

# What do Shannon-type Inequalities, Submodular Width, and Disjunctive Datalog have to do with one another?

Mahmoud Abo Khamis<sup>1</sup>   Hung Q. Ngo<sup>1</sup>   Dan Suciu<sup>1,2</sup>

<sup>1</sup>LogicBlox Inc.

<sup>2</sup>University of Washington

PODS 2017



# Contributions

- ▶ A Join Algorithm

# Contributions

- ▶ **A Join Algorithm**
  - ▶ first to meet the **Submodular Width** bound!

# Contributions

- ▶ **A Join Algorithm**
  - ▶ first to meet the **Submodular Width** bound!
  - ▶ works for and relies on **Disjunctive Datalog**.

# Contributions

- ▶ A **Join Algorithm**
  - ▶ first to meet the **Submodular Width** bound!
  - ▶ works for and relies on **Disjunctive Datalog**.
  - ▶ fully utilizes **Functional DEPs** and **Degree Bounds**.

# Contributions

- ▶ A **Join Algorithm**
  - ▶ first to meet the **Submodular Width** bound!
  - ▶ works for and relies on **Disjunctive Datalog**.
  - ▶ fully utilizes **Functional DEPs** and **Degree Bounds**.
- ▶ A **Unified Framework for Join Bounds**

# Contributions

- ▶ A **Join Algorithm**
  - ▶ first to meet the **Submodular Width** bound!
  - ▶ works for and relies on **Disjunctive Datalog**.
  - ▶ fully utilizes **Functional DEPs** and **Degree Bounds**.
- ▶ A Unified Framework for **Join Bounds**
  - ▶ **subsumes** most known bounds.

# Contributions

- ▶ A **Join Algorithm**
  - ▶ first to meet the **Submodular Width** bound!
  - ▶ works for and relies on **Disjunctive Datalog**.
  - ▶ fully utilizes **Functional DEPs** and **Degree Bounds**.
- ▶ A Unified Framework for **Join Bounds**
  - ▶ **subsumes** most known bounds.
  - ▶ **extends** them to **Functional DEPs** and **Degree Bounds**.



# Contributions

- ▶ A **Join Algorithm**
  - ▶ first to meet the **Submodular Width** bound!
  - ▶ works for and relies on **Disjunctive Datalog**.
  - ▶ fully utilizes **Functional DEPs** and **Degree Bounds**.
- ▶ A Unified Framework for **Join Bounds**
  - ▶ **subsumes** most known bounds.
  - ▶ **extends** them to Functional DEPs and Degree Bounds.
- ▶ **Results on Shannon-type Inequalities**

# Table of Contents

Size Bounds for Full Conjunctive Queries

Size Bounds for Disjunctive Datalog

Algorithms for Disjunctive Datalog

Algorithms for Conjunctive Queries

# Table of Contents

Size Bounds for Full Conjunctive Queries

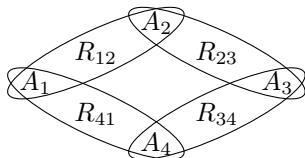
Size Bounds for Disjunctive Datalog

Algorithms for Disjunctive Datalog

Algorithms for Conjunctive Queries

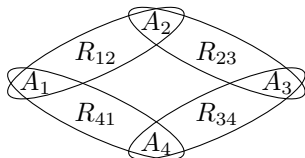
# Size Bounds for Full Conjunctive Queries

$Q(A_1, A_2, A_3, A_4) :-$   $R_{12}(A_1, A_2), R_{23}(A_2, A_3),$   
 $R_{34}(A_3, A_4), R_{41}(A_4, A_1).$



# Size Bounds for Full Conjunctive Queries

$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3),$   
 $R_{34}(A_3, A_4), R_{41}(A_4, A_1).$



$A_1$	$A_2$
a	1
b	1
b	2

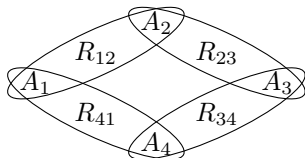
$A_2$	$A_3$
1	c
1	d
2	c

$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

# Size Bounds for Full Conjunctive Queries

$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3),$   
 $R_{34}(A_3, A_4), R_{41}(A_4, A_1).$



$A_1$	$A_2$	$A_3$	$A_4$
a	1	d	4
b	1	c	3
b	1	d	4
b	2	c	3

$A_1$	$A_2$
a	1
b	1
b	2

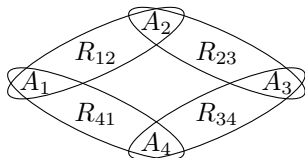
$A_2$	$A_3$
1	c
1	d
2	c

$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

# Size Bounds for Full Conjunctive Queries

$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3),$   
 $R_{34}(A_3, A_4), R_{41}(A_4, A_1).$



$A_1$	$A_2$	$A_3$	$A_4$	
a	1	d	4	1/4
b	1	c	3	1/4
b	1	d	4	1/4
b	2	c	3	1/4

$A_1$	$A_2$
a	1
b	1
b	2

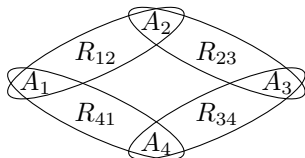
$A_2$	$A_3$
1	c
1	d
2	c

$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

# Size Bounds for Full Conjunctive Queries

$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3),$   
 $R_{34}(A_3, A_4), R_{41}(A_4, A_1).$



$A_1$	$A_2$	$A_3$	$A_4$	
a	1	d	4	1/4
b	1	c	3	1/4
b	1	d	4	1/4
b	2	c	3	1/4

$A_1$	$A_2$	
a	1	1/4
b	1	2/4
b	2	1/4

$A_2$	$A_3$	
1	c	1/4
1	d	2/4
2	c	1/4

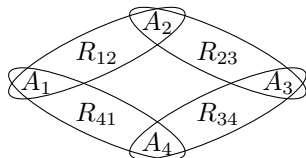
$A_3$	$A_4$	
c	3	2/4
d	4	2/4
d	5	0

$A_4$	$A_1$	
3	b	2/4
4	a	1/4
4	b	1/4



# Size Bounds for Full Conjunctive Queries

$$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



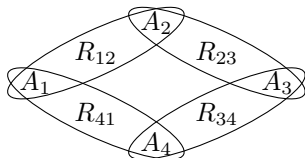
$A_1$	$A_2$	$A_3$	$A_4$	
a	1	d	4	1/4
b	1	c	3	1/4
b	1	d	4	1/4
b	2	c	3	1/4

$$h(A_1 A_2 A_3 A_4) = \log |Q|$$

$A_1$	$A_2$		$A_2$	$A_3$		$A_3$	$A_4$		$A_4$	$A_1$	
a	1	1/4	1	c	1/4	c	3	2/4	3	b	2/4
b	1	2/4	1	d	2/4	d	4	2/4	4	a	1/4
b	2	1/4	2	c	1/4	d	5	0	4	b	1/4

# Size Bounds for Full Conjunctive Queries

$$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$A_1$	$A_2$	$A_3$	$A_4$	
a	1	d	4	1/4
b	1	c	3	1/4
b	1	d	4	1/4
b	2	c	3	1/4

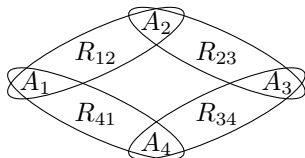
$$h(A_1 A_2 A_3 A_4) = \log |Q|$$

$A_1$	$A_2$		$A_2$	$A_3$		$A_3$	$A_4$		$A_4$	$A_1$	
a	1	1/4	1	c	1/4	c	3	2/4	3	b	2/4
b	1	2/4	1	d	2/4	d	4	2/4	4	a	1/4
b	2	1/4	2	c	1/4	d	5	0	4	b	1/4

$$h(A_1 A_2) \leq \log |R_{12}|, \quad h(A_2 A_3) \leq \log |R_{23}|, \quad h(A_3 A_4) \leq \log |R_{34}|, \quad \dots$$

