

Is Integer Linear Programming All You Need for Deletion Propagation?

A unified and practical approach to
Generalized Deletion Propagation

VLDB 2025 (London, Sept 4, 2025)



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

Neha Makhija

Northeastern University
(UMass Amherst from Sept 2025)



Wolfgang Gatterbauer

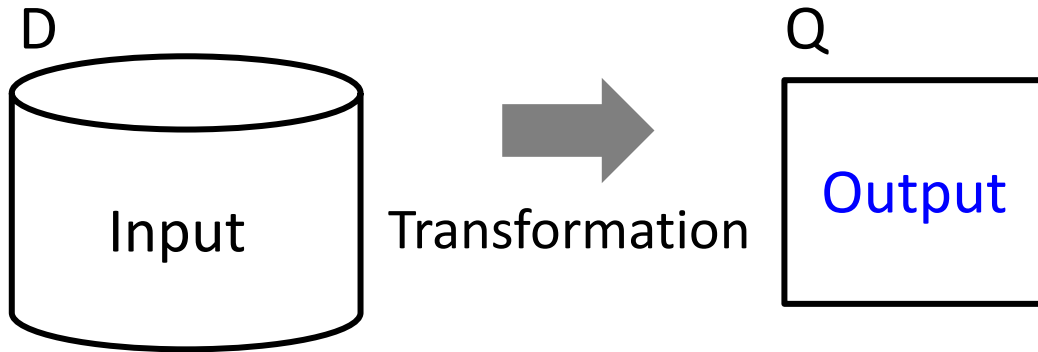
Northeastern University



Data Management

QUERY EVALUATION:

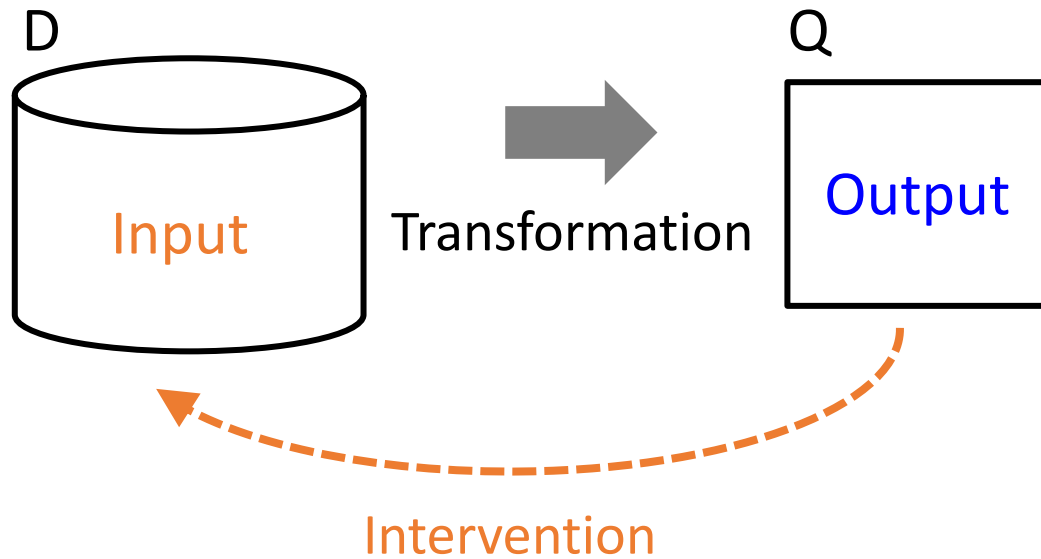
A transformation of the input to the output



Reverse Data Management (RDM)

QUERY EVALUATION:

A transformation of the input to the output



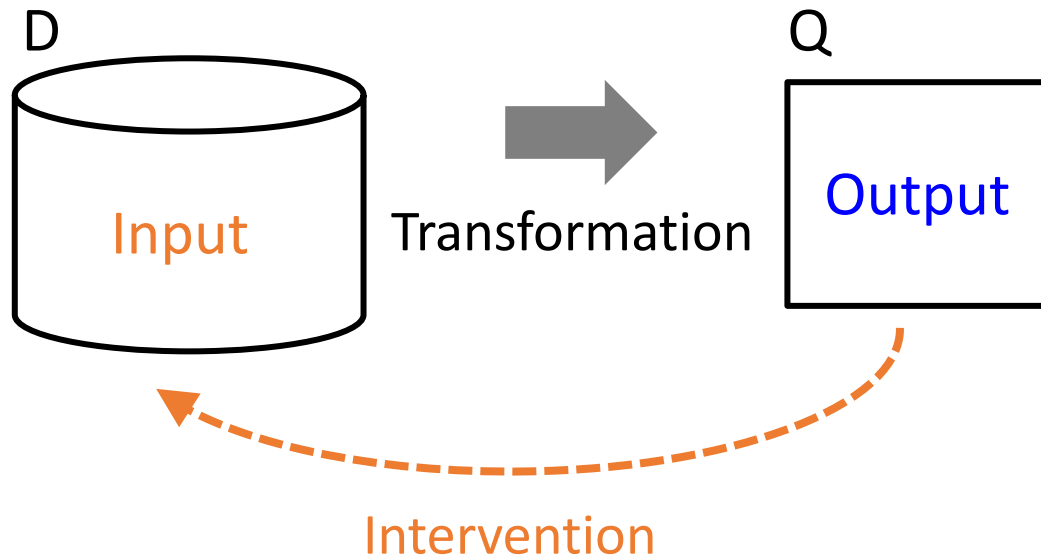
REVERSE DATA MANAGEMENT (RDM):

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

Reverse Data Management (RDM): 2 types of Explanations

QUERY EVALUATION:

A transformation of the input to the output



REVERSE DATA MANAGEMENT (RDM):

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

Two types of **explanations** in RDM, naming used from Explainable AI (XAI):

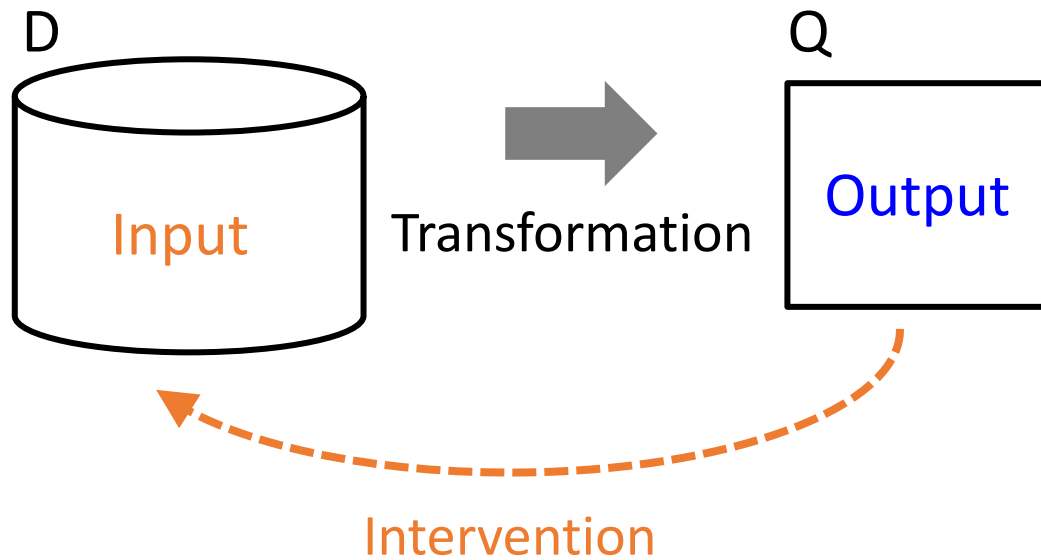
1. Contrastive (counterfactual):

2. Abductive (factual):

Reverse Data Management (RDM): 2 types of Explanations

QUERY EVALUATION:

A transformation of the input to the output



REVERSE DATA MANAGEMENT (RDM):

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

Two types of **explanations** in RDM, naming used from Explainable AI (XAI):

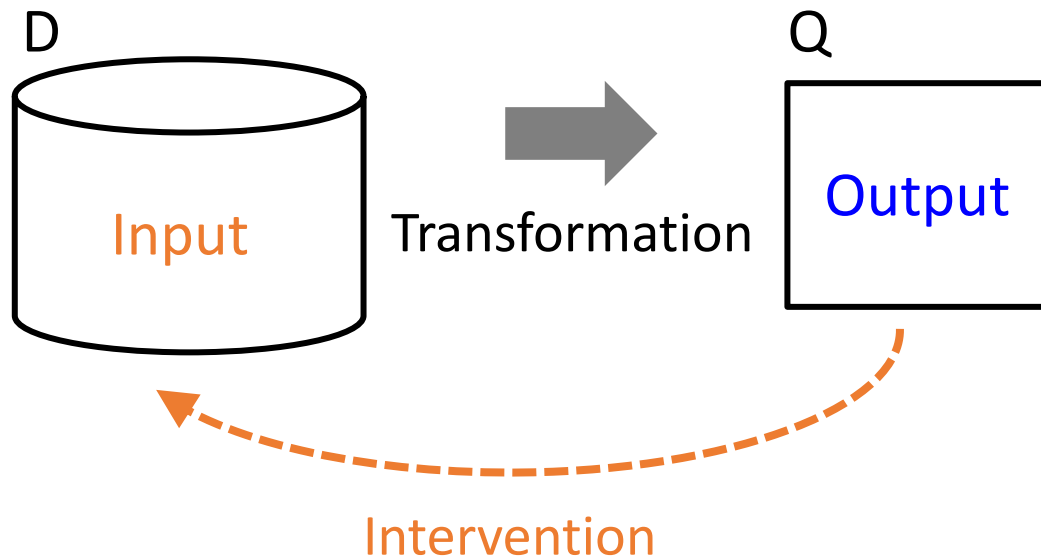
1. Contrastive (counterfactual): set of **tuples** (features) of **min size** that are sufficient to change an output (prediction)

2. Abductive (factual):

Reverse Data Management (RDM): 2 types of Explanations

QUERY EVALUATION:

A transformation of the input to the output



REVERSE DATA MANAGEMENT (RDM):

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

Two types of **explanations** in RDM, naming used from Explainable AI (XAI):

1. Contrastive (counterfactual): set of **tuples** (features) of **min size** that are sufficient to change an output (prediction)
= **min deletions of tuples**
(e.g., resilience, deletion propagation)

2. Abductive (factual): set of **tuples** (features) of **min size** that are sufficient for ensuring a certain output (prediction)
= **max deletions of tuples**
(e.g., smallest witness problem)

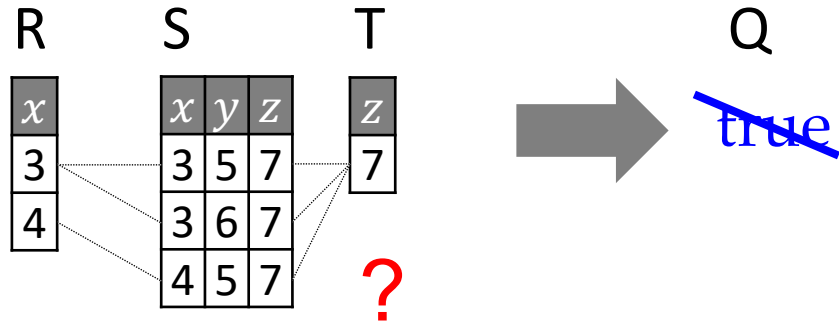
= **DELETION PROPAGATION**

Example Reverse Data Management Problems

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

1. Resilience

Delete min number of tuples to make Boolean Q false



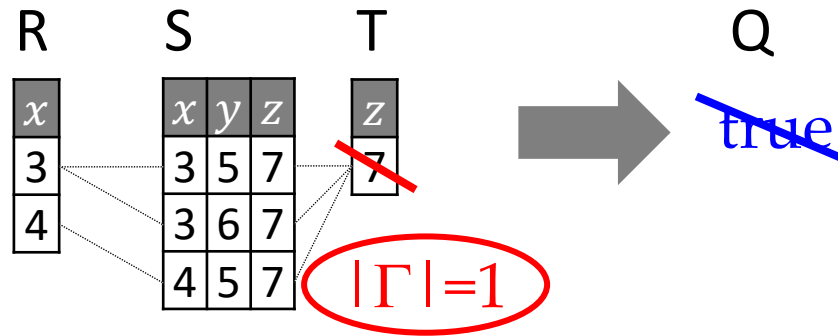
```
select exists(  
  select 1  
  from R, S, T  
  where R.x=S.x  
  and S.z=T.z)
```

Example Reverse Data Management Problems

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

1. Resilience

Delete min number of tuples to make Boolean Q false



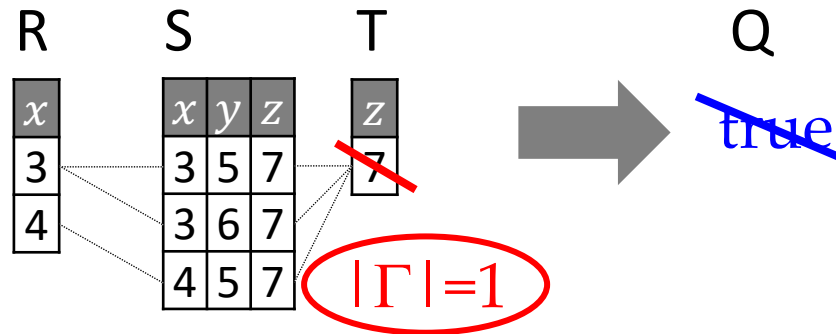
```
select exists(  
  select 1  
  from R, S, T  
  where R.x=S.x  
  and S.z=T.z)
```


Example Reverse Data Management Problems

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

1. Resilience

Delete min number of tuples to make Boolean Q false

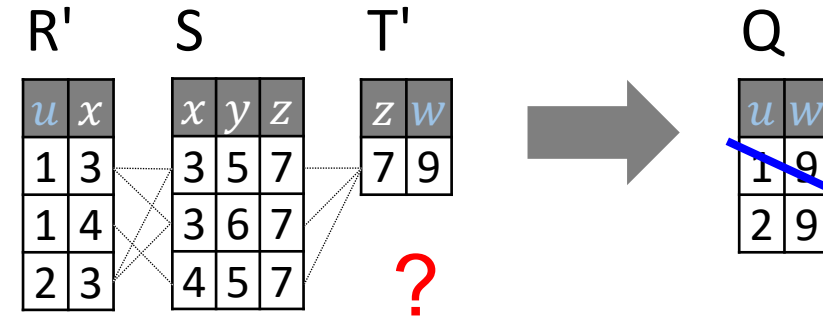


```
select exists(
  select 1
  from R, S, T
  where R.x=S.x
  and S.z=T.z)
```

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

2. Source side-effects (deletion propagation)

Delete min number of tuples to delete an output tuples



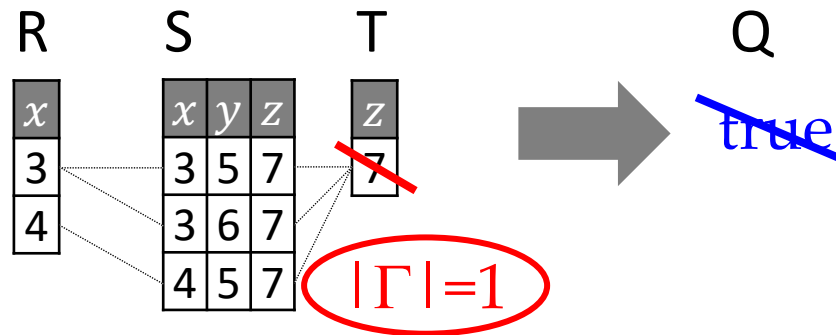
```
select R'.u, T'.w
from R', S, T'
where R'.x=S.x
and S.z=T'.z
```

Example Reverse Data Management Problems

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

1. Resilience

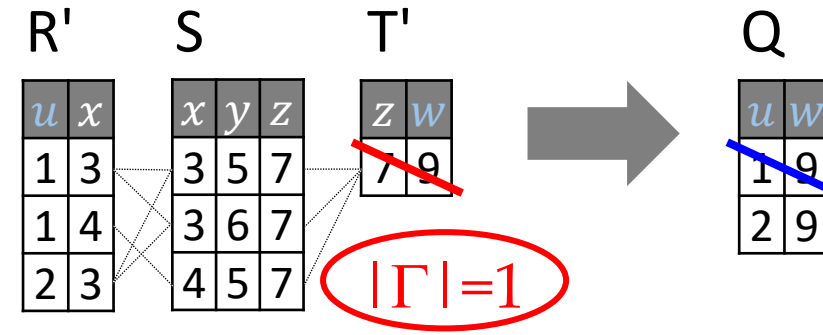
Delete min number of tuples to make Boolean Q false



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

2. Source side-effects (deletion propagation)

Delete min number of tuples to delete an output tuples



basically identical

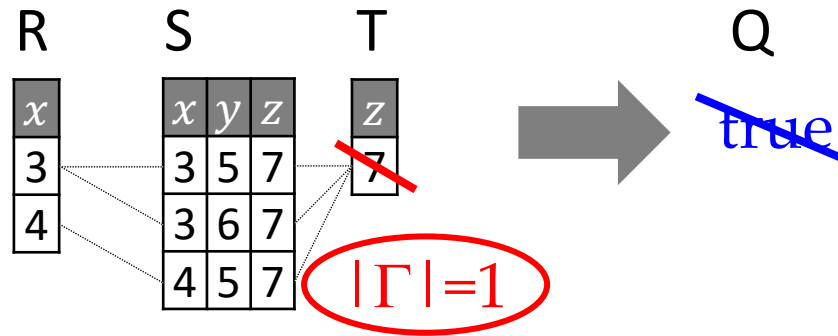
A long open problem: for which conjunctive queries is resilience in PTIME? (and for which NP-complete?). Only a partial classification known today

Example Reverse Data Management Problems

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

1. Resilience

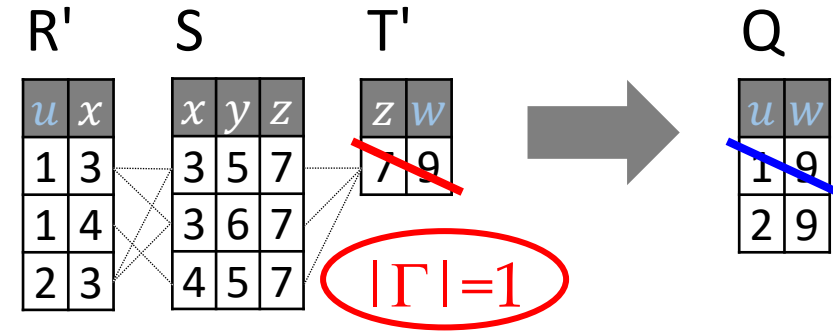
Delete min number of tuples to make Boolean Q false



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

2/3. (Aggregated) Source side-effects (d.p.)

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



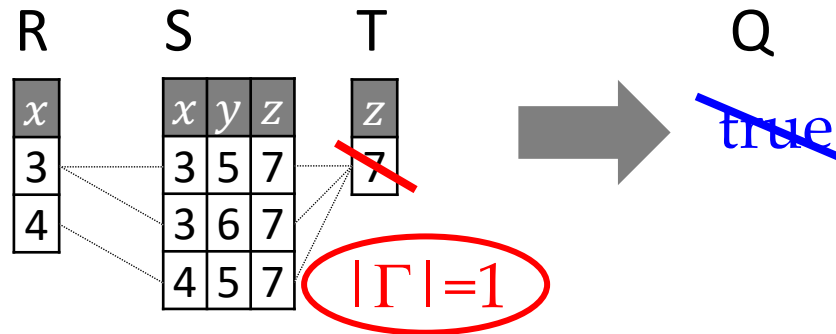
Different tractability results!

Example Reverse Data Management Problems

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

1. Resilience

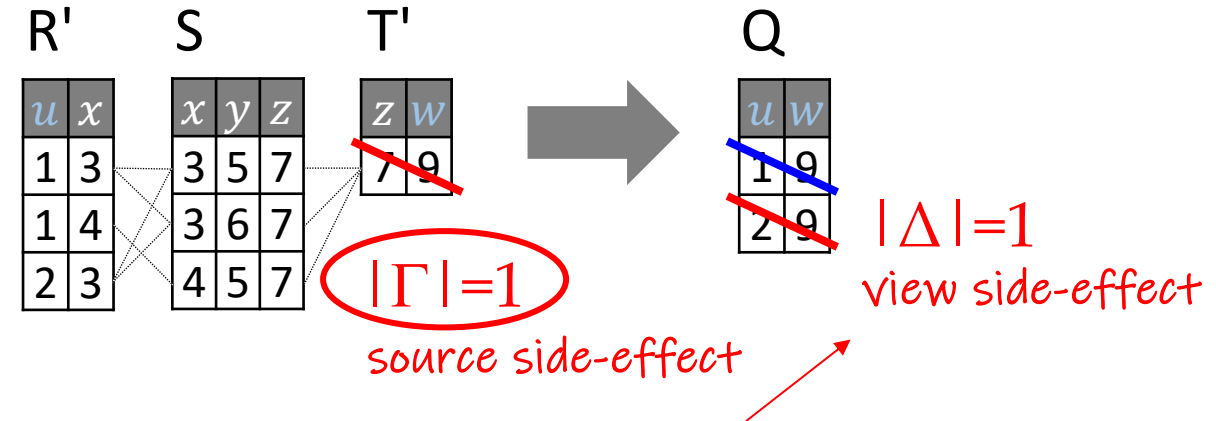
Delete min number of tuples to make Boolean Q false



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

2/3. (Aggregated) Source side-effects (d.p.)

