

# T3: Efficient query evaluation

## L15: Cyclic query evaluation

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp20)

<https://northeastern-datalab.github.io/cs7240/sp20/>

Version 2/25/2020

# Outline: Acyclic conjunctive queries

- Acyclic conjunctive queries
  - The semijoin operator
  - Join trees & Yannakakis algorithm
  - Query hypergraphs & GYO reduction
  - A detailed Yannakakis example
  - Full semijoin reductions
- Cyclic conjunctive queries

# Semijoin Reducer



$$Q(x,y,z) = R(x, y) \bowtie S(y, z) \bowtie T(z, w)$$

A full reducer is

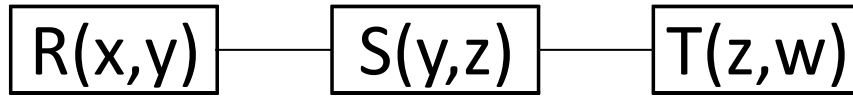
?

# Semijoin Reducer



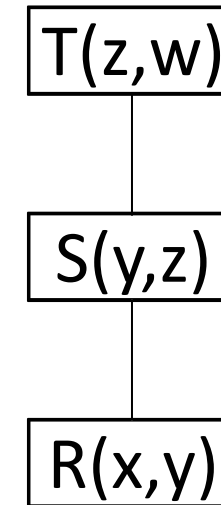
$$Q(x,y,z) = R(x,y) \bowtie S(y,z) \bowtie T(z,w)$$

A full reducer is



?

Join Tree for Q

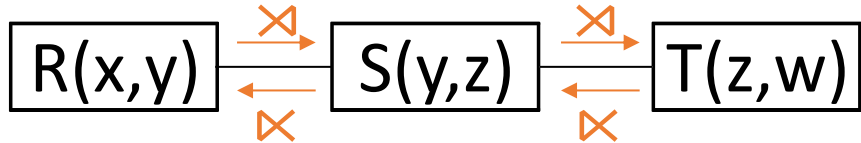


# Semijoin Reducer



$$Q(x,y,z) = R(x,y) \bowtie S(y,z) \bowtie T(z,w)$$

A full reducer is

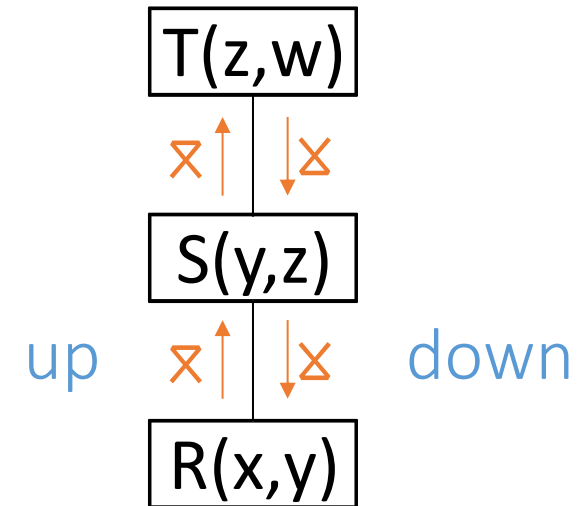


$$\begin{aligned} S_1(y,z) &= S(y,z) \bowtie R(x,y) \\ T_1(z,y) &= T(z,y) \bowtie S_1(y,z) \\ S_2(z,y) &= S_1(y,z) \bowtie T_1(z,y) \\ R_1(x,y) &= R(x,y) \bowtie S_2(y,z) \end{aligned}$$

The rewritten query is

?

Join Tree for Q

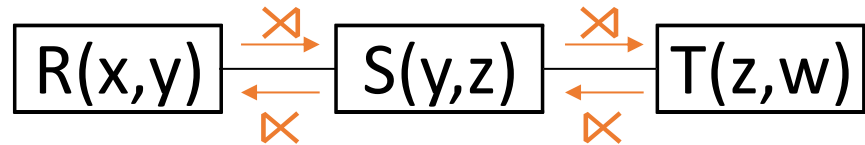


# Semijoin Reducer



$$Q(x,y,z) = R(x,y) \bowtie S(y,z) \bowtie T(z,w)$$

A full reducer is

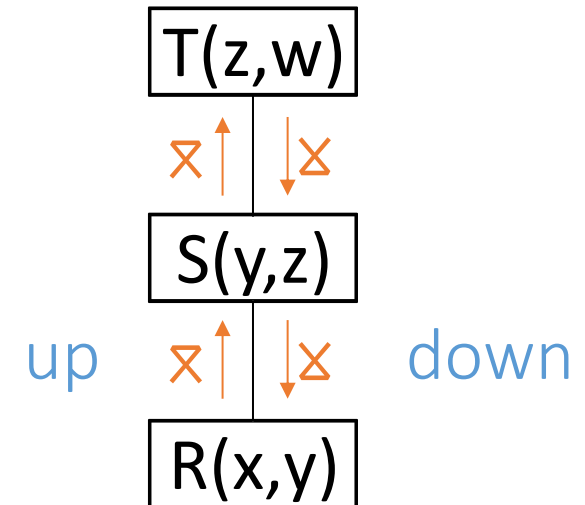


$$\begin{aligned} S_1(y,z) &= S(y,z) \bowtie R(x,y) \\ T_1(z,y) &= T(z,y) \bowtie S_1(y,z) \\ S_2(z,y) &= S_1(y,z) \bowtie T_1(z,y) \\ R_1(x,y) &= R(x,y) \bowtie S_2(y,z) \end{aligned}$$

The rewritten query is

$$Q(x,y,z) = R_1(x,y) \bowtie S_2(y,z) \bowtie T_1(z,w)$$

Join Tree for Q



# Semi-join reducers

## GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

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



Join tree

?

Query hypergraph

?

# Semi-join reducers

## GYO ear removal

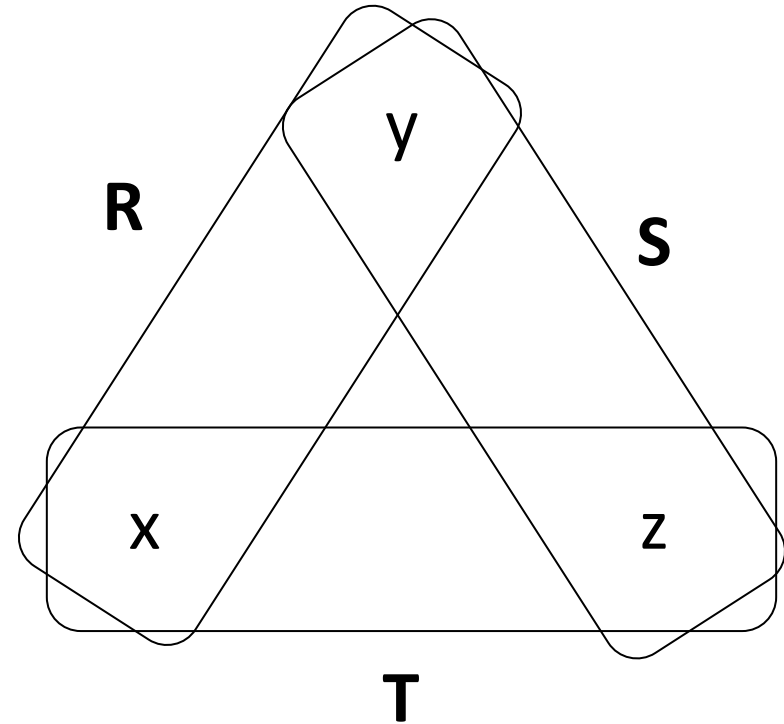
- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

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



Join tree

?



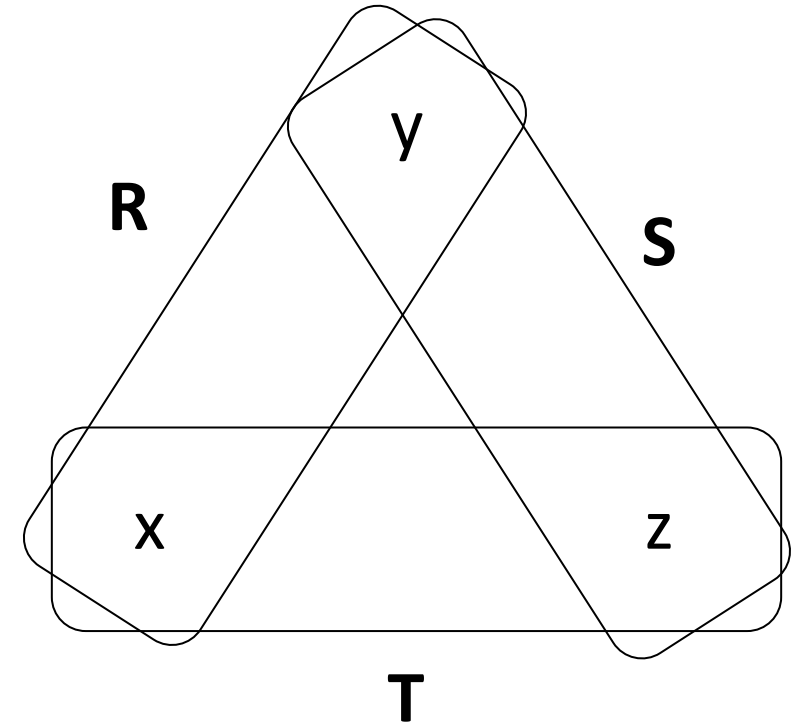
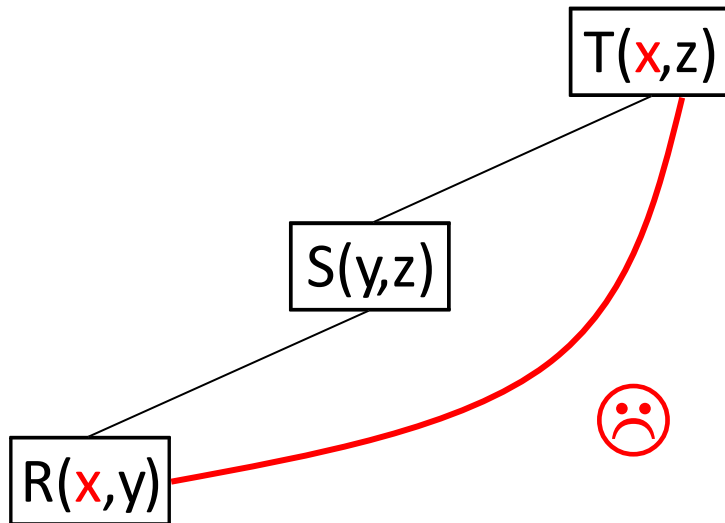


# Semi-join reducers

## GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

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



# Semi-join reducers

## GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

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



Join tree

?

Query hypergraph

?

# Semi-join reducers

## GYO ear removal

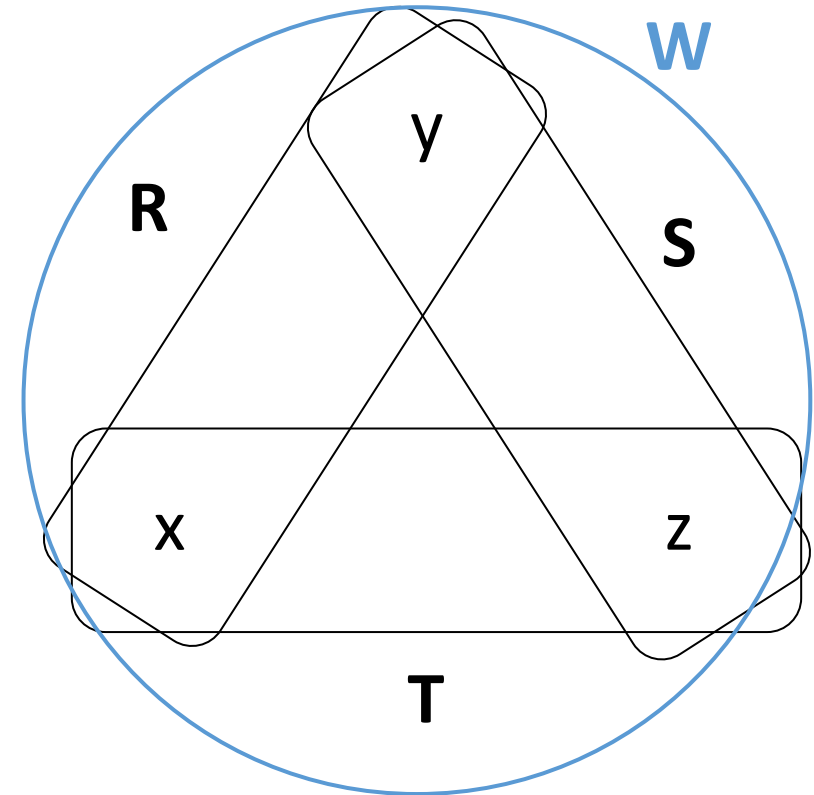
- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

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



Join tree

?