# Size Bounds for Full Conjunctive Queries

$$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$A_1$	$A_2$	$A_3$	$A_4$	
a	1	d	4	1/4
b	1	c	3	1/4
b	1	d	4	1/4
b	2	c	3	1/4

$$h(A_1 A_2 A_3 A_4) = \log |Q|$$

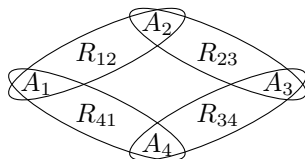
$A_1$	$A_2$		$A_2$	$A_3$		$A_3$	$A_4$		$A_4$	$A_1$	
a	1	1/4	1	c	1/4	c	3	2/4	3	b	2/4
b	1	2/4	1	d	2/4	d	4	2/4	4	a	1/4
b	2	1/4	2	c	1/4	d	5	0	4	b	1/4

$$h(A_1 A_2) \leq \log |R_{12}|, \quad h(A_2 A_3) \leq \log |R_{23}|, \quad h(A_3 A_4) \leq \log |R_{34}|, \quad \dots$$

$$h(A_2 | A_1 = 'a') \leq \log |\sigma_{A_1='a'} R_{12}|, \quad h(A_2 | A_1 = 'b') \leq \log |\sigma_{A_1='b'} R_{12}|, \quad \dots$$

# Size Bounds for Full Conjunctive Queries

$$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$A_1$	$A_2$	$A_3$	$A_4$	
a	1	d	4	1/4
b	1	c	3	1/4
b	1	d	4	1/4
b	2	c	3	1/4

$$h(A_1 A_2 A_3 A_4) = \log |Q|$$

$A_1$	$A_2$		$A_2$	$A_3$		$A_3$	$A_4$		$A_4$	$A_1$	
a	1	1/4	1	c	1/4	c	3	2/4	3	b	2/4
b	1	2/4	1	d	2/4	d	4	2/4	4	a	1/4
b	2	1/4	2	c	1/4	d	5	0	4	b	1/4

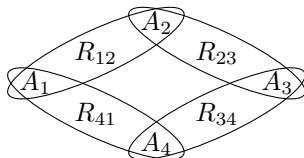
$$h(A_1 A_2) \leq \log |R_{12}|, \quad h(A_2 A_3) \leq \log |R_{23}|, \quad h(A_3 A_4) \leq \log |R_{34}|, \quad \dots$$

$$h(A_2 | A_1 = 'a') \leq \log |\sigma_{A_1='a'} R_{12}|, \quad h(A_2 | A_1 = 'b') \leq \log |\sigma_{A_1='b'} R_{12}|, \quad \dots$$

$$h(A_2 | A_1) \leq \log \max_x |\sigma_{A_1=x} R_{12}|$$

# Size Bounds for Full Conjunctive Queries

$$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$A_1$	$A_2$	$A_3$	$A_4$	
a	1	d	4	1/4
b	1	c	3	1/4
b	1	d	4	1/4
b	2	c	3	1/4

$$h(A_1 A_2 A_3 A_4) = \log |Q|$$

$A_1$	$A_2$		$A_2$	$A_3$		$A_3$	$A_4$		$A_4$	$A_1$	
a	1	1/4	1	c	1/4	c	3	2/4	3	b	2/4
b	1	2/4	1	d	2/4	d	4	2/4	4	a	1/4
b	2	1/4	2	c	1/4	d	5	0	4	b	1/4

$$h(A_1 A_2) \leq \log |R_{12}|, \quad h(A_2 A_3) \leq \log |R_{23}|, \quad h(A_3 A_4) \leq \log |R_{34}|, \quad \dots$$

$$h(A_2 | A_1 = 'a') \leq \log |\sigma_{A_1='a'} R_{12}|, \quad h(A_2 | A_1 = 'b') \leq \log |\sigma_{A_1='b'} R_{12}|, \quad \dots$$

$$h(A_2 | A_1) \leq \log \underbrace{\max_x |\sigma_{A_1=x} R_{12}|}_{\text{deg}_{R_{12}}(A_2 | A_1)}$$

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes
- ▶  $\mathbf{A}_F = \{A_f \mid f \in F \subseteq [n]\}$  is a subset



# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes
- ▶  $\mathbf{A}_F = \{A_f \mid f \in F \subseteq [n]\}$  is a subset

- ▶ Degree Constraints (DC):

$$\deg_F(\mathbf{A}_Y | \mathbf{A}_X) \leq N_{Y|X}, \quad X \subset Y \subseteq F \in \mathcal{E}$$

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes
- ▶  $\mathbf{A}_F = \{A_f \mid f \in F \subseteq [n]\}$  is a subset

- ▶ Degree Constraints (DC):

$$\deg_F(\mathbf{A}_Y | \mathbf{A}_X) \leq N_{Y|X}, \quad X \subset Y \subseteq F \in \mathcal{E}$$

- ▶ Cardinality Constraints (CC):  $|R_F| \leq N_{F|\emptyset}$

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes
- ▶  $\mathbf{A}_F = \{A_f \mid f \in F \subseteq [n]\}$  is a subset

- ▶ Degree Constraints (DC):

$$\text{deg}_F(\mathbf{A}_Y | \mathbf{A}_X) \leq N_{Y|X}, \quad X \subset Y \subseteq F \in \mathcal{E}$$

- ▶ Cardinality Constraints (CC):  $|R_F| \leq N_{F|\emptyset}$
- ▶ Functional Dependencies (FD):  $\mathbf{A}_X \rightarrow \mathbf{A}_Y$

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes
- ▶  $\mathbf{A}_F = \{A_f \mid f \in F \subseteq [n]\}$  is a subset

- ▶ Degree Constraints (DC):

$$\text{deg}_F(\mathbf{A}_Y | \mathbf{A}_X) \leq N_{Y|X}, \quad X \subset Y \subseteq F \in \mathcal{E}$$

- ▶ Cardinality Constraints (CC):  $|R_F| \leq N_{F|\emptyset}$
- ▶ Functional Dependencies (FD):  $\mathbf{A}_X \rightarrow \mathbf{A}_Y$

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes
- ▶  $\mathbf{A}_F = \{A_f \mid f \in F \subseteq [n]\}$  is a subset

- ▶ Degree Constraints (DC):

$$\deg_F(\mathbf{A}_Y | \mathbf{A}_X) \leq N_{Y|X}, \quad X \subset Y \subseteq F \in \mathcal{E}$$

- ▶ Cardinality Constraints (CC):  $|R_F| \leq N_{F|\emptyset}$
- ▶ Functional Dependencies (FD):  $\mathbf{A}_X \rightarrow \mathbf{A}_Y$

## Bound Idea

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes
- ▶  $\mathbf{A}_F = \{A_f \mid f \in F \subseteq [n]\}$  is a subset

- ▶ Degree Constraints (DC):

$$\deg_F(\mathbf{A}_Y | \mathbf{A}_X) \leq N_{Y|X}, \quad X \subset Y \subseteq F \in \mathcal{E}$$

- ▶ Cardinality Constraints (CC):  $|R_F| \leq N_{F|\emptyset}$
- ▶ Functional Dependencies (FD):  $\mathbf{A}_X \rightarrow \mathbf{A}_Y$

## Bound Idea

$$\log |Q| \leq \text{maximum } h(A_1, \dots, A_n)$$

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes
- ▶  $\mathbf{A}_F = \{A_f \mid f \in F \subseteq [n]\}$  is a subset

- ▶ Degree Constraints (DC):

$$\deg_F(\mathbf{A}_Y | \mathbf{A}_X) \leq N_{Y|X}, \quad X \subset Y \subseteq F \in \mathcal{E}$$

- ▶ Cardinality Constraints (CC):  $|R_F| \leq N_{F|\emptyset}$
- ▶ Functional Dependencies (FD):  $\mathbf{A}_X \rightarrow \mathbf{A}_Y$

## Bound Idea

$$\log |Q| \leq \text{maximum over all entropies } h(A_1, \dots, A_n)$$

# Size Bounds: Input

- ▶ A full conjunctive query

$$Q(\mathbf{A}_{[n]}) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

- ▶  $\mathbf{A}_{[n]} = \{A_1, \dots, A_n\}$  is the full set of attributes
- ▶  $\mathbf{A}_F = \{A_f \mid f \in F \subseteq [n]\}$  is a subset
- ▶ Degree Constraints (DC):

$$\deg_F(\mathbf{A}_Y \mid \mathbf{A}_X) \leq N_{Y \mid X}, \quad X \subset Y \subseteq F \in \mathcal{E}$$

- ▶ Cardinality Constraints (CC):  $|R_F| \leq N_{F \mid \emptyset}$
- ▶ Functional Dependencies (FD):  $\mathbf{A}_X \rightarrow \mathbf{A}_Y$

