

# A **Unified** Approach for **Resilience** and **Causal Responsibility** with Integer Linear Programming (ILP) and LP Relaxations

Neha Makhija

Northeastern University

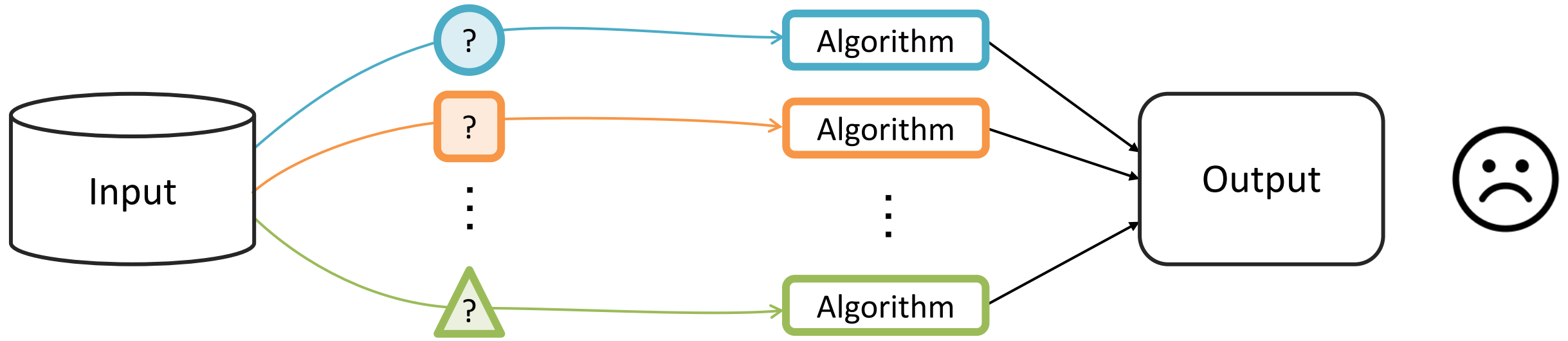
(Joint work with **Wolfgang Gatterbauer**)

SIGMOD 2024, Santiago



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

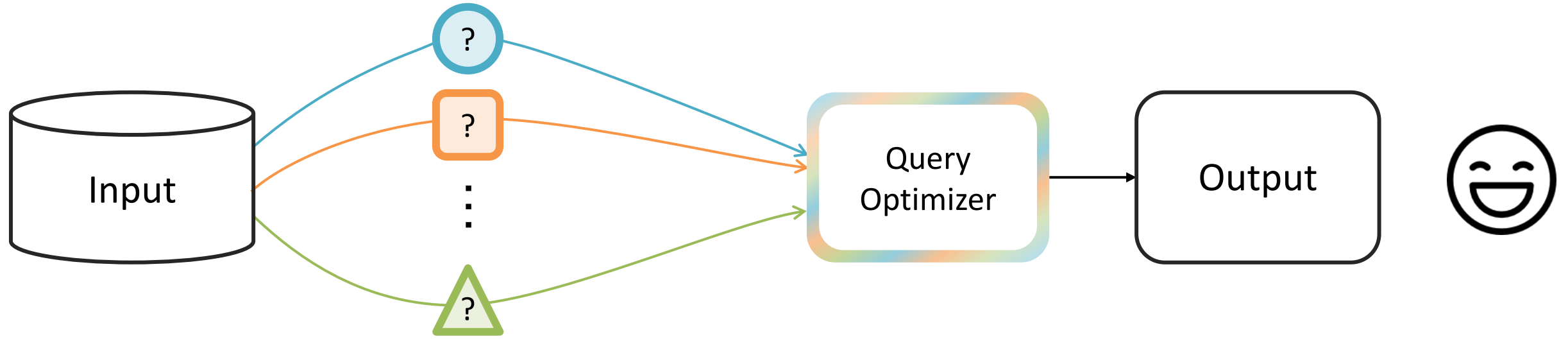
# Traditional Data Management (before query optimizers)



Many types of  
user questions!

