

# A unified and practical approach to Generalized Deletion Propagation

**Wolfgang Gatterbauer**, joint work with **Neha Makhija**

<https://northeastern-datalab.github.io/unified-reverse-data-management/>



Logic and Algorithms in Database Theory and AI Reunion

Jan 15, 2024

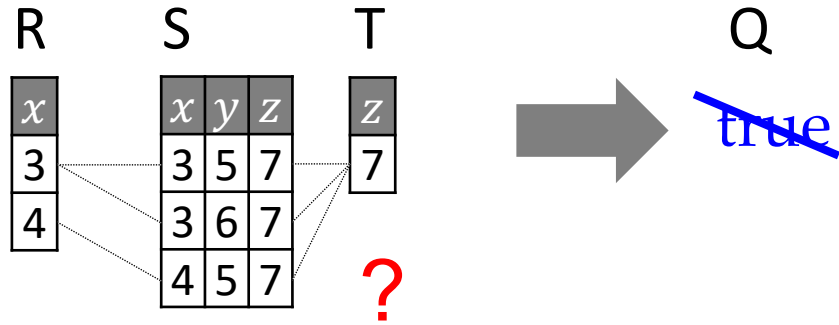
<https://simons.berkeley.edu/workshops/logic-algorithms-database-theory-ai-reunion>

# Example Reverse Data Management Problems

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

## 1. Resilience

Delete min number of tuples to make Q false

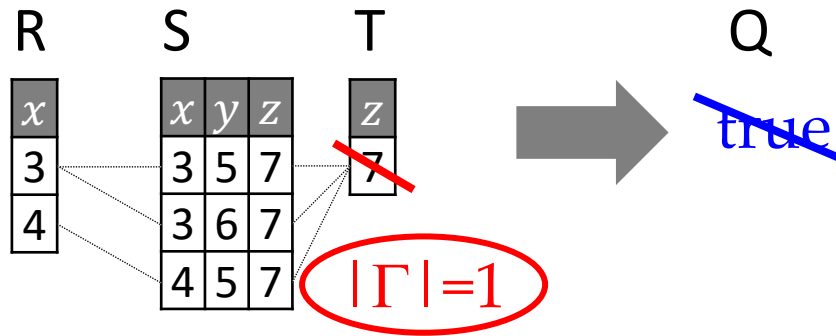


# Example Reverse Data Management Problems

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

## 1. Resilience

Delete min number of tuples to make Q false



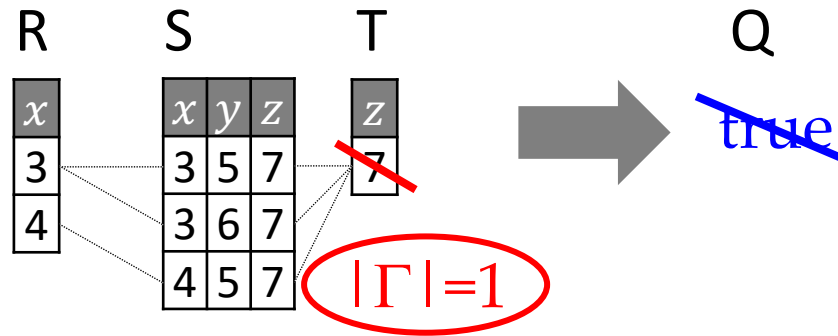
# Example Reverse Data Management Problems

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

Q(u,w):-R'(u,x),S(x,y,z),T'(z,w)

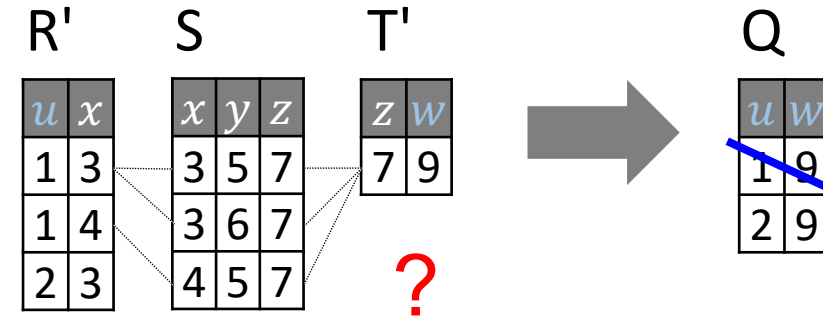
## 1. Resilience

Delete min number of tuples to make Q false



## 2. Source side-effects

Delete min number of tuples to delete an output tuples



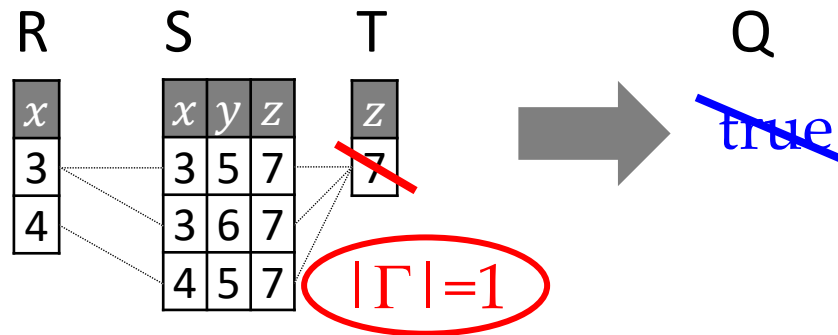
# Example Reverse Data Management Problems

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

Q(u,w):-R'(u,x),S(x,y,z),T'(z,w)

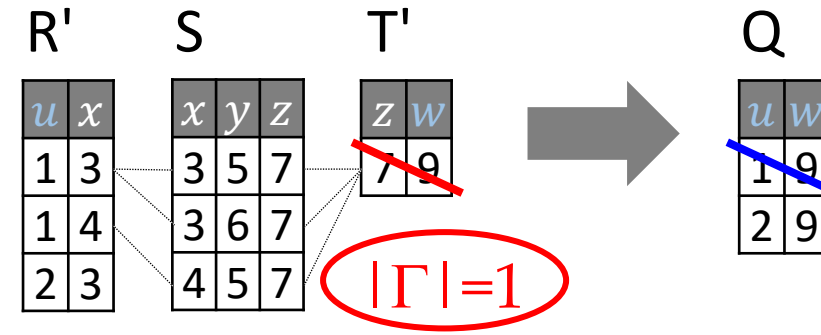
## 1. Resilience

Delete min number of tuples to make Q false



## 2. Source side-effects

Delete min number of tuples to delete an output tuples

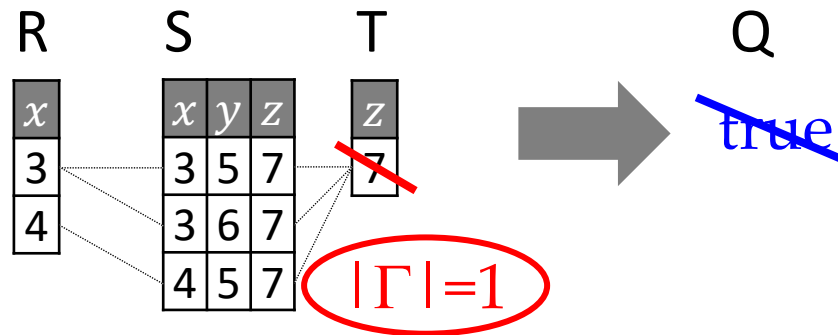


# Example Reverse Data Management Problems

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

## 1. Resilience

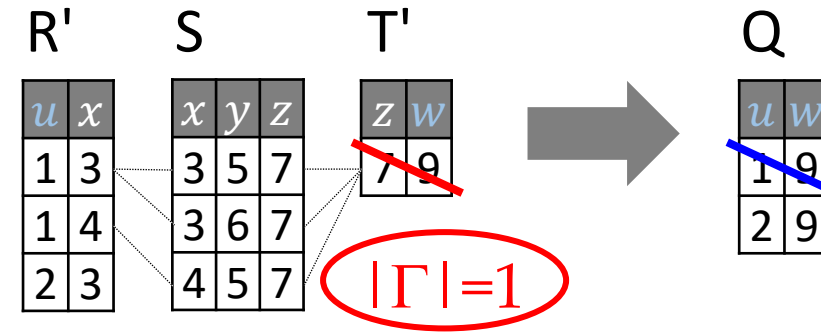
Delete min number of tuples to make Q false



Q(u,w):  $\neg R'(u,x), S(x,y,z), T'(z,w)$

## 2. Source side-effects

Delete min number of tuples to delete an output tuples



## 3. Aggregated deletion propagation

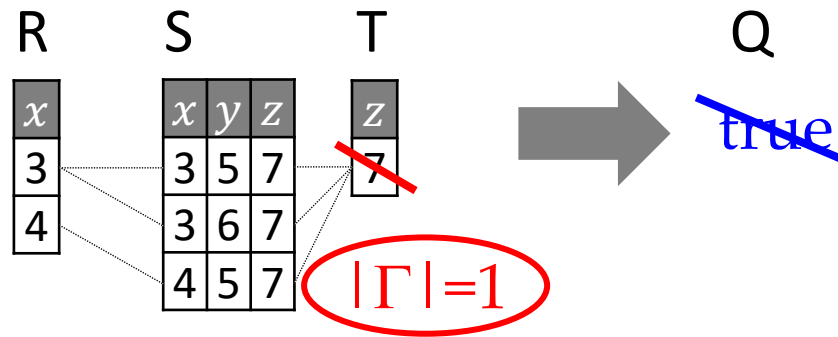
Delete min number of tuples to delete  $\geq k$  output tuples

# Example Reverse Data Management Problems

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

## 1. Resilience

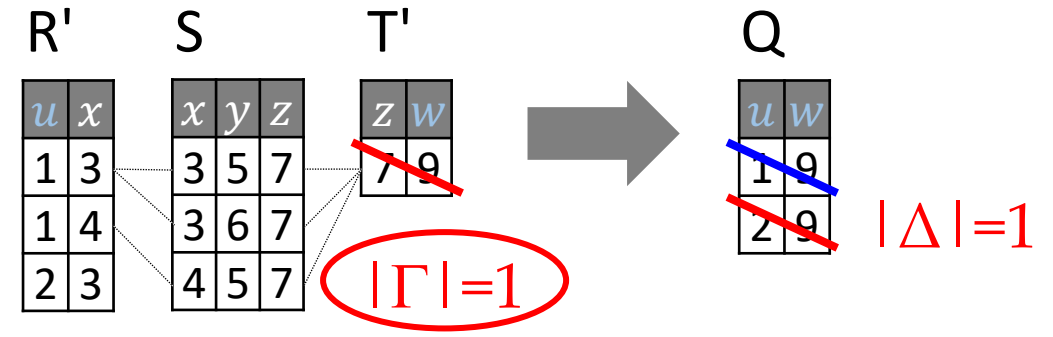
Delete min number of tuples to make Q false



$Q(u,w):-R'(u,x),S(x,y,z),T'(z,w)$

## 2/3. (Aggregated) Source side-effects

Delete min number of tuples to delete  $\geq k$  output tuples

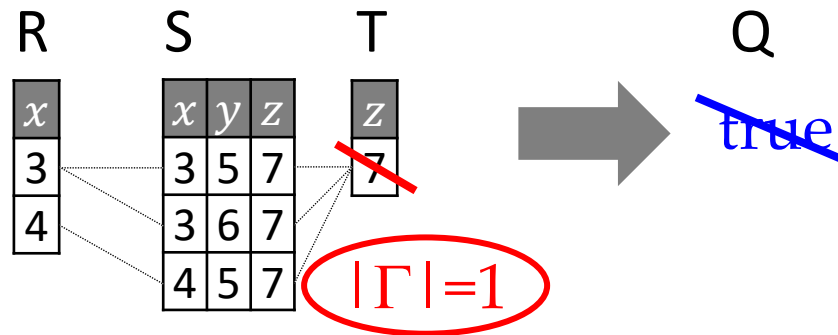


# Example Reverse Data Management Problems

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

## 1. Resilience

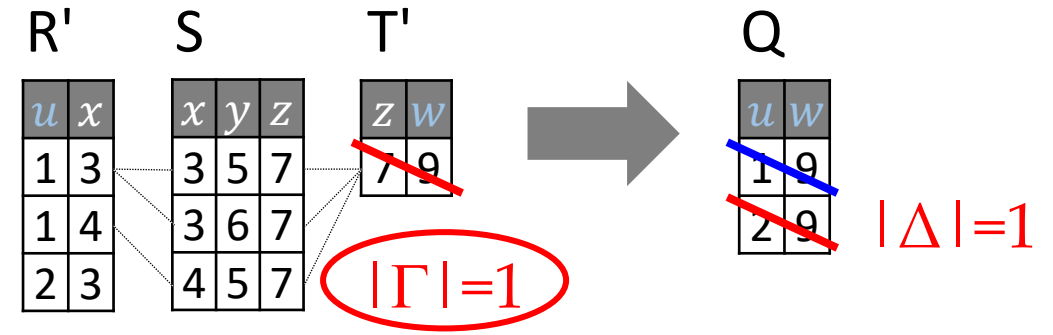
Delete min number of tuples to make Q false



Q(u,w):  $\neg R'(u,x), S(x,y,z), T'(z,w)$

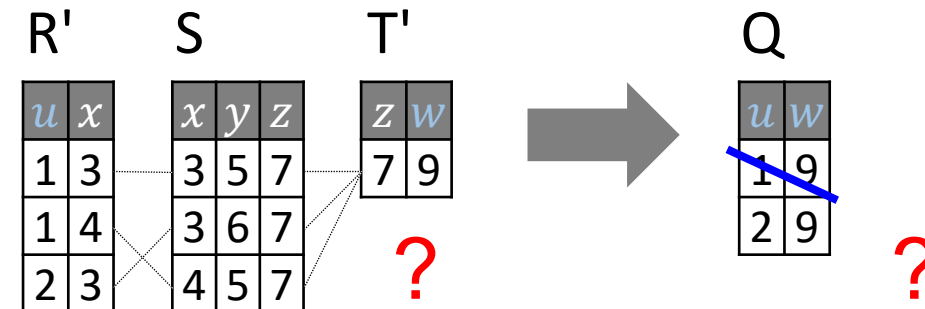
## 2/3. (Aggregated) Source side-effects

Delete min number of tuples to delete  $\geq k$  output tuples



## 4. View side-effects (in deletion propagation)

Delete tuples in order to delete an output tuple, while minimizing the other output tuples deleted



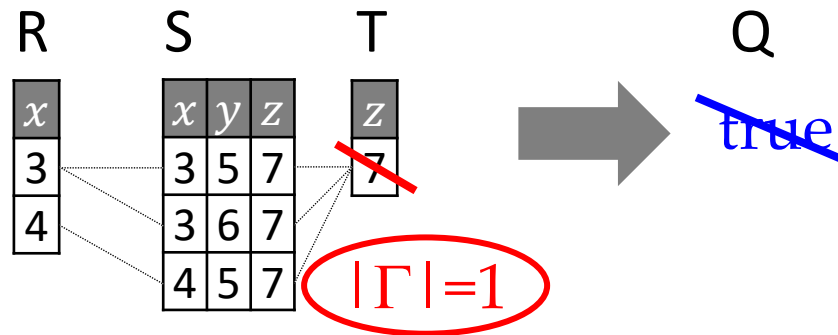


# Example Reverse Data Management Problems

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

## 1. Resilience

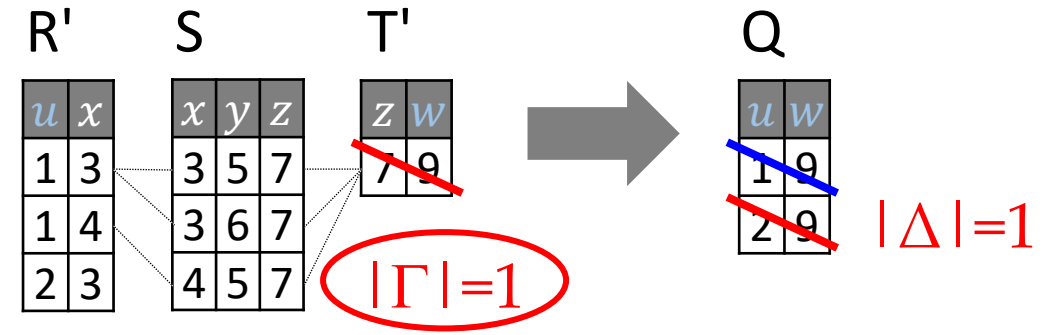
Delete min number of tuples to make Q false



