

FAQ-AI: Functional Aggregate Queries with Additive Inequalities

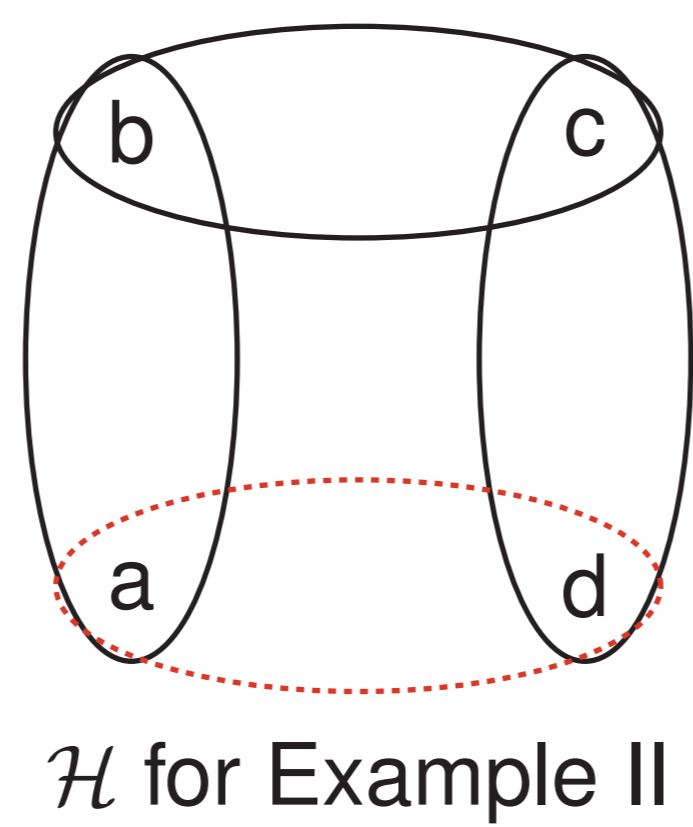
$$Q(\mathbf{x}_F) = \bigoplus_{\mathbf{x}_{V \setminus F}} \left(\bigotimes_{K \in \mathcal{E}_s} R_K(\mathbf{x}_K) \right) \otimes \left(\bigotimes_{K \in \mathcal{E}_l} \mathbf{1}_{\sum_{v \in K} \theta_v^K(x_v) \leq 0} \right)$$

$Q(a, b) = \sum_c R(a, b) \cdot S(b, c) \cdot T(c, d) \cdot \mathbf{1}_{a \leq d} \cdot \mathbf{1}_{\frac{b}{2} \leq c} \cdot \mathbf{1}_{a^2 + \frac{b}{2} \leq 5c}$ Sum-Product: $(\mathbb{R}, +, \times)$
 $Q(a, b) = \bigvee_c R(a, b) \wedge S(b, c) \wedge T(c, d) \wedge \mathbf{1}_{a \leq d} \wedge \mathbf{1}_{\frac{b}{2} \leq c} \wedge \mathbf{1}_{a^2 + \frac{b}{2} \leq 5c}$ Boolean: $(\{1, 0\}, \vee, \wedge)$
 $Q(a, b) = \bigotimes_c R(a, b) \oplus S(b, c) \oplus T(c, d) \oplus \mathbf{1}_{a \leq d} \oplus \mathbf{1}_{\frac{b}{2} \leq c} \oplus \mathbf{1}_{a^2 + \frac{b}{2} \leq 5c}$ Arbitrary: $(\mathbf{D}, \oplus, \otimes)$