Different questions  
→ Different Optimal Algorithms

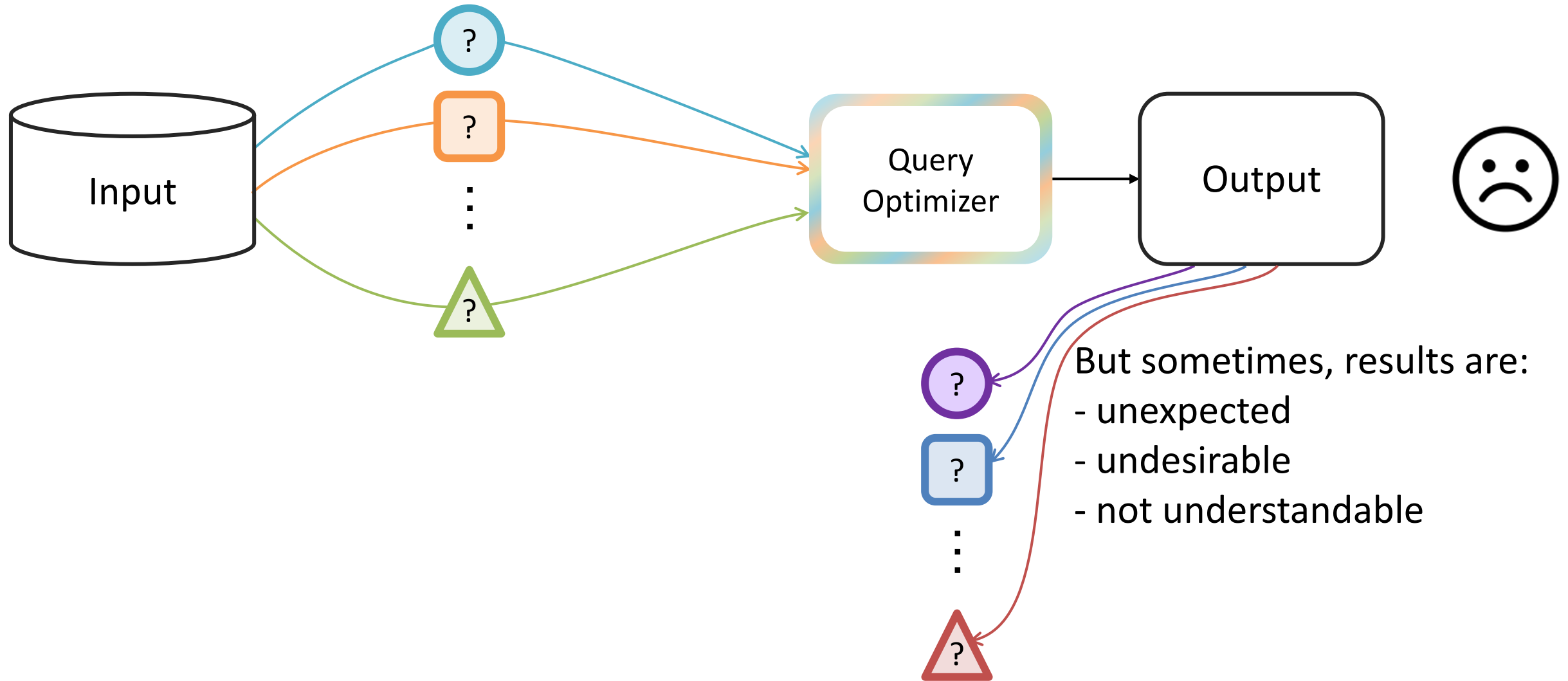
# Traditional Data Management



Many types of  
user questions!