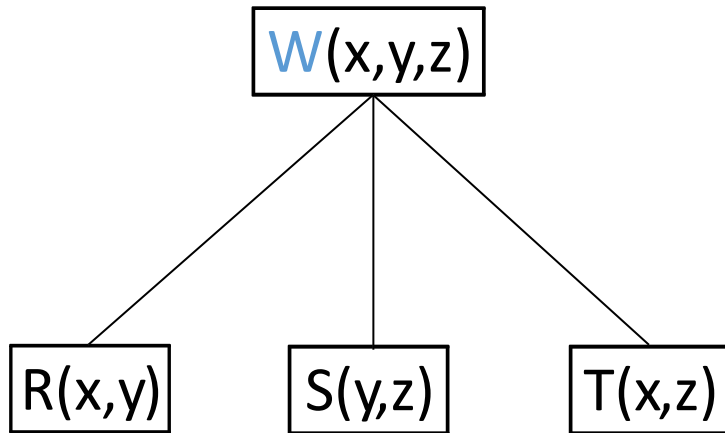


# Semi-join reducers

## GYO ear removal

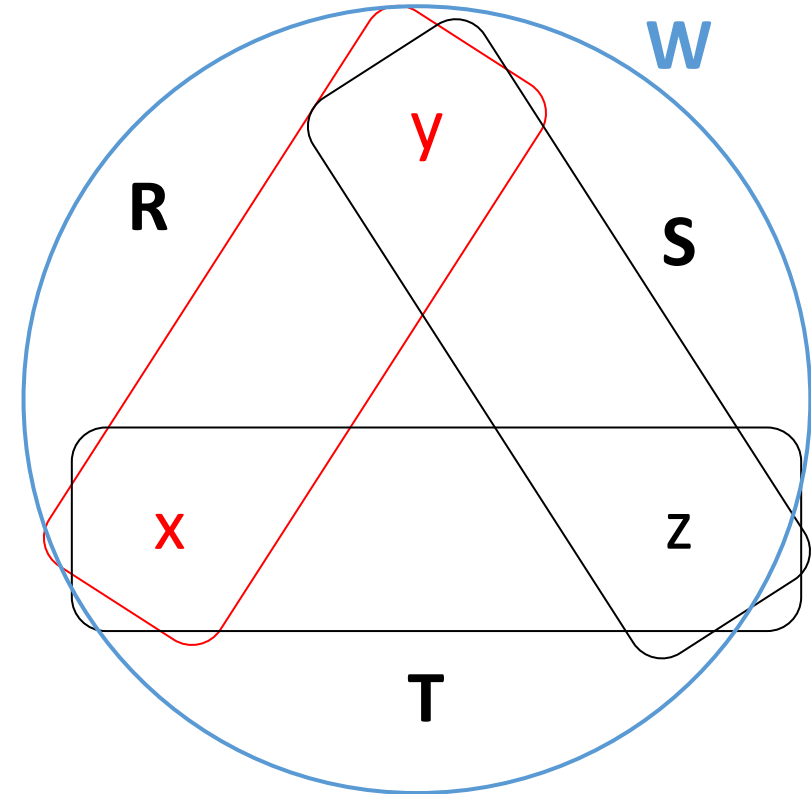
- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

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



Full reducer

?

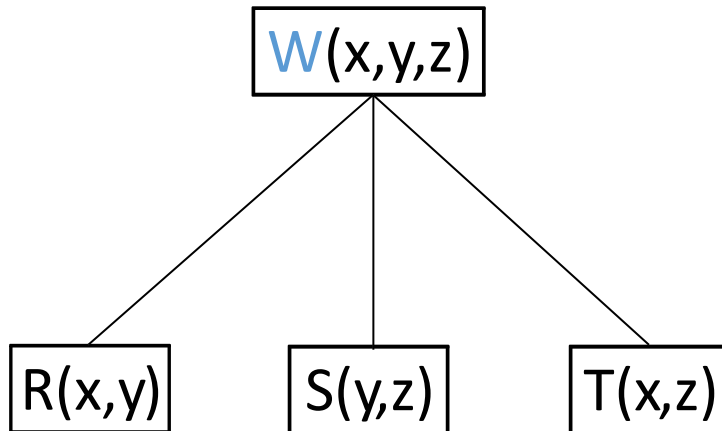


# Semi-join reducers

## GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

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



$$W_1(x, y, z) = W(x, y, z) \bowtie R(x, y)$$

$$W_2(x, y, z) = W_1(x, y, z) \bowtie S(y, z)$$

$$W_3(x, y, z) = W_2(x, y, z) \bowtie T(x, z)$$

$$R_1(x, y) = R(x, y) \bowtie W_3(x, y, z)$$

$$S_1(y, z) = S(y, z) \bowtie W_3(x, y, z)$$

$$T_1(x, z) = T(x, z) \bowtie W_3(x, y, z)$$

$$Q(x,y,z) = R_1(x, y) \bowtie S_1(y, z) \bowtie T_1(x, z) \bowtie W_3(x, y, z)$$

?

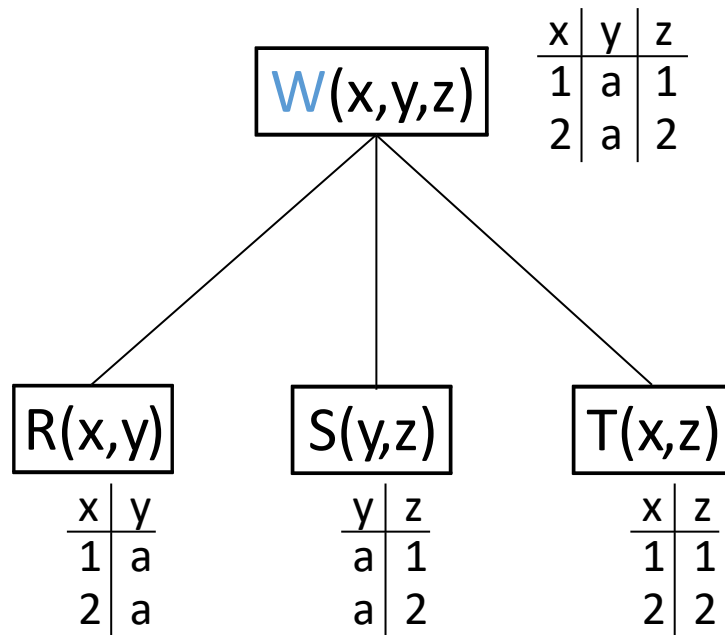
# Semi-join reducers

## GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)



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



$$W_1(x, y, z) = W(x, y, z) \bowtie R(x, y)$$

$$W_2(x, y, z) = W_1(x, y, z) \bowtie S(y, z)$$

$$W_3(x, y, z) = W_2(x, y, z) \bowtie T(x, z)$$

$$R_1(x, y) = R(x, y) \bowtie W_3(x, y, z)$$

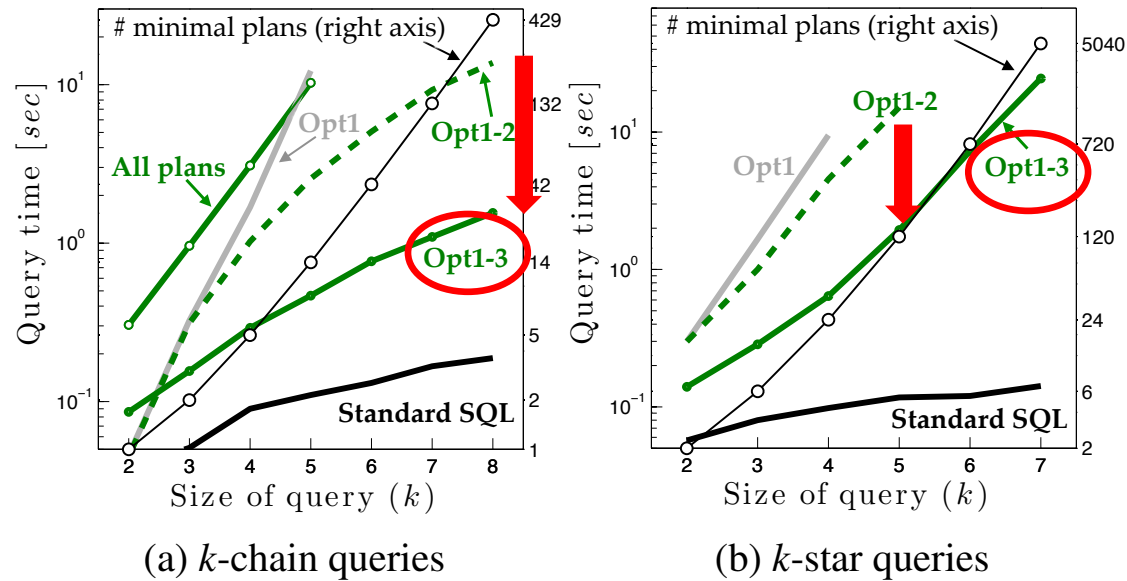
$$S_1(y, z) = S(y, z) \bowtie W_3(x, y, z)$$

$$T_1(x, z) = T(x, z) \bowtie W_3(x, y, z)$$

~~$$Q(x,y,z) = R_1(x, y) \bowtie S_1(y, z) \bowtie T_1(x, z) \bowtie W_3(x, y, z)$$~~

$$Q(x,y,z) = R_1(x, y) \bowtie W_3(x, y, z) \bowtie S_1(y, z) \bowtie T_1(x, z)$$

# Semi-join reductions can be extremely powerful



**Fig. 16** While the query complexity is exponential (number of minimal plans are shown on the right side), our optimizations can even evaluate a very large number of minimal plans (here shown up to 429 for a 8-chain query and 5040 (!) for a 7-star query).

## 6.3 Opt. 3: Deterministic semi-join reduction

The most expensive operations in probabilistic query plans are the group-bys for the probabilistic project operations. These are often applied early in the plans to tuples which are later pruned and do not contribute to the final query result. Our third optimization is to first apply a *full semi-join reduction on the input relations* before starting the probabilistic evaluation from these *reduced input relations*.

We like to draw here an important connection to [54], which introduces the idea of “lazy plans” and shows orders of magnitude performance improvements for safe plans by computing confidences not after each join and projection, but rather at the very end of the plan. We note that our semi-join reduction *serves the same purpose* with similar performance improvements and also apply for safe queries. The advantage of semi-join reductions, however, is that we *do not require any modifications to the query engine*.

# Outline: Cyclic conjunctive queries

- Acyclic conjunctive queries
- Cyclic conjunctive queries
  - 2SAT (a detour)
  - Tree decompositions
  - AGM bound (join processing of cyclic queries)
  - Duality in Linear programming (a quick primer)
  - Worst-case optimal joins
  - Hypertree & other decompositions
  - Optimal joins

*cycles make everything  
more complicated ☹️*



# Why cyclic queries (other than social networks)

```
Likes(person, drink)
Frequents(person, bar)
Serves(bar, drink, cost)
```

## 2. Specify or choose a Query

Supported grammar

104 Bars: Persons who frequent some bar that serves some drink they like.

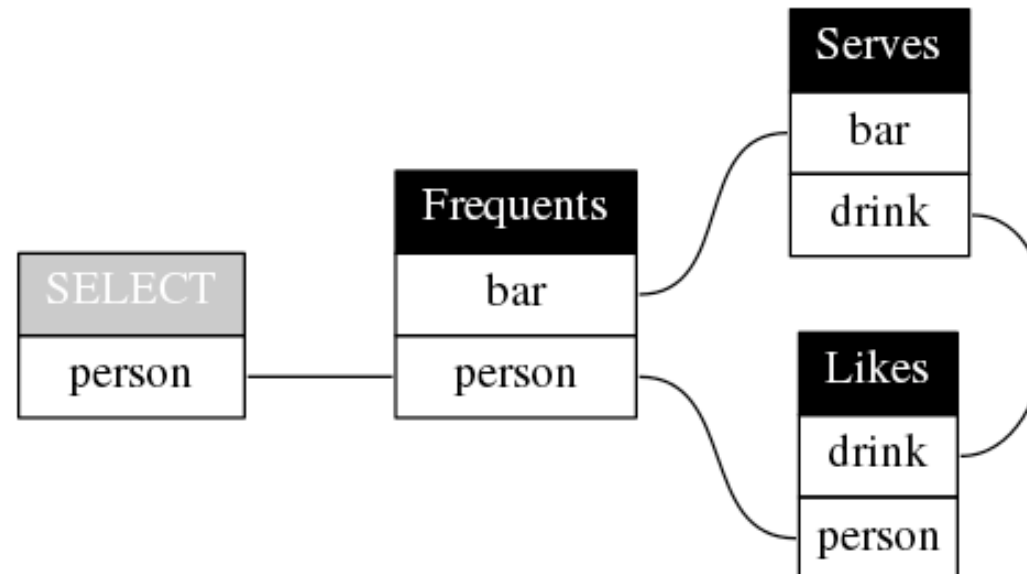
# Why cyclic queries (other than social networks)

```
Likes(person, drink)
Frequents(person, bar)
Serves(bar, drink, cost)
```