Query Hypergraph for FAQ-AIs

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

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



\mathcal{H} for Example II

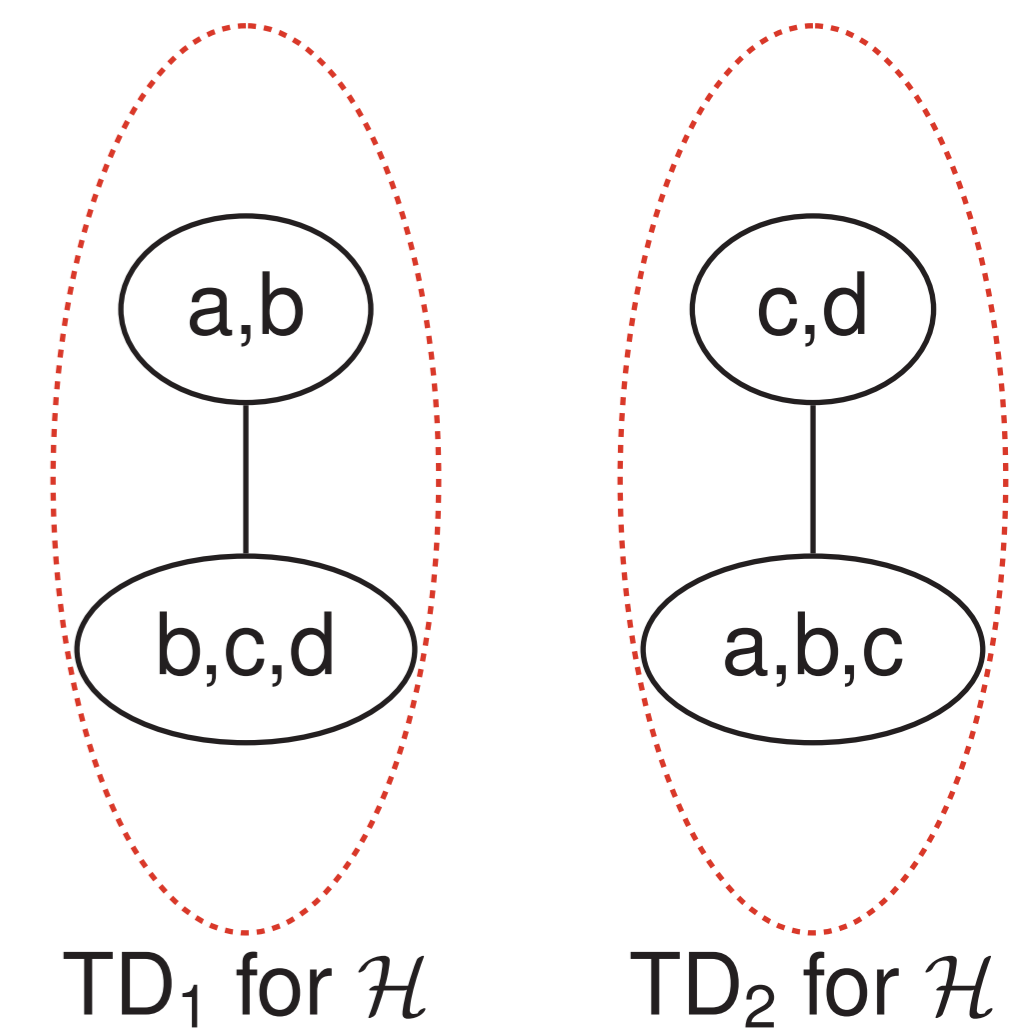
Relaxed Tree Decompositions

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

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

A TD satisfies:

- Running intersection property
- Containment property



TD₁ for \mathcal{H}

TD₂ for \mathcal{H}

Containment for Tree Decompositions:

- every hyperedge is covered by some bag

Containment for Relaxed Tree Decompositions:

- every skeleton hyperedge is covered by some bag
- every ligament hyperedge is covered by two adjacent bags

Example I

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

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

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

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

R	a	b	S	b	c	#
1	1	1	1	1	2	
1	2	1	1	2	1	
2	1	2	2	0	4	
2	2	2	2	2	3	find smallest $c \geq a$
3	2	2	2	3	2	
3	3	3	2	4	1	

lookup $b=2$
count = 3

- Step 1: Pre-count S for each b
 Step 2: Foreach $R(a, b)$, locate first $S(b, c)$ with $a \leq c$
 Step 3: Return pre-aggregated count

