

# EmptyHeaded: A Relational Engine for Graph Processing

Jason Priest  
2018-02-28  
6.886

# Key Contributions

- Join optimization based on *generalized hypertree decomposition*
- Data and algorithm optimization based on local graph skew



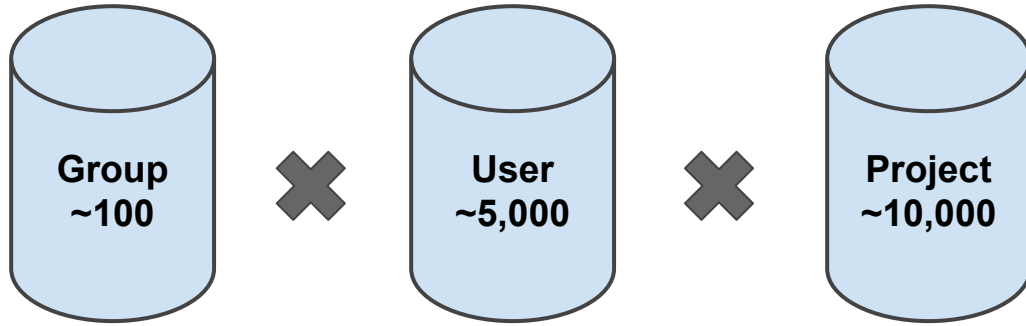
# Traditional Relational Query

How many github organizations (groups) have a C++ developer who registered this year?

```
select count(*)  
  from group  
  join user      on group.users contains user.id  
  join project  on project.owner is      user.id  
  where user.joined after "2018-01-01"  
  and  project.lang is "C++";
```



# Join Order Matters



# Cyclic Queries

$$Q_{\Delta} = R(A, B) \bowtie S(B, C) \bowtie T(A, C).$$

- Joins usually implemented pairwise, or between two sets at a time
- For cyclic queries such as above, this leads to suboptimal behavior

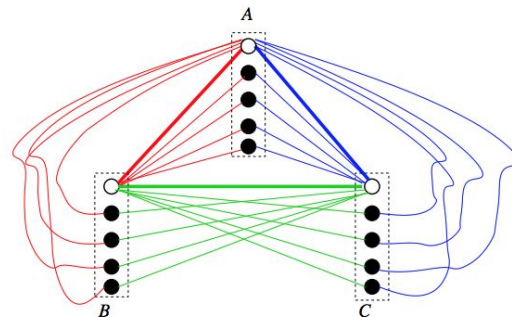


# Pairwise Joins Insufficient

$$R = \{a_0\} \times \{b_0, \dots, b_m\} \cup \{a_0, \dots, a_m\} \times \{b_0\}$$

$$S = \{b_0\} \times \{c_0, \dots, c_m\} \cup \{b_0, \dots, b_m\} \times \{c_0\}$$

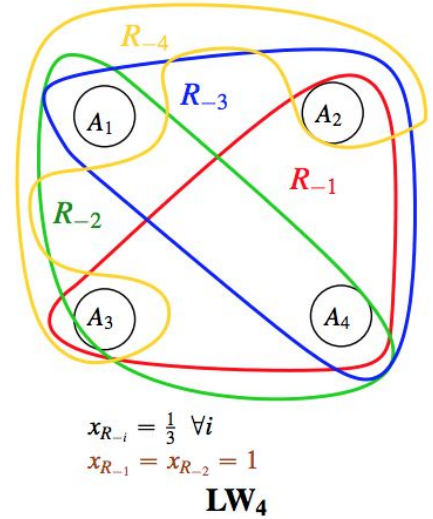
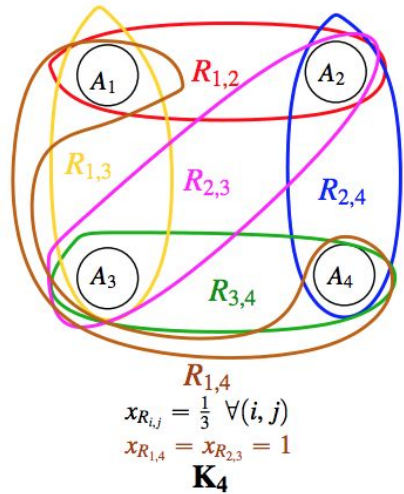
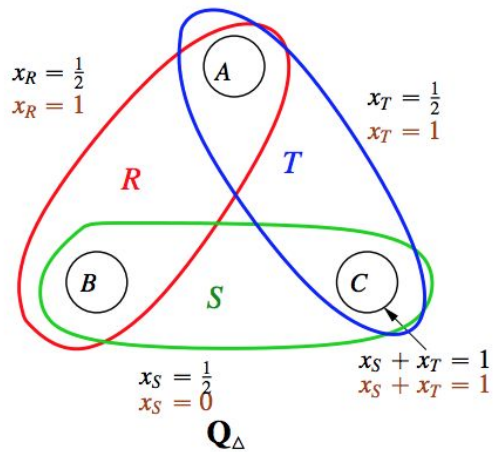
$$T = \{a_0\} \times \{c_0, \dots, c_m\} \cup \{a_0, \dots, a_m\} \times \{c_0\}$$



