### The Fine Grained Data Complexity of Conjunctive Query Evaluation

#### Dan Suciu

University of Washington

#### **Problem Definition**

Boolean conjunctive query  $Q() = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$ 

Given input database D, what is the time complexity to compute Q(D)?

Time expressed in terms of statistics on **D**.

#### **Problem Definition**

Boolean conjunctive query  $Q() = R_1(\boldsymbol{U}_1) \land \cdots \land R_m(\boldsymbol{U}_m)$ 

Given input database D, what is the time complexity to compute Q(D)?

Time expressed in terms of statistics on **D**.

In this talk: Combinatorial Algorithms: time  $O(|\mathbf{D}|^{\text{subw}(Q)})$  [Khamis et al., 2016a, Khamis et al., 2024c]. Some lower bound techniques [Fan et al., 2023].

Not in this talk: Using Fast Matrix Multiplication: time  $O(|D|^{\omega-\text{subw}(Q)})$ 

[Khamis et al., 2024a].

Assume the algorithm for Q(D) computes several intermediate results.

Assume the algorithm for Q(D) computes several intermediate results.

- If these intermediate outputs have size B, then the runtime is  $\Omega(B)$
- Thus, any algorithm with runtime O(B) is a proof that outputs  $\leq B$ .

Assume the algorithm for Q(D) computes several intermediate results.

- If these intermediate outputs have size B, then the runtime is  $\Omega(B)$
- Thus, any algorithm with runtime O(B) is a proof that outputs  $\leq B$ .

"Proofs to algorithms" paradigm is the converse:

(1) Prove: if **D** satisfies the given stats, then all outputs have size  $\leq B$ . (2) Convert this proof into an algorithm for  $Q(\mathbf{D})$ , with runtime  $\tilde{O}(B)$ 

3/34

Assume the algorithm for Q(D) computes several intermediate results.

- If these intermediate outputs have size B, then the runtime is  $\Omega(B)$
- Thus, any algorithm with runtime O(B) is a proof that outputs  $\leq B$ .

"Proofs to algorithms" paradigm is the converse:

(1) Prove: if **D** satisfies the given stats, then all outputs have size  $\leq B$ . (2) Convert this proof into an algorithm for  $Q(\mathbf{D})$ , with runtime  $\tilde{O}(B)$ 

To prove bounds on the output size: Information Inequalities

Overview 00	Information Inequalities	Output Bounds 000000	Proofs to Algorithms	PANDA 0000000	Lower Bounds 000000	

# Information Inequalities

Definition ([Shannon, 1948])

The entropy of a r.v. X with domain D is  $h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$ 

5/34

Definition ([Shannon, 1948])

The entropy of a r.v. X with domain D is  $h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$ 

Definition ([Shannon, 1948])

The entropy of a r.v. X with domain D is  $h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$ 

X	Y
а	р
a	q
b	q
a	m

Definition ([Shannon, 1948])

The entropy of a r.v. X with domain D is  $h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$ 



Definition ([Shannon, 1948])

The entropy of a r.v. X with domain D is  $h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$ 

X a a	Y p q	p 1/4 1/4	X a b	p   3/4   1/4	Y P q	p 1/4 2/4	Ø	р 1
b a h(X	$\begin{pmatrix} q \\ m \\ (Y) =$	1/4   1/4 = log 4	h(>	<) ≤ log 2	h( Y	$1/4$ () $\leq \log 3$	$h(\emptyset)$	) = <mark>0</mark>

Definition ([Shannon, 1948])

The entropy of a r.v. X with domain D is  $h(X) \stackrel{\text{def}}{=} -\sum_{x \in D} p(x) \log p(x)$ 

The entropic vector of *n* r.v.  $X_1, \ldots, X_n$  is  $\boldsymbol{h} = (h(X_\alpha))_{\alpha \subseteq [n]} \in \mathbb{R}^{2^n}_+$ .



 $\Gamma_n^* =$  set of entropic vectors (a.k.a. entropic functions)

#### **Basic Shannon Inequalities**

For all sets 
$$U, V \subseteq \{X_1, \dots, X_n\}$$
:  
 $h(\emptyset) = 0$   
 $h(U \cup V) \ge h(U)$  Monotonicity  
 $h(U) + h(V) \ge h(U \cup V) + h(U \cap V)$  Submodularity

### **Basic Shannon Inequalities**

For all sets 
$$U, V \subseteq \{X_1, \dots, X_n\}$$
:  
 $h(\emptyset) = 0$   
 $h(U \cup V) \ge h(U)$  Monotonicity  
 $h(U) + h(V) \ge h(U \cup V) + h(U \cap V)$  Submodularity

Setting  $h(V|U) \stackrel{\text{def}}{=} h(UV) - h(U)$ , the inequalities become:

$$h(V|U) \ge 0$$
  $h(V|U) \ge h(V|UW)$ 

### Basic Shannon Inequalities

For all sets 
$$U, V \subseteq \{X_1, \dots, X_n\}$$
:  
 $h(\emptyset) = 0$   
 $h(U \cup V) \ge h(U)$  Monotonicity  
 $h(U) + h(V) \ge h(U \cup V) + h(U \cap V)$  Submodularity

Setting  $h(V|U) \stackrel{\text{def}}{=} h(UV) - h(U)$ , the inequalities become:

$$h(V|U) \ge 0 \qquad \qquad h(V|U) \ge h(V|UW)$$

A Shannon inequality is a consequence of the basic inequalities.

A Polymatroid is a vector  $h \in \mathbb{R}^{2^n}_+$  that satisfies all Shannon inequalities

$$\Gamma_n = \text{set of polymatroids}$$
  $\Gamma_n^* \subseteq \Gamma_n$ 

Proofs to Algorit

s PANDA 0000000 Lower Bounds

Summary 0

#### Example: A Shannon Inequality

#### Example (Shearer) $h(XY) + h(YZ) + h(XZ) \ge 2h(XYZ)$

E

#### Example: A Shannon Inequality

xample (Shearer)  
$$h(XY) + h(YZ) + h(XZ) \ge 2h(XYZ)$$

h(XY) + h(YZ) + h(XZ)

Proofs to Algorithn 000 PANDA 0000000 Lower Bounds

Summary

#### Example: A Shannon Inequality

Example (Shearer)  
$$h(XY) + h(YZ) + h(XZ) \ge 2h(XYZ)$$

 $\underline{h(XY) + h(YZ)} + h(XZ)$ 

7/34

Proofs to Algorit

PANDA

Lower Bounds

Summary

Example (Shearer)  
$$h(XY) + h(YZ) + h(XZ) \ge 2h(XYZ)$$

$$\frac{h(XY) + h(YZ)}{\geq} h(XYZ) + h(Y) + h(XZ)$$

E

roofs to Algorithm

PANDA 0000000 Lower Bounds 000000 Summary

xample (Shearer)  
$$h(XY) + h(YZ) + h(XZ) \ge 2h(XYZ)$$

$$\frac{h(XY) + h(YZ)}{\geq} h(XYZ) + \frac{h(Y) + h(XZ)}{\leq}$$

Example (Shearer)  
$$h(XY) + h(YZ) + h(XZ) \ge 2h(XYZ)$$

$$\frac{h(XY) + h(YZ) + h(XZ)}{\geq h(XYZ) + \underline{h(Y) + h(XZ)}}$$
$$\geq 2h(XYZ) + h(\emptyset)$$