## 2. Specify or choose a Query

[Supported grammar](#)

104 Bars: Persons who frequent some bar that serves some drink they like.



# Why cyclic queries (other than social networks)

```
Likes(person, drink)
Frequents(person, bar)
Serves(bar, drink, cost)
```

## 2. Specify or choose a Query

[Supported grammar](#)

104 Bars: Persons who frequent some bar that serves some drink they like.

```
SELECT  F1.person
FROM    Frequents F1
WHERE   exists
        (SELECT *
         FROM    Serves S2
         WHERE   S2.bar = F1.bar
         AND     exists
                (SELECT *
                 FROM    Likes L3
                 WHERE   L3.person = F1.person
                 AND     S2.drink = L3.drink))
```



# Outline: Cyclic conjunctive queries

- Acyclic conjunctive queries
- Cyclic conjunctive queries
  - 2SAT (a detour)
  - Tree decompositions
  - AGM bound (join processing of cyclic queries)
  - Duality in Linear programming (a quick primer)
  - Worst-case optimal joins
  - Hypertree & other decompositions
  - Optimal joins

# 2SAT

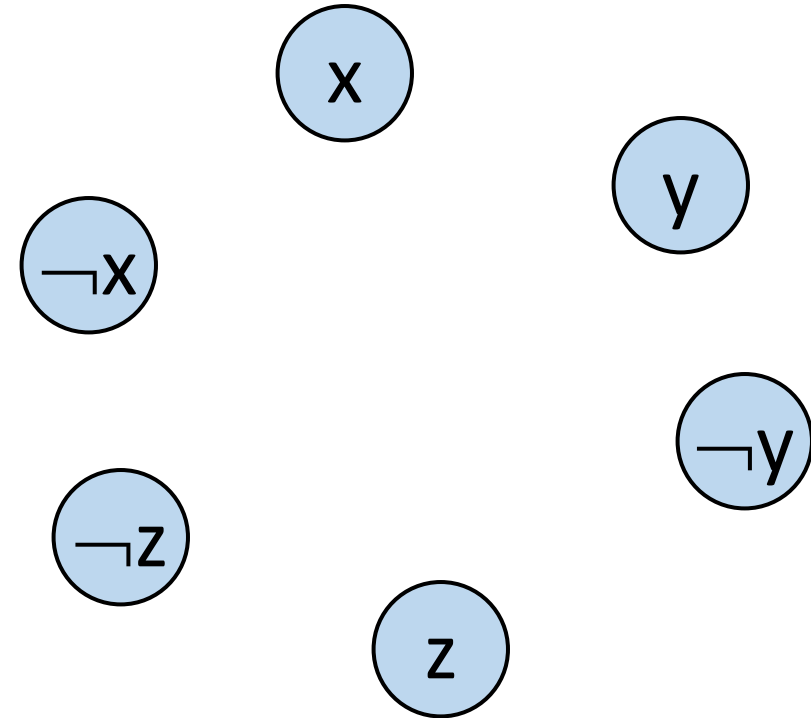
$$\varphi = (x \vee y) \wedge (\neg y \vee z) \wedge (\neg x \vee \neg z) \wedge (z \vee y)$$

- Instance: A 2-CNF formula  $\varphi$
- Problem: To decide if  $\varphi$  is satisfiable
- Theorem: 2SAT is polynomial-time decidable.
  - Proof: We'll show how to solve this problem efficiently using path searches in graphs...
- Background: Given a graph  $G=(V,E)$  and two vertices  $s,t \in V$ , finding if there is a path from  $s$  to  $t$  in  $G$  is polynomial-time decidable. Use some search algorithm (DFS/BFS).

## 2SAT: Graph Construction

$$\varphi = (x \vee y) \wedge (\neg y \vee z) \wedge (\neg x \vee \neg z) \wedge (z \vee y)$$

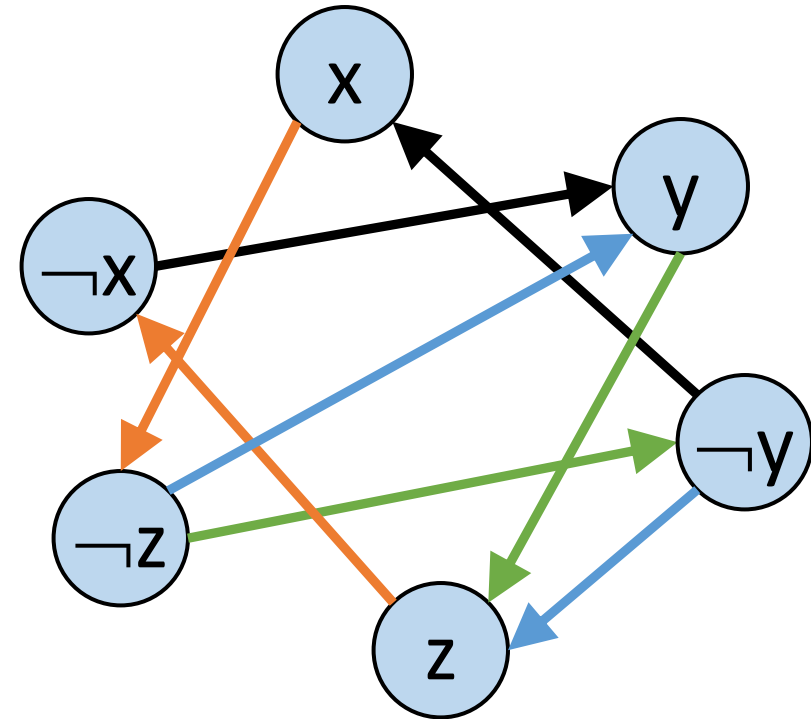
- Vertex for each variable and a negation of a variable



# 2SAT: Graph Construction

$$\varphi = \underbrace{(x \vee y)}_{\text{black}} \wedge \underbrace{(\neg y \vee z)}_{\text{green}} \wedge \underbrace{(\neg x \vee \neg z)}_{\text{orange}} \wedge \underbrace{(z \vee y)}_{\text{blue}}$$

- Vertex for each variable and a negation of a variable
- Edge  $(\neg x \rightarrow y)$  iff there exists a clause equivalent to  $(x \vee y)$ 
  - Recall  $(x \vee y)$  same as  $(\neg x \Rightarrow y)$  and  $(\neg y \Rightarrow x)$ , thus also  $(\neg y \rightarrow x)$

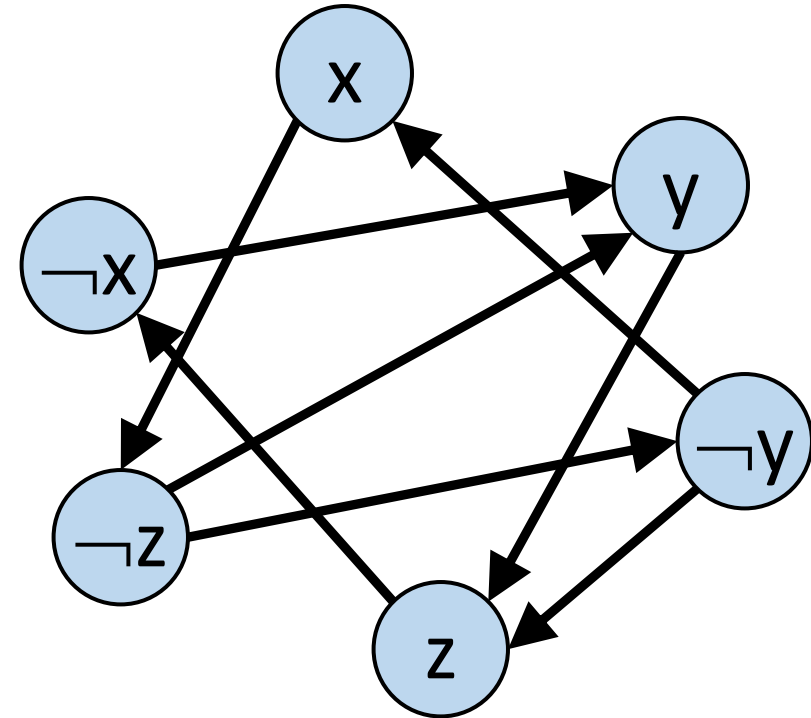




# 2SAT: Graph Construction

$$\varphi = (x \vee y) \wedge (\neg y \vee z) \wedge (\neg x \vee \neg z) \wedge (z \vee y)$$

- Claim: a 2-CNF formula  $\varphi$  is unsatisfiable iff there exists a variable  $x$ , such that:
  - there is a path from  $x$  to  $\neg x$  in the graph, and
  - there is a path from  $\neg x$  to  $x$  in the graph

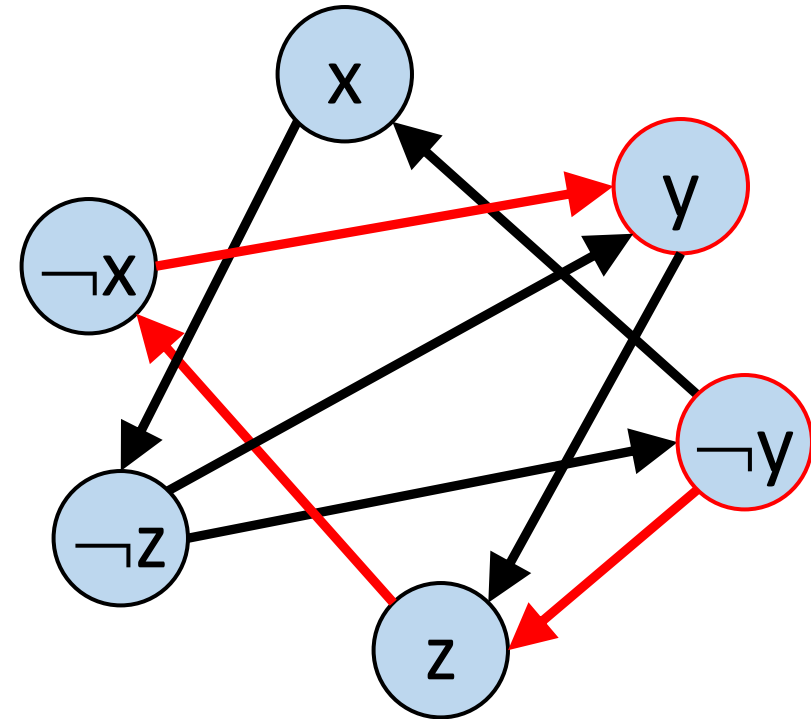


# 2SAT: Graph Construction

$$\varphi = (x \vee y) \wedge (\neg y \vee z) \wedge (\neg x \vee \neg z) \wedge (z \vee y)$$

- Claim: a 2-CNF formula  $\varphi$  is unsatisfiable iff there exists a variable  $x$ , such that:
  - there is a path from  $x$  to  $\neg x$  in the graph, and
  - there is a path from  $\neg x$  to  $x$  in the graph

*not enough,  
needs both directions!*

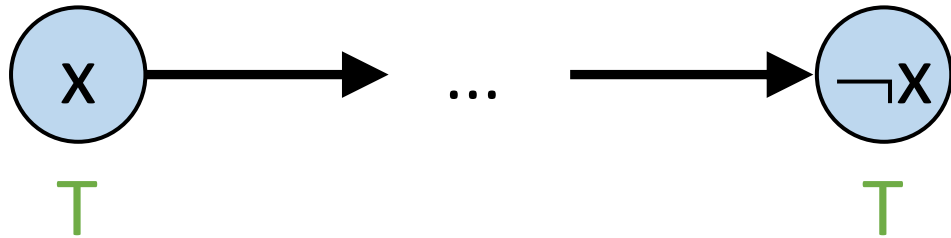


# Correctness (1)

$$\varphi = (x \vee y) \wedge (\neg y \vee z) \wedge (\neg x \vee \neg z) \wedge (z \vee y)$$

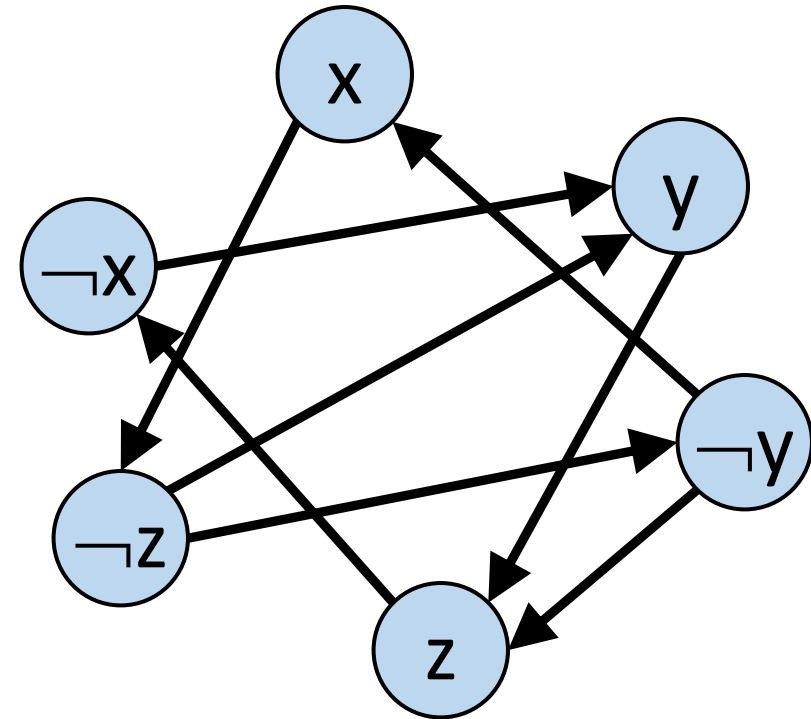
- Suppose there are paths  $x \dots \neg x$  and  $\neg x \dots x$  for some variable  $x$ , but there's also a satisfying assignment  $\rho$ .