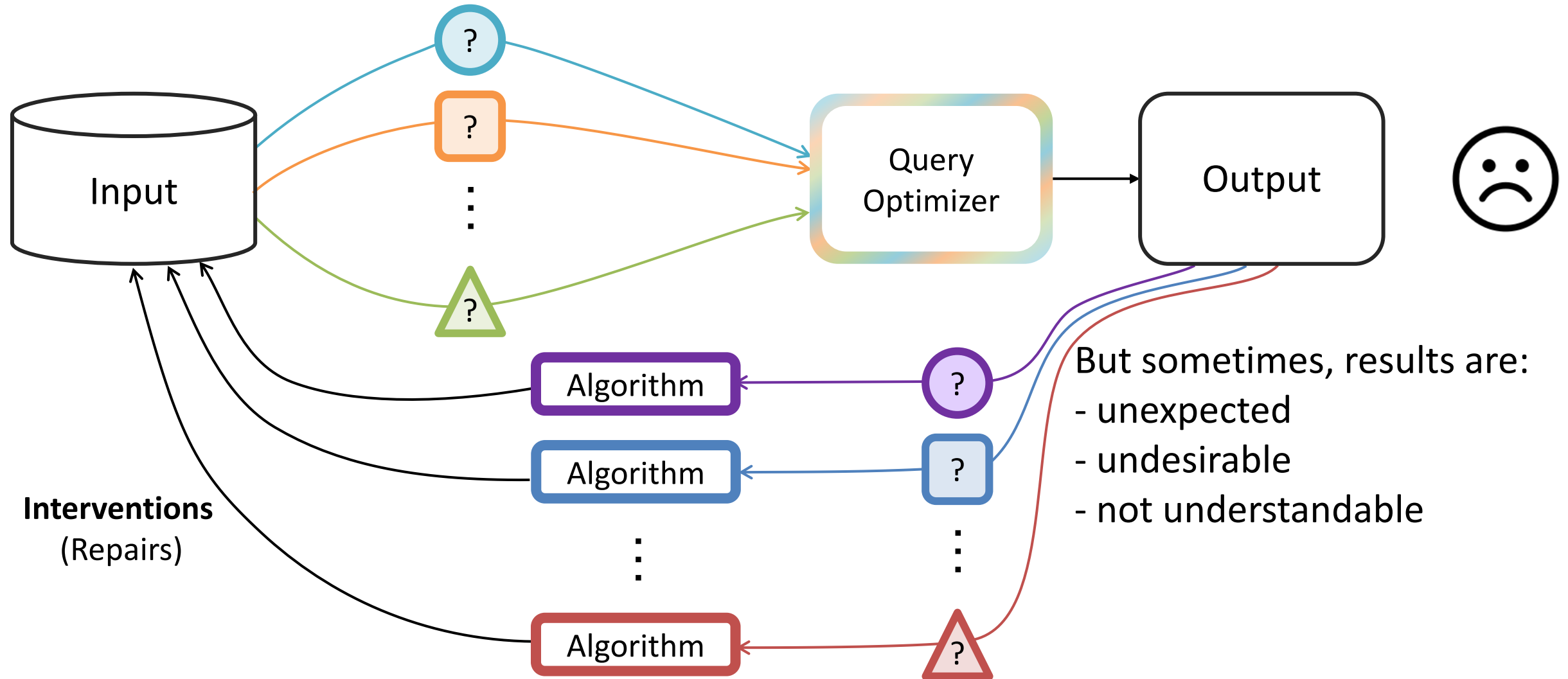
Different questions  
→ Query Optimizer solves  
everything optimally  
→ Choose between different access paths

