, then the aggregation step can be integrated into a cascade of two-way joins, making it more efficient than the other algorithm, and thus becomes the preferred solution.

## I. INTRODUCTION

The importance of joins cannot be overstated. Joins are used explicitly or implicitly when performing a multitude of everyday tasks, such as connecting two or more datasets based on common attributes, comparing tuples with other tuples in the same table using selfjoins, traversing graphs (c.f. [1], [12], [8], [10]), mining graphs and social networks (c.f. [17], [4], [18], [11]) computing graph statistics and cubes (c.f. [16], [19], [20]), and multiplying matrices (c.f. [13]), to name a few. Joins are so useful, they find application even in algorithms requiring the computation of intersection of large automata and transducers (c.f. [15], [14]), or probabilistic reasoning on graph databases (c.f. [9]). In fact, it is hard to imagine a domain where joins are *not* present, albeit sometimes in disguise. This is because of our innate need, in many fields of research, to always be able to connect different pieces of data or information together.

In many of the aforementioned settings, the datasets involved are very big. For example, Facebook, the most popular social network, contains data for over 900 million users and their relationships. To traverse the Facebook graph of friendships and compute statistics would involve a series of challenging multi-way joins. The main reason for the difficulty is that joins are expensive, even for datasets of moderate size. As such, devising distributed algorithms for computing joins on large data sets is of the utmost importance.

MapReduce (c.f. [5], [6], [7]) is a popular distributed computing framework that can work with thousands of machines in a fault-tolerant way. Unsurprisingly, joins have been one of the first candidates to be considered for implementation

in MapReduce. While joining two tables of data is easy to implement from an algorithmic point of view, joining three or more tables brings several challenges to overcome.

In this paper, we focus on three-way joins. Afrati and Ullman in [2], [3] give an elegant algorithm for computing three-way and multi-way joins in MapReduce. However, as we describe later, their main idea for computing the join is based on the assumption of having a limited number of processes (reducers) for processing intermediate results. This is a limiting factor for the scalability of a multi-way join for large clusters with thousands of machines. Also, more often than not, one is more interested in summarizing or aggregating the result of the join in some way. The algorithm proposed in [2], [3] needs to produce the entire join result before an aggregator can summarize it. In contrast, a simple cascade of two way joins may be a better choice, as it allows interleaving the aggregation with the computation of the intermediate result.

We outline the details of how to perform this optimization, and then perform a detailed experimental study on the performance of the algorithm of [2], [3] versus a simple cascade of two-way joins. Our results reveal two surprising facts:

- 1) For real data (such as those coming from edge-lists of social networks), the algorithm of [2], [3] can scale on clusters much bigger than the original papers suggest, and
- 2) When aggregation is required (rather than enumerating the raw join result) a cascade of two-way joins is the preferred choice exhibiting a significant gain in execution cost.

We pay special attention to applying joins for multiplication of large, sparse matrices because of the wide applicability in social network and web analysis. This is also in line with the prime use cases of [2], [3], which also come from the analysis of large social graphs. Such graphs are often represented as sparse matrices by listing the non-empty elements of their incidence matrix (the so-called edge list table).

The outline of the paper is as follows: In Section II, we formally describe joins, and then matrix multiplication and graph computations based on joins. In Section III we describe the MapReduce framework. In Section IV we discuss three-way join algorithms for MapReduce. In Section V we give aggregation algorithms for join results. In Section VI we present our experimental evaluation. Section VII concludes the paper.

## II. JOINS, MATRICES, GRAPHS

Let  $R(A, B, V)$  and  $S(B, C, W)$ , be two tables with attributes (columns)  $A, B, V$ , and  $B, C, W$ , respectively. The join  $R(A, B, V) \bowtie S(B, C, W)$  is table

$$J(A, B, C, V, W) = \{(a, b, c, v, w) : \\ (a, b, x) \in R(A, B, V) \text{ and} \\ (b, c, y) \in S(B, C, W)\}.$$

The definition is extended in the natural way when  $A, B, C, V$ , and  $W$  are sets of attributes, as opposed to single attributes.