– If  $\rho(x)=T$ :



– Similarly for  $\rho(x)=F$ ...

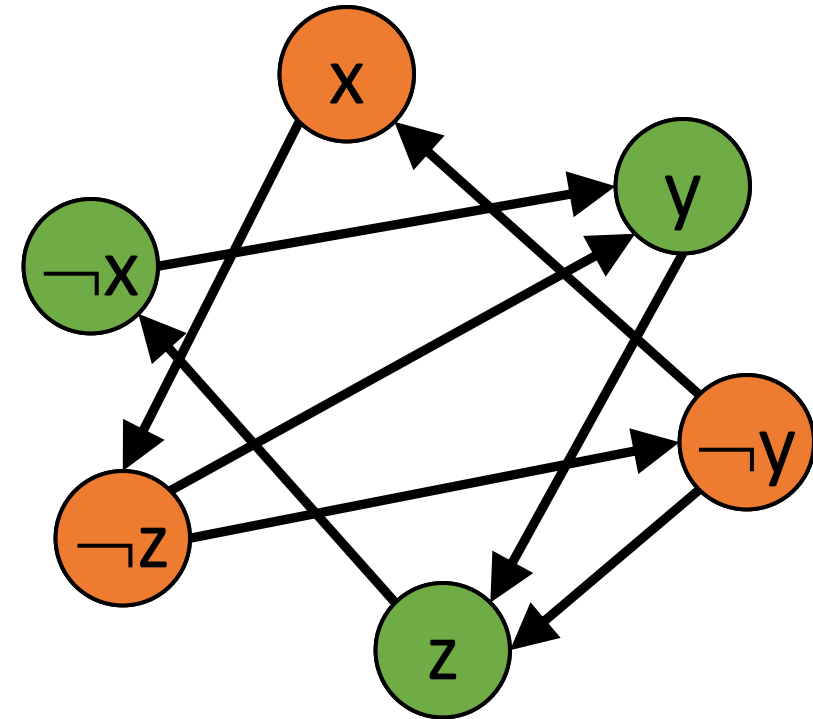
*recall, needs to hold in both directions!*



## Correctness (2)

$$\varphi = (x \vee y) \wedge (\neg y \vee z) \wedge (\neg x \vee \neg z) \wedge (z \vee y)$$

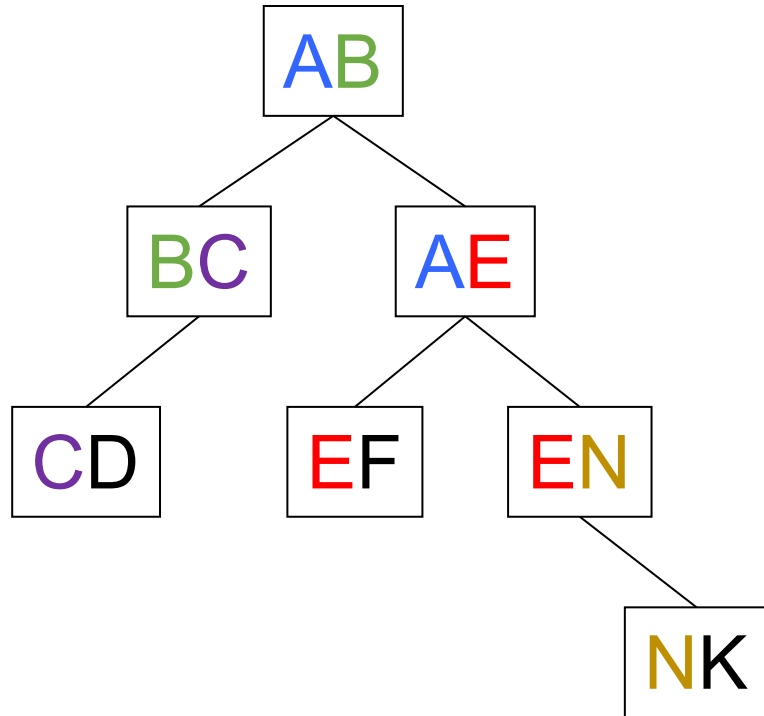
- Suppose there are no such paths.
- Construct an assignment as follows:
  1. pick an unassigned literal  $\alpha$ , with no path from  $\alpha$  to  $\neg\alpha$ , and assign it **T**
  2. assign **T** to all reachable vertices
  3. assign **F** to their negations
  4. Repeat until all vertices are assigned



# Outline: Cyclic conjunctive queries

- Acyclic conjunctive queries
- Cyclic conjunctive queries
  - 2SAT (a detour)
  - Tree decompositions
  - AGM bound (join processing of cyclic queries)
  - Duality in Linear programming (a quick primer)
  - Worst-case optimal joins
  - Hypertree & other decompositions
  - Optimal joins

# Definition of an *attribute-connected tree* (also running *intersection property* or *coherence*)



A tree is **attribute-connected** if the sub-tree induced by each attribute is connected

# Tree decomposition

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one

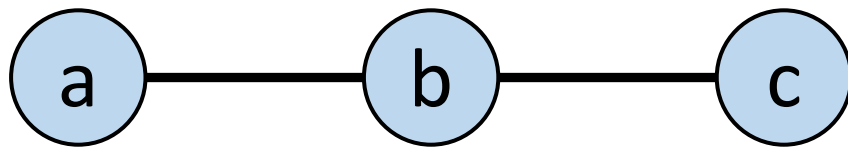


# Tree decomposition ... of a tree

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



*tree decomposition*





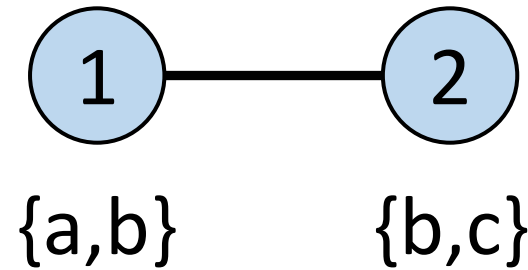
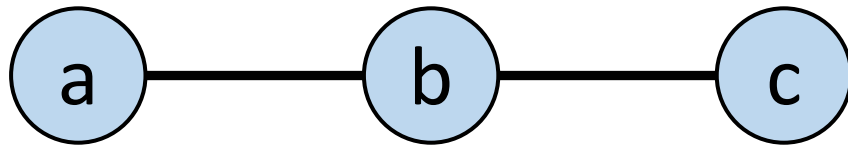


# Tree decomposition ... of a tree

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



That's why **treewidth** defined as max cardinality - 1

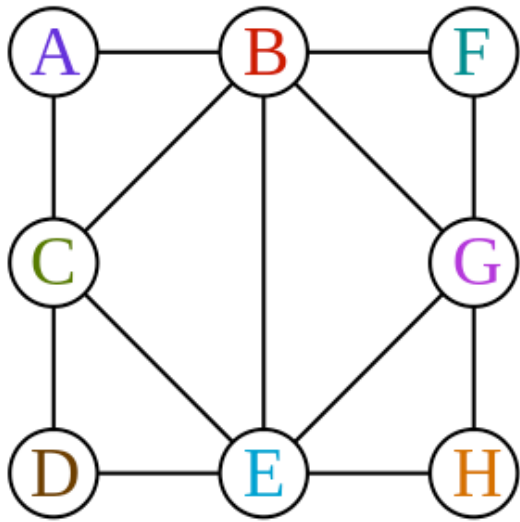


# Tree decomposition example

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



*tree decomposition*

?

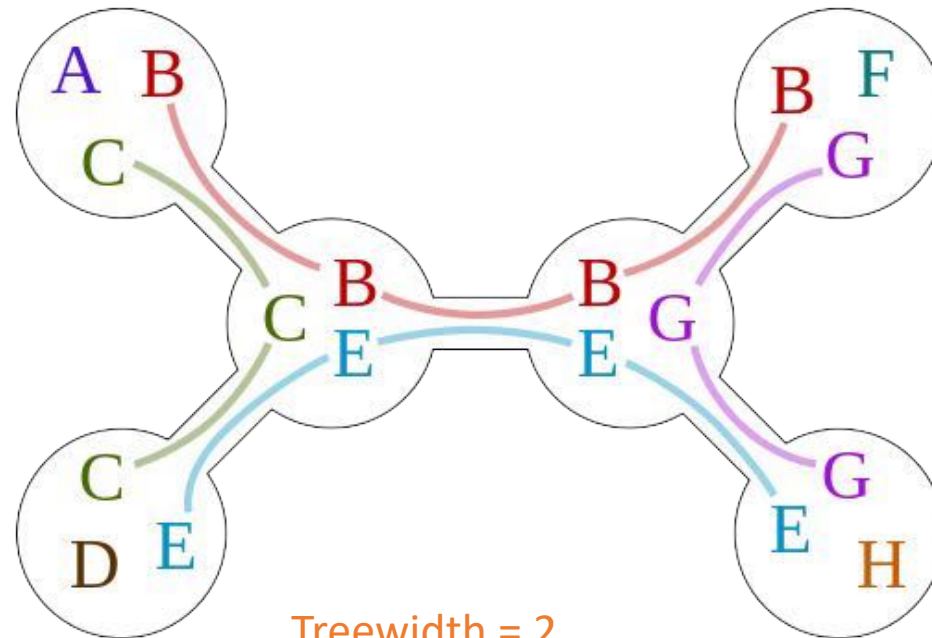
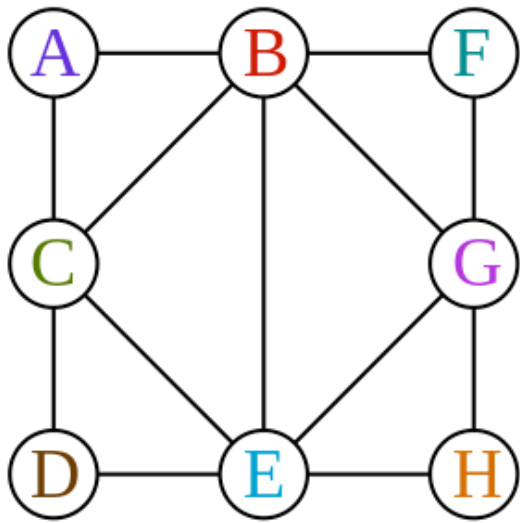


# Tree decomposition example

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



Treewidth = 2

Notice running intersection property

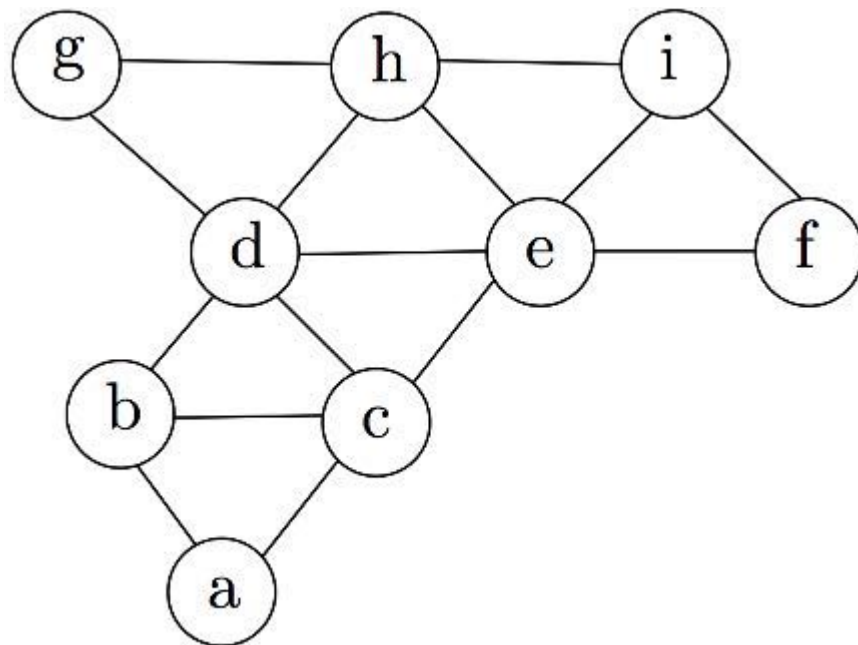


# Tree decomposition example

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



*tree decomposition*

?