## Bound Idea

$\log |Q| \leq$  maximum  $h(A_1, \dots, A_n)$   
over all *entropies*  $h$   
such that  $h$  satisfies degree constraints of  $Q$



# Size Bounds: Preliminaries

- ▶ **HDC** is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  *satisfying* the degree constraints

$$\text{HDC} \stackrel{\text{def}}{=} \{h \mid h(Y|X) \leq \log N_{Y|X}, \quad \forall (X, Y, N_{Y|X})\}$$

## Size Bounds: Preliminaries

- ▶ **HDC** is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  *satisfying* the degree constraints

$$\text{HDC} \stackrel{\text{def}}{=} \{h \mid h(Y|X) \leq \log N_{Y|X}, \quad \forall (X, Y, N_{Y|X})\}$$

- ▶  $\Gamma_n^*$  is the set of *entropic* functions

# Size Bounds: Preliminaries

- ▶ **HDC** is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  *satisfying* the degree constraints

$$\text{HDC} \stackrel{\text{def}}{=} \{h \mid h(Y|X) \leq \log N_{Y|X}, \quad \forall (X, Y, N_{Y|X})\}$$

- ▶  $\Gamma_n^*$  is the set of *entropic* functions
- ▶  $\overline{\Gamma}_n^*$  is the topological closure of  $\Gamma_n^*$

## Size Bounds: Preliminaries

- ▶ **HDC** is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  *satisfying* the degree constraints

$$\text{HDC} \stackrel{\text{def}}{=} \{h \mid h(Y|X) \leq \log N_{Y|X}, \quad \forall (X, Y, N_{Y|X})\}$$

- ▶  $\Gamma_n^*$  is the set of *entropic* functions
- ▶  $\overline{\Gamma}_n^*$  is the topological closure of  $\Gamma_n^*$
- ▶  $\Gamma_n$  is the set of *polymatroids*, i.e. functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  satisfying

$$\begin{aligned} h(X \cup Y) + h(X \cap Y) &\leq h(X) + h(Y), & X, Y &\subseteq [n] && \text{(submodularity)} \\ h(X) &\leq h(Y), & X &\subseteq Y \subseteq [n] && \text{(monotonicity)} \\ h(\emptyset) &= 0 &&&& \text{(strictness)} \end{aligned}$$

# Size Bounds: Preliminaries

- ▶ **HDC** is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  *satisfying* the degree constraints

$$\text{HDC} \stackrel{\text{def}}{=} \{h \mid h(Y|X) \leq \log N_{Y|X}, \quad \forall (X, Y, N_{Y|X})\}$$

- ▶  $\Gamma_n^*$  is the set of *entropic* functions
- ▶  $\bar{\Gamma}_n^*$  is the topological closure of  $\Gamma_n^*$
- ▶  $\Gamma_n$  is the set of *polymatroids*, i.e. functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  satisfying

$$\begin{aligned} h(X \cup Y) + h(X \cap Y) &\leq h(X) + h(Y), & X, Y \subseteq [n] & & \text{(submodularity)} \\ h(X) &\leq h(Y), & X \subseteq Y \subseteq [n] & & \text{(monotonicity)} \\ h(\emptyset) &= 0 & & & \text{(strictness)} \end{aligned}$$

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

# Size Bounds for Full Conjunctive Queries

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

# Size Bounds for Full Conjunctive Queries

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

Bound	Entropic Bound	Polymatroid Bound
Definition	$\log  Q  \leq \max_{h \in \bar{\Gamma}_n^* \cap \text{HDC}} h([n])$	$\log  Q  \leq \max_{h \in \Gamma_n \cap \text{HDC}} h([n])$

# Size Bounds for Full Conjunctive Queries

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

Bound	Entropic Bound	Polymatroid Bound
Definition	$\log  Q  \leq \max_{h \in \bar{\Gamma}_n^* \cap \text{HDC}} h([n])$	$\log  Q  \leq \max_{h \in \Gamma_n \cap \text{HDC}} h([n])$
CC only	AGM bound (Tight) [Atserias et al. FOCS'08]	AGM bound (Tight) [Atserias et al. FOCS'08]



# Size Bounds for Full Conjunctive Queries

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

Bound	Entropic Bound	Polymatroid Bound
Definition	$\log  Q  \leq \max_{h \in \bar{\Gamma}_n^* \cap \text{HDC}} h([n])$	$\log  Q  \leq \max_{h \in \Gamma_n \cap \text{HDC}} h([n])$
CC only	AGM bound (Tight) [Atserias et al. FOCS'08]	AGM bound (Tight) [Atserias et al. FOCS'08]
CC + FD only	<b>Entropic</b> Bound for FD [Gottlob et al. JACM'12] (Tight [Gogacz et al. ICDT'17])	<b>Polymatroid</b> Bound for FD [Gottlob et al. JACM'12] (Not tight [Our work] )

# Size Bounds for Full Conjunctive Queries

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

Bound	Entropic Bound	Polymatroid Bound
Definition	$\log  Q  \leq \max_{h \in \bar{\Gamma}_n^* \cap \text{HDC}} h([n])$	$\log  Q  \leq \max_{h \in \Gamma_n \cap \text{HDC}} h([n])$
CC only	AGM bound (Tight) [Atserias et al. FOCS'08]	AGM bound (Tight) [Atserias et al. FOCS'08]
CC + FD only	Entropic Bound for FD [Gottlob et al. JACM'12] (Tight [Gogacz et al. ICDT'17])	Polymatroid Bound for FD [Gottlob et al. JACM'12] (Not tight [Our work] )
DC	Entropic Bound for DC (Tight [Our work] )	Polymatroid Bound for DC (Not tight [Our work] )

# Table of Contents

Size Bounds for Full Conjunctive Queries

**Size Bounds for Disjunctive Datalog**

Algorithms for Disjunctive Datalog

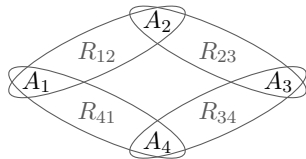
Algorithms for Conjunctive Queries

# Disjunctive Datalog

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

# Disjunctive Datalog

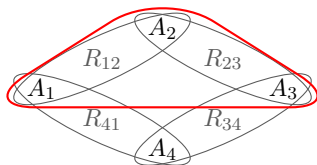
$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

# Disjunctive Datalog

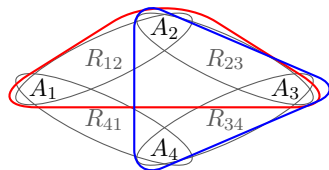
$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

# Disjunctive Datalog

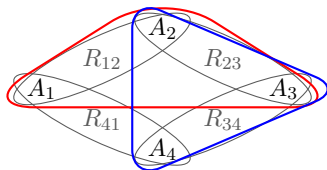
$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

# Disjunctive Datalog

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$A_1$	$A_2$
a	1
b	1
b	2

$A_2$	$A_3$
1	c
1	d
2	c

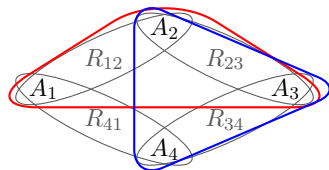
$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b



# Disjunctive Datalog

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$A_1$	$A_2$
a	1
b	1
b	2

$A_2$	$A_3$
1	c
1	d
2	c

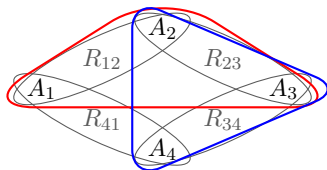
$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

$A_1$	$A_2$	$A_3$	$A_4$
a	1	d	4
b	1	c	3
b	1	d	4
b	2	c	3

# Disjunctive Datalog

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$A_1$	$A_2$
a	1
b	1
b	2

$A_2$	$A_3$
1	c
1	d
2	c

$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

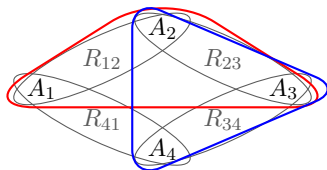
$A_1$	$A_2$	$A_3$	$A_4$
a	1	d	4
b	1	c	3
b	1	d	4
b	2	c	3

