

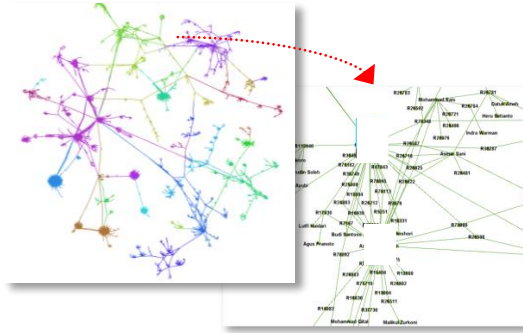
From NoSQL Accumulo to NewSQL Graphulo: Design and Utility of Graph Algorithms inside a BigTable Database

Dylan Hutchison Jeremy Kepner
Vijay Gadepally Bill Howe

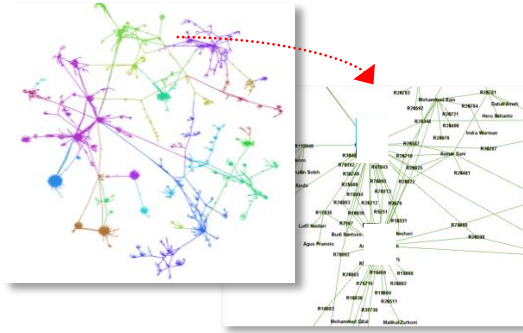


UNIVERSITY *of* WASHINGTON
Massachusetts Institute of Technology



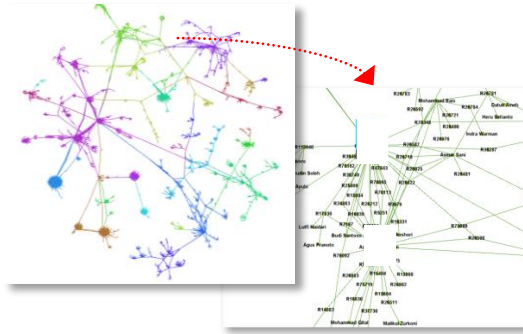


Computation ➔ Databases



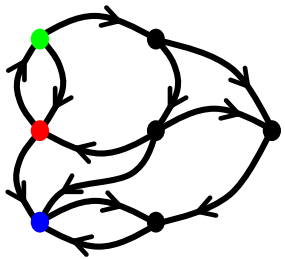
Computation → Databases

Graph Algorithms → BigTable

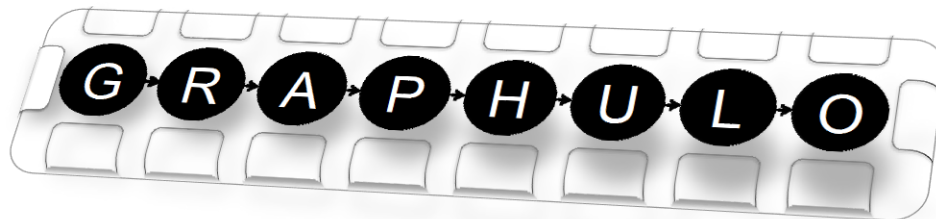


Computation → Databases

Graph Algorithms → BigTable



GraphBLAS →



Why push Compute into Accumulo?

- > Data Locality**
 - Save communication
- > Reuse infrastructure**
 - One less system to adopt and maintain
- > Database features for free**
 - Indexed access
 - Distributed execution

if the computation aligns with the DB's access path

1. *How* to do matrix computation in Accumulo?

2. Applications

- Jaccard coefficients
- k-Truss subgraph

3. *When* is this a good idea?

- Spoiler: Compare Memory and I/O

Adjacency Matrix Schema

| Key | | | | Value | |
|-----|--------|-----------|------------|-------|-----------|
| Row | Column | | | | Timestamp |
| | Family | Qualifier | Visibility | | |

> (Row, Column Qualifier, Value)
 = (v_1 , v_2 , weight)
 [Transpose: (v_2 , v_1 , weight)]

$$\begin{array}{c}
 1 \quad 10 \quad \dots \\
 1 \begin{bmatrix} 141 & 12 & \dots \\ 18 & & \vdots \\ \vdots & & \ddots \end{bmatrix} \\
 10 \\
 \vdots
 \end{array}$$

| Adjacency Table | | | |
|-----------------|---------|----|-----|
| row | :colq | -> | val |
| 1 | :1 [] | -> | 141 |
| 1 | :10 [] | -> | 12 |
| 1 | :101 [] | -> | 9 |
| 1 | :105 [] | -> | 3 |
| 1 | :11 [] | -> | 9 |
| 1 | :110 [] | -> | 3 |
| 1 | :111 [] | -> | 3 |
| 1 | :113 [] | -> | 12 |
| 10 | :1 [] | -> | 18 |
| 10 | :109 [] | -> | 2 |

Other schemas supported:

- Graph as Incidence Matrix
- Graph as Single Table with degrees

Adjacency Matrix Schema

| Key | | | | Value | |
|-----|--------|-----------|------------|-------|-----------|
| Row | Column | | | | Timestamp |
| | Family | Qualifier | Visibility | | |

- > (Row, Column Qualifier, Value)
= (v_1 , v_2 , weight)
[Transpose: (v_2 , v_1 , weight)]
- > Degree table: store
vertex degrees separately

| Degree Table | |
|--------------|---------|
| 1 :in [] | -> 1084 |
| 1 :out [] | -> 1027 |
| 10 :in [] | -> 118 |
| 10 :out [] | -> 94 |
| 100 :in [] | -> 8 |
| 100 :out [] | -> 10 |

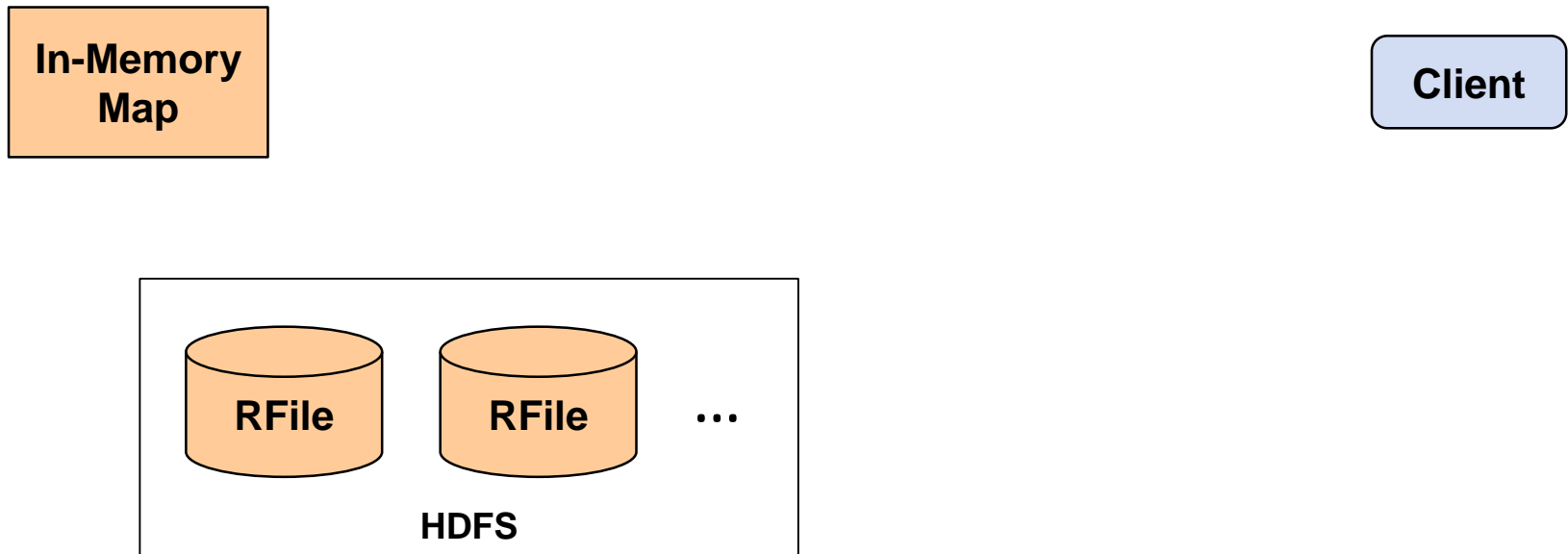
$$\begin{matrix} & 1 & 10 & \dots \\ \begin{matrix} 1 \\ 10 \\ \vdots \end{matrix} & \begin{bmatrix} 141 & 12 & \dots \\ 18 & & \vdots \\ \vdots & \dots & \ddots \end{bmatrix}
 \end{matrix}$$

Other schemas supported:

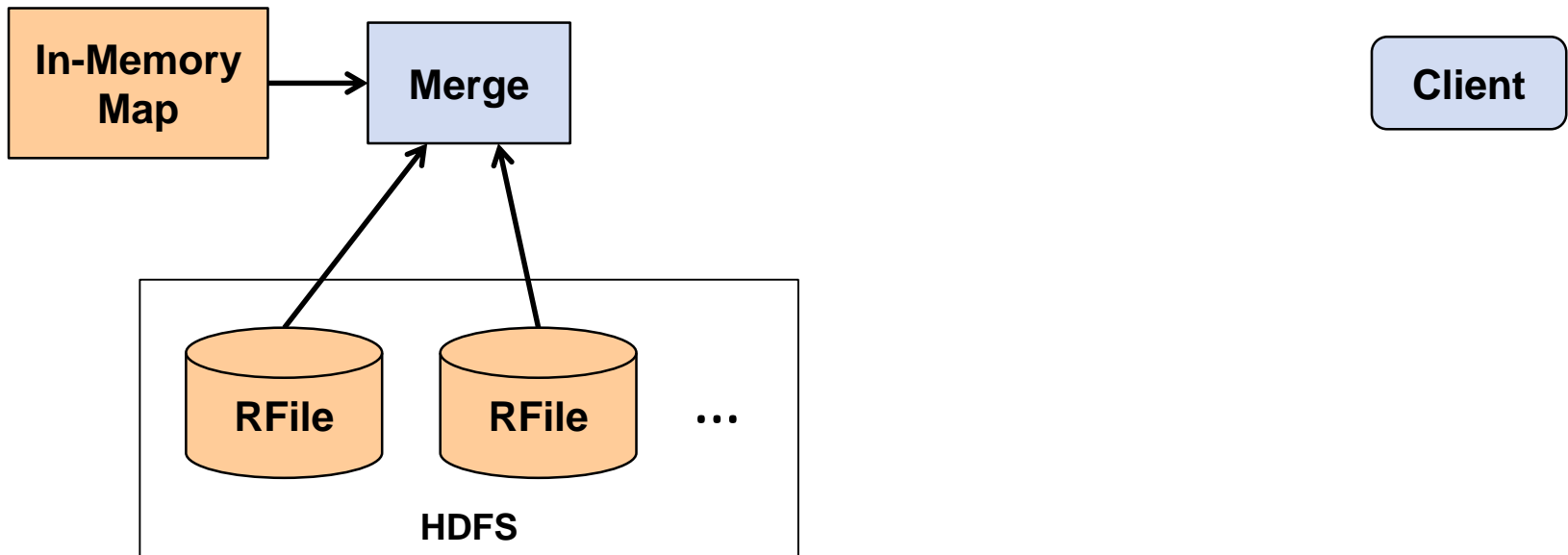
- Graph as Incidence Matrix
- Graph as Single Table with degrees

| Adjacency Table | |
|-----------------|--------|
| row :colq | ->val |
| 1 :1 [] | -> 141 |
| 1 :10 [] | -> 12 |
| 1 :101 [] | -> 9 |
| 1 :105 [] | -> 3 |
| 1 :11 [] | -> 9 |
| 1 :110 [] | -> 3 |
| 1 :111 [] | -> 3 |
| 1 :113 [] | -> 12 |
| 10 :1 [] | -> 18 |
| 10 :109 [] | -> 2 |