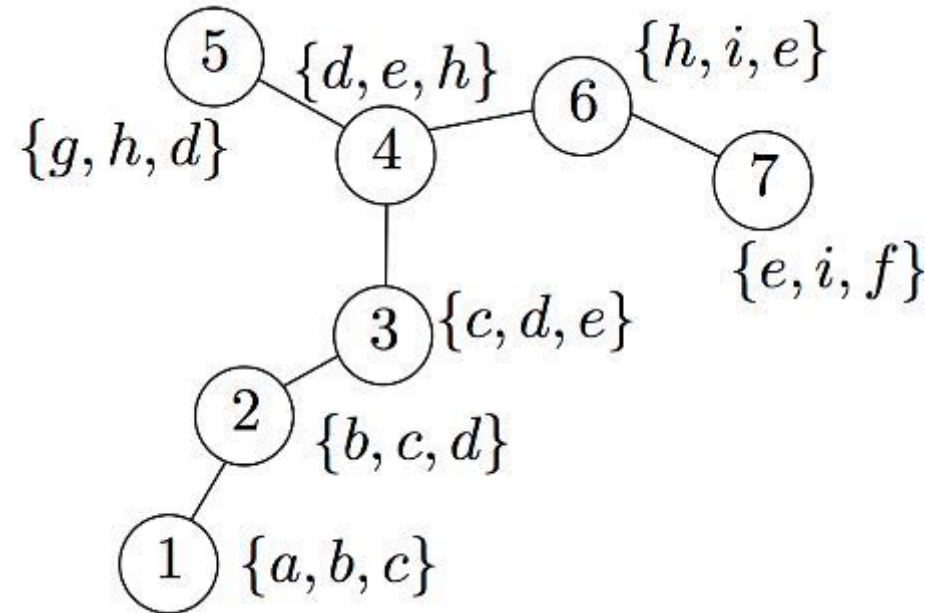
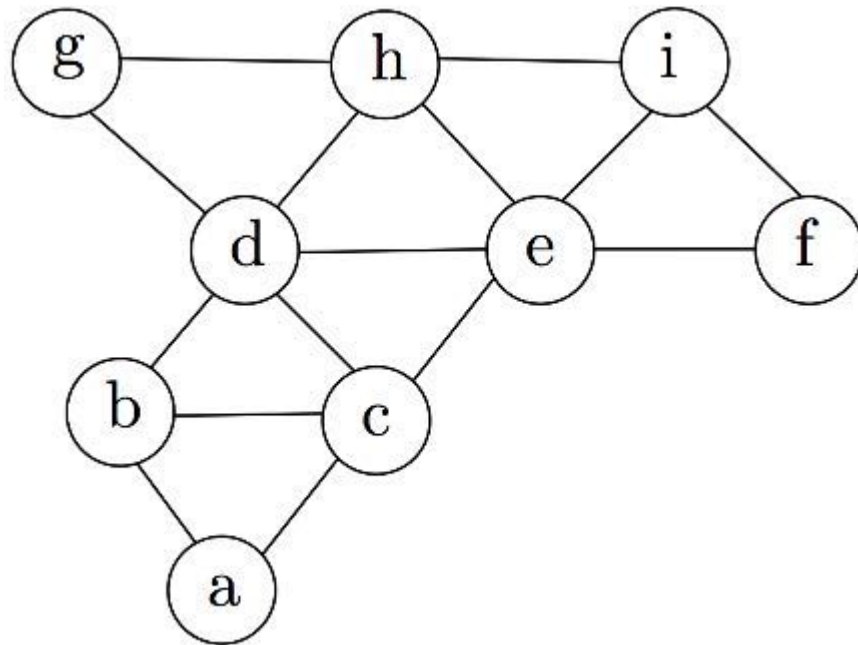


# Tree decomposition example

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



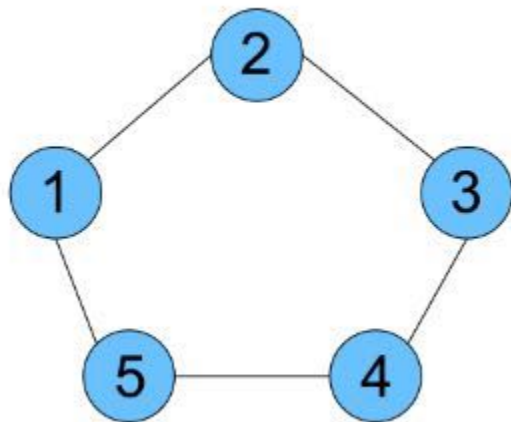


# Tree decomposition of a cycle

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



*tree decomposition*



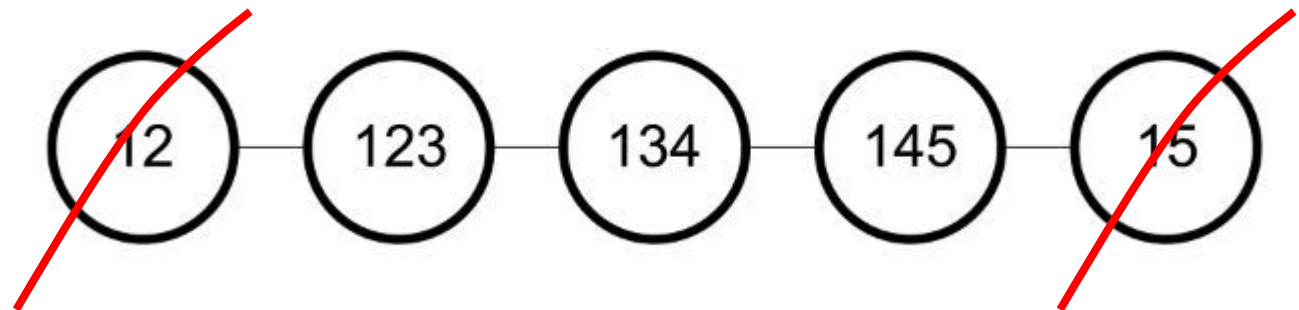
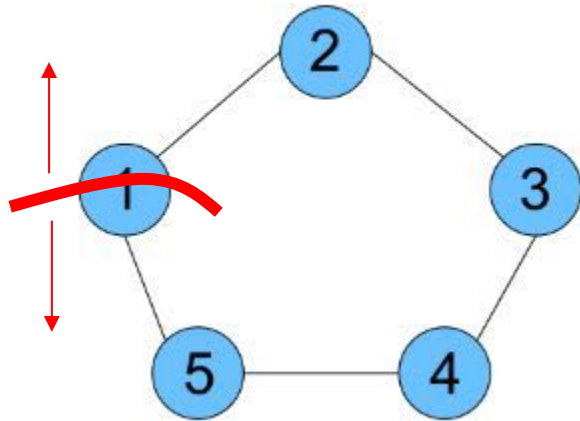


# Tree decomposition of a cycle

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



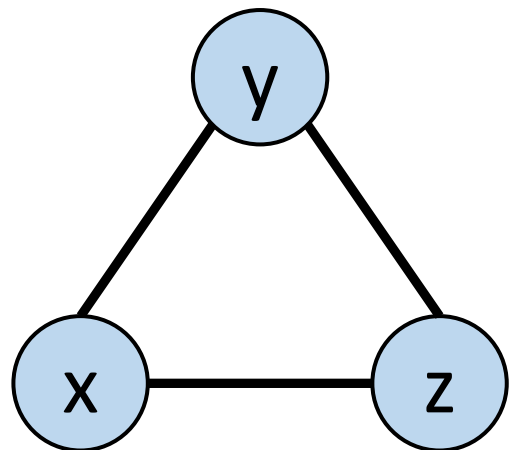


# Tree decomposition of a triangle

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



*tree decomposition*

?



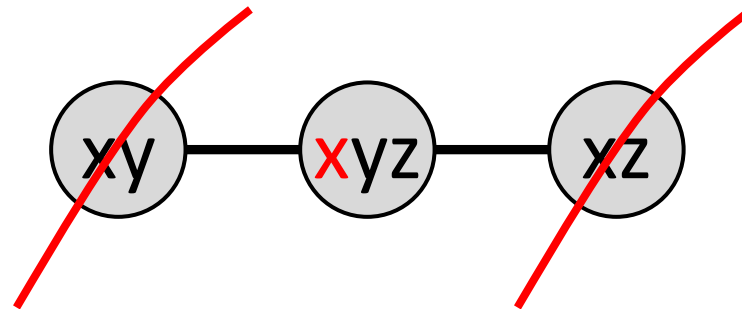
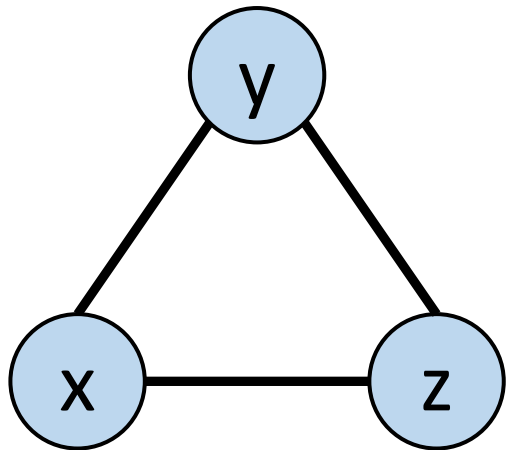


# Tree decomposition of a triangle

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



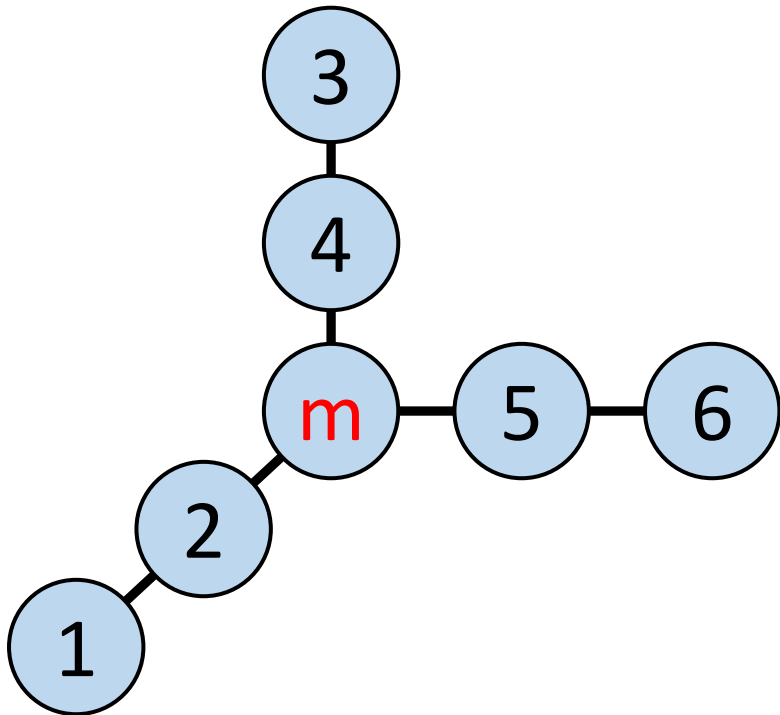


# Tree decomposition of a longer tree

A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

The **width of a tree decomposition** is the size of its largest set minus one



*tree decomposition*

?