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



Oops, we deleted more output tuples than needed.

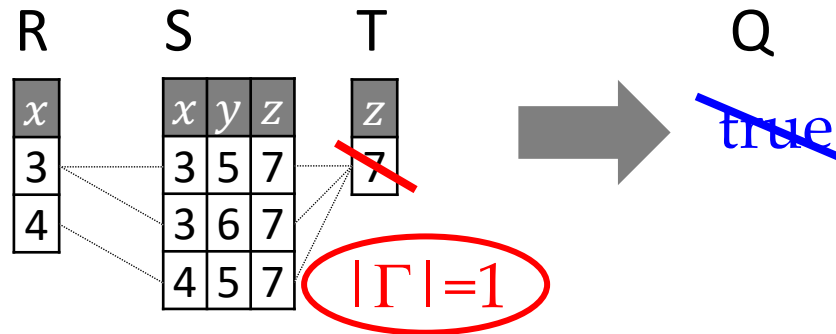
That is not a concern when minimizing source-side effects.
But when minimizing view-side effects!

Example Reverse Data Management Problems

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

1. Resilience

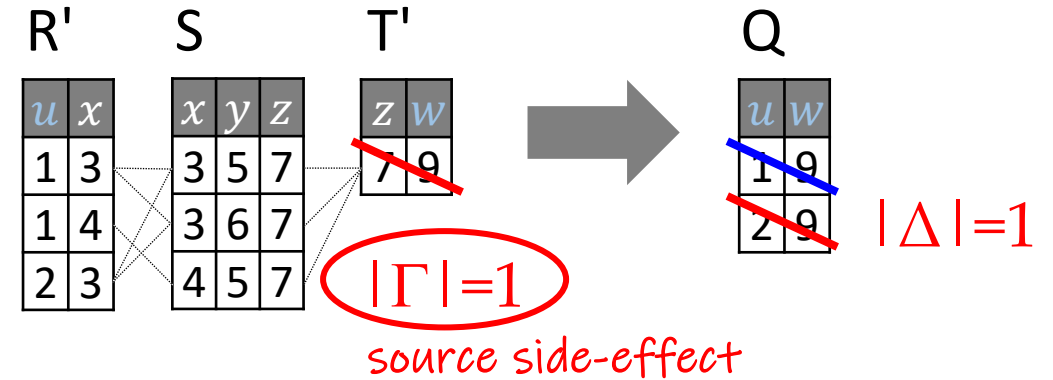
Delete min number of tuples to make Boolean Q false



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

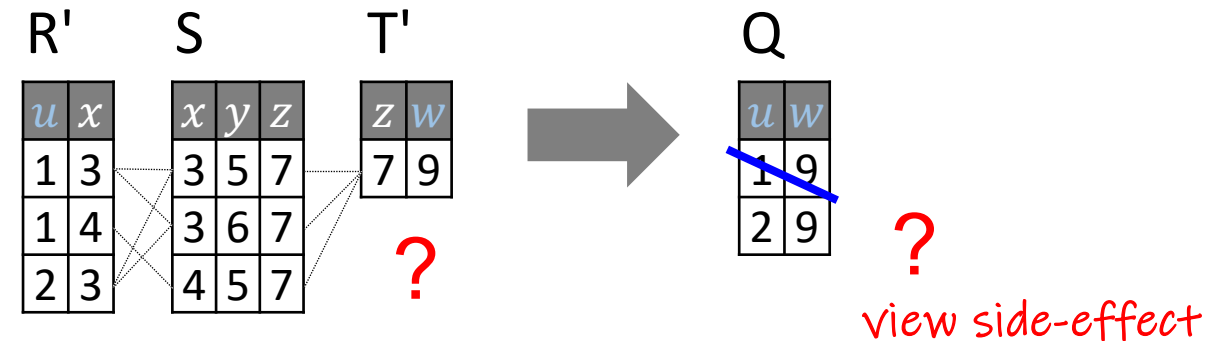
2/3. (Aggregated) Source side-effects (d.p.)

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



4. View side-effects (deletion propagation)

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

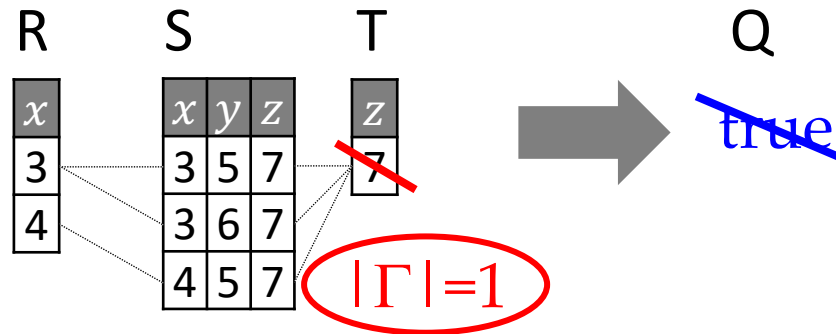


Example Reverse Data Management Problems

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

1. Resilience

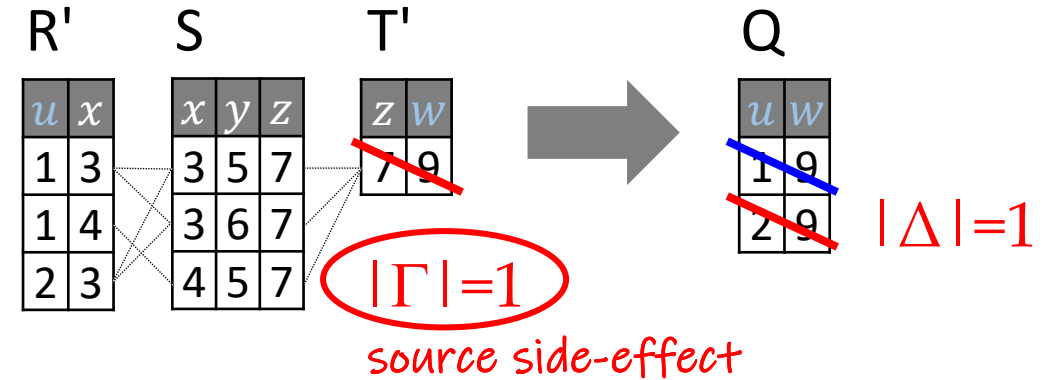
Delete min number of tuples to make Boolean Q false



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

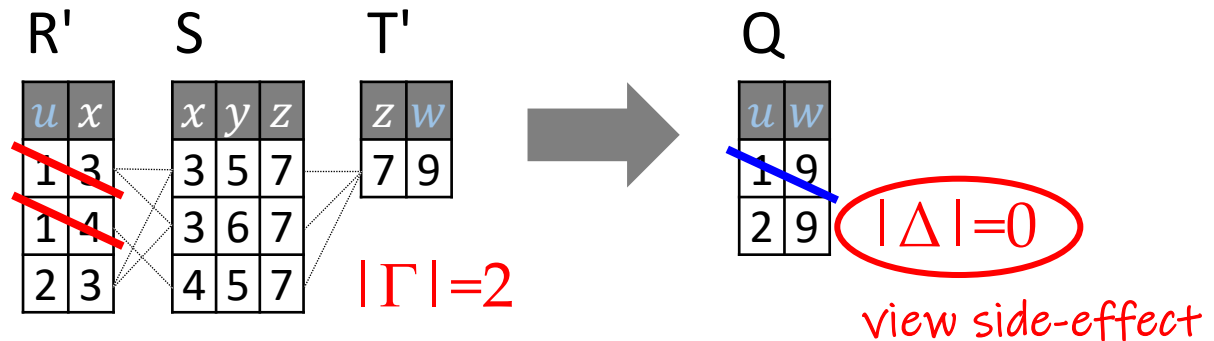
2/3. (Aggregated) Source side-effects (d.p.)

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



4. View side-effects (deletion propagation)

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



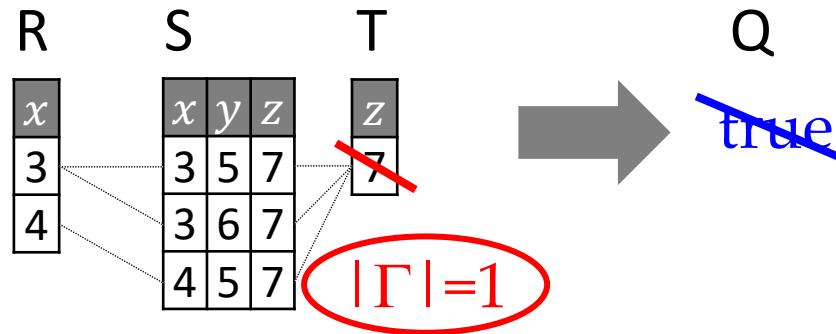
Again: Different tractability results!

Example Reverse Data Management Problems

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

1. Resilience

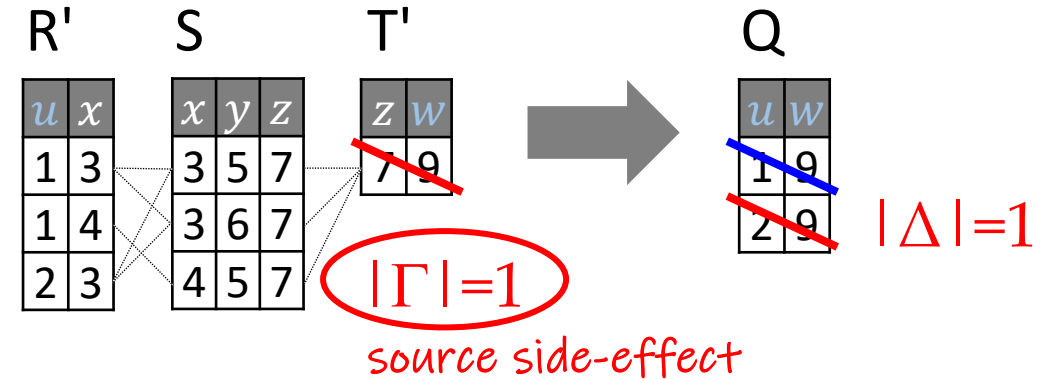
Delete min number of tuples to make Boolean Q false



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

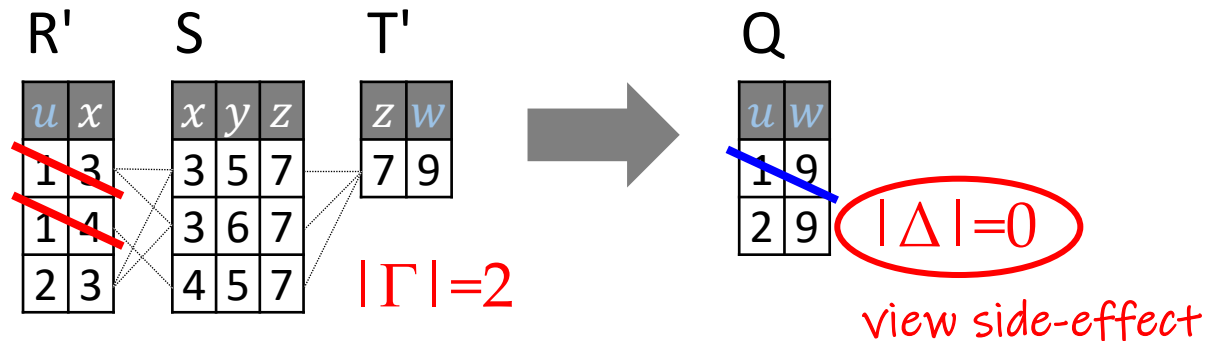
2/3. (Aggregated) Source side-effects (d.p.)

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



4/5. (Aggregated) View side-effects (d.p.)

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



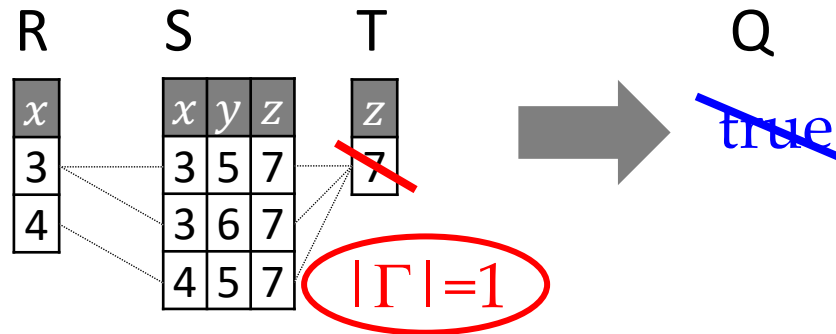
No prior work yet

Example Reverse Data Management Problems

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

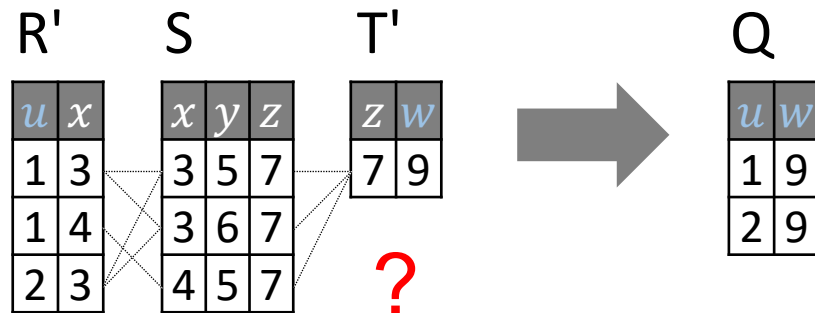
1. Resilience

Delete min number of tuples to make Boolean Q false



5. Smallest witness problem

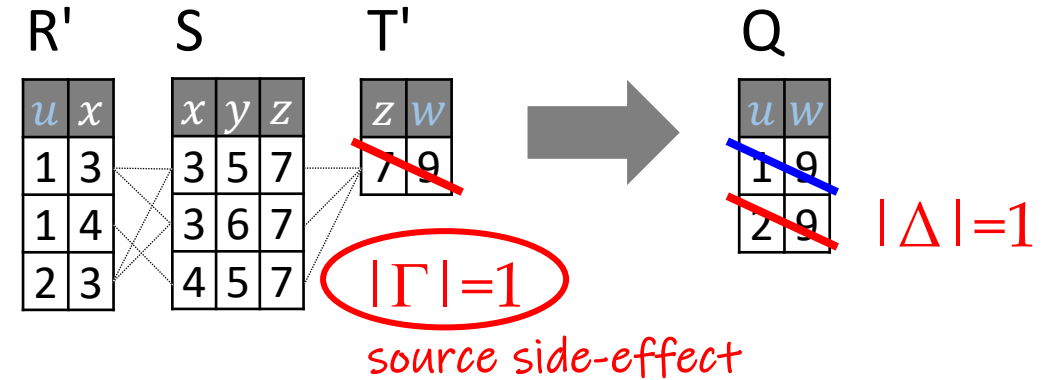
Delete max num. of tuples while keeping all output tuples



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

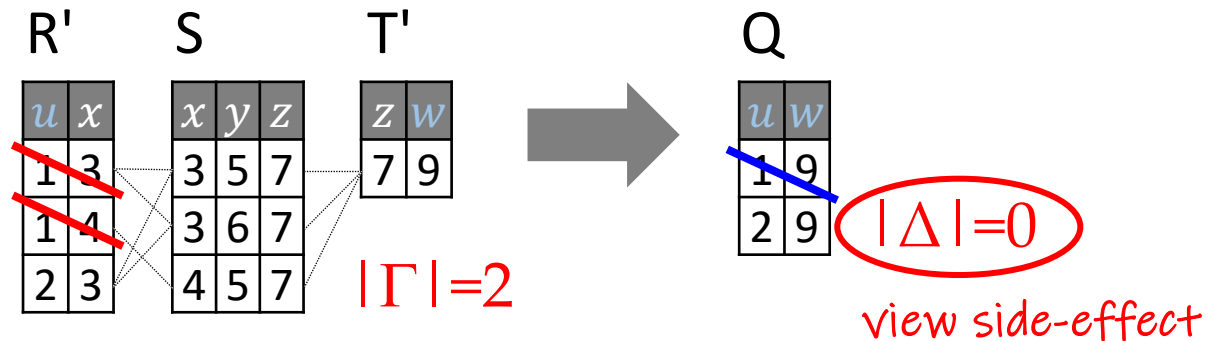
2/3. (Aggregated) Source side-effects (d.p.)

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



4/5. (Aggregated) View side-effects (d.p.)

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

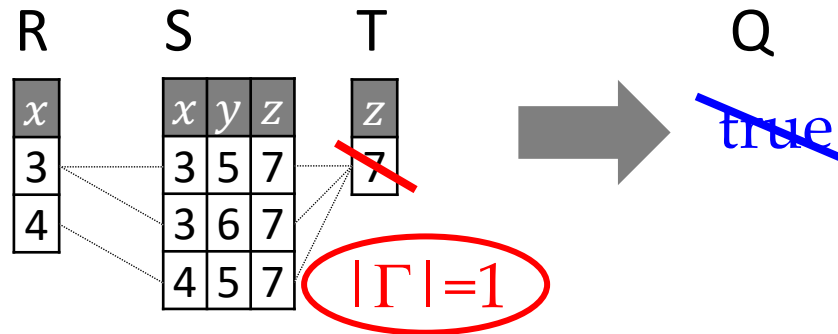


Example Reverse Data Management Problems

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

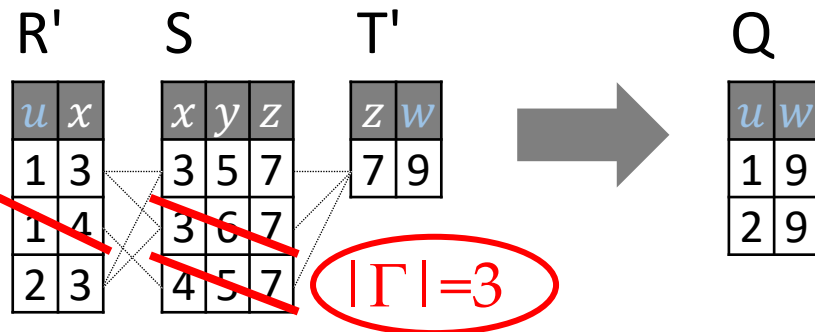
1. Resilience

Delete min number of tuples to make Boolean Q false



5. Smallest witness problem

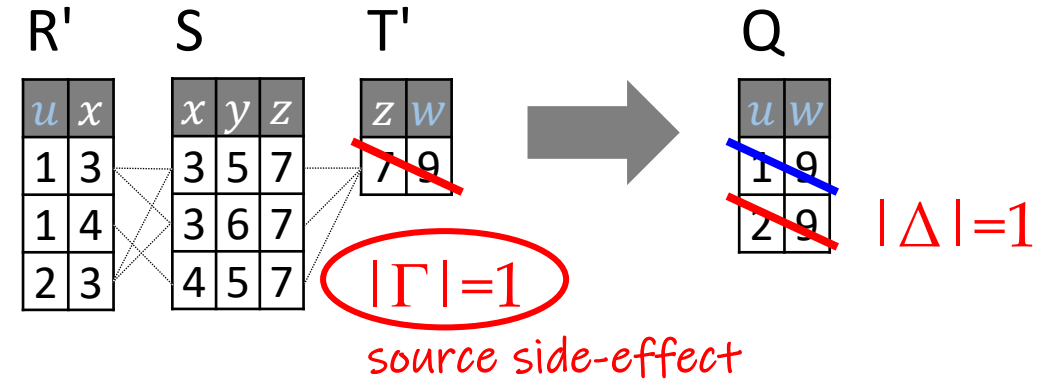
Delete max num. of tuples while keeping all output tuples



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

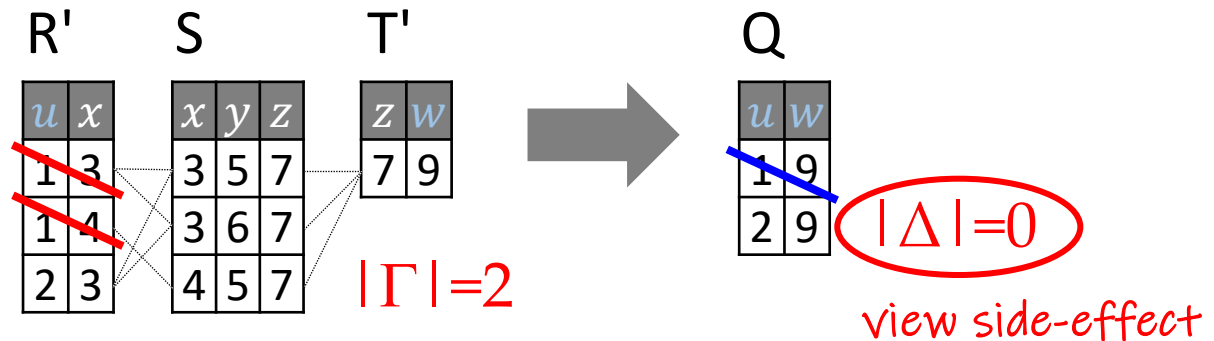
2/3. (Aggregated) Source side-effects (d.p.)