# Reverse Data Management



**Reverse user questions: Why, how to, what if...**

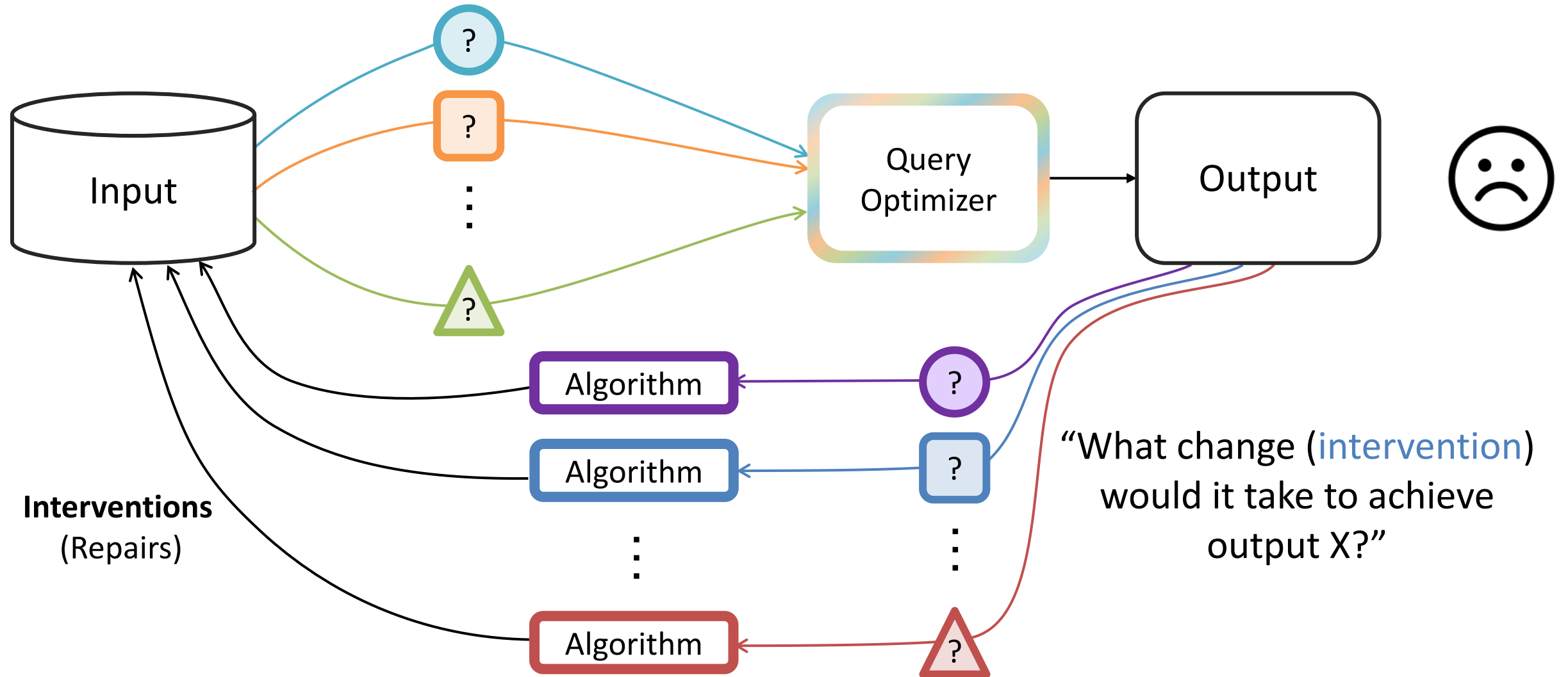
# Reverse Data Management



But sometimes, results are:  
- unexpected  
- undesirable  
- not understandable