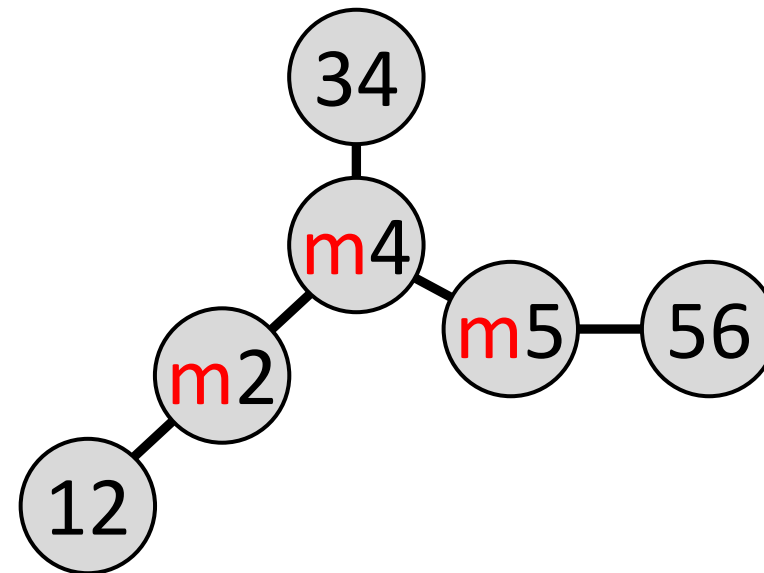
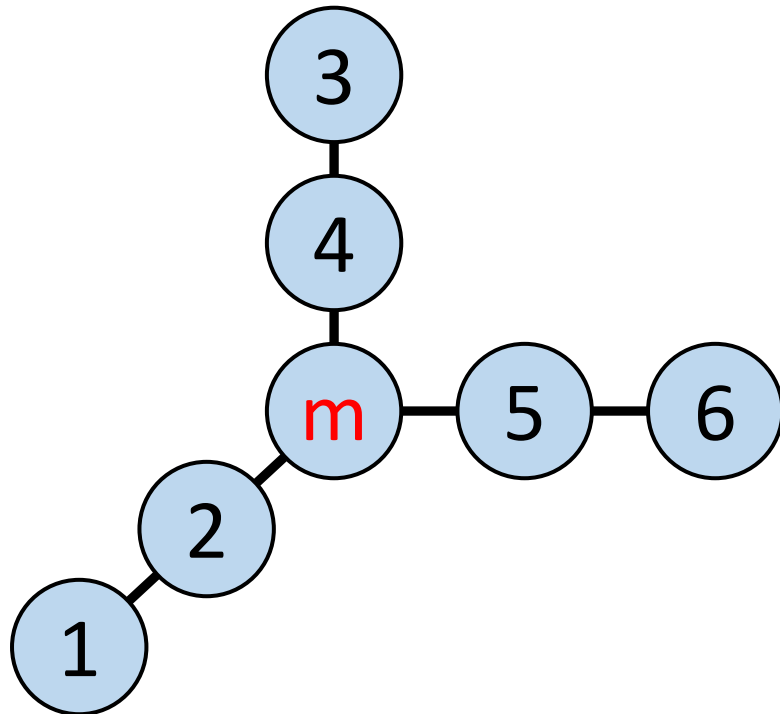


# Tree decomposition of a longer tree

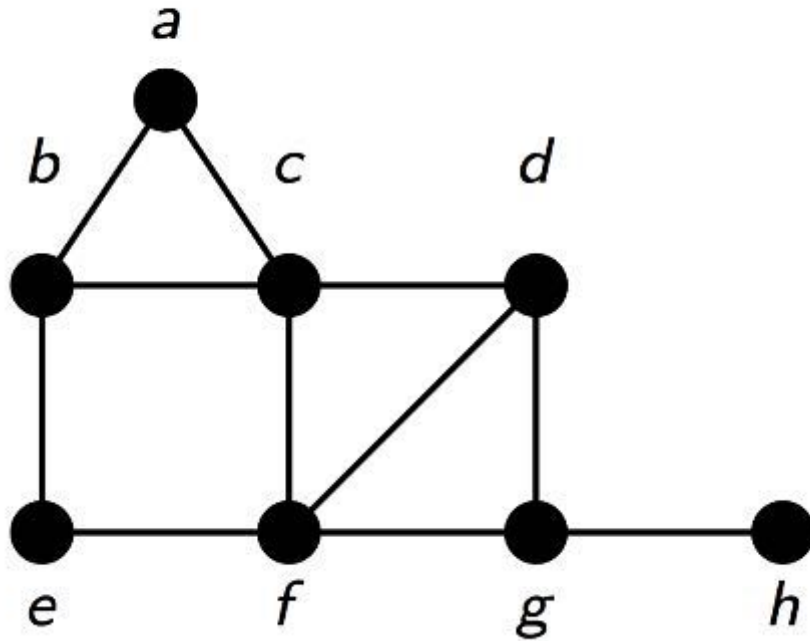
A **tree decomposition** of graph  $G(N, E)$  is a tree  $T(V, F)$  and a subset  $N_v \subseteq N$  assigned to each vertex  $v \in V$  s.t.:

- (1) **Node coverage**: Every vertex of  $G$  is assigned least one vertex in  $T$
- (2) **Edge coverage**: For every edge  $e$  of  $G$ , there is a vertex in  $T$  that contains both ends of  $e$
- (3) **Coherence**: The tree is "attribute-connected"

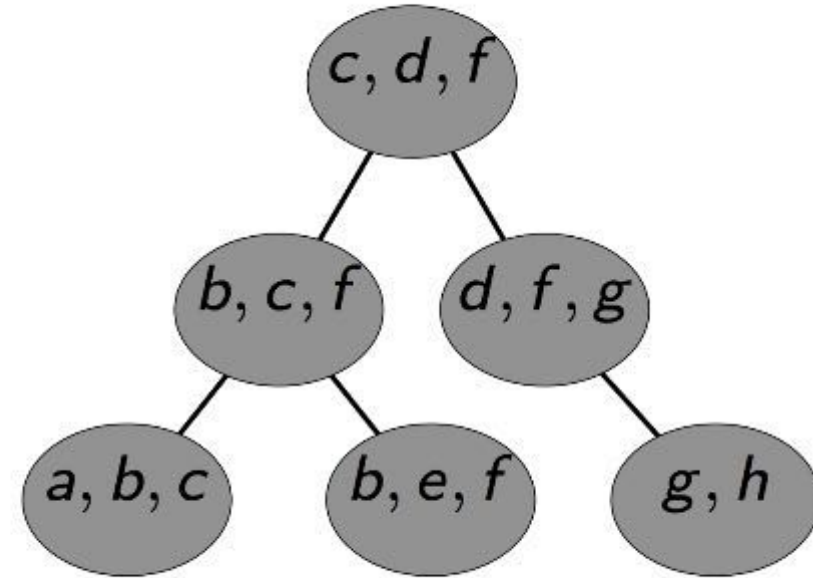
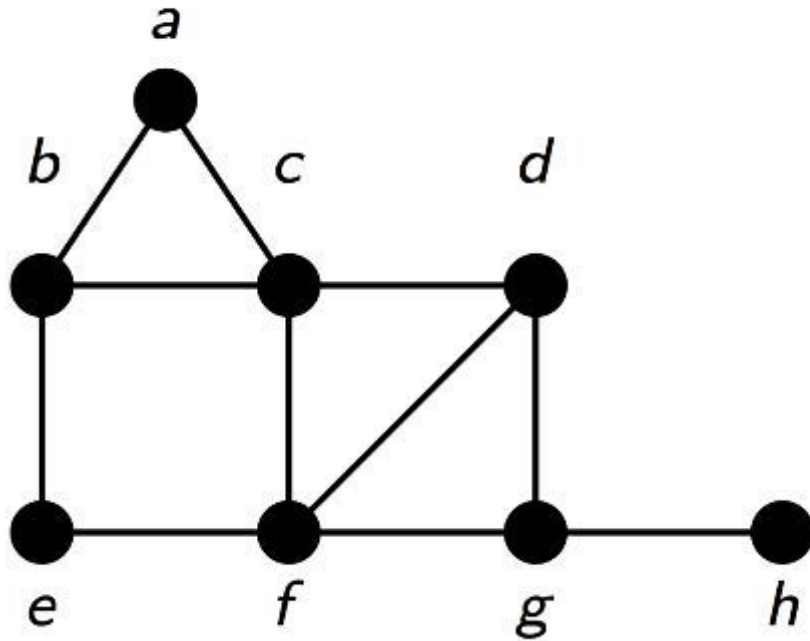
The **width of a tree decomposition** is the size of its largest set minus one



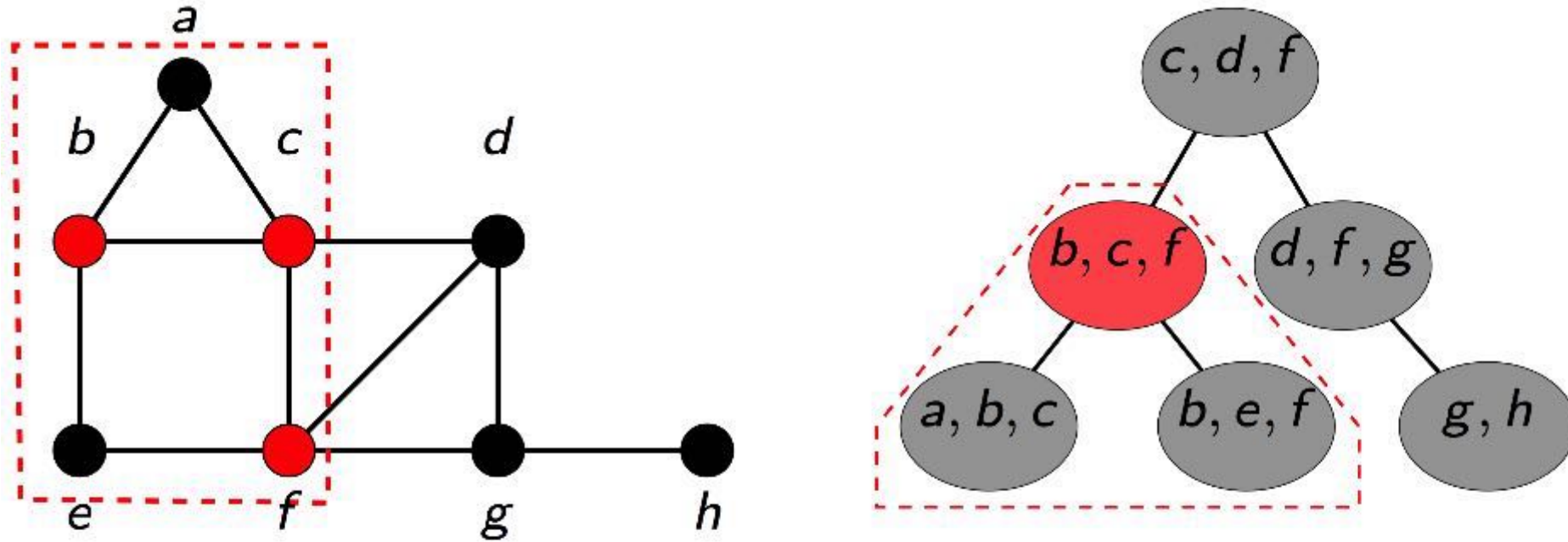
# Tree decomposition



# Tree decomposition



# Tree decomposition

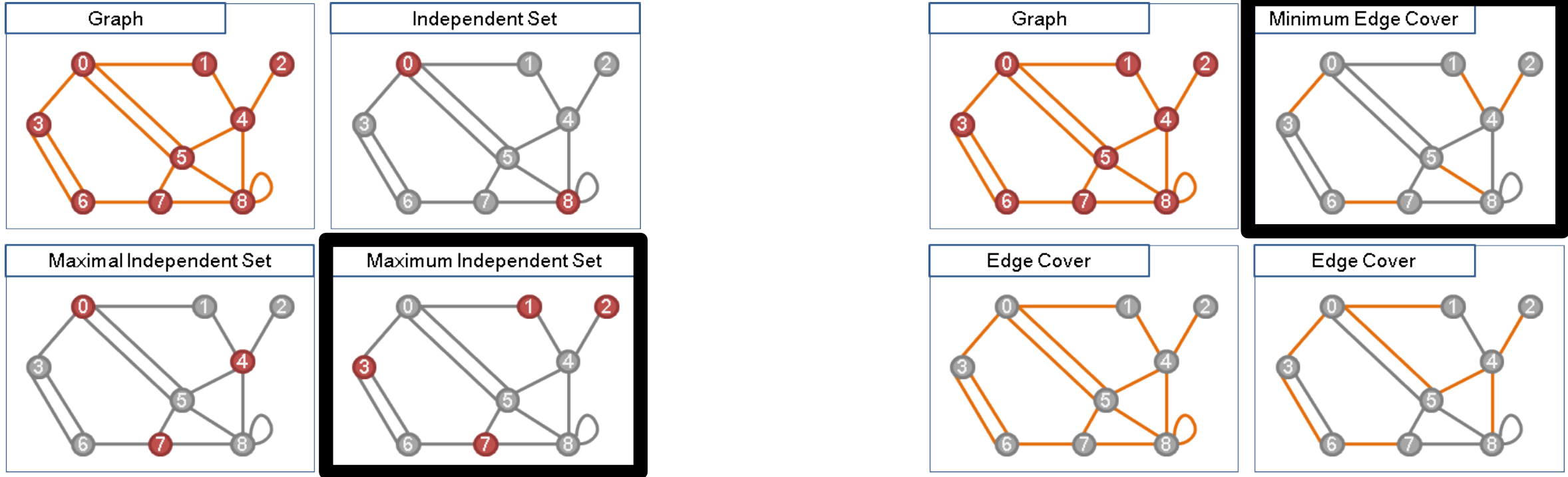


A subtree communicates with the outside world only via the root of the subtree.