Example (Shearer)  
$$h(XY) + h(YZ) + h(XZ) \ge 2h(XYZ)$$

$$\frac{h(XY) + h(YZ)}{\geq h(XYZ) + h(Y) + h(XZ)}$$
$$\geq 2h(XYZ) + h(\emptyset)$$
$$= 2h(XYZ)$$

Information Inequalities

00000

- Shannon inequality: one satisfied by all polymatroids Γ<sub>n</sub> Information inequality:<sup>1</sup> one satisfied by all entropic vectors Γ<sub>n</sub><sup>\*</sup>
- Pippenger [Pippenger, 1986]: Are Shannon inequalities complete? Breakthrough [Zhang and Yeung, 1998]: NO for  $n \ge 4$ . It follows  $\Gamma_n^* \subsetneq \Gamma_n$ .
- $\Gamma_n^*$  is not convex;  $\overline{\Gamma}_n^*$  is convex [Yeung, 2008], but is not a polytope [Matús, 2007].
- The characterization of  $\overline{\Gamma}_n^*$  is open to date.

<sup>&</sup>lt;sup>1</sup>A.k.a. entropic inequality

Overview 00	Information Inequalities	Output Bounds ●00000	Proofs to Algorithms	PANDA 0000000	Lower Bounds 000000	

## Output Size Bounds

Proofs to Algorithm

PANDA Lowe

Upper Bounds on the Output Size

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$ 

Given statistics Stats, compute  $\boldsymbol{B} \stackrel{\text{def}}{=} \max_{\boldsymbol{D}:\boldsymbol{D} \models \boldsymbol{Stats}} Q(|\boldsymbol{D}|).$ 

Proofs to Algorithm

PANDA Lowe

Upper Bounds on the Output Size

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$ 

Given statistics Stats, compute  $B \stackrel{\text{def}}{=} \max_{D:D \models Stats} Q(|D|)$ .

• Cardinality constraints <sup>2</sup>  $|R| \leq N$ :



<sup>&</sup>lt;sup>2</sup>Assume R is ternary, R(XYZ).

Proofs to Algorithm

PANDA Lowe

Upper Bounds on the Output Size

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$ 

Given statistics Stats, compute  $B \stackrel{\text{def}}{=} \max_{D:D \models Stats} Q(|D|)$ .

- Cardinality constraints <sup>2</sup>  $|R| \leq N$ :
- Functional dependency  $X \to Y$ :
- Degree constraint  $||\deg_R(YZ|X)||_{\infty} \leq P$ :



h(Y|X)=0

$$h(YZ|X) \leq \log P$$

<sup>&</sup>lt;sup>2</sup>Assume *R* is ternary, R(XYZ).

Upper Bounds on the Output Size

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$ 

Given statistics Stats, compute  $B \stackrel{\text{def}}{=} \max_{D:D \models Stats} Q(|D|)$ .

- Cardinality constraints <sup>2</sup>  $|R| \leq N$ :
- Functional dependency  $X \to Y$ :
- Degree constraint  $|| \deg_R(YZ|X) ||_{\infty} \leq P$ :
- $\ell_p$ -constraint  $|| \deg_R(YZ|X) ||_p \leq S$ :



h(Y|X)=0

 $h(YZ|X) \leq \log P$ 

$$\left| \frac{1}{p} h(X) + h(YZ|X) \le \log S \right|$$

<sup>&</sup>lt;sup>2</sup>Assume R is ternary, R(XYZ).

Overview 00	Information Inequalities	Output Bounds	Proofs to Algorithms 000	PANDA 0000000	Lower Bounds 000000	Summary O		
Exam	ple							
$Full \ CQ : \ Q_{\Delta}(XYZ) = R(XY) \land S(YZ) \land T(ZX)$								
			R	$\leq N_1,  S $	$  \leq N_2,  T $	$\leq N_3$ .		

Overview 00	Information Inequalities	Output Bounds	Proofs to Algorithms 000	PANDA 0000000	Lower Bounds 000000			
Exam	ple							
Full CQ: $Q_{\Lambda}(XYZ) = R(XY) \land S(YZ) \land T(ZX)$								
	· · · · · · · · · · · · · · · · · · ·			≤ <b>№</b> 1,   <i>S</i>	$ \leq N_2,  T $	$\leq N_3$ .		

Example ([Chung et al., 1986] [Grohe and Marx, 2014] [Atserias et al., 2013])

 $|Q_{\Delta}(\boldsymbol{D})| \leq (N_1 \cdot N_2 \cdot N_3)^{\frac{1}{2}}$ 

Overview 00	Information Inequalities	Output Bounds	Proofs to Algorithms 000	PANDA 0000000	Lower Bounds 000000	Summary O
Exam	ple					
	(YYZ) =		$(\mathbf{V}\mathbf{Z}) \wedge \mathbf{T}(\mathbf{Z}\mathbf{V})$			
Full C	$Q. Q\Delta(\Lambda TZ) = 1$	n( <i>x</i> i ) // 3(	R	≤ <b>№</b> 1,   <i>S</i>	$ \leq N_2,  T $	≤ <b>№</b> 3.

Example ([Chung et al., 1986] [Grohe and Marx, 2014] [Atserias et al., 2013])

Let  $p: Q_{\Delta}(\boldsymbol{D}) \rightarrow [0,1]$  be uniform and  $\boldsymbol{h}$  its entropy.

Dan Suciu

 $|Q_{\Delta}(\boldsymbol{D})| \leq (\boldsymbol{N}_1 \cdot \boldsymbol{N}_2 \cdot \boldsymbol{N}_3)^{\frac{1}{2}}$ 

 $|R| \leq N_1, |S| \leq N_2, |T| \leq N_3.$ 

Example ([Chung et al., 1986] [Grohe and Marx, 2014] [Atserias et al., 2013])

 $|Q_{\Delta}(\boldsymbol{D})| \leq (\boldsymbol{N}_1 \cdot \boldsymbol{N}_2 \cdot \boldsymbol{N}_3)^{\frac{1}{2}}$ 

Let  $p: Q_{\Delta}(\mathbf{D}) \rightarrow [0, 1]$  be uniform and  $\mathbf{h}$  its entropy.