Q(u,w):  $\neg R'(u,x), S(x,y,z), T'(z,w)$

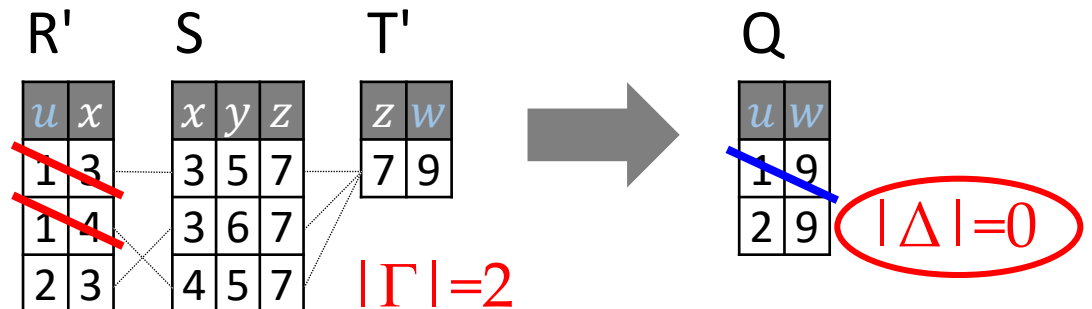
## 2/3. (Aggregated) Source side-effects

Delete min number of tuples to delete  $\geq k$  output tuples



## 4. View side-effects (in deletion propagation)

Delete tuples in order to delete an output tuple, while minimizing the other output tuples deleted

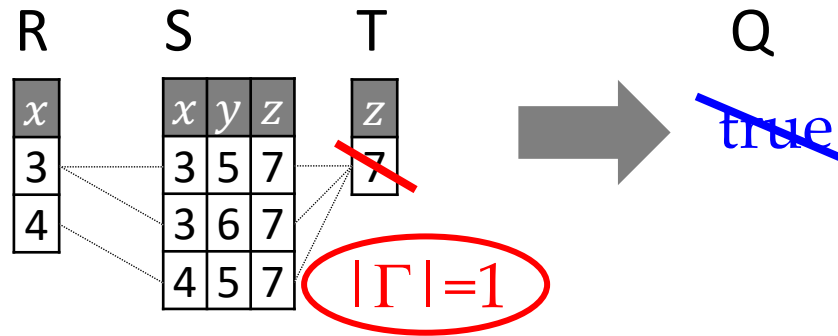


# Example Reverse Data Management Problems

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

## 1. Resilience

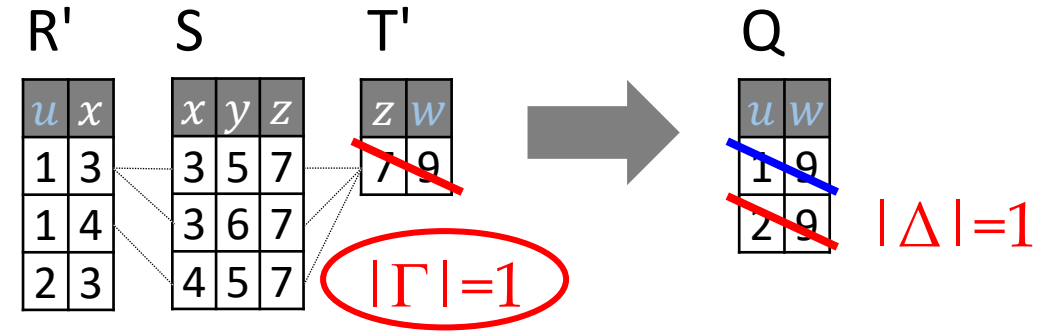
Delete min number of tuples to make Q false



Q(u,w):-R'(u,x),S(x,y,z),T'(z,w)

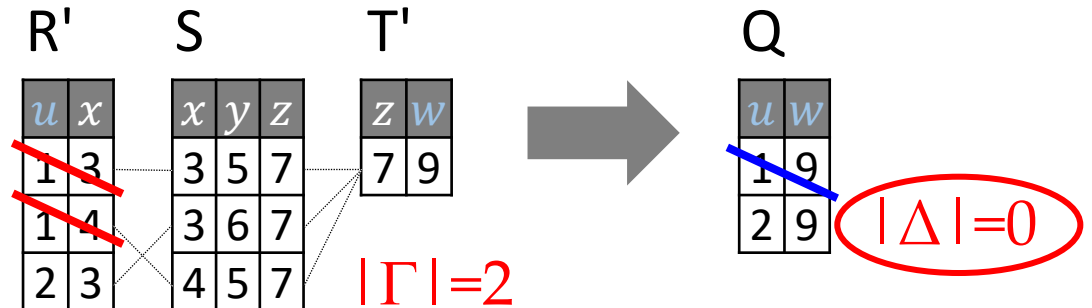
## 2/3. (Aggregated) Source side-effects

Delete min number of tuples to delete  $\geq k$  output tuples



## 4/5. (Aggregated) View side-effects

Delete tuples in order to delete  $\geq k$  output tuple, while minimizing the other output tuples deleted

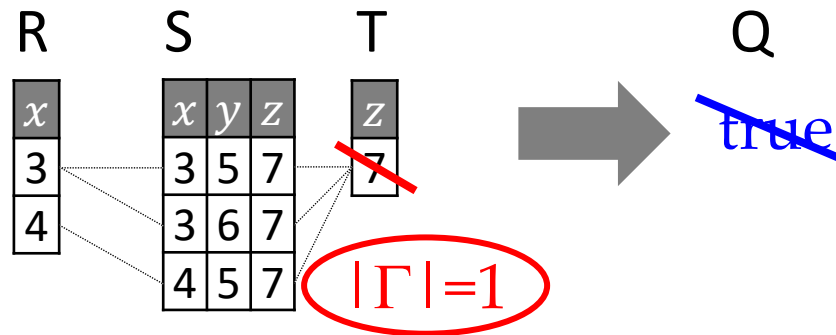


# Example Reverse Data Management Problems

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

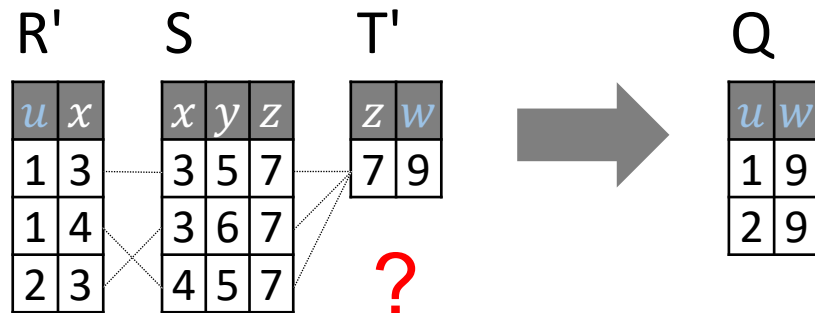
## 1. Resilience

Delete min number of tuples to make  $Q$  false



## 5. Smallest witness problem

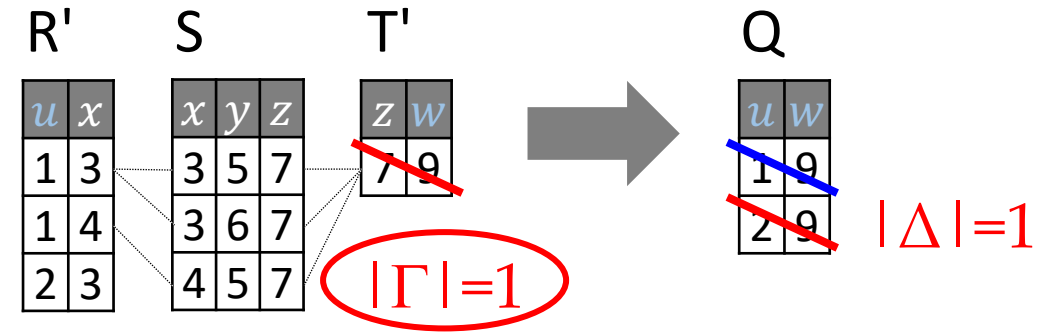
Delete max num. of tuples while keeping all output tuples



$Q(u,w):-R'(u,x),S(x,y,z),T'(z,w)$

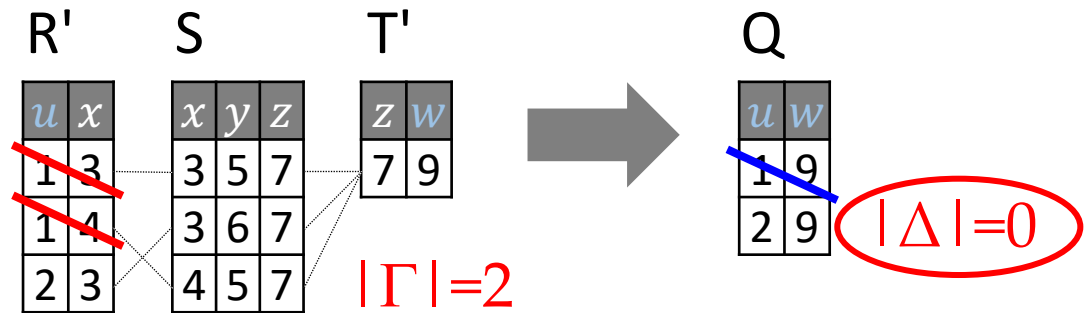
## 2/3. (Aggregated) Source side-effects

Delete min number of tuples to delete  $\geq k$  output tuples



## 4/5. (Aggregated) View side-effects

Delete tuples in order to delete  $\geq k$  output tuple, while minimizing the other output tuples deleted

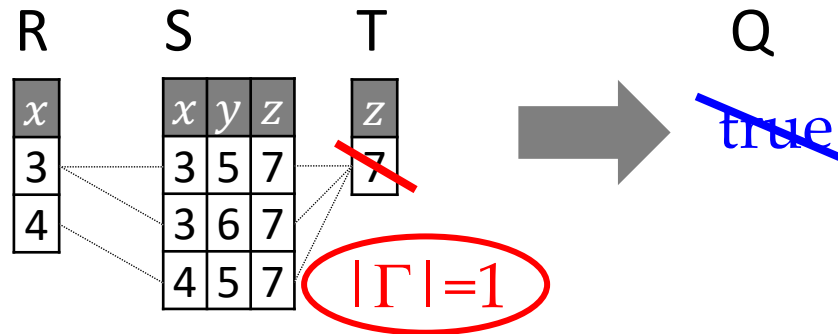


# Example Reverse Data Management Problems

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

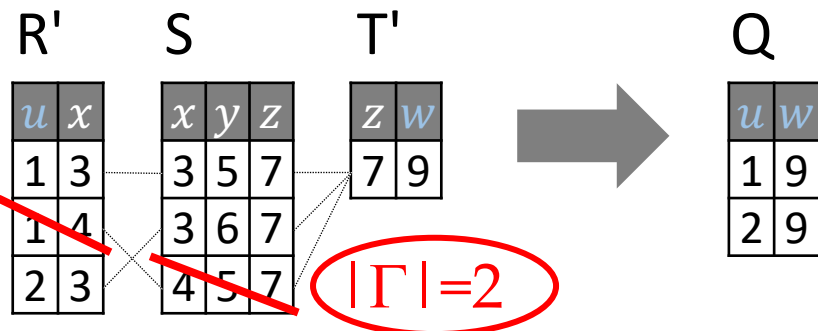
## 1. Resilience

Delete min number of tuples to make Q false



## 5. Smallest witness problem

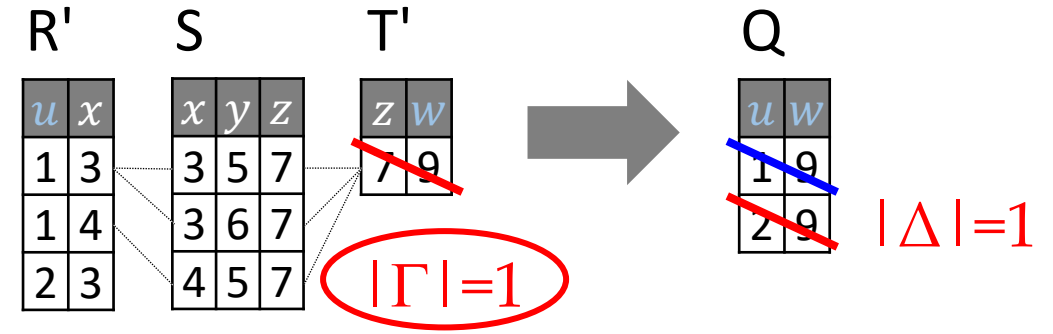
Delete max num. of tuples while keeping all output tuples



Q(u,w):-R'(u,x),S(x,y,z),T'(z,w)

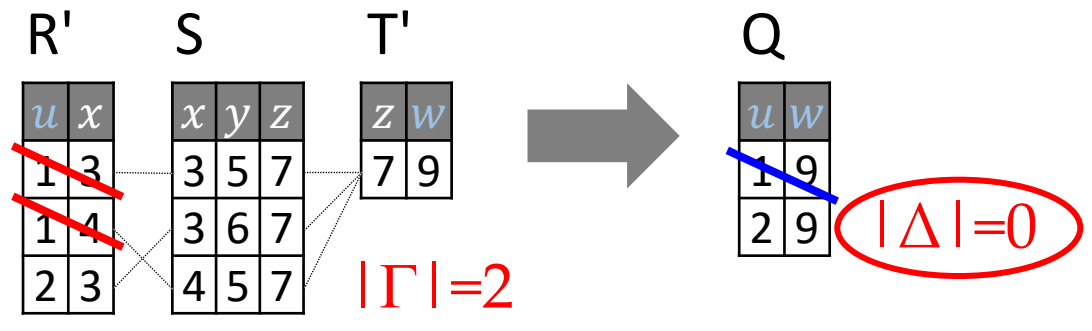
## 2/3. (Aggregated) Source side-effects

Delete min number of tuples to delete  $\geq k$  output tuples



## 4/5. (Aggregated) View side-effects

Delete tuples in order to delete  $\geq k$  output tuple, while minimizing the other output tuples deleted

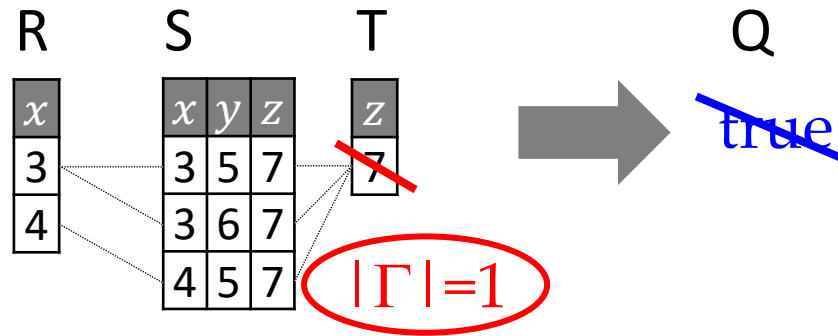


# Example Reverse Data Management Problems

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

## 1. Resilience

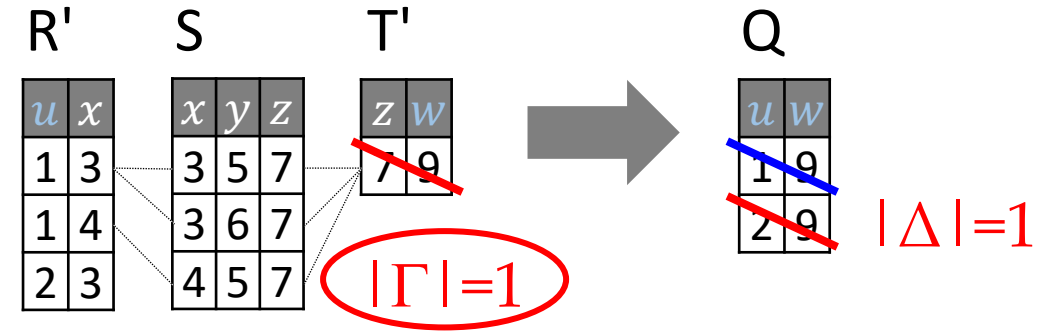
Delete min number of tuples to make Q false