# Outline: Cyclic conjunctive queries

- Acyclic conjunctive queries
- Cyclic conjunctive queries
  - 2SAT (a detour)
  - Tree decompositions
  - AGM bound (join processing of cyclic queries)
  - Duality in Linear programming (a quick primer)
  - Worst-case optimal joins
  - Hypertree & other decompositions
  - Optimal joins

# Background: MAX independent (vertex) set $\leq$ MIN edge cover



- Assume graph  $G$  is connected. Thus, every vertex has at least one edge (unless just one vertex)
- Suppose  $S$  is an independent set and  $E$  is an edge cover.
- Then for each vertex  $v \in S$  there exists at least one edge  $e \in E$  incident with  $v$ .
- By definition of independent set no two  $u, v \in S$ , have a common edge in  $E$ .
- Therefore  $|S| \leq |E|$

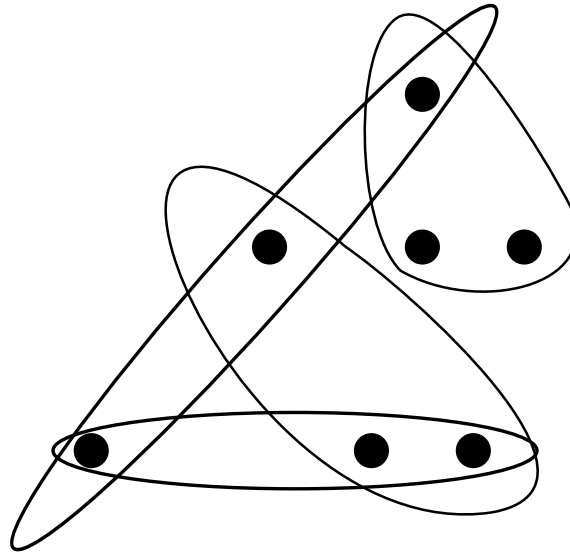


What do we know about  
bounding the size of the  
answer?

(...and enumerating all solutions)

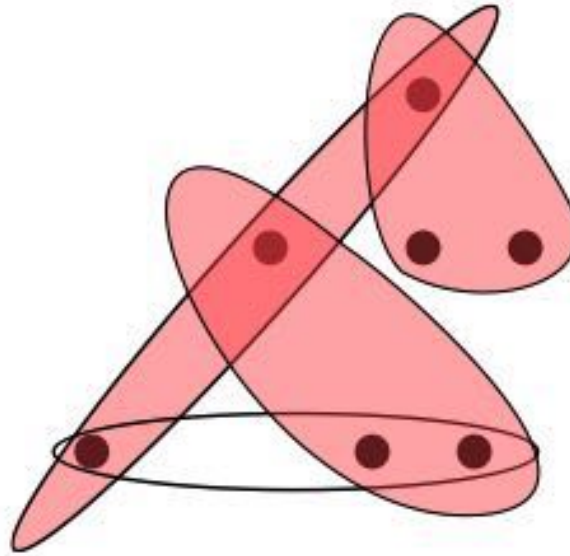
## Upper bound

**Observation:** If the hypergraph has edge cover number  $\rho$  and every relation has size at most  $N$ , then there are at most  $N^\rho$  tuples in the answer.



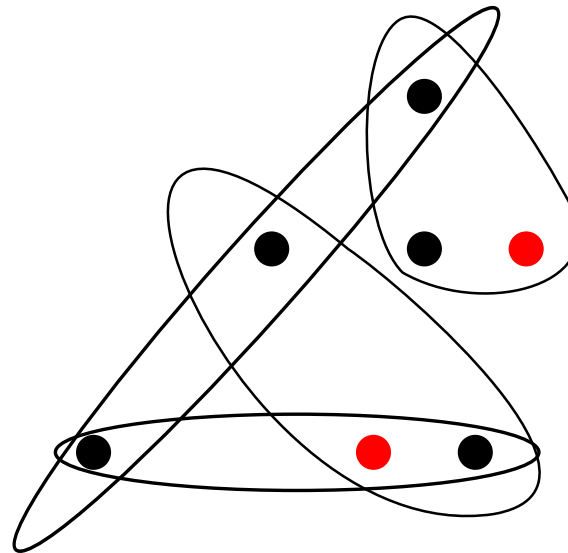
## Upper bound

**Observation:** If the hypergraph has edge cover number  $\rho$  and every relation has size at most  $N$ , then there are at most  $N^\rho$  tuples in the answer.



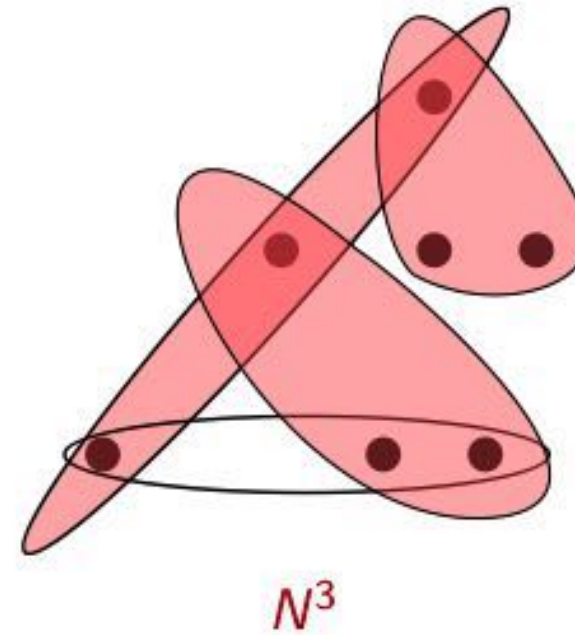
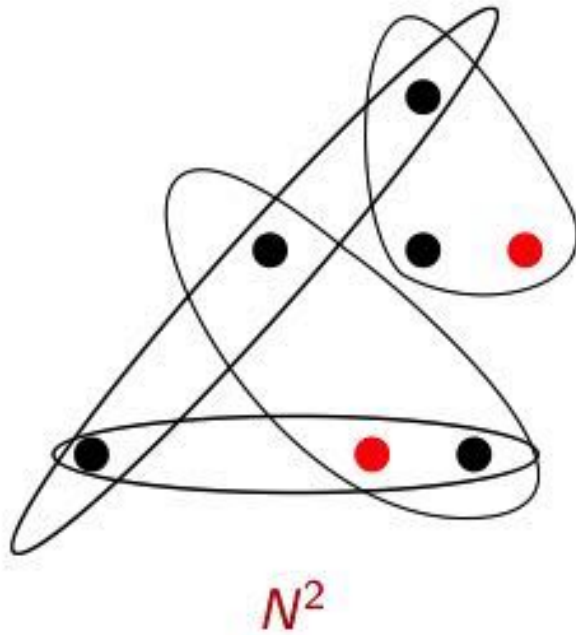
## Lower bound

**Observation:** If the hypergraph has independence number  $\alpha$ , then one can construct an instance where every relation has size  $N$  at the answer has size  $N^\alpha$ .



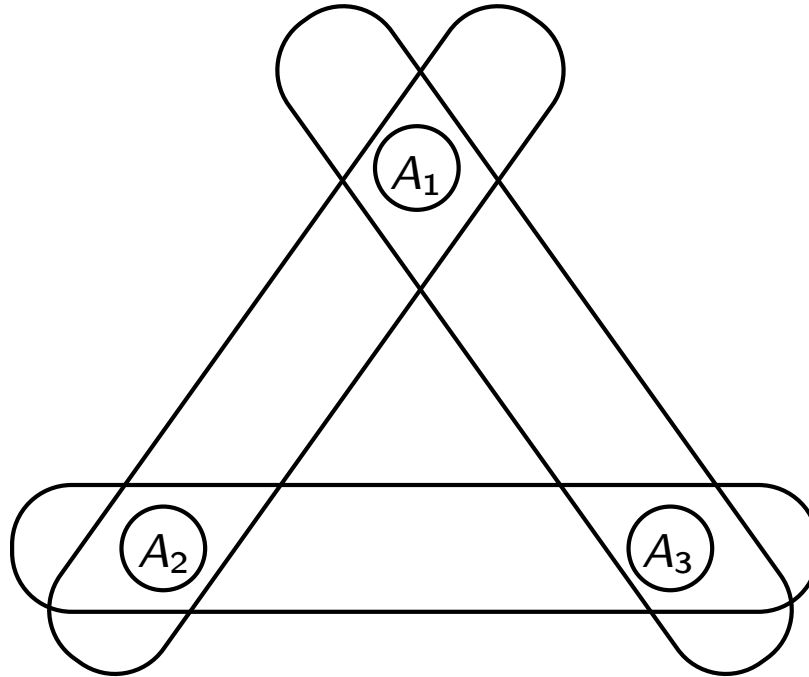
Definition of the relations:

- If variable  $A$  is in the independent set, then it can take any value in  $[N]$ .
- Otherwise it is forced to 1.



Which is tight: the upper bound or the lower bound?

## Example: triangles

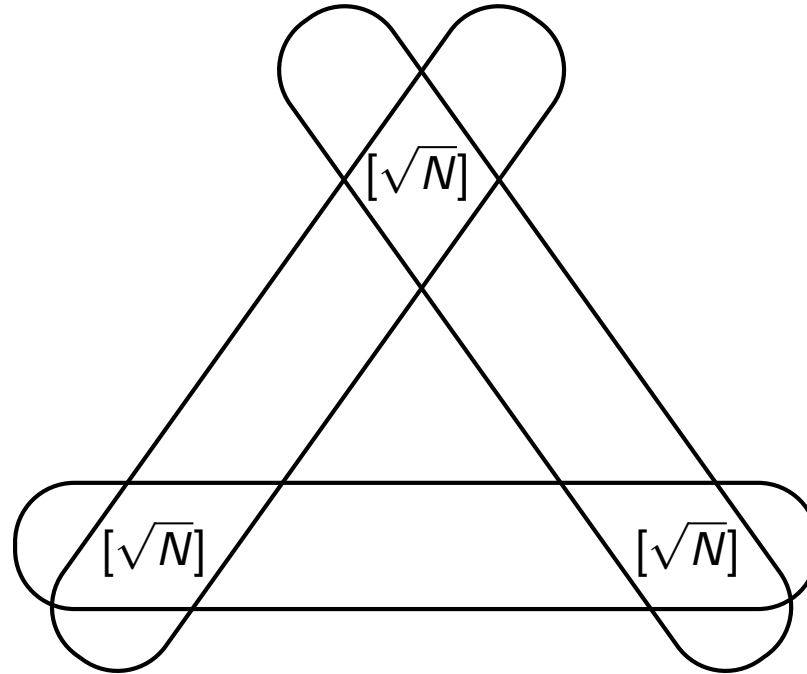


### Upper bound

Two kind of values for  $A_1$ :

- **Light:** can be extended to at most  $\sqrt{N}$  ways to  $A_2$ .  
 $\Rightarrow \leq N \cdot \sqrt{N}$  answers with light  $A_1$
  - **Heavy:** can be extended to at least  $\sqrt{N}$  ways to  $A_2$ .  
 $\Rightarrow \leq \sqrt{N}$  heavy values  $\Rightarrow \leq \sqrt{N} \cdot N$  answers with heavy  $A_1$
- $\Rightarrow$  At most  $2 \cdot N^{3/2}$  answers.