 $\log N_1 + \log N_2 + \log N_3 \ge h(XY) + h(YZ) + h(XZ)$ 

Overview  
OcInformation InequalitiesOutput Bounds  
ooleoolProofs to Algorithms  
ooleoolPANDA  
cooleoolLower Bounds  
ooleoolSummary  
ooleoolExampleFull CQ:
$$Q_{\Delta}(XYZ) = R(XY) \land S(YZ) \land T(ZX)$$

 $\begin{aligned} \text{Full CQ: } Q_{\Delta}(XYZ) &= R(XY) \land S(YZ) \land T(ZX) \\ & |R| \leq N_1, |S| \leq N_2, |T| \leq N_3. \end{aligned}$ 

Example ([Chung et al., 1986] [Grohe and Marx, 2014] [Atserias et al., 2013])

 $|Q_{\Delta}(\boldsymbol{D})| \leq (N_1 \cdot N_2 \cdot N_3)^{\frac{1}{2}}$ 

Let  $p: Q_{\Delta}(\boldsymbol{D}) \rightarrow [0,1]$  be uniform and  $\boldsymbol{h}$  its entropy.

$$\log N_1 + \log N_2 + \log N_3 \ge h(XY) + h(YZ) + h(XZ)$$
$$\ge 2h(XYZ)$$

Overview  
OccorrectionInformation InequalitiesOutput Bounds  
occorrectionProofs to Algorithms  
occorrectionPANDA  
cocorrectionLower Bounds  
occorrectionSummary  
occorrectionExampleFull CQ:
$$Q_{\Lambda}(XYZ) = R(XY) \land S(YZ) \land T(ZX)$$

 $\begin{aligned} \text{Full CQ: } Q_{\Delta}(XYZ) &= R(XY) \land S(YZ) \land T(ZX) \\ & |R| \leq N_1, |S| \leq N_2, |T| \leq N_3. \end{aligned}$ 

Example ([Chung et al., 1986] [Grohe and Marx, 2014] [Atserias et al., 2013])

 $|Q_{\Delta}(\boldsymbol{D})| \leq (N_1 \cdot N_2 \cdot N_3)^{\frac{1}{2}}$ 

Let  $p: Q_{\Delta}(\boldsymbol{D}) \rightarrow [0,1]$  be uniform and  $\boldsymbol{h}$  its entropy.

$$\log N_1 + \log N_2 + \log N_3 \ge h(XY) + h(YZ) + h(XZ)$$
$$\ge 2h(XYZ) = 2\log |Q(\mathbf{D})|$$

11/34
Overview  
OccoreInformation InequalitiesOutput Bounds  
ooleoooProofs to Algorithms  
occorePANDA  
CococoreLower Bounds  
occoreSummary  
occoreExampleFull CQ:
$$Q_{\Lambda}(XYZ) = R(XY) \land S(YZ) \land T(ZX)$$

 $\begin{aligned} \text{Full CQ: } Q_{\Delta}(XYZ) &= R(XY) \land S(YZ) \land T(ZX) \\ & |R| \leq N_1, |S| \leq N_2, |T| \leq N_3. \end{aligned}$ 

Example ([Chung et al., 1986] [Grohe and Marx, 2014] [Atserias et al., 2013])

 $|Q_{\Delta}(\boldsymbol{D})| \leq (N_1 \cdot N_2 \cdot N_3)^{\frac{1}{2}}$ 

Let  $p: Q_{\Delta}(\boldsymbol{D}) \rightarrow [0,1]$  be uniform and  $\boldsymbol{h}$  its entropy.

$$\log N_1 + \log N_2 + \log N_3 \ge h(XY) + h(YZ) + h(XZ)$$
$$\ge 2h(XYZ) = 2\log |Q(\mathbf{D})|$$

Other inequalities exists, e.g.  $|Q_{\Delta}| \leq N_1 \cdot N_2$ .

## The Polymatroid Bound

Full CQ  $Q(\mathbf{X})$ , statistics Stats.

The polymatroid bound is:  $b(Q, \text{Stats}) \stackrel{\text{def}}{=} \max_{h \in \Gamma_n, h \models \text{stats}} h(X)$ 

Fact: if  $\boldsymbol{D} \models Stats$ , then  $|Q(\boldsymbol{D})| \leq 2^{b(Q,Stats)}$ .

Equal to the best bound one can derive using inequalities (by strong duality).

### Example Revisited

$$Q_{\Delta}(XYZ) = R(XY) \land S(YZ) \land T(ZX)$$

b(Q, Stats) can be computed in one of two ways

13/34

### Example Revisited

$$Q_{\Delta}(XYZ) = R(XY) \land S(YZ) \land T(ZX)$$

b(Q, Stats) can be computed in one of two ways

```
Primal program:

Maximize h(XYZ) where

\begin{cases}
h(XY) \le \log |N_1| \\
h(YZ) \le \log |N_2| \\
h(XZ) \le \log |N_3| \\
h \in \Gamma_3
\end{cases}
```

13/34

Example Revisited

$$Q_{\Delta}(XYZ) = R(XY) \land S(YZ) \land T(ZX)$$

b(Q, Stats) can be computed in one of two ways

```
Primal program:

Maximize h(XYZ) where

\begin{cases}
h(XY) \le \log |N_1| \\
h(YZ) \le \log |N_2| \\
h(XZ) \le \log |N_3| \\
h \in \Gamma_3
\end{cases}
```

Dual program: Minimize  $w_1 \log |N_1| + w_2 \log |N_2| + w_3 \log |N_3|$ where  $w_1, w_2, w_3 \ge 0$  s.t.  $\forall \mathbf{h} \in \Gamma_3$ :  $w_1 h(XY) + w_2 h(YZ) + w_3 h(XZ) \ge h(XYZ)$ 

# Short History of Output Size Bounds

AGM bound [Atserias et al., 2013]

Add FDs [Gottlob et al., 2012, Khamis et al., 2016b, Gogacz and Torunczyk, 2017, Ngo, 2022]

Add max degrees [Khamis et al., 2017].

Add degree sequences [Deeds et al., 2023a, Deeds et al., 2023b, Khamis et al., 2024b, Zhang et al., 2025]

PANDA Low

# Short History of Output Size Bounds

AGM bound [Atserias et al., 2013]

### $|Q_{\Delta}| \leq (|R| \cdot |S| \cdot |T|)^{1/2}$

Add FDs [Gottlob et al., 2012, Khamis et al., 2016b, Gogacz and Torunczyk, 2017, Ngo, 2022]

Add max degrees [Khamis et al., 2017].

Add degree sequences [Deeds et al., 2023a, Deeds et al., 2023b, Khamis et al., 2024b, Zhang et al., 2025]

PANDA Low

# Short History of Output Size Bounds

AGM bound [Atserias et al., 2013]

 $|Q_{\Delta}| \leq \left(|R| \cdot |S| \cdot |T|\right)^{1/2}$ 

Add FDs [Gottlob et al., 2012, Khamis et al., 2016b, Gogacz and Torunczyk, 2017, Ngo, 2022]

 $|Q_{\Delta}| \leq |R|$  when S.Y is a key

Add max degrees [Khamis et al., 2017].

 $|Q_{\Delta}| \leq |R| \cdot ||\deg_{S}(Z|Y)||_{\infty}$ 

Add degree sequences [Deeds et al., 2023a, Deeds et al., 2023b, Khamis et al., 2024b, Zhang et al., 2025]

# Short History of Output Size Bounds

AGM bound [Atserias et al., 2013]

 $|Q_{\Delta}| \leq \left(|R| \cdot |S| \cdot |T|\right)^{1/2}$ 

Add FDs [Gottlob et al., 2012, Khamis et al., 2016b, Gogacz and Torunczyk, 2017, Ngo, 2022]

 $|Q_{\Delta}| \leq |R|$  when S.Y is a key

Add max degrees [Khamis et al., 2017].

