

Pessimistic Query Optimization: Tighter Upper Bounds for Intermediate Join Cardinalities

Walter Cai Magdalena Balazinska Dan Suciu

University of Washington

[walter,magda,suciu}@cs.washington.edu

July 19th, 2019

Systematic Underestimation

Query optimizers assume:

- ▶ Uniformity
- ▶ Independence

1 Background: Cardinality Bounds

2 Tightened Cardinality Bounds

3 Evaluation

1 Background: Cardinality Bounds

2 Tightened Cardinality Bounds

3 Evaluation

Example Query (SQL)

```
SELECT
    *
FROM
    pseudonym ,
    cast ,
    movie_companies ,
    company_name
WHERE
    pseudonym.person_id = cast.person_id AND
    cast.movie_id_id = movie_companies.movie_id AND
    movie_companies.company_id = company_name.id;
```

Example Query (Join Graph & Datalog)

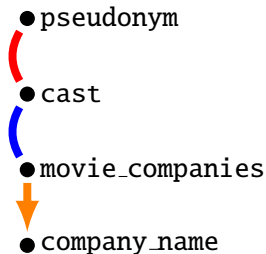


Figure: Join Graph.

$$Q(x, y, z, w) :- \text{pseudo}(x, y),$$

$$\text{cast}(y, z),$$

$$\text{mc}(z, w),$$

$$\text{cn}(w)$$

$y \mapsto \text{person}$

$z \mapsto \text{movie}$

$w \mapsto \text{company}$

Review: Entropy

Take random variable X :

$$h(X) = - \sum_a \mathbb{P}(X = a) \cdot \log(\mathbb{P}(X = a))$$

Multiple variables:

$$h(X, Y) = - \sum_{a,b} \mathbb{P}(X = a, Y = b) \cdot \log(\mathbb{P}(X = a, Y = b))$$

Conditional entropy:

$$h(X|Y) = - \sum_{a,b} \mathbb{P}(X = a, Y = b) \cdot \log\left(\frac{\mathbb{P}(X = a, Y = b)}{\mathbb{P}(Y = b)}\right)$$

Review: Entropy

$$X \sim P(x_1, \dots, x_n)$$

- ▶ Fact: $h(X) \leq \log(n)$
- ▶ $h(X) = \log(n)$ iff P is uniform

Connection to Entropy

$$Q(x, y, z, w) := \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$$

- ▶ Create random variable for each attribute.

$$x \rightarrow X, \quad y \rightarrow Y, \quad z \rightarrow Z, \quad w \rightarrow W$$

- ▶ Let (X, Y, Z, W) be uniformly distributed over true output of Q .

$$h(X, Y, Z, W) = \log |Q(x, y, z, w)|$$

$$\exp(h(X, Y, Z, W)) = |Q(x, y, z, w)|$$

Entropic Bound

$$|Q(x, y, z, w)| = \exp(h(X, Y, Z, W))$$

- ▶ Suffices to bound $h(X, Y, Z, W)$.

Entropic Bound

$$h(X, Y, Z, W) \leq h(X|Y) + h(Y, Z) + h(W|Z)$$

Entropic Bound

$$|Q(x, y, z, w)| = \exp(h(X, Y, Z, W)) \leq \exp(h(X|Y) + h(Y, Z) + h(W|Z))$$

$$h(Y, Z) \leq \log(\text{count}(\text{cast}))$$

$$h(X|Y) \leq \log(\max \text{ degree}(\text{pseudonym}))$$

$$h(W|Z) \leq \log(\max \text{ degree}(\text{movie_companies}))$$

Entropic Bound

$$|Q(x, y, z, w)| = \exp(h(X, Y, Z, W)) \leq \exp(h(X|Y) + h(Y, Z) + h(W|Z))$$

$$h(Y, Z) \leq \log c_{\text{cast}}$$

$$h(X|Y) \leq \log d_{\text{pseudo}}^y$$

$$h(W|Z) \leq \log d_{\text{mc}}^z$$

Cardinality Bound

$$\begin{aligned}
 |Q(x, y, z, w)| &= \exp(h(X, Y, Z, W)) \\
 &\leq \exp(\underbrace{h(X|Y)}_{\leq \log d_{\text{pseudo}}^y} + \underbrace{h(Y, Z)}_{\leq \log c_{\text{cast}}} + \underbrace{h(W|Z)}_{\leq \log d_{\text{mc}}^z}) \\
 &\leq d_{\text{pseudo}}^y \cdot c_{\text{cast}} \cdot d_{\text{mc}}^z
 \end{aligned}$$

