

# On Functional Aggregate Queries with Additive Inequalities

[fdbresearch.github.io](https://fdbresearch.github.io)

[relational.ai](https://relational.ai)



**relationalAI**  
AI for the enterprise

**Maximilian Schleich**

University of Oxford

**Mahmoud Abo Khamis**, [relationalAI](https://relational.ai)

**Ryan R. Curtin**, [relationalAI](https://relational.ai)

**Benjamin Moseley**, CMU

**Hung Q. Ngo**, [relationalAI](https://relational.ai)

**XuanLong Nguyen**, University of Michigan

**Dan Olteanu**, University of Oxford

**Carnegie Mellon University**

ACM PODS

June, 2019

## Example I

**Given:** Relations  $R, S$  of size  $O(N)$ .

**Task:** Count number of tuples that satisfy

$$R(a, b) \wedge S(b, c) \wedge a \leq c$$

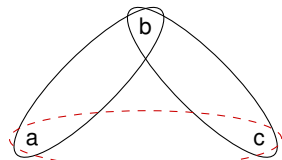
■ Existing approaches take  $O(N^2)$  time.

1. Join  $R$  and  $S$
2. Count number of tuples that satisfy  $a \leq c$

■ Our approach takes  $O(N \log N)$  time.

<b>R</b>	<b>a</b>	<b>b</b>
	1	1
	1	2
	2	1
	2	2
	3	2
	3	3

<b>S</b>	<b>b</b>	<b>c</b>
	1	1
	1	2
	2	0
	2	2
	2	3
	2	4



Hypergraph  $\mathcal{H}$

$R$  and  $S$  are sorted

## Example I

**Given:** Relations  $R, S$  of size  $O(N)$ .

**Task:** Count number of tuples that satisfy

$$R(a, b) \wedge S(b, c) \wedge a \leq c$$

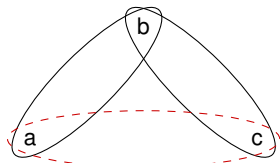
■ Existing approaches take  $O(N^2)$  time.

1. Join  $R$  and  $S$
2. Count number of tuples that satisfy  $a \leq c$

■ Our approach takes  $O(N \log N)$  time.

<b>R</b>	<b>a</b>	<b>b</b>
	1	1
	1	2
	2	1
	2	2
	3	2
	3	3

<b>S</b>	<b>b</b>	<b>c</b>	<b>#</b>
	1	1	2
	1	2	1
	2	0	4
	2	2	3
	2	3	2
	2	4	1



Hypergraph  $\mathcal{H}$

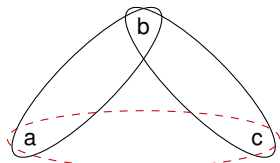
**Step 1:** Pre-aggregate  $S$

## Example I

**Given:** Relations  $R, S$  of size  $O(N)$ .

**Task:** Count number of tuples that satisfy

$$R(a, b) \wedge S(b, c) \wedge a \leq c$$



Hypergraph  $\mathcal{H}$

■ Existing approaches take  $O(N^2)$  time.

1. Join  $R$  and  $S$
2. Count number of tuples that satisfy  $a \leq c$

■ Our approach takes  $O(N \log N)$  time.

R	a	b	S	b	c	#
	1	1		1	1	2
	1	2		1	2	1
	2	1		2	0	4
	2	2		2	2	3
	3	2		2	3	2
	3	3		2	4	1

lookup  $b = 2$

find smallest  $c \geq 1$

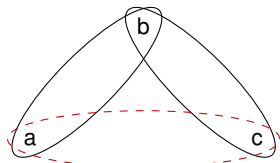
**Step 2:** For each  $R(a, b)$ , locate  $S(b, c)$  with  $c \geq a$

## Example I

**Given:** Relations  $R, S$  of size  $O(N)$ .

**Task:** Count number of tuples that satisfy

$$R(a, b) \wedge S(b, c) \wedge a \leq c$$

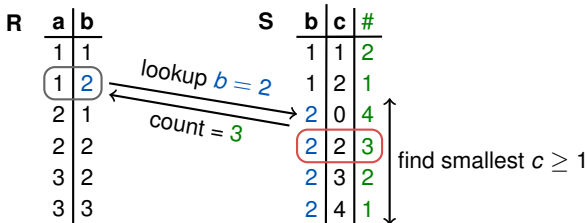


Hypergraph  $\mathcal{H}$

■ Existing approaches take  $O(N^2)$  time.