Q(u,w):-R'(u,x),S(x,y,z),T'(z,w)

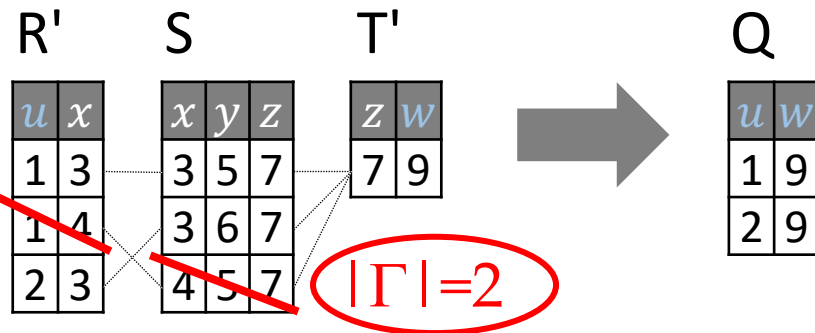
## 2/3. (Aggregated) Source side-effects

Delete min number of tuples to delete  $\geq k$  output tuples



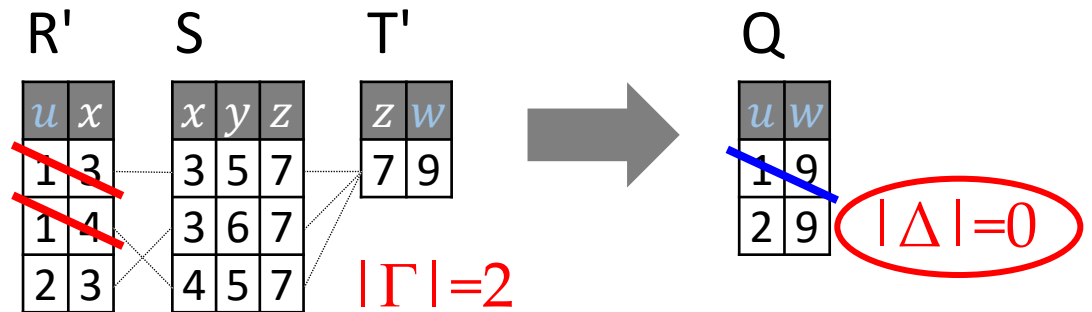
## 6./7 (Aggregated) Smallest witness problem

Delete max num. of tuples while keeping  $\geq k$  output tuples



## 4/5. (Aggregated) View side-effects

Delete tuples in order to delete  $\geq k$  output tuple, while minimizing the other output tuples deleted

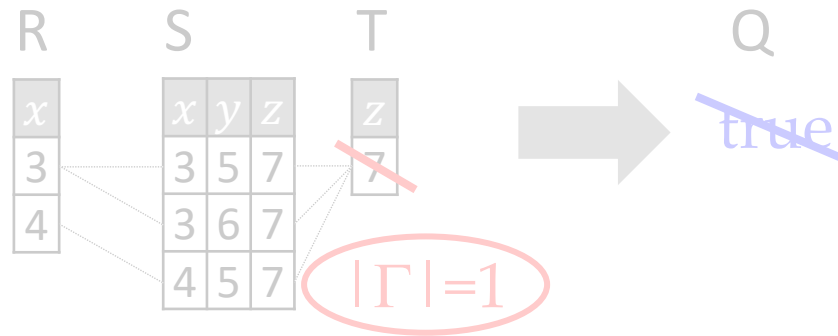


# Example Reverse Data Management Problems

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

## 1. Resilience

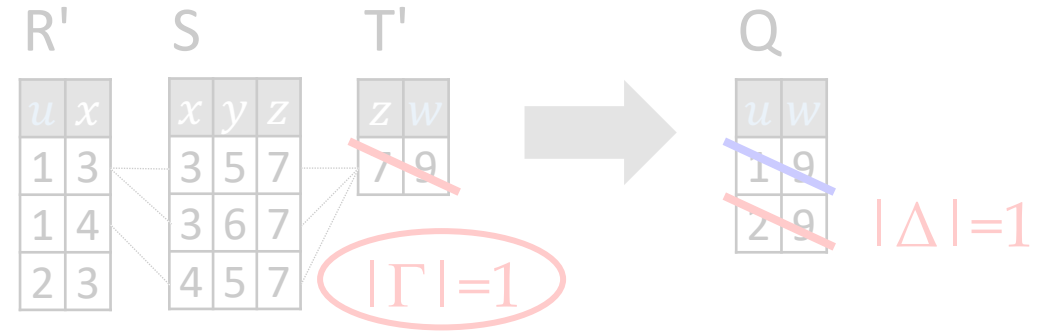
Delete min number of tuples to make  $Q$  false



$Q(u,w):-R'(u,x),S(x,y,z),T'(z,w)$

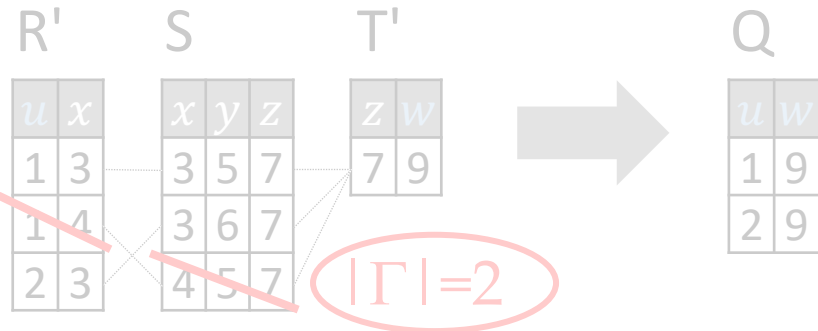
## 2/3. (Aggregated) Source side-effects

Delete min number of tuples to delete  $\geq k$  output tuples



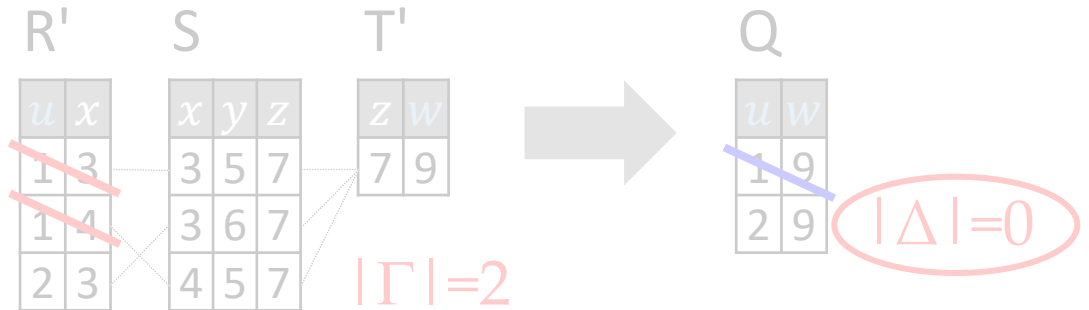
## 6./7 (Aggregated) Smallest witness problem

Delete max num. of tuples while keeping  $\geq k$  output tuples



## 4/5. (Aggregated) View side-effects

Delete tuples in order to delete  $\geq k$  output tuple, while minimizing the other output tuples deleted



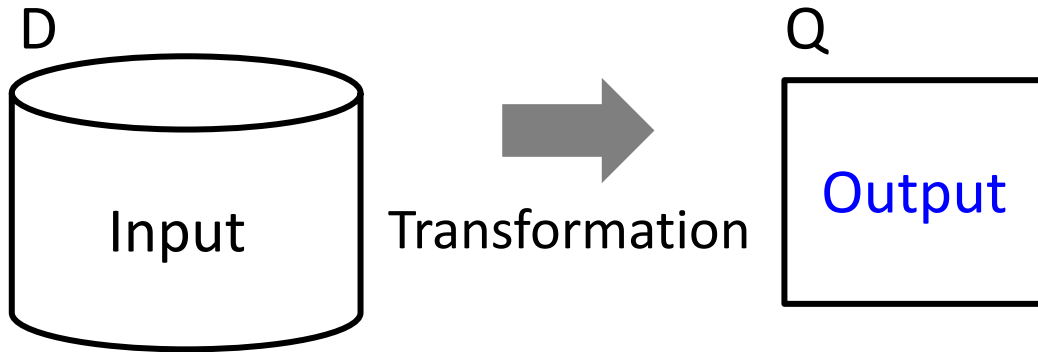
# Our topic today: Generalized Deletion Propagation

1. Unify and generalize these problems as instances of "generalized deletion propagation"
  - also allows new problem variants
2. Propose one algorithm to solve all these problems
  - including difficult cases, such as self-joins, or bag semantics
3. The algorithm terminates in PTIME for all known PTIME cases
  - including all known PTIME cases for self-join free queries for the problems of: resilience [VLDB'15], aggregated deletion propagation [VLDB'20], view-side effects [PODS'11], smallest witness problem [ICDT'24]

# Reverse Data Management for Conjunctive Queries

QUERY EVALUATION: *(Talks today in the morning)*

A transformation of the input to the output

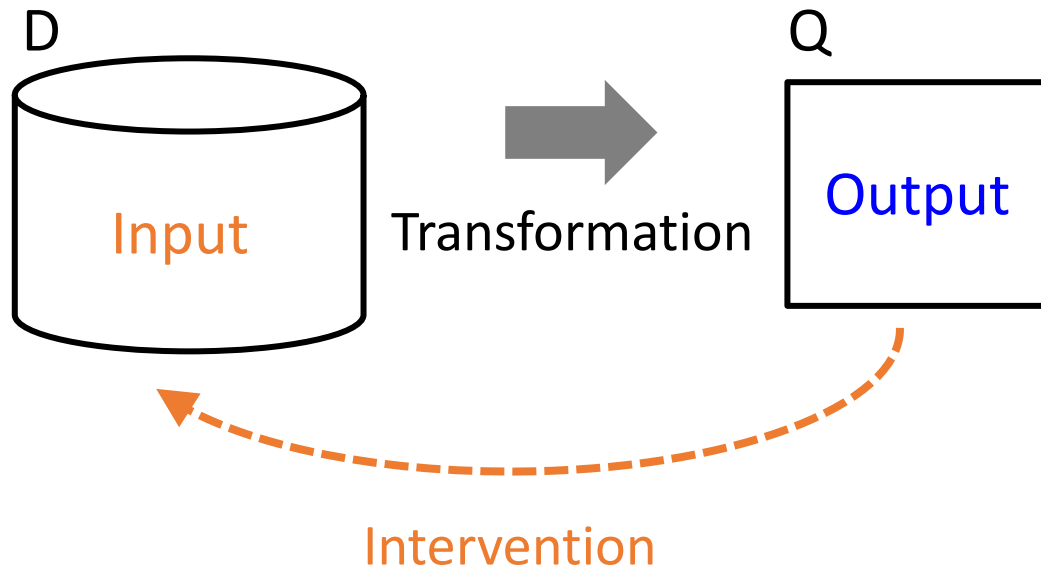




# Reverse Data Management

## QUERY EVALUATION:

A transformation of the input to the output



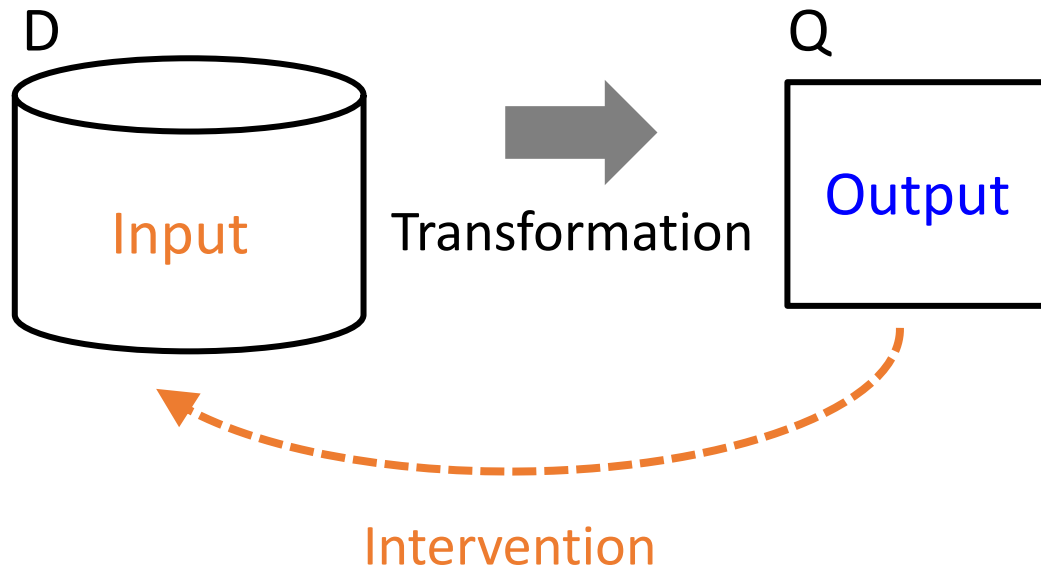
## REVERSE DATA MANAGEMENT:

What are the required changes to the input, in order to achieve a desired output?

# Reverse Data Management for Conjunctive Queries

## QUERY EVALUATION:

A transformation of the input to the output



## REVERSE DATA MANAGEMENT:

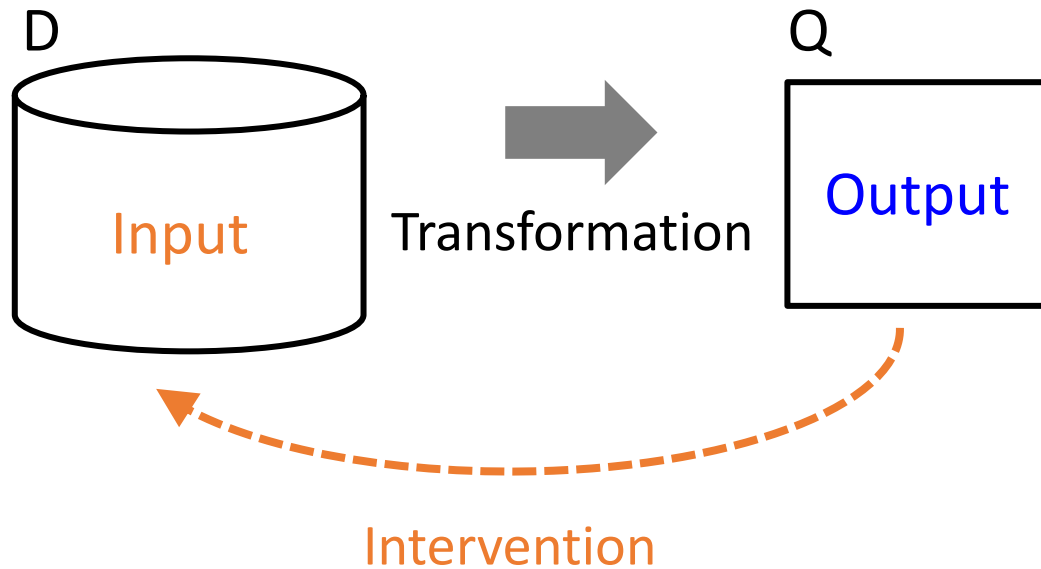
What are the required changes to the input, in order to achieve a desired output?

**1. Contrastive (CXp):** subset-minimal set of features (tuples) that are sufficient for changing a prediction (from true to false).