Many Entropic Bounds

$$h(X, Y, Z, W) \leq \dots$$

$$h(X, Y) + h(Z|Y) + h(W|Z)$$

$$h(X, Y) + h(Z|Y) + h(W)$$

$$h(X, Y) + h(Z, W)$$

$$h(X, Y) + h(Z|W) + h(W)$$

$$h(X|Y) + h(Y, Z) + h(W|Z)$$

$$h(X|Y) + h(Y|Z) + h(Z, W)$$

$$h(X|Y) + h(Y|Z) + h(Z|W) + h(Z)$$

Entropic Bounds

$$Q(x, y, z, w) := \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$$

$$|Q(x, y, z, w)| \leq \min \left\{ \begin{array}{l} c_{\text{pseudo}} \cdot d_{\text{cast}}^y \cdot d_{\text{mc}}^z \\ c_{\text{pseudo}} \cdot d_{\text{cast}}^y \cdot c_{\text{cn}} \\ c_{\text{pseudo}} \cdot c_{\text{mc}} \\ c_{\text{pseudo}} \cdot d_{\text{mc}}^w \cdot c_{\text{cn}} \\ d_{\text{pseudo}}^y \cdot c_{\text{cast}} \cdot d_{\text{mc}}^z \\ d_{\text{pseudo}}^y \cdot c_{\text{cast}} \cdot c_{\text{cn}} \\ d_{\text{pseudo}}^y \cdot d_{\text{cast}}^z \cdot c_{\text{mc}} \\ d_{\text{pseudo}}^y \cdot d_{\text{cast}}^z \cdot d_{\text{mc}}^w \cdot c_{\text{cn}} \end{array} \right.$$

Entropic Bounds

Neat! But is it useful?

- ▶ Short answer: 'No'. (Not yet, anyway)
 - ▶ Bounds still too loose (overestimation)
 - ▶ Need to tighten
- ▶ How to tighten? Partitioning

- 1 Background: Cardinality Bounds
- 2 Tightened Cardinality Bounds**
- 3 Evaluation

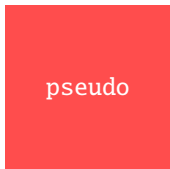
$$Q(x, y, z, w) :- \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$$

pseudonym

cast

movie_companies

company_name

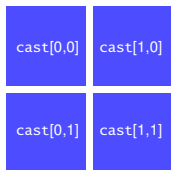


$$Q(x, y, z, w) := \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$$

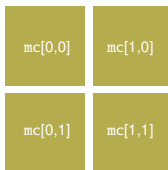
pseudonym



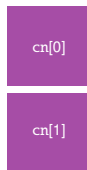
cast



movie_companies

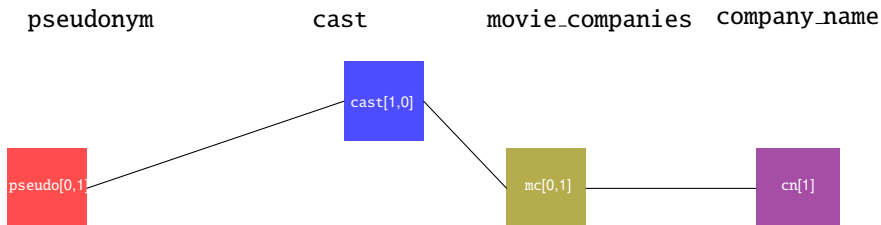


company_name



- ▶ Value based hashing
- ▶ Analagous to hash join

$$Q(x, y, z, w) := \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$$



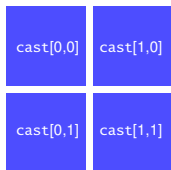
- ▶ Value based hashing
- ▶ Analagous to hash join

$$Q(x, y, z, w) := \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$$

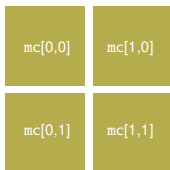
pseudonym



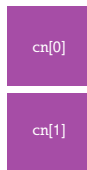
cast



movie_companies



company_name



- ▶ $Q(D)$: query evaluated on database D
- ▶ $Q(D[J])$: query evaluated on partition $D[J]$