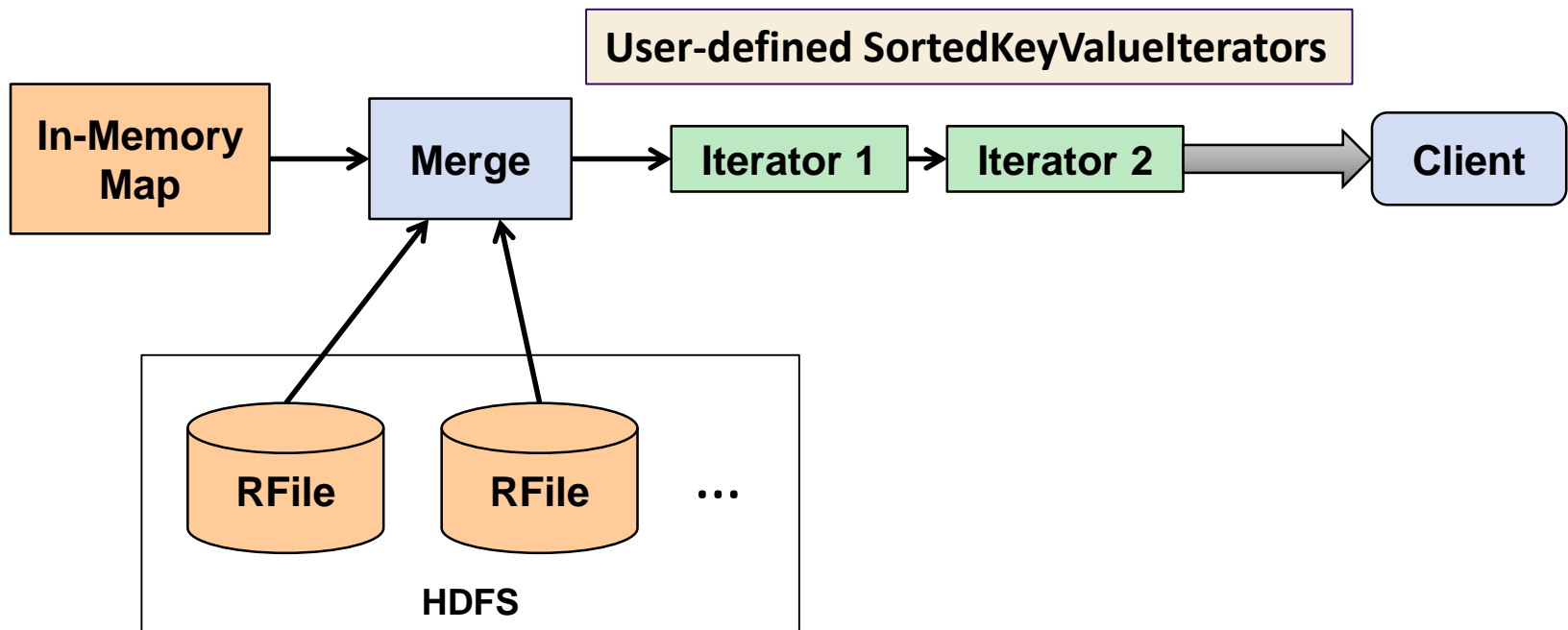
Accumulo Scan Iterator Pipeline



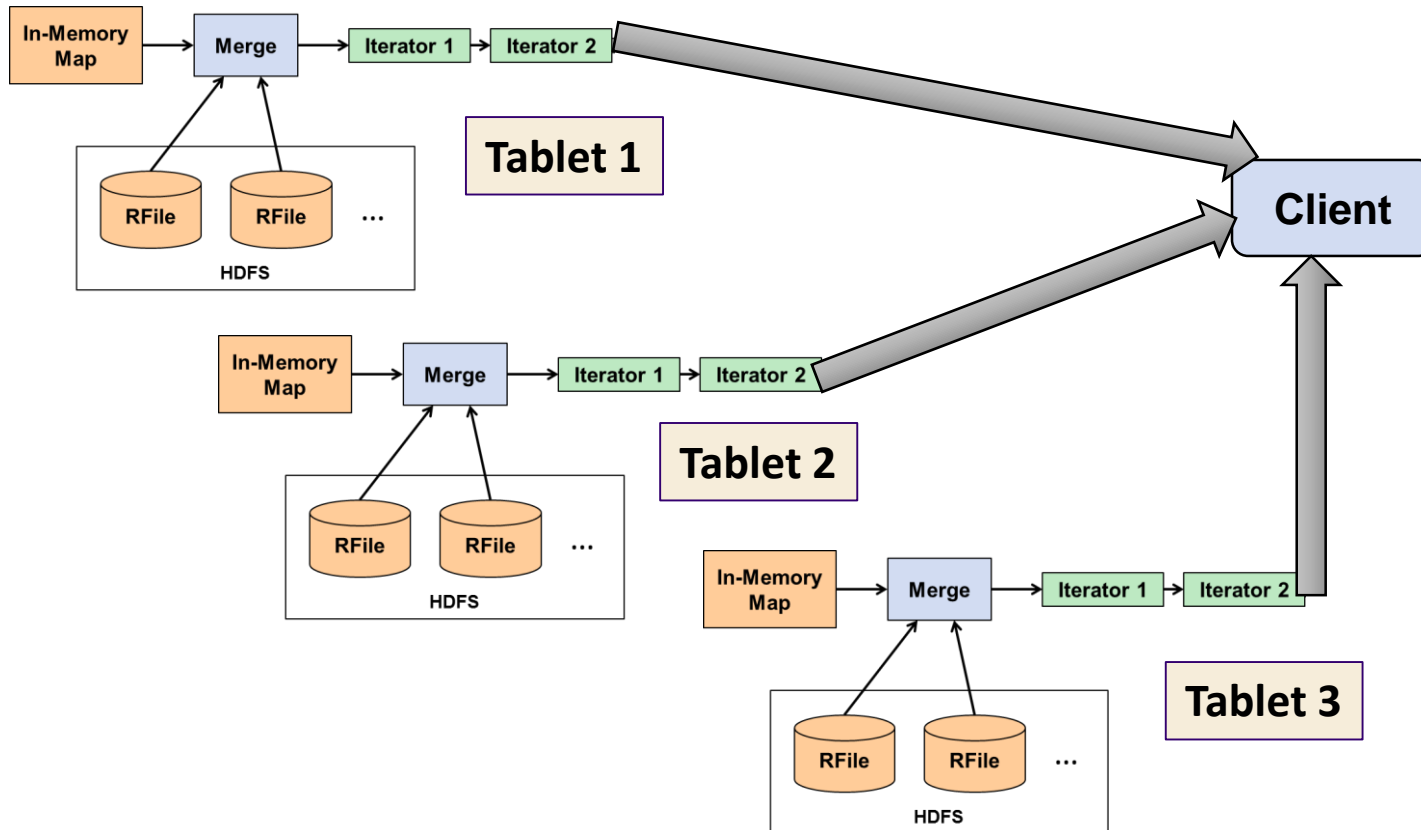
Accumulo Scan Iterator Pipeline



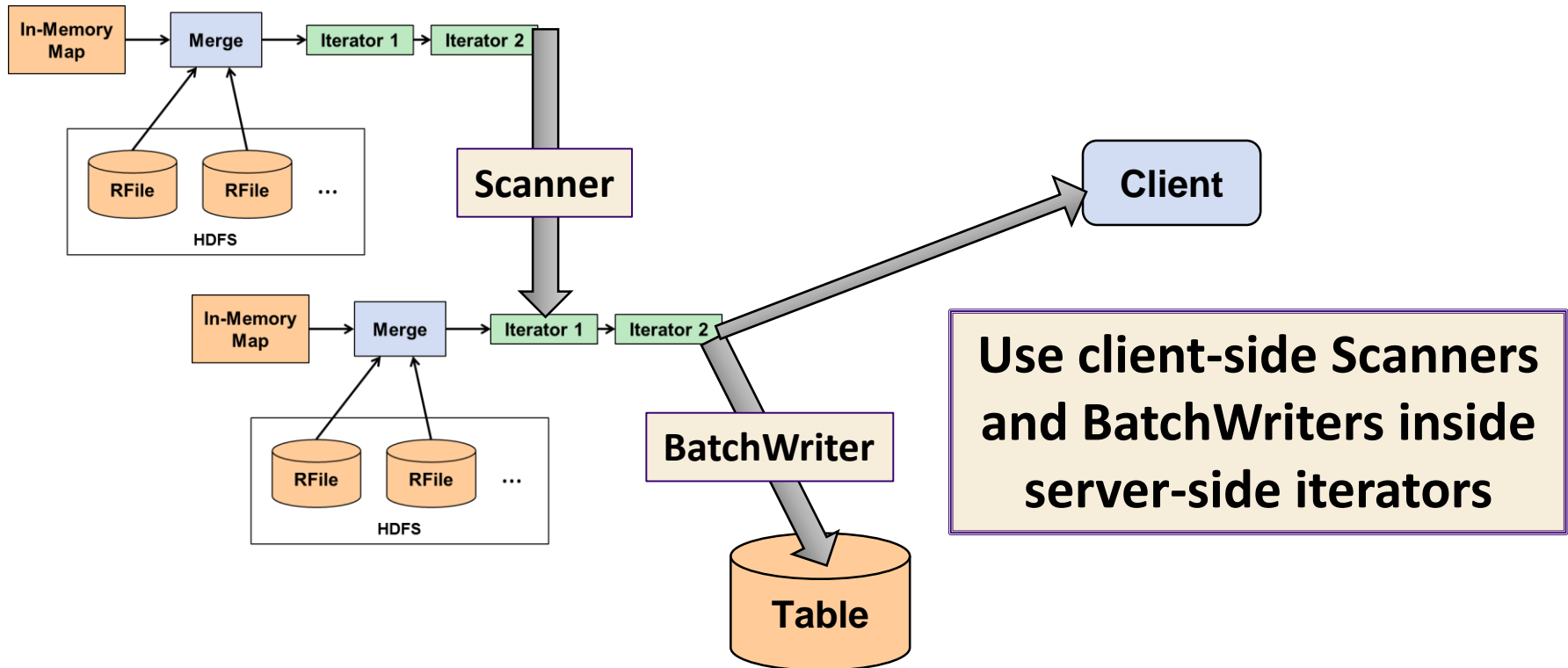
Accumulo Scan Iterator Pipeline



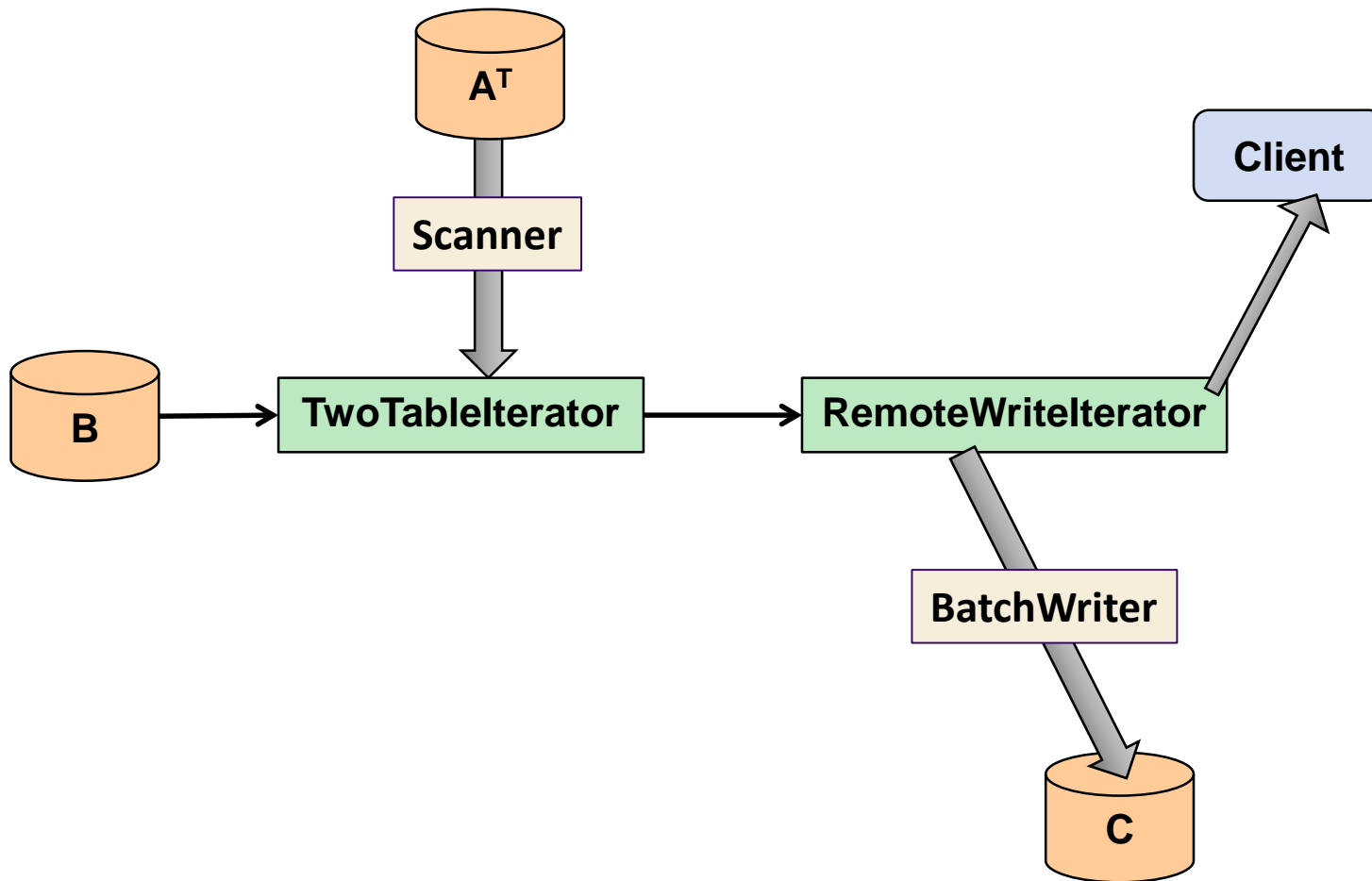
Accumulo BatchScan Iterator Pipeline



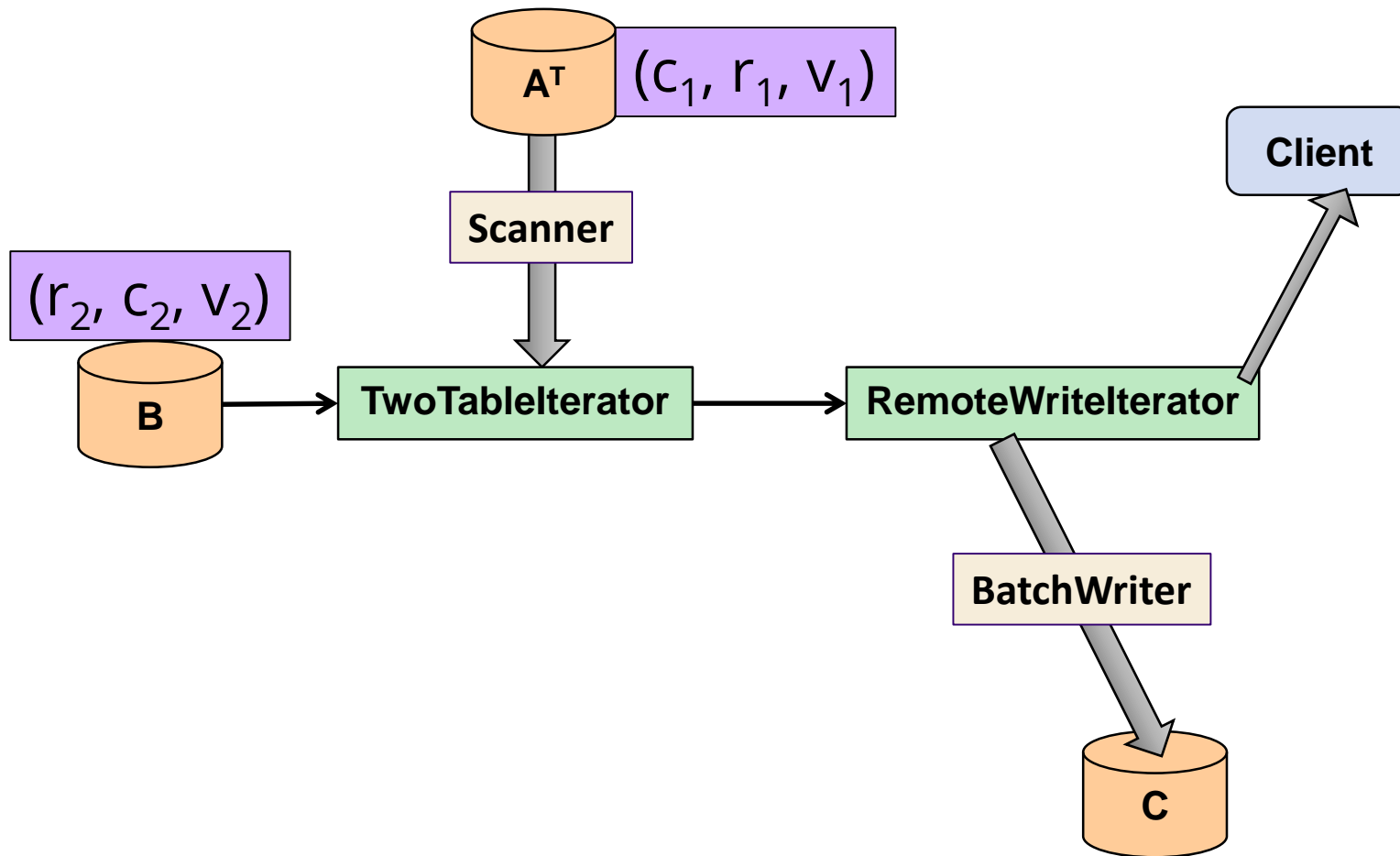
Graphulo addition to Iterator Pipeline



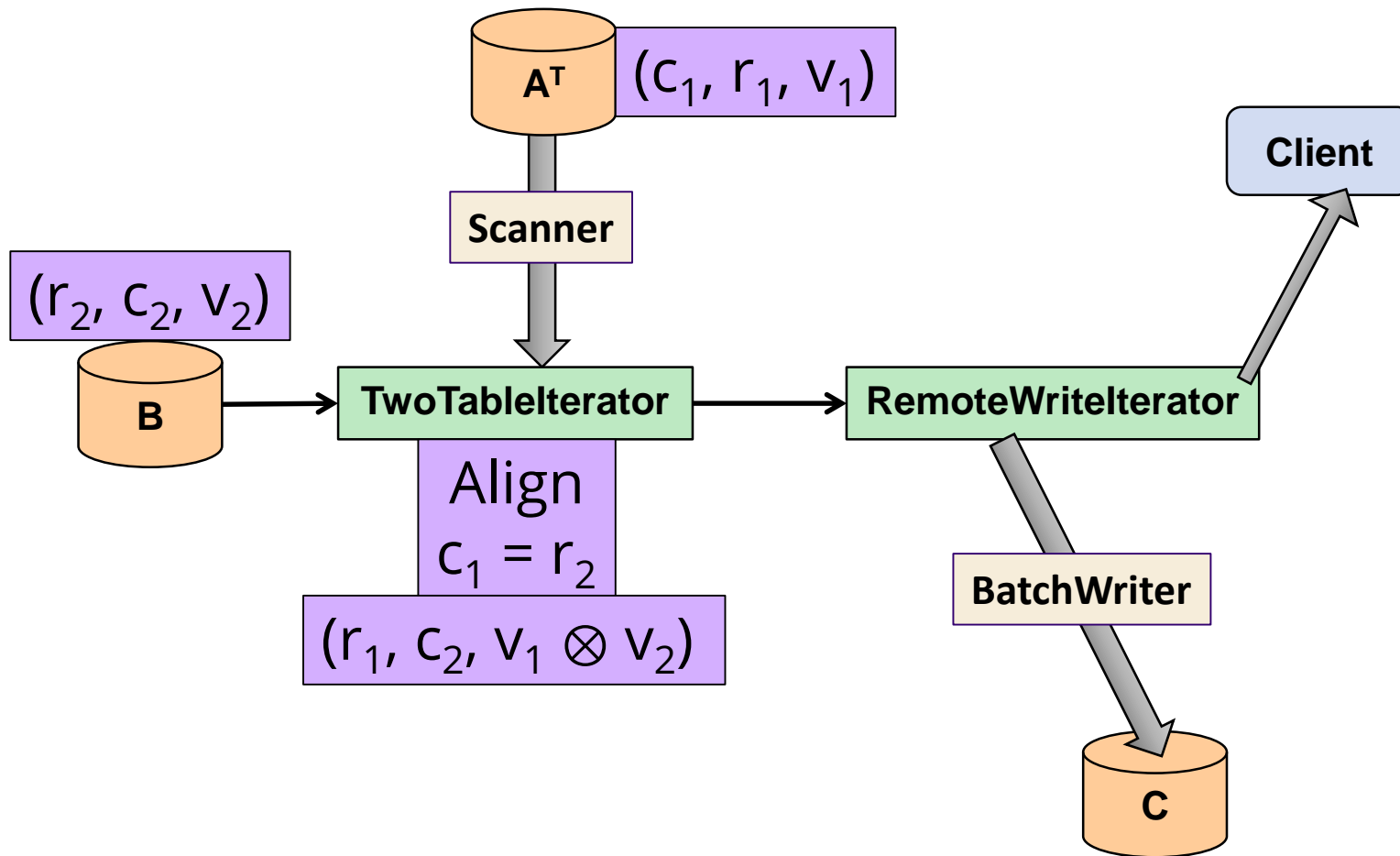