# Reverse Data Management for Conjunctive Queries

## QUERY EVALUATION:

A transformation of the input to the output



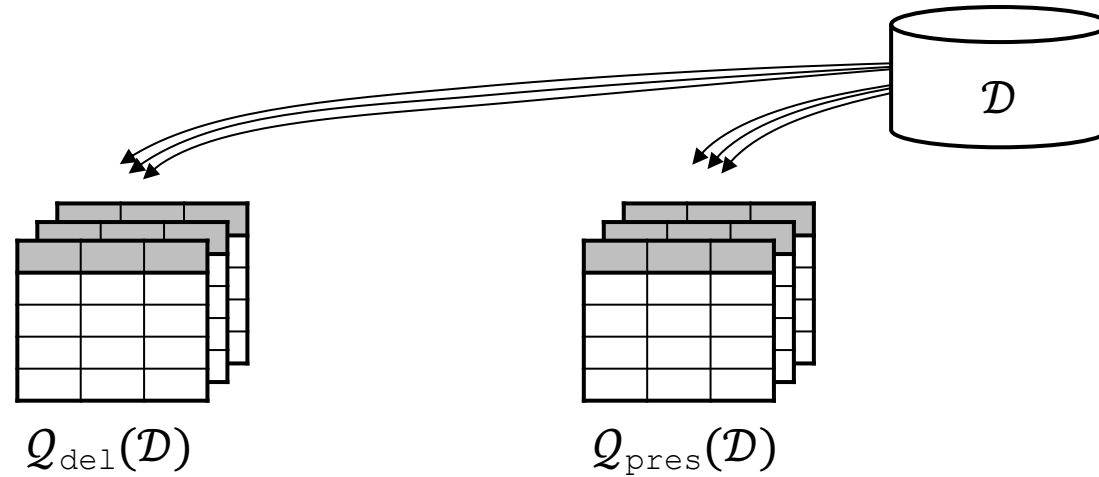
## REVERSE DATA MANAGEMENT:

What are the required changes to the input, in order to achieve a desired output?

**2. Abductive (AXp):** a subset-minimal set of features (tuples) that are sufficient for ensuring a certain prediction (i.e. fixing these tuples guarantees the output)  
= maximum deletions of tuples

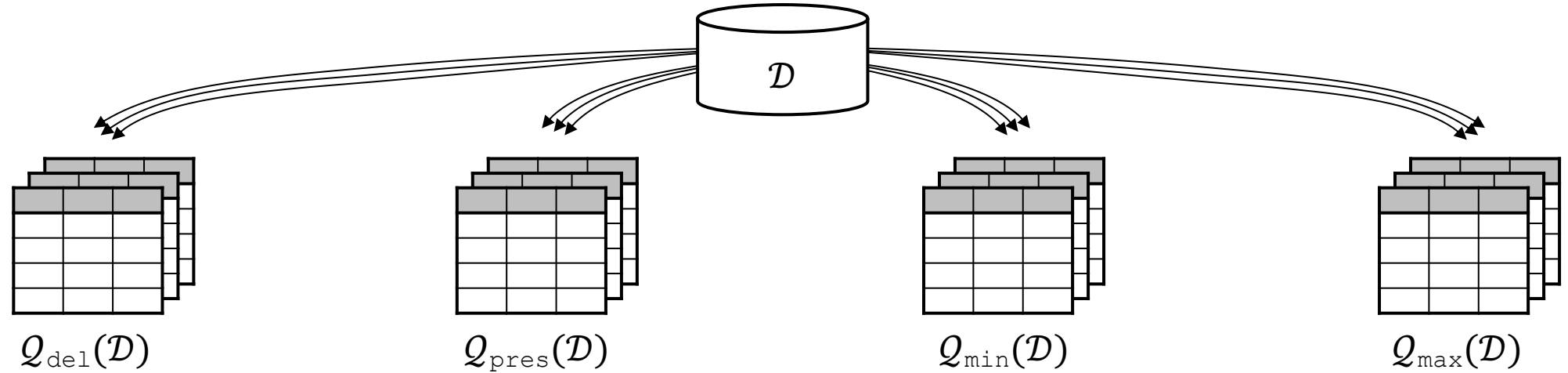
**1. Contrastive (CXp):** subset-minimal set of features (tuples) that are sufficient for changing a prediction (from true to false).  
= minimum deletions of tuples

# Generalized Deletion Propagation



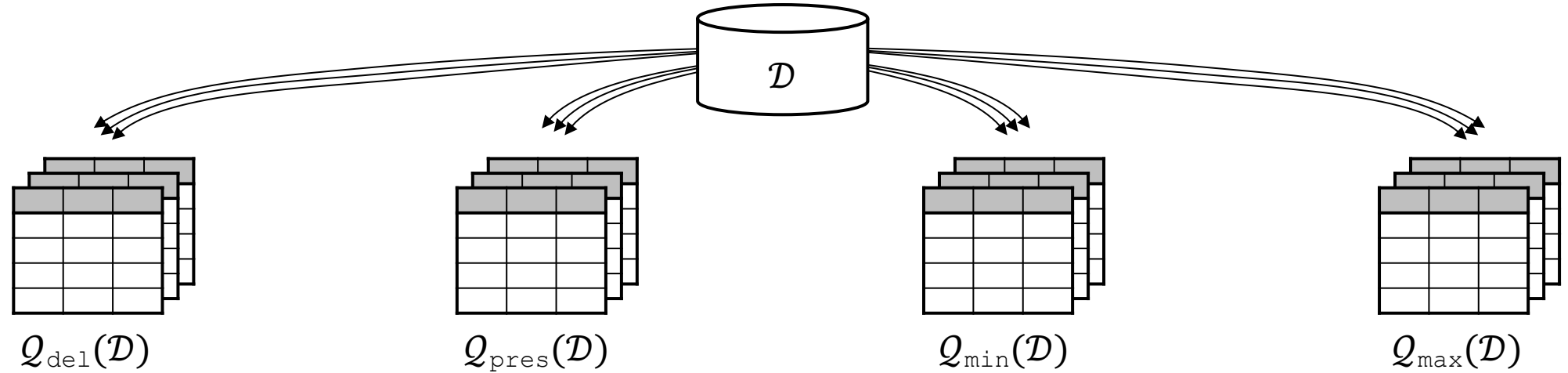
Propagate Requirements (hard constraints)	
<b>Delete</b> $\geq k_{\text{del}}^i$ tuples in $Q_{\text{del}}^i(\mathcal{D})$	<b>Preserve</b> $\geq k_{\text{pres}}^i$ tuples in $Q_{\text{pres}}^i(\mathcal{D})$

# Generalized Deletion Propagation



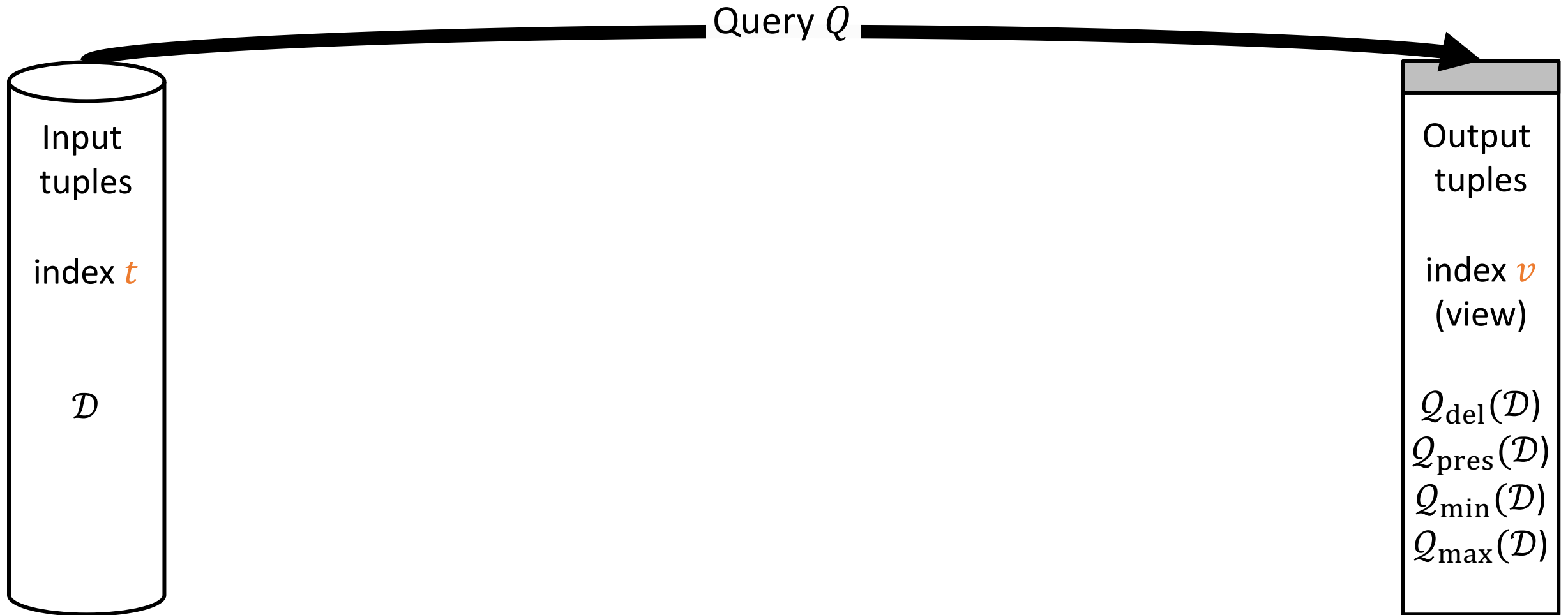
Propagate Requirements (hard constraints)		Optimize Side-Effect (soft constraints)	
<b>Delete</b>	<b>Preserve</b>	<b>Minimize</b>	<b>Maximize</b>
$\geq k_{\text{del}}^i$ tuples in $Q_{\text{del}}^i(\mathcal{D})$	$\geq k_{\text{pres}}^i$ tuples in $Q_{\text{pres}}^i(\mathcal{D})$	Deletions in $Q_{\text{min}}$	Deletions in $Q_{\text{max}}$

# Generalized Deletion Propagation



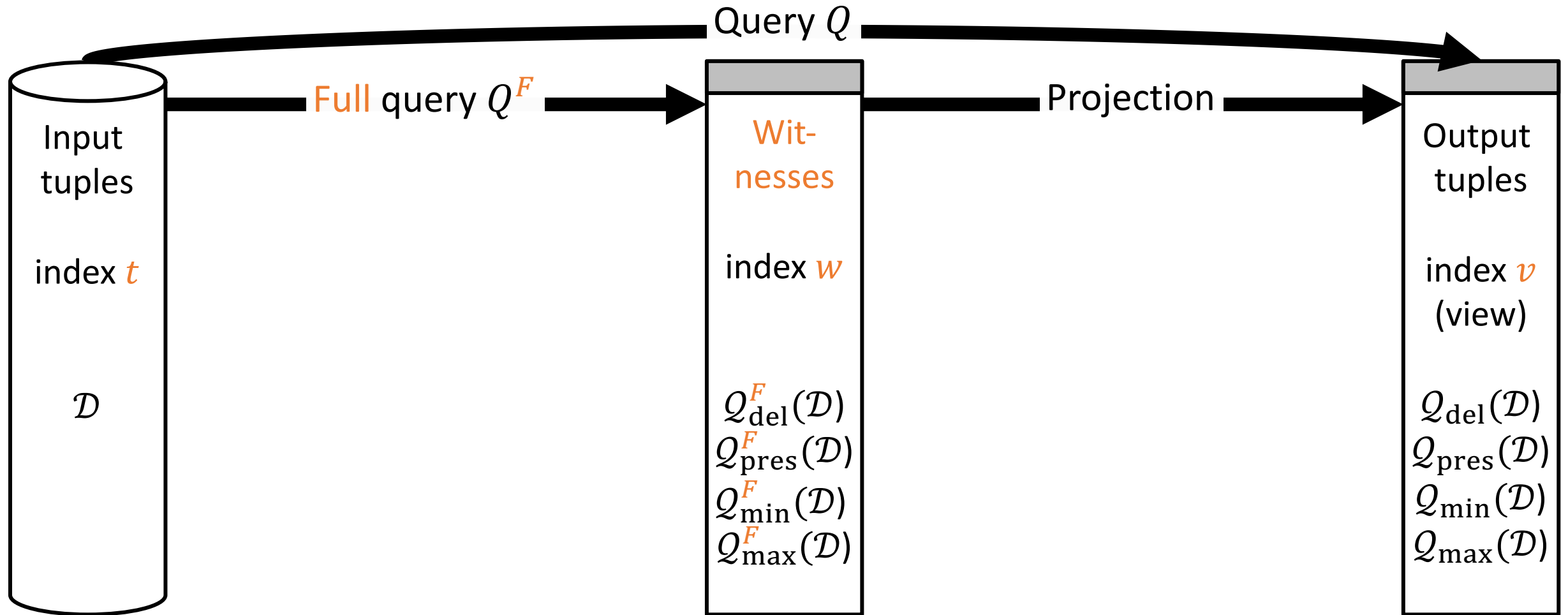
	Propagate Requirements (hard constraints)		Optimize Side-Effect (soft constraints)	
	Delete $\geq k_{\text{del}}^i$ tuples in $Q_{\text{del}}^i(\mathcal{D})$	Preserve $\geq k_{\text{pres}}^i$ tuples in $Q_{\text{pres}}^i(\mathcal{D})$	Minimize Deletions in $Q_{\text{min}}$	Maximize Deletions in $Q_{\text{max}}$
DP-SS	All tuples in $Q_{\text{del}}(\mathcal{D})$	-	Deletions in $Q_{\text{min}} = \{R \in \mathcal{D}\}$	-
ADP-SS	$\geq k$ tuples in $Q_{\text{del}}(\mathcal{D})$	-	Deletions in $Q_{\text{min}} = \{R \in \mathcal{D}\}$	-
DP-VS	All tuples in $Q_{\text{del}}(\mathcal{D})$	-	Deletions in $Q_{\text{min}}(\mathcal{D}) \supseteq Q_{\text{del}}(\mathcal{D})$	-
SWP	-	All tuples in $Q_{\text{pres}}(\mathcal{D})$	-	Deletions in $Q_{\text{max}} = \{R \in \mathcal{D}\}$

# Propagation Constraints in the Language of ILPs



*$X[t] = 1$  means that  $t$  is deleted*

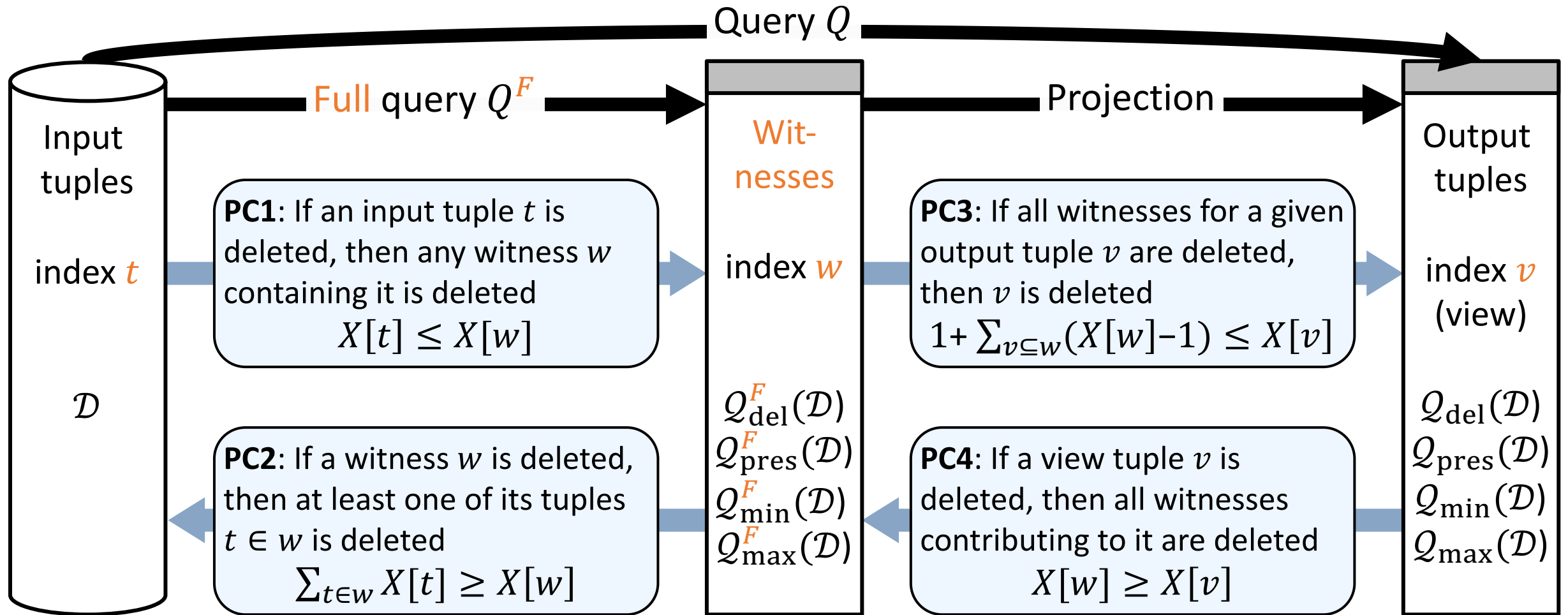
# Propagation Constraints in the Language of ILPs