1. Join  $R$  and  $S$
2. Count number of tuples that satisfy  $a \leq c$

■ Our approach takes  $O(N \log N)$  time.



**Step 4:** Return pre-aggregated count

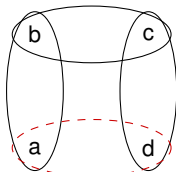
## Example II

**Given:** Relations  $R, S, T$  of size  $O(N)$ .

**Task:** Count number of tuples that satisfy

$$R(a, b) \wedge S(b, c) \wedge T(c, d) \wedge a \leq d$$

- Existing approaches take  $O(N^2)$  time.
- Our approach takes  $O(N^{1.5} \log N)$  time.



Hypergraph  $\mathcal{H}$

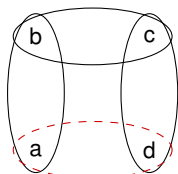
## Example II

**Given:** Relations  $R, S, T$  of size  $O(N)$ .

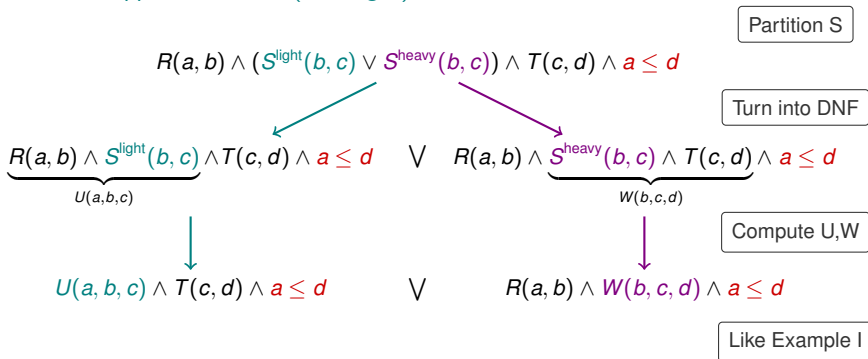
**Task:** Count number of tuples that satisfy

$$R(a, b) \wedge S(b, c) \wedge T(c, d) \wedge a \leq d$$

- Existing approaches take  $O(N^2)$  time.
- Our approach takes  $O(N^{1.5} \log N)$  time.



Hypergraph  $\mathcal{H}$

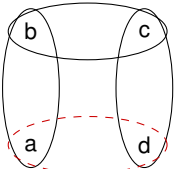


# Example II

**Given:** Relations  $R, S, T$  of size  $O(N)$ .

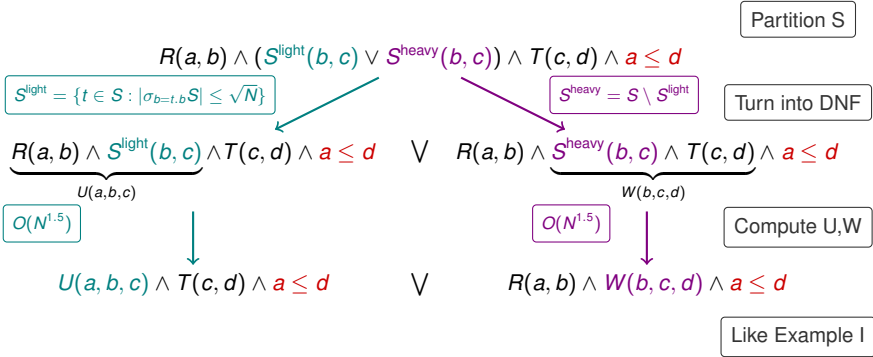
**Task:** Count number of tuples that satisfy

$$R(a, b) \wedge S(b, c) \wedge T(c, d) \wedge a \leq d$$



Hypergraph  $\mathcal{H}$

- Existing approaches take  $O(N^2)$  time.
- Our approach takes  $O(N^{1.5} \log N)$  time.