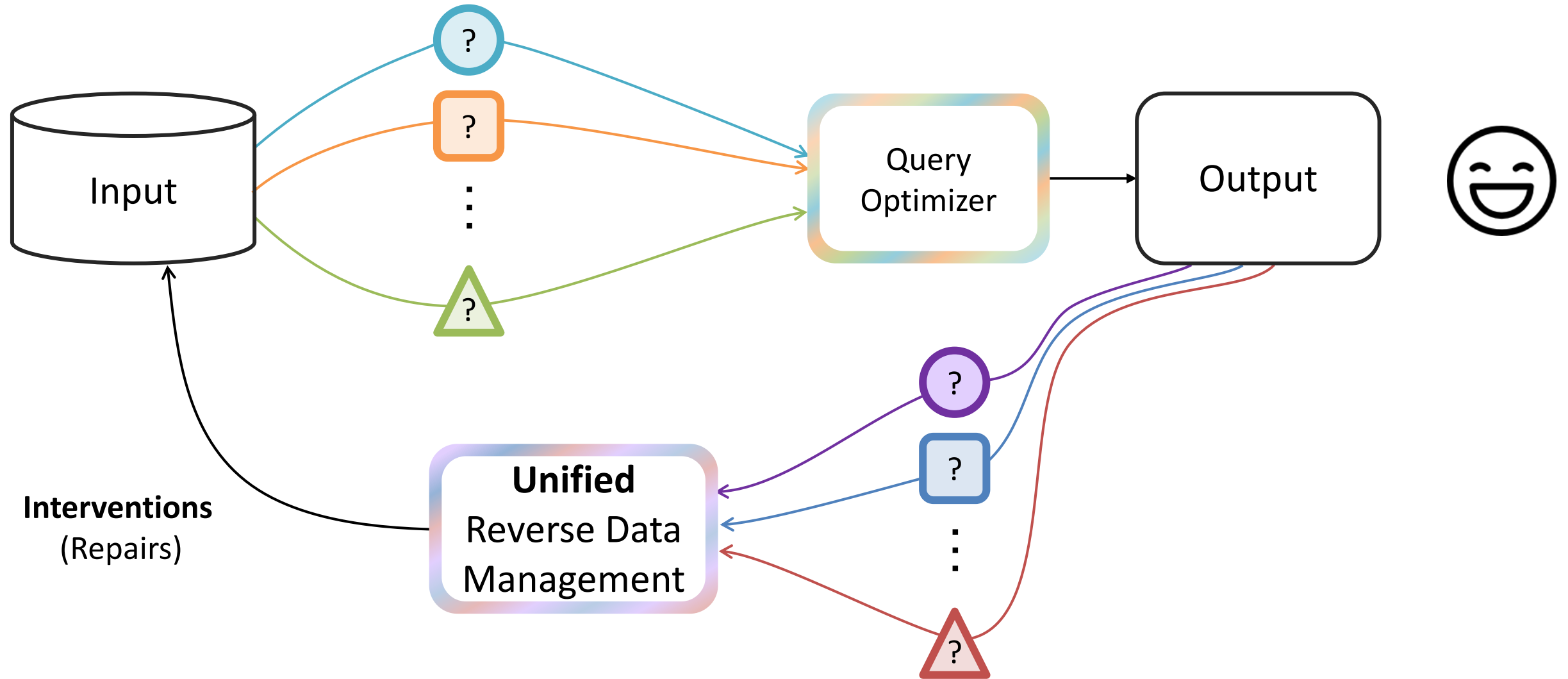
**Reverse** user questions: Why, how to, what if...