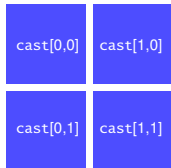
$$Q(D) = \bigcup_J Q(D[J])$$

$$Q(x, y, z, w) := \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$$

pseudonym



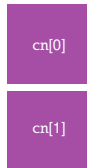
cast



movie_companies



company_name



- ▶ Bound each partition $D[J]$
- ▶ Sum will be bound on full database

$$Q(D) = \bigcup_J Q(D[J])$$

$$|Q(D)| \leq \sum_J \text{bound}(Q(D[J]))$$

Partition Bounding

$$|Q(D)| \leq \sum_{J \in \{0,1\}^4} \min \left\{ \begin{array}{l} C_{\text{pseudo}[J]} \cdot d_{\text{cast}[J]}^y \cdot d_{\text{mc}[J]}^z \\ C_{\text{pseudo}[J]} \cdot d_{\text{cast}[J]}^y \cdot C_{\text{cn}[J]} \\ C_{\text{pseudo}[J]} \cdot C_{\text{mc}[J]} \\ C_{\text{pseudo}[J]} \cdot d_{\text{mc}[J]}^w \cdot C_{\text{cn}[J]} \\ d_{\text{pseudo}[J]}^y \cdot C_{\text{cast}[J]} \cdot d_{\text{mc}[J]}^z \\ d_{\text{pseudo}[J]}^y \cdot C_{\text{cast}[J]} \cdot C_{\text{cn}[J]} \\ d_{\text{pseudo}[J]}^y \cdot d_{\text{cast}[J]}^z \cdot C_{\text{mc}[J]} \\ d_{\text{pseudo}[J]}^y \cdot d_{\text{cast}[J]}^z \cdot d_{\text{mc}[J]}^w \cdot C_{\text{cn}[J]} \end{array} \right.$$

Optimizations

- ▶ Bound Formula Generation
- ▶ Partition Budgeting
 - ▶ Combats exponential runtime w.r.t. hash size
 - ▶ Non-monotonic behaviour
- ▶ Filter Predicates

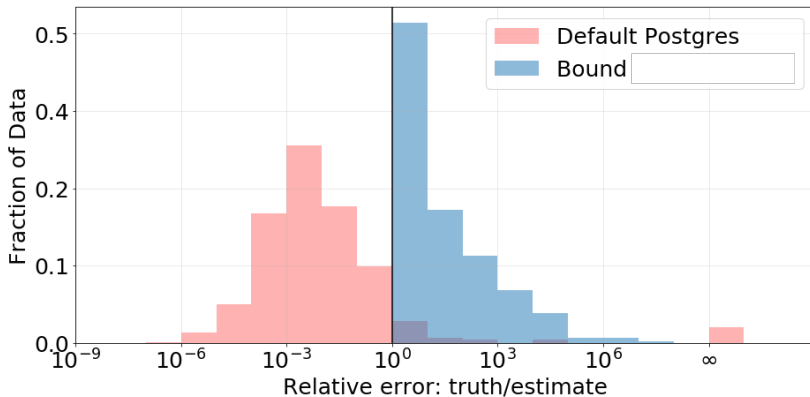
- 1 Background: Cardinality Bounds
- 2 Tightened Cardinality Bounds
- 3 Evaluation**