## Example III: Linear SVM over Databases

**Task:** Compute  $J(\beta)$  over dataset  $\mathbf{D}$  defined by query  $Q$  over database  $\mathcal{I}$ .

$$\begin{aligned} J(\beta) &= \sum_{(\mathbf{x}, y) \in \mathbf{D}} \underbrace{\max\{0, 1 - y \cdot f_{\beta}(\mathbf{x})\}}_{\text{Hinge Loss}} \\ &= \sum_{(\mathbf{x}, y) \in \mathbf{D}} \begin{cases} 1 & \text{if } y \cdot f_{\beta}(\mathbf{x}) < 1 \\ 0 & \text{otherwise} \end{cases} \\ &= \underbrace{\sum_{(\mathbf{x}, y) \in \mathbf{D}} (1 - y \cdot f_{\beta}(\mathbf{x})) \cdot \mathbf{1}_{y \cdot f_{\beta}(\mathbf{x}) \leq 1}}_{\text{Aggregate with inequality}} \end{aligned}$$

Existing approaches:

1. materialize  $\mathbf{D}$
2. learn model in favorite ML tool

Our Approach:

1. avoid materialization of  $\mathbf{D}$
2. learn model using aggregates with inequalities

We can learn the SVM model in time **sublinear** in the size of  $\mathbf{D}$ .

# Relational Machine Learning over Databases

We can express range of models as aggregates with inequalities, including:

- Support Vector Machines (SVM)
- k-Means Clustering
- Robust Regression with Huber Loss
- Boolean Principle Component Analysis (PCA)
- Low Rank Matrix Factorization
- ... and several other models trained with Non-Polynomial Loss functions

# Functional Aggregate Queries with Additive Inequalities

FAQ-Als encode:

1. Relations as factors  $R_K : \prod_{i \in K} \text{Dom}(X_i) \rightarrow \mathbf{D}$
2. Additive Inequality  $E$  as Kronecker delta  $\mathbf{1}_E$
3. Sum-Product Operations over Semiring

Examples for Sum-Product Semiring  $(\mathbb{R}, +, \times)$ :

$$Q() = \sum_{a,b,c} R(a,b) \cdot S(b,c) \cdot \mathbf{1}_{a \leq c} \quad (\text{Example 1})$$

$$Q() = \sum_{a,b,c} R(a,b) \cdot S(b,c) \cdot \mathbf{1}_{a \leq c} \cdot \mathbf{1}_{\frac{b}{2} \leq c} \cdot \mathbf{1}_{a^2 + \frac{b}{2} + 5c \leq 0} \quad (\text{more Als})$$

$$Q(a,b) = \sum_c R(a,b) \cdot S(b,c) \cdot \mathbf{1}_{a \leq c} \cdot \mathbf{1}_{\frac{b}{2} \leq c} \cdot \mathbf{1}_{a^2 + \frac{b}{2} + 5c \leq 0} \quad (\text{with free variables})$$

# Functional Aggregate Queries with Additive Inequalities

FAQ-Als encode:

1. Relations as factors  $R_K : \prod_{i \in K} \text{Dom}(X_i) \rightarrow \mathbf{D}$
2. **Additive Inequality**  $E$  as Kronecker delta  $\mathbf{1}_E$
3. Sum-Product Operations over Semiring

Examples for Sum-Product Semiring  $(\mathbb{R}, +, \times)$ :

$$Q() = \sum_{a,b,c} R(a,b) \cdot S(b,c) \cdot \mathbf{1}_{a \leq c} \quad (\text{Example 1})$$

$$Q() = \sum_{a,b,c} R(a,b) \cdot S(b,c) \cdot \mathbf{1}_{a \leq c} \cdot \mathbf{1}_{\frac{b}{2} \leq c} \cdot \mathbf{1}_{a^2 + \frac{b}{2} + 5c \leq 0} \quad (\text{more Als})$$

$$Q(a,b) = \sum_c R(a,b) \cdot S(b,c) \cdot \mathbf{1}_{a \leq c} \cdot \mathbf{1}_{\frac{b}{2} \leq c} \cdot \mathbf{1}_{a^2 + \frac{b}{2} + 5c \leq 0} \quad (\text{with free variables})$$