The above is also called a “two-way” join because *two* tables are joined. If three tables are joined, we refer to the join as being “three-way”. Join is an associative operation, i.e.  $(R(A, B, V) \bowtie S(B, C, W)) \bowtie T(C, D, X) = R(A, B, V) \bowtie (S(B, C, W) \bowtie T(C, D, X))$ .

With minimal work, the concept of a join can be extended to perform matrix multiplication. A table  $R(A, B, V)$  represents a sparse matrix, with each tuple  $(a, b, v)$  in  $R$  representing the presence of  $v$  in the matrix, at row  $a$  and column  $b$ .

When two matrix tables  $R(A, B, V)$  and  $S(B, C, W)$  are “joined” the effect is to perform the first step of matrix multiplication.  $R$  and  $S$  are joined on their common attribute,  $B$ , and the  $V$  and  $W$  values are multiplied. Call the result  $J(A, C, P)$ . Observe that in the result we only keep  $A$  and  $C$ , and multiply the values of  $V$  and  $W$  rather than just output them as in the join definition. In effect, the row vectors of the matrix represented by  $R$  have been multiplied by the column vectors of the matrix represented by  $S$ .

The next step is to perform summation of the  $p = v \cdot w$  values of  $J(A, C, P)$  tuples that agree on  $A$  and  $C$ , and produce one tuple  $(a, c, s_{a,c})$  for each existing  $a, c$  combination in  $J(A, C, P)$ . That is, we group by  $a, c$  and aggregate using sum. This is the equivalent of summing the intermediate results of matrix multiplication to yield the finished matrix.

Now, graphs can be represented as incidence matrices, which are often quite sparse for real graphs. We can use join-based matrix multiplication to multiply graph matrices with themselves  $n$  times and thus obtain the number (or the weight) of paths of length  $n + 1$  between the starting and ending nodes listed in the final output. This is important in the friend-of-friend analysis of social networks.

Also, by considering the diagonal of the result (those  $(a, c, s_{a,c})$  tuples with  $a = c$ ) for binary incidence matrices, we can obtain the number of triangles in the graph. Namely, the number of triangles is the sum of all  $s_{a,c}$ , with  $a = c$ , divided by three.

As the matrix multiplication is a simple extension of the join followed by a group by and aggregation, we will first focus on MapReduce algorithms for join, then modify them to handle matrix multiplication.

## III. MAPREDUCE

MapReduce is used to describe both a distributed computing system and an algorithmic paradigm, used to process large sets of data. The most popular incarnation of MapReduce as a

distributed system is Hadoop; an open-source Java implementation based on the Hadoop Distributed File System (HDFS).

From an algorithmic point of view, MapReduce simplifies distributed computing. All a programmer needs to do is implement a Map and a Reduce function, without having to deal with low level details of machine communication, data transfer, scheduling, and fault-tolerance. Depending on the number of machines and configuration of the system, the Map and Reduce functions can run in many machines at once.

There are always two phases in a MapReduce job, the Map phase and the Reduce phase. The latter starts once the former completes. The processes running the map function are called *mappers*, and the processes running the reduce function are called *reducers*.

Both the map and reduce functions take *key-value pairs* (KVPs) as input. They also emit key-value pairs; the exact construction of any KVP depends on the individual map or reduce function used. The records emitted by the mappers are sorted and shuffled by the system before being sent to the reducers. The guarantee of the system is that all the emitted KVPs with the same key are sent to the same reducer. A reducer can receive KVPs with many different keys, but if it receives one KVPs with a specific key, it is certain to receive all other KVPs with that same key.