 $|Q_{\Delta}| \leq |R| \cdot ||\deg_{\mathcal{S}}(Z|Y)||_{\infty}$ 

 $\begin{array}{l} \text{Add degree sequences }_{[\text{Deeds et al., 2023a, Deeds et al., 2023b, Khamis et al., 2024b, Zhang et al., 2025]} \\ |Q_{\Delta}| \leq \left( ||\text{deg}_{\mathcal{R}}(\mathcal{Y}|\mathcal{X})||_{2}^{2} \cdot ||\text{deg}_{\mathcal{S}}(\mathcal{Z}|\mathcal{Y})||_{2}^{2} \cdot ||\text{deg}_{\mathcal{T}}(\mathcal{X}|\mathcal{Z})||_{2}^{2} \right)^{1/3} \end{array}$ 

# Short History of Output Size Bounds

AGM bound [Atserias et al., 2013]

 $|Q_{\Delta}| \leq (|R| \cdot |S| \cdot |T|)^{1/2}$ 

Add FDs [Gottlob et al., 2012, Khamis et al., 2016b, Gogacz and Torunczyk, 2017, Ngo, 2022]

 $|Q_{\Delta}| \leq |R|$  when S.Y is a key

Add max degrees [Khamis et al., 2017].

 $|Q_{\Delta}| \leq |R| \cdot ||\deg_{\mathcal{S}}(Z|Y)||_{\infty}$ 

 $\begin{array}{l} \text{Add degree sequences } _{\text{[Deeds et al., 2023a, Deeds et al., 2023b, Khamis et al., 2024b, Zhang et al., 2025]} \\ |Q_{\Delta}| \leq \left( ||\text{deg}_{R}(Y|X)||_{2}^{2} \cdot ||\text{deg}_{S}(Z|Y)||_{2}^{2} \cdot ||\text{deg}_{T}(X|Z)||_{2}^{2} \right)^{1/3} \\ |Q_{\Delta}| \leq \left( ||\text{deg}_{R}(Y|X)||_{3}^{3} \cdot ||\text{deg}_{S}(Y|Z)||_{3}^{3} \cdot |T|^{5} \right)^{1/6} \end{array}$ 

Overview 00	Information Inequalities	Output Bounds 000000	Proofs to Algorithms ●00	PANDA 0000000	Lower Bounds 000000	

# From Proofs to Algorithms

PANDA Lo

Lower Bounds

Summary O

# Data Partitioning: Example

### Example

 $Q_{\Delta}(XYZ) = R(XY) \wedge S(YZ) \wedge T(ZX),$ 

 $|\mathcal{Q}_{\Delta}| \leq (\mathcal{N}_1 \mathcal{N}_2 \mathcal{N}_3)^{\frac{1}{2}}$ 

PANDA Lov

Lower Bounds

# Data Partitioning: Example

### Example

 $Q_{\Delta}(XYZ) = R(XY) \wedge S(YZ) \wedge T(ZX),$ 

$$|Q_{\Delta}| \leq (N_1 N_2 N_3)^{\frac{1}{2}}$$



PANDA L

Lower Bounds 000000

### Data Partitioning: Example

#### Example

$$Q_{\Delta}(XYZ) = R(XY) \wedge S(YZ) \wedge T(ZX),$$

$$|Q_{\Delta}| \leq (N_1 N_2 N_3)^{\frac{1}{2}}$$





PANDA 0000000

Lower Bounds 000000

## Data Partitioning: Example

#### Example

 $Q_{\Delta}(XYZ) = R(XY) \wedge S(YZ) \wedge T(ZX),$ 





The light/heavy thresholds are given by  $h^*$ :

Light:
$$|| \deg_{S_{\text{light}}}(Z|Y)||_{\infty} \le 2^{h^*(Z|Y)}$$
Heavy: $|S_{\text{heavy}}(Y)| \le 2^{h^*(Y)}$ 

### Data Partitioning

Theorem ([Khamis et al., 2024c])

 $Q(\mathbf{D})$  can be computed in time  $\tilde{O}(2^{b(Q,Stats)})$  by data partitioning & joins.

• Light/heavy with threshold  $h^*(Z|Y)$ :

We can only use this partitioning when inequality is divergent: i.e. a sum of two inequalities each with h(Y) and h(Z|Y) respectively.

### Data Partitioning

Theorem ([Khamis et al., 2024c])

 $Q(\mathbf{D})$  can be computed in time  $\tilde{O}(2^{b(Q,Stats)})$  by data partitioning & joins.

### • Light/heavy with threshold $h^*(Z|Y)$ :

We can only use this partitioning when inequality is divergent: i.e. a sum of two inequalities each with h(Y) and h(Z|Y) respectively.

 Uniformization partition into log N buckets Bucket i holds tuples with deg<sub>R</sub>(Z|y) ∈ [2<sup>i−1</sup>, 2<sup>i</sup>].

Repeated uniformizations lead to poly-log factor.

Overview 00	Information Inequalities	Output Bounds 000000	Proofs to Algorithms 000	PANDA ●ooooooo	Lower Bounds 000000	

#### Proof-Assisted eNtropic Degree-Aware

# The Fractional Tree Width

Boolean CQ  $Q() = \bigwedge_j R(U_j)$ 

Denote a tree decomposition by  $\boldsymbol{T} = (T, \chi)$ 

Fractional tree width

$$\mathsf{ftw}(Q) \stackrel{\mathsf{def}}{=} \mathsf{min}_{\mathcal{T}} \mathsf{max}_{\boldsymbol{h}\models\mathsf{Stats}} \mathsf{max}_{t\in \mathtt{Nodes}(\mathcal{T})} h(\chi(t))$$





 $Q(\mathbf{D})$  can be computed in time  $\tilde{O}(2^{ftw(Q, Stats)})$ 

19/34

# The Submodular Width

$$\mathsf{ftw}(Q, \mathsf{Stats}) \stackrel{\mathsf{def}}{=} \mathsf{min}_{\mathcal{T}} \mathsf{max}_{\mathbf{h} \models \mathsf{Stats}} \mathsf{max}_{t \in \mathsf{Nodes}(\mathcal{T})} h(\chi(t))$$

Marx [Marx, 2013] defines the submodular width using all tree decompositions:

 $\mathsf{subw}(Q, \mathsf{Stats}) \stackrel{\mathsf{def}}{=} \mathsf{max}_{\pmb{h} \models \mathsf{Stats}} \mathsf{min}_{\pmb{\tau}} \mathsf{max}_{t \in \mathtt{Nodes}(\mathcal{T})} h(\chi(t))$ 

# The Submodular Width

$$\mathsf{ftw}(Q, \mathsf{Stats}) \stackrel{\mathsf{def}}{=} \mathsf{min}_{\mathcal{T}} \mathsf{max}_{h \models \mathsf{Stats}} \mathsf{max}_{t \in \mathtt{Nodes}(\mathcal{T})} h(\chi(t))$$

Marx [Marx, 2013] defines the submodular width using all tree decompositions:

 $\mathsf{subw}(Q, \mathsf{Stats}) \stackrel{\mathsf{def}}{=} \mathsf{max}_{h \models \mathsf{Stats}} \mathsf{min}_{\mathcal{T}} \mathsf{max}_{t \in \mathtt{Nodes}(\mathcal{T})} h(\chi(t))$ 

Theorem ([Khamis et al., 2016a, Khamis et al., 2024c])

PANDA computes a Boolean query Q in time  $\tilde{O}(2^{subw(Q, Stats)})$ 

Key technique: disjunctive datalog rules & data partitioning

A Disjunctive Datalog Rule, DDR, is:

$$\Sigma: \left[ Q_1(\boldsymbol{V}_1) \vee \cdots \vee Q_k(\boldsymbol{V}_k) \leftarrow R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m) \right]$$

A Disjunctive Datalog Rule, DDR, is:

$$\Sigma: \left| Q_1(\boldsymbol{V}_1) \vee \cdots \vee Q_k(\boldsymbol{V}_k) \leftarrow R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m) \right|$$

A model is a tuple of relations  $\Sigma = (Q_1, \ldots, Q_k)$  that satisfies the rule.