Change the semiring to get different aggregates:

$$Q(a,b) = \bigvee_c R(a,b) \wedge S(b,c) \wedge \mathbf{1}_{a \leq c} \wedge \mathbf{1}_{\frac{b}{2} \leq c} \wedge \mathbf{1}_{a^2 + \frac{b}{2} + 5c \leq 0} \quad \text{Boolean: } (\{1, 0\}, \vee, \wedge)$$

$$Q(a,b) = \max_c R(a,b) \cdot S(b,c) \cdot \mathbf{1}_{a \leq c} \cdot \mathbf{1}_{\frac{b}{2} \leq c} \cdot \mathbf{1}_{a^2 + \frac{b}{2} + 5c \leq 0} \quad \text{Max-Product: } (\mathbb{R}, \max, \times)$$

$$Q(a,b) = \bigoplus_c R(a,b) \otimes S(b,c) \otimes \mathbf{1}_{a \leq c} \otimes \mathbf{1}_{\frac{b}{2} \leq c} \otimes \mathbf{1}_{a^2 + \frac{b}{2} + 5c \leq 0} \quad \text{Arbitrary: } (\mathbf{D}, \oplus, \otimes)$$

# Functional Aggregate Queries with Additive Inequalities

$$Q(\mathbf{x}_F) = \bigoplus_{\mathbf{x}_{\mathcal{V} \setminus F}} \left( \bigotimes_{K \in \mathcal{E}_s} R_K(\mathbf{x}_K) \right) \otimes \left( \bigotimes_{K \in \mathcal{E}_\ell} \mathbf{1}_{\sum_{v \in K} \theta_v^K(x_v) \leq 0} \right)$$

Query Hypergraph  $\mathcal{H} = (\mathcal{V}, \mathcal{E} = \mathcal{E}_s \cup \mathcal{E}_\ell)$

- Set of variables  $\mathcal{V} = \{X_1, \dots, X_n\}$
- Set of “**skeleton**” hyperedges  $\mathcal{E}_s$ 
  - ▶ Each hyperedge is defined by a **factor**  $R_K(\mathbf{x}_K)$
- Set of “**ligament**” hyperedges  $\mathcal{E}_\ell$ 
  - ▶ Each hyperedge is defined by **sum of univariate functions**

$X_i$  - variable

$x_i$  - value in  $Dom(X_i)$

$\mathbf{x}_K$  - tuple of values in  $\prod_{i \in K} Dom(X_i)$

For this talk:

Queries without free variables.

The paper:

Generalizes all results via notion of FAQ-width.

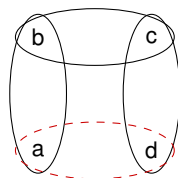
# Generalized Hypertree Decompositions

Generalized Hypertree Decomposition (TD) for  $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ :

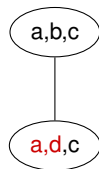
- Tree  $T = (V(T), E(T))$
- Bag  $\chi(t) \subseteq \mathcal{V}$  for each tree-node  $t \in V(T)$

A TD must satisfy:

1. Running intersection property
2. Containment property
  - ▶ every hyperedge in  $\mathcal{E}$  is covered by some bag  $\chi(t)$ .

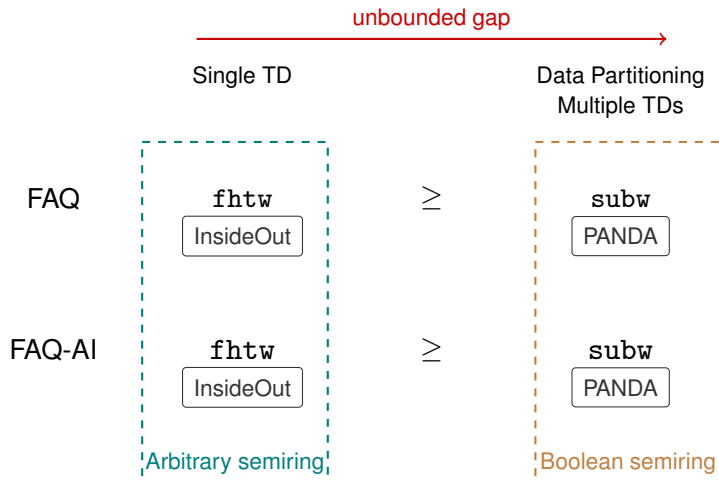


Hypergraph  $\mathcal{H}$



TD for  $\mathcal{H}$

# State-of-the-art for evaluating FAQs and FAQ-AIs



We can do better for FAQ-AI! Using Relaxed Tree Decompositions.