# Reverse Data Management



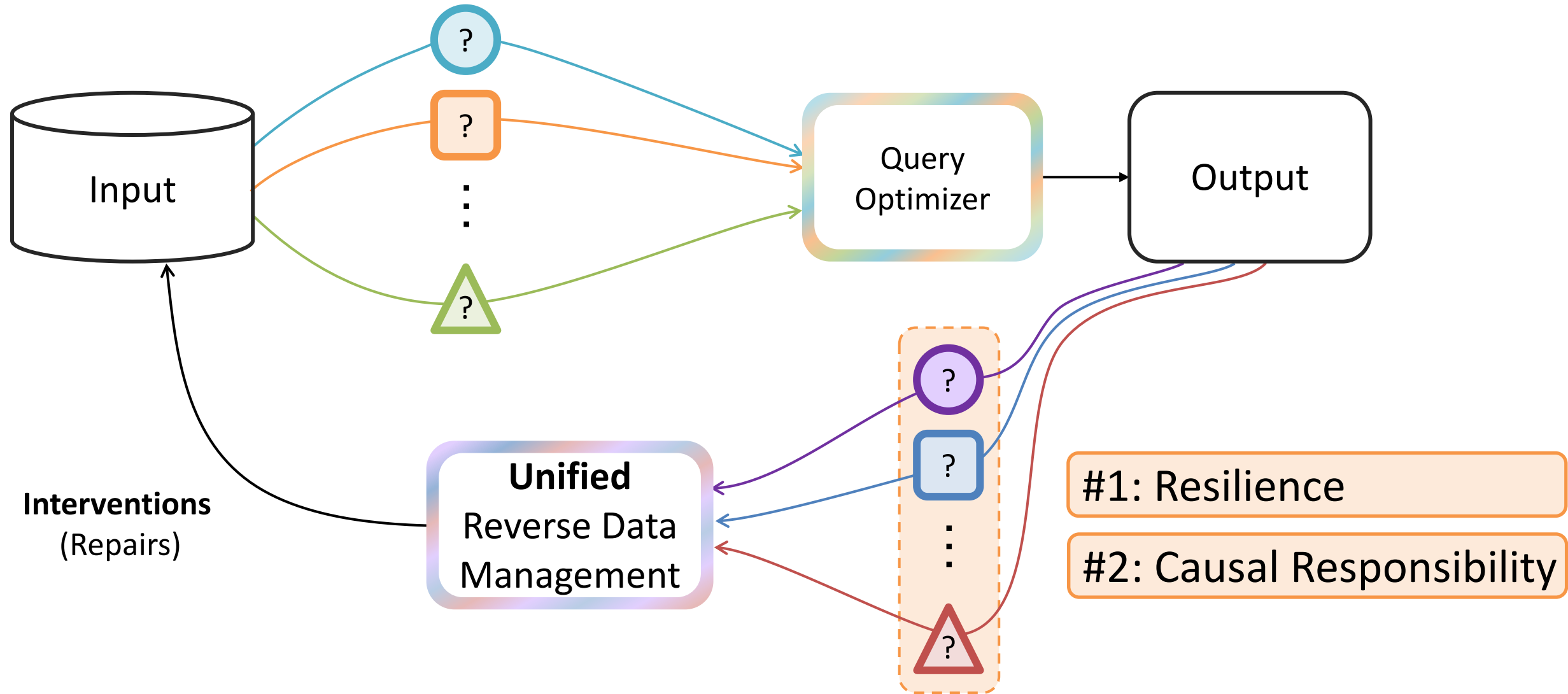
**Reverse** user questions: Why, how to, what if...

