



# A Unified Approach For Reverse Data Management

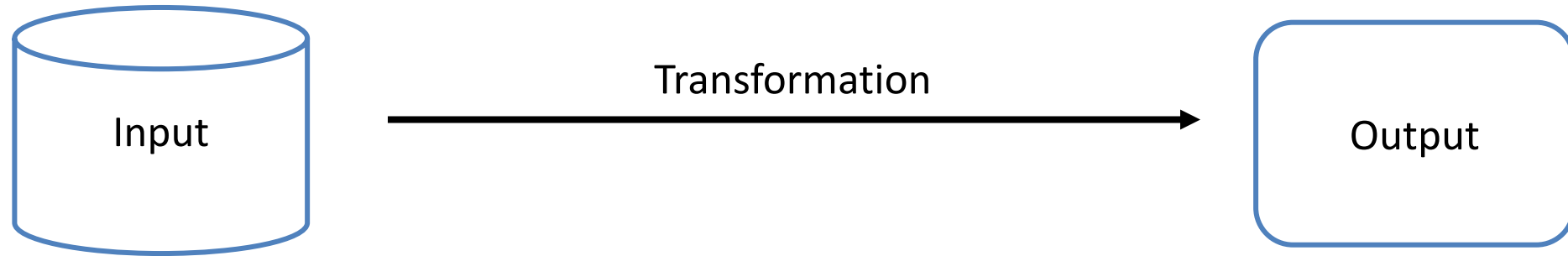
Neha Makhija

Northeastern University

(Joint work with Wolfgang Gatterbauer)

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

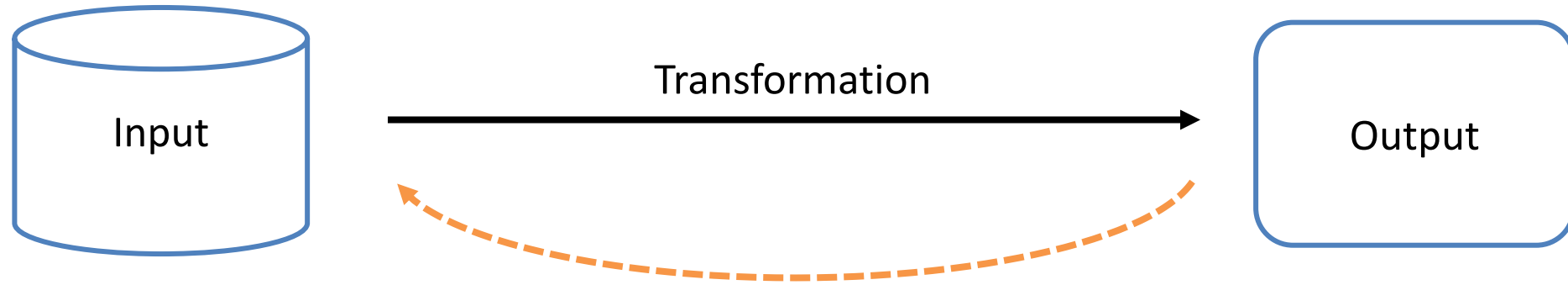
# Data Management



But sometimes, results are

- unexpected
- undesirable
- not understandable

# Reverse Data Management

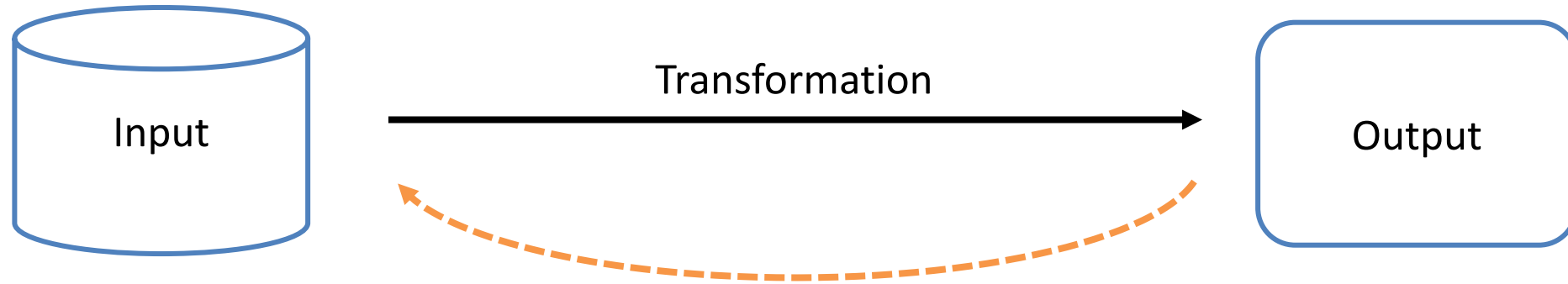


But sometimes, results are

- unexpected
- undesirable
- not understandable

Then we need to reason about the **input** in terms of the **output**!

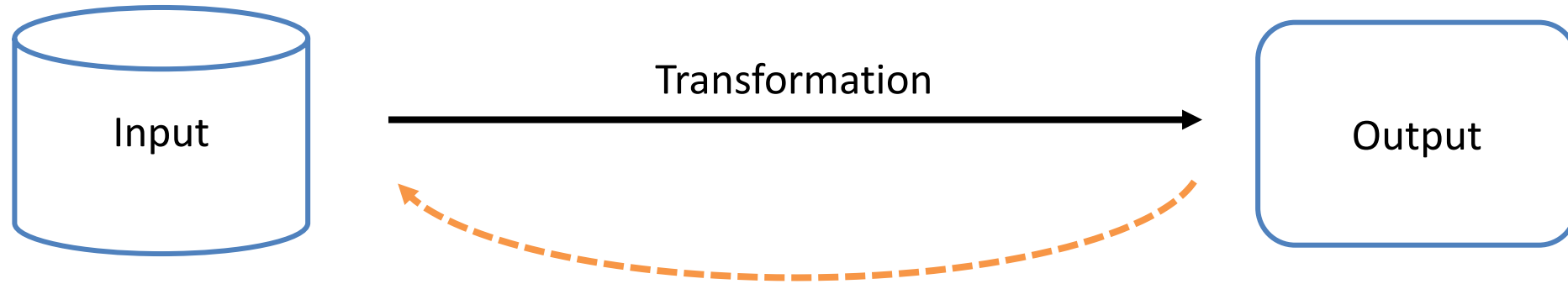
# Reverse Data Management



“What change would it take to achieve output X?”

“The handling of transformations that perform **actions on the input data**, **on behalf of desired outcomes** in the output data”

# Reverse Data Management: Example

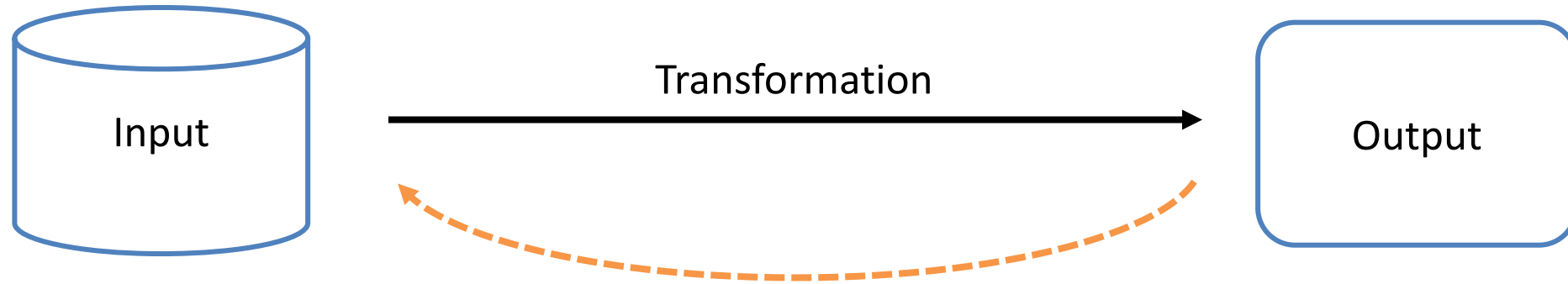


“What change would it take to **delete** some values from the output?”

## Deletion Propagation

- Minimize Source Side Effect
- Minimize View Side Effect

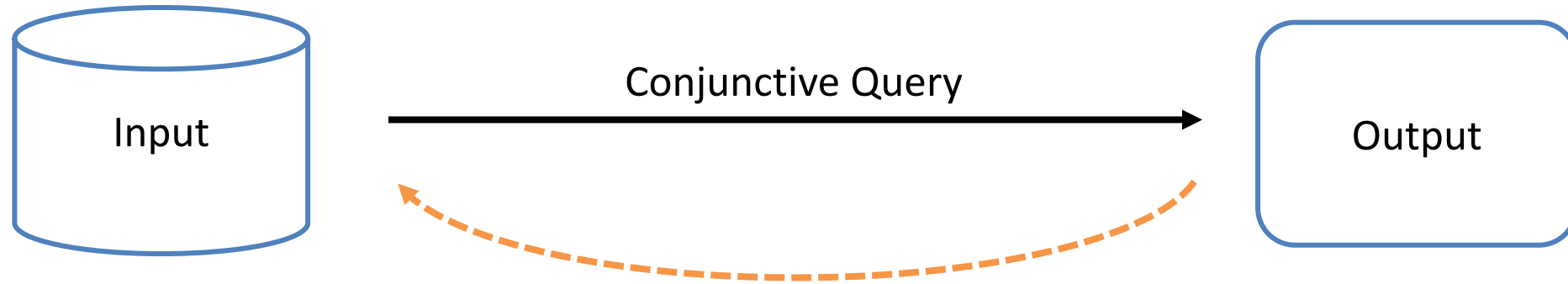
# Reverse Data Management: Example



“What minimum change would it take to ensure output is **“fair”**?”

Algorithmic Fairness

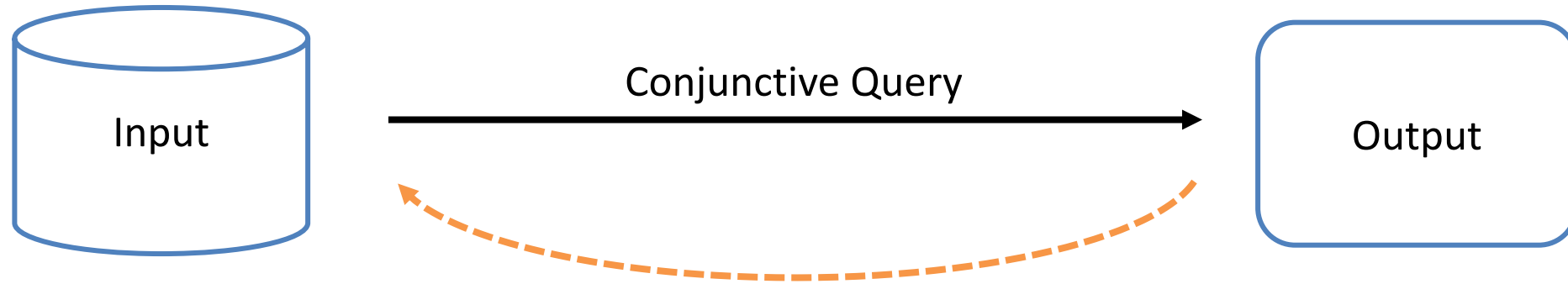
# Reverse Data Management: Example



“What change would it take to  
minimize the generated **provenance** expression?”