# Relaxed Tree Decompositions

Containment for Tree Decompositions:

1. every hyperedge is covered by some bag

Containment for Relaxed Tree Decompositions:

1. every 'skeleton' hyperedge is covered by some bag
2. every 'ligament' hyperedge is covered by two adjacent bags

# Relaxed Tree Decompositions

Containment for Tree Decompositions:

1. every hyperedge is covered by some bag

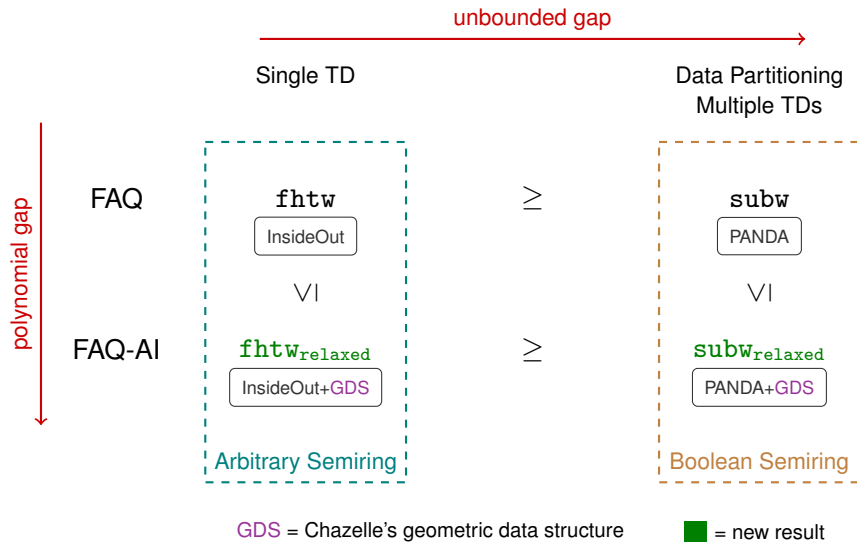
Containment for Relaxed Tree Decompositions:

1. every 'skeleton' hyperedge is covered by some bag
2. every 'ligament' hyperedge is covered by two adjacent bags

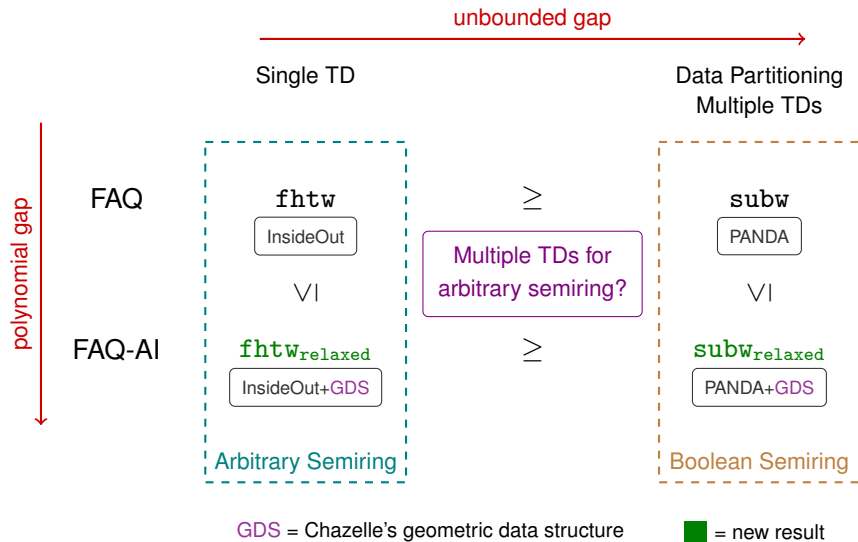
Evaluation of ligament hyperedges over two adjacent bags based on