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



4/5. (Aggregated) View side-effects (d.p.)

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

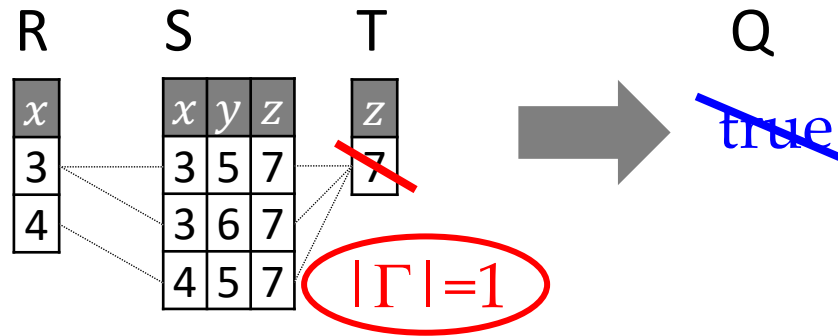


Example Reverse Data Management Problems

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

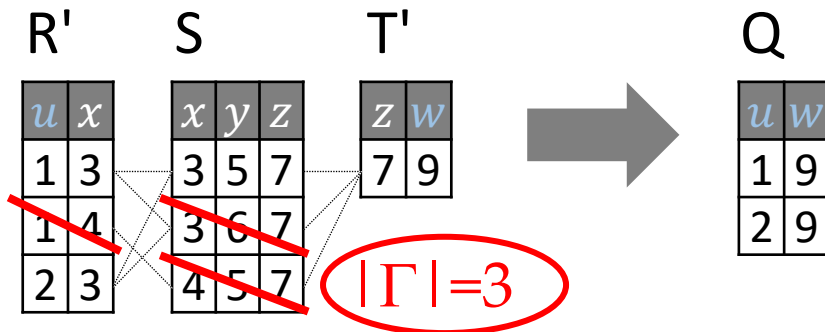
1. Resilience

Delete min number of tuples to make Boolean Q false



6./7 (Aggregated) Smallest witness problem

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

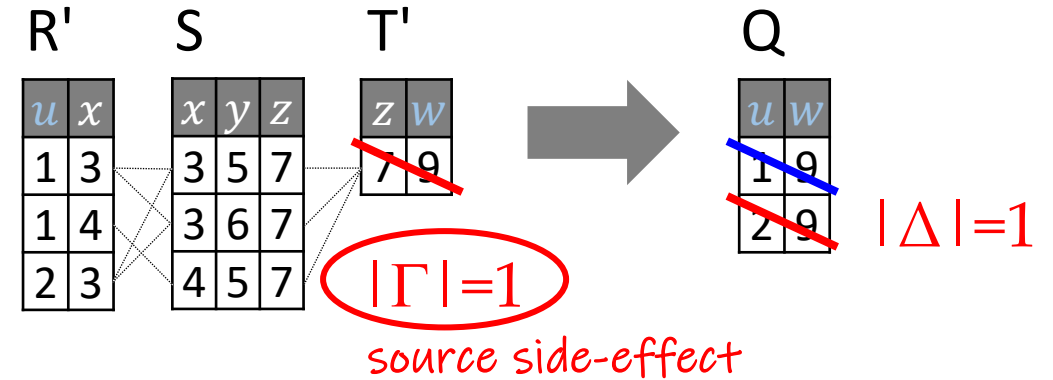


No prior work yet

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

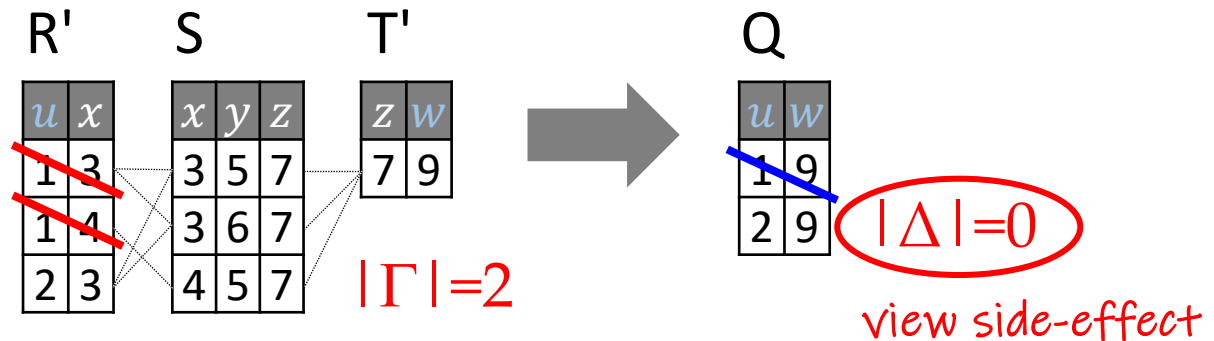
2/3. (Aggregated) Source side-effects (d.p.)

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



4/5. (Aggregated) View side-effects (d.p.)

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

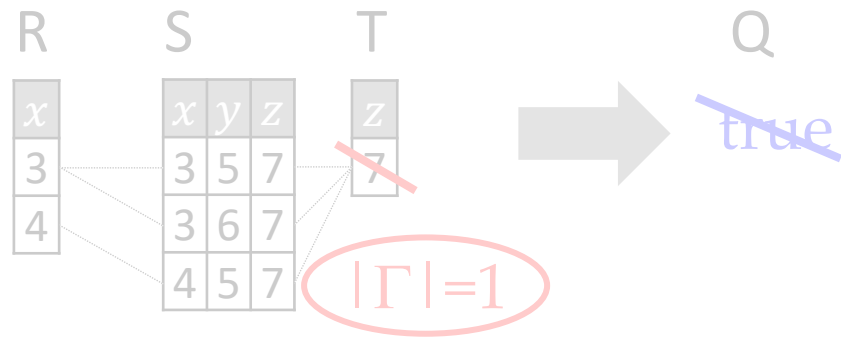


A Plethora of Reverse Data Management Problems ...

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

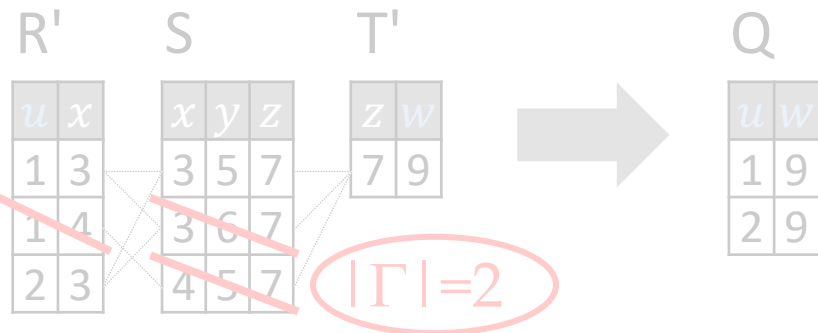
1. Resilience

Delete min number of tuples to make Boolean Q false



6./7 (Aggregated) Smallest witness problem

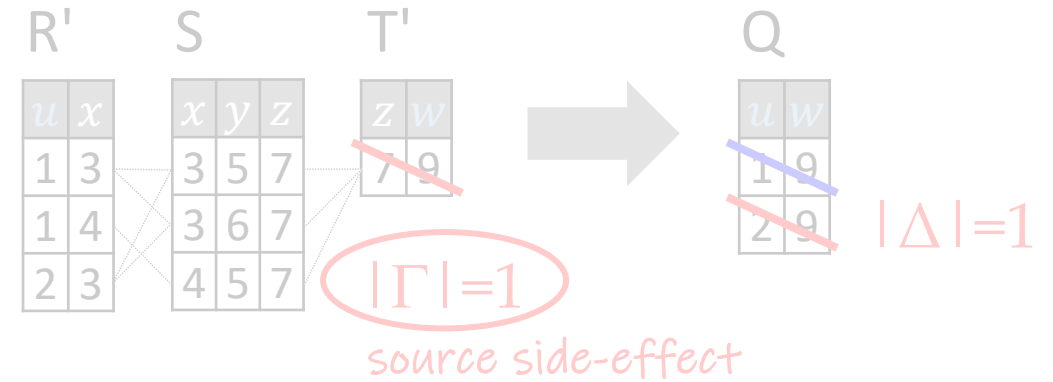
Delete max num. of tuples while keeping $\geq k$ output tuples



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

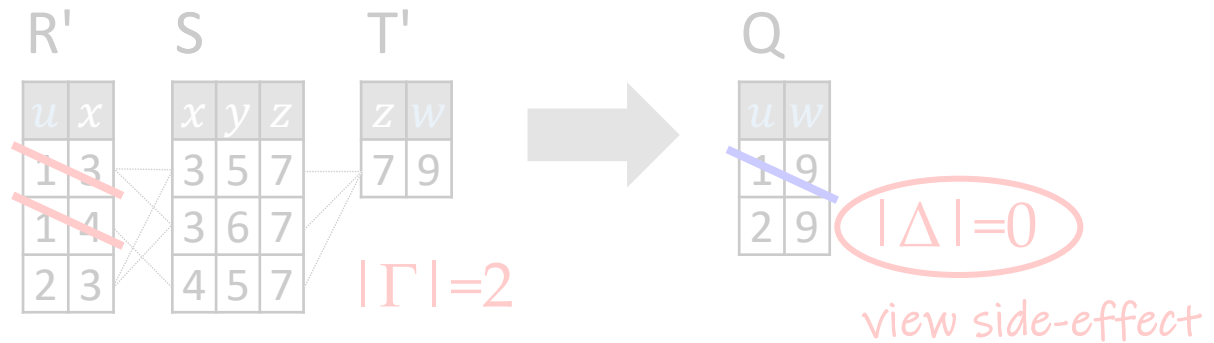
2/3. (Aggregated) Source side-effects (d.p.)

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



4/5. (Aggregated) View side-effects (d.p.)

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 **ILP encoding** to solve all these problems
 - including difficult cases, such as self-joins, or bag semantics
3. The ILP formulation **is solvable in PTIME for all known PTIME cases**
 - including all known PTIME cases for the problems of: **resilience** [VLDB'15], **aggregated deletion propagation** [VLDB'20], **view-side effects** [PODS'11], **smallest witness problem** [ICDT'24], under functional dependencies, and both set and bag semantics (where results are known)

Outline

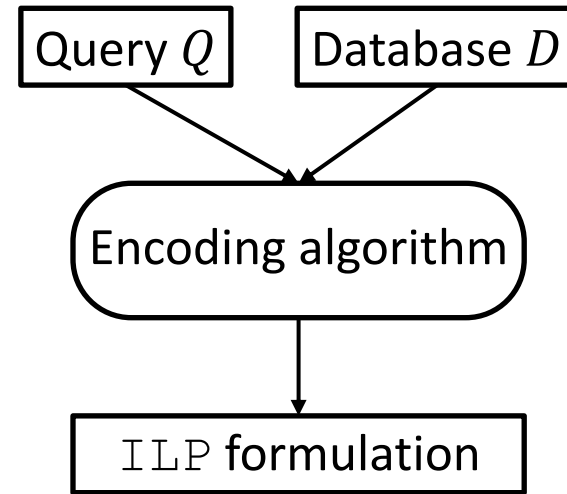
1. Reverse Data Management (RDM)
2. A magical ILP formulation
3. Take-aways

(An illustrated example: only if time remains)

Unified Algorithms for Reverse Data Management

ILP formulation:

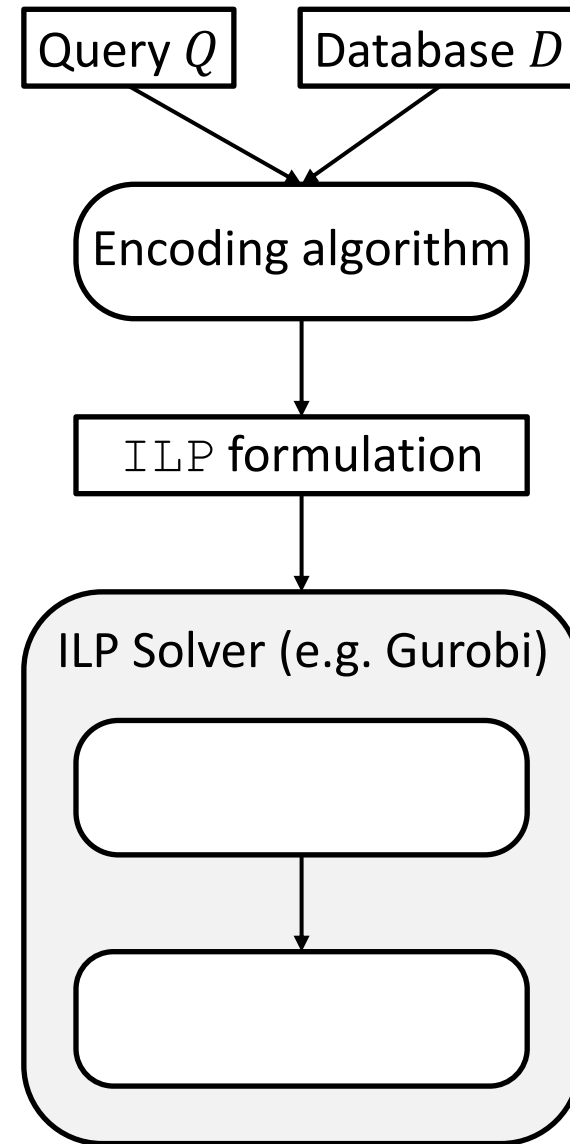
$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n \end{aligned}$$



Unified Algorithms for Reverse Data Management

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n \end{aligned}$$



Unified Algorithms for Reverse Data Management

ILP formulation:

$$f^* = \min[\mathbf{c} \cdot \mathbf{x}]$$

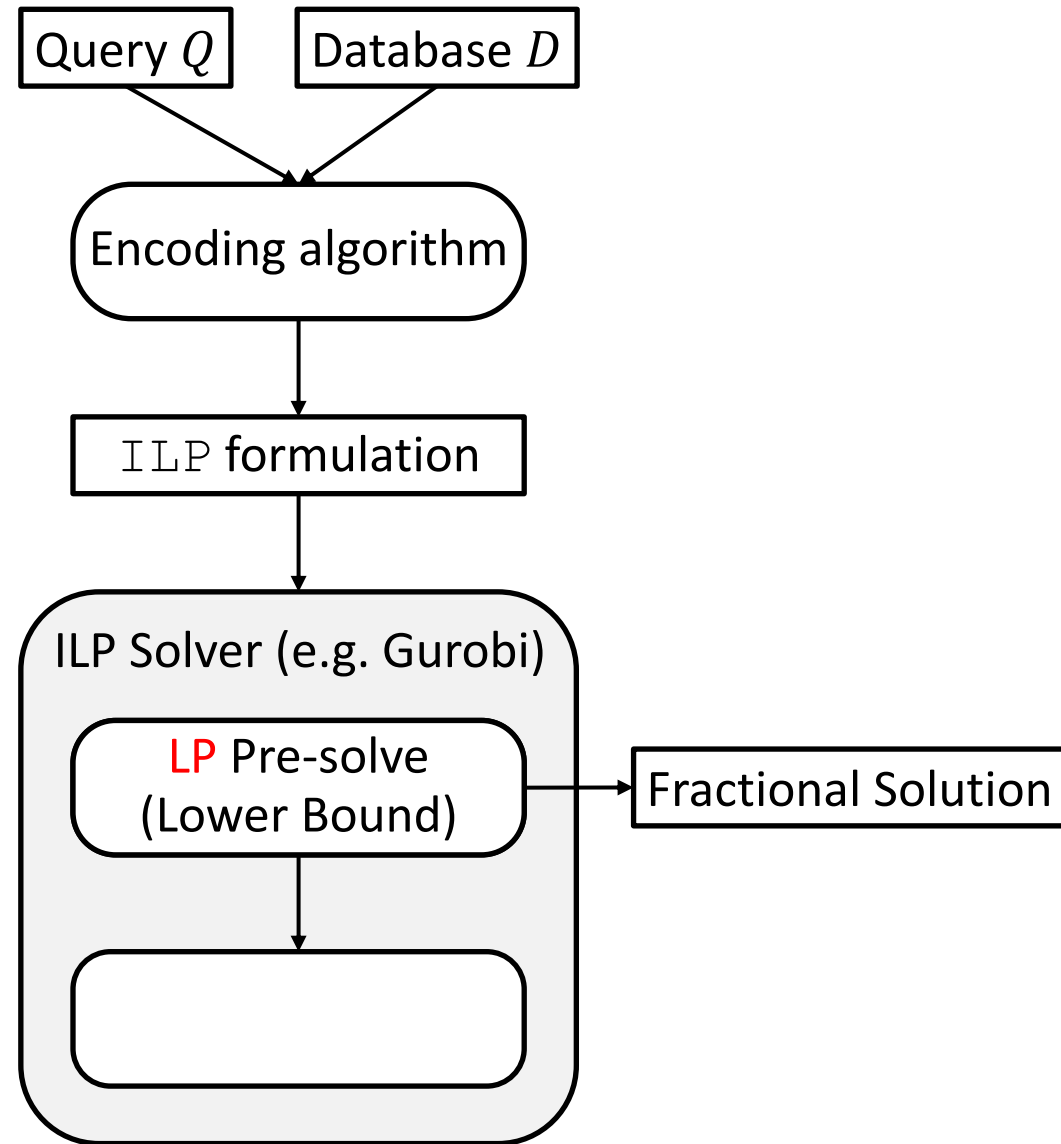
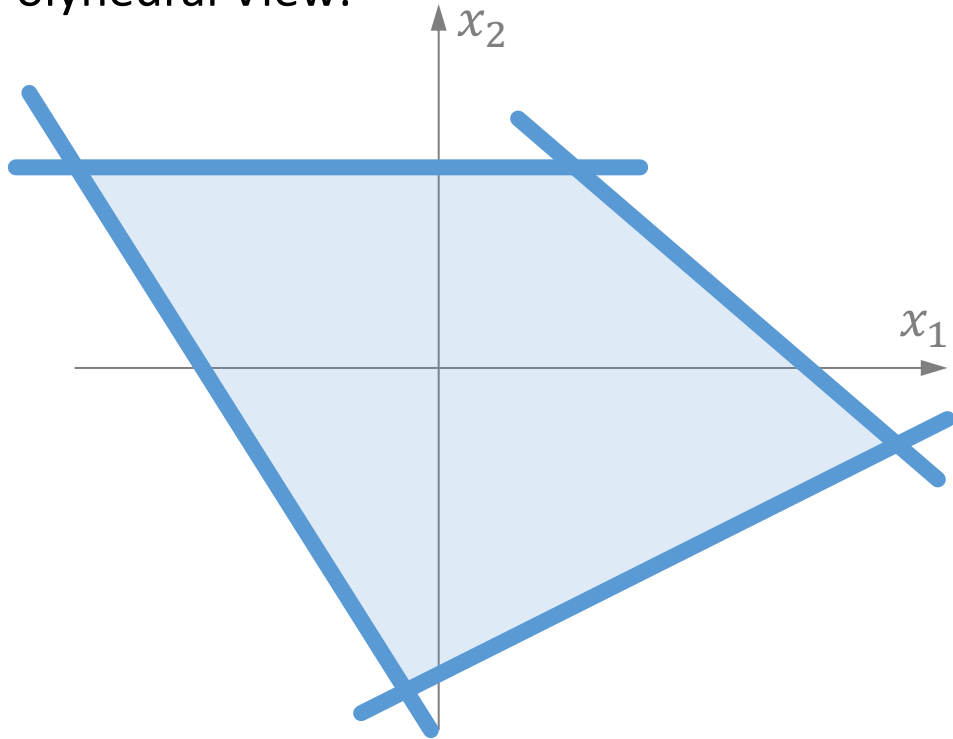
constraint vector

$$\text{s.t. } \mathbf{A} \cdot \mathbf{x} \geq \mathbf{b}$$

constraint matrix

$$\mathbf{x} \in \mathbb{N}^n / \mathbb{R}^n$$

Polyhedral View:



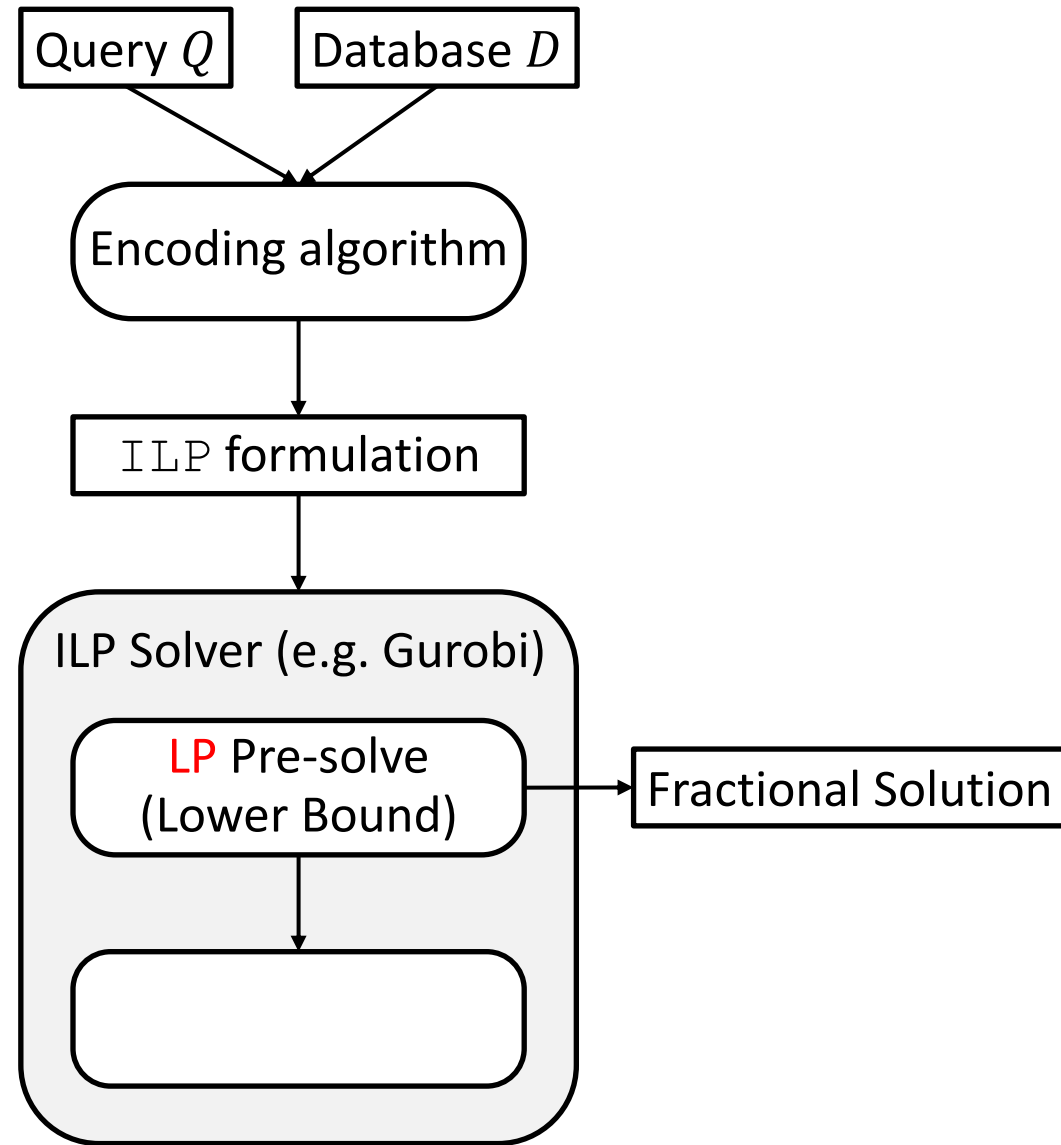
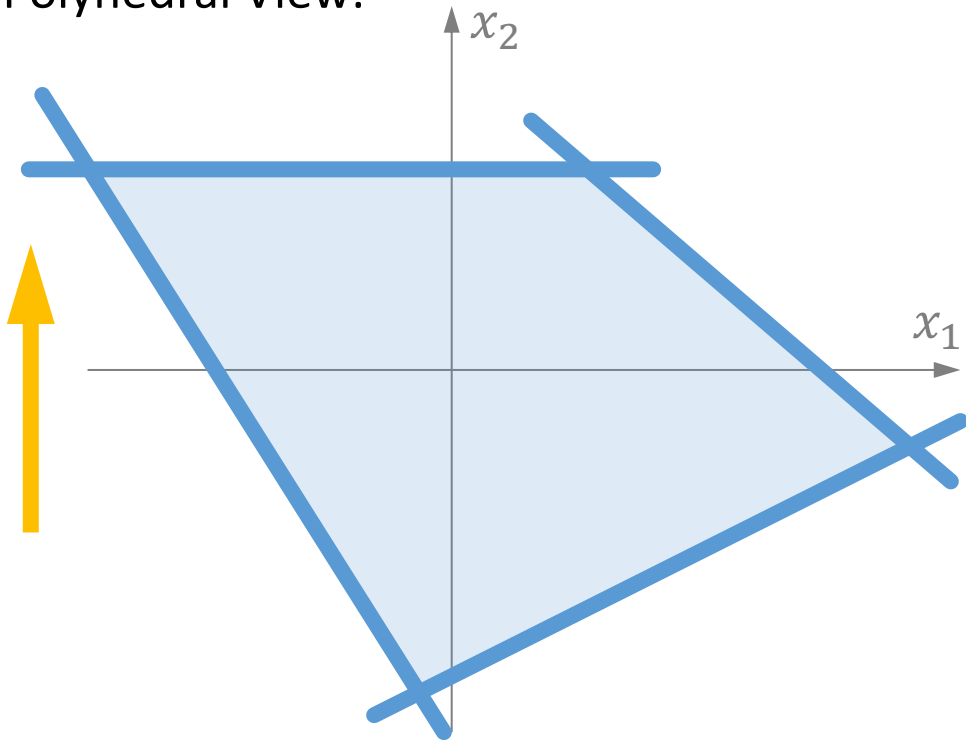
Unified Algorithms for Reverse Data Management

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

objective vector (pointing to \mathbf{c})
constraint vector (pointing to \mathbf{b})
constraint matrix (pointing to \mathbf{A})

Polyhedral View:



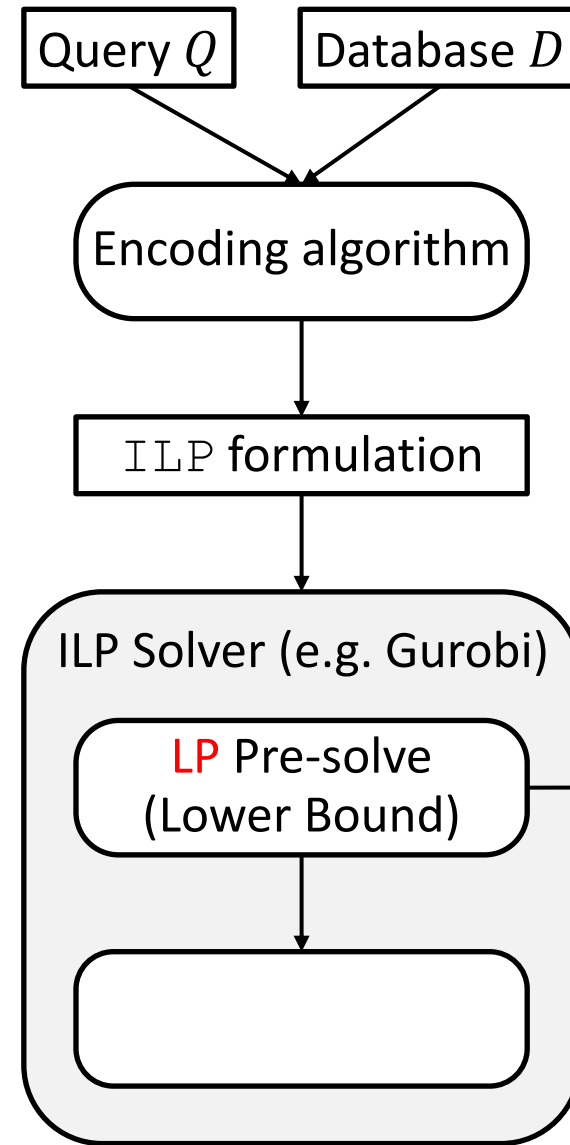
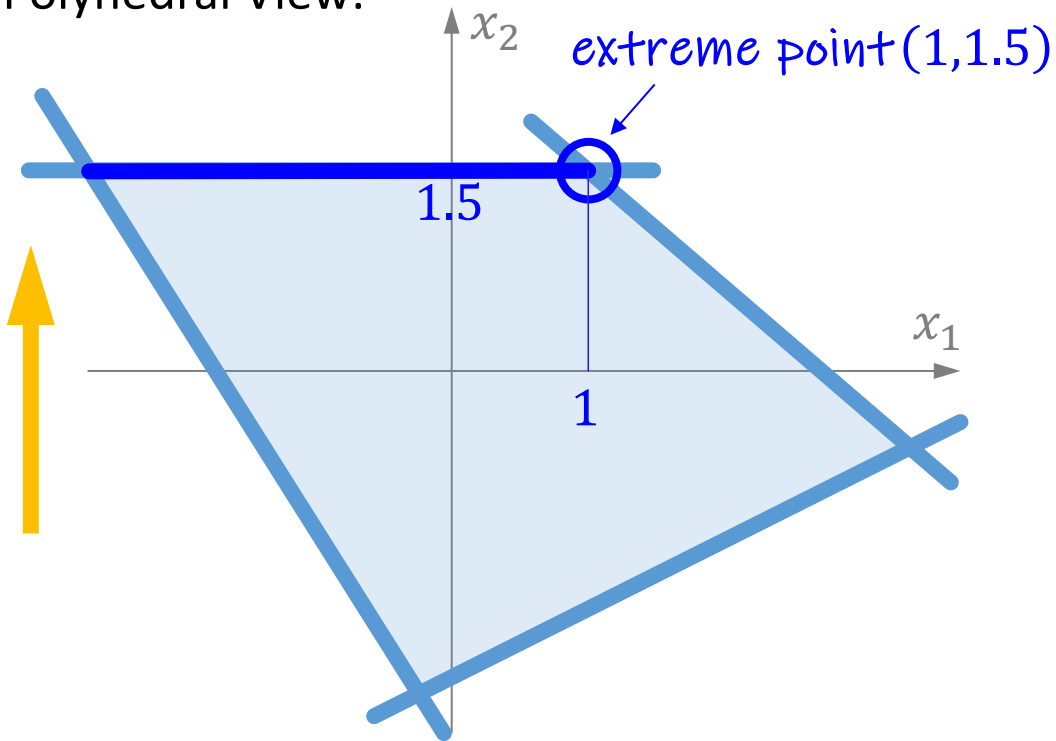
Unified Algorithms for Reverse Data Management

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

objective vector (pointing to \mathbf{c})
constraint vector (pointing to \mathbf{b})
constraint matrix (pointing to \mathbf{A})

Polyhedral View:



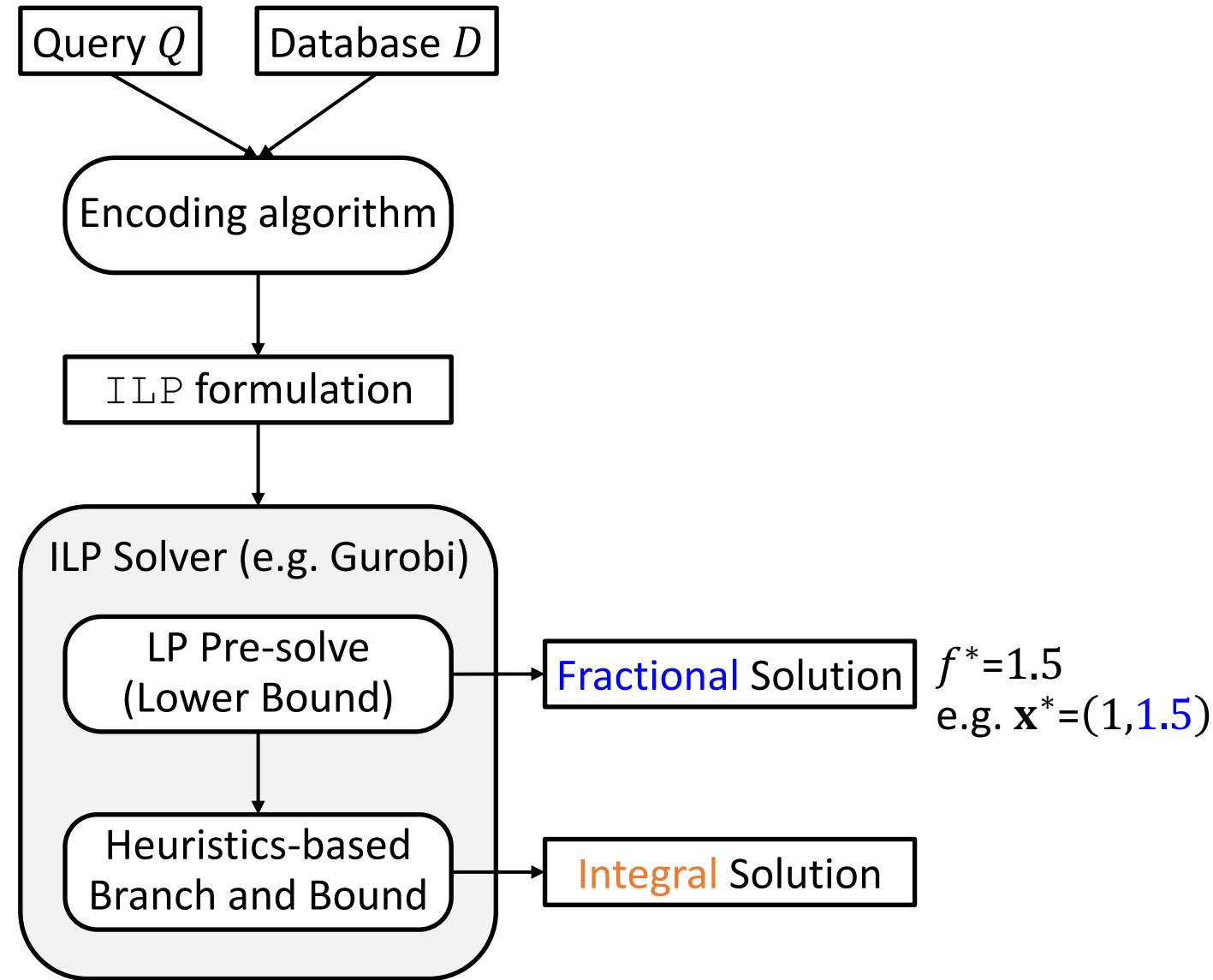
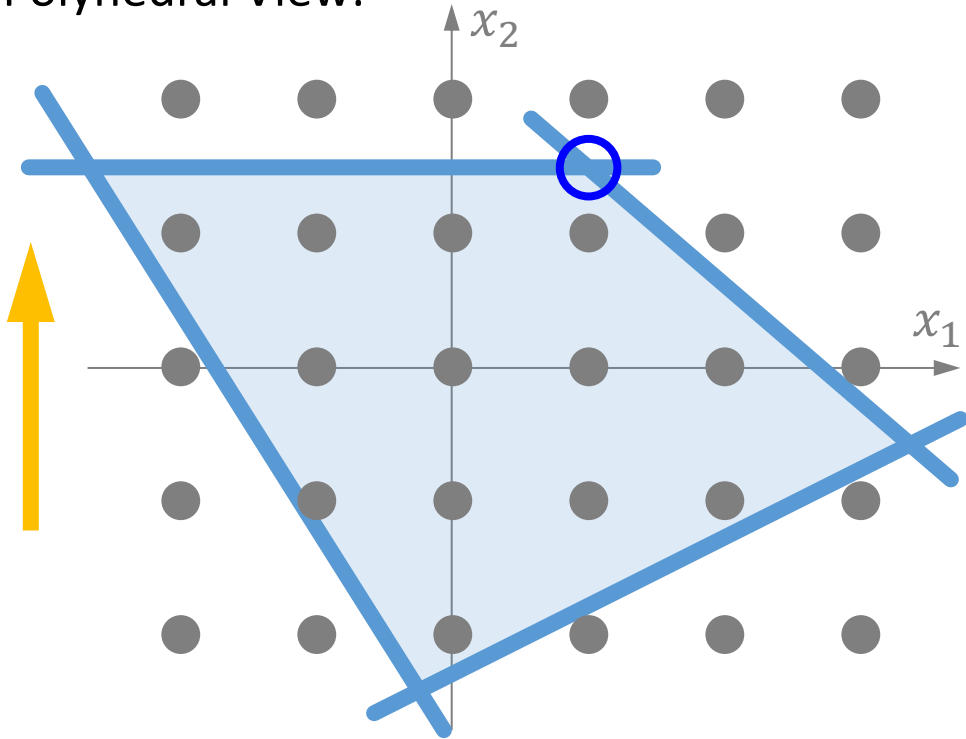
Fractional Solution $f^*=1.5$
e.g. $\mathbf{x}^*=(1,1.5)$

Unified Algorithms for Reverse Data Management

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Polyhedral View:

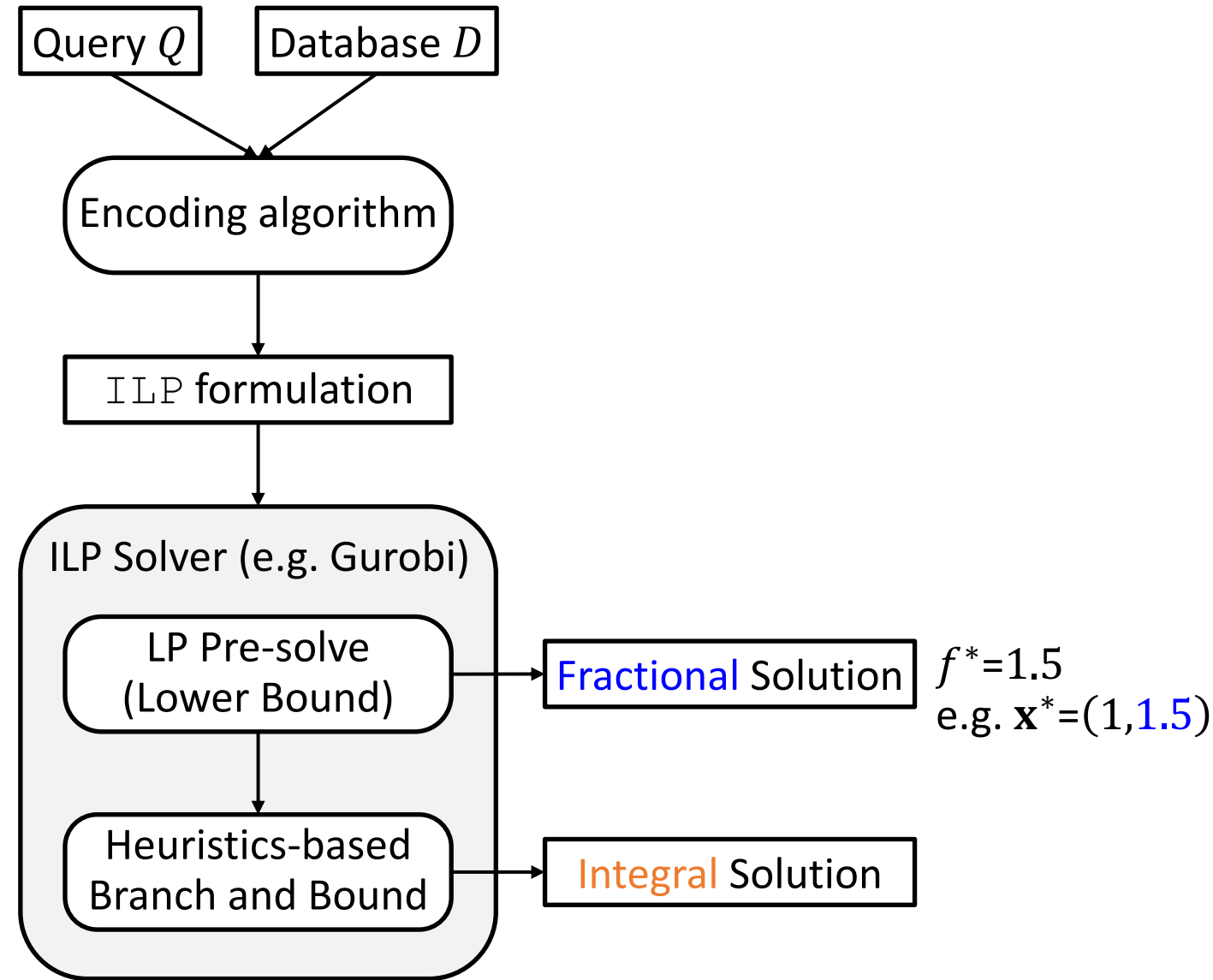
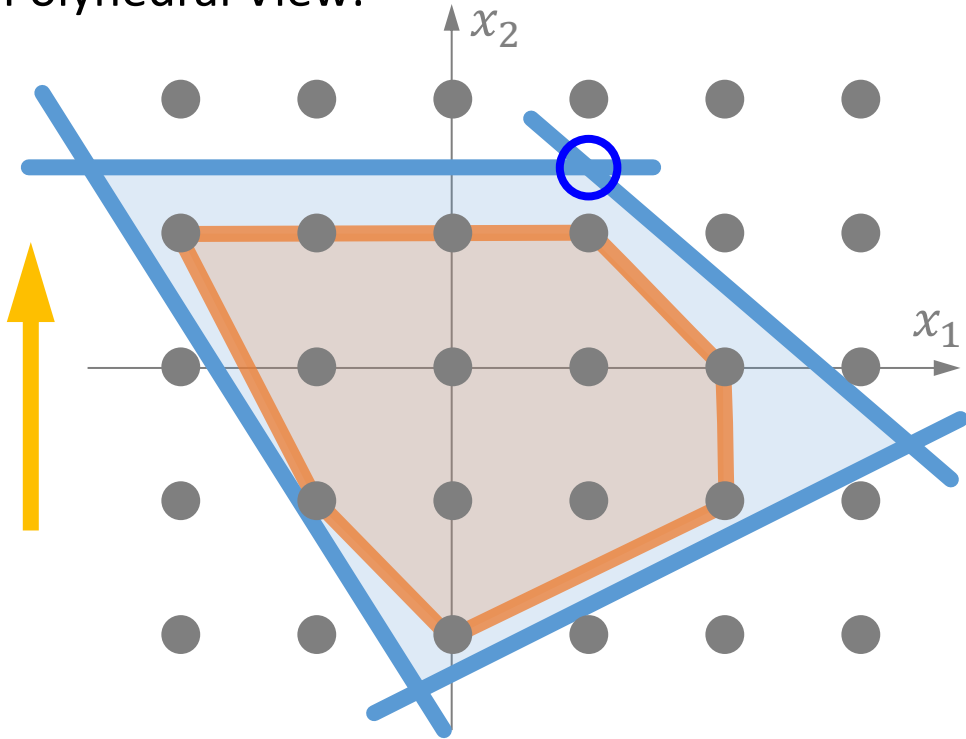