$A_1$	$A_2$	$A_3$
b	1	c
b	2	c

$A_2$	$A_3$	$A_4$
1	d	4
2	c	3
2	d	4

# Disjunctive Datalog

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$A_1$	$A_2$
a	1
b	1
b	2

$A_2$	$A_3$
1	c
1	d
2	c

$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

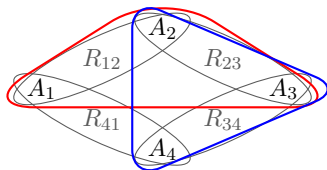
$A_1$	$A_2$	$A_3$	$A_4$
a	1	d	4
b	1	c	3
b	1	d	4
b	2	c	3

$A_1$	$A_2$	$A_3$
b	1	c
b	2	c

$A_2$	$A_3$	$A_4$
1	d	4
2	c	3
2	d	4

# Disjunctive Datalog

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$A_1$	$A_2$
a	1
b	1
b	2

$A_2$	$A_3$
1	c
1	d
2	c

$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

$A_1$	$A_2$	$A_3$	$A_4$
a	1	d	4
b	1	c	3
b	1	d	4
b	2	c	3

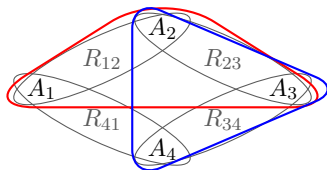
$A_1$	$A_2$	$A_3$
b	1	c
b	2	c

$A_2$	$A_3$	$A_4$
1	d	4
2	c	3
2	d	4

Model size is  $\max(|T_{123}|, |T_{234}|) = 3$

# Disjunctive Datalog

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$A_1$	$A_2$
a	1
b	1
b	2

$A_2$	$A_3$
1	c
1	d
2	c

$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

$A_1$	$A_2$	$A_3$	$A_4$
a	1	d	4
b	1	c	3
b	1	d	4
b	2	c	3

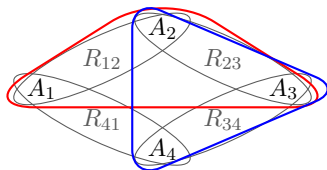
$A_1$	$A_2$	$A_3$
b	1	c
b	2	c

$A_2$	$A_3$	$A_4$
1	d	4
2	c	3
2	d	4

Output size is the *minimum* over all models

# Disjunctive Datalog

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$A_1$	$A_2$
a	1
b	1
b	2

$A_2$	$A_3$
1	c
1	d
2	c

$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

$A_1$	$A_2$	$A_3$	$A_4$
a	1	d	4
b	1	c	3
b	1	d	4
b	2	c	3

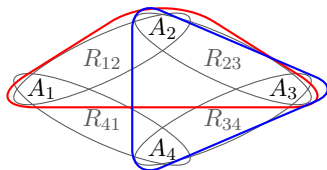
$A_1$	$A_2$	$A_3$
b	1	c
b	2	c

$A_2$	$A_3$	$A_4$
1	d	4

A minimum-sized model of size 2

# Disjunctive Datalog

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$



$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$A_1$	$A_2$
a	1
b	1
b	2

$A_2$	$A_3$
1	c
1	d
2	c

$A_3$	$A_4$
c	3
d	4
d	5

$A_4$	$A_1$
3	b
4	a
4	b

$A_1$	$A_2$	$A_3$	$A_4$
a	1	d	4
b	1	c	3
b	1	d	4
b	2	c	3

$A_1$	$A_2$	$A_3$
b	1	c
b	2	c

$A_2$	$A_3$	$A_4$
1	d	4

A minimum-sized model of size 2  
 $\Rightarrow$  Output size is 2

# Disjunctive Datalog: Output Size

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

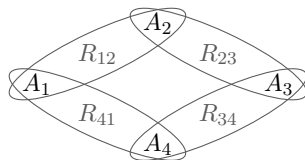
$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T} : \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$



# Disjunctive Datalog: Output Size

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

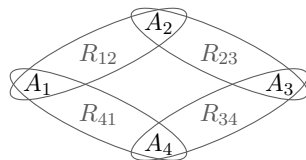
$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T}: \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$



# Disjunctive Datalog: Output Size

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T} : \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$

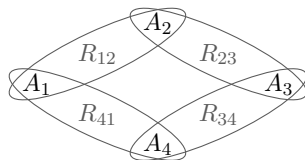


$$\mathbf{D} : |R_{12}| \leq N, \quad |R_{23}| \leq N, \quad |R_{34}| \leq N, \quad |R_{41}| \leq N.$$

# Disjunctive Datalog: Output Size

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T} : \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$



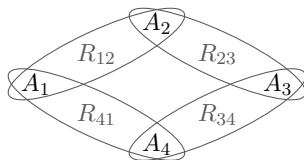
$$\mathbf{D} : |R_{12}| \leq N, \quad |R_{23}| \leq N, \quad |R_{34}| \leq N, \quad |R_{41}| \leq N.$$

- $P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :-$   
 $R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$

# Disjunctive Datalog: Output Size

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T} : \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$



$$\mathbf{D} : |R_{12}| \leq N, \quad |R_{23}| \leq N, \quad |R_{34}| \leq N, \quad |R_{41}| \leq N.$$

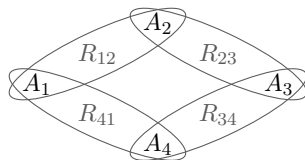
- $P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :-$   
 $R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$

$$|P(\mathbf{D})| \leq N^{3/2}, \text{ for all } \mathbf{D}$$

# Disjunctive Datalog: Output Size

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T} : \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$



$$\mathbf{D} : |R_{12}| \leq N, \quad |R_{23}| \leq N, \quad |R_{34}| \leq N, \quad |R_{41}| \leq N.$$

- ▶  $P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :-$   
 $R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$

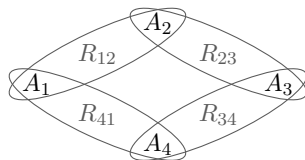
$$|P(\mathbf{D})| \leq N^{3/2}, \text{ for all } \mathbf{D}$$

- ▶  $P' : T_{123}(A_1, A_2, A_3) :-$   
 $R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$

# Disjunctive Datalog: Output Size

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T} : \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$



$$\mathbf{D} : |R_{12}| \leq N, \quad |R_{23}| \leq N, \quad |R_{34}| \leq N, \quad |R_{41}| \leq N.$$

- ▶  $P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :-$   
 $R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$

$$|P(\mathbf{D})| \leq N^{3/2}, \text{ for all } \mathbf{D}$$

- ▶  $P' : T_{123}(A_1, A_2, A_3) :-$   
 $R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$

$$|P'(\mathbf{D})| = N^2, \text{ for some } \mathbf{D}$$

# Disjunctive Datalog: Size Bounds

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T}: \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$

# Disjunctive Datalog: Size Bounds

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T}: \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$

Recall that:

- ▶ HDC is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  satisfying the degree constraints

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$



# Disjunctive Datalog: Size Bounds

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T}: \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$

Recall that:

- ▶ HDC is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  satisfying the degree constraints

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

## Theorem

$$\log |P(\mathbf{D})|$$

# Disjunctive Datalog: Size Bounds

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T}: \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$

Recall that:

- ▶ HDC is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  satisfying the degree constraints

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

## Theorem

$$\log |P(\mathbf{D})| \leq \underbrace{\max_{h \in \bar{\Gamma}_n^* \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)}_{\text{entropic bound}}$$

# Disjunctive Datalog: Size Bounds

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T}: \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$

Recall that:

- ▶ HDC is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  satisfying the degree constraints

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

## Theorem

$$\log |P(\mathbf{D})| \leq \underbrace{\max_{h \in \bar{\Gamma}_n^* \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)}_{\text{entropic bound}} \leq \underbrace{\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)}_{\text{polymatroid bound}}$$

# Disjunctive Datalog: Size Bounds

$$P : \bigvee_{B \in \mathcal{B}} T_B(\mathbf{A}_B) :- \bigwedge_{F \in \mathcal{E}} R_F(\mathbf{A}_F)$$

$$|P(\mathbf{D})| \stackrel{\text{def}}{=} \min_{\mathbf{T}: \mathbf{T} \models P} \max_{B \in \mathcal{B}} |T_B|$$

Recall that:

- ▶ HDC is the set of functions  $h : 2^{[n]} \rightarrow \mathbb{R}_+$  satisfying the degree constraints

$$\underbrace{\Gamma_n^*}_{\text{entropic functions}} \subset \underbrace{\bar{\Gamma}_n^*}_{\text{topological closure of } \Gamma_n^*} \subset \underbrace{\Gamma_n}_{\text{polymatroids}}$$