Join Order Benchmark¹

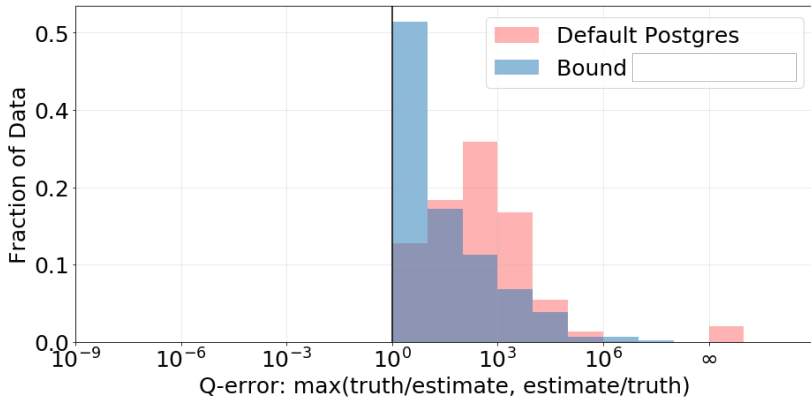
- ▶ Built on the IMDb dataset
 - ▶ 113 queries
 - ▶ 33 unique topologies
 - ▶ Skew!
 - ▶ Correlation!
 - ▶ Complex selection predicates!

¹ *How Good Are Query Optimizers, Really?* Leis et al. VLDB 2015.

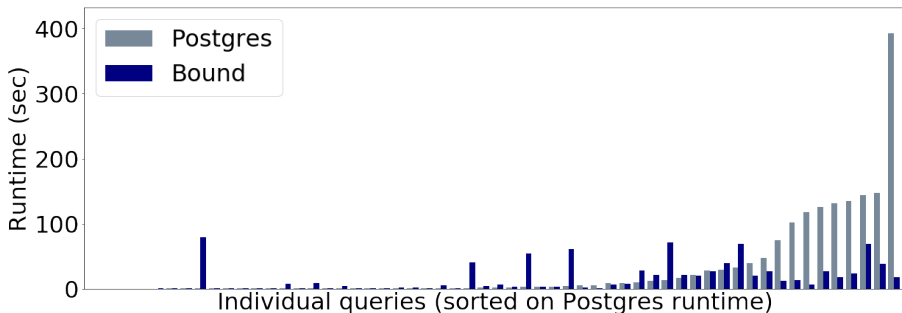
Bound Relative Error Versus Postgres Relative Error



Bound Q-Error Versus Postgres Q-Error



Plan Execution Runtime (With Foreign Keys Indexes)



Plan Execution Runtime (With Foreign Keys Indexes)

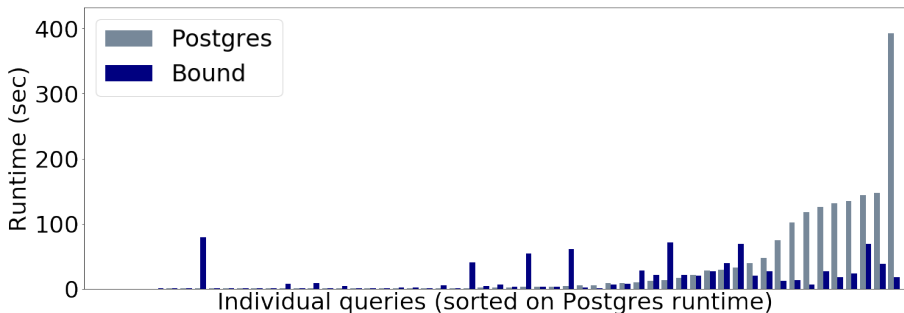


Figure: Linear scale runtime improvements over JOB queries.

- ▶ Total runtime:
 - ▶ Postgres: 3,190 seconds
 - ▶ Tightened Bound: 1,832 seconds

Plan Execution Runtime (No Foreign Key Indexes)

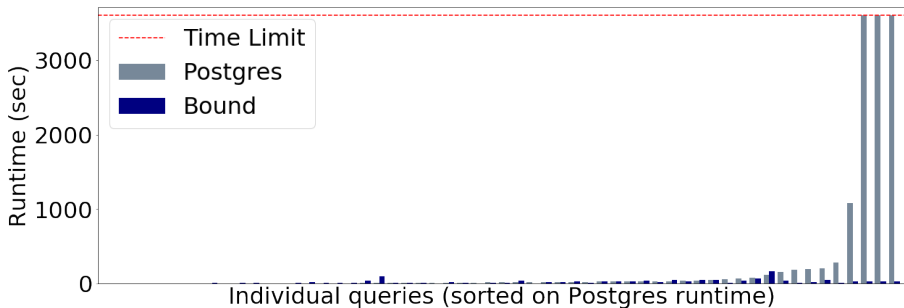


Figure: Linear scale plan execution time over JOB queries.