Its size is  $\|\Sigma\| = \max_i |Q_i|$ .

A Disjunctive Datalog Rule, DDR, is:

$$\Sigma: \left| Q_1(\boldsymbol{V}_1) \vee \cdots \vee Q_k(\boldsymbol{V}_k) \leftarrow R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m) \right|$$

A model is a tuple of relations  $\Sigma = (Q_1, \ldots, Q_k)$  that satisfies the rule.

Its size is  $\|\Sigma\| = \max_i |Q_i|$ .

Examples:

```
Q(XZ) \leftarrow R(XY) \land S(YZ)
```

A Disjunctive Datalog Rule, DDR, is:

$$\Sigma: \ Q_1(\boldsymbol{V}_1) \lor \cdots \lor Q_k(\boldsymbol{V}_k) \leftarrow R_1(\boldsymbol{U}_1) \land \cdots \land R_m(\boldsymbol{U}_m)$$

A model is a tuple of relations  $\Sigma = (Q_1, \ldots, Q_k)$  that satisfies the rule.

Its size is  $\|\Sigma\| = \max_i |Q_i|$ .

Examples:

$$Q(XZ) \leftarrow R(XY) \land S(YZ)$$
  
 $A(X) \lor B(Y) \leftarrow R(XY)$ 

A Disjunctive Datalog Rule, DDR, is:

$$\Sigma: \left[ Q_1(\boldsymbol{V}_1) \vee \cdots \vee Q_k(\boldsymbol{V}_k) \leftarrow R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m) \right]$$

A model is a tuple of relations  $\Sigma = (Q_1, \ldots, Q_k)$  that satisfies the rule.

Its size is  $\|\Sigma\| = \max_i |Q_i|$ .

Examples:

$$egin{aligned} Q(XZ) \leftarrow R(XY) \wedge S(YZ) \ A(X) \lor B(Y) \leftarrow R(XY) \ A(XYZ) \lor B(YZU) \leftarrow R(XY) \wedge S(YZ) \wedge T(ZU) \end{aligned}$$

$$\Sigma: \quad Q_1(\boldsymbol{V}_1) \vee \cdots \vee Q_k(\boldsymbol{V}_k) \leftarrow R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$$

Polymatroid bound 
$$b(\Sigma, Stats) \stackrel{\text{def}}{=} \max_{\boldsymbol{h} \in \Gamma_n, \boldsymbol{h} \models \text{Stats}} \min(h(\boldsymbol{V}_1), \dots, h(\boldsymbol{V}_k))$$

A linear optimization problem which has a matching Shannon inequality.

Theorem ([Khamis et al., 2024c])  

$$\exists model \Sigma \text{ of size} \leq 2^{b(\Sigma, \text{Stats})}; \text{ it can be computed in time } \tilde{O}(2^{b(\Sigma, \text{Stats})})$$

Output Bound

Proofs to Algorithı 000 PANDA 00000●0

Lower Bounds

Summary 0

### PANDA

# Boolean CQ $Q() = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$

Boolean CQ 
$$Q() = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$$

$$\mathsf{subw}(Q,\mathsf{Stats}) = \max_{\boldsymbol{h} \models \mathsf{Stats}} \min_{\boldsymbol{T} \in \{\boldsymbol{T}_1, \dots, \boldsymbol{T}_k\}} \max_{t \in \mathsf{Nodes}(T)} h(\chi(t))$$

Boolean CQ 
$$Q() = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$$

$$\begin{aligned} \mathsf{subw}(Q,\mathsf{Stats}) &= \max_{\boldsymbol{h} \models \mathsf{Stats}} \min_{\boldsymbol{T} \in \{\boldsymbol{T}_1, \dots, \boldsymbol{T}_k\}} \max_{t \in \mathsf{Nodes}(\mathcal{T})} h(\chi(t)) \\ &= \max_{t_1 \in \mathsf{Nodes}(\mathcal{T}_1), \dots, t_k \in \mathsf{Nodes}(\mathcal{T}_k)} \max_{\boldsymbol{h} \models \mathsf{Stats}} (\min(h(\chi_1(t_1)), \dots, h(\chi_k(t_k)))) \end{aligned}$$

Boolean CQ 
$$Q() = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$$

$$subw(Q, Stats) = \max_{\substack{h \models Stats}} \min_{\substack{T \in \{T_1, \dots, T_k\}}} \max_{t \in Nodes(T)} h(\chi(t))$$
$$= \max_{\substack{t_1 \in Nodes(T_1), \dots, t_k \in Nodes(T_k)}} \max_{\substack{h \models Stats}} (\min(h(\chi_1(t_1)), \dots, h(\chi_k(t_k))))$$

Boolean CQ 
$$Q() = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$$

$$subw(Q, Stats) = \max_{\substack{h \models Stats}} \min_{T \in \{T_1, \dots, T_k\}} \max_{t \in Nodes(T)} h(\chi(t))$$
$$= \max_{\substack{t_1 \in Nodes(T_1), \dots, t_k \in Nodes(T_k)}} \max_{\substack{h \models Stats}} (\min(h(\chi_1(t_1)), \dots, h(\chi_k(t_k))))$$
$$\underbrace{b(\Sigma, Stats)}$$

DDR: 
$$\Sigma : B_1(\chi(t_1)) \vee \cdots \vee B_k(\chi(t_k)) \leftarrow R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$$

Boolean CQ 
$$Q() = R_1(\boldsymbol{U}_1) \land \cdots \land R_m(\boldsymbol{U}_m)$$

$$subw(Q, Stats) = \max_{\substack{h \models Stats}} \min_{T \in \{T_1, \dots, T_k\}} \max_{t \in Nodes(T)} h(\chi(t))$$
$$= \max_{\substack{t_1 \in Nodes(T_1), \dots, t_k \in Nodes(T_k)}} \max_{\substack{h \models Stats}} (\min(h(\chi_1(t_1)), \dots, h(\chi_k(t_k))))$$
$$\underbrace{h(\Sigma, Stats)}$$

DDR: 
$$\Sigma : B_1(\chi(t_1)) \vee \cdots \vee B_k(\chi(t_k)) \leftarrow R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$$

Compute bags  $B_1, \ldots, B_k$  in time  $\tilde{O}(2^{b(\Sigma, \text{Stats})})$ , repeat for all  $(t_1, \ldots, t_k)$ . Run Yannakakis' algorithm on each tree  $T_i$ . Runtime:  $\tilde{O}(2^{\text{subw}(Q, \text{Stats})})$ 

## PANDA: Discussion

Positives:

- Runtime  $\tilde{O}(2^{subw})$
- Steps of the algorithm entirely guided by solving an LP system.
- Generalizes to  $\omega$ -subw [Khamis et al., 2024a].