A two-way join  $R(A, B, V) \bowtie S(B, C, W)$  can be implemented in MapReduce as follows. Initially, each mapper is assigned a chunk of data, which can contain tuples from  $R(A, B, V)$ ,  $S(B, C, W)$ , or both. The specification for the Map and Reduce functions is as follows.

### Map function.

For each pair  $(tid, (a, b, v))$ , where  $(a, b, v) \in R(A, B, V)$ , emit  $(b, (a, v, R))$ .

For each pair  $(tid, (b, c, w))$ , where  $(b, c, w) \in S(B, C, W)$ , emit  $(b, (c, w, S))$ .

### Reduce function.

Join all  $(a, b, v)$ ’s from  $(b, (a, v, R))$

with all  $(b, c, w)$ ’s from  $(b, (c, w, S))$ , matching on  $b$ .

Emit  $((a, b, c, v, w), \dagger)$ , where  $\dagger$ ’s value is unimportant. Some attributes may optionally be omitted from the output.

Because of the system guarantee that all KVPs with the same key are sent to the same reducer, the reducer receiving the  $(b, (\_, \_, \_))$  pairs has all the information it needs to compute the section of the join related to value  $b$ .

What really matters for the efficiency of a MapReduce algorithms is the total amount of I/O performed by all the processes. Emitted data is considered I/O; this is because data emission and transmission is realized using HDFS. The total amount of I/O is called the *communication cost*. However, as typical in databases when we compare algorithms, we do not count the cost of writing the final output. The reasoning behind this is that the output is either small enough to be consumed by human users (in which case the cost can be safely ignored), or if not, it will be pipelined into another process (possibly another MapReduce round) which will summarize or aggregate it in some way. In this case, the size of the output would be counted in the communication cost of the next process, and so is not included in this cost estimate.



### Reduce function.

For all tuples  $((a, c), p)$ , sum all the tuples'  $p$  values, returning  $((a, c), s_{a,c})$ .

The same aggregator is also used after the second join.

The aggregator employed after computing the three-way join using 1,3J is as follows.

### Map function.

For all tuples  $((a, b, c, d, v, w, x), \dagger)$ , emit  $((a, d), p)$ , where  $p = v \cdot w \cdot x$ .

### Reduce function.

For all tuples  $((a, d), p)$ , sum all the tuples'  $p$  values, returning  $((a, d), s_{a,d})$ .

We call 1,3J followed by this aggregator 1,3JA.

2,3JA computes  $R(A, B, V) \bowtie S(B, C, W)$  [of size  $r'$ ], then aggregates the result, yielding  $Agg(R(A, B, V) \bowtie S(B, C, W))$  [of size  $r''$ ]. This is then joined with  $T(C, D, X)$ . The communication cost is similar to that of 2,3J, but with the addition of  $2r''$ , for a total of  $6r + 2r' + 2r''$  tuples. The exact  $r''$  depends on how well the aggregator can reduce  $R(A, B, V) \bowtie S(B, C, W)$ , but it is always equal to or less than (usually much less than)  $r'$ .

Unlike 2,3J and 2,3JA, the communication cost of 1,3J for 1,3JA rises with the number of reducers. Recall, 1,3J's communication cost is  $4r + 2r\sqrt{k}$  tuples, where  $k$  is the number of reducers. The computation cost of 1,3JA is  $4r + 2r\sqrt{k} + 2r'''$ , where  $r'''$  is the size of the raw three-way join.

## VI. EVALUATION

Our experiments ran on a 33-node (4 cores per node, maximum 132 MapReduce instances running at any one time) cluster, using Apache Hadoop 1.2.1. For data, we acquired seven datasets from the Stanford Large Network Dataset Collection (<http://snap.stanford.edu/data/>). Each dataset represented a directed graph. Of the datasets, five (Slashdot, Twitter, Wikitalk, Pokec, and LiveJournal) were social in nature, representing user relationships (Twitter, Pokec, Livejournal) or interactions (voting on other users for Slashdot, comments on talk pages for Wikitalk). The remaining two datasets each represented different data: the Amazon set was derived from Amazon's "customers who bought item A also bought item B" database, while the Google Web dataset is a small chunk of internet structure (with edges representing links between pages), which was released by Google as part of a programming competition.