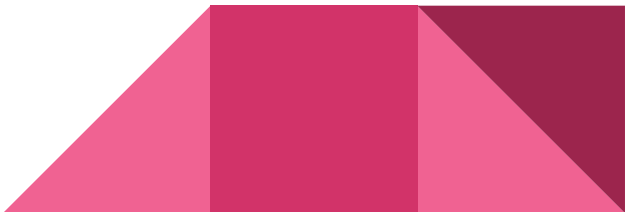
**Figure 2: Counter-example for join-project only plans for the triangles (left) and an illustration for  $m = 4$  (right). The pairs connected by the red/green/blue edges form the tuples in the relations  $R/S/T$  respectively. Note that in this case each relation has  $N = 2m + 1 = 9$  tuples and there are  $3m + 1 = 13$  output tuples in  $Q_{\Delta}$ . Any pair-wise join however has size  $m^2 + m = 20$ .**

# Worst-Case Optimal Join Algorithm

# Graph Covers

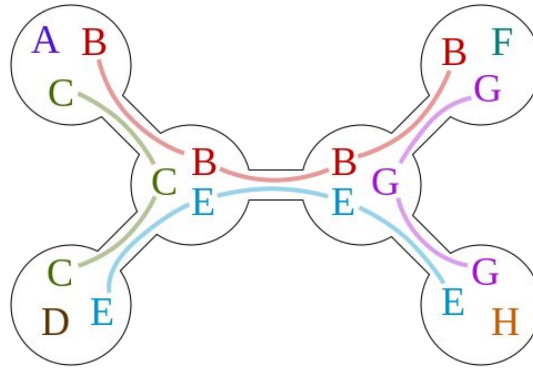
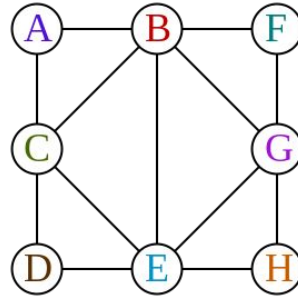


$$|Q| = |\times_{F \in \mathcal{E}} R_F| \leq \prod_{F \in \mathcal{E}} |R_F|^{x_F}. \quad (6)$$





# Graph Decomposition



[https://en.wikipedia.org/wiki/Tree\\_decomposition](https://en.wikipedia.org/wiki/Tree_decomposition)

# Worst-Case Optimal Join Algorithm

- Brute force all decompositions of relations
- Find minimum width decomposition
- Use decomposition to inform much order of joins, order of comparing fields, etc
- Resulting plan for joins is optimal intermediate output sizes

---

## ALGORITHM 2: Enumerating all GHDs via Brute Force Search

---

```
1 // Input: Hypergraph H=(V,E) and a set of parent edges P.
2 // Output: A list of GHDs in the form of (GHD node, subtree) pairs.
3 GHD-Enumeration(V,E,P):
4   GHDs = []
5   // Iterate over all subset combinations of edges.
6   for [C | C ⊆ E] do
7     // The remaining edges, not in C.
8     R = E - C
9     // If the running intersection property is broken, the GHD is
10    // not valid. The check makes sure that all attributes in the
11    // parent and subtree of a specified GHD node also appear
12    // within the specified GHD node. Here we use ∪ on a set
13    // of edges to indicate the union of their attributes.
14    if not ( (∪P) ∩ (∪R) ⊆ ∪C ) then continue
15    // Consider each subgraph of the remaining edges. For each
16    // subgraph, recursively enumerate all possible GHDs.
17    PartitionChildren = []
18    for Pr in Partition(R) do
19      PartitionChildren += [GHD-Enumeration(∪Pr,Pr,C)]
20      // Consider all possible combinations of subtrees by calling the
21      // method below. For each, construct a GHD with C as the root.
22    for Ch in Subtree-Combinations(PartitionChildren) do
23      GHDs += [(C,Ch)]
24    return GHDs
25
26 // Input: A list of lists of GHDs: for each partition of a hypergraph
27 // (the outer list), all possible decompositions for that partition
28 // (the inner lists).
29 // Output: A list where each member of this list is a list that
30 // contains one subtree from each partition.
31 Subtree-Combinations(PartitionChildren):
32   ChildrenCombinations = []
33   if |PartitionChildren| > 0 then
34     if |PartitionChildren| = 1 then
35       // If there is only one partition, for each of the possible GHDs
36       // of this partition, add a combination with just this GHD.
37       for Ch in PartitionChildren[0] do
38         ChildrenCombinations += [[Ch]]
39     else
40       // Recursively generate combinations for the partitions after
41       // the first one.
42       RemainingCombinations = Subtree-Combinations(PartitionChildren[1..])
43       // If there is more than one partition, each subtree in the
44       // first partition is combined with each list of subtrees in
45       // the recursively generated combinations for the remaining
46       // partitions.
47       for Ch in PartitionChildren[0] do
48         for C in RemainingCombinations do
49           FinalCombination = [Ch] + C
50           ChildrenCombinations += [FinalCombination]
51   return ChildrenCombinations
```