## PANDA: Discussion

Positives:

- Runtime  $\tilde{O}(2^{subw})$
- Steps of the algorithm entirely guided by solving an LP system.
- Generalizes to  $\omega$ -subw [Khamis et al., 2024a].

Negatives:

- Polylog factor. But is it needed???
- Computing subw is really hard!!
- Counting queries, sum-products: require restriction [Khamis et al., 2019]

Overview 00	Information Inequalities	Output Bounds 000000	Proofs to Algorithms	PANDA 0000000	Lower Bounds ●00000	

# Lower Bound Techniques
### Overview

Is  $O(2^{\text{subw}})$  optimal? Two questions:

• Is the polymatroid bound of a full CQ tight?

Do some accepted hypothesis imply a runtime of Ω(2<sup>subw</sup>)?

Neither question has a definitive answer, but we have some insights.

Proofs to Algor

PANDA Lower Bo

Lower Bounds Su

## Types of Polymatroids

Polymatroids with n variables X:

- Polymatroids Γ<sub>n</sub>: Shannon inequalities
- Almost entropic  $\overline{\Gamma}_n^*$ : closure of  $\Gamma_n^*$
- Entropic  $\Gamma_n^*$ : joint random variables
- Normal  $N_n$ :  $h(\boldsymbol{U}) = \mu(f(\boldsymbol{U}))$  where Set function:  $f : \boldsymbol{X} \to 2^{\Omega}$ Set measure:  $\mu : \Omega \to \mathbb{R}_p$
- Modular  $M_n$ :  $\Omega = X$ , f =identity.



## Examples





### Modular M<sub>3</sub>

## Examples









Modular M<sub>3</sub>

Normal N<sub>3</sub>

utput Bounds

Proofs to Algorith

PANDA 0000000 Lower Bounds

## Examples









X 0 1 1 Y 0 1 0 3 Z



Modular M<sub>3</sub>

Normal N<sub>3</sub>

Entropic  $\Gamma_3^*$ 

## Tightness of the Polymatroid Bound for Full CQ

• Cardinality stats:  $h^* \in M_n$  and the AGM bound is tight.  $\exists D \text{ s.t. } |Q(D)| \ge \frac{1}{2^n} AGM(Q).$  [Atserias et al., 2013]

• Add FDs: there is a query for which  $h^* \in \Gamma_n - \overline{\Gamma}_n^*$ . Not tight.

[Khamis et al., 2017, Gogacz and Torunczyk, 2017, Suciu, 2023]

• Add "simple" degree constraints:  $h^* \in N_n$ : bound is tight.  $\exists D, |Q(D)| \ge \frac{1}{2^{2^n-1}} 2^{h^*(X)} \text{ [Khamis et al., 2021, Suciu, 2023]:}$ 

### Lower Bounds [Fan et al., 2023]

Conjecture (Boolean Clique Conjecture)

 $\forall \varepsilon > 0$ , no algorithm can check for a *k*-clique in time  $O(n^{1-\varepsilon})$ .

 ${}^{3}\mu$  is a constant function and (1)  $\forall \omega \in \Omega$ ,  $\{Z \mid \omega \in f(Z)\}$  is connected in Q, (2)  $\forall \omega_{1}, \omega_{2} \in \Omega$ ,  $\exists j, \exists Y_{1}, Y_{2} \in U_{j}$  s.t.  $\omega_{i} \in f(Y_{i})$ , i = 1, 2.

## Conjecture (Boolean Clique Conjecture)

 $\forall \varepsilon > 0$ , no algorithm can check for a k-clique in time  $O(n^{1-\varepsilon})$ .

Boolean CQ 
$$Q() = \bigwedge_j R_j(\boldsymbol{U}_j)$$

An embedding of Q is  $h = \mu \circ f$ , satisfying some syntactic conditions<sup>3</sup>

Embedding power:  $|\operatorname{emb}(Q) \stackrel{\text{def}}{=} \operatorname{subw}(Q)$ -restricted to an embedding

 $^{3}\mu$  is a constant function and (1)  $\forall \omega \in \Omega$ ,  $\{Z \mid \omega \in f(Z)\}$  is connected in Q, (2)  $\forall \omega_{1}, \omega_{2} \in \Omega$ ,  $\exists j, \exists Y_{1}, Y_{2} \in U_{j}$  s.t.  $\omega_{i} \in f(Y_{i}), i = 1, 2$ .

# Lower Bounds [Fan et al., 2023]

Conjecture (Boolean Clique Conjecture)

 $\forall \varepsilon > 0$ , no algorithm can check for a *k*-clique in time  $O(n^{1-\varepsilon})$ .

Boolean CQ 
$$Q() = \bigwedge_j R_j(U_j)$$

An embedding of Q is  $h = \mu \circ f$ , satisfying some syntactic conditions<sup>3</sup>

Embedding power:  $emb(Q) \stackrel{\text{def}}{=} subw(Q)$ -restricted to an embedding

Theorem ([Fan et al., 2023])

 $\forall \varepsilon > 0$ , no algorithm can compute  $Q(\mathbf{D})$  in time  $O(2^{emb(Q)(1-\varepsilon)})$ .

 $^{3}\mu$  is a constant function and (1)  $\forall \omega \in \Omega$ ,  $\{Z \mid \omega \in f(Z)\}$  is connected in Q, (2)  $\forall \omega_{1}, \omega_{2} \in \Omega$ ,  $\exists j, \exists Y_{1}, Y_{2} \in U_{j}$  s.t.  $\omega_{i} \in f(Y_{i}), i = 1, 2$ .

## Discussion

Combinatorial Algorithms:  $Q(\mathbf{D})$  can be computed in time  $\tilde{O}(2^{\text{subw}})$ 

• Proofs to algorithms, Disjunctive Datalog Rules

Open problem Characterize divergent inequalities.

• Lower bounds: embeddings [Fan et al., 2023]; circuit lower bounds [Fan et al., 2024].

Summary

## Discussion

Combinatorial Algorithms:  $Q(\mathbf{D})$  can be computed in time  $\tilde{O}(2^{\text{subw}})$ 

• Proofs to algorithms, Disjunctive Datalog Rules

Open problem Characterize divergent inequalities.

• Lower bounds: embeddings [Fan et al., 2023]; circuit lower bounds [Fan et al., 2024].