Graphulo MxM: $A^T (\oplus, \otimes) B$



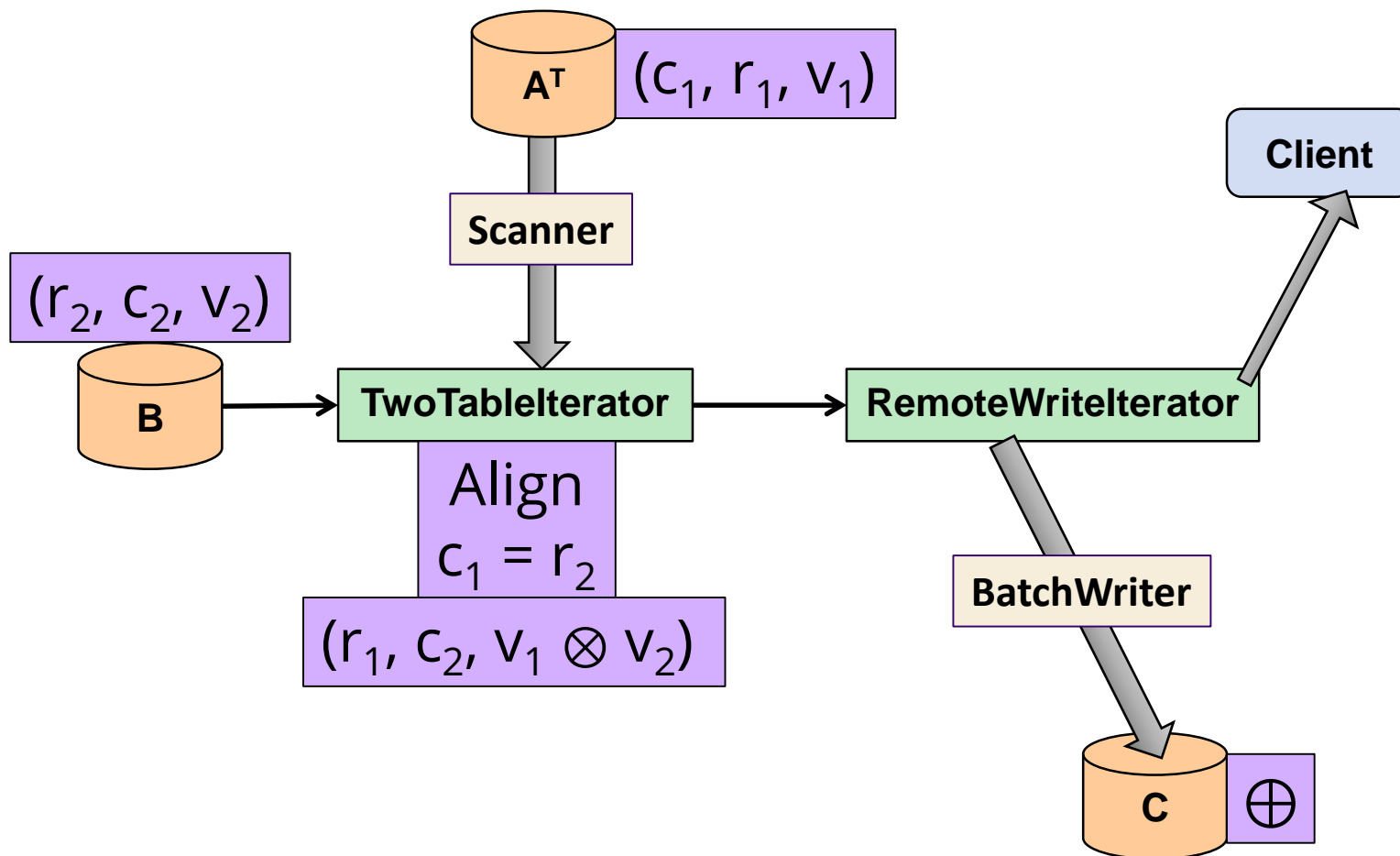
Graphulo MxM: $A^T (\oplus.\otimes) B$



Graphulo MxM: $A^T (\oplus.\otimes) B$



Graphulo MxM: $A^T (\oplus.\otimes) B$



Graphulo Client Functions

```
long TableMult(String Atable, String Btable, String Ctable)
```

```
long SpEwiseX(String Atable, String Btable, String Ctable)
```

```
long SpEwiseSum(String Atable, String Btable, String Ctable)
```

```
...
```