## Theorem

$$\log |P(\mathbf{D})| \leq \underbrace{\max_{h \in \bar{\Gamma}_n^* \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)}_{\substack{\text{entropic bound} \\ \text{(asymptotically tight!)}}} \leq \underbrace{\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)}_{\text{polymatroid bound}}$$

# Table of Contents

Size Bounds for Full Conjunctive Queries

Size Bounds for Disjunctive Datalog

**Algorithms for Disjunctive Datalog**

Algorithms for Conjunctive Queries

# PANDA (**P**roof-**A**ssisted **e**Ntropic **D**egree-**A**ware)

- ▶ An algorithm for disjunctive datalog

# PANDA (**P**roof-**A**ssisted **e**Ntropic **D**egree-**A**ware)

- ▶ An algorithm for **disjunctive datalog**
  - ▶ computes *a model*

# PANDA (Proof-Assisted eNtropic Degree-Aware)

- ▶ An algorithm for disjunctive datalog
  - ▶ computes *a model*
  - ▶ within the polymatroid bound:



# PANDA (Proof-Assisted eNtropic Degree-Aware)

- ▶ An algorithm for **disjunctive datalog**
  - ▶ computes *a model*
  - ▶ within the **polymatroid bound**:
    - ▶ the worst-case size of the *minimum model*.

# PANDA (Proof-Assisted eNtropic Degree-Aware)

- ▶ An algorithm for **disjunctive datalog**
  - ▶ computes *a model*
  - ▶ within the **polymatroid bound**:
    - ▶ the worst-case size of the *minimum model*.
- ▶ Outline

# PANDA (**P**roof-**A**ssisted **e**Ntropic **D**egree-**A**ware)

- ▶ An algorithm for **disjunctive datalog**
  - ▶ computes *a model*
  - ▶ within the **polymatroid bound**:
    - ▶ the worst-case size of the *minimum model*.
- ▶ Outline
  - ▶ Construct a **Proof Sequence** for the bound.

# PANDA (**P**roof-**A**ssisted **e**Ntropic **D**egree-**A**ware)

- ▶ An algorithm for **disjunctive datalog**
  - ▶ computes *a model*
  - ▶ within the **polymatroid bound**:
    - ▶ the worst-case size of the *minimum model*.
- ▶ Outline
  - ▶ Construct a **Proof Sequence** for the bound.
  - ▶ Interpret each proof step as an algorithmic step.

# PANDA

- ▶ Polymatroid bound:  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)$

# PANDA

- ▶ Polymatroid bound:  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)$
- ▶  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B) = \max_{h \in \Gamma_n \cap \text{HDC}} \sum_{B \in \mathcal{B}} \lambda_B h(B)$

# PANDA

- ▶ Polymatroid bound:  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)$
- ▶  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B) = \max_{h \in \Gamma_n \cap \text{HDC}} \sum_{B \in \mathcal{B}} \lambda_B h(B)$
- ▶  $\sum_{B \in \mathcal{B}} \lambda_B \cdot h(B) \leq \sum_{(X, Y, N_{Y|X})} \delta_{Y|X} \cdot \underbrace{h(Y|X)}_{\substack{|\wedge \\ \log N_{Y|X}}}$

# PANDA

- ▶ Polymatroid bound:  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)$
- ▶  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B) = \max_{h \in \Gamma_n \cap \text{HDC}} \sum_{B \in \mathcal{B}} \lambda_B h(B)$
- ▶  $\sum_{B \in \mathcal{B}} \lambda_B \cdot h(B) \leq \sum_{(X, Y, N_{Y|X})} \delta_{Y|X} \cdot \underbrace{h(Y|X)}_{\substack{|\wedge \\ \log N_{Y|X}}}$

## ▶ Proof Sequence

Given  $X \subseteq Y$ :

$$h(X) + h(Y|X) \rightarrow h(Y)$$

$$h(Y) \rightarrow h(X) + h(Y|X)$$

$$h(Y) \rightarrow h(X)$$

$$h(Y|X) \rightarrow h(Y \cup Z|X \cup Z)$$



# PANDA

- ▶ Polymatroid bound:  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)$
- ▶  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B) = \max_{h \in \Gamma_n \cap \text{HDC}} \sum_{B \in \mathcal{B}} \lambda_B h(B)$
- ▶  $\sum_{B \in \mathcal{B}} \lambda_B \cdot h(B) \leq \sum_{(X, Y, N_{Y|X})} \delta_{Y|X} \cdot \underbrace{h(Y|X)}_{\substack{|\wedge \\ \log N_{Y|X}}}$

- ▶ Proof Sequence

Given  $X \subseteq Y$ :

$$h(X) + h(Y|X) \rightarrow h(Y) \quad (\text{join})$$

$$h(Y) \rightarrow h(X) + h(Y|X) \quad (\text{data partition})$$

$$h(Y) \rightarrow h(X) \quad (\text{projection})$$

$$h(Y|X) \rightarrow h(Y \cup Z|X \cup Z) \quad (\text{nothing})$$

- ▶ Algorithmic Sequence

# PANDA

- ▶ Polymatroid bound:  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B)$
- ▶  $\max_{h \in \Gamma_n \cap \text{HDC}} \min_{B \in \mathcal{B}} h(B) = \max_{h \in \Gamma_n \cap \text{HDC}} \sum_{B \in \mathcal{B}} \lambda_B h(B)$
- ▶  $\sum_{B \in \mathcal{B}} \lambda_B \cdot h(B) \leq \sum_{(X, Y, N_{Y|X})} \delta_{Y|X} \cdot \underbrace{h(Y|X)}_{\substack{|\wedge \\ \log N_{Y|X}}}$

- ▶ Proof Sequence

Given  $X \subseteq Y$ :

$$h(X) + h(Y|X) \rightarrow h(Y) \quad (\text{join})$$

$$h(Y) \rightarrow h(X) + h(Y|X) \quad (\text{data partition})$$

$$h(Y) \rightarrow h(X) \quad (\text{projection})$$

$$h(Y|X) \rightarrow h(Y \cup Z|X \cup Z) \quad (\text{nothing})$$

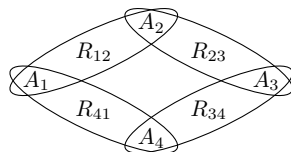
- ▶ Algorithmic Sequence

## Theorem

*PANDA solves any disjunctive datalog rule  $P$  in time within the polymatroid bound of  $P$ .*

# PANDA: Example

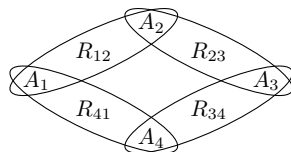
$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



# PANDA: Example

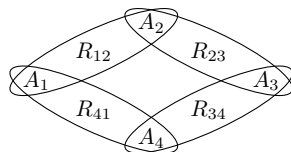
$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N$$



# PANDA: Example

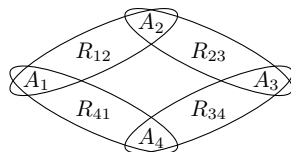
$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

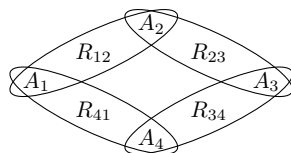


$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4))$$

# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

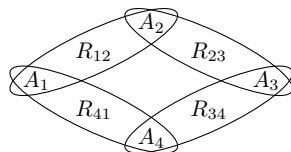


$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4))$$

# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



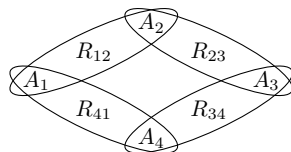
$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4))$$



# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

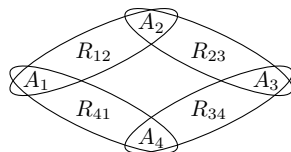


$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



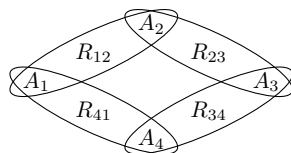
$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



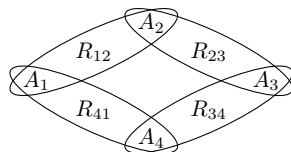
$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3))$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