## Example: triangles



### Lower bound

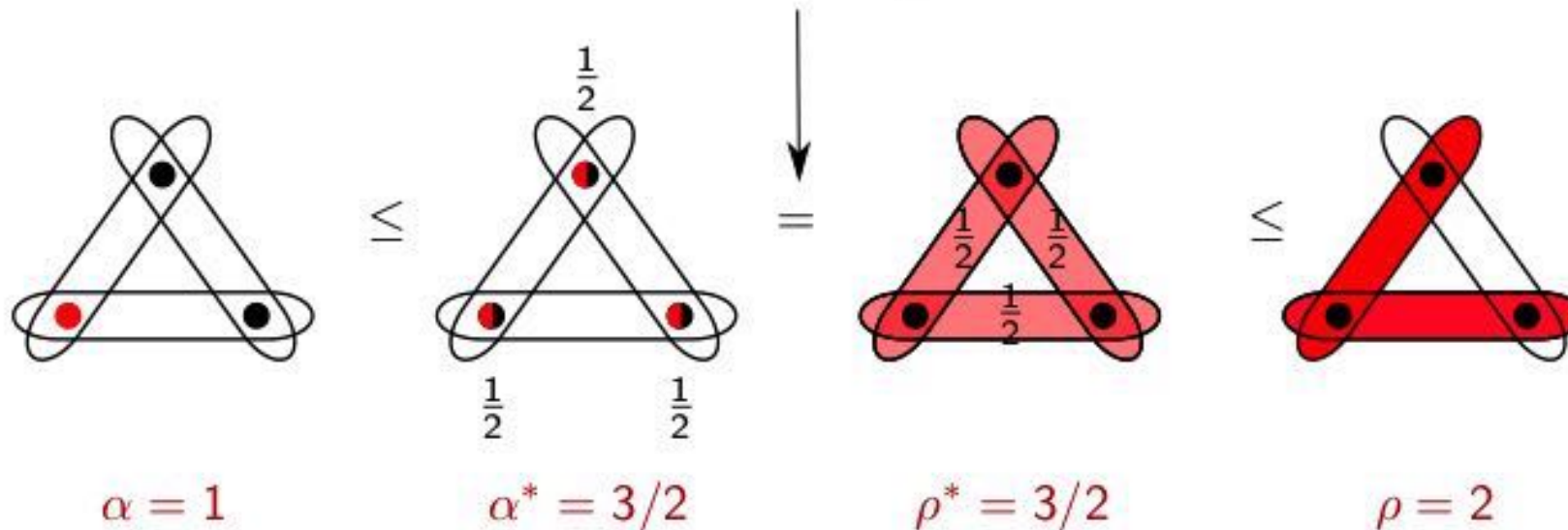
Allow every variable to be any value from  $[\sqrt{N}] \Rightarrow N^{3/2}$  answers.

The correct bound  $N^{3/2}$  is between  
 $N^\alpha = N^1$  and  $N^\rho = N^2$ .

# Fractional values

- $\alpha$ : independence number
- $\alpha^*$ : fractional independence number  
(max. weight of vertices s.t. each edge contains weight  $\leq 1$ )
- $\rho^*$ : fractional edge cover number  
(min. weight of edges s.t. each vertex receives weight  $\geq 1$ )
- $\rho$ : edge cover number

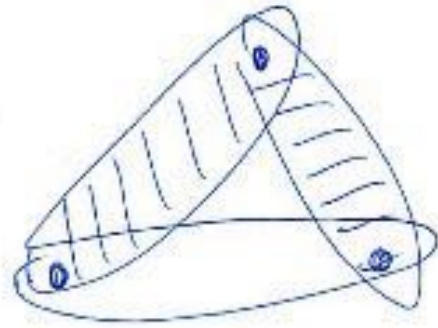
LP duality!



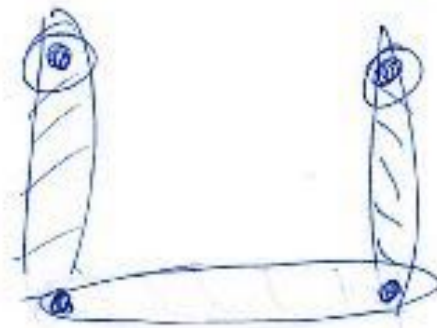


# Examples

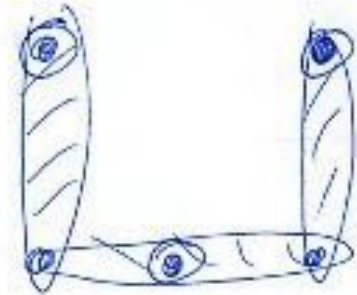
poly  
Cover



2

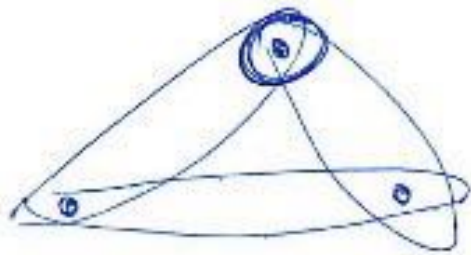


2

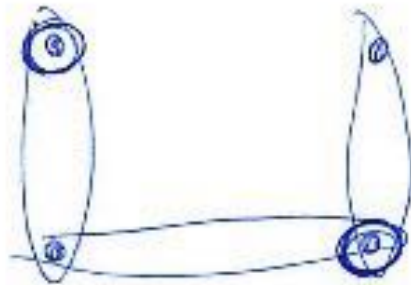


3

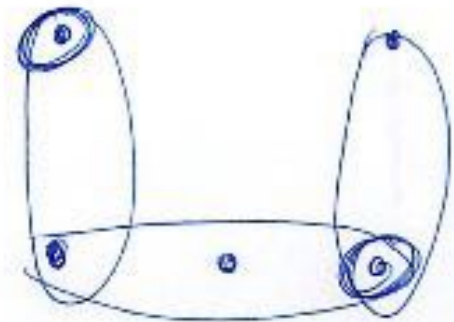
incl.  
SM



1



2



2

# Outline: Cyclic conjunctive queries

- Acyclic conjunctive queries
- Cyclic conjunctive queries
  - 2SAT (a detour)
  - Tree decompositions
  - AGM bound (join processing of cyclic queries)
  - **Duality in Linear programming (a quick primer)**
  - Worst-case optimal joins
  - Hypertree & other decompositions
  - Optimal joins

# Dual Optimization Problem

- A maximization problem **M** and a minimization problem **N**, defined on the same instances (such as graphs) such that:
  1. for every candidate solution  $M$  to **M** and every candidate solution  $N$  to **N**, the value of  $M$  is less than or equal to the value of  $N$
  2. obtaining candidate solutions  $M$  and  $N$  that have the same value proves that  $M$  and  $N$  are optimal solutions for that instance.

# A quick primer on Duality in Linear Programming

$$\max x_1 + 6x_2$$

$$x_1 \leq 200$$

$$x_2 \leq 300$$

$$x_1 + x_2 \leq 400$$

$$x_1, x_2 \geq 0$$

optimum solution to be  $(x_1, x_2) = (100, 300)$

objective value 1900

# A quick primer on Duality in Linear Programming

$$\max x_1 + 6x_2$$

$$x_1 \leq 200 \quad \times 1$$

$$x_2 \leq 300 \quad \times 6$$

$$x_1 + x_2 \leq 400 \quad \times 0$$

$$x_1, x_2 \geq 0$$

$$x_1 + 6x_2 \leq 2000$$

upper bound!

optimum solution to be  $(x_1, x_2) = (100, 300)$

objective value 1900

# A quick primer on Duality in Linear Programming

$$\max x_1 + 6x_2$$

$$x_1 \leq 200 \quad \times 0$$

$$x_2 \leq 300 \quad \times 5$$

$$x_1 + x_2 \leq 400 \quad \times 1$$

$$x_1, x_2 \geq 0$$

$$x_1 + 6x_2 \leq 1900$$

$(0,5,1)$ ... certificate of optimality

optimum solution to be  $(x_1, x_2) = (100, 300)$

objective value 1900

# A quick primer on Duality in Linear Programming

non-negative!

		Multiplier	Inequality
$\max x_1 + 6x_2$			
$x_1 \leq 200$	<b>x 0</b>	$y_1$	$x_1 \leq 200$
$x_2 \leq 300$	<b>x 5</b>	$y_2$	$x_2 \leq 300$
$x_1 + x_2 \leq 400$	<b>x 1</b>	$y_3$	$x_1 + x_2 \leq 400$
$x_1, x_2 \geq 0$			
			$x_1 + 6x_2 \leq 1900$

**(0,5,1)...** certificate of optimality

# A quick primer on Duality in Linear Programming

non-negative!

$$\max x_1 + 6x_2$$

$$x_1 \leq 200$$

$$x_2 \leq 300$$

$$x_1 + x_2 \leq 400$$

$$x_1, x_2 \geq 0$$

Multiplier

Inequality

$$y_1$$

$$x_1 \leq 200$$

$$y_2$$

$$x_2 \leq 300$$

$$y_3$$

$$x_1 + x_2 \leq 400$$

$$(y_1 + y_3)x_1 + (y_2 + y_3)x_2 \leq 200y_1 + 300y_2 + 400y_3$$

right side upper bound  
for  $(x_1 + 6x_2)$  if

$$y_1 + y_3 \geq 1$$

$$y_2 + y_3 \geq 6$$



# A quick primer on Duality in Linear Programming

$$\text{Primal : } (x_1, x_2) = (100, 300); \quad \text{Dual : } (y_1, y_2, y_3) = (0, 5, 1)$$

$$\max x_1 + 6x_2$$

$$x_1 \leq 200$$

$$x_2 \leq 300$$

$$x_1 + x_2 \leq 400$$

$$x_1, x_2 \geq 0$$

$$\min 200y_1 + 300y_2 + 400y_3$$

$$y_1 + y_3 \geq 1$$

$$y_2 + y_3 \geq 6$$

$$y_1, y_2, y_3 \geq 0$$

$$(y_1 + y_3)x_1 + (y_2 + y_3)x_2 \leq 200y_1 + 300y_2 + 400y_3$$

right side upper bound  
for  $(x_1 + 6x_2)$  if

$$y_1 + y_3 \geq 1$$

$$y_2 + y_3 \geq 6$$

# A quick primer on Duality in Linear Programming

**Figure 7.10** A generic primal LP in matrix-vector form, and its dual.

Primal LP:

$$\begin{aligned} \max \quad & \mathbf{c}^T \mathbf{x} \\ \mathbf{Ax} \leq & \mathbf{b} \\ \mathbf{x} \geq & 0 \end{aligned}$$

Dual LP:

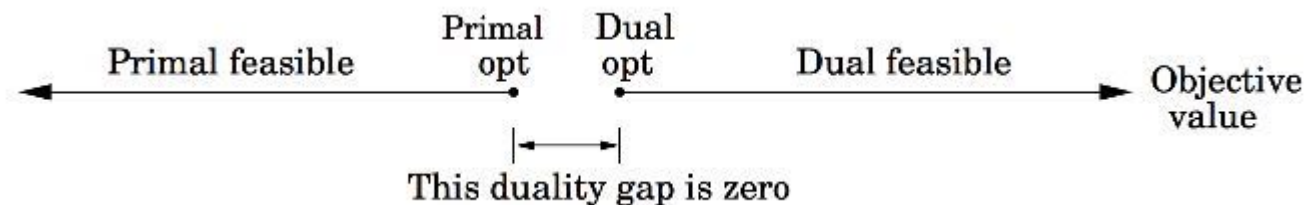
$$\begin{aligned} \min \quad & \mathbf{y}^T \mathbf{b} \\ \mathbf{y}^T \mathbf{A} \geq & \mathbf{c}^T \\ \mathbf{y} \geq & 0 \end{aligned}$$

Primal LP:

$$\begin{aligned} \max \quad & c_1x_1 + \cdots + c_nx_n \\ a_{i1}x_1 + \cdots + a_{in}x_n \leq & b_i \quad \text{for } i \in I \\ a_{i1}x_1 + \cdots + a_{in}x_n = & b_i \quad \text{for } i \in E \\ x_j \geq & 0 \quad \text{for } j \in N \end{aligned}$$

Dual LP:

$$\begin{aligned} \min \quad & b_1y_1 + \cdots + b_my_m \\ a_{1j}y_1 + \cdots + a_{mj}y_m \geq & c_j \quad \text{for } j \in N \\ a_{1j}y_1 + \cdots + a_{mj}y_m = & c_j \quad \text{for } j \notin N \\ y_i \geq & 0 \quad \text{for } i \in I \end{aligned}$$



# Pointers to related work

- "AGM bound": Atserias, Grohe, Marx. *Size bounds and query plans for relational joins*. SIAM J. Comput. 2013. <https://doi.org/10.1137/110859440> (also FOCS 2008)
- "Worst-Case Optimal (WCO) joins": Ngo, Porat, Re, Rudra. *Worst-case optimal join algorithms*. JACM 2018. <https://doi.org/10.1145/3180143> (also PODS 2012)
- "FAQ paper": Khamis, Ngo, Rudra. *FAQ: Questions Asked Frequently*. PODS 2016. <https://doi.org/10.1145/2902251.2902280> (see also SIGMOD record 2017).
- Khamis, Ngo, Suciu. *What do Shannon-type inequalities, submodular width, and disjunctive Datalog have to do with one another?* PODS 2017. <https://doi.org/10.1145/3034786.3056105>
- Robertson, Seymour. *Graph minors. II. Algorithmic aspects of tree-width*. Journal of Algorithms. 1986. [https://doi.org/10.1016/0196-6774\(86\)90023-4](https://doi.org/10.1016/0196-6774(86)90023-4)