Minimal Factorization

# Reverse Data Management: Example



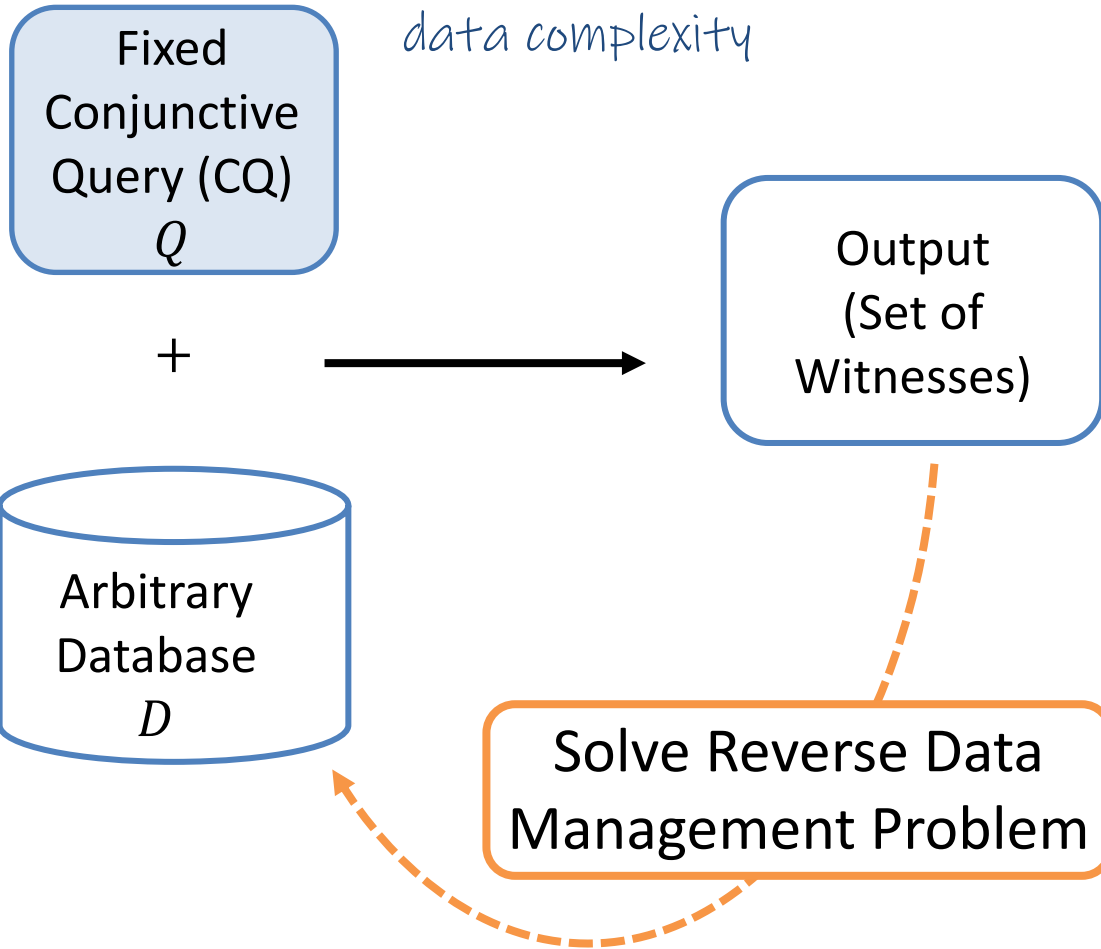
“What minimum change would it take to  
**make tuple X counterfactual?**”

Causal Responsibility



# Our Goal

*In other words,  
data complexity*



If hard

- Solve exactly
- Approximate in **PTIME**

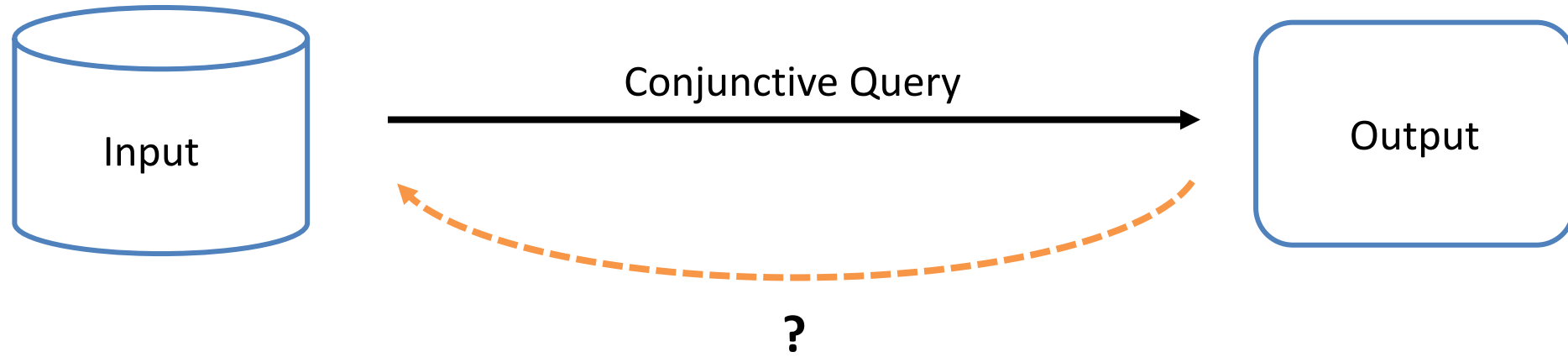
Find tractability border

Exploit **structure** in data

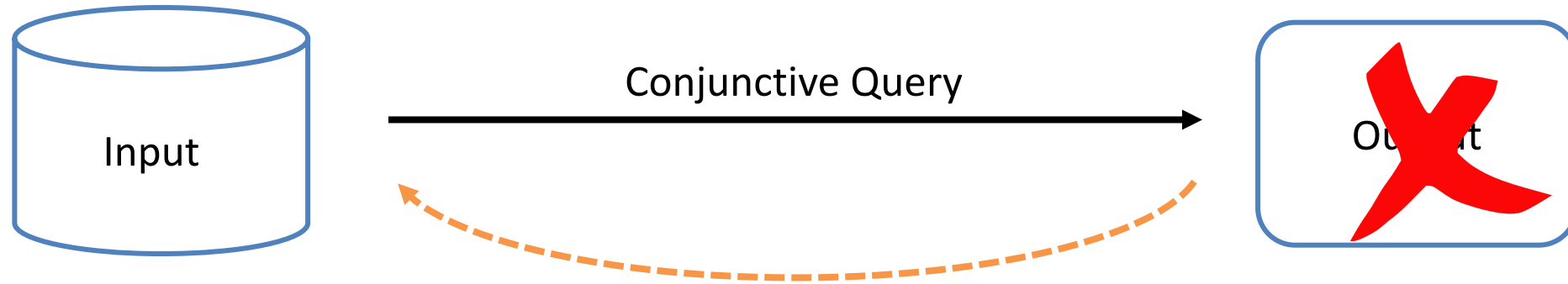
- Solve easy subclasses in PTIME

*One unified  
algorithm*

# Simplest Reverse Data Management Problem?



# Simplest Reverse Data Management Problem?



“What minimal change would it take to **delete the output?**”

Resilience

- Diagnose Points of Failure
- Equivalent to Deletion Propagation with Source Side-Effects

- Reverse Data Management Problems
- **Our Focus: Resilience**
- Results
- Our Unified Approach
- Takeaways + Open Questions

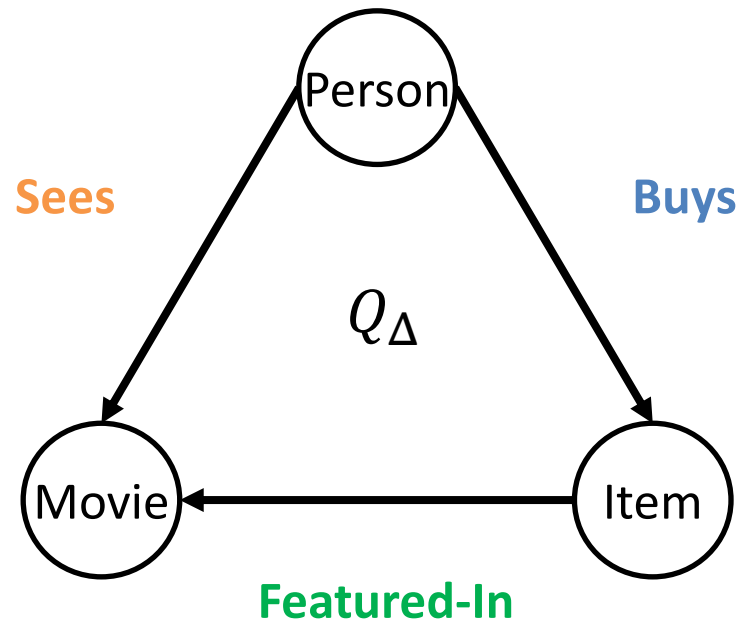
# Resilience: Example

**Sees**(person, movie) **Buys**(person, item) **Featured-In**(item, movie)

# Resilience: Example

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

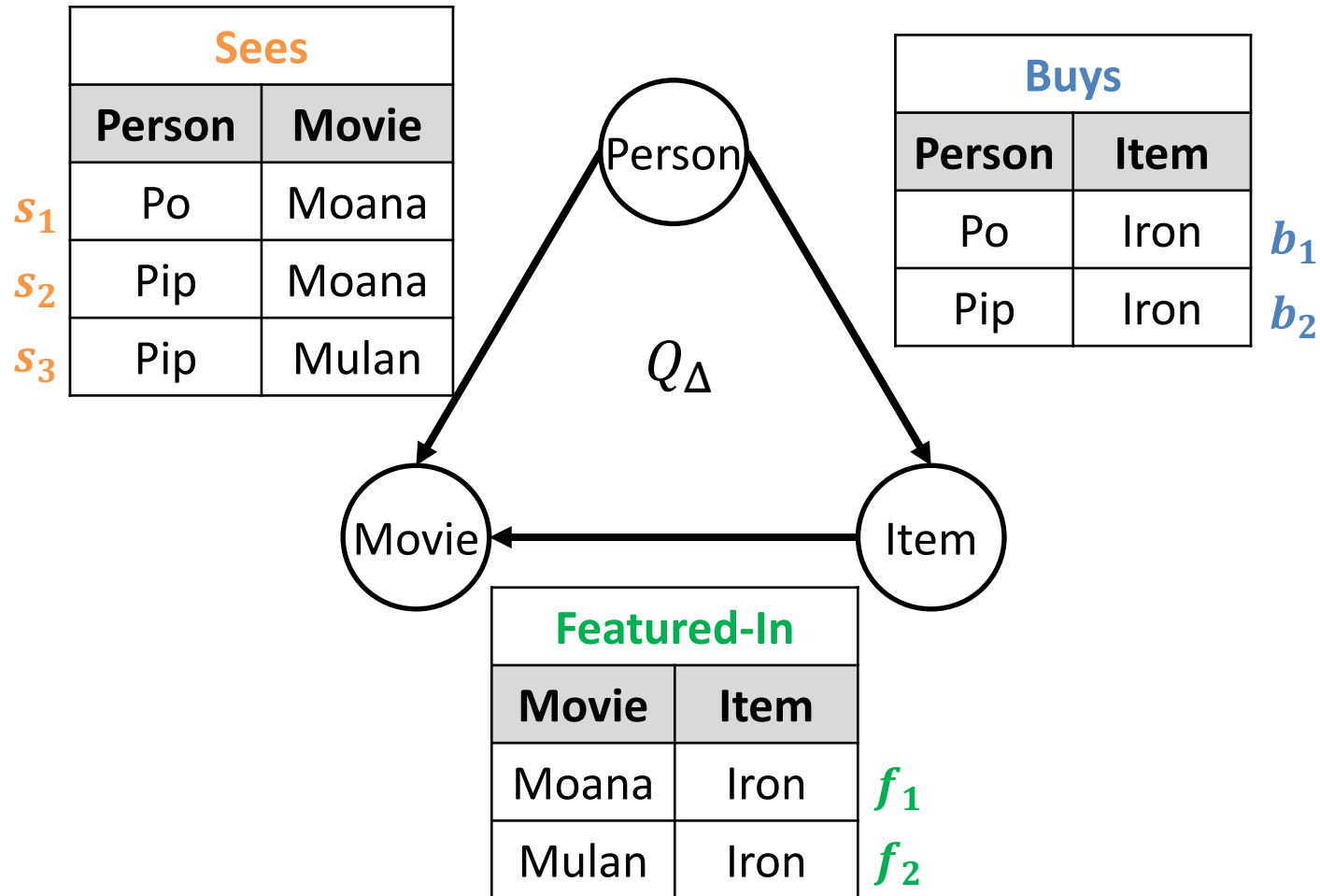
Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)



# Resilience: Example

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

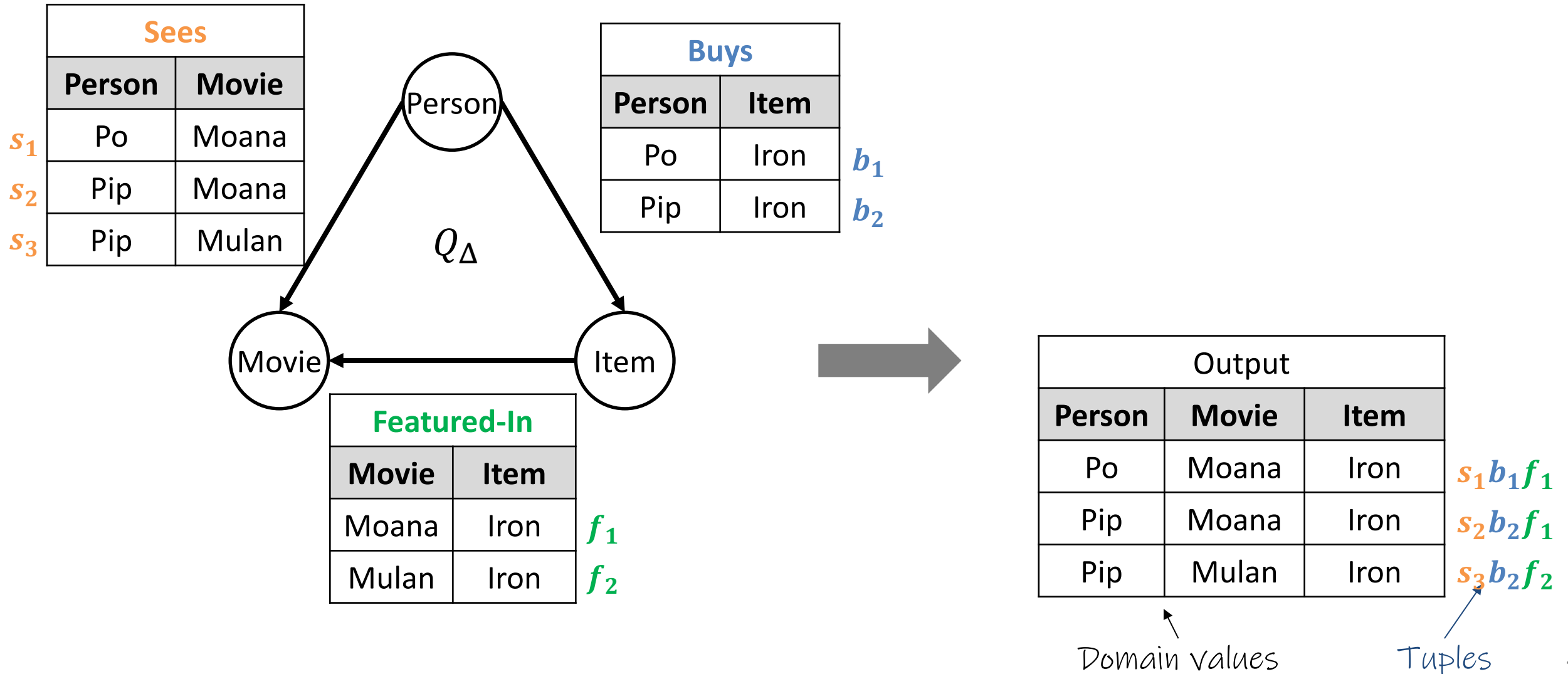
Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)