- ▶ Total runtime (including 1 hour cutoff for default Postgres):
 - ▶ Postgres: 21,125 seconds
 - ▶ Tightened Bound: 2,216 seconds

Plan Execution Runtime (No Foreign Key Indexes)

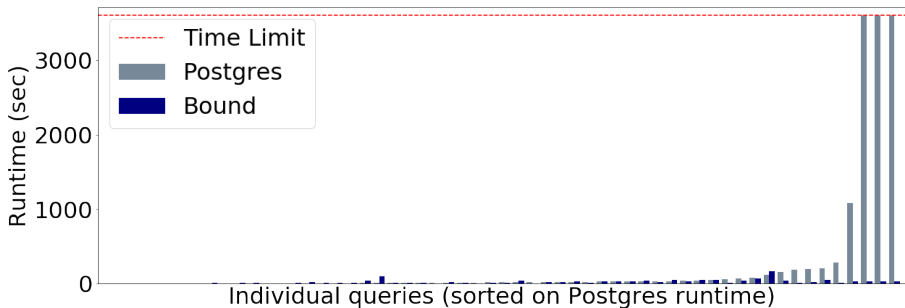


Figure: Linear scale plan execution time over JOB queries.

- ▶ Total runtime (including 1 hour cutoff for default Postgres):
 - ▶ Postgres: 21,125 seconds
 - ▶ Tightened Bound: 2,216 seconds

Takeaways

- ▶ Gains for very slow queries.
- ▶ On par for fast queries.
- ▶ Pessimistic but robust query optimization.

Acknowledgements

- ▶ Thank you to Jenny, Tomer, Laurel, Brandon, Jingjing, Tobin, Leilani, and Guna!
- ▶ This research is supported by NSF grant AITF 1535565 and III 1614738.

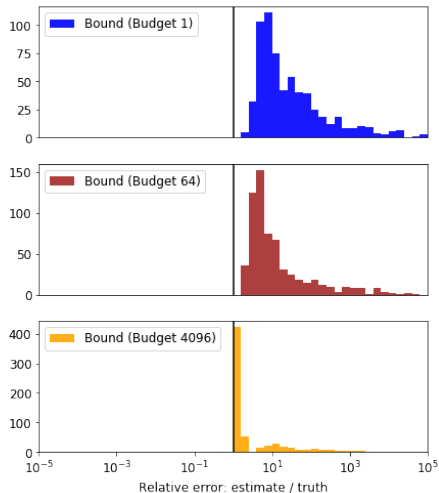
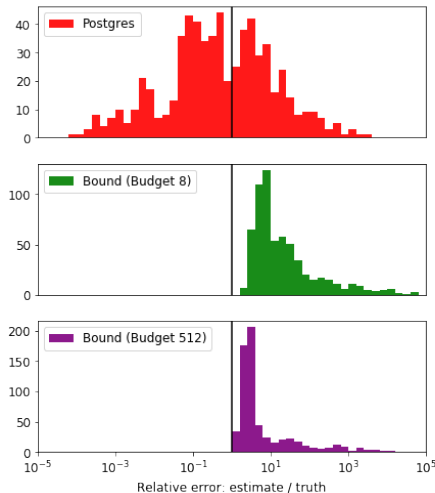
Googleplus Microbenchmark Examples

```
SELECT COUNT(*)  
FROM  
  community_44 AS t0,  
  community_44 AS t1,  
  community_44 AS t2,  
  community_44 AS t3  
WHERE  
  t0.object = t1.subject AND  
  t1.object = t2.subject AND  
  t2.object = t3.subject AND  
  t0.subject % 512 = 89 AND  
  t3.object % 512 = 174;
```

Googleplus Microbenchmark Examples

```
SELECT COUNT(*)
FROM
  community_30 AS t0,
  community_30 AS t1,
  community_30 AS t2,
  community_30 AS t3,
  community_30 AS t4
WHERE
  t0.object = t1.subject AND
  t0.object = t2.subject AND
  t0.object = t3.subject AND
  t3.object = t4.subject AND
  t0.subject % 256 = 49 AND
  t1.object % 256 = 213 AND
  t2.object % 256 = 152 AND
  t4.object % 256 = 248;
AND ci.movie_id = mc.movie_id;
```