Chazelle's geometric data structure (GDS)

# Width Measures for FAQs and FAQ-AIs



# Width Measures for FAQs and FAQ-AIs



# #PANDA: A PANDA Variant for Arbitrary Semirings

PANDA decomposes the query into several sub-queries

- Based on information theoretic inequalities

Each sub-query is:

- Computed over partitions of factors corresponding to skeleton hyperedges
- Defined by a different tree decomposition

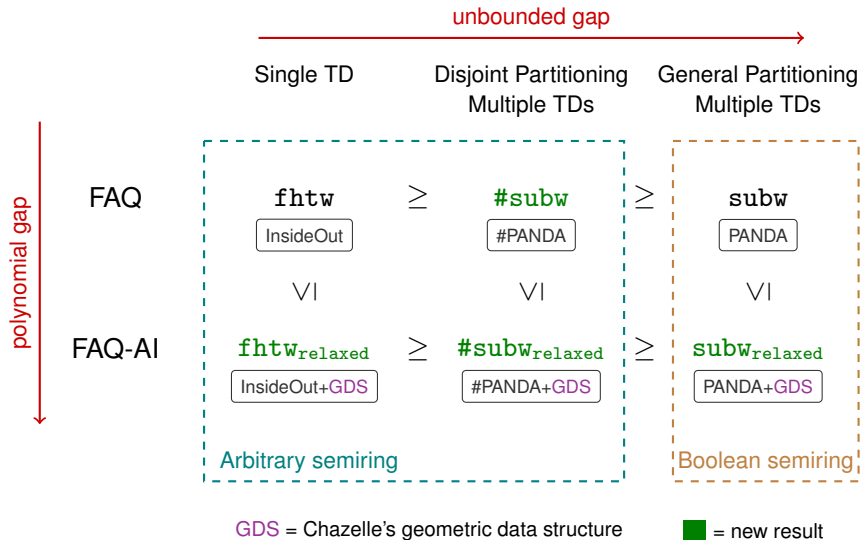
**Challenge:** The results of sub-queries may overlap

- Boolean semiring: OK (✓)
- Arbitrary semiring: Not OK (✗)

#PANDA ensures that the results of sub-queries are disjoint

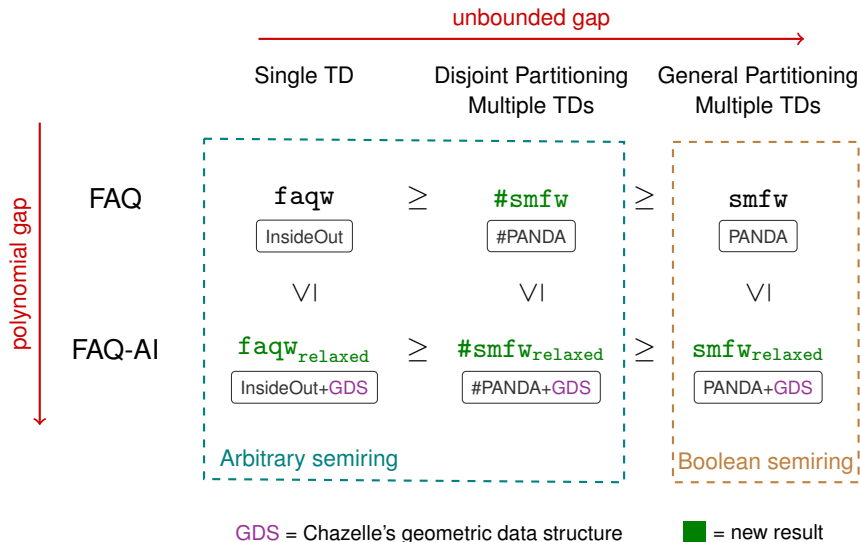
- Allows for aggregate computation with Arbitrary Semirings
- Recall Example II

# Width Measures for FAQs and FAQ-AIs



# Appendix

# Width Measures for FAQs and FAQ-AIs + free variables





# What are we hiding?