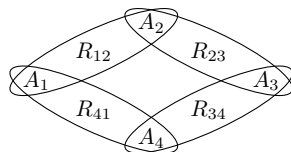
$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4|A_3) + h(A_3)) \\ h(A_1 A_2) + h(A_2 A_3) + h(A_4|A_3) + h(A_3)$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



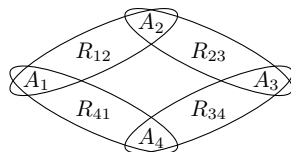
$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4|A_3) + h(A_3)) \\ h(A_1 A_2) + h(A_2 A_3) + h(A_4|A_3) + h(A_3) \rightarrow (h(A_4|A_3) \rightarrow h(A_4|A_2 A_3))$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

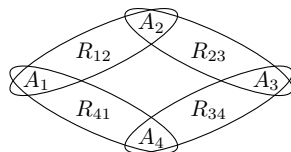
$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4|A_3) + h(A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4|A_3) + h(A_3) \rightarrow (h(A_4|A_3) \rightarrow h(A_4|A_2 A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4|A_2 A_3) + h(A_3)$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

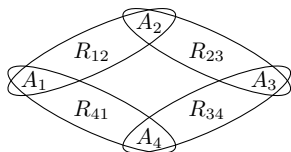
$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4 | A_3) + h(A_3) \rightarrow (h(A_4 | A_3) \rightarrow h(A_4 | A_2 A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4 | A_2 A_3) + h(A_3) \rightarrow (h(A_2 A_3) + h(A_4 | A_2 A_3) \rightarrow h(A_2 A_3 A_4))$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4 | A_3) + h(A_3) \rightarrow (h(A_4 | A_3) \rightarrow h(A_4 | A_2 A_3))$$

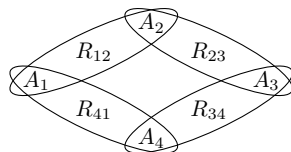
$$h(A_1 A_2) + h(A_2 A_3) + h(A_4 | A_2 A_3) + h(A_3) \rightarrow (h(A_2 A_3) + h(A_4 | A_2 A_3) \rightarrow h(A_2 A_3 A_4))$$

$$h(A_1 A_2) + h(A_2 A_3 A_4) + h(A_3)$$



# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3))$$

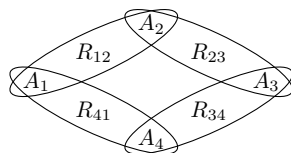
$$h(A_1 A_2) + h(A_2 A_3) + h(A_4 | A_3) + h(A_3) \rightarrow (h(A_4 | A_3) \rightarrow h(A_4 | A_2 A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4 | A_2 A_3) + h(A_3) \rightarrow (h(A_2 A_3) + h(A_4 | A_2 A_3) \rightarrow h(A_2 A_3 A_4))$$

$$h(A_1 A_2) + h(A_2 A_3 A_4) + h(A_3) \rightarrow (h(A_1 A_2) \rightarrow h(A_1 A_2 | A_3))$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4 | A_3) + h(A_3) \rightarrow (h(A_4 | A_3) \rightarrow h(A_4 | A_2 A_3))$$

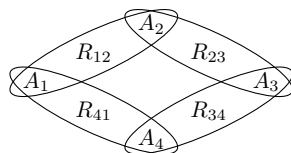
$$h(A_1 A_2) + h(A_2 A_3) + h(A_4 | A_2 A_3) + h(A_3) \rightarrow (h(A_2 A_3) + h(A_4 | A_2 A_3) \rightarrow h(A_2 A_3 A_4))$$

$$h(A_1 A_2) + h(A_2 A_3 A_4) + h(A_3) \rightarrow (h(A_1 A_2) \rightarrow h(A_1 A_2 | A_3))$$

$$h(A_1 A_2 | A_3) + h(A_2 A_3 A_4) + h(A_3)$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4|A_3) + h(A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4|A_3) + h(A_3) \rightarrow (h(A_4|A_3) \rightarrow h(A_4|A_2 A_3))$$

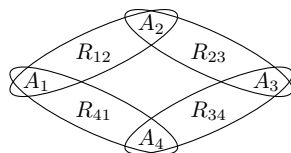
$$h(A_1 A_2) + h(A_2 A_3) + h(A_4|A_2 A_3) + h(A_3) \rightarrow (h(A_2 A_3) + h(A_4|A_2 A_3) \rightarrow h(A_2 A_3 A_4))$$

$$h(A_1 A_2) + h(A_2 A_3 A_4) + h(A_3) \rightarrow (h(A_1 A_2) \rightarrow h(A_1 A_2|A_3))$$

$$h(A_1 A_2|A_3) + h(A_2 A_3 A_4) + h(A_3) \rightarrow (h(A_1 A_2|A_3) + h(A_3) \rightarrow h(A_1 A_2 A_3))$$

# PANDA: Example

$$P : T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$\log |P| \leq \min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq \frac{1}{2} (h(A_1 A_2 A_3) + h(A_2 A_3 A_4)) \\ \leq \frac{1}{2} (h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4)) \leq \frac{3}{2} \log N$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_3 A_4) \rightarrow (h(A_3 A_4) \rightarrow h(A_4|A_3) + h(A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4|A_3) + h(A_3) \rightarrow (h(A_4|A_3) \rightarrow h(A_4|A_2 A_3))$$

$$h(A_1 A_2) + h(A_2 A_3) + h(A_4|A_2 A_3) + h(A_3) \rightarrow (h(A_2 A_3) + h(A_4|A_2 A_3) \rightarrow h(A_2 A_3 A_4))$$

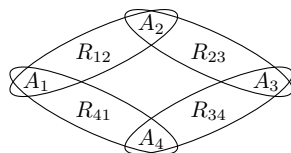
$$h(A_1 A_2) + h(A_2 A_3 A_4) + h(A_3) \rightarrow (h(A_1 A_2) \rightarrow h(A_1 A_2|A_3))$$

$$h(A_1 A_2|A_3) + h(A_2 A_3 A_4) + h(A_3) \rightarrow (h(A_1 A_2|A_3) + h(A_3) \rightarrow h(A_1 A_2 A_3))$$

$$h(A_1 A_2 A_3) + h(A_2 A_3 A_4)$$

# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3)$$

$$h(A_4 | A_3) \rightarrow h(A_4 | A_2 A_3)$$

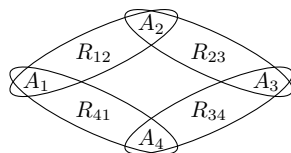
$$h(A_2 A_3) + h(A_4 | A_2 A_3) \rightarrow h(A_2 A_3 A_4)$$

$$h(A_1 A_2) \rightarrow h(A_1 A_2 | A_3)$$

$$h(A_1 A_2 | A_3) + h(A_3) \rightarrow h(A_1 A_2 A_3)$$

# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3)$$

$$R_{34}(A_3, A_4) \rightarrow R_{34}^{(l)}(A_3, A_4), R_3^{(h)}(A_3)$$

$$h(A_4 | A_3) \rightarrow h(A_4 | A_2 A_3)$$

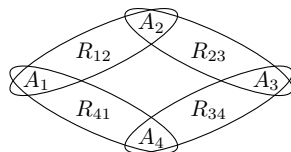
$$h(A_2 A_3) + h(A_4 | A_2 A_3) \rightarrow h(A_2 A_3 A_4)$$

$$h(A_1 A_2) \rightarrow h(A_1 A_2 | A_3)$$

$$h(A_1 A_2 | A_3) + h(A_3) \rightarrow h(A_1 A_2 A_3)$$

# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3)$$

$$R_{34}(A_3, A_4) \rightarrow R_{34}^{(l)}(A_3, A_4), R_3^{(h)}(A_3)$$

$$h(A_4 | A_3) \rightarrow h(A_4 | A_2 A_3)$$

$$R_{34}^{(l)}(A_3, A_4) \rightarrow R_{34}^{(l)}(A_3, A_4)$$

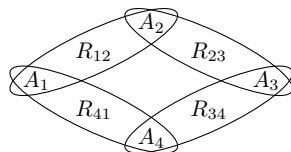
$$h(A_2 A_3) + h(A_4 | A_2 A_3) \rightarrow h(A_2 A_3 A_4)$$

$$h(A_1 A_2) \rightarrow h(A_1 A_2 | A_3)$$

$$h(A_1 A_2 | A_3) + h(A_3) \rightarrow h(A_1 A_2 A_3)$$

# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3)$$

$$h(A_4 | A_3) \rightarrow h(A_4 | A_2 A_3)$$

$$h(A_2 A_3) + h(A_4 | A_2 A_3) \rightarrow h(A_2 A_3 A_4)$$

$$h(A_1 A_2) \rightarrow h(A_1 A_2 | A_3)$$

$$h(A_1 A_2 | A_3) + h(A_3) \rightarrow h(A_1 A_2 A_3)$$

$$R_{34}(A_3, A_4) \rightarrow R_{34}^{(l)}(A_3, A_4), R_3^{(h)}(A_3)$$

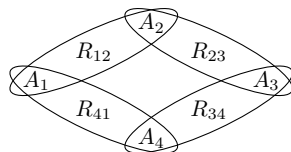
$$R_{34}^{(l)}(A_3, A_4) \rightarrow R_{34}^{(l)}(A_3, A_4)$$

$$R_{23}(A_2, A_3) \bowtie R_{34}^{(l)}(A_3, A_4) \rightarrow T_{234}(A_2, A_3, A_4)$$



# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$h(A_3 A_4) \rightarrow h(A_4 | A_3) + h(A_3)$$

$$h(A_4 | A_3) \rightarrow h(A_4 | A_2 A_3)$$

$$h(A_2 A_3) + h(A_4 | A_2 A_3) \rightarrow h(A_2 A_3 A_4)$$

$$h(A_1 A_2) \rightarrow h(A_1 A_2 | A_3)$$

$$h(A_1 A_2 | A_3) + h(A_3) \rightarrow h(A_1 A_2 A_3)$$

$$R_{34}(A_3, A_4) \rightarrow R_{34}^{(l)}(A_3, A_4), R_3^{(h)}(A_3)$$

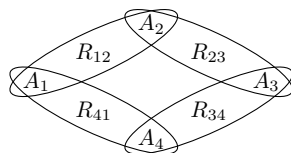
$$R_{34}^{(l)}(A_3, A_4) \rightarrow R_{34}^{(l)}(A_3, A_4)$$

$$R_{23}(A_2, A_3) \bowtie R_{34}^{(l)}(A_3, A_4) \rightarrow T_{234}(A_2, A_3, A_4)$$

$$R_{12}(A_1, A_2) \rightarrow R_{12}(A_1, A_2)$$

# PANDA: Example

$$P : \quad T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :- \\ R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$|R_{12}|, |R_{23}|, |R_{34}|, |R_{41}| \leq N \quad \Rightarrow \quad |P| \leq N^{3/2}$$

$$h(A_3 A_4) \rightarrow h(A_4|A_3) + h(A_3)$$

$$h(A_4|A_3) \rightarrow h(A_4|A_2 A_3)$$

$$h(A_2 A_3) + h(A_4|A_2 A_3) \rightarrow h(A_2 A_3 A_4)$$

$$h(A_1 A_2) \rightarrow h(A_1 A_2|A_3)$$

$$h(A_1 A_2|A_3) + h(A_3) \rightarrow h(A_1 A_2 A_3)$$

$$R_{34}(A_3, A_4) \rightarrow R_{34}^{(l)}(A_3, A_4), R_3^{(h)}(A_3)$$

$$R_{34}^{(l)}(A_3, A_4) \rightarrow R_{34}^{(l)}(A_3, A_4)$$

$$R_{23}(A_2, A_3) \bowtie R_{34}^{(l)}(A_3, A_4) \rightarrow T_{234}(A_2, A_3, A_4)$$

$$R_{12}(A_1, A_2) \rightarrow R_{12}(A_1, A_2)$$

$$R_{12}(A_1, A_2) \bowtie R_3^{(h)}(A_3) \rightarrow T_{123}(A_1, A_2, A_3)$$

# Table of Contents

Size Bounds for Full Conjunctive Queries

Size Bounds for Disjunctive Datalog

Algorithms for Disjunctive Datalog

Algorithms for Conjunctive Queries

# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

$$\tilde{O}\left( N^d + |\text{output}| \right)$$

# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

$$\tilde{O}\left( \underbrace{N^d}_{\text{Intrinsic Cost}} + |\text{output}| \right)$$

# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

$$\tilde{O}\left( \begin{array}{c} N^d \\ \text{Intrinsic Cost} \end{array} + \begin{array}{c} |\text{output}| \\ \text{Output Cost} \end{array} \right)$$

# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

$$\tilde{O}\left( \underbrace{N^d}_{\text{Intrinsic Cost}} + \underbrace{|\text{output}|}_{\text{Output Cost}} \right)$$

- ▶ **Submodular width** as a candidate for  $d$



# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

$$\tilde{O}\left( \underbrace{N^d}_{\text{Intrinsic Cost}} + \underbrace{|\text{output}|}_{\text{Output Cost}} \right)$$

- ▶ **Submodular width** as a candidate for  $d$  [Marx JACM'13]

# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

$$\tilde{O}\left( \underbrace{N^d}_{\text{Intrinsic Cost}} + \underbrace{|\text{output}|}_{\text{Output Cost}} \right)$$

- ▶ **Submodular width** as a candidate for  $d$  [Marx JACM'13]

$$\mathbf{FPT} \iff \text{Bounded subw}(Q)$$

# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

$$\tilde{O}\left( \underbrace{N^d}_{\text{Intrinsic Cost}} + \underbrace{|\text{output}|}_{\text{Output Cost}} \right)$$

- ▶ **Submodular width** as a candidate for  $d$  [Marx JACM'13]

**FPT**  $\Leftrightarrow$  Bounded  $\text{subw}(Q)$

$$\text{Boolean } Q \Rightarrow \tilde{O}\left(N^{\text{subw}(Q)} \times c\right)$$

# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

$$\tilde{O}\left( \underbrace{N^d}_{\text{Intrinsic Cost}} + \underbrace{|\text{output}|}_{\text{Output Cost}} \right)$$

- ▶ **Submodular width** as a candidate for  $d$  [Marx JACM'13]

**FPT**  $\Leftrightarrow$  Bounded  $\text{subw}(Q)$

$$\text{Boolean } Q \Rightarrow \tilde{O}\left(N^{\text{subw}(Q)} \times c\right)$$

- ▶ **Our goals**

# Beyond Worst-case Optimality

- ▶ **Output-sensitive** algorithms

$$\tilde{O}\left( \underbrace{N^d}_{\text{Intrinsic Cost}} + \underbrace{|\text{output}|}_{\text{Output Cost}} \right)$$

- ▶ **Submodular width** as a candidate for  $d$  [Marx JACM'13]

**FPT**  $\Leftrightarrow$  Bounded  $\text{subw}(Q)$

$$\text{Boolean } Q \Rightarrow \tilde{O}\left(N^{\text{subw}(Q)} \times c\right)$$

- ▶ **Our goals**

$$\text{Any } Q \Rightarrow \tilde{O}\left(N^{\text{da-subw}(Q)} \times 1 + |\text{output}|\right)$$

# Submodular Width

$$\text{fhtw}(Q) \stackrel{\text{def}}{=} \min_{(T, \chi)} \max_{t \in V(T)} \rho^*(\chi(t))$$

# Submodular Width

$$\begin{aligned} \text{fhtw}(Q) &\stackrel{\text{def}}{=} \min_{(T,\chi)} \max_{t \in V(T)} \rho^*(\chi(t)) \\ &= \min_{(T,\chi)} \max_{t \in V(T)} \max_{h \in \text{ED} \cap \Gamma_n} h(\chi(t)) \end{aligned}$$

# Submodular Width

$$\begin{aligned} \text{fhtw}(Q) &\stackrel{\text{def}}{=} \min_{(T,\chi)} \max_{t \in V(T)} \rho^*(\chi(t)) \\ &= \min_{(T,\chi)} \max_{t \in V(T)} \max_{h \in \text{ED} \cap \Gamma_n} h(\chi(t)) \\ &= \min_{(T,\chi)} \max_{h \in \text{ED} \cap \Gamma_n} \max_{t \in V(T)} h(\chi(t)) \end{aligned}$$



# Submodular Width

$$\begin{aligned} \text{fhtw}(Q) &\stackrel{\text{def}}{=} \min_{(T,\chi)} \max_{t \in V(T)} \rho^*(\chi(t)) \\ &= \min_{(T,\chi)} \max_{t \in V(T)} \max_{h \in \text{ED} \cap \Gamma_n} h(\chi(t)) \\ &= \min_{(T,\chi)} \max_{h \in \text{ED} \cap \Gamma_n} \max_{t \in V(T)} h(\chi(t)) \\ \text{subw}(Q) &\stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T,\chi)} \max_{t \in V(T)} h(\chi(t)) \end{aligned}$$