# Resilience: Example

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)

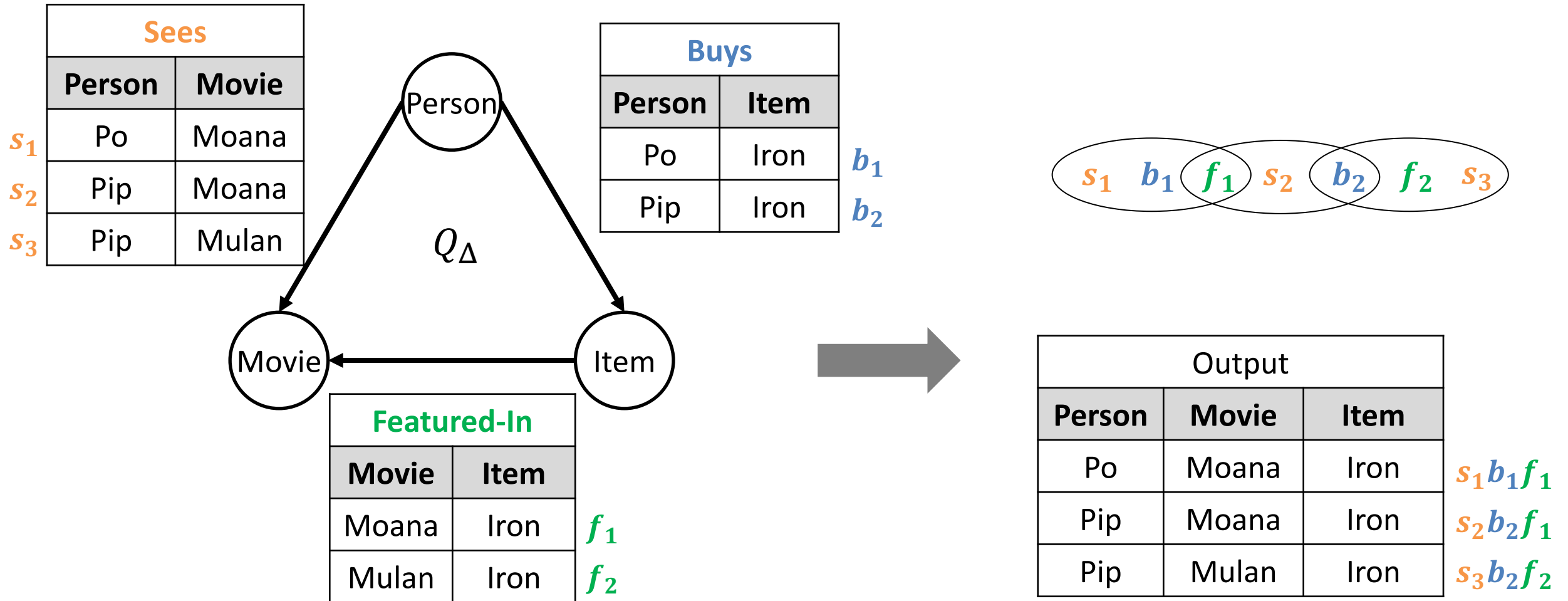




# Resilience

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

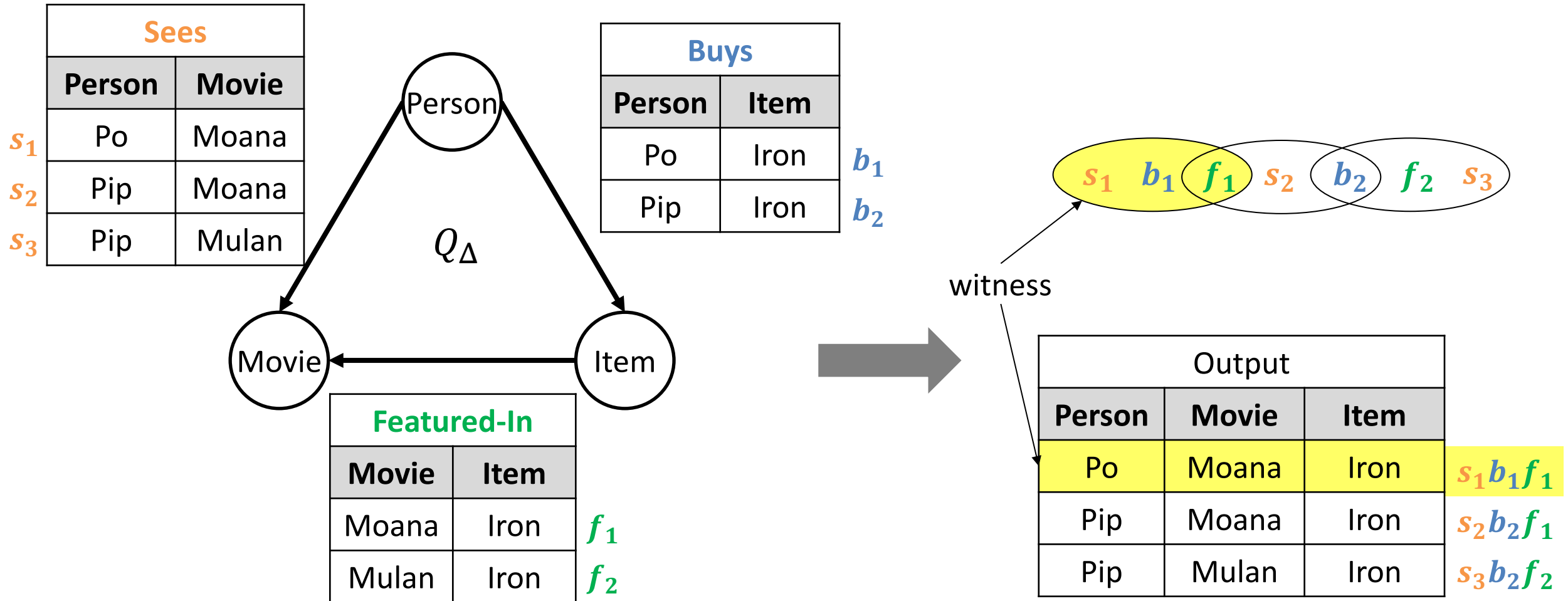
Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)



# Resilience

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

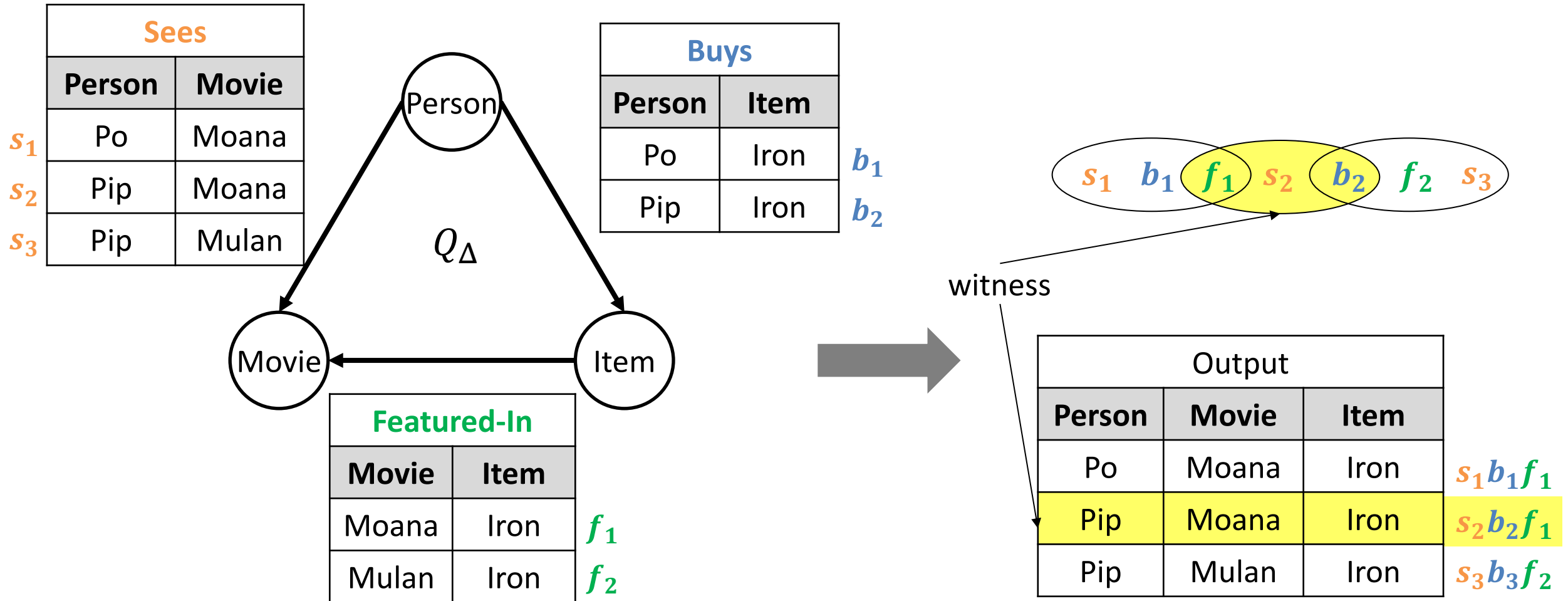
Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)



# Resilience

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

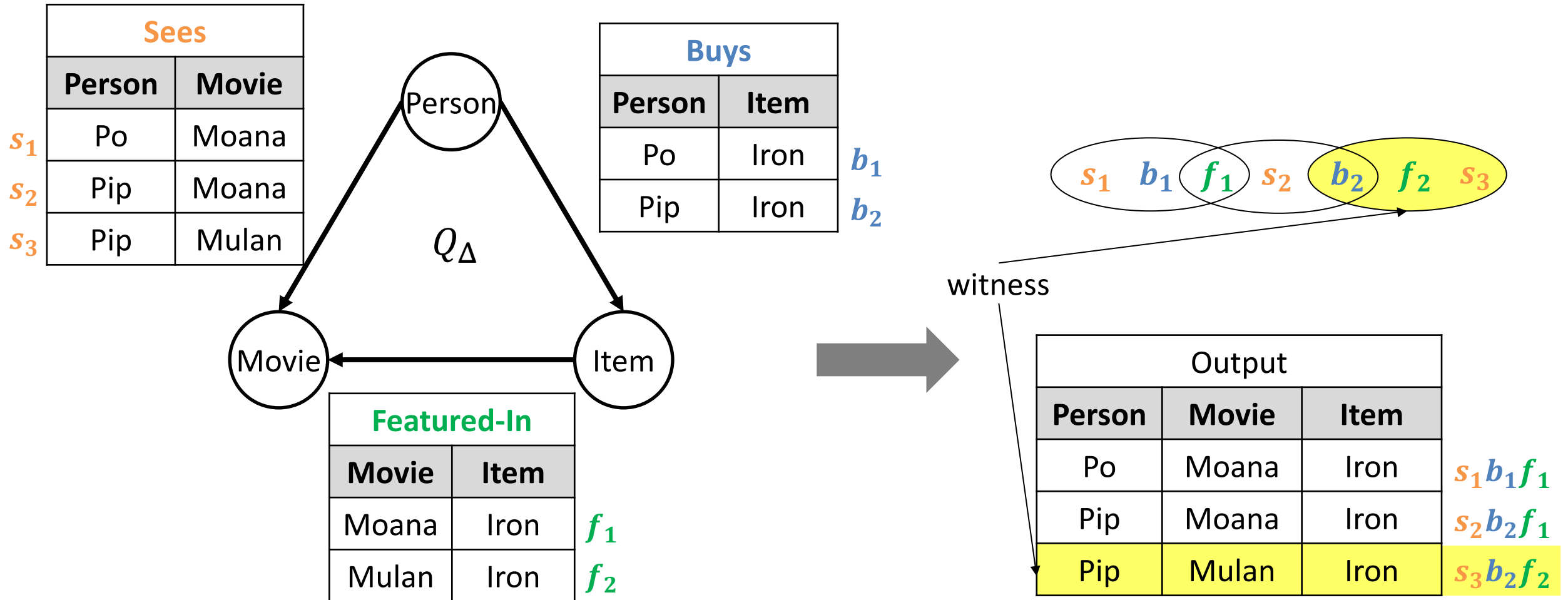
Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)