*$X[t] = 1$  means that  $t$  is deleted*

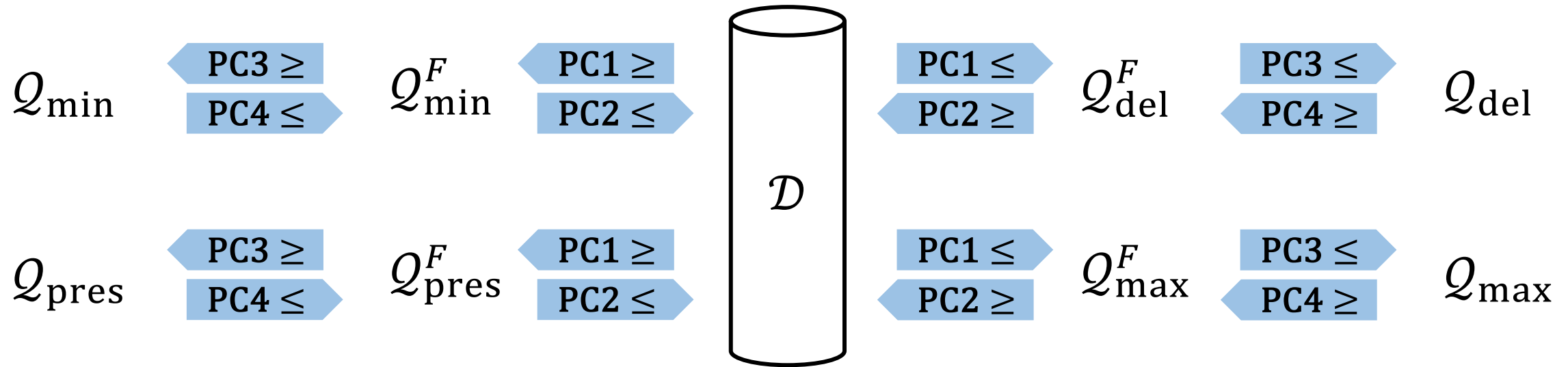


# Propagation Constraints in the Language of ILPs

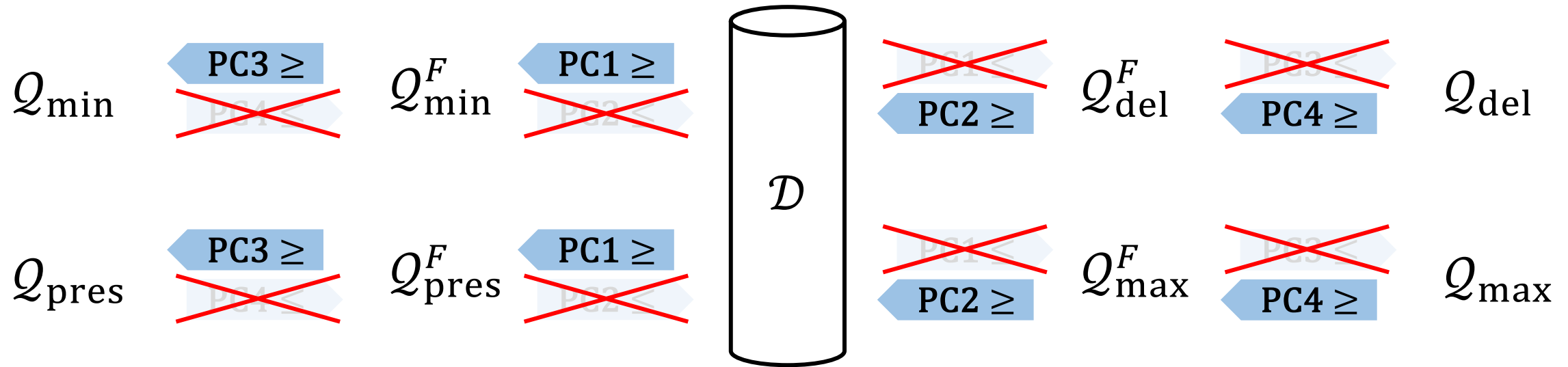


$X[t] = 1$  means that  $t$  is deleted

# Propagation constraints: Naive ILP (1/3)



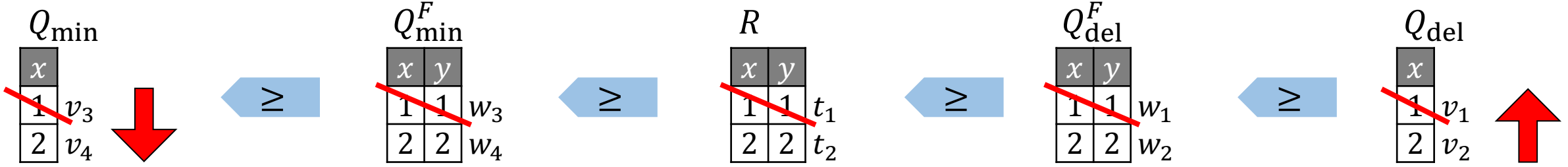
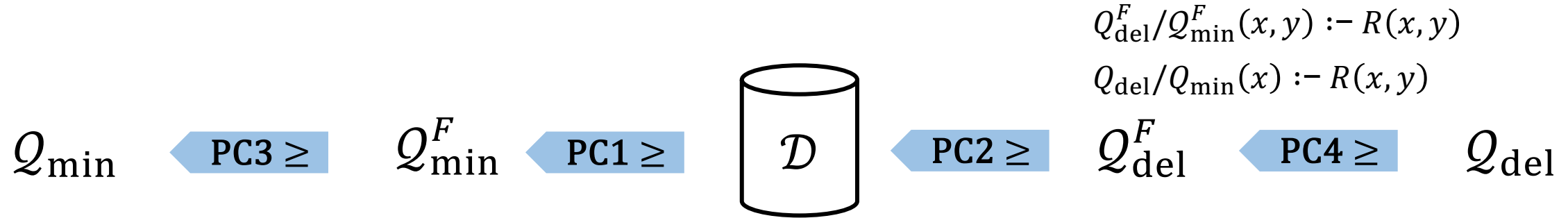
# Propagation constraints: Wildcard ILP (2/3)



Certain constraints don't need to be explicitly enforced for a given type of view.

Removing the unnecessary constraints \*increases\* the solution space (from the point of view of the ILP formulation) but keeps all the optimal solutions (projected on the actual interventions on  $\mathcal{D}$  that we care about).

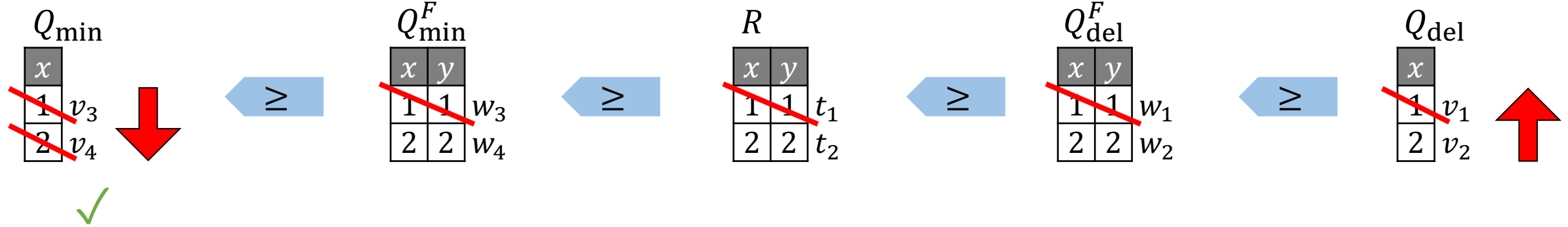
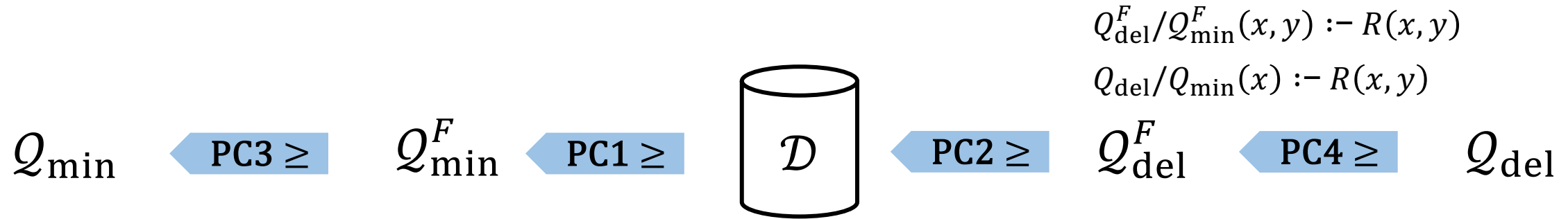
# Propagation constraints: Wildcard ILP (2/3)



Soft constraint: we want to minimize the number of deletions from  $Q_{\min}$

User constraint: we are required to delete a given number of tuples from  $Q_{\text{del}}$

# Propagation constraints: Wildcard ILP (2/3)



Soft constraint: we want to minimize the number of deletions from  $Q_{\min}$

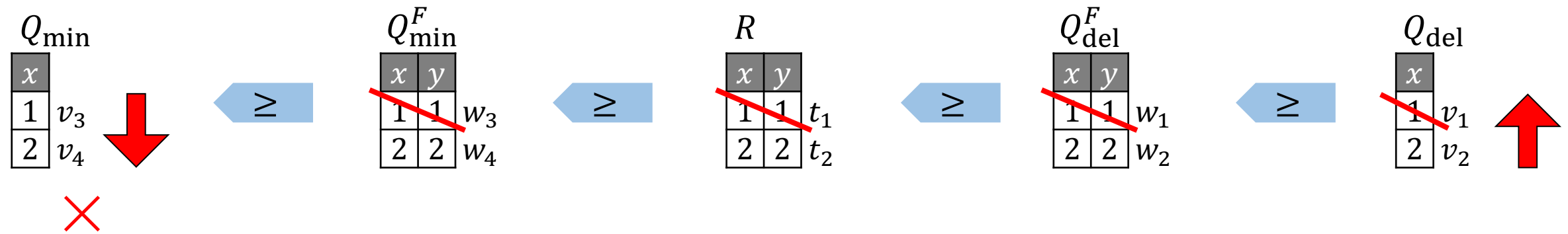
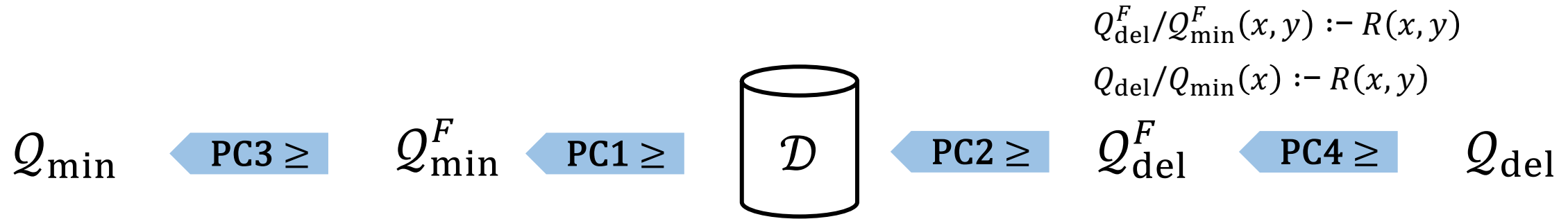
User constraint: we are required to delete a given number of tuples from  $Q_{\min}$

The variables are an upper bound and preservations (X=0) need to get propagated from left to right.

We can set more variables than necessary to 1, but that just penalizes our solution.

The variables are a lower bound and deletions (X=1) need to get propagated from right to left.

# Propagation constraints: Wildcard ILP (2/3)



Soft constraint: we want to minimize the number of deletions from  $Q_{\min}$

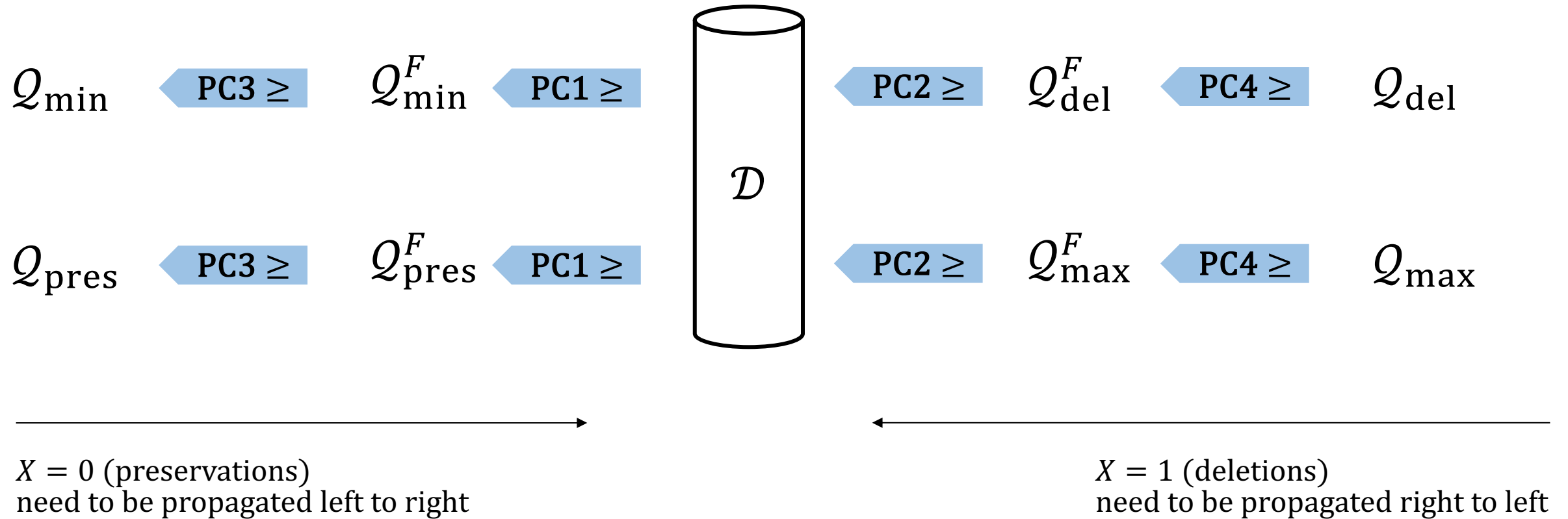
User constraint: we are required to delete a given number of tuples from  $Q_{\min}$

The variables are an upper bound and preservations ( $X=0$ ) need to get propagated from left to right.

We can set more variables than necessary to 1, but that just penalizes our solution.

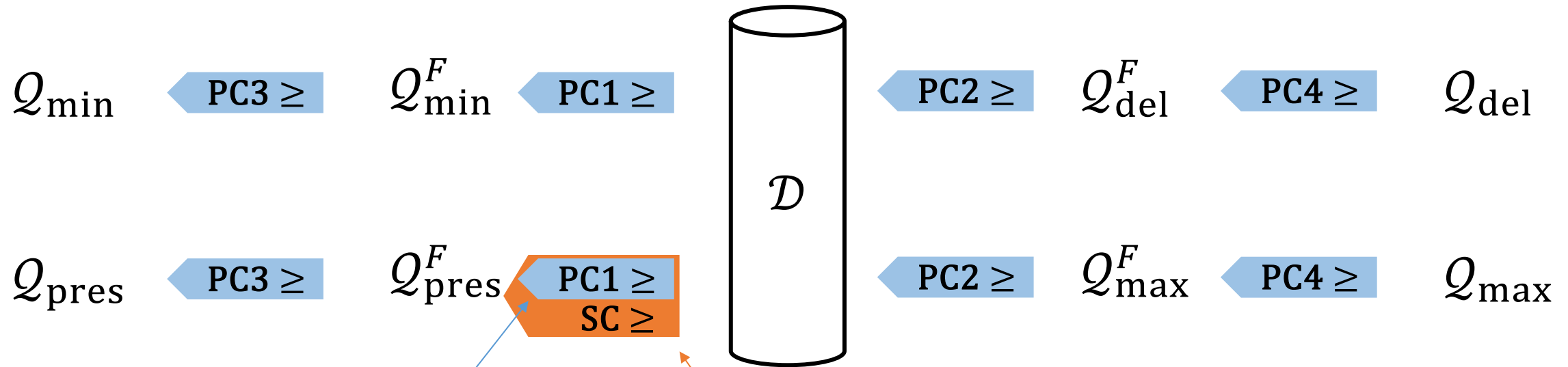
The variables are a lower bound and deletions ( $X=1$ ) need to get propagated from right to left.