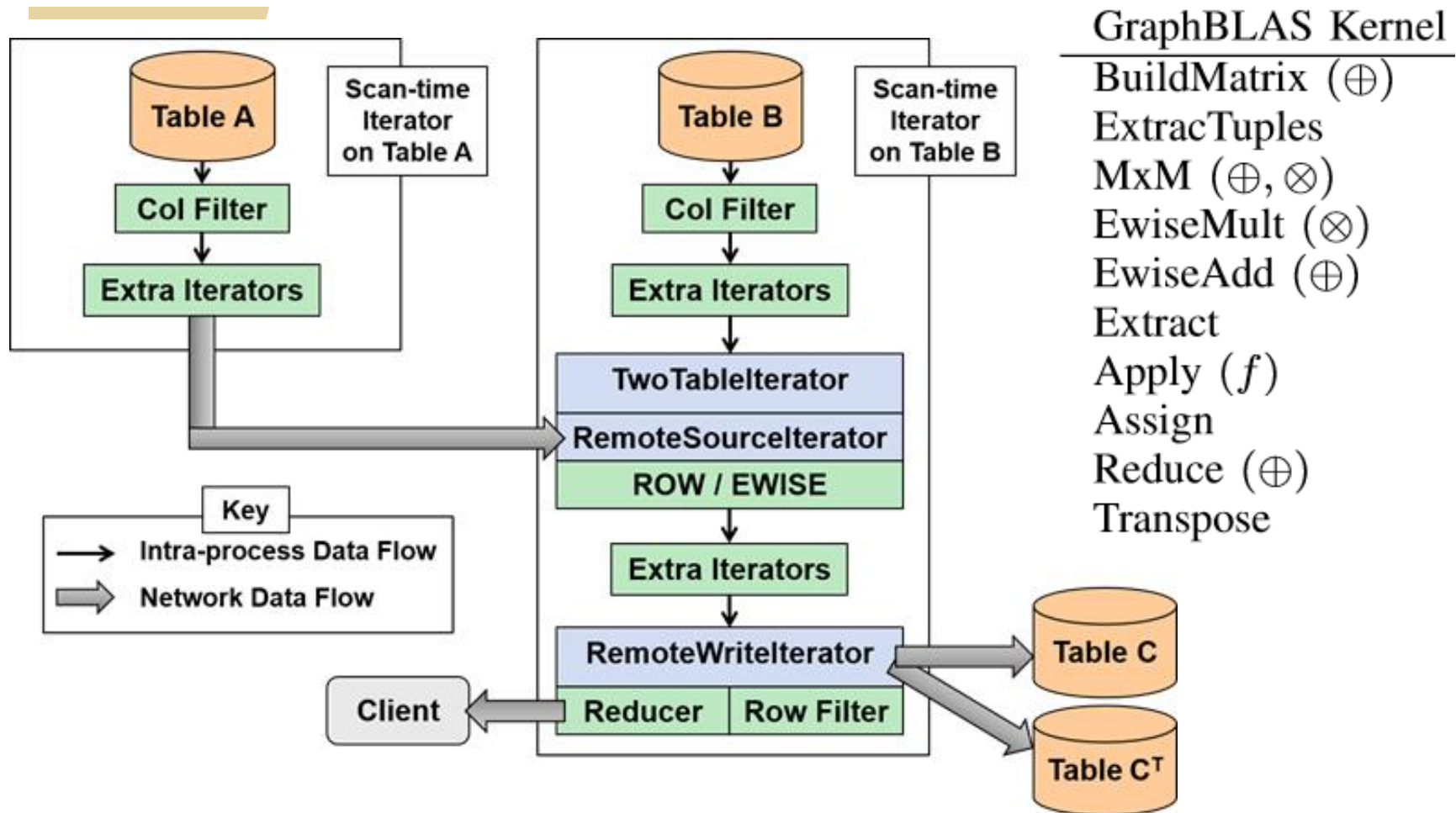
**Simple API abstracts
the iterator pipeline**

Graphulo Client Functions

```
long TwoTable(  
    String ATtable, String Btable, String Ctable, String CTtable,  
    int BScanIteratorPriority, TwoTableIterator.DOTMODE dotmode,  
    Map<String, String> optsTT, IteratorSetting plusOp,  
    String rowFilter, String colFilterAT, String colFilterB,  
    boolean emitNoMatchA, boolean emitNoMatchB,  
    List<IteratorSetting> iteratorsBeforeA,  
    List<IteratorSetting> iteratorsBeforeB,  
    List<IteratorSetting> iteratorsAfterTwoTable,  
    Reducer reducer, Map<String, String> reducerOpts,  
    int numEntriesCheckpoint, Authorizations ATauthorizations,  
    Authorizations Bauthorizations, int batchWriterThreads  
)
```

**Full control
when you need it**

Graphulo's TwoTable Pipeline

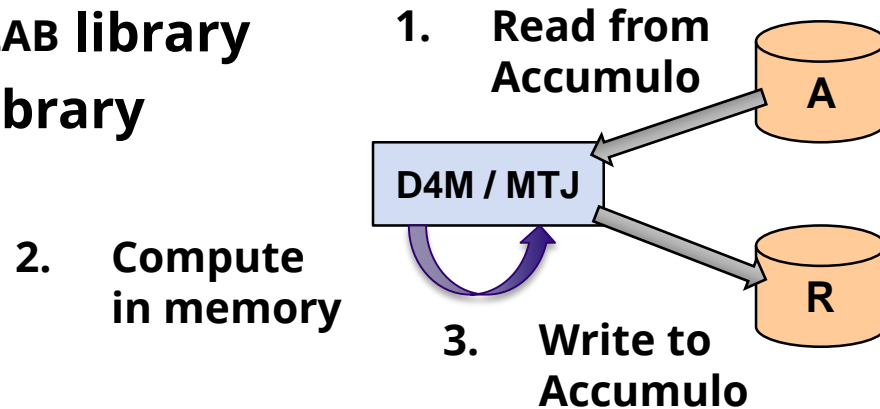


Applications and Performance

Applications and Performance

- > Algorithms: Jaccard and k-Truss
- > Comparison to main-memory alternatives

- Sparse: D4M MATLAB library
- Dense: MTJ Java library



- > Node: 16GB RAM, 8 i7 cores, Accumulo 1.8.0
- > Data: Unpermuted power law graph generator
 - 2^{SCALE} nodes, 16 edges/node
- > Accumulo Threads: 1-2 tablets

Alg 1: Jaccard Coefficients

- > Neighborhood overlap of two vertices
- > Vertex similarity measure

$$J_{ij} = \frac{|N(v_i) \cap N(v_j)|}{|N(v_i) \cup N(v_j)|}$$

```
long Jaccard(String Aorig, String ADeg, String Rfinal)
```

Graphulo: Linear algebra graph kernels for NoSQL databases, Gadepally et al, IEEE IPDPSW 2015