# Unified Reverse Data Management: Our Vision



**Reverse** user questions: Why, how to, what if...

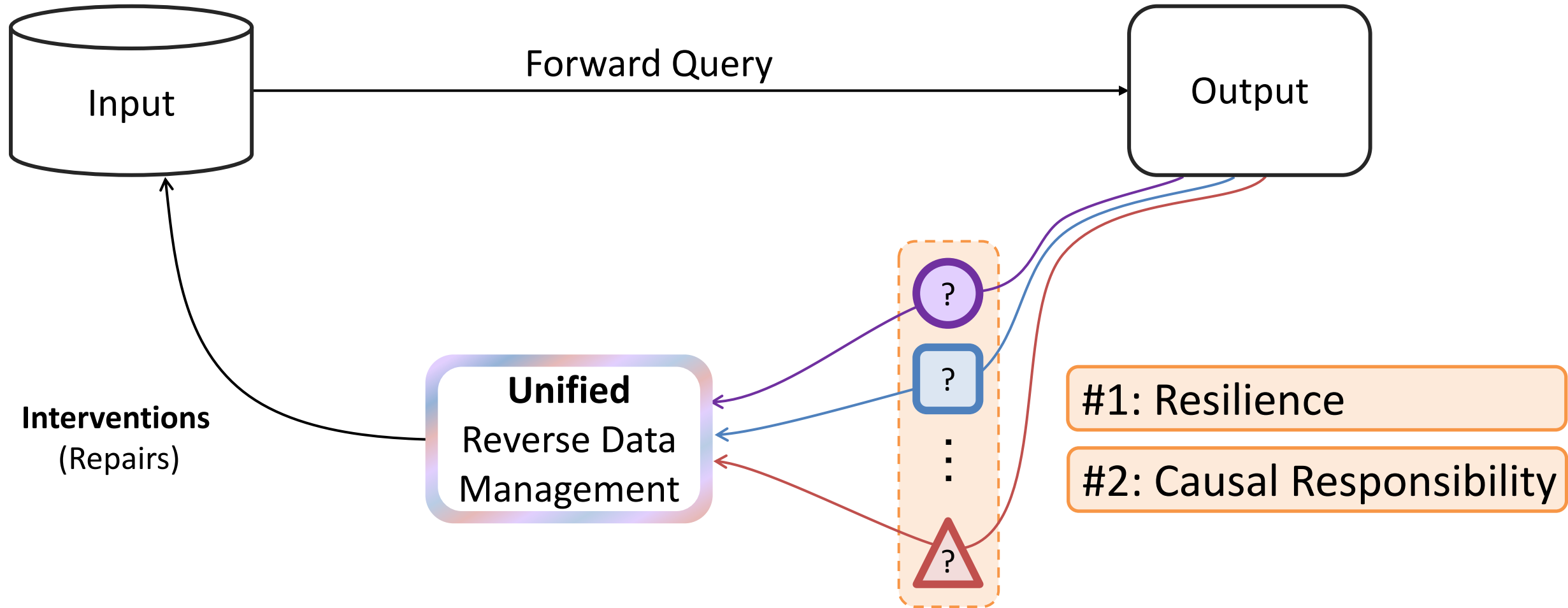
# Reverse Data Management: Our Focus



**Reverse** user questions: Why, how to, what if...

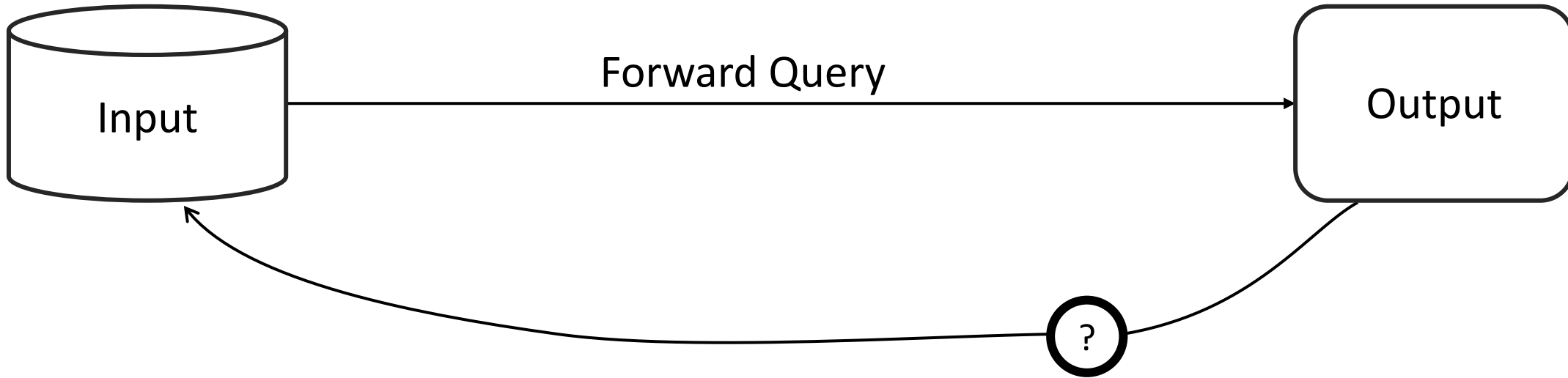


# Reverse Data Management: **Our Focus**



**Reverse** user questions: Why, how to, what if...

# Reverse Data Management: Resilience

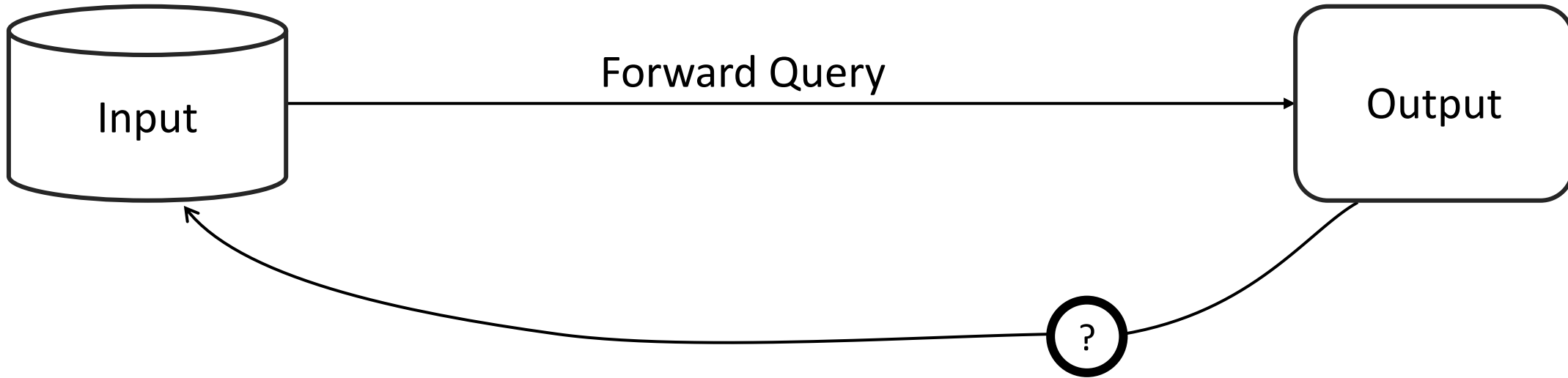


“What change (**intervention**) would it take to achieve **no output**?”

#1: Resilience

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