# Propagation constraints: Wildcard ILP (2/3)



Removing the unnecessary constraints \*increases\* the solution space (from the point of view of the ILP formulation) but keeps all the optimal solutions (projected on the actual interventions on  $\mathcal{D}$  that we care about).

# Propagation constraints: Smoothened ILP (3/3)



**PC1:** If an input tuple  $t$  is deleted, then any witness  $w$  containing it is deleted  
 $X[t] \leq X[w]$

**SC:** If an input tuple  $t$  is not deleted, then at least one of its witnesses is preserved, for each view tuple  $v$   

$$X[t] \leq 1 + \sum_{\substack{w: t \in w \\ w \supseteq v}} (X[w]-1)$$



# Example: Smoothing for Smallest Witness Problem DETAILS

$$Q_{\text{pres}}(x) := R(x, y), S(x)$$

$Q_{\text{pres}}$
$x$
$1$

 $v_1$ 

$R$	
$x$	$y$
$1$	$1$
$1$	$2$

 $t_1$   
 $t_2$ 

$S$
$x$
$1$

 $t_3$ 

$$Q_{\text{pres}}^F(x, y) := R(x, y), S(x)$$

$x$	$y$
$1$	$1$
$1$	$2$

 $w_1$   
 $w_2$ 

Optimal integral solution:

$$f^* = -1, \mathbf{X}^*[t_1, t_2, t_3] = (1, 0, 0) \text{ or } (0, 1, 0)$$

# Example: Smoothing for Smallest Witness Problem DETAILS

$$Q_{\text{pres}}(x) := R(x, y), S(x)$$

$Q_{\text{pres}}$
$x$
$1$

 $v_1$ 

$R$	
$x$	$y$
$1$	$1$
$1$	$2$

 $t_1$   
 $t_2$ 

$S$
$x$
$1$

 $t_3$ 

$$Q_{\text{pres}}^F(x, y) := R(x, y), S(x)$$

$x$	$y$	
$1$	$1$	$w_1$
$1$	$2$	$w_2$

Optimal integral solution:

$$f^* = -1, \mathbf{X}^*[t_1, t_2, t_3] = (1, 0, 0) \text{ or } (0, 1, 0)$$

# Example: Smoothing for Smallest Witness Problem DETAILS

$$Q_{\text{pres}}(x) := R(x, y), S(x)$$

$Q_{\text{pres}}$	$R$	$S$
$x$	$x$ $y$	$x$
$1$ $v_1$	$1$ $1$ $t_1$	$1$ $t_3$
	$1$ $2$ $t_2$	

$$Q_{\text{pres}}^F(x, y) := R(x, y), S(x)$$

$x$	$y$	
$1$	$1$	$w_1$
$1$	$2$	$w_2$

## Naive ILP formulation

Minimization objective

$$f(\mathbf{X}) = -X[t_1] - X[t_2] - X[t_3] \leftarrow \text{Delete as many tuples as possible...}$$

User constraint

$$X[w_1] + X[w_2] \leq 1 \leftarrow \text{... but preserve at least one witness.}$$

Propagation Constraints (PCs)\*

$$X[t_1] \leq X[w_1]$$

$$X[t_2] \leq X[w_2]$$

$$X[t_3] \leq X[w_1]$$

$$X[t_3] \leq X[w_2]$$

$X[t] = 1$  means that  $t$  is deleted

$\leftarrow$  Delete as many tuples as possible...

$\leftarrow$  ... but preserve at least one witness.

$\leftarrow$  If a tuple is deleted, then so is the corresponding witness

Optimal integral solution:

$$f^* = -1, \mathbf{X}^*[t_1, t_2, t_3] = (1, 0, 0) \text{ or } (0, 1, 0)$$

\*To simplify, ILP formulation is slightly simplified from the actual formulation (PC2, PC3, PC4 are moved into the user constraint). See example 3 in paper for details.

# Example: Smoothing for Smallest Witness Problem DETAILS

$$Q_{\text{pres}}(x) := R(x, y), S(x)$$

x
1

 $v_1$ 

x	y
1	1
1	2

 $t_1$   
 $t_2$ 

x
1

 $t_3$ 

$$Q_{\text{pres}}^F(x, y) := R(x, y), S(x)$$

x	y
1	1
1	2

 $w_1$   
 $w_2$ 

## Naive ILP formulation

$X[t] = 1$  means that  $t$  is deleted

Minimization objective

$$f(\mathbf{X}) = -X[t_1] - X[t_2] - X[t_3] \leftarrow \text{Delete as many tuples as possible...}$$

0.5      0.5      0.5

User constraint

$$X[w_1] + X[w_2] \leq 1 \leftarrow \text{... but preserve at least one witness.}$$

0.5      0.5

Propagation Constraints (PCs)\*

$$\begin{aligned} 0.5 X[t_1] &\leq X[w_1] & 0.5 \\ 0.5 X[t_2] &\leq X[w_2] & 0.5 \\ 0.5 X[t_3] &\leq X[w_1] & 0.5 \\ 0.5 X[t_3] &\leq X[w_2] & 0.5 \end{aligned}$$

$\leftarrow$  If a tuple is deleted, then so is the corresponding witness

Optimal integral solution (ILP)

$$f^* = -1, \mathbf{X}^*[t_1, t_2, t_3] = (1, 0, 0) \text{ or } (0, 1, 0)$$

Optimal fractional solution (LP relaxation)

$$f^* = -1.5, \mathbf{X}^*[t_1, t_2, t_3] = (0.5, 0.5, 0.5)$$



# Example: Smoothing for Smallest Witness Problem DETAILS

$$Q_{\text{pres}}(x) := R(x, y), S(x)$$

$x$
1

 $v_1$ 

$x$	$y$
1	1
1	2

 $t_1$   
 $t_2$ 

$x$
1

 $t_3$ 

$$Q_{\text{pres}}^F(x, y) := R(x, y), S(x)$$

$x$	$y$
1	1
1	2

 $w_1$   
 $w_2$ 

## Smoothened ILP formulation

$X[t] = 1$  means that  $t$  is deleted

Minimization objective

$$f(\mathbf{X}) = -X[t_1] - X[t_2] - X[t_3] \leftarrow \text{Delete as many tuples as possible...}$$

User constraint

$$X[w_1] + X[w_2] \leq 1 \leftarrow \dots \text{ but preserve at least one witness.}$$

Propagation Constraints (PCs)\*

$$X[t_1] \leq X[w_1]$$

$$X[t_2] \leq X[w_2]$$

$$X[t_3] \leq X[w_1]$$

$$X[t_3] \leq X[w_2]$$

If a tuple is deleted, then so is the corresponding witness

Smoothing constraint (SCs)

$$X[t_3] \leq X[w_1] + X[w_2] - 1$$

If a tuple is not deleted, then at least one of its witnesses is preserved

$$0.5 \not\leq 0.5 + 0.5 - 1$$

The smoothing constraint eliminates the fractional solution with  $f^* = -1.5$

Optimal integral solution (ILP)

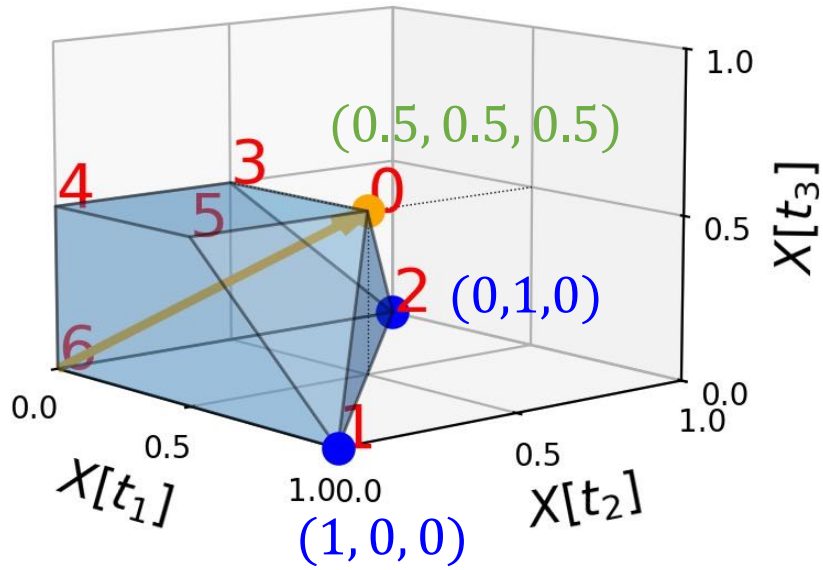
= optimal fractional solution (LP relaxation)

$$f^* = -1, \mathbf{X}^*[t_1, t_2, t_3] = (1, 0, 0) \text{ or } (0, 1, 0)$$

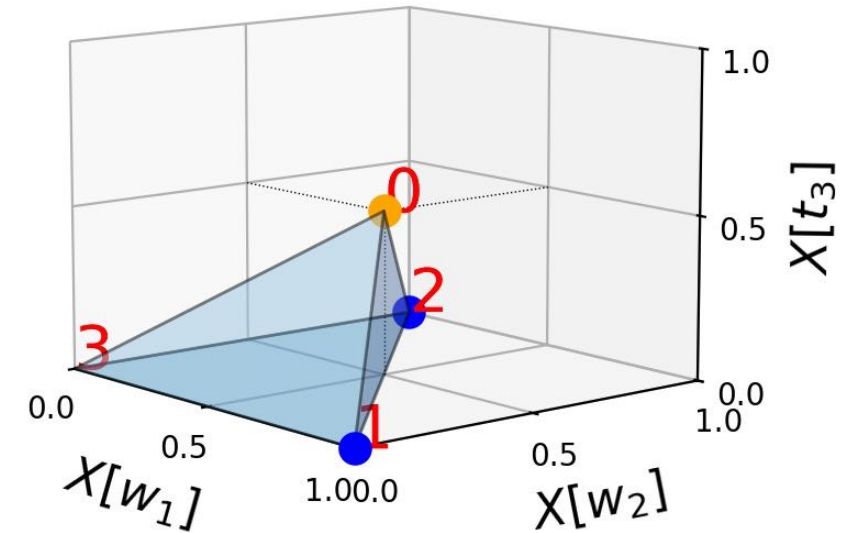


\*ILP formulation is slightly simplified from the actual formulation (PC2, PC3, PC4 are moved into the user constraint). See example 3 in paper for details.

# Example: Smoothing for Smallest Witness Problem DETAILS



- 0:  $\mathbf{X}[\mathbf{t}] = (0.5, 0.5, 0.5)$
- 1:  $\mathbf{X}[\mathbf{t}] = (1, 0, 0)$
- 2:  $\mathbf{X}[\mathbf{t}] = (0, 1, 0)$

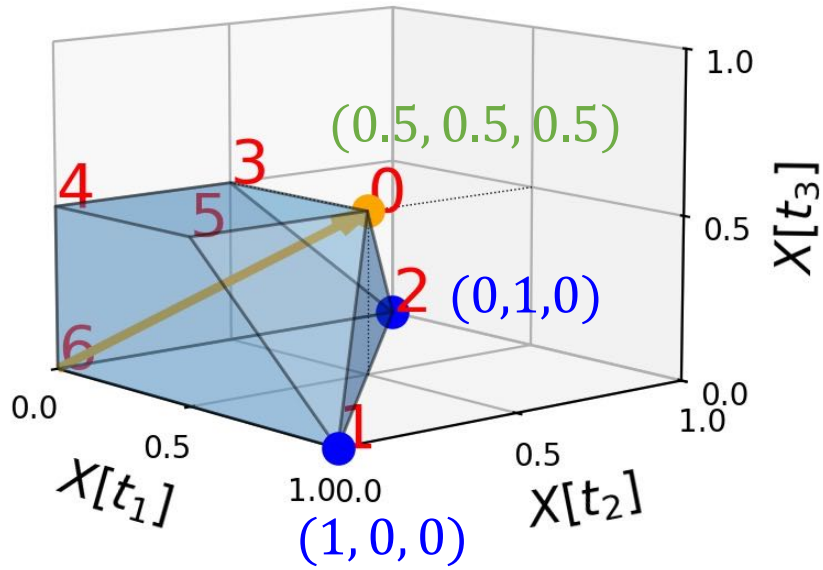


Smoothing constraint (SCs)

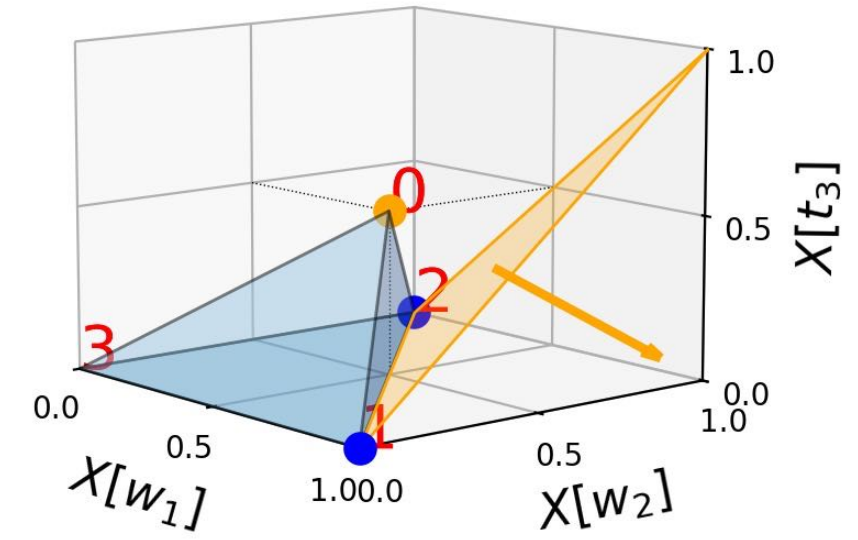
$$X[t_3] \leq X[w_1] + X[w_2] - 1$$

$$0.5 \not\leq 0.5 + 0.5 - 1$$

# Example: Smoothing for Smallest Witness Problem DETAILS



- 0:  $\mathbf{x}[t] = (0.5, 0.5, 0.5)$
- 1:  $\mathbf{x}[t] = (1, 0, 0)$
- 2:  $\mathbf{x}[t] = (0, 1, 0)$