Alg 1: Jaccard Coefficients

Input: Unweighted, undirected adjacency matrix \mathbf{A}

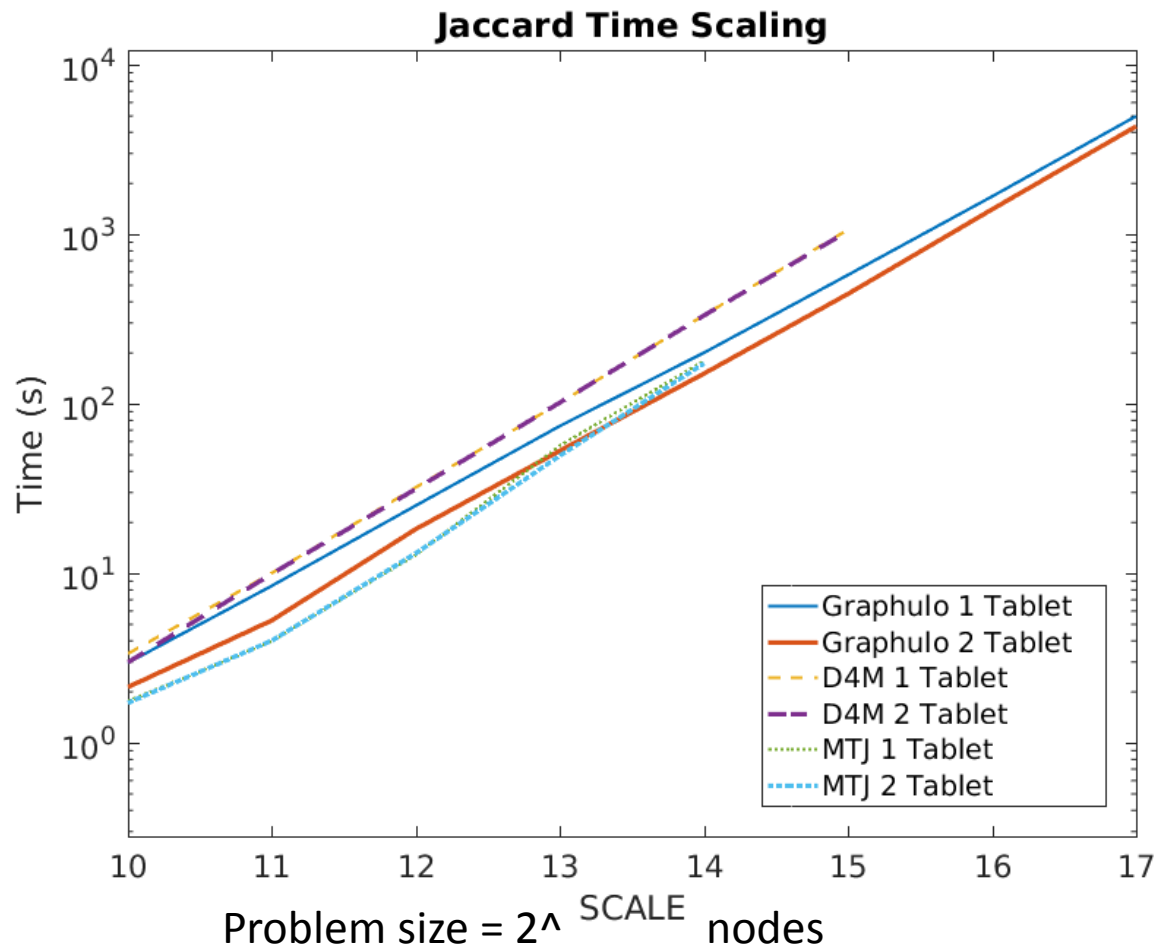
Output: Upper triangle of Jaccard coefficients \mathbf{J}

```
1  $\mathbf{d} = \text{sum}(\mathbf{A})$  // pre-computed in degree table
2  $\mathbf{U} = \text{triu}(\mathbf{A}, 1)$  // strict upper triangle filter
3  $\mathbf{J} = \text{triu}(\mathbf{U}\mathbf{U} + \mathbf{U}\mathbf{U}^\top + \mathbf{U}^\top\mathbf{U}, 1)$  // fused  $M \times M$ 
4 foreach nonzero entry  $\mathbf{J}_{ij} \in \mathbf{J}$  do
5   |  $\mathbf{J}_{ij} = \mathbf{J}_{ij} / (\mathbf{d}_i + \mathbf{d}_j - \mathbf{J}_{ij})$  // stateful Apply on  $\mathbf{J}$ 
6 end
```

Algorithm 1: Jaccard

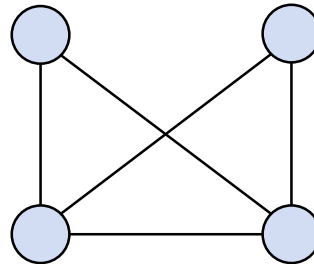
Fusion:
A single Graphulo
TwoTable pass!

Jaccard Performance



Alg 2: k-Truss Subgraph

- > A graph is a k-Truss if each edge is part of at least $k-2$ triangles
 - May be the empty graph
- > Indicates a graph's "core community"



Example 3-truss

```
long kTrussAdj(String Aorig, String Rfinal, int k)
```

Alg 2: k-Truss Subgraph

Input: Unweighted, undirected adjacency matrix \mathbf{A}_0 ,
integer k

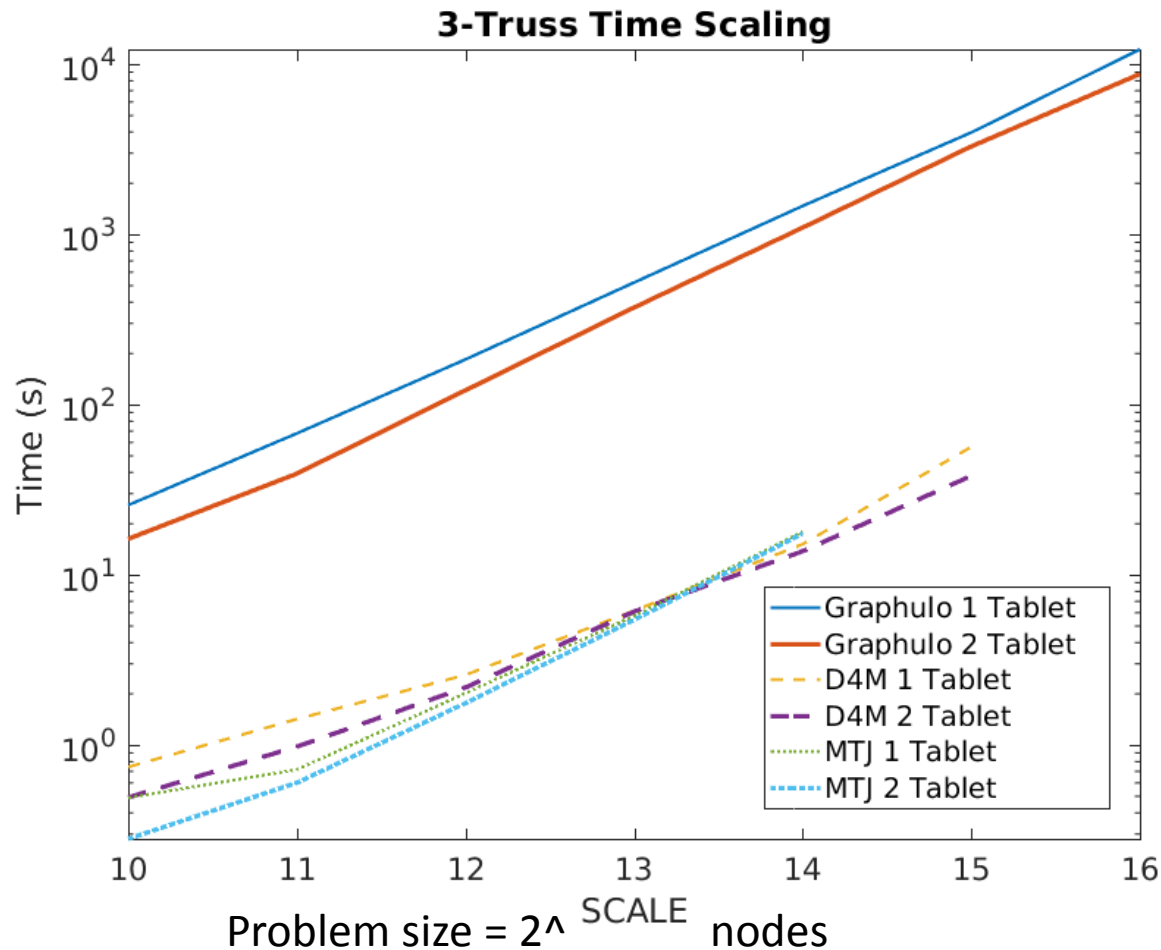
Output: Adjacency matrix of k -truss subgraph \mathbf{A}

```
1  $z' = \infty, \mathbf{A} = \mathbf{A}_0$  // table clone
2 repeat
3    $z = z'$ 
4    $\mathbf{B} = \mathbf{A}$  // table clone
5    $\mathbf{B} = \mathbf{B} + 2\mathbf{A}\mathbf{A}$  //  $M \times M$  with  $a \otimes b = 2$  if  $a, b \neq 0$ 
6    $\mathbf{B}(\mathbf{B} \% 2 == 0) = 0$  // filter on  $\mathbf{B}$ 
7    $\mathbf{B}((\mathbf{B} - 1)/2 < k - 2) = 0$  // filter on  $\mathbf{B}$ 
8    $\mathbf{A} = |\mathbf{B}|_0$  // Apply on  $\mathbf{B}$ ; switch  $\mathbf{A} \leftrightarrow \mathbf{B}$ 
9    $z' = \text{nnz}(\mathbf{A})$  // Reduce, gathering nnz at client
10 until  $z == z'$  // client controls iteration
```