Unified Algorithms for Reverse Data Management

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Polyhedral View:

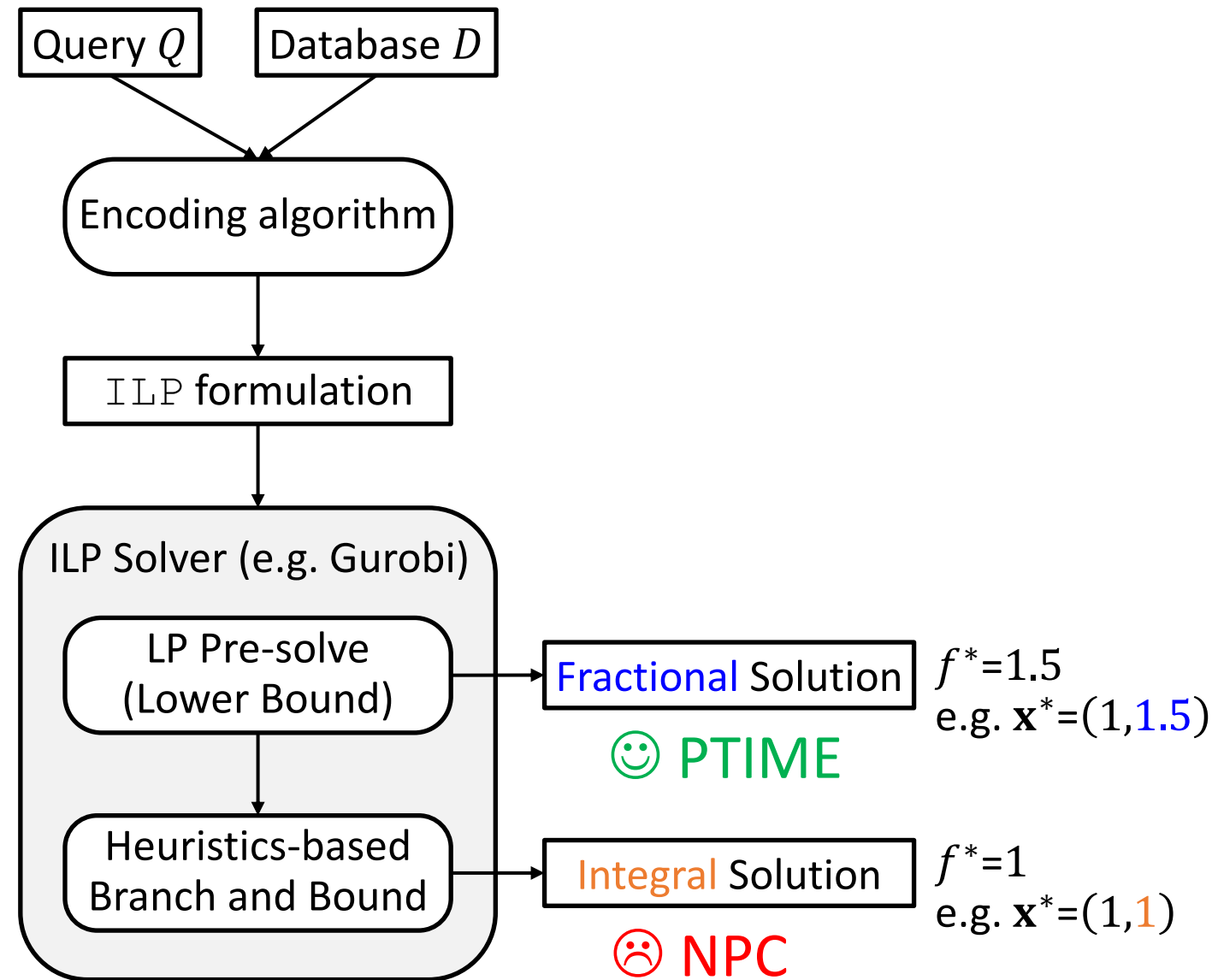
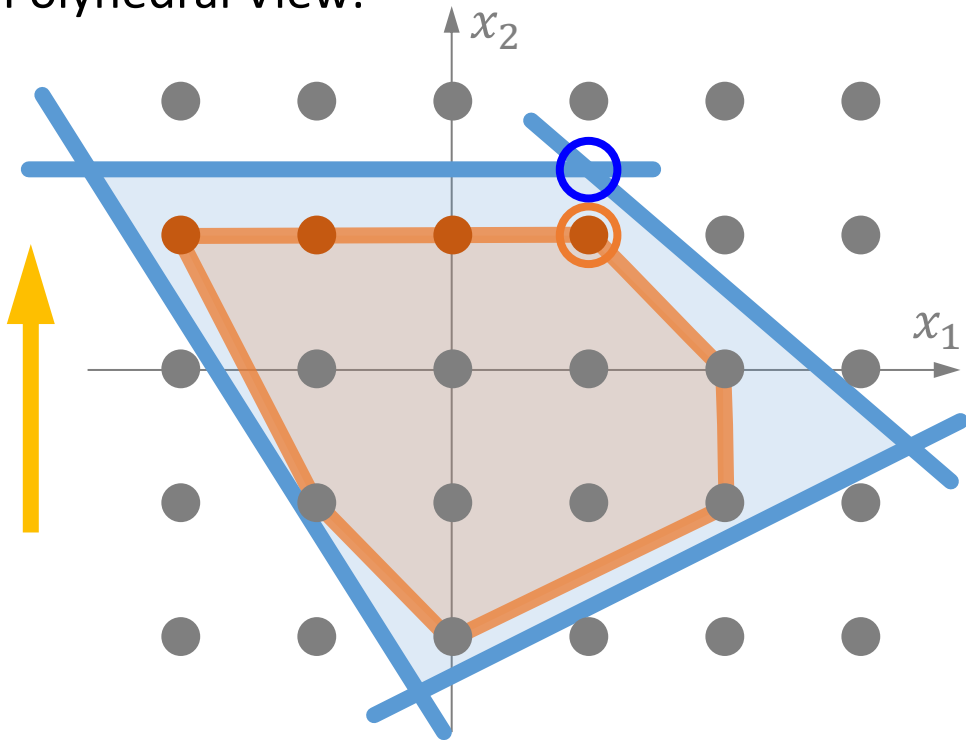


Unified Algorithms for Reverse Data Management

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Polyhedral View:

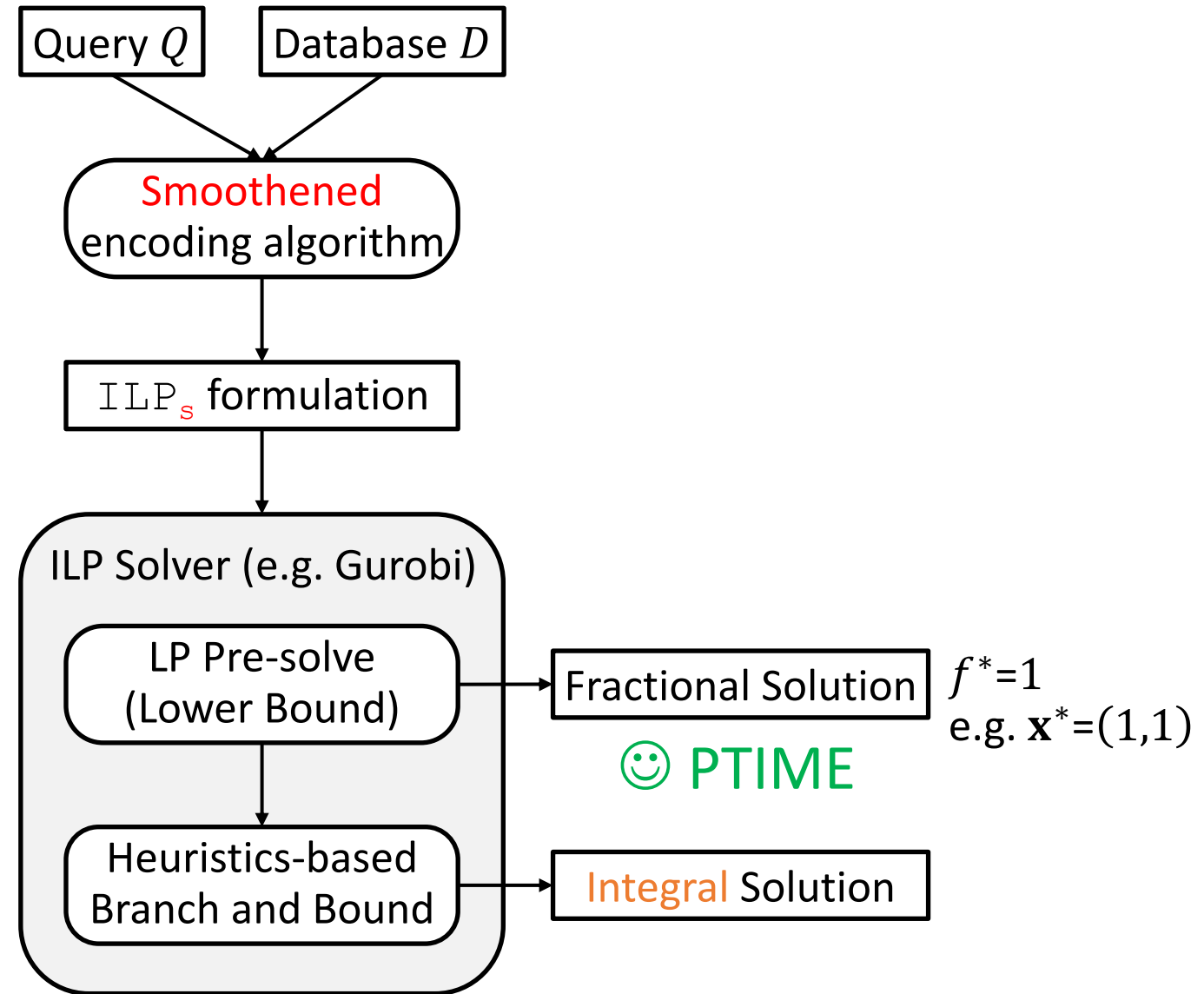
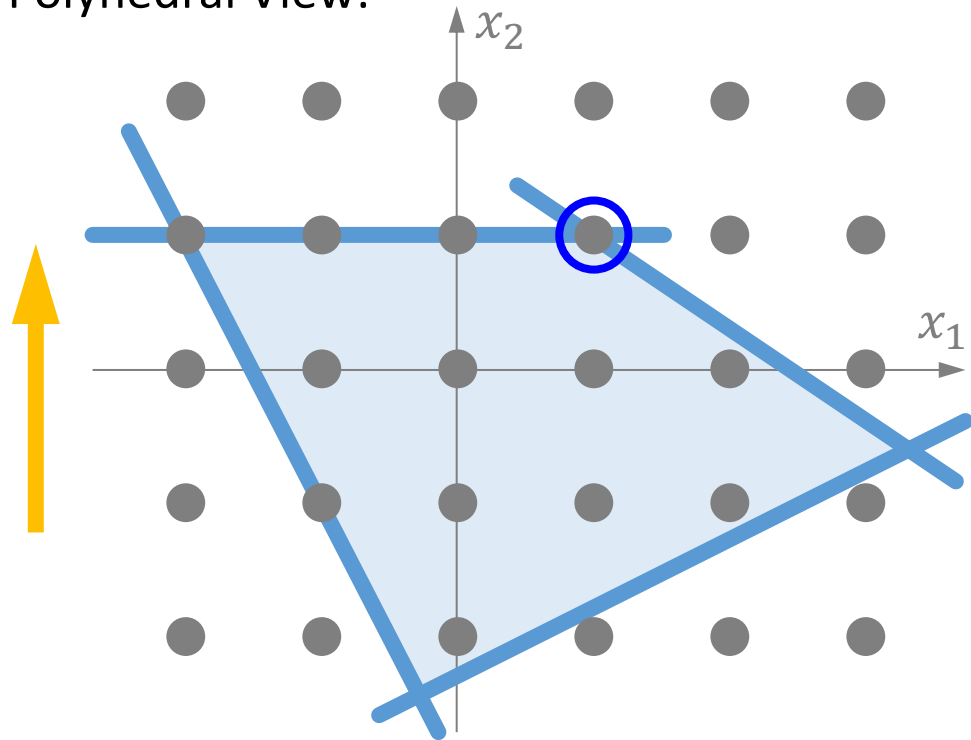


Unified Algorithms for Reverse Data Management

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Polyhedral View:

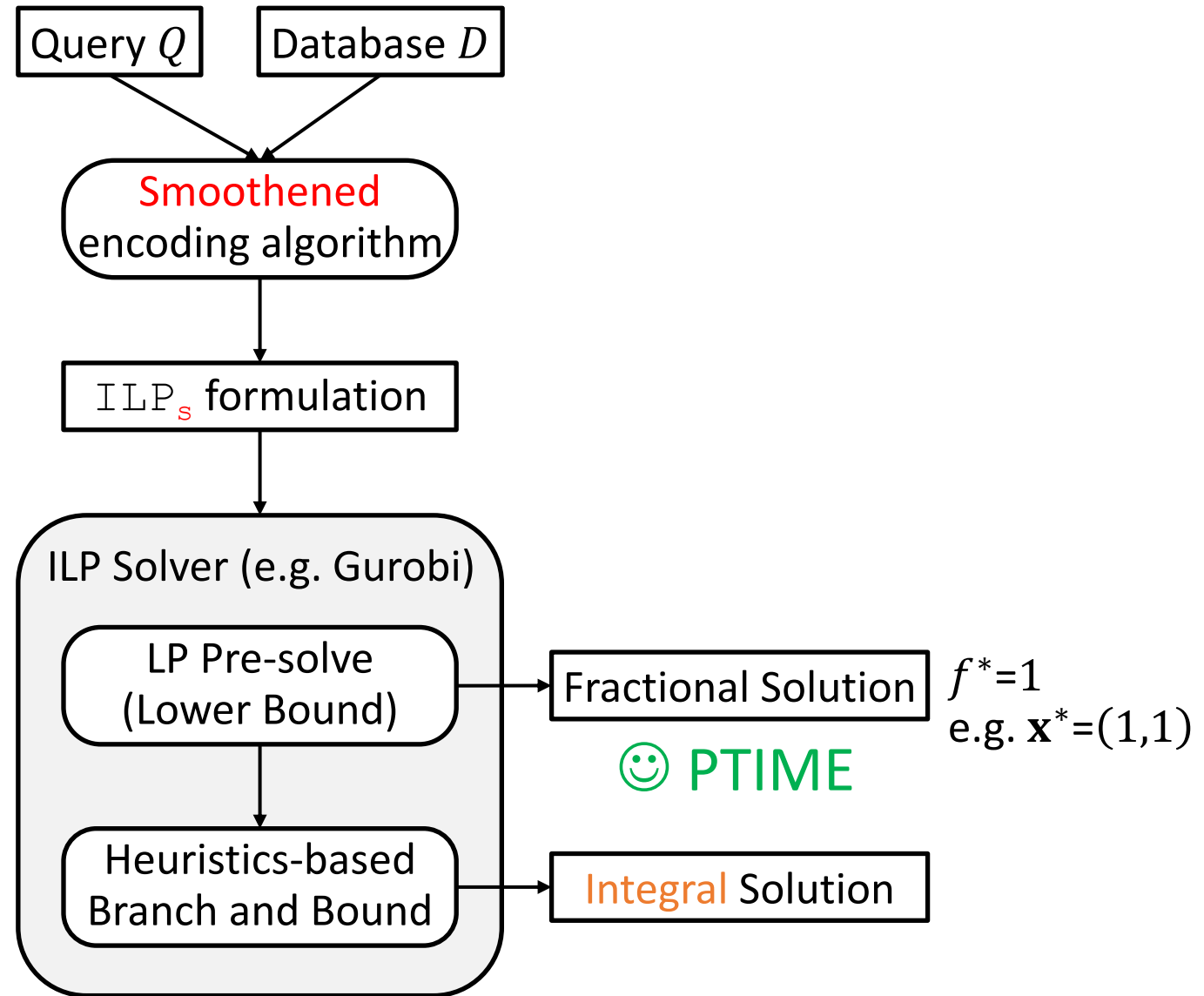
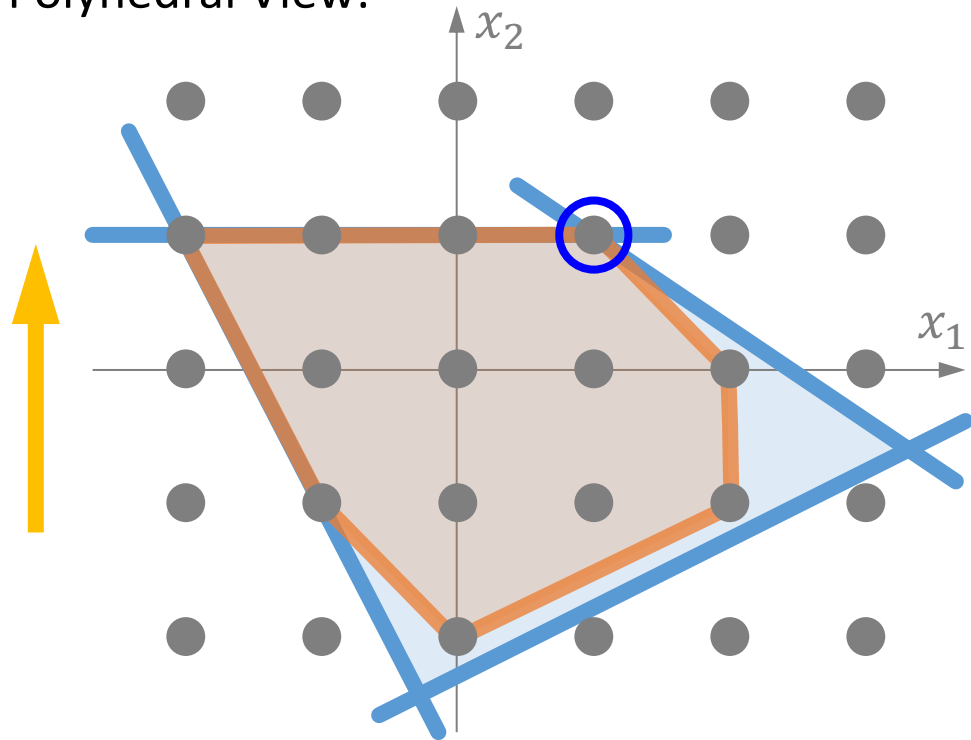