In all experiments, each dataset was joined to itself twice, creating a three-way selfjoin. The three copies of the dataset will still be referred to as  $R$ ,  $S$ , and  $T$ .

### A. Results

For the first set of tests, we ran our implementation of 2,3J and 1,3J on each dataset. We measured the communication cost as defined above. The results are shown in Figure 2.

For every dataset, 1,3J had a lower communication cost than 2,3J for a large number of reducers. After a certain reducer threshold, the 2,3J cost became lower than the 1,3J cost, but the

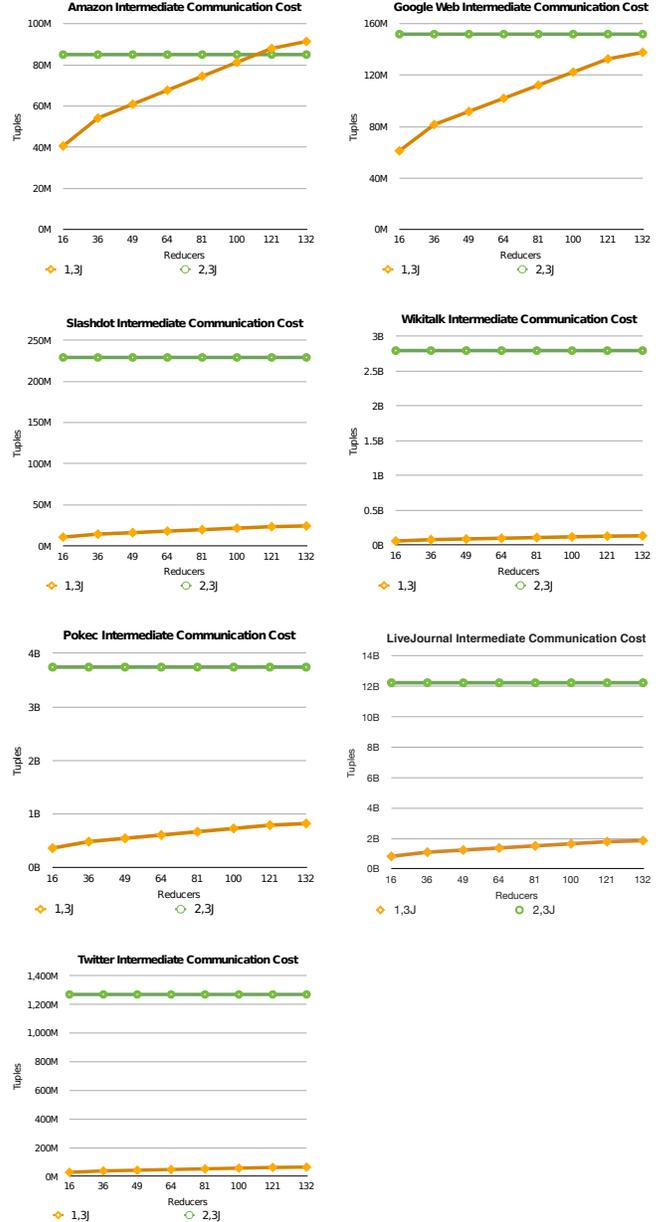


Fig. 2. Sizes of intermediate 1,3J and 2,3J communication cost, compared. Top row: Amazon (left), Google Web (right). 2nd row: Slashdot (left), Wikitalk (right). 3rd row: Pokec (left), LiveJournal (right). Bottom left: Twitter.