Googleplus Progressive Bound Tightness



Example of Non-monotonic Behavior

$Q(x, y, z, w) :- R(z, y), S(y, z), T(z, w)$

x	y		y	z		z	w		x	y	z	w
0	0		0	0		0	0		0	0	0	0
0	1	⊗	1	0	⊗	1	1	=	0	1	0	0
1	0		2	1		2	2		1	0	0	0
1	1		3	1		3	3		1	1	0	0
R			S			T			Q			

Example of Non-monotonic Behavior

x	y		y	z		z	w		x	y	z	w
0	0		0	0		0	0		0	0	0	0
0	1	⋈	1	0	⋈	1	1	=	0	1	0	0
1	0		2	1		2	2		1	0	0	0
1	1		3	1		3	3		1	1	0	0
R			S			T			Q			

$$|Q(x, y, z, w)| \leq \min \begin{cases} c_R \cdot d_S^y \cdot d_T^z \\ d_R^y \cdot c_S \cdot d_T^z \\ d_R^y \cdot d_S^z \cdot c_T \\ c_R \cdot c_T \end{cases}$$

$$c_{R(0)} \cdot d_{S(0,0)}^y \cdot d_{T(0)}^z = 4 \cdot 1 \cdot 1 = 4$$

Example of Non-monotonic Behavior

x	y		y	z		z	w		x	y	z	w
0	0		0	0		0	0		0	0	0	0
0	1	⋈	1	0	⋈	1	1	=	0	1	0	0
1	0		2	1		2	2		1	0	0	0
1	1		3	1		3	3		1	1	0	0
<i>R</i>			<i>S</i>			<i>T</i>			<i>Q</i>			

$$|Q(x, y, z, w)| \leq \min \begin{cases} c_R \cdot d_S^y \cdot d_T^z \\ d_R^y \cdot c_S \cdot d_T^z \\ d_R^y \cdot d_S^z \cdot c_T \\ c_R \cdot c_T \end{cases}$$

$$c_{R(0)} \cdot d_{S(0,0)}^y \cdot d_{T(0)}^z = 4 \cdot 1 \cdot 1 = 4$$

Example of Non-monotonic Behavior

x	y		y	z		z	w		x	y	z	w
0	0		0	0		0	0		0	0	0	0
0	1	⋈	1	0	⋈	1	1	=	0	1	0	0
1	0		2	1		2	2		1	0	0	0
1	1		3	1		3	3		1	1	0	0
R			S			T			Q			

$$|Q(x, y, z, w)| \leq \min \begin{cases} c_R \cdot d_S^y \cdot d_T^z \\ d_R^y \cdot c_S \cdot d_T^z \\ d_R^y \cdot d_S^z \cdot c_T \\ c_R \cdot c_T \end{cases}$$

$$c_{R(0)} \cdot d_{S(0,0)}^y \cdot d_{T(0)}^z = 4 \cdot 1 \cdot 1 = 4$$

Example of Non-monotonic Behavior

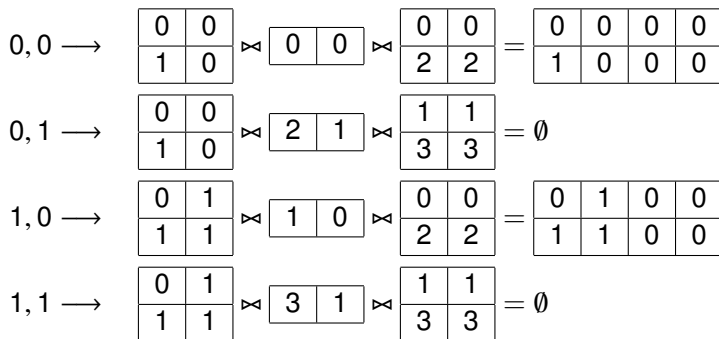
- ▶ Define hash function $\text{hash}(u_i) = i\%2$.