---

# Execution Engine

# Skew

- Density skew
  - some values are much more common
  - some relations are much more selective
- Cardinality skew
  - some tables are much larger
  - some nodes have much greater degree



# Set Layout

- bitsets: bitvectors with offsets to first element
- pshort: nearby values may have similar prefix, thus store repeated high 16 bits
- varint: difference encoding with continue bits
- uint: sorted array with binary search



# Set Intersection Algorithm

- uint & uint
  - 5 SIMD techniques, chosen based on density skew and cardinality
- bitset & bitset
  - Just SIMD AND comparison



# Evaluation

# Simplicity

Name	Query Syntax
Triangle	$\text{Triangle}(x, y, z): -R(x, y), S(y, z), T(x, z).$
4-Clique	$4\text{Clique}(x, y, z, w): -R(x, y), S(y, z), T(x, z), U(x, w), V(y, w), Q(z, w).$
Lollipop	$\text{Lollipop}(x, y, z, w): -R(x, y), S(y, z), T(x, z), U(x, w).$
Barbell	$\text{Barbell}(x, y, z, x', y', z'): -R(x, y), S(y, z), T(x, z), U(x, x'),$ $R'(x', y'), S'(y', z'), T'(x', z').$
Count Triangle	$\text{CntTriangle}(; w: \text{long}): -R(x, y), S(x, z), T(x, z); w = \ll \text{COUNT}(\ast) \gg.$
4-Clique-Selection	$S4\text{Clique}(x, y, z, w): -R(x, y), S(y, z), T(x, z), U(x, w),$ $V(y, w), Q(z, w), P(x, \text{'node'}).$
Barbell-Selection	$S\text{Barbell}(x, y, z, x', y', z'): -R(x, y), S(y, z), T(x, z), U(x, \text{'node'}),$ $V(\text{'node'}, x'), R'(x', y'), S'(y', z'), T'(x', z').$
PageRank	$N(; w: \text{int}): -\text{Edge}(x, y); w = \ll \text{COUNT}(x) \gg.$ $\text{PageRank}(x; y: \text{float}): -\text{Edge}(x, z); y = 1/N.$ $\text{PageRank}(x; y: \text{float}) * [i=5]: -\text{Edge}(x, z), \text{PageRank}(z), \text{InvDeg}(z);$ $y = 0.15 + 0.85 * \ll \text{SUM}(z) \gg.$
SSSP	$\text{SSSP}(x; y: \text{int}): -\text{Edge}(\text{'start'}, x); y = 1.$ $\text{SSSP}(x; y: \text{int}) * : -\text{Edge}(w, x), \text{SSSP}(w); y = \ll \text{MIN}(w) \gg + 1.$



# Performance

Table 9. Triangle Counting Runtime (in Seconds) for EmptyHeaded and Relative Slowdown for Other Engines Including PowerGraph, a Commercial Graph Tool (CGT-X), Snap-Ringo, Socialite, and LogicBlox

Dataset	EmptyHeaded	Low-Level			High-Level	
		PowerGraph	CGT-X	Snap-Ringo	Socialite	LogicBlox
Google+	<b>0.31</b>	8.40×	62.19×	4.18×	1390.75×	83.74×
Higgs	<b>0.15</b>	3.25×	57.96×	5.84×	387.41×	29.13×
LiveJournal	<b>0.48</b>	5.17×	3.85×	10.72×	225.97×	23.53×
Orkut	<b>2.36</b>	2.94×	-	4.09×	191.84×	19.24×
Patents	<b>0.14</b>	10.20×	7.45×	22.14×	49.12×	27.82×
Twitter	<b>56.81</b>	4.40×	-	2.22×	t/o	30.60×

48 threads used for all engines. “-” indicates the engine does not process over 70 million edges. “t/o” indicates the engine ran for over 30 minutes.

# Performance

Table 11. SSSP Runtime (in Seconds) Using 48 Threads for All Engines

Dataset	EmptyHeaded	Low-Level			High-Level	
		Galois	PowerGraph	CGT-X	SociaLite	LogicBlox
Google+	0.024	<b>0.008</b>	0.22	0.51	0.27	41.81
Higgs	0.035	<b>0.017</b>	0.34	0.91	0.85	58.68
LiveJournal	0.19	<b>0.062</b>	1.80	-	3.40	102.83
Orkut	0.24	<b>0.079</b>	2.30	-	7.33	215.25
Patents	0.15	<b>0.054</b>	1.40	4.70	3.97	159.12
Twitter	7.87	<b>2.52</b>	36.90	-	x	379.16

“-” indicates the engine does not process over 70 million edges. The other engines include Galois, PowerGraph, a commercial graph tool (CGT-X), SociaLite, and LogicBlox. “x” indicates the engine did not compute the query properly.