# Resilience

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

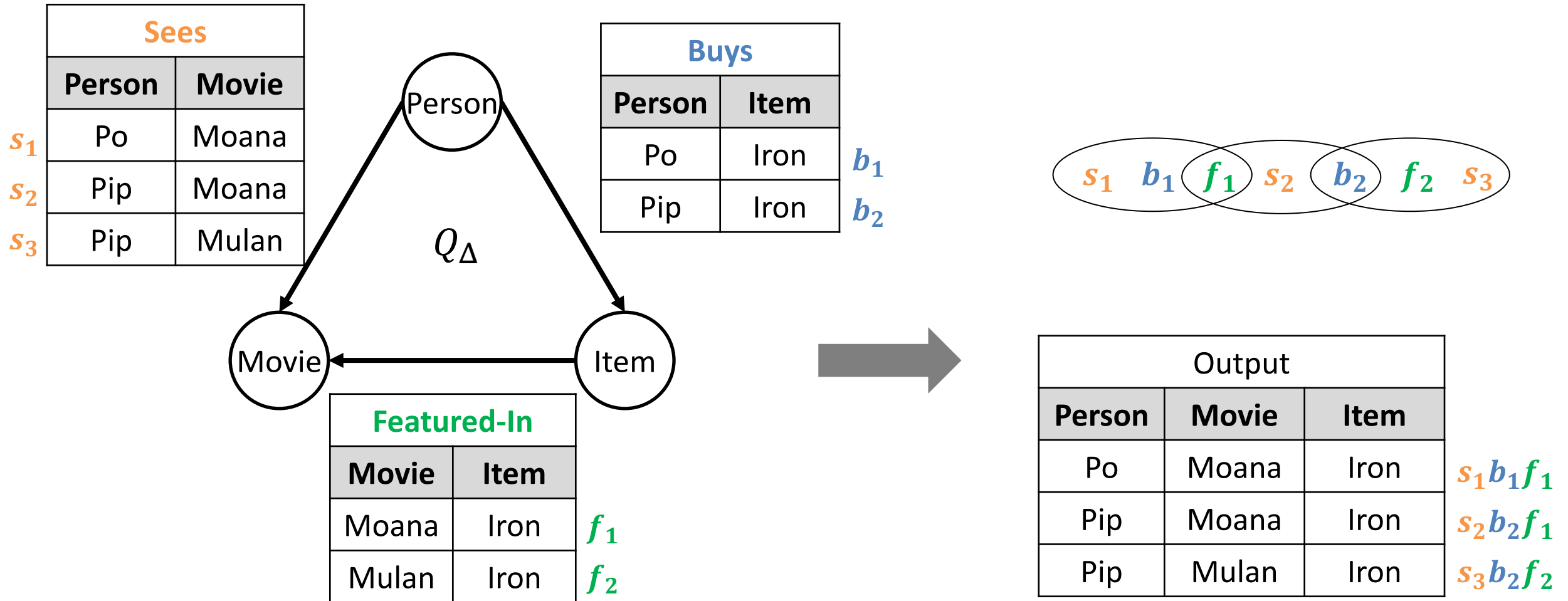
Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)



# Resilience

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)

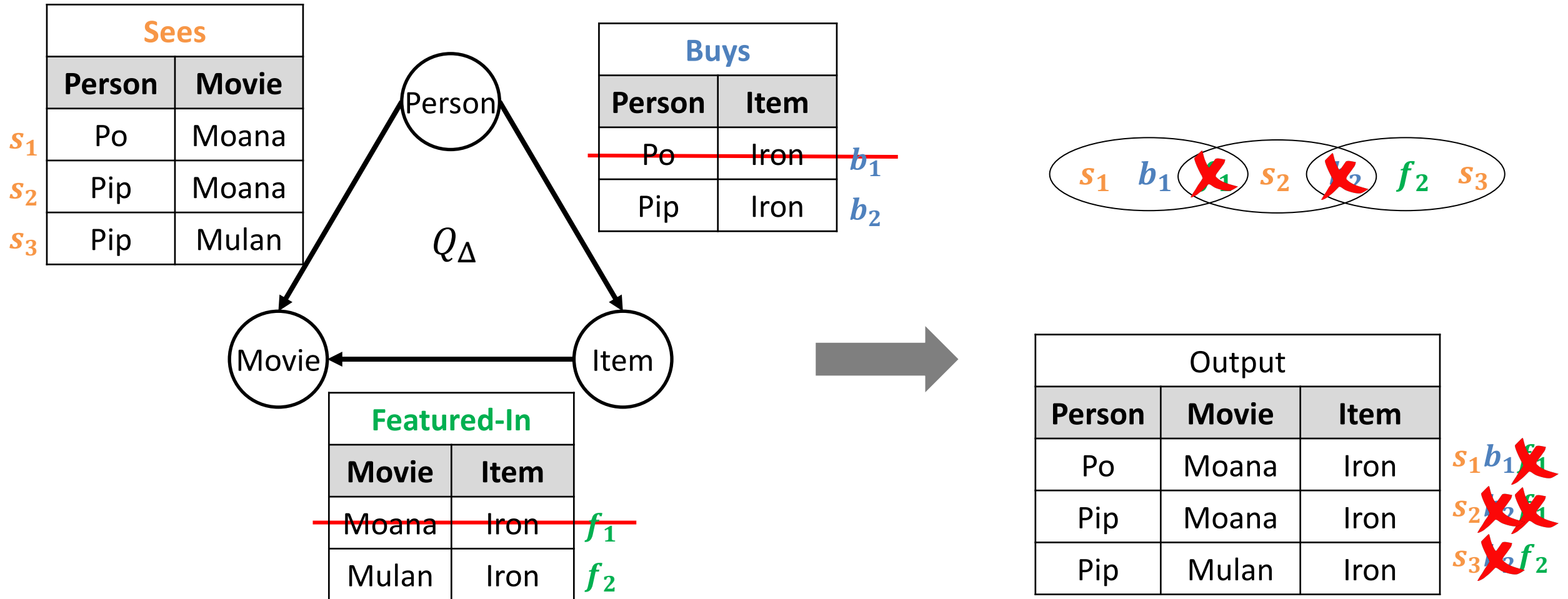


Recall: Resilience = What minimal change would it take to **delete the output?**

# Resilience

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)

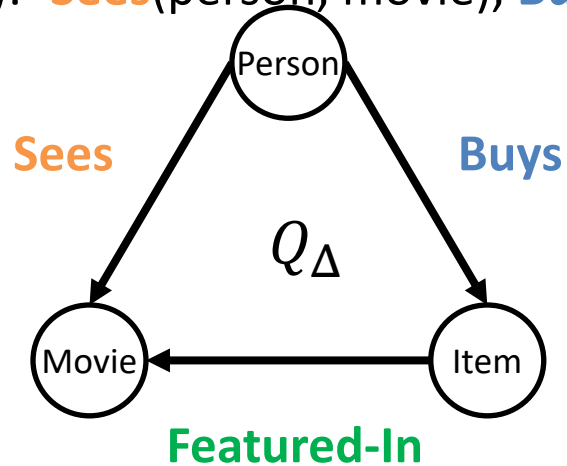


Recall: Resilience = What minimal change would it take to **delete the output?**

# Resilience Complexity

Query:- What person *sees* a movie and *buys* an item *featured in* the movie?

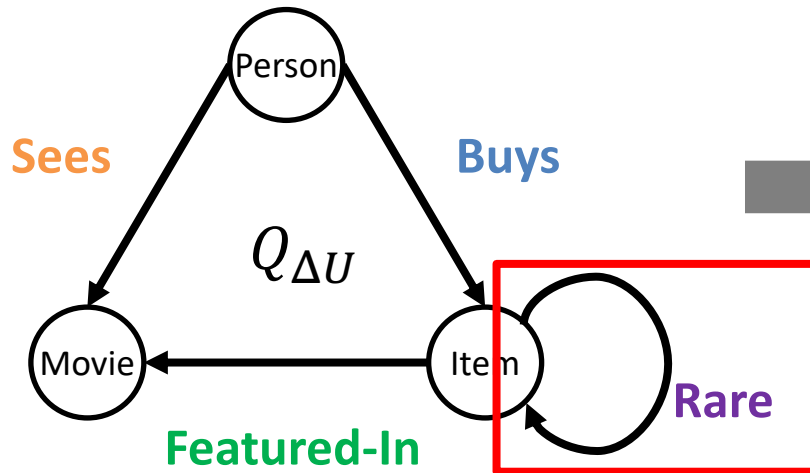
Q(person, movie, item):- *Sees*(person, movie), *Buys*(person, item), *Featured-In*(item, movie)



**NP-Complete**

Query:- What person *sees* a movie and *buys* an item *featured in* the movie that is *rare*?

Q(person, movie, item):- *Sees*(person, movie), *Buys*(person, item), *Featured-In*(item, movie), **Rare**(Item)

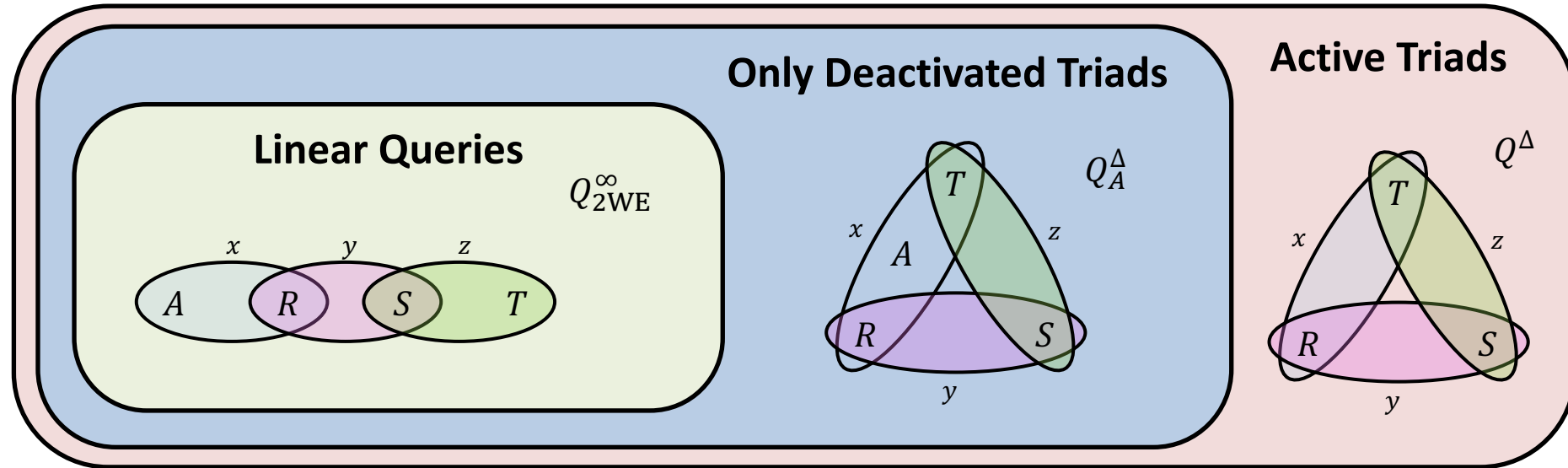


**PTIME**

- Reverse Data Management Problems
- Our Focus: Resilience
- **Results**
- Our Unified Approach
- Takeaways + Open Questions



# Resilience Complexity: Self-Join Free Queries



Resilience Set Semantics:

Project-Join Queries are <b>NP-Hard</b>		
<b>PTIME</b>	<b>PTIME</b>	<b>NPC</b>

Buneman+02

Freire+15

Resilience **Bag Semantics:**

<b>PTIME</b>	<b>NPC</b>	<b>NPC</b>
--------------	------------	------------

Makhija+24

## Complexity Results for Self-Join Free Queries

Buneman, Khanna, Tan. On propagation of deletions and annotations through views, PODS 2002 <https://doi.org/10.1145/543613.543633>

Freire, Gatterbauer, Immerman, Meliou. The complexity of resilience and responsibility for self-join-free conjunctive queries, PVLDB 2015 <https://dl.acm.org/doi/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://arxiv.org/abs/2212.08898>

# Unified Algorithm (Easy or Hard)

## Prior Work (RES+RESP)

Self-Join Free **PTIME** Cases

Flow based encodings  
specific to query  
[Freire+15]

Self-Join Free **NPC** Cases:  
Exact Evaluation

No prior algorithm:  
brute force trivially

Self-Join Free **NPC** Cases:  
Approximations

No prior algorithm

Queries with Self-Joins

Flow algorithms for some  
known PTIME queries  
[Freire+20]

## Our Approach (RES+RESP)

### One **Unified** Algorithm

- Solves all cases
- All known PTIME cases terminate in PTIME

# Unified Algorithm (Across Different Settings)

## Prior Work

## Our Work

Semantics

Set Semantics

Set + **Bag** Semantics

Queries

Self-Join Free Queries

+

Some Restricted SJ Queries

**All** UCQs

Functional  
Dependencies

Can leverage known FDs

Take advantage of  
**unspecified FDs** in data

- Reverse Data Management Problems
- Our Focus: Resilience
- Results
- **Our Unified Approach**
  - Unified Algorithm
  - Unified Hardness Criterion
- Takeaways + Open Questions

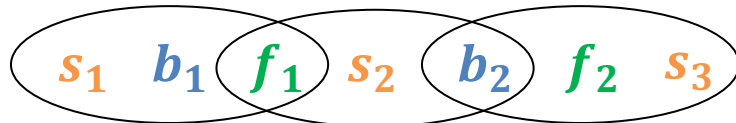
# Resilience as an Integer Linear Program (ILP)

Weights if applied can go in the objective

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x} \\ \text{constraint matrix} \quad & \mathbf{A} \mathbf{x} \geq \mathbf{b} \\ \text{constraint vector} \quad & \mathbf{x} \in \{0,1\}^n \end{aligned}$$

$$\min \left[ \begin{aligned} & x[s_1] + x[b_1] + x[f_1] + x[s_2] \\ & + x[b_2] + x[f_2] + x[s_3] \end{aligned} \right]$$

$Q_{\Delta}$  Example:



Output		
Person	Movie	Item
Po	Moana	Iron
Pip	Moana	Iron
Pip	Mulan	Iron

$s_1 b_1 f_1$

$s_2 b_2 f_1$

$s_3 b_2 f_2$

Queries with self-joins are modelled the same way. Not all constraints must have the same arity.

$$\begin{aligned} x[s_1] + x[b_1] + x[f_1] &\geq 1 \\ x[f_1] + x[s_2] + x[b_2] &\geq 1 \\ x[b_2] + x[f_2] + x[s_3] &\geq 1 \\ \forall \text{ tuples } \mathbf{t}: \quad &\mathbf{x}[\mathbf{t}] \in \{0,1\} \end{aligned}$$

Setting  $X[f_1] = X[b_2] = 1$ , and all other variables to 0 satisfies the ILP.