Using FMM: Q(D) can be computed in time  $\tilde{O}(2^{\omega-\text{subw}})$  [Khamis et al., 2024a]. • Triangle:  $O(N^{\frac{2\omega}{\omega+1}})$ 

- 4-Clique:  $O(N^{\frac{\omega+1}{2}})$
- 4-Cycle:  $O(N^{\frac{4\omega-1}{2\omega+1}})$
- k-pyramid:  $\tilde{O}(N^{2-\frac{2}{\omega(k-1)-k+3}})$  new result

Summary



Finding and counting given length cycles. *Algorithmica*, 17(3):209–223.



Atserias, A., Grohe, M., and Marx, D. (2013).

Size bounds and query plans for relational joins. SIAM J. Comput., 42(4):1737–1767.



Chung, F. R. K., Graham, R. L., Frankl, P., and Shearer, J. B. (1986).

Some intersection theorems for ordered sets and graphs. J. Comb. Theory, Ser. A, 43(1):23–37.



#### Degree sequence bound for join cardinality estimation.

In Geerts, F. and Vandevoort, B., editors, 26th International Conference on Database Theory, ICDT 2023, March 28-31, 2023, Ioannina, Greece, volume 255 of LIPIcs, pages 8:1–8:18. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.

Deeds, K. B., Suciu, D., and Balazinska, M. (2023b).

Safebound: A practical system for generating cardinality bounds. *Proc. ACM Manag. Data*, 1(1):53:1–53:26.



Fan, A. Z., Koutris, P., and Zhao, H. (2023).

The fine-grained complexity of boolean conjunctive queries and sum-product problems.

In Etessami, K., Feige, U., and Puppis, G., editors, 50th International Colloquium on Automata, Languages, and Programming, ICALP 2023, July 10-14, 2023, Paderborn, Germany, volume 261 of LIPIcs, pages 127:1–127:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.



Fan, A. Z., Koutris, P., and Zhao, H. (2024).

Tight bounds of circuits for sum-product queries. Proc. ACM Manag. Data, 2(2):87.

#### Gogacz, T. and Torunczyk, S. (2017).

Entropy bounds for conjunctive queries with functional dependencies.

In Benedikt, M. and Orsi, G., editors, 20th International Conference on Database Theory, ICDT 2017, March 21-24, 2017, Venice, Italy, volume 68 of LIPIcs, pages 15:1–15:17. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.



Gottlob, G., Lee, S. T., Valiant, G., and Valiant, P. (2012).

Size and treewidth bounds for conjunctive queries. J. ACM, 59(3):16:1–16:35.



Grohe, M. and Marx, D. (2014).

Constraint solving via fractional edge covers. ACM Trans. Algorithms, 11(1):4:1-4:20.



Khamis, M. A., Curtin, R. R., Moseley, B., Ngo, H. Q., Nguyen, X., Olteanu, D., and Schleich, M. (2019).

On functional aggregate queries with additive inequalities.

In Suciu, D., Skritek, S., and Koch, C., editors, Proceedings of the 38th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2019, Amsterdam, The Netherlands, June 30 - July 5, 2019, pages 414–431. ACM.



Khamis, M. A., Hu, X., and Suciu, D. (2024a).

Fast matrix multiplication meets the subdmodular width. *CoRR*, abs/2412.06189.



Khamis, M. A., Kolaitis, P. G., Ngo, H. Q., and Suciu, D. (2021).

Bag query containment and information theory. *ACM Trans. Database Syst.*, 46(3):12:1–12:39.



Khamis, M. A., Nakos, V., Olteanu, D., and Suciu, D. (2024b).

Join size bounds using I<sub>p</sub>-norms on degree sequences. *Proc. ACM Manag. Data*, 2(2):96.



Khamis, M. A., Ngo, H. Q., and Rudra, A. (2016a).

#### FAQ: questions asked frequently.

In Milo, T. and Tan, W., editors, Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2016, San Francisco, CA, USA, June 26 - July 01, 2016, pages 13–28. ACM.

#### Khamis, M. A., Ngo, H. Q., and Suciu, D. (2016b).

Computing join queries with functional dependencies.

In Milo, T. and Tan, W., editors, Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2016, San Francisco, CA, USA, June 26 - July 01, 2016, pages 327–342. ACM.

#### Khamis, M. A., Ngo, H. Q., and Suciu, D. (2017).

What do shannon-type inequalities, submodular width, and disjunctive datalog have to do with one another? In Sallinger, E., den Bussche, J. V., and Geerts, F., editors, *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2017, Chicago, IL, USA, May 14-19, 2017*, pages 429–444. ACM. Extended version available at http://arxiv.org/abs/1612.02603.



Khamis, M. A., Ngo, H. Q., and Suciu, D. (2024c).

PANDA: query evaluation in submodular width. *CoRR*, abs/2402.02001.



Marx, D. (2013).

Tractable hypergraph properties for constraint satisfaction and conjunctive queries. J. ACM, 60(6):42:1-42:51.



#### Matús, F. (2007).

Infinitely many information inequalities.

In IEEE International Symposium on Information Theory, ISIT 2007, Nice, France, June 24-29, 2007, pages 41-44. IEEE.



#### Ngo, H. Q. (2022).

#### On an information theoretic approach to cardinality estimation (invited talk).

In Olteanu, D. and Vortmeier, N., editors, 25th International Conference on Database Theory, ICDT 2022, March 29 to April 1, 2022, Edinburgh, UK (Virtual Conference), volume 220 of LIPIcs, pages 1:1–1:21. Schloss Dagstuhl -Leibniz-Zentrum für Informatik.



Pippenger, N. (1986).

What are the laws of information theory.

In 1986 Special Problems on Communication and Computation Conference, pages 3–5.



Shannon, C. E. (1948).

#### A mathematical theory of communication.

Bell Syst. Tech. J., 27(3):379-423.



#### Suciu, D. (2023).

Applications of information inequalities to database theory problems. In LICS, pages 1–30.



Yeung, R. W. (2008).

Information Theory and Network Coding. Springer Publishing Company, Incorporated, 1 edition.



Zhang, H., Mayer, C., Khamis, M. A., Olteanu, D., and Suciu, D. (2025).

Lpbound: Pessimistic cardinality estimation using p-norms of degree sequences. *CoRR*, abs/2502.05912.



Zhang, Z. and Yeung, R. W. (1998).

On characterization of entropy function via information inequalities. *IEEE Transactions on Information Theory*, 44(4):1440–1452.

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \wedge \cdots \wedge R_m(\boldsymbol{U}_m)$ ,

Statistics are cardinality constraints:  $|R_j| \leq N_j$ , j = 1, m.

Fix any fractional edge cover  $w_1, \ldots, w_m$ .

Theorem (The AGM Bound [Atserias et al., 2013])

If **D** satisfies the cardinality constraints, then  $|Q(\mathbf{D})| \leq \prod_i |N_j|^{w_j}$ .

We will prove the AGM bound, then derive an algorithm from its proof.

32 / 34

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \land \cdots \land R_m(\boldsymbol{U}_m)$ , frac. edge cover  $w_1, \ldots, w_m$ .

Shearer's Inequality:  $\sum_{j} w_{j}h(\boldsymbol{U}_{j}) \geq h(\boldsymbol{X})$ 

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \land \cdots \land R_m(\boldsymbol{U}_m)$ , frac. edge cover  $w_1, \ldots, w_m$ .