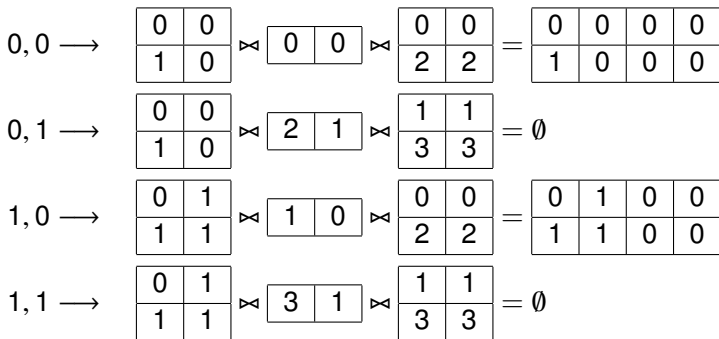
$$\text{hash}(0) = \text{hash}(2) = 0$$

$$\text{hash}(1) = \text{hash}(3) = 1$$

Partitioned Relations

$\text{hash}(y), \text{hash}(z) = \dots$





$$\sum_{i,j \in \{0,1\}} \min \begin{cases} c_{R(i)} \cdot d_{S(i,j)}^y \cdot d_{T(i)}^z \\ d_{R(i)}^y \cdot c_{S(i,j)} \cdot d_{T(i)}^z \\ d_{R(i)}^y \cdot d_{S(i,j)}^z \cdot c_{T(i)} \\ c_{R(i)} \cdot c_{T(i)} \end{cases} = \sum_{i,j \in \{0,1\}} \min \begin{cases} 2 \cdot 1 \cdot 1 \\ 2 \cdot 1 \cdot 1 \\ 2 \cdot 1 \cdot 2 \\ 2 \cdot 2 \end{cases} = \sum_{i,j \in \{0,1\}} 2 = 8$$

$$\begin{array}{l}
0,0 \rightarrow \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 1 & 0 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 0 & 0 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 2 & 2 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 \\ \hline \end{array} \\
0,1 \rightarrow \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 1 & 0 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 2 & 1 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 3 & 3 \\ \hline \end{array} = \emptyset \\
1,0 \rightarrow \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 1 & 1 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 1 & 0 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 2 & 2 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 1 & 0 & 0 \\ \hline 1 & 1 & 0 & 0 \\ \hline \end{array} \\
1,1 \rightarrow \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 1 & 1 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 3 & 1 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 3 & 3 \\ \hline \end{array} = \emptyset
\end{array}$$

$$\sum_{i,j \in \{0,1\}} \min \begin{cases} c_{R(i)} \cdot d_{S(i,j)}^y \cdot d_{T(j)}^z \\ d_{R(i)}^y \cdot c_{S(i,j)} \cdot d_{T(j)}^z \\ d_{R(i)}^y \cdot d_{S(i,j)}^z \cdot c_{T(j)} \\ c_{R(i)} \cdot c_{T(j)} \end{cases} = \sum_{i,j \in \{0,1\}} \min \begin{cases} 2 \cdot 1 \cdot 1 \\ 2 \cdot 1 \cdot 1 \\ 2 \cdot 1 \cdot 2 \\ 2 \cdot 2 \end{cases} = \sum_{i,j \in \{0,1\}} 2 = 8$$

$$\begin{array}{l}
0,0 \rightarrow \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 1 & 0 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 0 & 0 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 2 & 2 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 \\ \hline \end{array} \\
0,1 \rightarrow \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 1 & 0 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 2 & 1 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 3 & 3 \\ \hline \end{array} = \emptyset \\
1,0 \rightarrow \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 1 & 1 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 1 & 0 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 2 & 2 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 1 & 0 & 0 \\ \hline 1 & 1 & 0 & 0 \\ \hline \end{array} \\
1,1 \rightarrow \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 1 & 1 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 3 & 1 \\ \hline \end{array} \bowtie \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 3 & 3 \\ \hline \end{array} = \emptyset
\end{array}$$

$$\sum_{i,j \in \{0,1\}} \min \begin{cases} c_{R(i)} \cdot d_{S(i,j)}^y \cdot d_{T(j)}^z \\ d_{R(i)}^y \cdot c_{S(i,j)} \cdot d_{T(j)}^z \\ d_{R(i)}^y \cdot d_{S(i,j)}^z \cdot c_{T(j)} \\ c_{R(i)} \cdot c_{T(j)} \end{cases} = \sum_{i,j \in \{0,1\}} \min \begin{cases} 2 \cdot 1 \cdot 1 \\ 2 \cdot 1 \cdot 1 \\ 2 \cdot 1 \cdot 2 \\ 2 \cdot 2 \end{cases} = \sum_{i,j \in \{0,1\}} 2 = 8$$