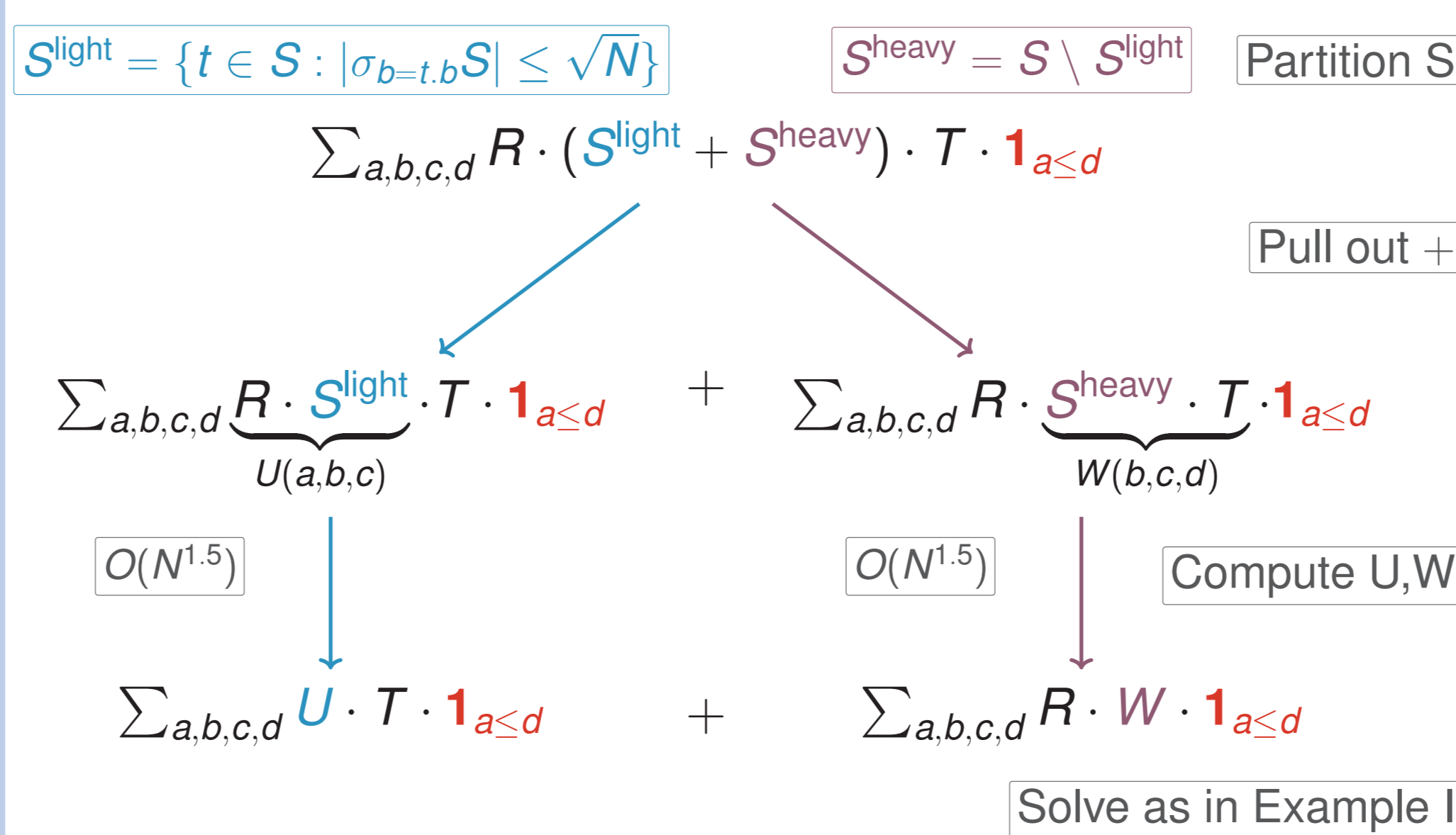
Example II

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

Task: $Q() = \sum_{a,b,c,d} R(a, b) \cdot S(b, c) \cdot T(c, d) \cdot \mathbf{1}_{a \leq d}$

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

Our approach takes $O(N^{1.5} \log N)$ time



Example III: Learning SVM over Databases

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

$$J(\beta) = \sum_{(x,y) \in \mathbf{D}} \max\{0, 1 - y \cdot f_\beta(x)\} = \sum_{(x,y) \in \mathbf{D}} (1 - y \cdot f_\beta(x)) \cdot \mathbf{1}_{y \cdot f_\beta(x) \leq 1}$$

Hinge Loss
FAQ-AI

Existing approaches:

- materialize \mathbf{D}
- learn model using favorite ML tool

Our approach:

- avoid materialization of \mathbf{D}
- express learning using FAQ-AIs

SVM models can be learned in time **sublinear** in $|\mathbf{D}|$.

Further ML models that can benefit from FAQ-AI:

- k-Means Clustering
- Robust Regression with Huber Loss
- Boolean Principle Component Analysis
- Other models trained with non-polynomial loss

Width Measures for FAQs and FAQ-AIs without Free Variables

polynomial gap ↓		unbounded gap →			
		Single TD	Disjoint Partitioning Multiple TDs	General Partitioning Multiple TDs	
FAQ	fhtw InsideOut	≥	#subw #PANDA	≥	subw PANDA
FAQ-AI	fhtw _{relaxed} InsideOut+GDS	≥	#subw _{relaxed} #PANDA+GDS	≥	subw _{relaxed} PANDA+GDS
			Arbitrary semiring		Boolean semiring
			fhtw = fractional hypertree width		subw = submodular width
			GDS = Chazelle's geometric data structure		■ = new result

Width Measures for FAQs and FAQ-AIs with Free Variables

polynomial gap ↓		unbounded gap →			
		Single TD	Disjoint Partitioning Multiple TDs	General Partitioning Multiple TDs	
FAQ	faqw InsideOut	≥	#smfw #PANDA	≥	smfw PANDA
FAQ-AI	faqw _{relaxed} InsideOut+GDS	≥	#smfw _{relaxed} #PANDA+GDS	≥	smfw _{relaxed} PANDA+GDS
			Arbitrary semiring		Boolean semiring
			faqw = FAQ width		smfw = submodular FAQ width
			GDS = Chazelle's geometric data structure		■ = new result