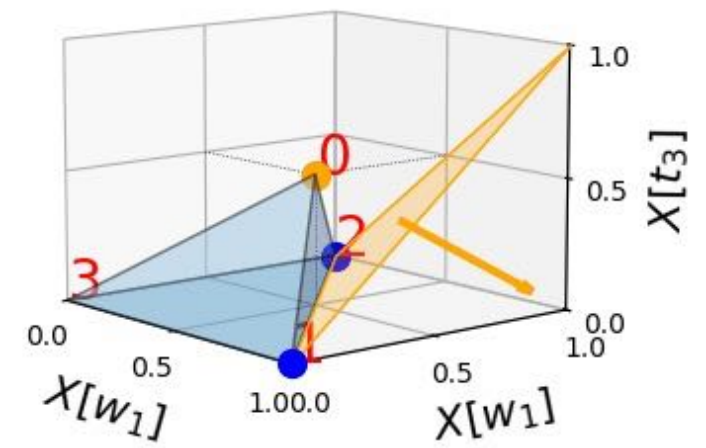


The smoothing constraint eliminates the fractional solution with  $f^* = -1.5$

Smoothing constraint (SCs)

$$X[t_3] \leq X[w_1] + X[w_2] - 1$$

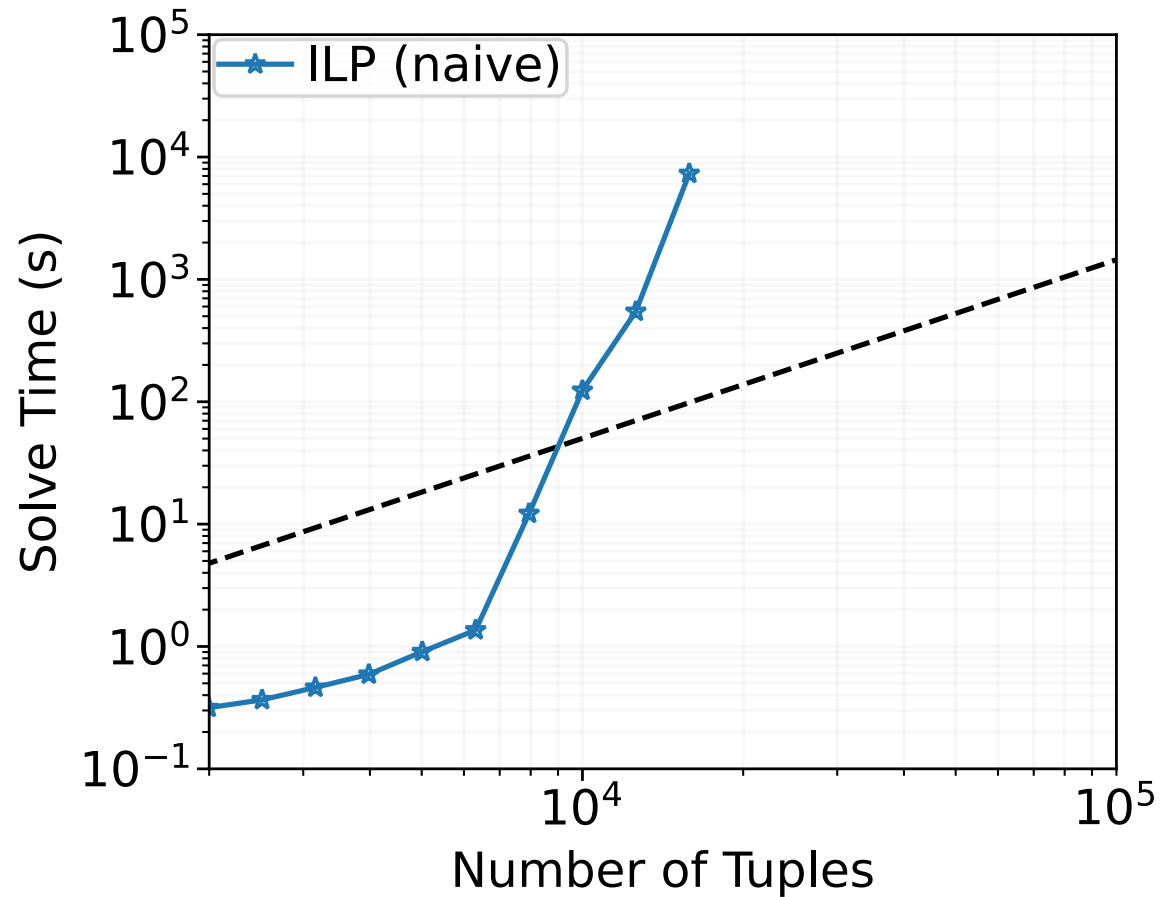
$$0.5 \not\leq 0.5 + 0.5 - 1$$



# Scalability of Naive vs. Smoothened ILP for PTIME query

Finding an optimal solution for the Smallest Witness Problem

$Q_{5\text{star}}(x): -R(x, a), S(x, b), T(x, c), U(x, d), V(x, e), A(x)$

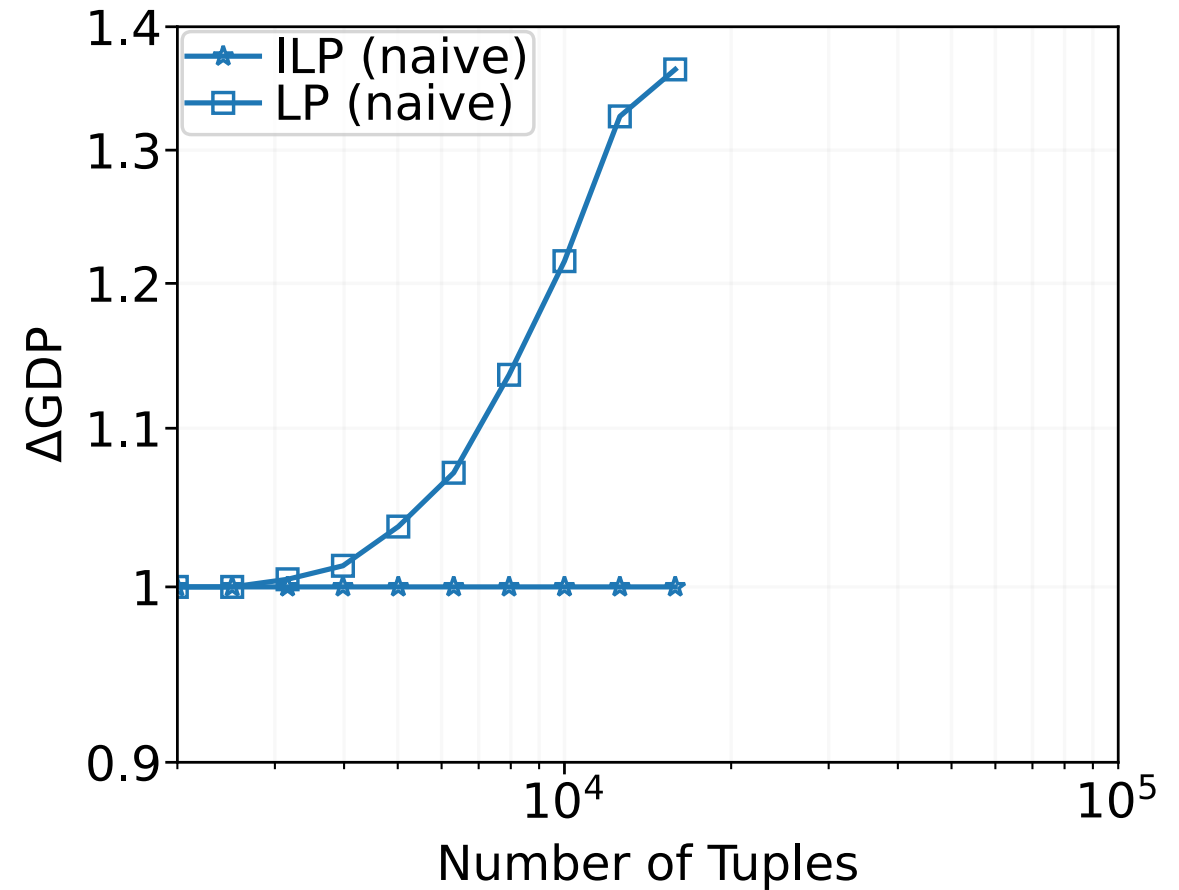
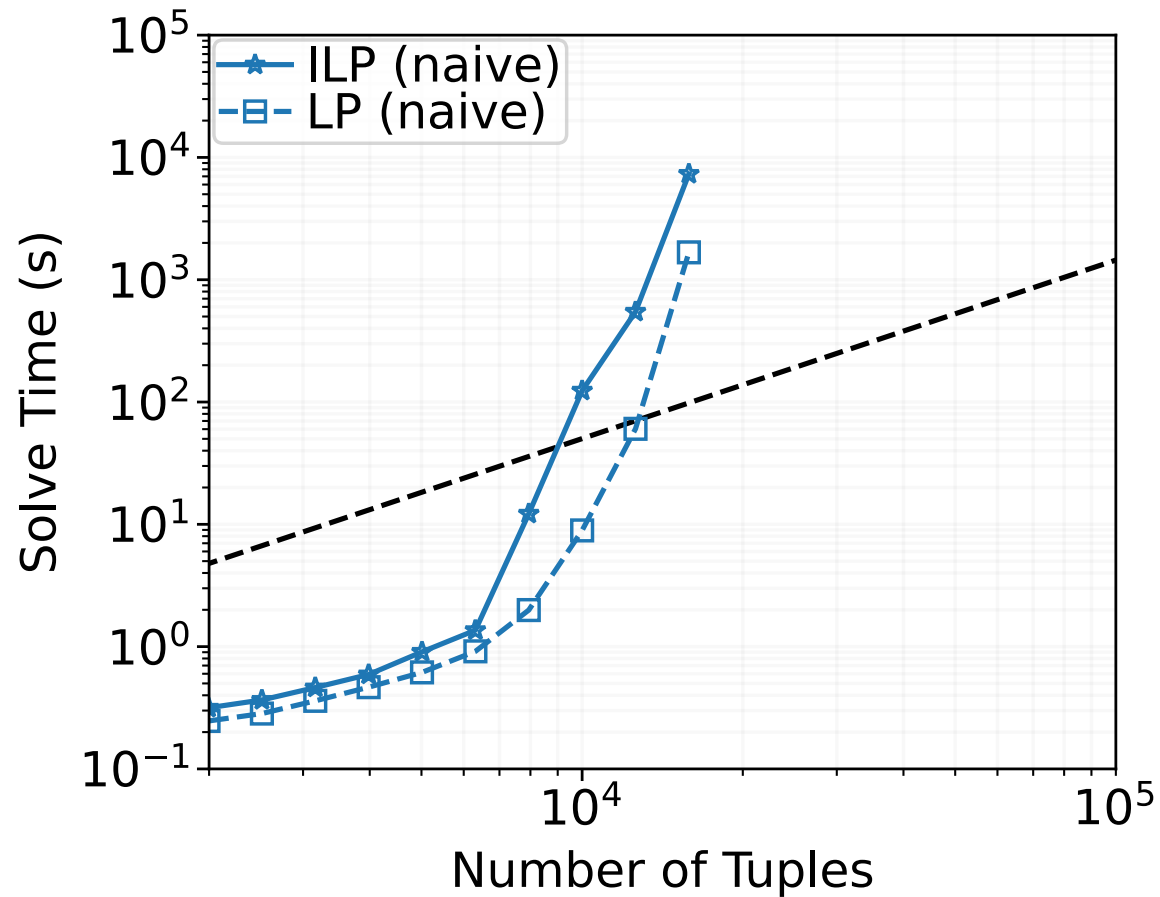




# Scalability of Naive vs. Smoothened ILP for PTIME query

Finding an optimal solution for the Smallest Witness Problem

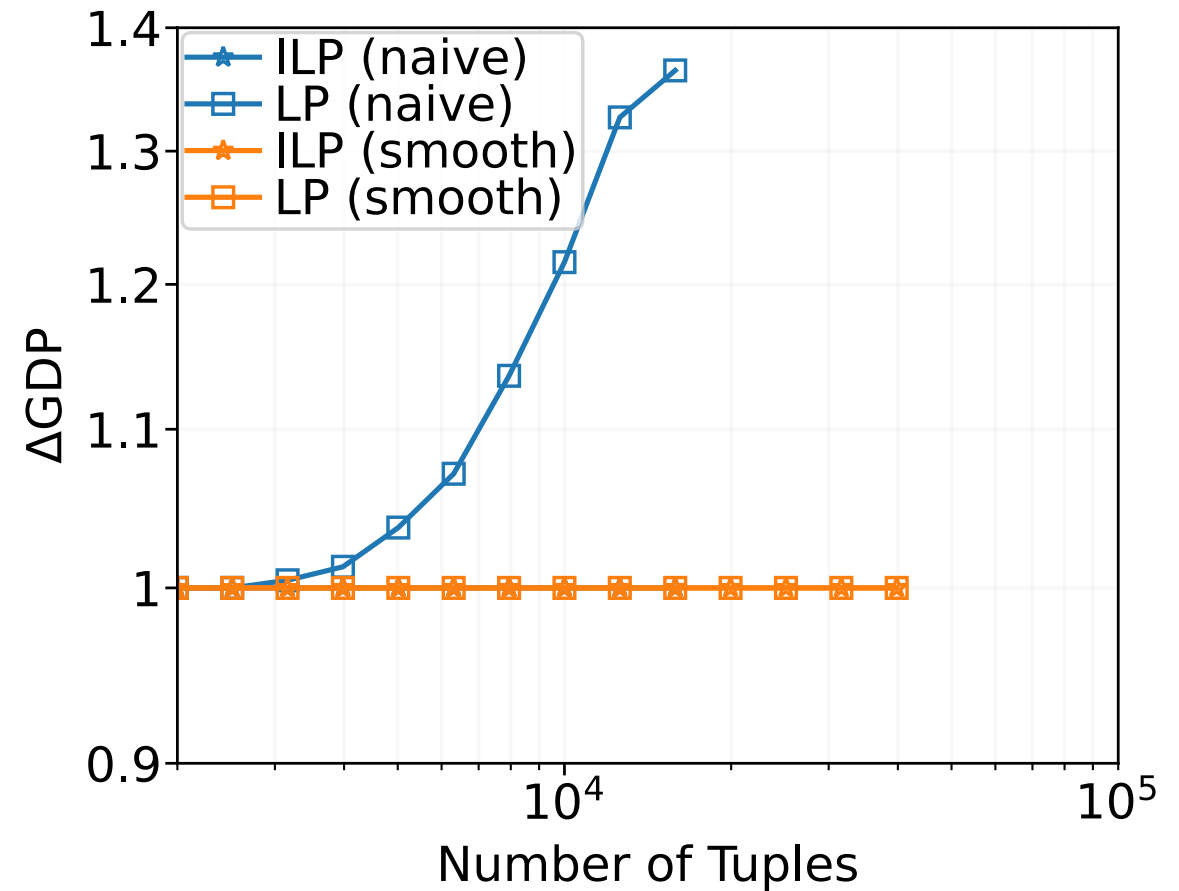
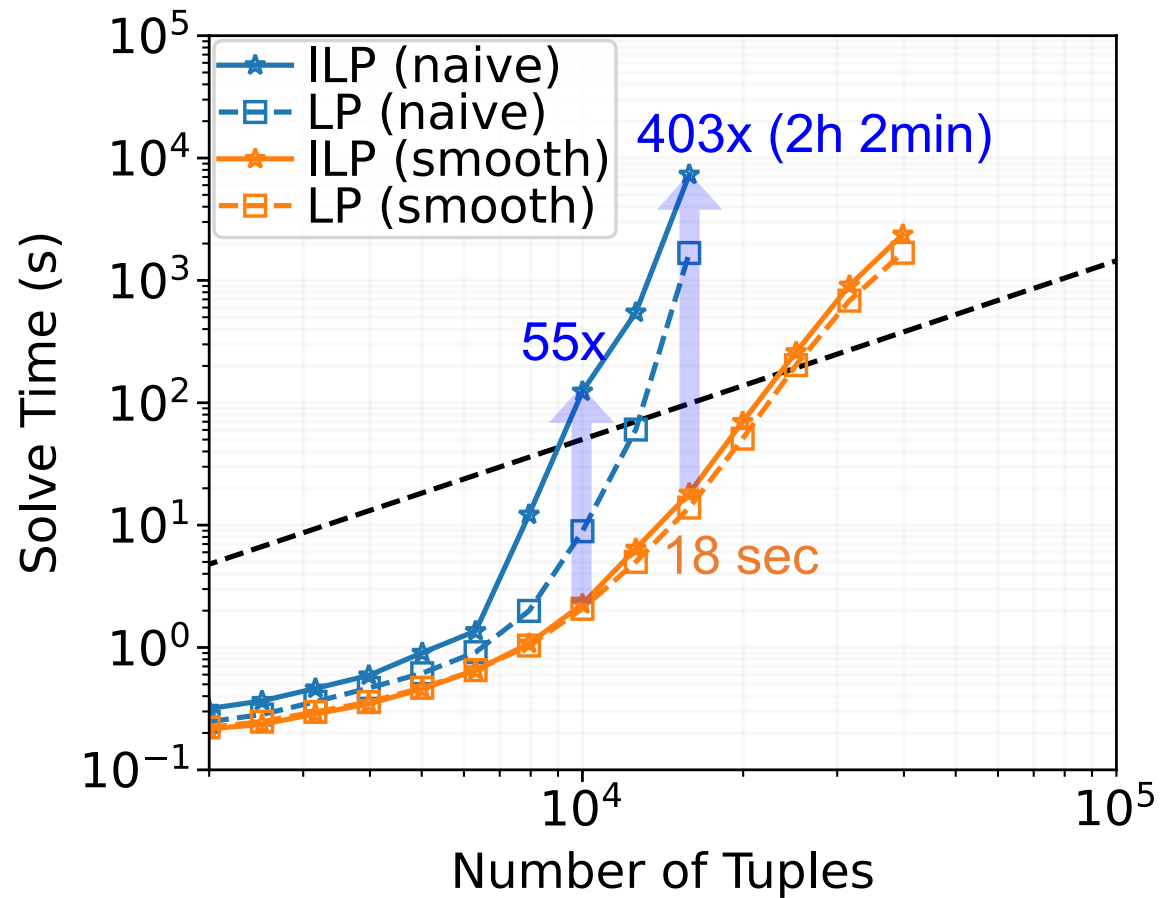
$$Q_{5\text{star}}(x): -R(x, a), S(x, b), T(x, c), U(x, d), V(x, e), A(x)$$