Algorithm 2: kTruss

Iterations of Graphulo TwoTables

3-Truss Performance



Performance Diff

| SCALE | nnz(A) | nnz(Jaccard(A)) | Partial Products | Graphulo Overhead | Graphulo 1 Tablet | Graphulo 2 Tablets | D4M 1 Tablet | D4M 2 Tablets | MTJ 1 Tablet | MTJ 2 Tablets |
|-------|--------------------|--------------------|--------------------|-------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| 10 | 2.10×10^4 | 2.15×10^5 | 1.01×10^6 | 4.7x | 2.97 | 2.14 | 3.36 | 2.99 | 1.76 | 1.72 |
| 11 | 4.52×10^4 | 7.07×10^5 | 3.10×10^6 | 4.4x | 8.46 | 5.29 | 1.01×10^1 | 9.96 | 3.99 | 4.01 |
| 12 | 9.67×10^4 | 2.18×10^6 | 9.29×10^6 | 4.3x | 2.52×10^1 | 1.83×10^1 | 3.22×10^1 | 3.16×10^1 | 1.29×10^1 | 1.32×10^1 |
| 13 | 2.04×10^5 | 6.75×10^6 | 2.71×10^7 | 4.0x | 7.42×10^1 | 5.30×10^1 | 1.02×10^2 | 1.02×10^2 | 5.68×10^1 | 4.96×10^1 |
| 14 | 4.26×10^5 | 2.02×10^7 | 7.77×10^7 | 3.8x | 2.01×10^2 | 1.51×10^2 | 3.34×10^2 | 3.33×10^2 | 1.79×10^2 | 1.73×10^2 |
| 15 | 8.83×10^5 | 6.07×10^7 | 2.22×10^8 | 3.7x | 5.77×10^2 | 4.46×10^2 | 1.07×10^3 | 1.05×10^3 | | |
| 16 | 1.82×10^6 | 1.77×10^8 | 6.20×10^8 | 3.5x | 1.68×10^3 | 1.41×10^3 | | | | |
| 17 | 3.73×10^6 | 5.16×10^8 | 1.72×10^9 | 3.3x | 4.99×10^3 | 4.34×10^3 | | | | |

TABLE II: Jaccard experiment statistics. Graphulo is competitive and better scales due to low overhead.

| SCALE | nnz(A) | nnz(3Truss(A)) | Partial Products | Graphulo Overhead | Graphulo 1 Tablet | Graphulo 2 Tablets | D4M 1 Tablet | D4M 2 Tablets | MTJ 1 Tablet | MTJ 2 Tablets |
|-------|--------------------|--------------------|--------------------|-------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| 10 | 2.10×10^4 | 2.03×10^4 | 5.94×10^6 | 293.3x | 2.57×10^1 | 1.63×10^1 | 0.74 | 0.49 | 0.49 | 0.28 |
| 11 | 4.52×10^4 | 4.35×10^4 | 1.22×10^7 | 280.7x | 6.78×10^1 | 3.93×10^1 | 1.42 | 0.98 | 0.72 | 0.60 |
| 12 | 9.67×10^4 | 9.20×10^4 | 5.45×10^7 | 592.7x | 1.84×10^2 | 1.21×10^2 | 2.58 | 2.18 | 2.02 | 1.76 |
| 13 | 2.04×10^5 | 1.93×10^5 | 1.59×10^8 | 825.5x | 5.22×10^2 | 3.72×10^2 | 6.16 | 6.09 | 5.74 | 5.44 |
| 14 | 4.26×10^5 | 3.99×10^5 | 4.55×10^8 | 1140.6x | 1.47×10^3 | 1.10×10^3 | 1.52×10^1 | 1.38×10^1 | 1.79×10^1 | 1.75×10^1 |
| 15 | 8.83×10^5 | 8.20×10^5 | 1.30×10^9 | 1582.5x | 3.97×10^3 | 3.29×10^3 | 5.65×10^1 | 3.82×10^1 | | |
| 16 | 1.82×10^6 | 1.67×10^6 | 3.62×10^9 | 2167.0x | 1.22×10^4 | 8.77×10^3 | | | | |

TABLE III: 3Truss experiment statistics. D4M and MTJ execute faster, assuming sufficient memory, due to high overhead. Graphulo overhead is 1.6x faster than D4M and 1.5x faster than MTJ. Graphulo overhead is listed in seconds.

Graphulo Overhead: How many more entries Graphulo writes to Accumulo than the main-memory systems

(Overhead is due to Graphulo writing all partial products whereas main-memory systems pre-sum)

**Comparable I/O;
savings in communication
outweighs I/O overhead**

| SCALE | nnz(A) | nnz(3Truss(A)) | Partial Products | Graphulo Overhead | Graphulo 1 Tablet | Graphulo 2 Tablets | D4M 1 Tablet | D4M 2 Tablets | MTJ 1 Tablet | MTJ 2 Tablets |
|-------|--------------------|-----------------------------|---------------------|----------------------|----------------------|-----------------------|--------------------|--------------------|--------------------|--------------------|
| 10 | 2.10×10^4 | 2.03×10^4 | 5.94×10^6 | 293.3x | 2.57×10^1 | 1.63×10^1 | 0.74 | 0.49 | 0.49 | 0.28 |
| 11 | 4.52×10^4 | 4.35×10^4 | 1.22×10^7 | 280.7x | 6.78×10^1 | 3.93×10^1 | 1.42 | 0.98 | 0.72 | 0.60 |
| 12 | 9.67×10^4 | 9.20×10^4 | 5.45×10^7 | 592.7x | 1.84×10^2 | 1.21×10^2 | 2.58 | 2.18 | 2.02 | 1.76 |
| 13 | 2.04×10^5 | 1.93×10^5 | 1.59×10^8 | 825.5x | 5.22×10^2 | 3.72×10^2 | 6.16 | 6.09 | 5.74 | 5.44 |
| 14 | 4.26×10^5 | 3.99×10^5 | 4.55×10^8 | 1140.6x | 1.47×10^3 | 1.10×10^3 | 1.52×10^1 | 1.38×10^1 | 1.79×10^1 | 1.75×10^1 |
| 15 | 8.83×10^5 | 8.20×10^5 | 1.30×10^9 | 1582.5x | 3.97×10^3 | 3.29×10^3 | 5.65×10^1 | 3.82×10^1 | | |
| 16 | 1.82×10^6 | 1.67×10^6 | 3.62×10^9 | 2167.0x | 1.22×10^4 | 8.77×10^3 | | | | |

Graphulo Overhead: How many more entries Graphulo writes to Accumulo than the main-memory systems

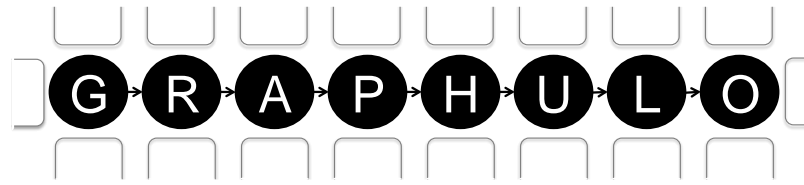
(Overhead is due to Graphulo writing all partial products whereas main-memory systems pre-sum)

**Comparable I/O;
savings in communication
outweighs I/O overhead**

**Vastly different I/O;
main-memory systems cache
temporary tables whereas
Graphulo writes all to Accumulo**

**Vastly different I/O;
main-memory systems cache
temporary tables whereas
Graphulo writes all to Accumulo**

Graphulo Overhead: How many more entries Graphulo writes to Accumulo than the main-memory systems



<http://graphulo.mit.edu>

- > ***How to do matrix math in Accumulo?***
 - The TwoTable iterator pipeline
- > **Jaccard & k-Truss**
- > ***When to do matrix math in Accumulo?***
 - Memory requirements
 - Compare in-database I/O vs. alternatives
- > **Future Work: *Multi-Node*,
expand to *Relational Algebra*,
use an *Optimizer* to choose the best plan**