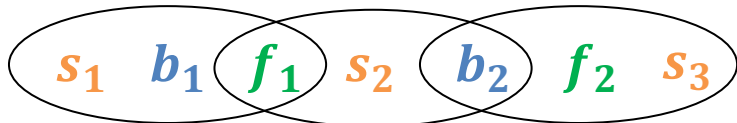
Optimal value = 2

# Resilience as a Linear Program (LP)

$$\begin{array}{l}
 \min \mathbf{c}^T \mathbf{x} \\
 \mathbf{A} \mathbf{x} \geq \mathbf{b} \\
 \mathbf{x} \in \{0,1\}^n \\
 \mathbf{x} = [0,1]^n
 \end{array}$$

*objective vector* →  $\mathbf{c}^T$   
*constraint matrix* →  $\mathbf{A}$   
*constraint vector* →  $\mathbf{b}$

$Q_\Delta$  Example:



Output			
Po	Moana	Iron	$s_1 b_1 f_1$
Pip	Moana	Iron	$s_2 b_2 f_1$
Pip	Mulan	Iron	$s_3 b_2 f_2$

$$\min \left[ \begin{array}{l} x[s_1] + x[b_1] + x[f_1] + x[s_2] \\ + x[b_2] + x[f_2] + x[s_3] \end{array} \right]$$

$$x[s_1] + x[b_1] + x[f_1] \geq 1$$

$$x[f_1] + x[s_2] + x[b_2] \geq 1$$

$$x[b_2] + x[f_2] + x[s_3] \geq 1$$

~~$$\forall \text{ tuples } \mathbf{t}: \mathbf{x}[\mathbf{t}] \in \{0,1\}$$~~

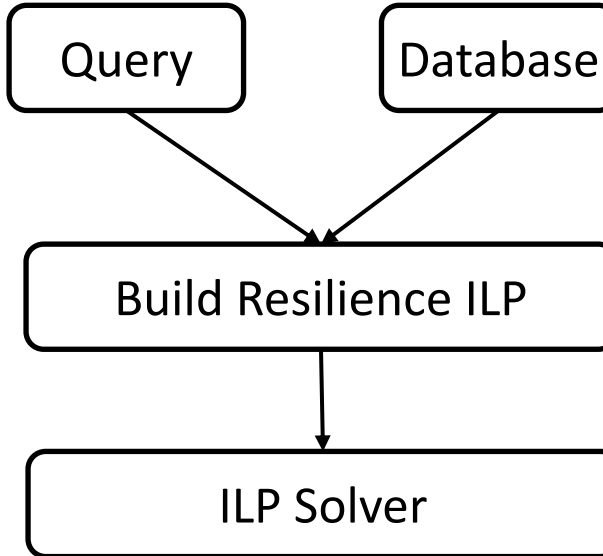
$$\forall \text{ tuples } \mathbf{t}: 0 \leq \mathbf{x}[\mathbf{t}] \leq 1$$

Linear programs allow fractional solutions  
 We can get a lower bound in PTIME.

# Unified Algorithm for Resilience

Theorem.

For all known PTIME cases of resilience, the **objective** of the LP relaxation is identical to the ILP.



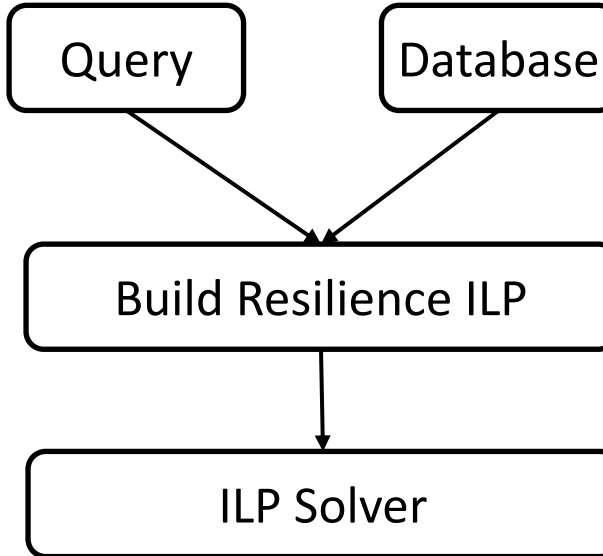
*But ILPs are NP-Hard!*

# Unified Algorithm for Resilience



Theorem.

For all known PTIME cases of Resilience, **LP=ILP**



*But ILPs are NP-Hard!*



# Unified Algorithm for Resilience



Theorem.

For all known PTIME cases of Resilience, **LP=ILP**

Query

Database

Build Resilience ILP

ILP Solver

LP Pre-solve  
(Lower Bound)

Fractional Solution

Heuristics based  
Branch and Bound

=LP?

Integral Solution

Exhaust search space

Integral Solution



If **LP=ILP**, solvers can recover a solution efficiently!

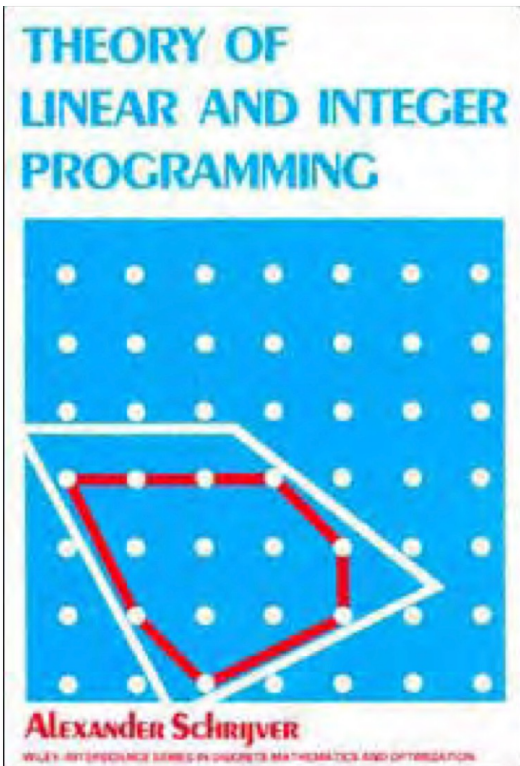
# When does LP = ILP?

Theorem.

For all known PTIME cases of Resilience, **LP=ILP**

*Our PTIME constraint matrixes need not be balanced or Totally Unimodular*

*The PTIME cases go beyond these known criterion!*



Alexander Schrijver

## Combinatorial Optimization

Polyhedra and Efficiency

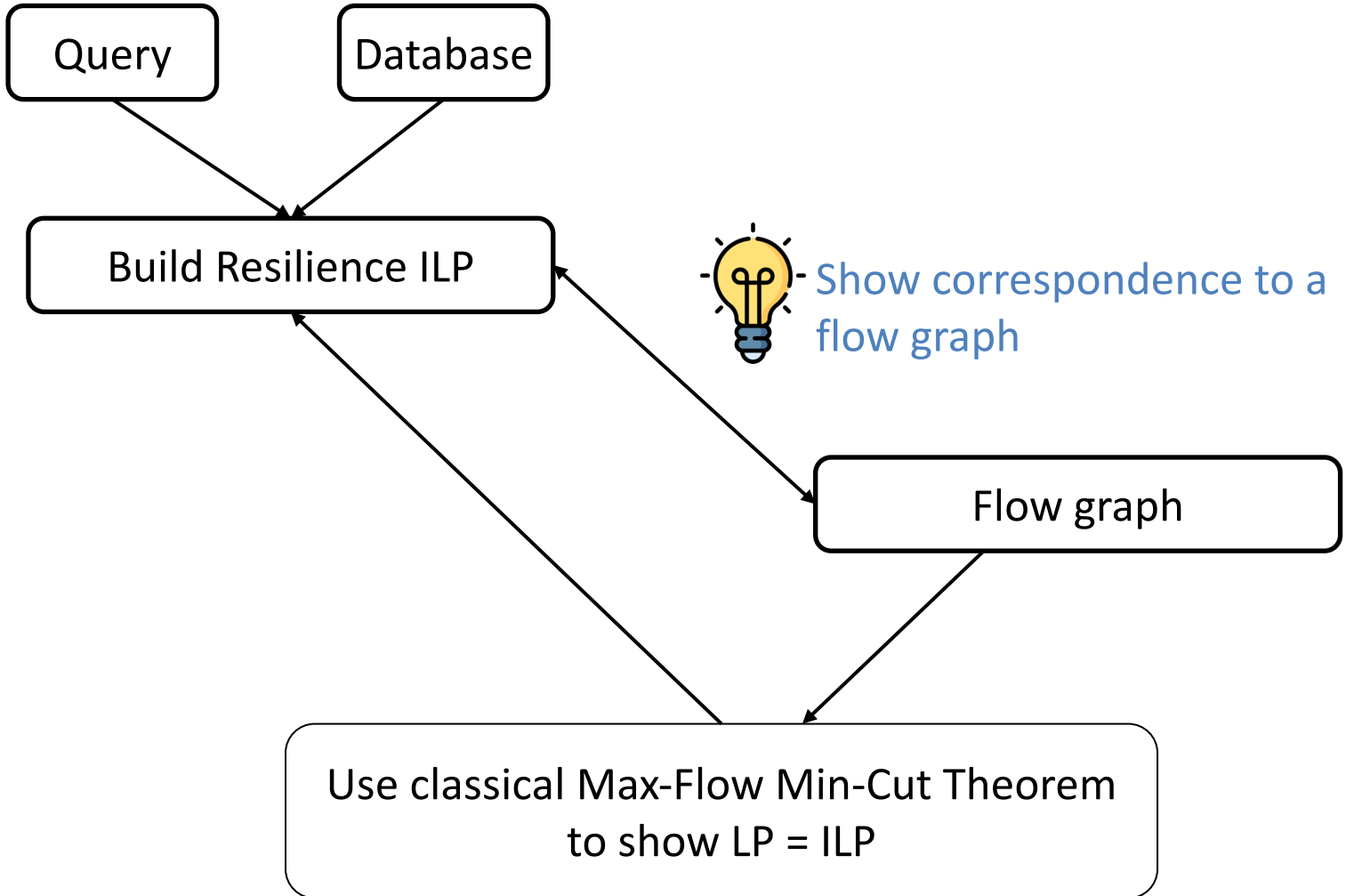
Volume A-C

<b>83</b>	<b>Balanced and unimodular hypergraphs</b>	1439
83.1	Balanced hypergraphs	1439
83.2	Characterization 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

# When does LP = ILP?

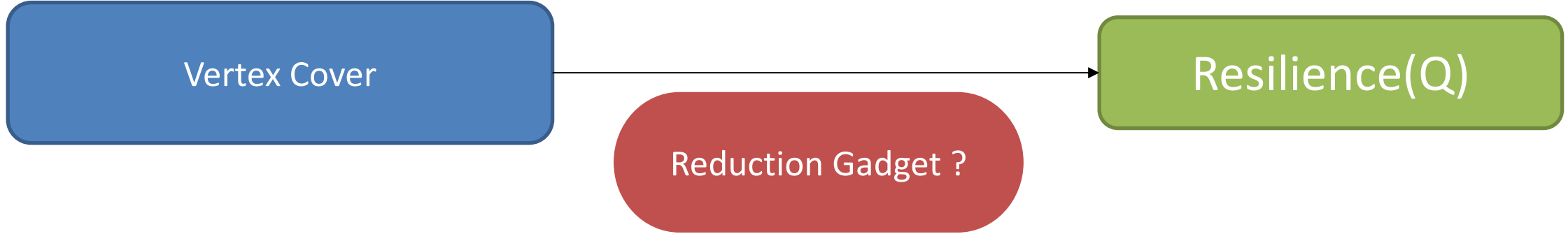
Theorem.