# Submodular Width

$$\begin{aligned} \text{fhtw}(Q) &\stackrel{\text{def}}{=} \min_{(T,\chi)} \max_{t \in V(T)} \rho^*(\chi(t)) \\ &= \min_{(T,\chi)} \max_{t \in V(T)} \max_{h \in \text{ED} \cap \Gamma_n} h(\chi(t)) \\ &= \min_{(T,\chi)} \max_{h \in \text{ED} \cap \Gamma_n} \max_{t \in V(T)} h(\chi(t)) \\ \text{subw}(Q) &\stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T,\chi)} \max_{t \in V(T)} h(\chi(t)) \end{aligned}$$

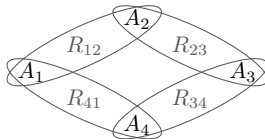
$$\text{subw}(Q) \leq \text{fhtw}(Q)$$

# Submodular Width: Example

$$\text{fhtw}(Q) \stackrel{\text{def}}{=} \min_{(T, \chi)} \max_{t \in V(T)} \rho^*(\chi(t)),$$

$$\text{subw}(Q) \stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T, \chi)} \max_{t \in V(T)} h(\chi(t))$$

$$Q(A_1, A_2, A_3, A_4) :- \quad R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

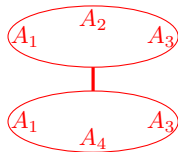
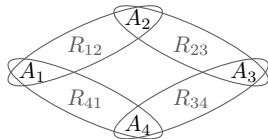


# Submodular Width: Example

$$\text{fhtw}(Q) \stackrel{\text{def}}{=} \min_{(T, \chi)} \max_{t \in V(T)} \rho^*(\chi(t)),$$

$$\text{subw}(Q) \stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T, \chi)} \max_{t \in V(T)} h(\chi(t))$$

$$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

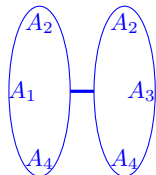
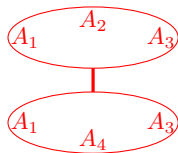
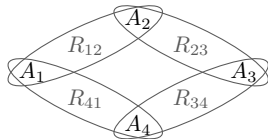


# Submodular Width: Example

$$\text{fhtw}(Q) \stackrel{\text{def}}{=} \min_{(T, \chi)} \max_{t \in V(T)} \rho^*(\chi(t)),$$

$$\text{subw}(Q) \stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T, \chi)} \max_{t \in V(T)} h(\chi(t))$$

$$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



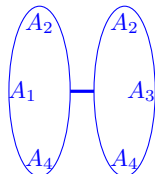
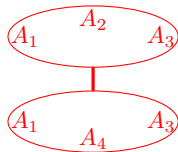
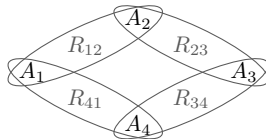
# Submodular Width: Example

$$\text{fhtw}(Q) \stackrel{\text{def}}{=} \min_{(T, \chi)} \max_{t \in V(T)} \rho^*(\chi(t)),$$

$$\text{subw}(Q) \stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T, \chi)} \max_{t \in V(T)} h(\chi(t))$$

$$Q(A_1, A_2, A_3, A_4) :- \quad R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$$\text{fhtw}(Q) = 2$$



# Submodular Width: Example

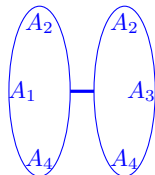
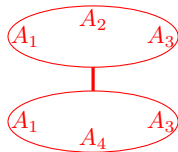
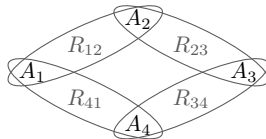
$$\text{fhtw}(Q) \stackrel{\text{def}}{=} \min_{(T, \chi)} \max_{t \in V(T)} \rho^*(\chi(t)),$$

$$\text{subw}(Q) \stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T, \chi)} \max_{t \in V(T)} h(\chi(t))$$

$$Q(A_1, A_2, A_3, A_4) :- \quad R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$$\text{fhtw}(Q) = 2$$

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



# Submodular Width: Example

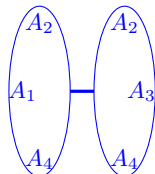
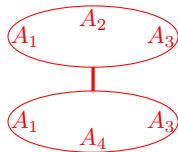
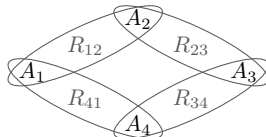
$$\text{fhtw}(Q) \stackrel{\text{def}}{=} \min_{(T, \chi)} \max_{t \in V(T)} \rho^*(\chi(t)),$$

$$\text{subw}(Q) \stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T, \chi)} \max_{t \in V(T)} h(\chi(t))$$

$$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$$\begin{aligned} \text{fhtw}(Q) &= 2 \\ \text{subw}(Q) &= 3/2 \end{aligned}$$

$$\begin{aligned} \min(\max(h(A_1 A_2 A_3), h(A_3 A_4 A_1)), \\ \max(h(A_4 A_1 A_2), h(A_2 A_3 A_4))) \leq 3/2 \end{aligned}$$





# Submodular Width: Example

$$\text{fhtw}(Q) \stackrel{\text{def}}{=} \min_{(T, \chi)} \max_{t \in V(T)} \rho^*(\chi(t)),$$

$$\text{subw}(Q) \stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T, \chi)} \max_{t \in V(T)} h(\chi(t))$$

$$Q(A_1, A_2, A_3, A_4) :- \quad R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$$\text{fhtw}(Q) = 2$$

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

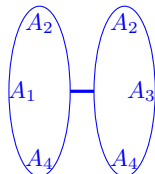
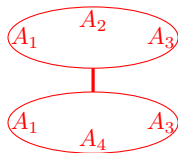
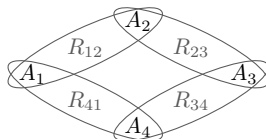
$$\min(\max(h(A_1 A_2 A_3), h(A_3 A_4 A_1)), \\ \max(h(A_4 A_1 A_2), h(A_2 A_3 A_4))) \leq 3/2$$

$$\min(h(A_1 A_2 A_3), h(A_4 A_1 A_2)) \leq 3/2$$

$$\min(h(A_1 A_2 A_3), h(A_2 A_3 A_4)) \leq 3/2$$

$$\min(h(A_3 A_4 A_1), h(A_4 A_1 A_2)) \leq 3/2$$

$$\min(h(A_3 A_4 A_1), h(A_2 A_3 A_4)) \leq 3/2$$

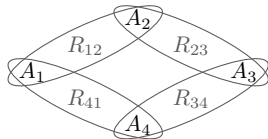


# Submodular Width: Example

$$\text{fhtw}(Q) \stackrel{\text{def}}{=} \min_{(T, \chi)} \max_{t \in V(T)} \rho^*(\chi(t)),$$

$$\text{subw}(Q) \stackrel{\text{def}}{=} \max_{h \in \text{ED} \cap \Gamma_n} \min_{(T, \chi)} \max_{t \in V(T)} h(\chi(t))$$

$$Q(A_1, A_2, A_3, A_4) :- R_{12}(A_1, A_2), R_{23}(A_2, A_3), \\ R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



$$T_{123}(A_1, A_2, A_3) \vee T_{412}(A_4, A_1, A_2) :-$$

$$R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$$T_{123}(A_1, A_2, A_3) \vee T_{234}(A_2, A_3, A_4) :-$$

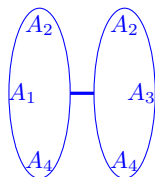
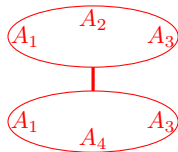
$$R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$$T_{341}(A_3, A_4, A_1) \vee T_{412}(A_4, A_1, A_2) :-$$

$$R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$

$$T_{341}(A_3, A_4, A_1) \vee T_{234}(A_2, A_3, A_4) :-$$

$$R_{12}(A_1, A_2), R_{23}(A_2, A_3), R_{34}(A_3, A_4), R_{41}(A_4, A_1).$$



# Summary of Bounds