Shearer's Inequality:  $\sum_{j} w_{j}h(\boldsymbol{U}_{j}) \geq h(\boldsymbol{X})$ 

If  $X_1 \in \boldsymbol{U}_j$ ,  $h(\boldsymbol{U}_j) = h(X_1) + h(\boldsymbol{U}_j|X_1)$ ; If  $X_1 \notin \boldsymbol{U}_j$ ,  $h(\boldsymbol{U}_j) \ge h(\boldsymbol{U}_j|X_1)$ .

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \land \cdots \land R_m(\boldsymbol{U}_m)$ , frac. edge cover  $w_1, \ldots, w_m$ .

Shearer's Inequality:  $\sum_{j} w_{j}h(\boldsymbol{U}_{j}) \geq h(\boldsymbol{X})$ 

If  $X_1 \in \boldsymbol{U}_j$ ,  $h(\boldsymbol{U}_j) = h(X_1) + h(\boldsymbol{U}_j|X_1)$ ; If  $X_1 \notin \boldsymbol{U}_j$ ,  $h(\boldsymbol{U}_j) \ge h(\boldsymbol{U}_j|X_1)$ .

 $\sum w_j h(oldsymbol{U}_j) \geq (\sum_{j:X_1 \in oldsymbol{U}_j} w_j) h(X_1) + \sum_j w_j h(oldsymbol{U}_j|X_1) \geq h(X_1) + h(oldsymbol{X}|X_1) = h(oldsymbol{X})$ 

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \land \cdots \land R_m(\boldsymbol{U}_m)$ , frac. edge cover  $w_1, \ldots, w_m$ .

Shearer's Inequality:  $\sum_{j} w_{j} h(\boldsymbol{U}_{j}) \geq h(\boldsymbol{X})$ 

If  $X_1 \in \boldsymbol{U}_j$ ,  $h(\boldsymbol{U}_j) = h(X_1) + h(\boldsymbol{U}_j|X_1)$ ; If  $X_1 \notin \boldsymbol{U}_j$ ,  $h(\boldsymbol{U}_j) \ge h(\boldsymbol{U}_j|X_1)$ .

 $\sum w_j h(oldsymbol{U}_j) \geq (\sum_{j:X_1 \in oldsymbol{U}_j} w_j) h(X_1) + \sum_j w_j h(oldsymbol{U}_j | X_1) \geq h(X_1) + h(oldsymbol{X} | X_1) = h(oldsymbol{X})$ 

Generic Join

for  $x_1 \in \mathsf{Dom}(X_1)$ compute  $Q[\boldsymbol{X}|X_1 := x_1]$ 

Runtime:  $\tilde{O}(\prod_j |R_j|^{w_j})$ 

Full CQ:  $Q(\boldsymbol{X}) = R_1(\boldsymbol{U}_1) \land \cdots \land R_m(\boldsymbol{U}_m)$ , frac. edge cover  $w_1, \ldots, w_m$ .

Shearer's Inequality:  $\sum_{j} w_{j} h(\boldsymbol{U}_{j}) \geq h(\boldsymbol{X})$ 

If  $X_1 \in \boldsymbol{U}_j$ ,  $h(\boldsymbol{U}_j) = h(X_1) + h(\boldsymbol{U}_j|X_1)$ ; If  $X_1 \notin \boldsymbol{U}_j$ ,  $h(\boldsymbol{U}_j) \ge h(\boldsymbol{U}_j|X_1)$ .

 $\sum w_j h(oldsymbol{U}_j) \geq (\sum_{j:X_1 \in oldsymbol{U}_j} w_j) h(X_1) + \sum_j w_j h(oldsymbol{U}_j | X_1) \geq h(X_1) + h(oldsymbol{X} | X_1) = h(oldsymbol{X})$ 

Generic JoinE.g.  $Q_{\Delta}(XYZ) = R(XY) \land S(YZ) \land T(ZX)$ :for  $x_1 \in \text{Dom}(X_1)$ <br/>compute  $Q[\boldsymbol{X}|X_1 := x_1]$ for  $x \in R[X] \cap T[X]$ <br/>for  $y \in R[Y|X] \cap S[Y]$ <br/>for  $z \in S[Z|y] \cap T[Z|X]$ <br/>output(x, y, z)

 $\overline{Q() = R(XY) \land S(YZ) \land T}(ZU) \land K(UX)$ 

Cardinality constraints  $|R|, |S|, |T|, |K| \leq N$ .

$$Q() = R(XY) \land S(YZ) \land T(ZU) \land K(UX)$$

Cardinality constraints  $|R|, |S|, |T|, |K| \leq N$ .



$$Q() = R(XY) \land S(YZ) \land T(ZU) \land K(UX)$$

Cardinality constraints  $|R|, |S|, |T|, |K| \leq N$ .

Claim:

$$subw(Q) = 3/2$$



$$Q() = R(XY) \land S(YZ) \land T(ZU) \land K(UX)$$

Cardinality constraints  $|R|, |S|, |T|, |K| \leq N$ .

Claim:

$$subw(Q) = 3/2$$



$$Q() = R(XY) \land S(YZ) \land T(ZU) \land K(UX)$$

Cardinality constraints  $|R|, |S|, |T|, |K| \leq N$ .

Claim:

$$subw(Q) = 3/2$$



 $3 \geq$ 

### $2\min(h(XYZ), h(YZU))$

$$Q() = R(XY) \land S(YZ) \land T(ZU) \land K(UX)$$

Cardinality constraints  $|R|, |S|, |T|, |K| \leq N$ .

Claim:

$$\mathsf{subw}(Q) = 3/2$$



 $3 \ge h(XY) + h(YZ) + h(ZU) = h(XY) + h(Z|Y) + h(Y) + h(ZU)$  $\ge h(XYZ) + h(YZU) \ge 2\min(h(XYZ), h(YZU))$ 

$$Q() = R(XY) \land S(YZ) \land T(ZU) \land K(UX)$$

Cardinality constraints  $|R|, |S|, |T|, |K| \leq N$ .

Claim:

$$\mathsf{subw}(Q) = 3/2$$



 $3 \ge h(XY) + h(YZ) + h(ZU) = h(XY) + h(Z|Y) + h(Y) + h(ZU)$  $\ge h(XYZ) + h(YZU) \ge 2\min(h(XYZ), h(YZU))$ 

 $B_1(XYZ) \lor B_3(YZU) \leftarrow R(XY) \land S(YZ) \land T(ZU)$ 

$$Q() = R(XY) \land S(YZ) \land T(ZU) \land K(UX)$$

Cardinality constraints  $|R|, |S|, |T|, |K| \leq N$ .

Claim:

$$subw(Q) = 3/2$$



 $3 \ge h(XY) + h(YZ) + h(ZU) = h(XY) + h(Z|Y) + h(Y) + h(ZU)$  $\ge h(XYZ) + h(YZU) \ge 2\min(h(XYZ), h(YZU))$ 

 $B_1(XYZ) \lor B_3(YZU) \leftarrow R(XY) \land S(YZ) \land T(ZU)$ 

Compute in time  $\tilde{O}(N^{3/2})$  (light/heavy), repeat for 4 DDRs.