Exponential Growth

- ▶ Sketch size (number of buckets) exponential in hash size.
 - ▶ Exponent = number of attributes in relation.
- ▶ Number of elements to sum up exponential in hash size.
 - ▶ Exponent = number of attributes in entire query.
- ▶ Non-monotonic bound behavior

Tuning Bucket Allocation

- ▶ Larger hash size \implies more information \implies tighter bounds, right?
 - ▶ Partitioning *unconditionally* covered attributes: yes.
 - ▶ Partitioning *conditionally* covered attributes: no.
 - ▶ Non-monotonic tradeoff space.

Non-Linearity of Degree Statistic

- ▶ Count is linear with respect to disjoint union!

$$\text{count}(A) + \text{count}(B) = \text{count}(A \cup B)$$

- ▶ Degree is not...

$$\text{degree}(A) + \text{degree}(B) \geq \text{degree}(A \cup B)$$

$Q(x, y, z, w) :- pseudo(x, y), cast(y, z), mc(z, w), cn(w)$

pseudo

cast

movie_companies

company_name



pseudo



cast



mc



cn

$Q(x, y, z, w) :- pseudo(x, y), cast(y, z), mc(z, w), cn(w)$

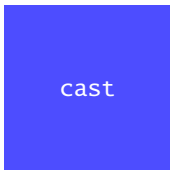
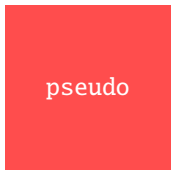
$$C_{pseudo} \cdot d_{cast}^y \cdot d_{mc}^z$$

pseudo

cast

movie_companies

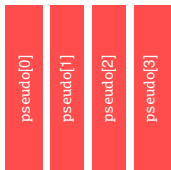
company_name



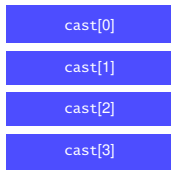
$Q(x, y, z, w) := \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$

$$C_{\text{pseudo}} \cdot d_{\text{cast}}^y \cdot d_{\text{mc}}^z$$

pseudo



cast



movie_companies



company_name



$Q(x, y, z, w) := \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$

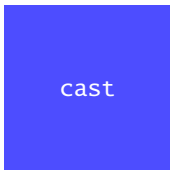
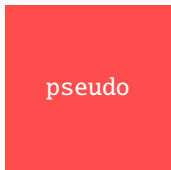
$$d_{\text{pseudo}}^y \cdot C_{\text{cast}} \cdot d_{\text{mc}}^z$$

pseudo

cast

movie_companies

company_name



$Q(x, y, z, w) := \text{pseudo}(x, y), \text{cast}(y, z), \text{mc}(z, w), \text{cn}(w)$

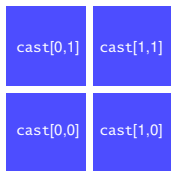
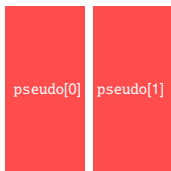
$$d_{\text{pseudo}}^y \cdot C_{\text{cast}} \cdot d_{\text{mc}}^z$$

pseudo

cast

movie_companies

company_name



Reformulated Bound Formula

▶ Old:

$$|Q(D)| \leq \sum_{J \in \text{partition indexes}} \left(\min_{b \in \text{bounding formulas}} b(Q(D[J])) \right)$$

▶ New:

$$|Q(D)| \leq \min_{b \in \text{bounding formulas}} \left(\sum_{\substack{J \in \\ \text{partition indexes} \\ \text{w.r.t. } b}} b(Q(D[J])) \right)$$