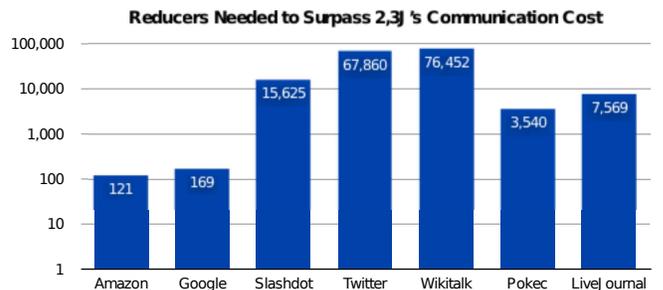


Fig. 3. For the 1,3J's intermediate communication cost to surpass the 2,3J's on a specific dataset, it would have to use the listed number of reducers.

number of reducers the 1,3J would need to use to cost more than the 2,3J is typically very large, as shown in Figure 3. This is a surprising fact that shows that 1,3J can scale much more than what [2], [3] suggest (only 960 reducers for a large hypothetical social network).

As shown in Figure 2, the Twitter dataset’s communication cost was far lower when running the 1,3J algorithm. A cluster running the two algorithms would need to have 67,860 reducers (a 260x261 reducer array, or about 8,400 8-core machines) before the 1,3J’s communication cost would be larger than the 2,3J’s. Similarly, the LiveJournal dataset would have to be run on a cluster with 7,569 reducers (an 87x87 reducer array, about 950 8-core machines) before the cost of running the 1,3J algorithm upon it would grow greater than the cost of running the 2,3J on the same.

For the second set of tests, we compared 2,3JA and 1,3JA on each dataset, and the results are shown in Figure 6. On the graphs of the larger datasets, the 1,3JA line still has a slope (the 2,3JA line is flat), but it cannot be easily observed due to the graph’s scale; Figure 6 [bottom right] illustrates the actual slope of one such line. The 2,3JAs cost does not change as the cluster size increases, while the 1,3JAs cost only gets larger.

For each dataset, the 2,3JA algorithm’s communication cost was far less than the 1,3JA algorithm’s, a fact due entirely to the reduced output size. If the intermediate aggregator combined some number of KVPs ( $n$ ) into one, and that KVP produced  $m$  tuples in the aggregated final output, then  $n * m$  tuples ( $n$  identical sets of  $m$  tuples) would be output in the unaggregated result.

The benefits of this (as shown in Figure 4) were transferred over to the 2,3JA’s second join round, producing a reduction in output size. The exact reduction in size is shown in Figure 5. As an example, the output from the primary aggregation round on the Pokec dataset is 76.4% of the size of first two-way join, a little larger than the average. This benefit carries over to the algorithm’s output: the 2,3JA’s Pokec output is 69.1% the size of the 1,3J’s output on the same dataset. In comparison, the LiveJournal dataset’s intermediate aggregated result is 56.9% the size of the first two-way join; the final 2,3JA output is 42.2% the size of the 1,3J output.

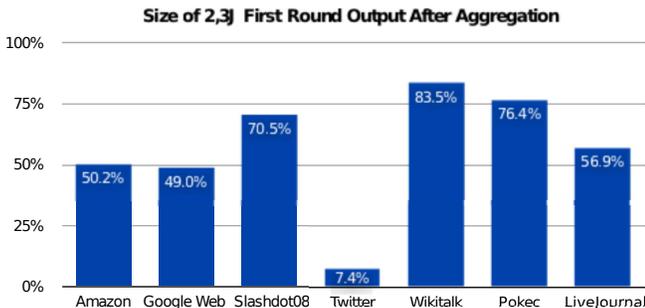


Fig. 4. This chart shows the size of the intermediate aggregation of the first round of the 2,3JA, as a percentage of the size of the first two-way join.