For all known PTIME cases of Resilience, **LP=ILP**

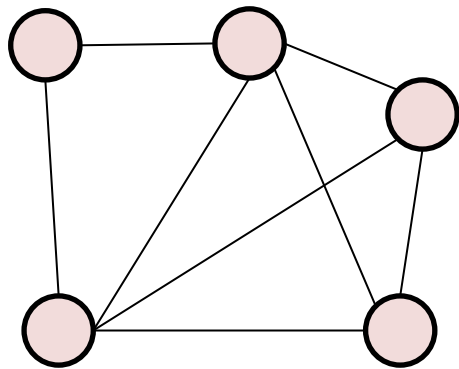


- Reverse Data Management Problems
- Our Focus: Resilience
- Results
- **Our Unified Approach**
  - Unified Algorithm
  - **Unified Hardness Criterion**
- Takeaways + Open Questions

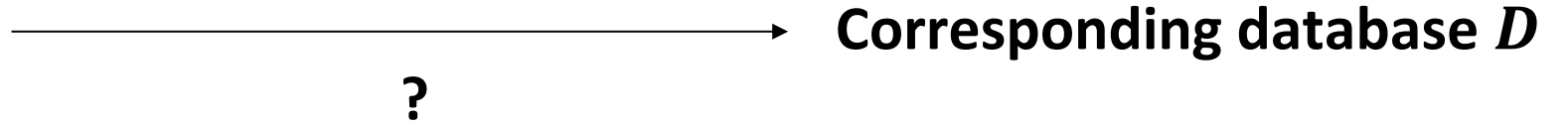
# Unified Approach to prove Hardness



Arbitrary graph  $G$



$$VC(G) = k$$

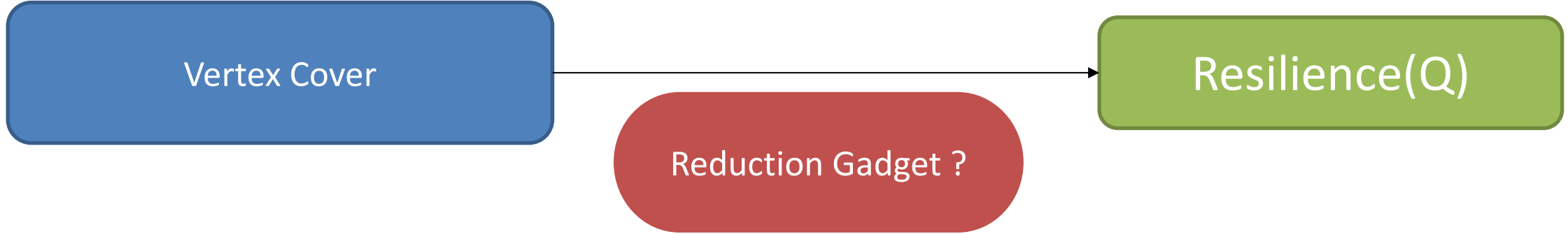


Corresponding database  $D$

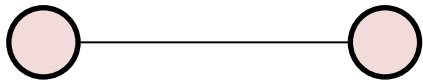
Sees		Featured-In		Buys	
$s_1$	...	$f_1$	...	$b_1$	...
$s_2$	...	$f_2$	...	$b_2$	...
$s_3$	...	$f_3$	...	$b_3$	...

$$\text{Resilience}(Q, D) = f(k)$$

# Unified Approach to prove Hardness



Single Edge



$$VC(G) = k$$

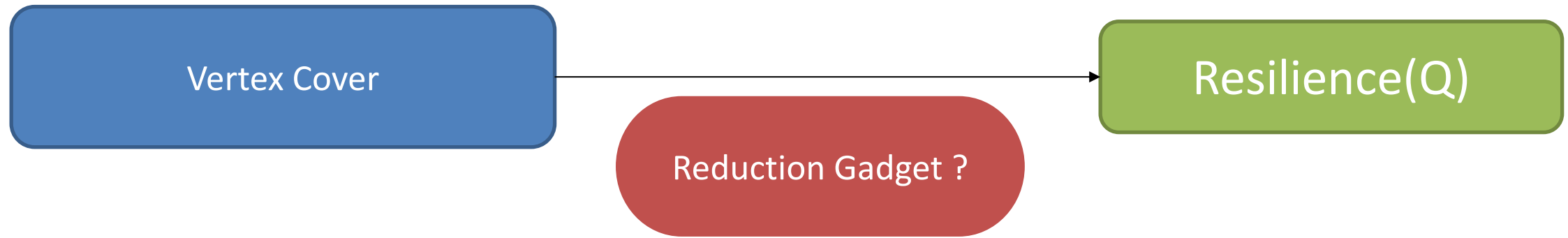
Edge Gadget

Corresponding database  $D$

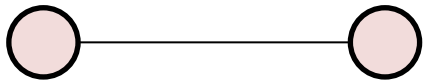
	Sees		Featured-In		Buys			
$s_1$	...	...	$f_1$	...	...	$b_1$	...	...
$s_2$	...	...	$f_2$	...	...	$b_2$	...	...
$s_3$	...	...	$f_3$	...	...	$b_3$	...	...

$$\text{Resilience}(Q, D) = f(k)$$

# Unified Approach to prove Hardness



Single Edge



$$VC(G) = k$$

Edge Gadget

Independent Join Path

**We prove:** If a query forms an “Independent Join Path” → it is NPC (conjecture Freire+20)

Corresponding database  $D$



Data Hypergraph

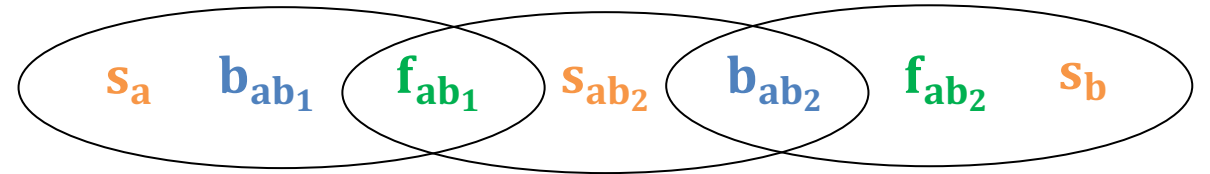
$$\text{Resilience}(Q, D) = f(k)$$

# Unified Approach to prove Hardness

## What is an Independent Join Path?

Database under query  $Q$ , with endpoints, with 5 *testable* properties:

1. Data hypergraph is connected



Data Hypergraph

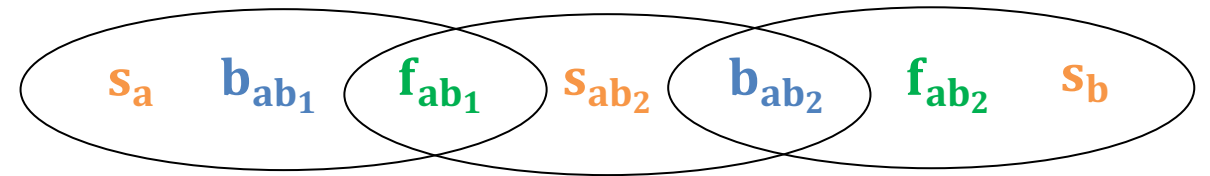


# Unified Approach to prove Hardness

## What is an Independent Join Path?

Database under query  $Q$ , with endpoints, with 5 *testable* properties:

1. Data hypergraph is connected
2. Database is reduced



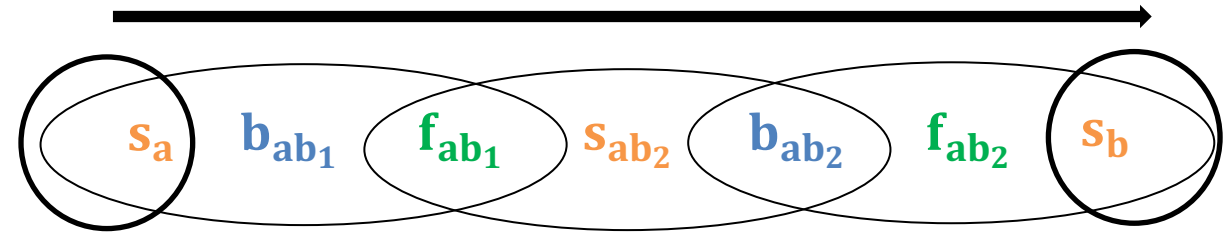
Data Hypergraph

# Unified Approach to prove Hardness

## What is an Independent Join Path?

Database under query  $Q$ , with endpoints, with 5 *testable* properties:

1. Data hypergraph is connected
2. Database is reduced
3. Endpoints are “valid”



Data Hypergraph

# Unified Approach to prove Hardness

## What is an Independent Join Path?

Database under query  $Q$ , with endpoints, with 5 *testable* properties:

1. Data hypergraph is connected

2. Database is reduced

3. Endpoints are “valid”

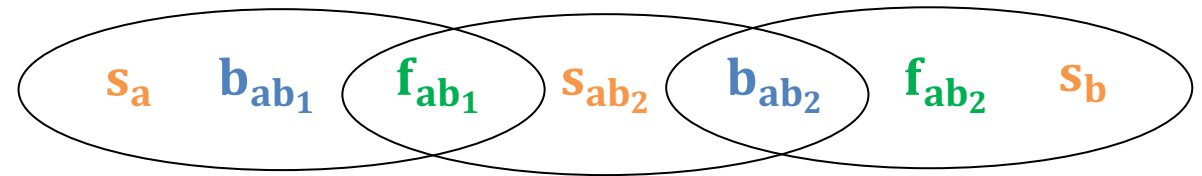
4. OR – property

5. Composability

“Key” properties:

- Semantically defined

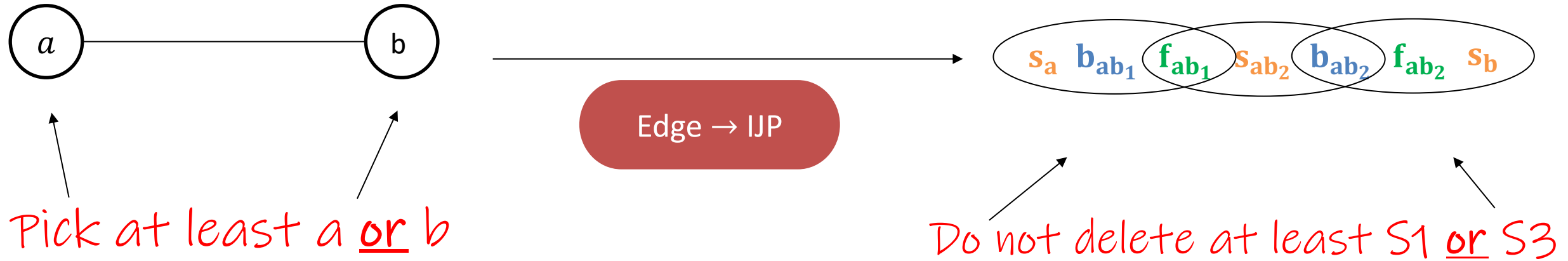
- I will just show intuition



Data Hypergraph

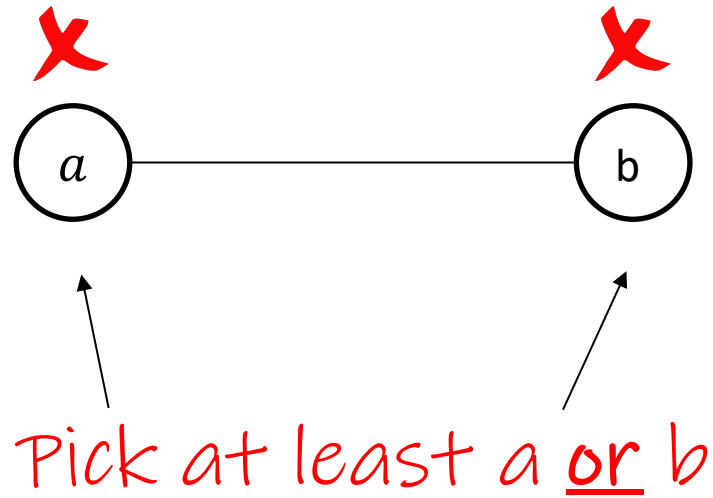
# Unified Approach to prove Hardness

## Key Property #1: OR property



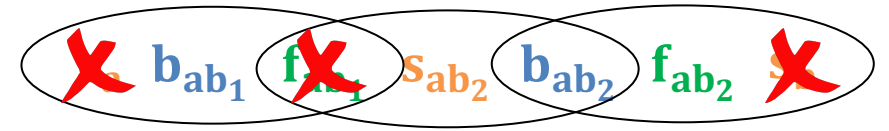
# Unified Approach to prove Hardness

## Key Property #1: OR property



VC = ??