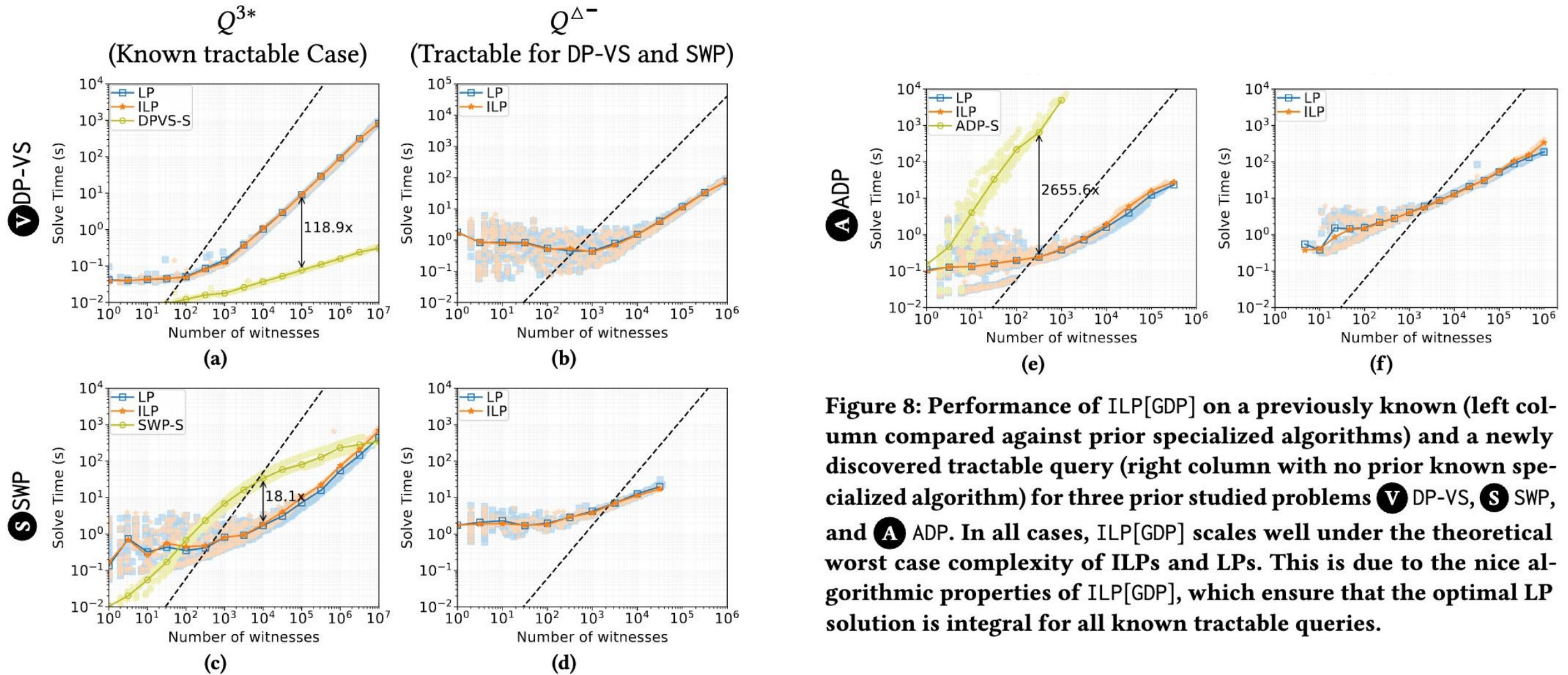
# Scalability of Naive vs. Smoothed ILP for PTIME query

Finding an optimal solution for the Smallest Witness Problem

$$Q_{5\text{star}}(x): -R(x, a), S(x, b), T(x, c), U(x, d), V(x, e), A(x)$$



# Comparison with specialized algorithms for PTIME cases



**Figure 8: Performance of ILP[GDP] on a previously known (left column compared against prior specialized algorithms) and a newly discovered tractable query (right column with no prior known specialized algorithm) for three prior studied problems  $\mathbf{V}$  DP-VS,  $\mathbf{S}$  SWP, and  $\mathbf{A}$  ADP. In all cases, ILP[GDP] scales well under the theoretical worst case complexity of ILPs and LPs. This is due to the nice algorithmic properties of ILP[GDP], which ensure that the optimal LP solution is integral for all known tractable queries.**

# A possible shift of focus of algorithm design in database theory?

## Current focus: discrete algorithms

Identify tractable cases for a class of problems that can be solved with some dedicated discrete algorithm (like dynamic programming or greedy) or a reduction to flow

For the hard cases:

- prove hardness via some dedicated reduction from some NPC problem.
- optionally design a separate dedicated approximation algorithm

Partial solutions: Often, the algorithm (the dichotomy) does not extend to all types of queries like CQs with self-joins, or problems under bag semantics

## Future: polyhedral algorithms

Design one "appropriate" ILP program to solve all problems

- "appropriate" here means that their natural LP relaxation has the same optimal objective for all PTIME cases ("LP=ILP"), which proves the ILP can be solved in PTIME.

**All cases** are covered (including the hard ones) ✓

- also approximation algorithms, just stop evaluation early, **anytime** algorithm comes for free ✓

Complete solutions: All problem types are covered ✓ (including self-joins or bag semantics)

# A possible shift of focus of algorithm design in database theory?

## Current focus: discrete algorithms

Identify tractable cases for a class of problems that can be solved with some dedicated discrete algorithm (like dynamic programming or greedy) or a reduction to flow

For the hard cases:

- prove hardness via some dedicated reduction from some NPC problem.
- optionally design a separate dedicated approximation algorithm

Partial solutions: Often, the algorithm (the dichotomy) does not extend to all types of queries like CQs with self-joins, or problems under bag semantics

## Future: polyhedral algorithms

Design one "appropriate" ILP program to solve all problems

- "appropriate" here means that their natural LP relaxation has the same optimal objective for all PTIME cases ("LP=ILP"), which proves the ILP can be solved in PTIME.

**All cases** are covered (including the hard ones) ✓

- also approximation algorithms, just stop evaluation early, **anytime** algorithm comes for free ✓

Complete solutions: All problem types are covered ✓ (including self-joins or bag semantics)

*An anonymous concern: "...the ILP constructed is not a simple mathematical object... Since the construction given is .... not a simple mathematical object, it is not clear to me how deep one can push this further by analyzing it."*

# A possible shift of focus of algorithm design in database theory?

## Current focus: discrete algorithms

Identify tractable cases for a class of problems that can be solved with some dedicated discrete algorithm (like dynamic programming or greedy) or a reduction to flow

For the hard cases:

- prove hardness via some dedicated reduction from some NPC problem.
- optionally design a separate dedicated approximation algorithm

Partial solutions: Often, the algorithm (the dichotomy) does not extend to all types of queries like CQs with self-joins, or problems under bag semantics

**Practical aspects**: usually only some problem cases are solved, hard cases often not treated, the practical nature of the algorithms is not always clear

## Future: polyhedral algorithms

Design one "appropriate" ILP program to solve all problems

- "appropriate" here means that their natural LP relaxation has the same optimal objective for all PTIME cases ("LP=ILP"), which proves the ILP can be solved in PTIME.

**All cases** are covered (including the hard ones) ✓

- also approximation algorithms, just stop evaluation early, **anytime** algorithm comes for free ✓

**Complete solutions**: All problem types are covered ✓ (including self-joins or bag semantics)

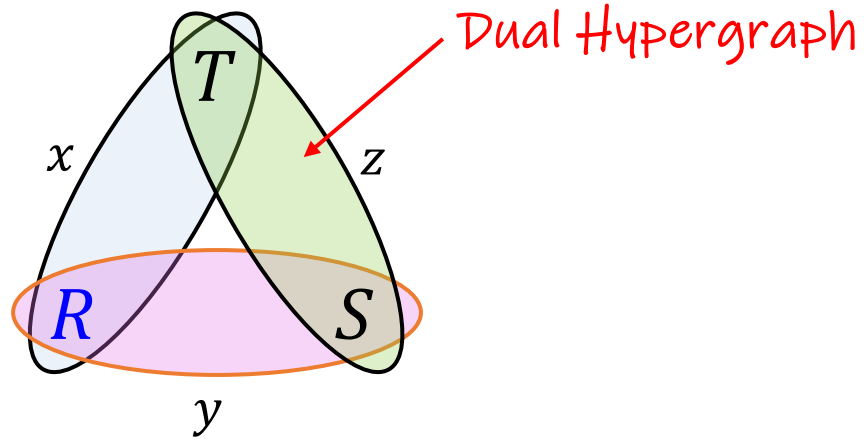
**Practical aspects**: it works from day one ✓

Theory

# Example complicated landscape for resilience

Triangle query

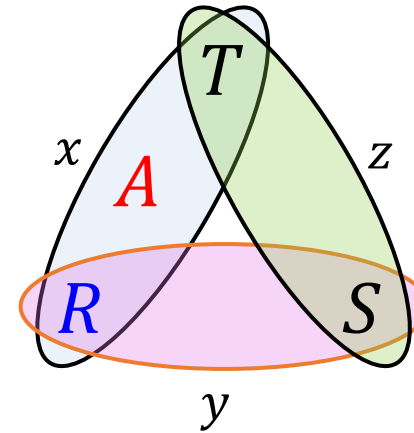
$$Q^\Delta: \neg R(x, y), S(y, z), T(x, z)$$



**NPC**

Triangle unary

$$Q_A^\Delta: \neg R(x, y), S(y, z), T(x, z), A(x)$$

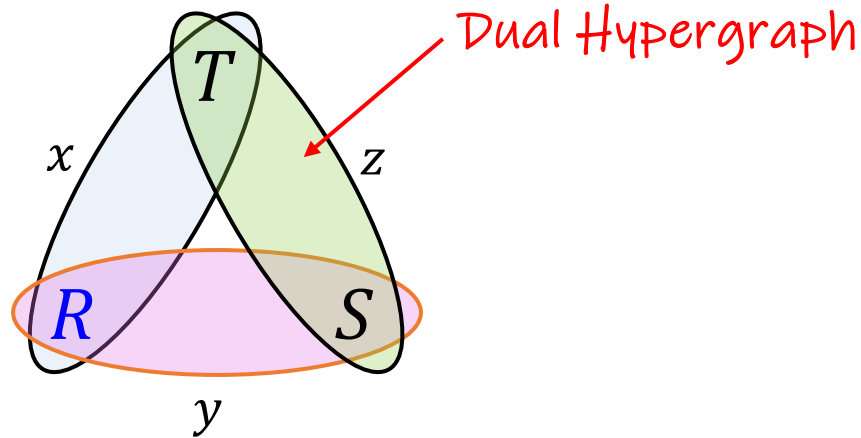


**PTIME**

# Example complicated landscape for resilience

Triangle query

$$Q^\Delta: \neg R(x, y), S(y, z), T(x, z)$$



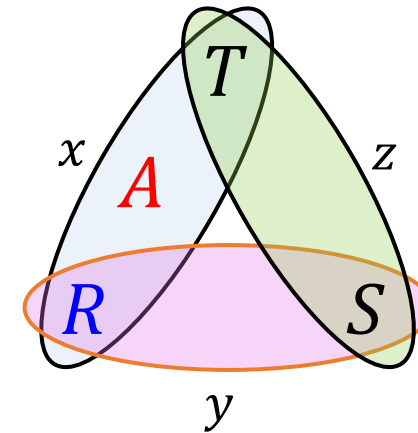
**NPC**

**PTIME** for FD  $x \rightarrow y$

**PTIME** if provenance happens to be read-once

Triangle unary

$$Q_A^\Delta: \neg R(x, y), S(y, z), T(x, z), A(x)$$



**PTIME**

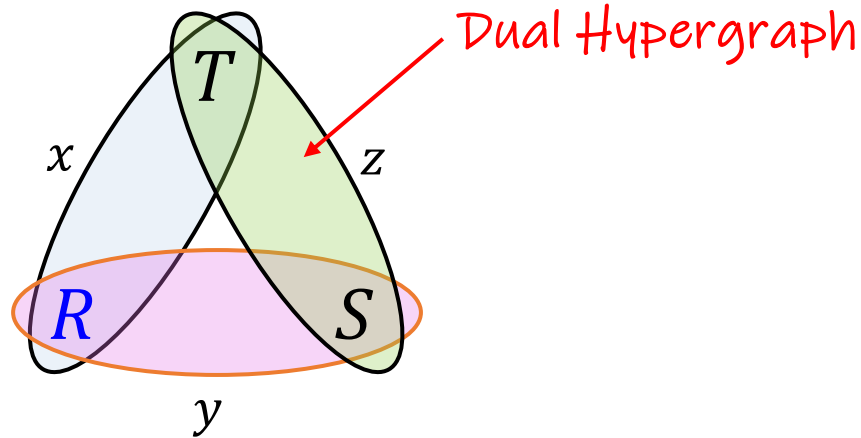
**NPC** under bag semantics



# Example complicated landscape for resilience

Triangle query

$$Q^\Delta: \neg R(x, y), S(y, z), T(x, z)$$



**NPC**

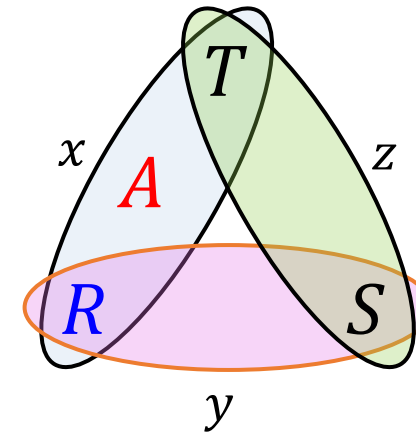
**PTIME** for FD  $x \rightarrow y$

**PTIME** if provenance happens to be read-once

*"Coarse-grained instance-optimal" algorithm*

Triangle unary

$$Q_A^\Delta: \neg R(x, y), S(y, z), T(x, z), A(x)$$



**PTIME**

**NPC** under bag semantics

Freire, Gatterbauer, Immerman, Meliou. The complexity of resilience and responsibility for self-join-free conjunctive queries, VLDB 2015. <https://doi.org/10.14778/2850583.2850592>

Makhija, Gatterbauer. A Unified Approach for Resilience and Causal Responsibility with Integer Linear Programming (ILP) and LP Relaxations, SIGMOD 2024. <https://doi.org/10.1145/3626715>

Makhija, Gatterbauer: A Unified and Practical Approach for Generalized Deletion Propagation. <https://arxiv.org/pdf/2411.17603>

# Take-aways Generalized Deletion Propagation

- Polyhedral theory solves many database theory problems "out of the box".
  - Shift in focus: instead of trying to find a dedicated PTIME algorithm for PTIME cases, start with a general formulation and prove it finishes in PTIME for PTIME cases.
  - There is some magic in getting the "right" formulation (ILP = LP for PTIME), we don't yet have "the" recipe
  - The proofs for LP=ILP go beyond standard optimization literature. Polyhedral theory alone does not help.
- The overall philosophy is way more general than reverse data management.
  - Makhija, Gatterbauer. *Minimally Factorizing the Provenance of Self-Join Free Conjunctive Queries*, PODS 2024.  
<https://doi.org/10.1145/3651605>
  - What about consistent query answering? And even more general logic optimization problems?
- More concretely open: Unifying deletion and insertion propagation ("change propagation"), basically positive and negative provenance / Why or Why not?
  - Meliou, Gatterbauer, Moore, Suciu. *Why so? or Why no? Functional causality for explaining query answers*. MUD 2010.  
<https://arxiv.org/pdf/0912.5340>
- Please talk to Neha 😊



Thank you 😊