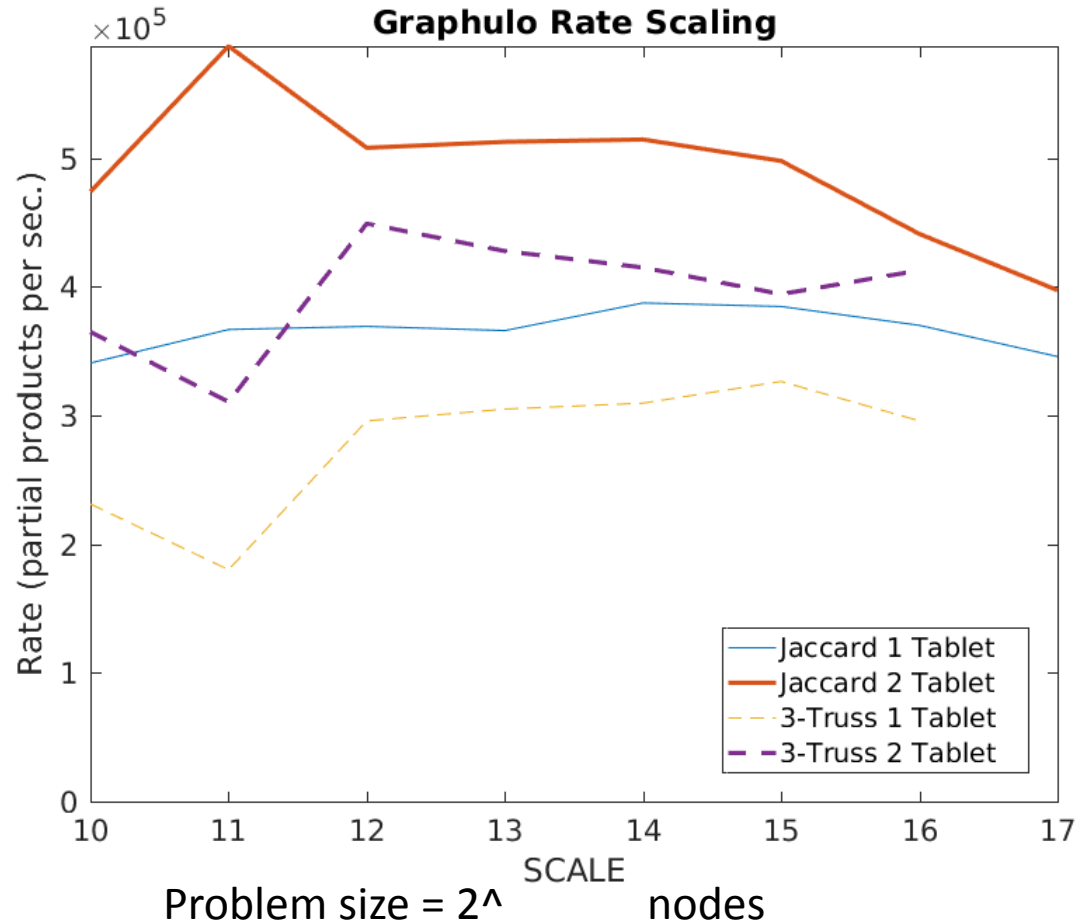


Contact: Dylan Hutchison
dhutchis@cs.washington.edu



Backup

Overall Graphulo Performance



Alg 1: Jaccard Coefficients

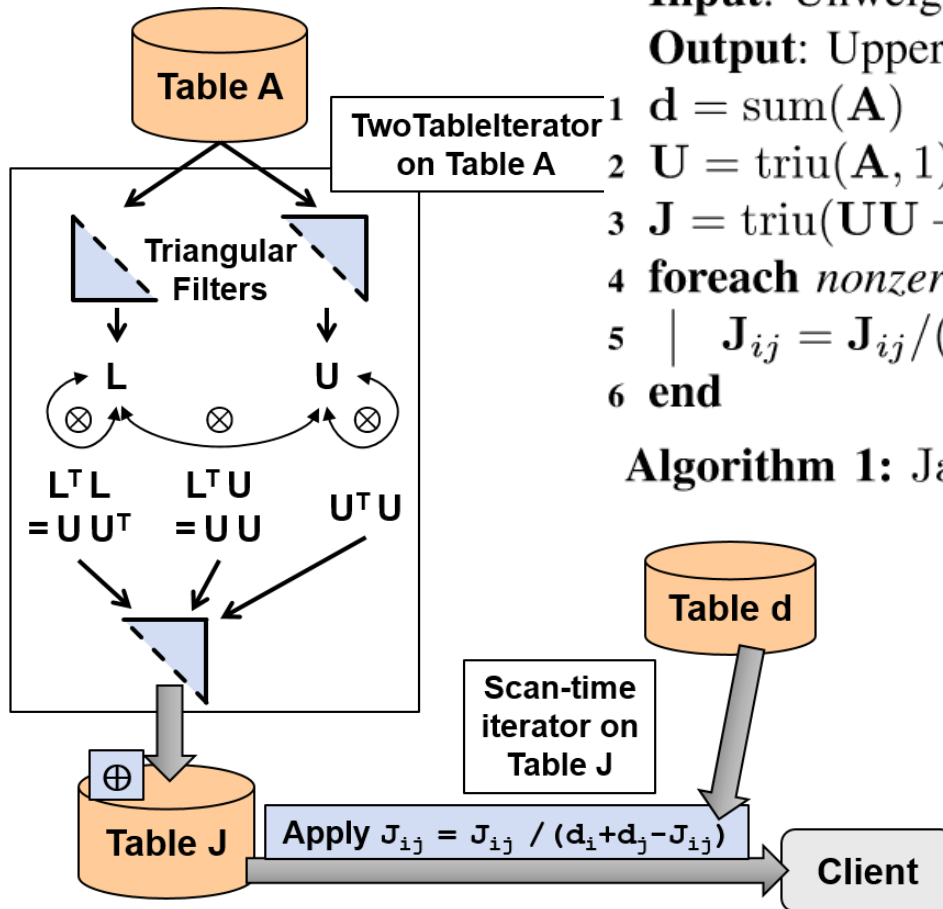
Input: Unweighted, undirected adjacency matrix A

Output: Upper triangle of Jaccard coefficients J

```

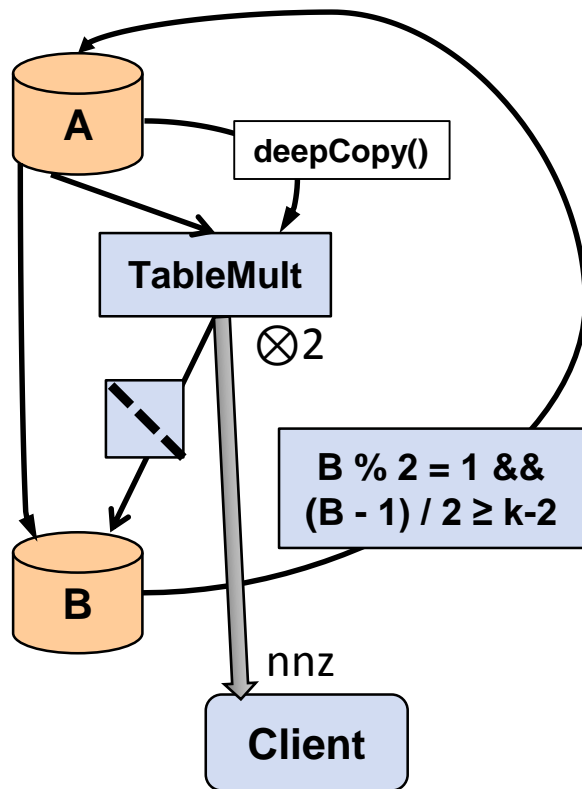
1  $d = \text{sum}(A)$  // pre-computed in degree table
2  $U = \text{triu}(A, 1)$  // strict upper triangle filter
3  $J = \text{triu}(UU + UU^T + U^T U, 1)$  // fused MxM
4 foreach nonzero entry  $J_{ij} \in J$  do
5    $J_{ij} = J_{ij} / (d_i + d_j - J_{ij})$  // stateful Apply on J
6 end
    
```

Algorithm 1: Jaccard



Fusion:
A single Graphulo
TwoTable pass!

Alg 2: k-Truss Subgraph



Input: Unweighted, undirected adjacency matrix A_0 , integer k

Output: Adjacency matrix of k -truss subgraph A

```

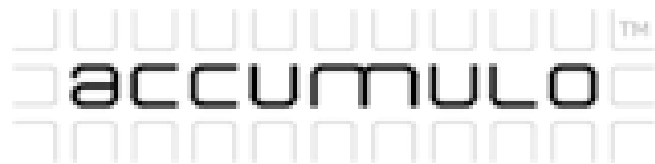
1  $z' = \infty, A = A_0$  // table clone
2 repeat
3    $z = z'$ 
4    $B = A$  // table clone
5    $B = B + 2AA$  //  $M \times M$  with  $a \otimes b = 2$  if  $a, b \neq 0$ 
6    $B(B \% 2 == 0) = 0$  // filter on B
7    $B((B - 1)/2 < k - 2) = 0$  // filter on B
8    $A = |B|_0$  // Apply on B; switch  $A \leftrightarrow B$ 
9    $z' = \text{nnz}(A)$  // Reduce, gathering nnz at client
10 until  $z == z'$  // client controls iteration
  
```

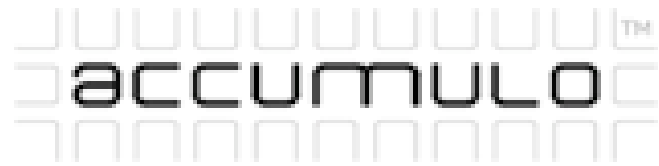
Algorithm 2: kTruss

Iterations of Graphulo TwoTables

| Key | | | | | Value |
|-----|--------|-----------|------------|-----------|-------|
| Row | Column | | | Timestamp | |
| | Family | Qualifier | Visibility | | |

Best for:

- > **Large, de-normalized tables; no schema necessary**
 - Unlimited columns; un-interpreted values; everything is a byte[]
 - > **TBs to PBs of data; robust horizontal scaling**
 - > **Hadoop HDFS / Java ecosystem**
 - > **Cell-level visibility**
 - > **Row store by default**
 - Scan over rows for $O(\log n)$ lookup & sorted order
 - Use Transpose Tables for column indexing¹
 - > **Iterator processing framework**
- 
- ¹D4M 2.0 Schema: A General Purpose High



¹*D4M 2.0 Schema: A General Purpose High Performance Schema for the Accumulo Database*, Kepner et al, IEEE HPEC 2013

GraphBLAS Operations

| GraphBLAS Kernel | Graphulo Implementation |
|---------------------------|---|
| BuildMatrix (\oplus) | Accumulo BatchWriter |
| ExtractTuples | Accumulo BatchScanner |
| MxM (\oplus, \otimes) | TwoTableIterator ROW mode, performing $A^T B$ |
| EwiseMult (\otimes) | TwoTableIterator EWISE mode |
| EwiseAdd (\oplus) | Similar to EwiseMult, with non-matching entries |
| Extract | Row and column filtering |
| Apply (f) | Extra Iterators |
| Assign | Apply with a key-transforming function |
| Reduce (\oplus) | Reducer module on RemoteWriteIterator |
| Transpose | Transpose option on RemoteWriteIterator |