Unified Algorithms for Reverse Data Management

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Polyhedral View:

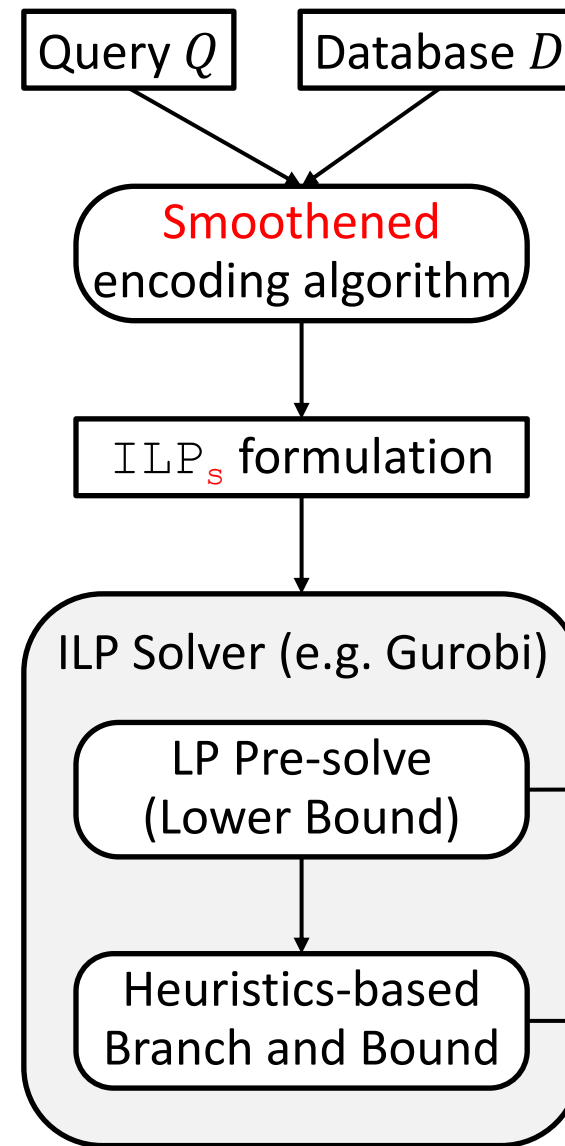
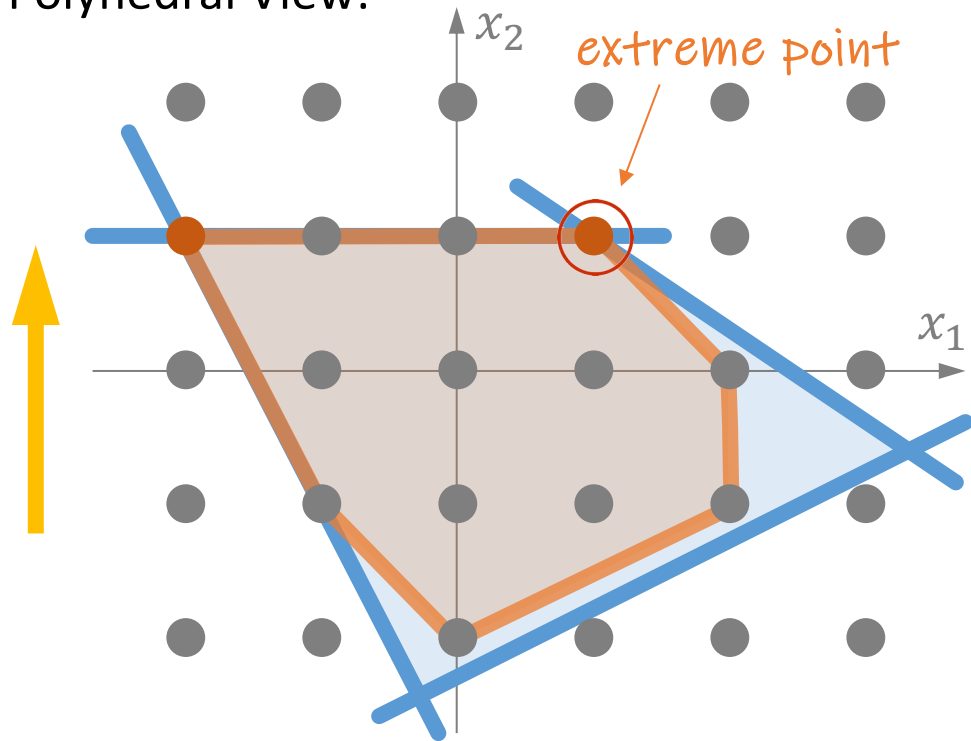


Unified Algorithms for Reverse Data Management

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Polyhedral View:



objective value f^*



If **LP=ILP**, then solvers can find an opt. integral solution efficiently!

Fractional Solution $f^*=1$
e.g. $\mathbf{x}^*=(1, \mathbf{1})$
☺ PTIME
identical !!!

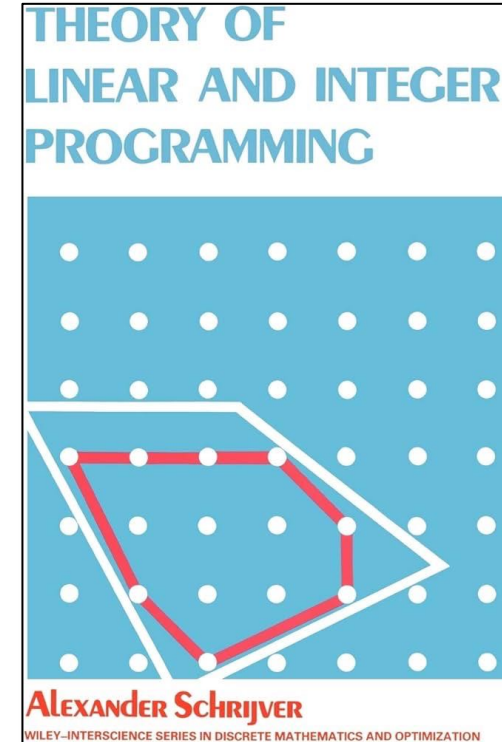
Integral Solution $f^*=1$
e.g. $\mathbf{x}^*=(1, \mathbf{1})$
☺ PTIME

When is LP=ILP according to the literature?

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

83	Balanced and unimodular hypergraphs	1439
83.1	Balanced hypergraphs	1439
83.2	Characterizations of balanced hypergraphs	1440
83.2a	Totally balanced matrices	1444
83.2b	Examples of balanced hypergraphs	1447
83.2c	Balanced $0, \pm 1$ matrices	1447
83.3	Unimodular hypergraphs	1448
83.3a	Further notes	1450



Alexander Schrijver

Combinatorial Optimization

Polyhedra and Efficiency

Volume A-C

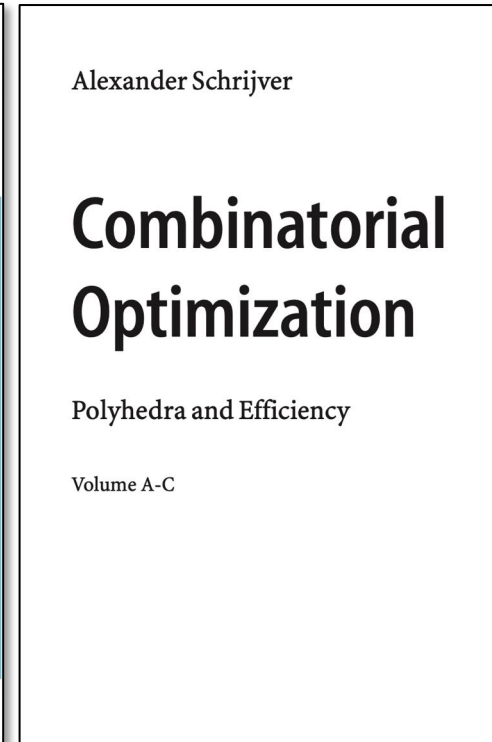
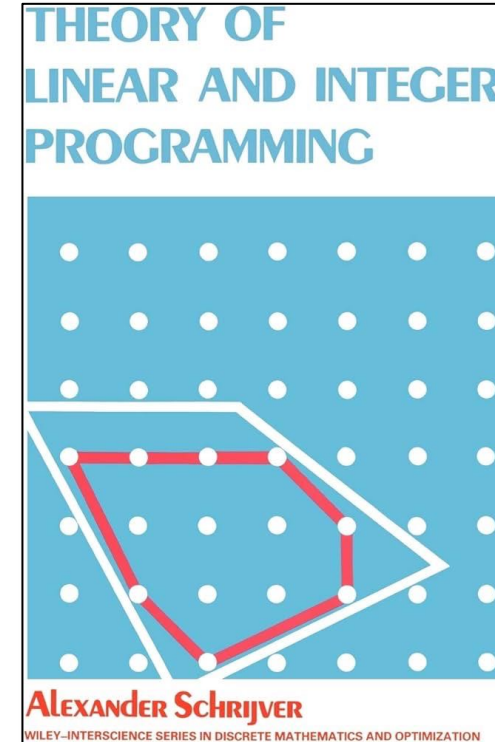
When is LP=ILP according to the literature: not useful 😞

ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Annotations:
 - \mathbf{c} : objective vector
 - \mathbf{A} : constraint matrix
 - \mathbf{b} : constraint vector

83	Balanced and unimodular hypergraphs	1439	😞
83.1	Balanced hypergraphs	1439	
83.2	Characterizations of balanced hypergraphs	1440	
83.2a	Totally balanced matrices	1444	
83.2b	Examples of balanced hypergraphs	1447	
83.2c	Balanced $0, \pm 1$ matrices	1447	
83.3	Unimodular hypergraphs	1448	😞
83.3a	Further notes	1450	



- Focus of polyhedral theory mainly on constraint matrix \mathbf{A} . But our PTIME constraint matrixes need not be balanced, nor Totally Unimodular, etc.

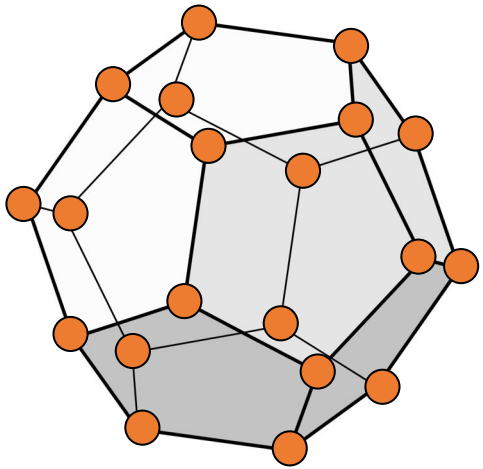
1. Our complexity results take into account the objective vector \mathbf{c} !
2. This gives us a separation between the problem under set vs. bag semantics!
3. We use an indirect proof via problem-specific MFMC encodings 😊

So what do we do to show ILP=LP for PTIME cases?

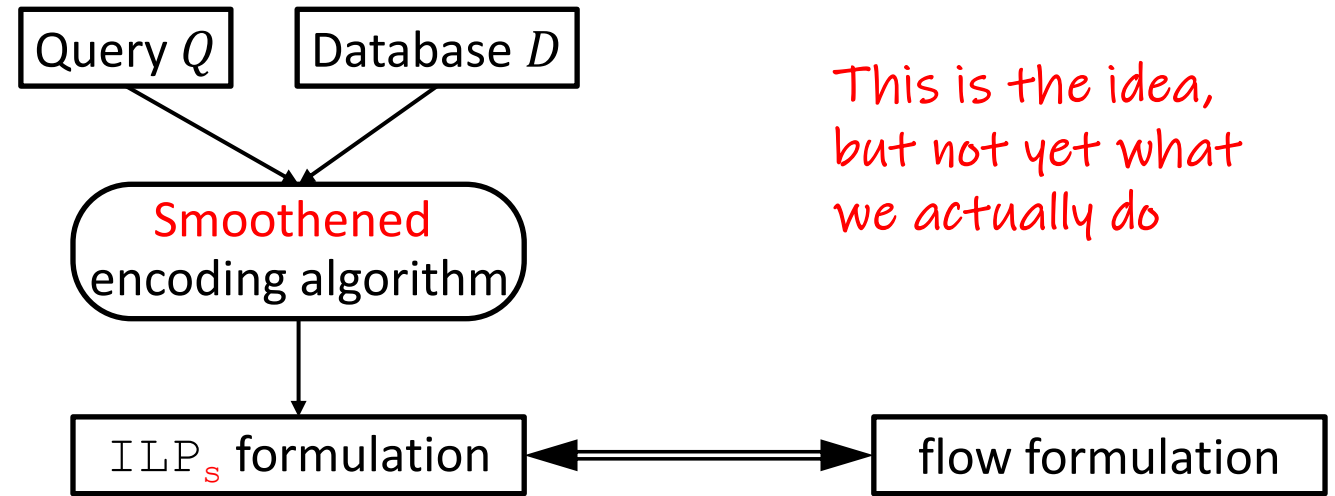
ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Polyhedral View



- For an "ideal" **constraint matrix**, the vertices (extreme points) of its polytope are all **integral**



This is the idea, but not yet what we actually do



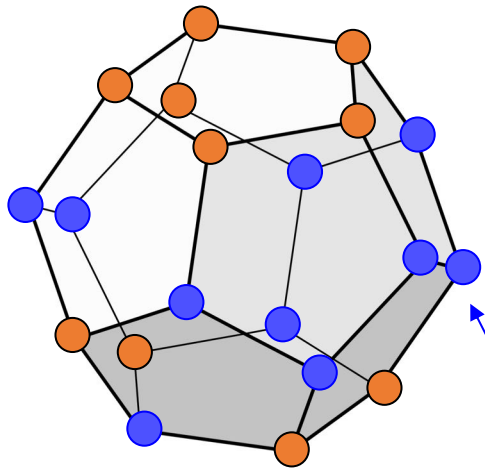
1. Show correspondence to a flow formulation
2. **Constraint matrix** represented by a flow graph is an "ideal matrix"
3. COROLLARY: LP = ILP, **PTIME** 😊

So what do we do to show ILP=LP for PTIME cases?

ILP formulation:

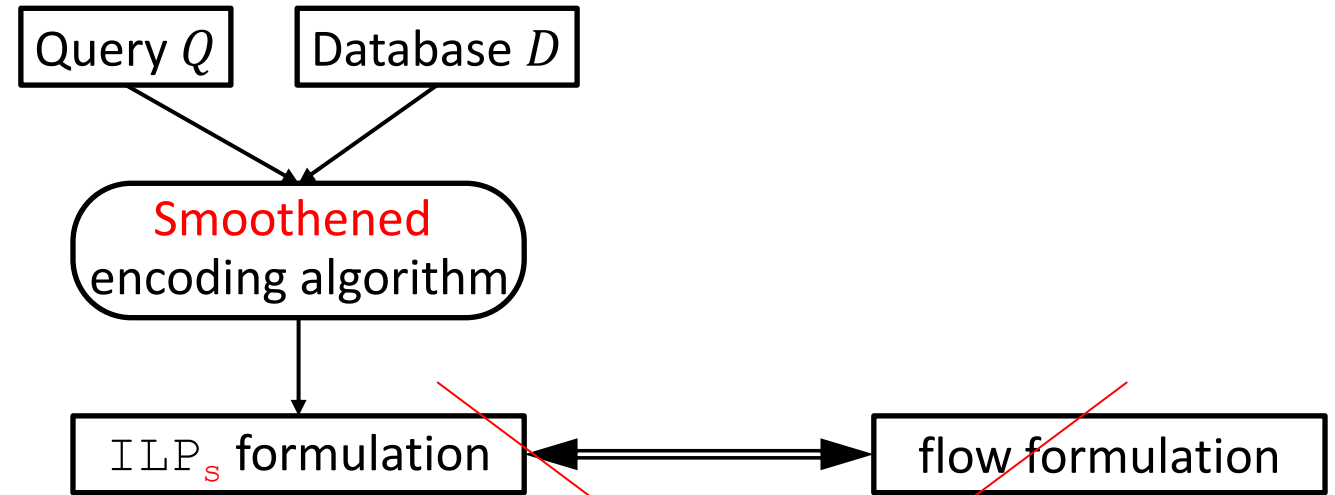
$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Polyhedral View



- For an "ideal" **constraint matrix**, the vertices (extreme points) of its polytope are all integral

fractional points



1. Show correspondence to a flow formulation
2. **Constraint matrix** represented by a flow graph is an "ideal matrix"
3. COROLLARY: LP = ILP, **PTIME** 😊

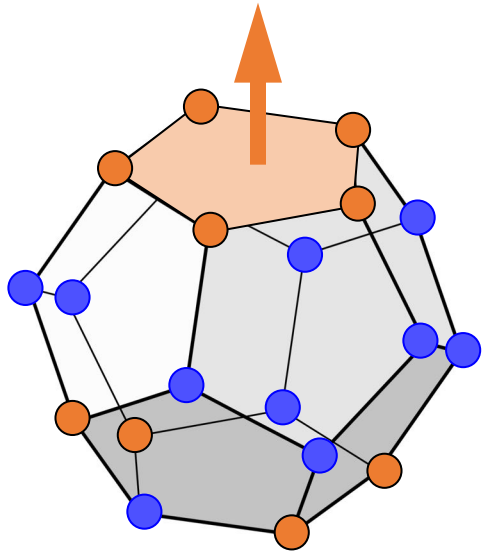
But it is not as easy!
We have fractional extremal points ☹️

So what do we do to show ILP=LP for PTIME cases?