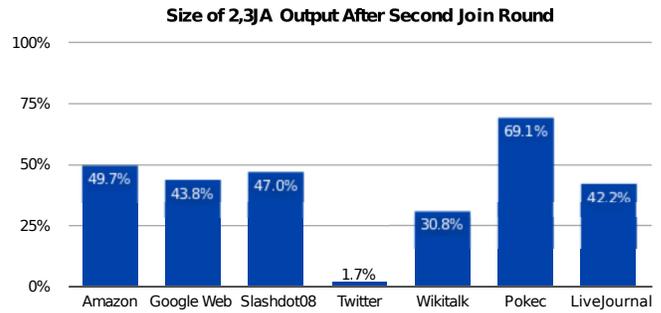


Fig. 5. This chart shows the size of the output of the 2,3JA, measured as a percentage of the output size of the 1,3J, where aggregation is not used.

## VII. CONCLUSIONS

We focused on three-way joins for MapReduce. Such joins are especially useful for friends-of-friends analysis and triangle computation in social networks. We considered the algorithm of [2], [3] (1,3J) versus a simple cascade of two-way joins (2,3J). The communication cost of 1,3J is dependent on the number of reducers; the more the number of reducers used in the computation, the bigger the communication cost becomes. On the other hand, the communication cost of 2,3J does not change when the number of reducers changes. We showed that 1,3J can scale much better than what was suggested in [2], [3], often by one or two orders of magnitude. However, when the result of the join needs to be aggregated in some way as in the case of matrix multiplication, then a cascade of two-way joins (2,3JA) is preferable to 1,3JA as the aggregation can be pushed to the intermediate results of 2,3JA significantly reducing the communication cost incurred.

## REFERENCES

- [1] Gephi. <http://gephi.org/>.
- [2] F. N. Afrati and J. D. Ullman. Optimizing joins in a map-reduce environment. In *EDBT*, pages 99–110, 2010.
- [3] F. N. Afrati and J. D. Ullman. Optimizing multiway joins in a map-reduce environment. *IEEE Trans. Knowl. Data Eng.*, 23(9):1282–1298, 2011.
- [4] D. Chakrabarti and C. Faloutsos. Graph mining: Laws, generators, and algorithms. *ACM Comput. Surv.*, 38(1), 2006.
- [5] J. Dean and S. Ghemawat. Mapreduce: Simplified data processing on large clusters. In *OSDI*, pages 137–150, 2004.
- [6] J. Dean and S. Ghemawat. Mapreduce: simplified data processing on large clusters. *Commun. ACM*, 51(1):107–113, 2008.
- [7] J. Dean and S. Ghemawat. Mapreduce: a flexible data processing tool. *Commun. ACM*, 53(1):72–77, 2010.
- [8] G. Grahne and A. Thomo. Regular path queries under approximate semantics. *Ann. Math. Artif. Intell.*, 46(1-2):165–190, 2006.
- [9] N. Hassanlou, M. Shoaran, and A. Thomo. Probabilistic graph summarization. In *WAIM*, pages 545–556, 2013.
- [10] N. Koochakzadeh, A. Sarraf, K. Kianmehr, J. G. Rokne, and R. Al-hajj. Netdriller: A powerful social network analysis tool. In *ICDM Workshops*, pages 1235–1238, 2011.
- [11] N. Korovaiko and A. Thomo. Trust prediction from user-item ratings. *Social Netw. Analys. Mining*, 3(3):749–759, 2013.
- [12] Z. Miao, D. C. Stefanescu, and A. Thomo. Grid-aware evaluation of regular path queries on spatial networks. In *AINA*, pages 158–165, 2007.
- [13] J. Myung and S. goo Lee. Matrix chain multiplication via multi-way join algorithms in mapreduce. In *ICUIMC*, page 53, 2012.