Violates constraints



Do not delete at least  $S_1$  or  $S_3$

RES = 3

Violates minimality

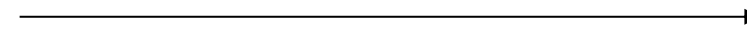
# Unified Approach to prove Hardness

## Key Property #1: OR property

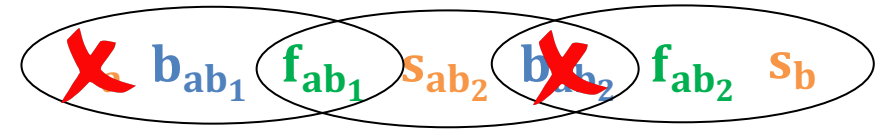


Pick at least a or b

VC = 1



Edge  $\rightarrow$  IJP



Do not delete at least  $S1$  or  $S3$

RES = 2

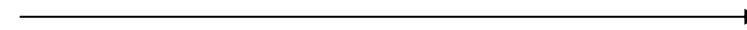
# Unified Approach to prove Hardness

## Key Property #1: OR property

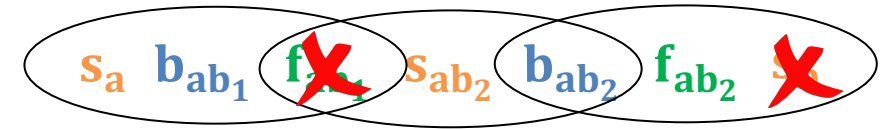


Pick at least a or b

VC = 1



Edge  $\rightarrow$  IJP



Do not delete at least  $S1$  or  $S3$

RES = 2

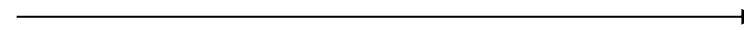
# Unified Approach to prove Hardness

## Key Property #1: OR property

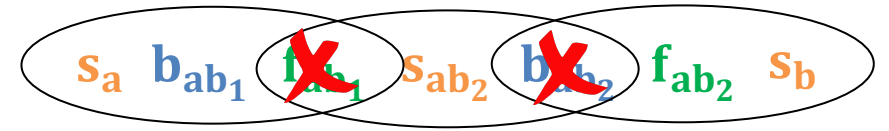


Pick at least a or b

VC = 2



Edge  $\rightarrow$  IJP



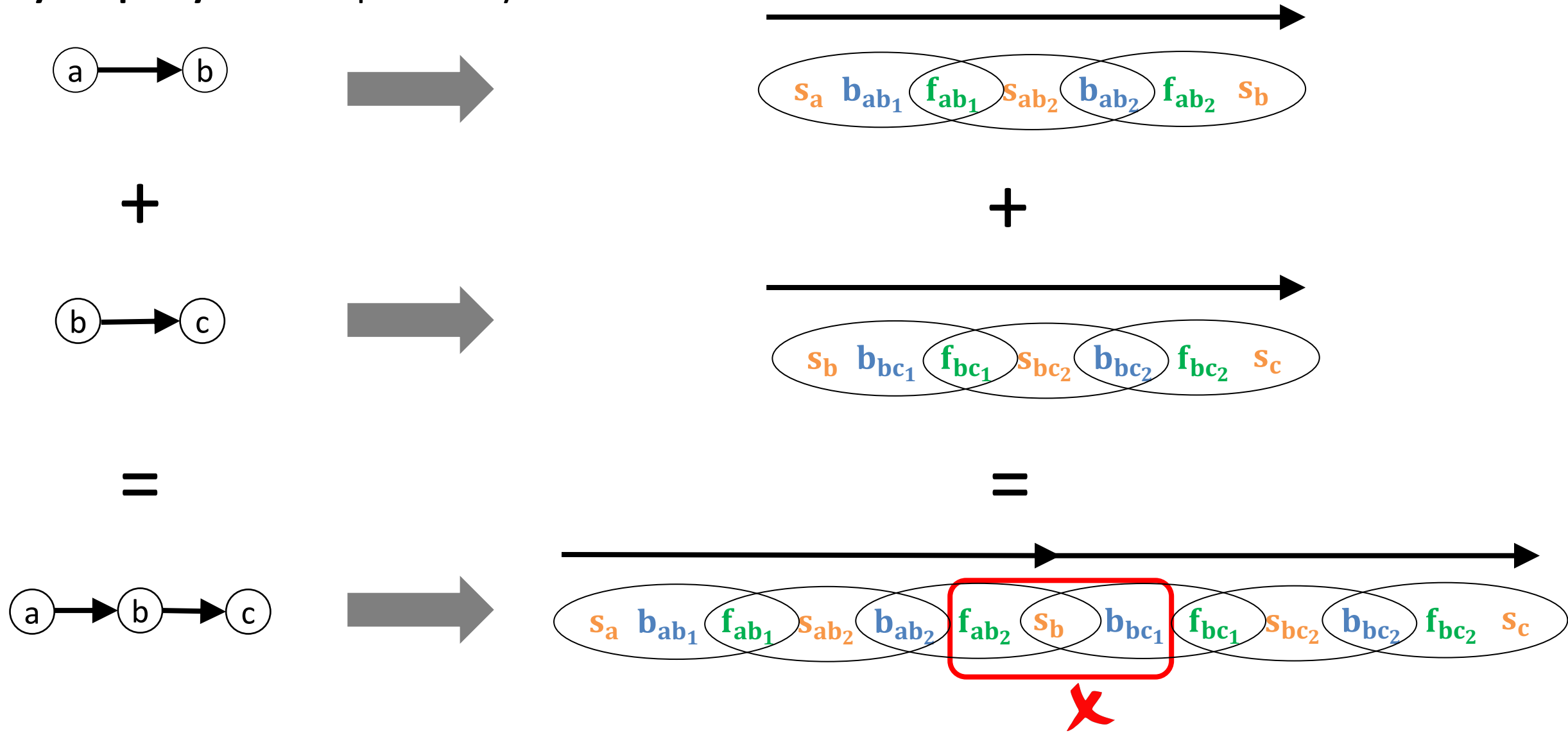
Do not delete at least  $s_1$  or  $s_3$

RES = 2



# Unified Approach to prove Hardness

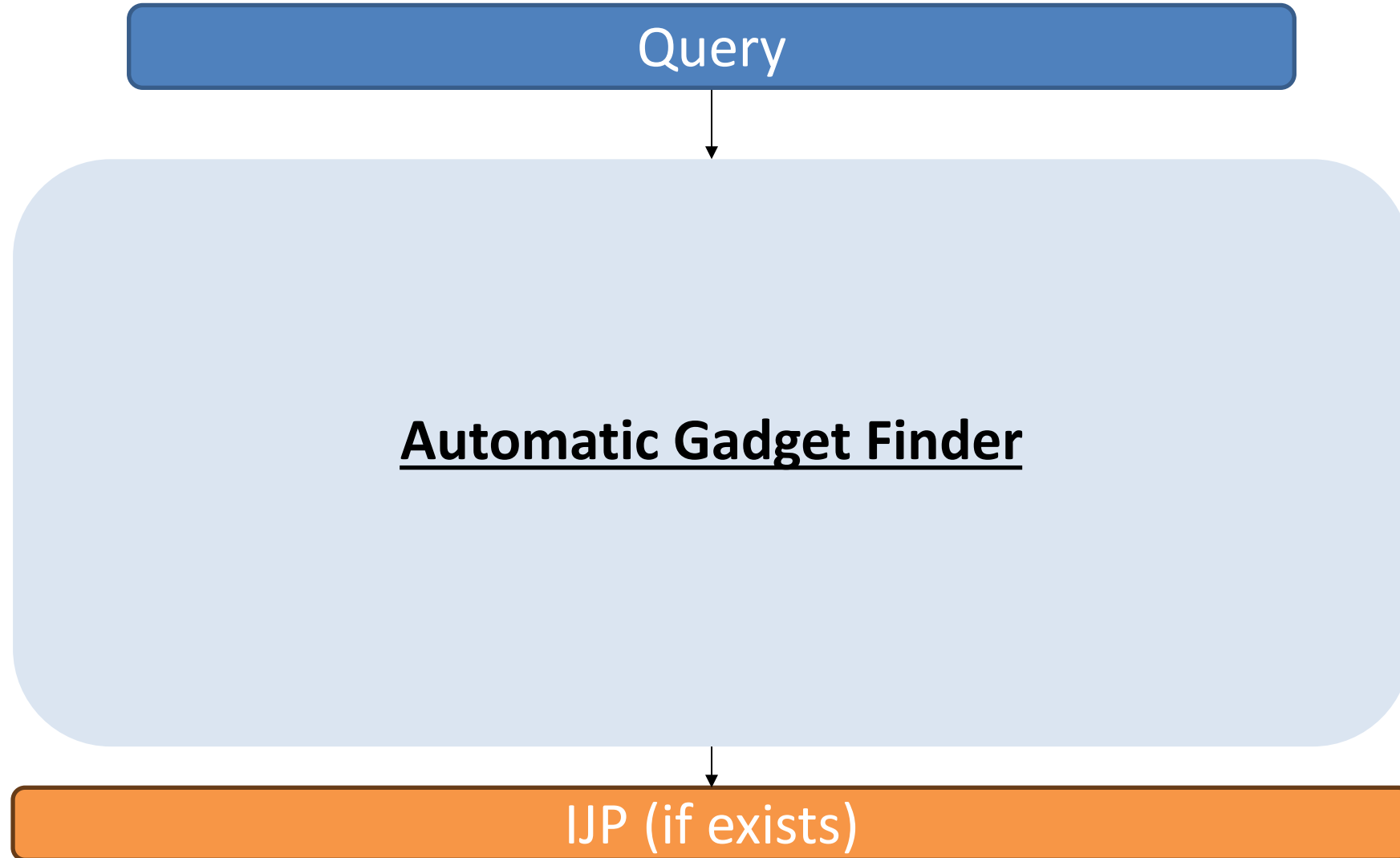
## Key Property #2: Composability of IJPs



Composing two IJPs should not lead to additional witnesses

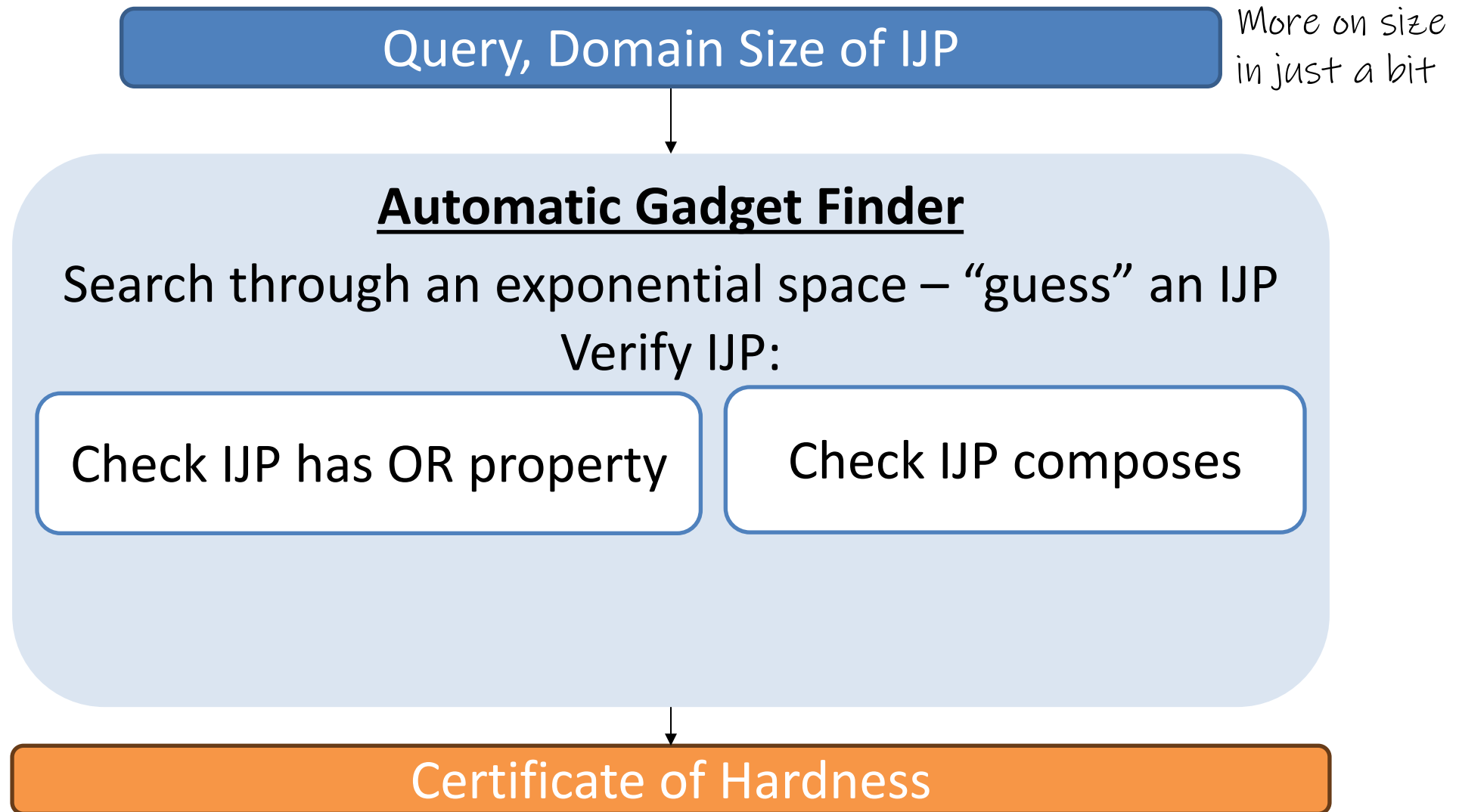
# Our Goal

Using the complete criterion for IJP, can we build a **principled** way to find IJPs?