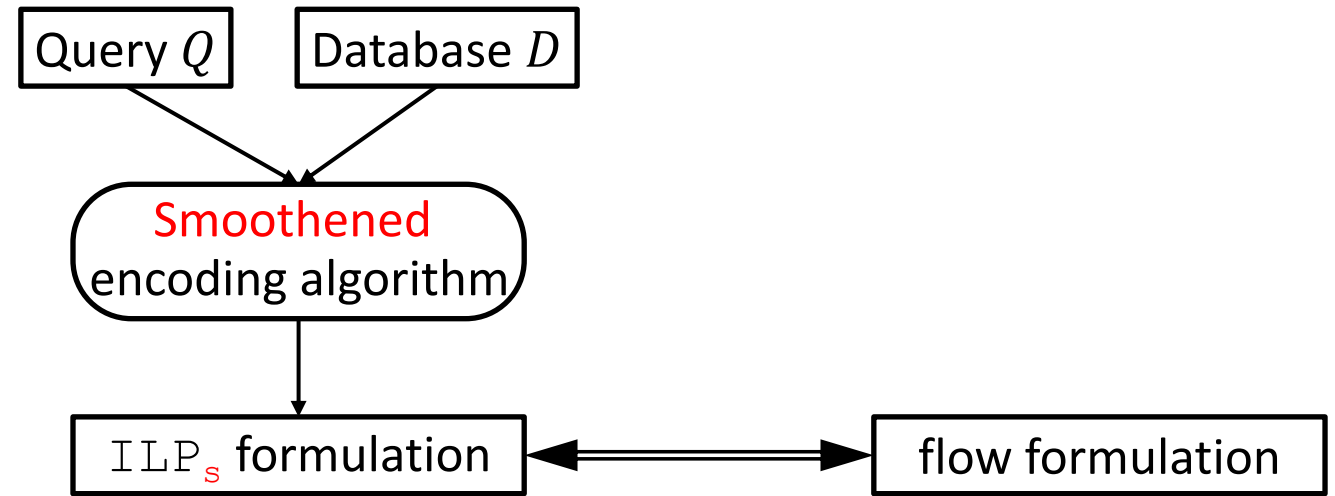
ILP formulation:

$$\begin{aligned} f^* &= \min[\mathbf{c} \cdot \mathbf{x}] \\ \text{s.t. } \mathbf{A} \cdot \mathbf{x} &\geq \mathbf{b} \\ \mathbf{x} &\in \mathbb{N}^n / \mathbb{R}^n \end{aligned}$$

Polyhedral View



- The **optimal solution face** contains only **integral vertices**
- Polytope may have **non-integral vertices**



1. Show correspondence *after processing* that takes the **objective vector** into account
2. *After processing*, constraint matrix is ideal
3. COROLLARY: LP = ILP, **PTIME** 😊

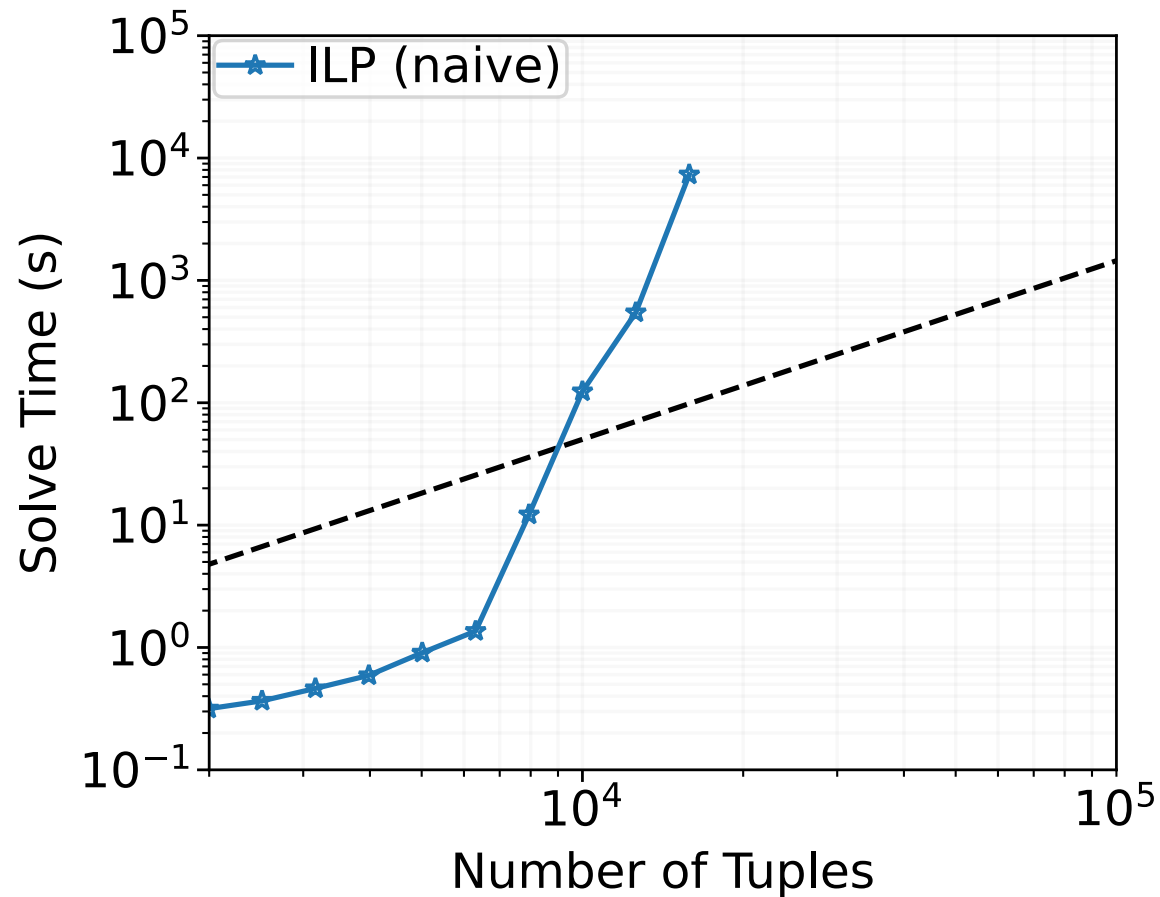


- Showing such correspondences for all PTIME cases in all scenarios is non-trivial
- This moves the challenge from algorithm development to proofs!

Scalability of Naive vs. Smoothened ILP for PTIME query

Finding an optimal solution for the Smallest Witness Problem ([SIGMOD'19], [ICDT'24])

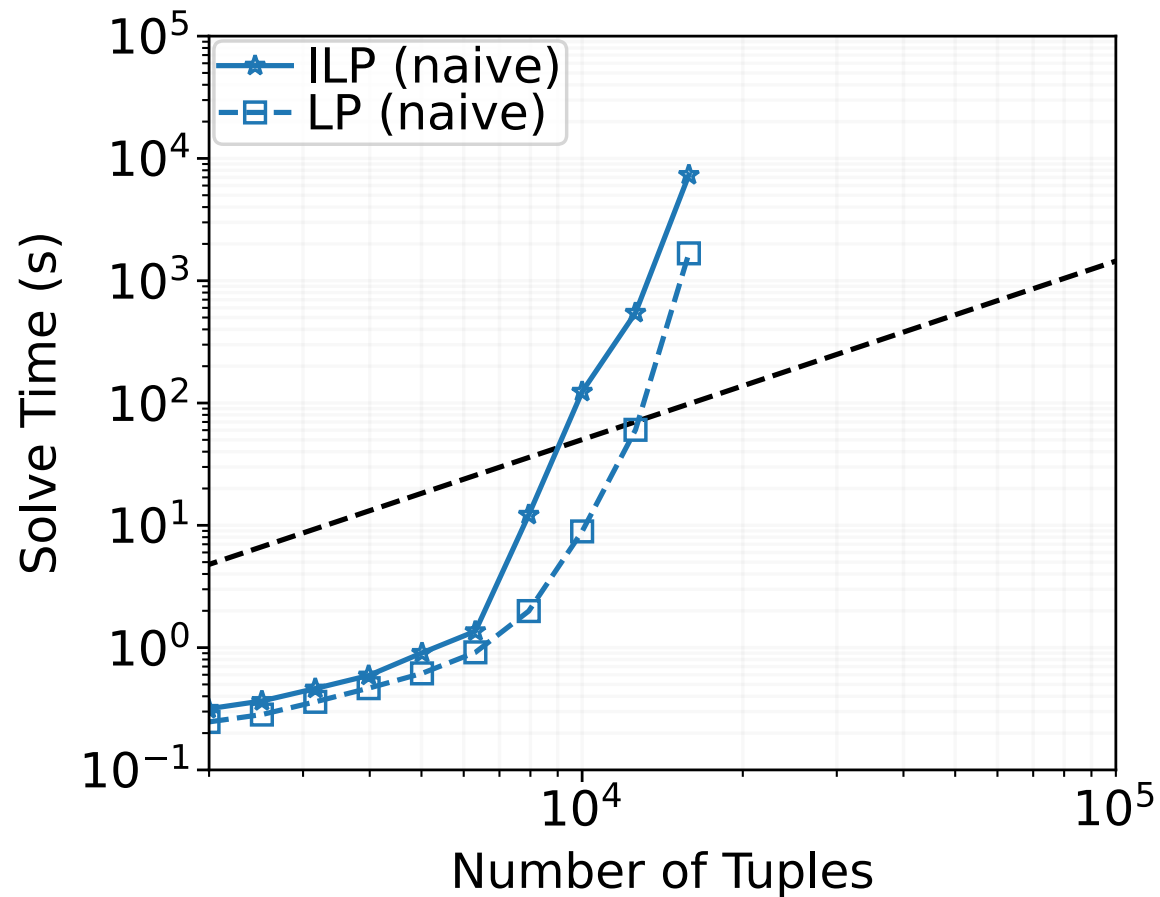
$Q_{5\text{star}}(x): \neg R(x, a), S(x, b), T(x, c), U(x, d), V(x, e), A(x)$



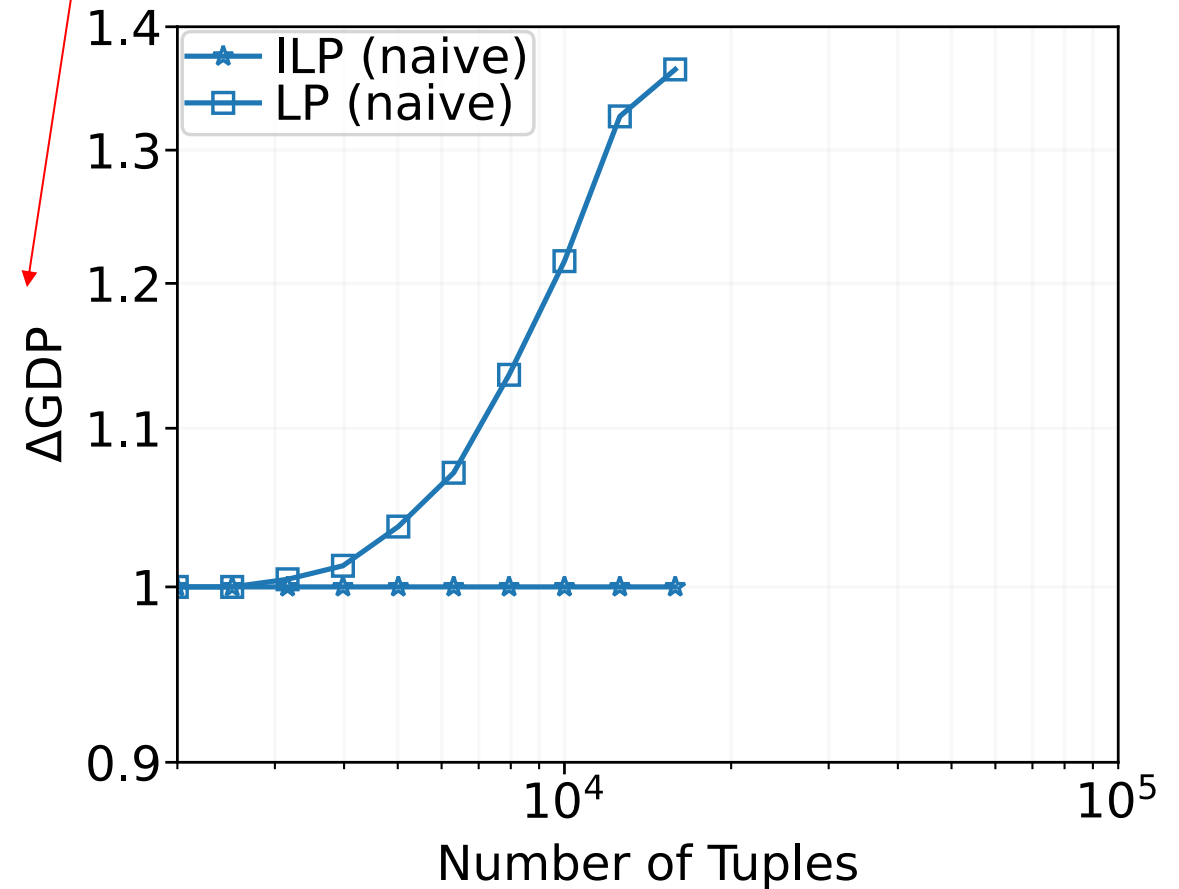
Scalability of Naive vs. Smoothened ILP for PTIME query

Finding an optimal solution for the Smallest Witness Problem ([SIGMOD'19], [ICDT'24])

$Q_{5\text{star}}(x): -R(x, y), S(x, z), T(x, u), U(x, v), V(x, w), W(x)$



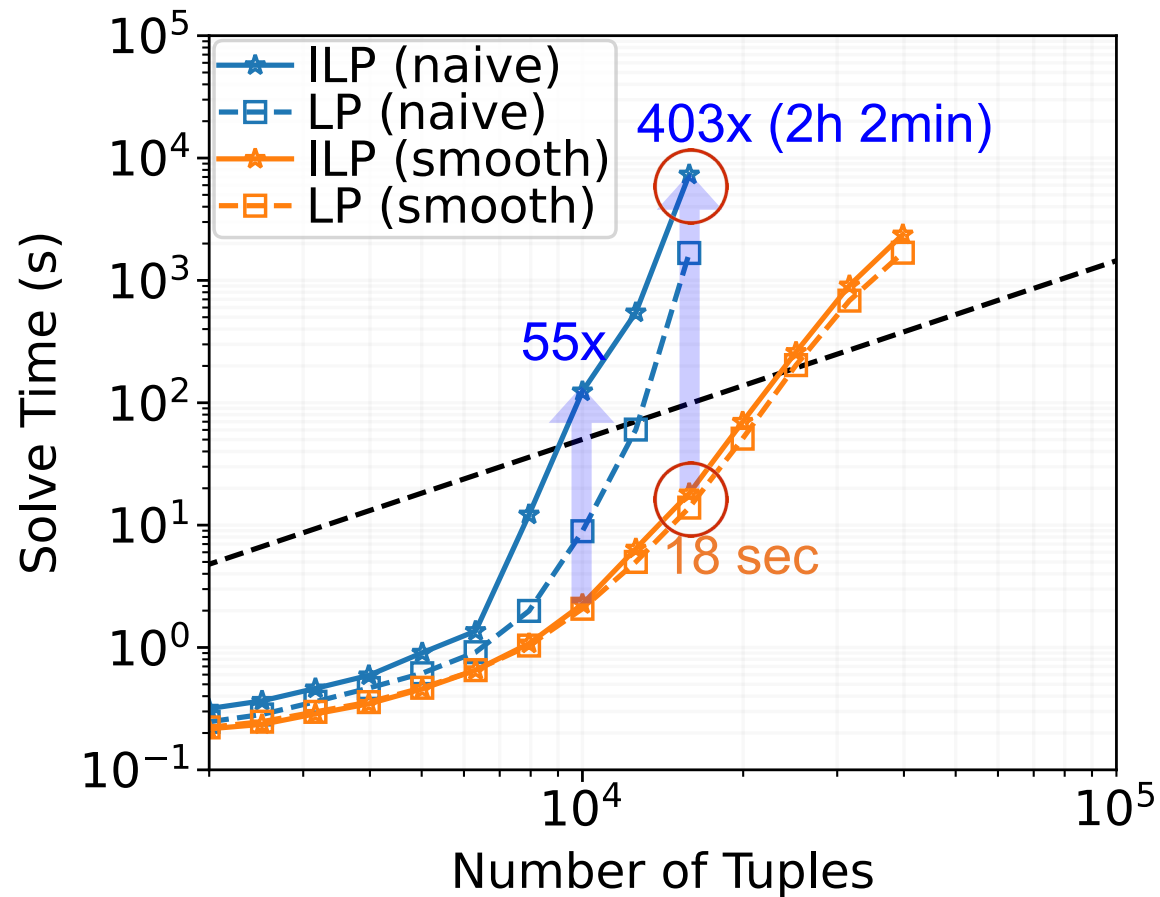
How different are the optimal values for the ILP formulation and its LP relaxation



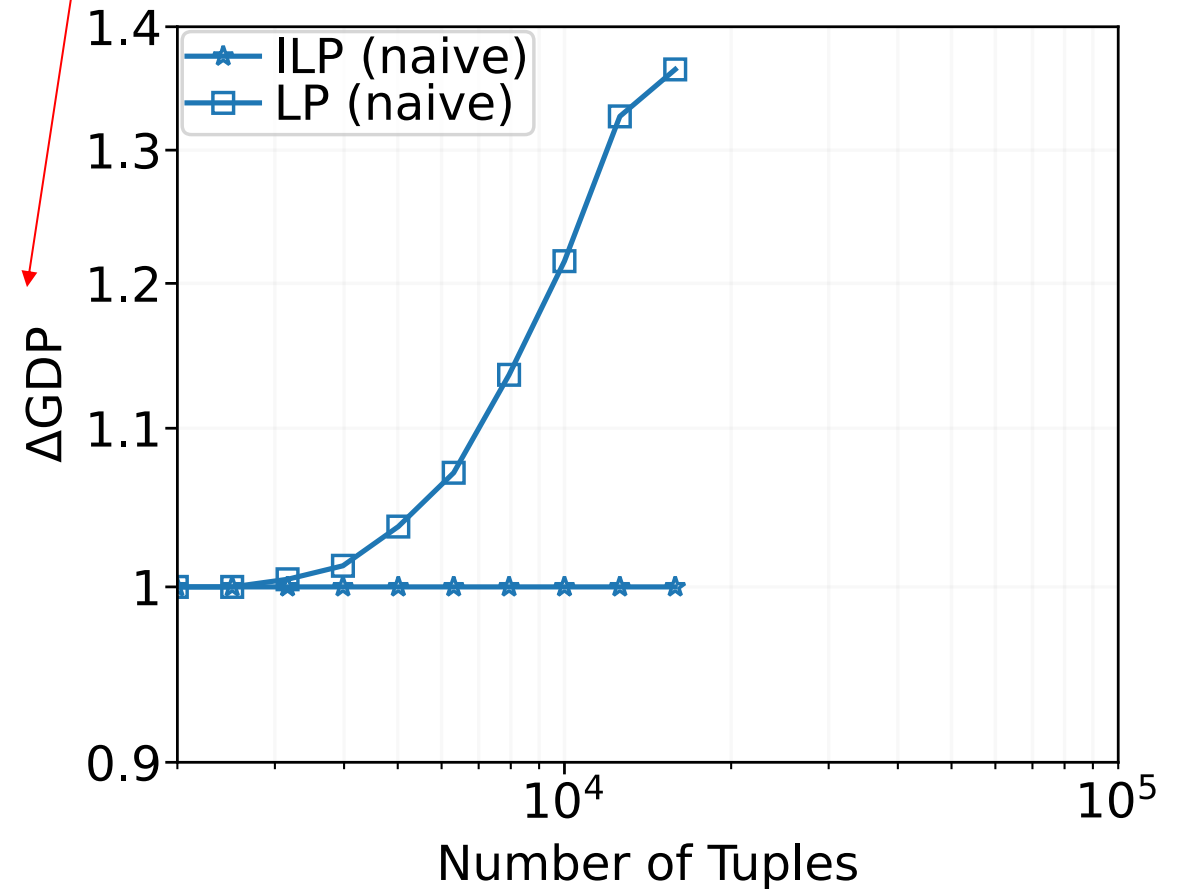
Scalability of Naive vs. Smoothed ILP for PTIME query

Finding an optimal solution for the Smallest Witness Problem ([SIGMOD'19], [ICDT'24])

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



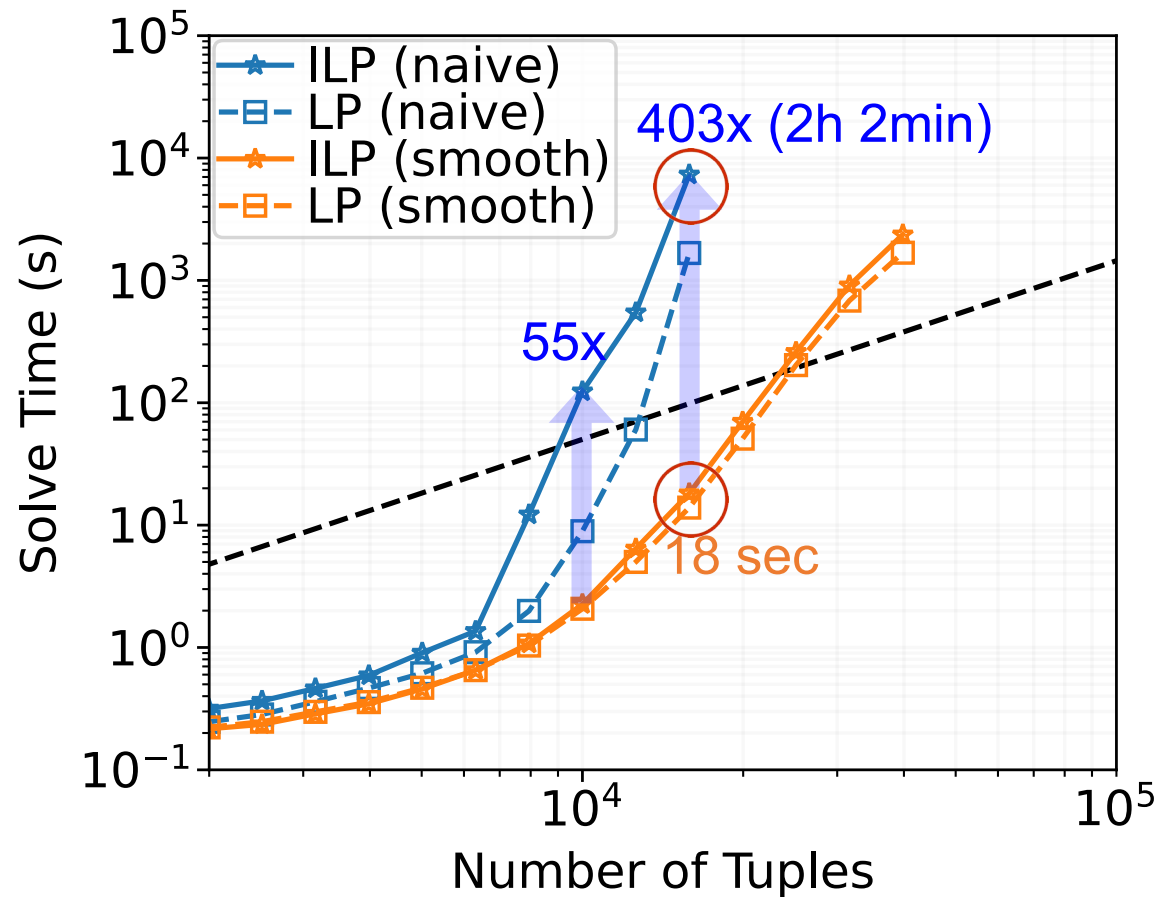
How different are the optimal values for the ILP formulation and its LP relaxation