- [14] M. Shoaran and A. Thomo. Evolving schemas for streaming xml. *Theor. Comput. Sci.*, 412(35):4545–4557, 2011.
- [15] A. Thomo, S. Venkatesh, and Y. Y. Ye. Visibly pushdown transducers for approximate validation of streaming xml. In *FoIKS*, pages 219–238, 2008.
- [16] Y. Tian, R. A. Hankins, and J. M. Patel. Efficient aggregation for graph summarization. In *SIGMOD Conference*, pages 567–580, 2008.
- [17] W. Wang, C. Wang, Y. Zhu, B. Shi, J. Pei, X. Yan, and J. Han. Graphminer: a structural pattern-mining system for large disk-based graph databases and its applications. In *SIGMOD Conference*, pages 879–881, 2005.
- [18] X. Xu, N. Yuruk, Z. Feng, and T. A. J. Schweiger. Scan: a structural clustering algorithm for networks. In *KDD*, pages 824–833, 2007.
- [19] N. Zhang, Y. Tian, and J. M. Patel. Discovery-driven graph summarization. In *ICDE*, pages 880–891, 2010.
- [20] P. Zhao, X. Li, D. Xin, and J. Han. Graph cube: on warehousing and olap multidimensional networks. In *SIGMOD Conference*, pages 853–864, 2011.

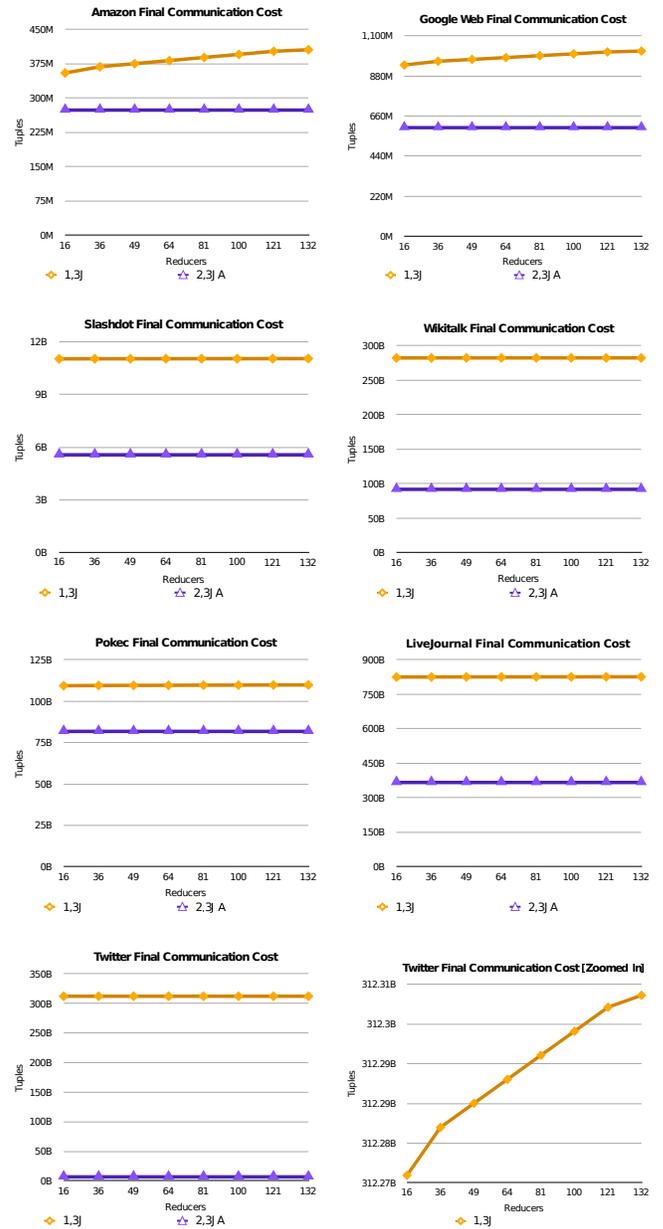


Fig. 6. Sizes of communication cost of 1,3JA and 2,3JA. Top row: Amazon (left), Google Web (right). 2nd row: Slashdot (left), Wikitalk (right). 3rd row: Pokec (left), LiveJournal (right). Bottom: Twitter (left), Twitter graph scaled to show slope of 1,3J line (right).