# Unified Approach to prove Hardness

Using the complete criterion for IJP, we can build a **principled** way to find IJPs



# Unified Approach to prove Hardness

Using the complete criterion for IJP, we can build a **principled** way to find IJPs

Query, Domain Size of IJP

More on size  
in just a bit

## Automatic Gadget Finder

Search through an exponential space – “guess” an IJP

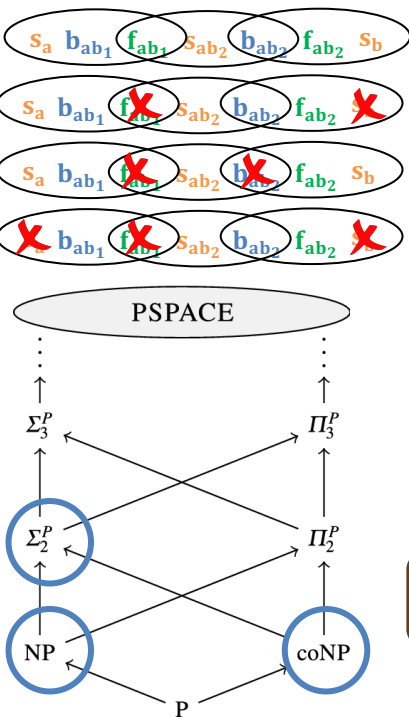
Verify IJP:

Check IJP has OR property

- Find size  $k$  resilience set (**NP**)
- Find no size  $k-1$  resilience set (**Co-NP**)

Check IJP composes

Certificate of Hardness



# Unified Approach to prove Hardness

Using the complete criterion for IJP, we can build a **principled** way to find IJPs

Query, Domain Size of IJP

More on size  
in just a bit

## Automatic Gadget Finder


Search through an exponential space – “guess” an IJP

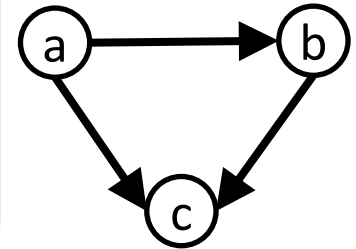
Verify IJP:

Check IJP has OR property

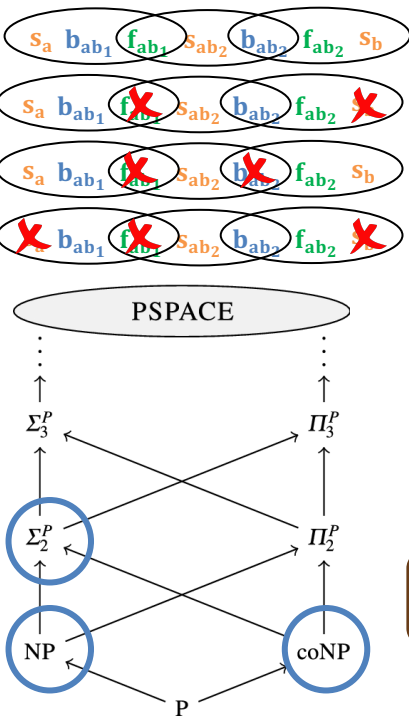
- Find size  $k$  resilience set (**NP**)
- Find no size  $k-1$  resilience set (**Co-NP**)

Check IJP composes

Theorem.   
It suffices to check just  
3 join paths compose

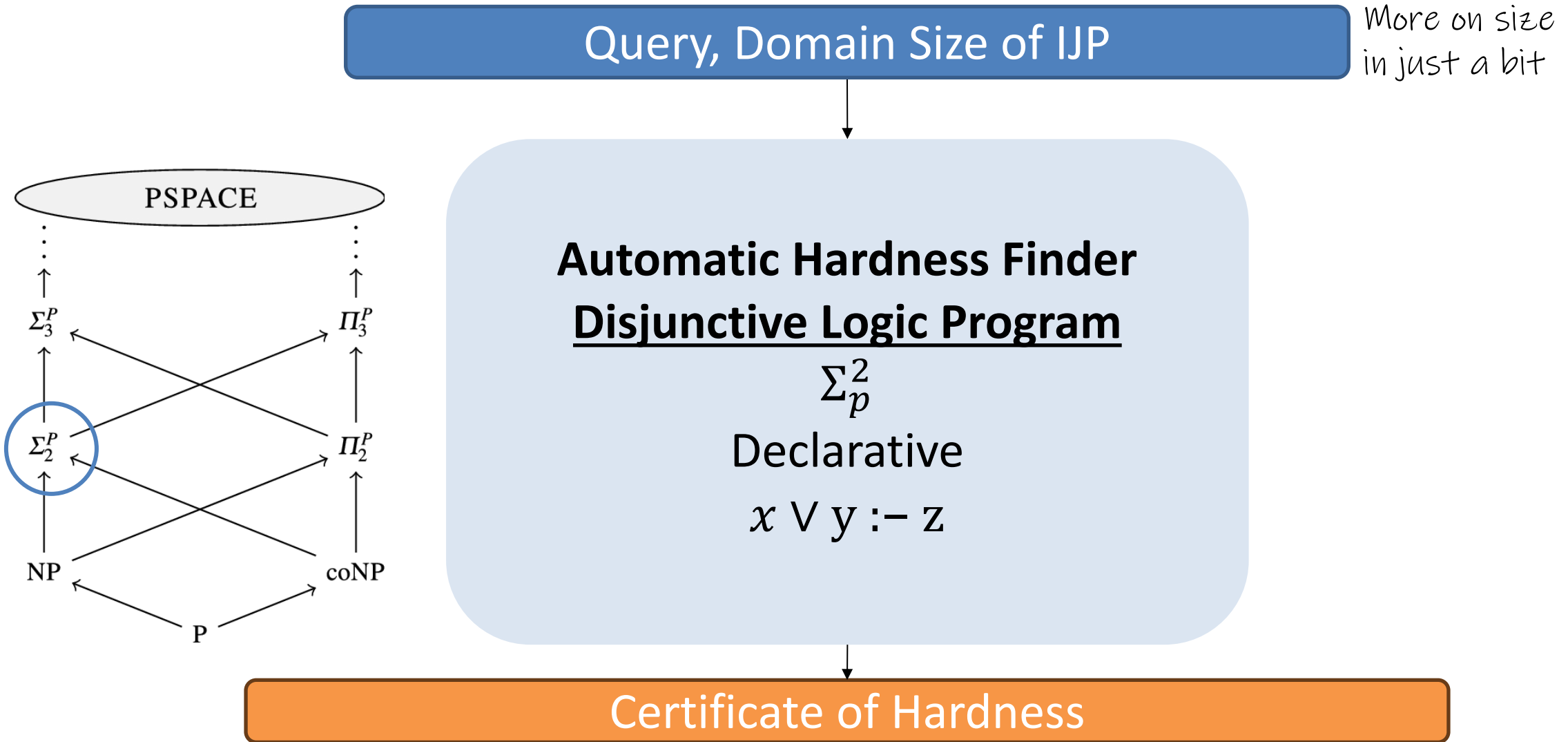


Certificate of Hardness



# Unified Approach to prove Hardness

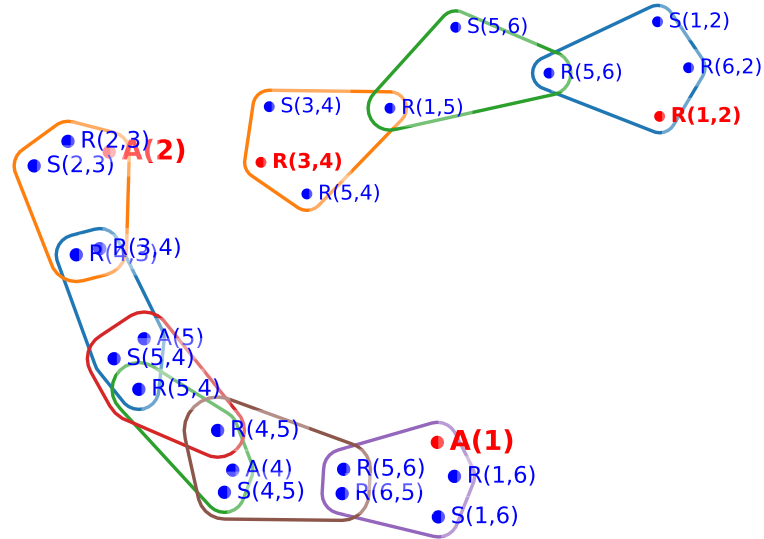
Using the complete criterion for IJP, we can build a **principled** way to find IJPs



# Finding IJPs with DLP: 5 New Hardness Gadgets

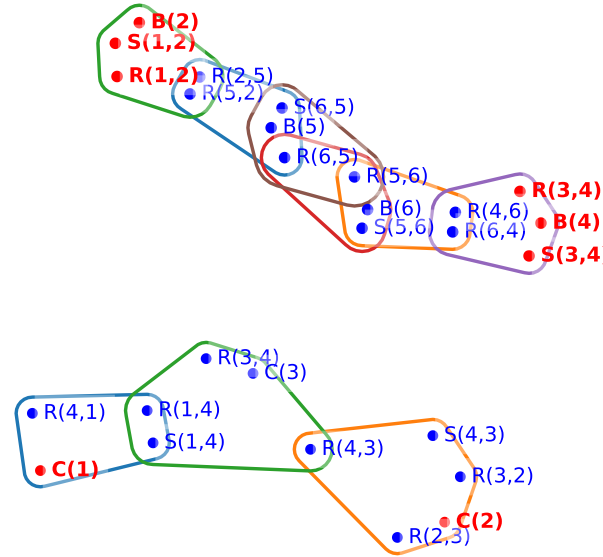
- Using the **Automatic IJP Generator** we proved 5 queries hard (out of 7 previously **open** from Freire+20)

$$q_{3cc}^S :- R(x,y), R(y,z), R(w,z), S(w,z)$$

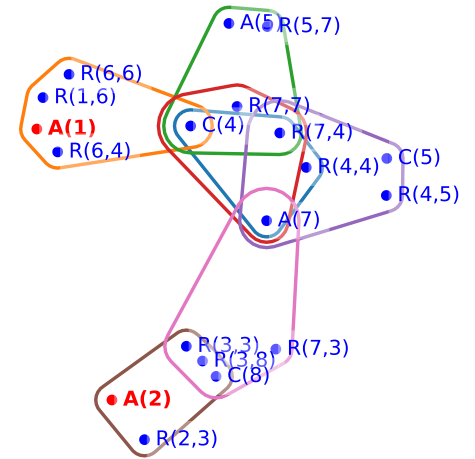


$$q_{3perm-R}^{AS_{xy}} :- A(x), S(x,y), R(x,y), R(y,z), R(z,y)$$

$$q_{3perm-R}^{S_{xy}B} :- S(x,y), R(x,y), B(y), R(y,z), R(z,y)$$



$$q_{3perm-R}^{S_{xy}C} :- S(x,y), R(x,y), R(y,z), R(z,y), C(z)$$



$$z_6 :- A(x), R(x,y), R(y,y), R(y,z), C(z)$$

- Can recover all previous hardness results + find new ones!

# Dichotomy Conjectures for Resilience

Theorem.

$IJP \rightarrow NPC$

Conjecture. [Hardness]

$IJP \leftrightarrow NPC$

Conjecture. [Hardness]

$IJP \rightarrow IJP$  of domain size  $\leq 7 * var(Q)$

Conjecture. [PTIME]

$\nexists IJP \rightarrow LP = ILP$

Theorem.

$IJP \leftrightarrow NPC$  for SJ-Free Queries

Conjecture. [Hardness] Corollary

$NPC \rightarrow$  DLP finds a hardness proof

Conjecture. [PTIME]

$\nexists IJP \rightarrow$  There is a flow graph that encodes resilience



# Takeaways

- One **unified** algorithm, only need to prove PTIME
- One **unified** hardness criterion
  - Automatic search



# Open Problems

- Which RDM problems can we solve with this unified approach?
  - Resilience
  - Causal Responsibility
  - Minimal Factorization of Provenance of CQs
  - ..... **Claim: many more**

Many more details, proofs, experiments, approximations:

- <https://northeastern-datalab.github.io/unified-reverse-data-management/>
- Makhija, Gatterbauer. A Unified Approach for Resilience and Causal Responsibility with Integer Linear Programming (ILP) and LP Relaxations, *To appear SIGMOD 2024* ( <https://arxiv.org/abs/2212.08898> )
- Makhija, Gatterbauer. Towards a Dichotomy for Minimally Factorizing the Provenance of Self-Join Free Conjunctive Queries, Arxiv 2021 ( <https://arxiv.org/abs/2105.14307> )