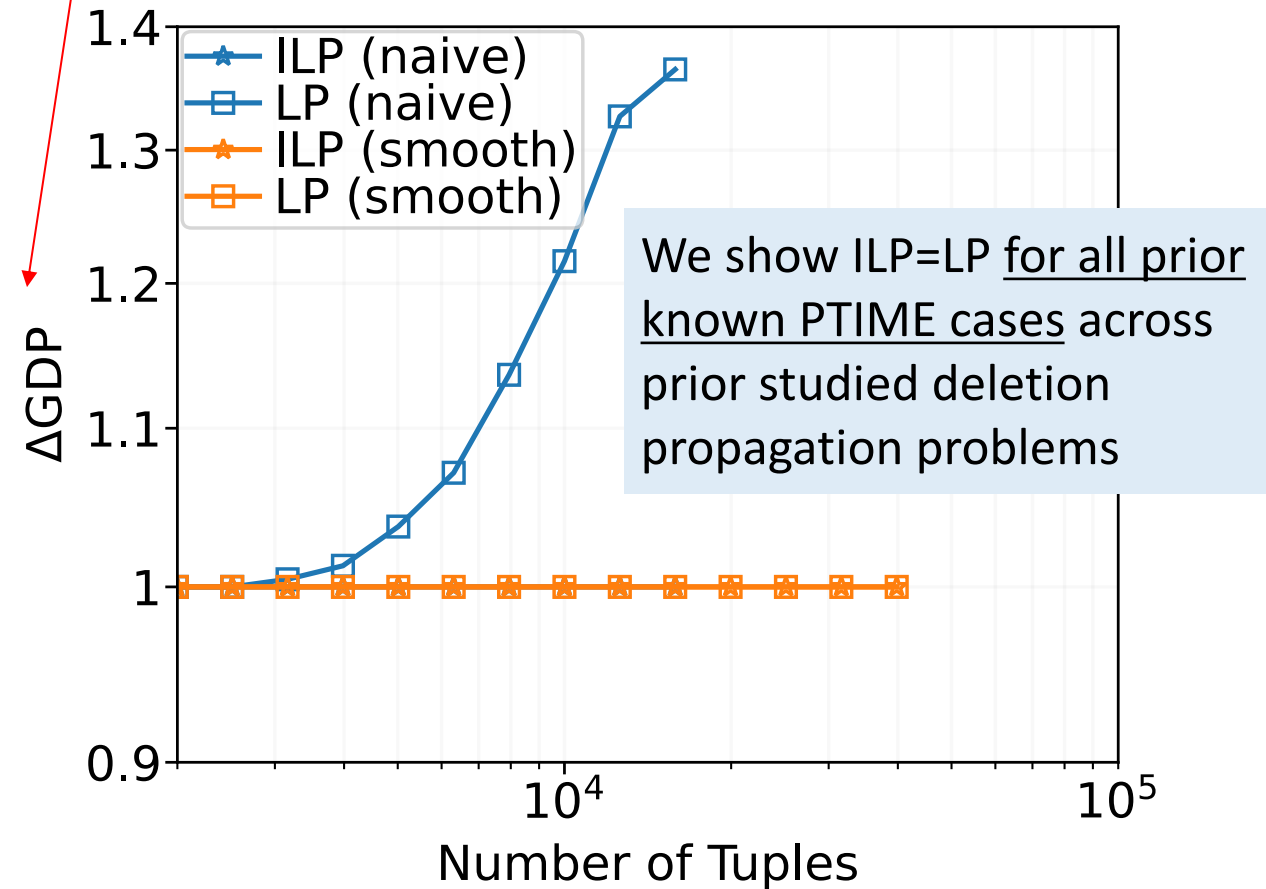
Scalability of Naive vs. Smoothed ILP for PTIME query

Finding an optimal solution for the Smallest Witness Problem ([SIGMOD'19], [ICDT'24])

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



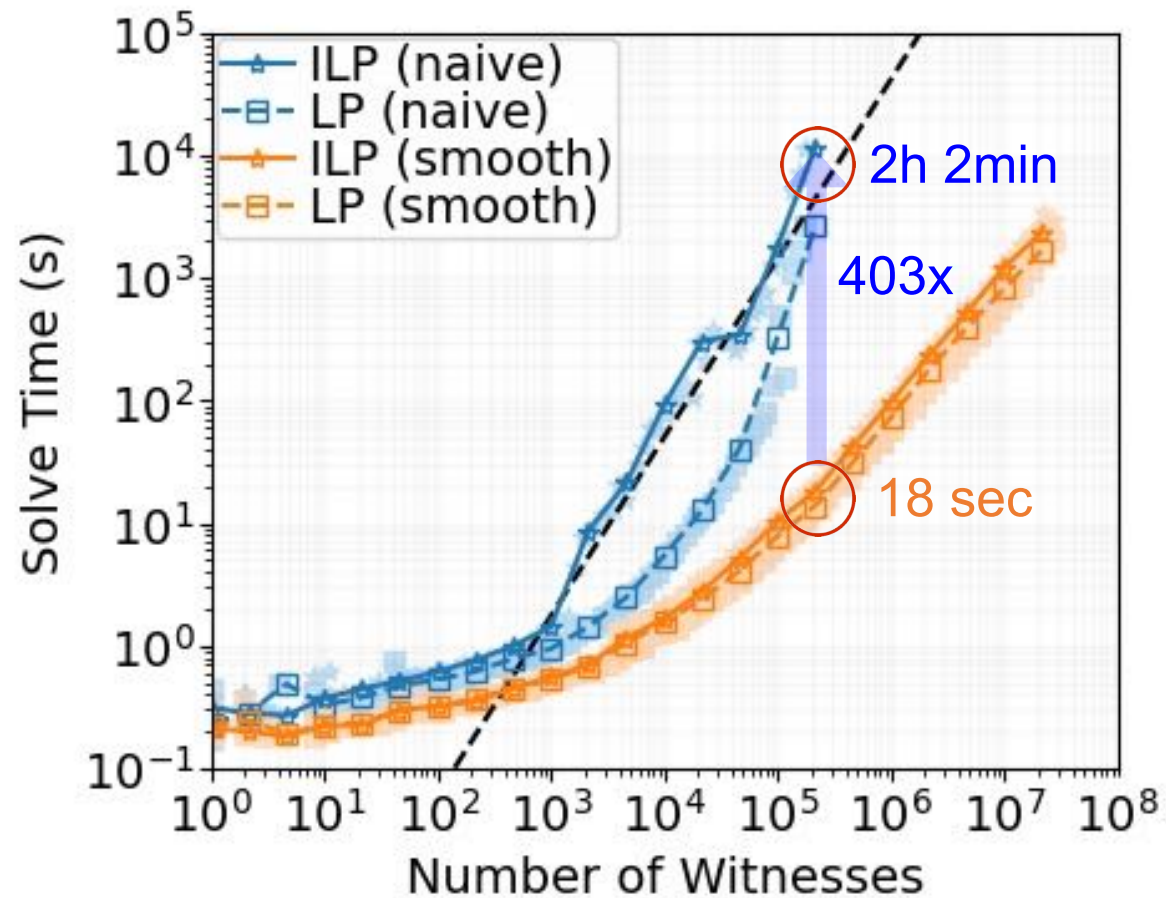
How different are the optimal values for the ILP formulation and its LP relaxation



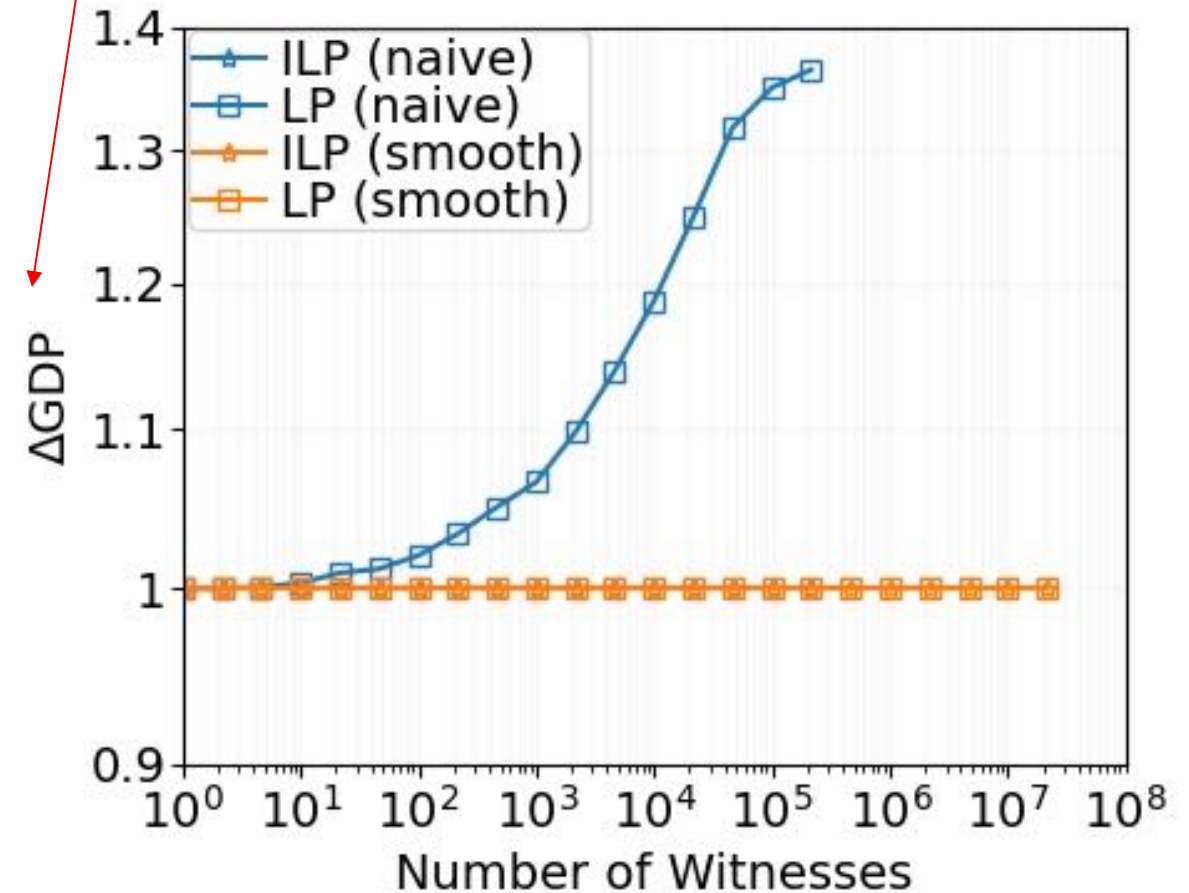
Scalability of Naive vs. Smoothed ILP for PTIME query

Finding an optimal solution for the Smallest Witness Problem ([SIGMOD'19], [ICDT'24])

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



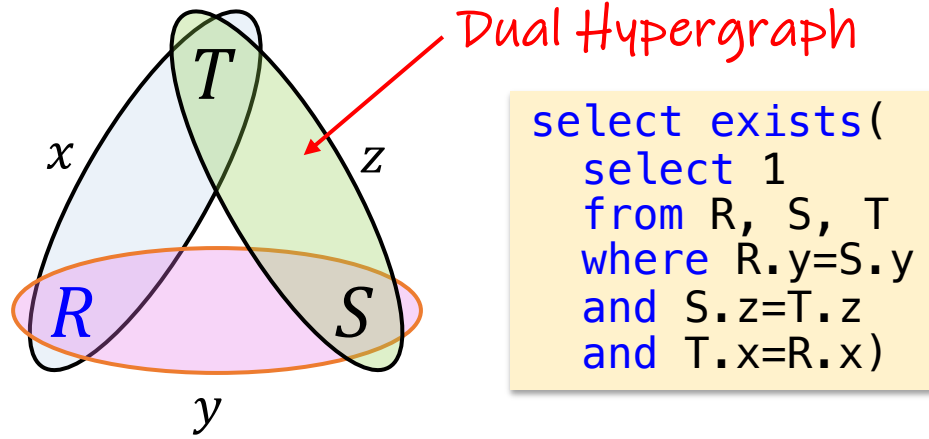
How different are the optimal values for the ILP formulation and its LP relaxation



Example complicated landscape for resilience

Triangle query

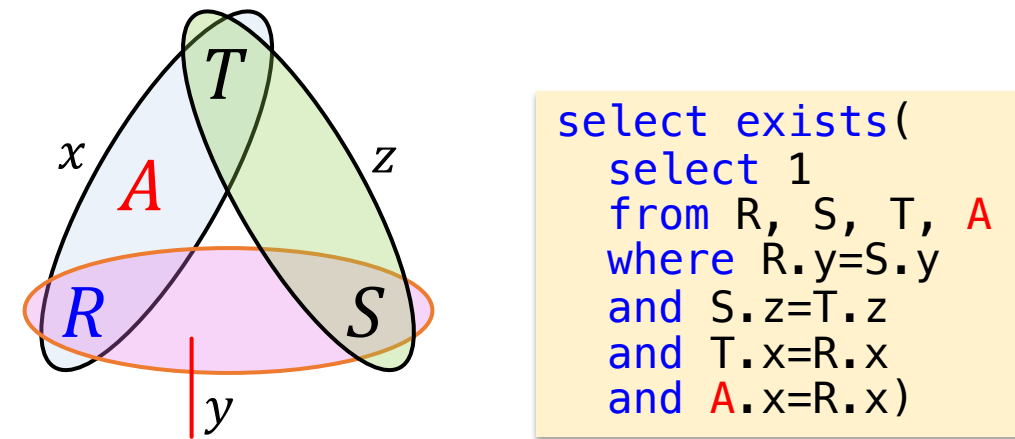
$Q^\Delta: \neg R(x, y), S(y, z), T(x, z)$



NPC

Triangle unary

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

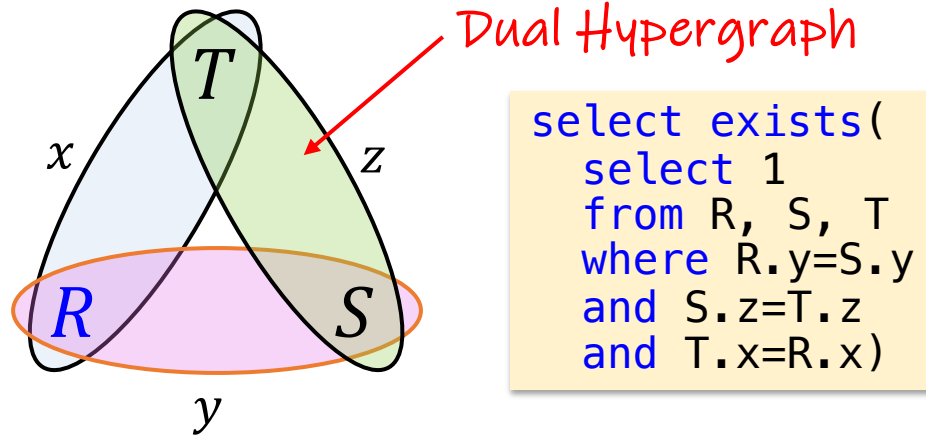


PTIME

Example complicated landscape for resilience

Triangle query

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



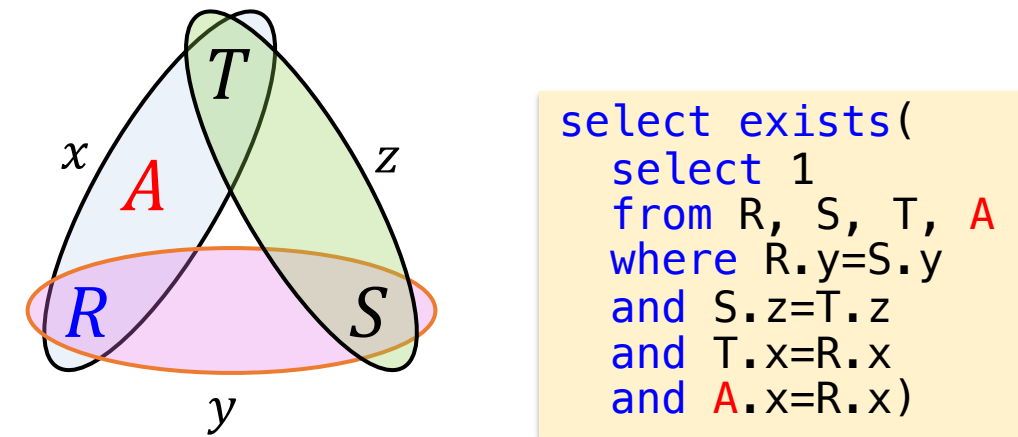
NPC

PTIME for FD $x \rightarrow y$

PTIME if provenance happens to be read-once

Triangle unary

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



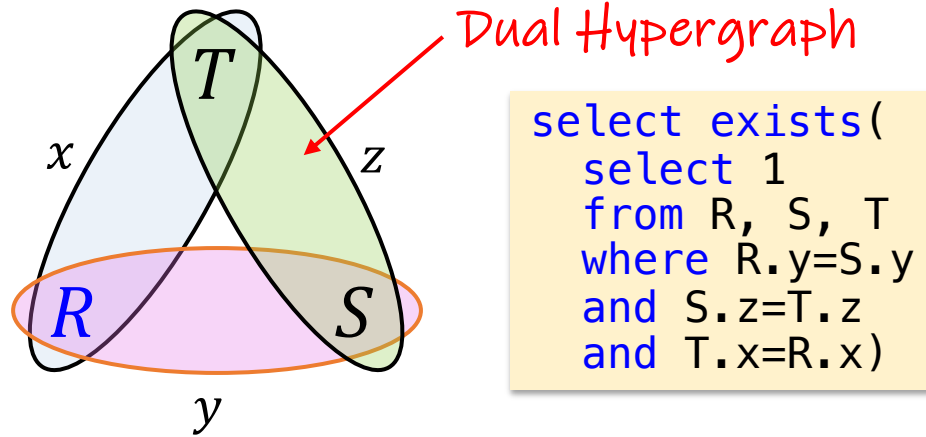
PTIME

NPC under bag semantics

Example complicated landscape for resilience

Triangle query

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



NPC

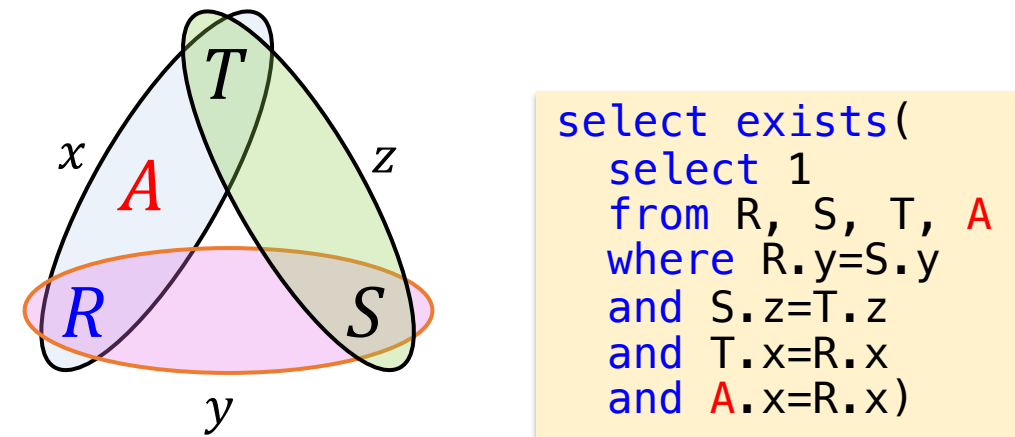
PTIME for FD $x \rightarrow y$

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

Outline

1. Reverse Data Management (RDM)
2. A magical ILP formulation
3. Take-aways

(An illustrated example: only if time remains)

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 ✓

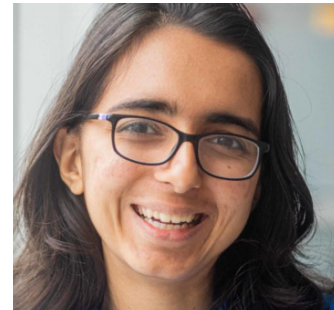
Complete solutions: All problem types are covered ✓ (including self-joins or bag semantics)

Practical aspects: it works from day one ✓

Theory

Take-aways: "Is ILP all you need? ..."

- 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 also talk to Neha 😊
Faculty at UMass Amherst from Fall'25



Thank you 😊