# Reverse Data Management: Resilience



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

#2: Causal Responsibility

- Halpern-Pearl Framework of Counterfactual Causality adapted to Conjunctive Queries

- Unified Reverse Data Management Framework
  - Example: Resilience
- Insight #1: How to build **Unified Algorithms**
- Insight #2: How to build **Automatic Hardness Provers**
- What else is in the paper?
- Takeaways + Open Questions

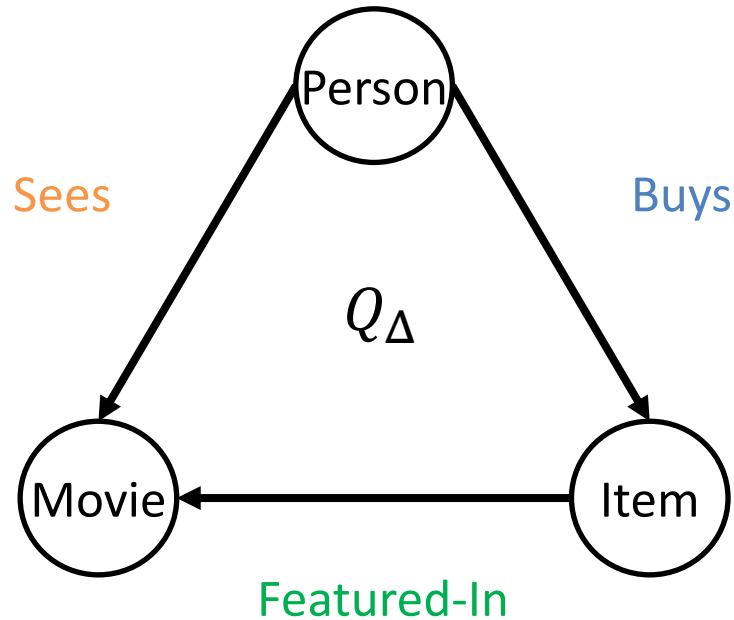
# Resilience: Example

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?

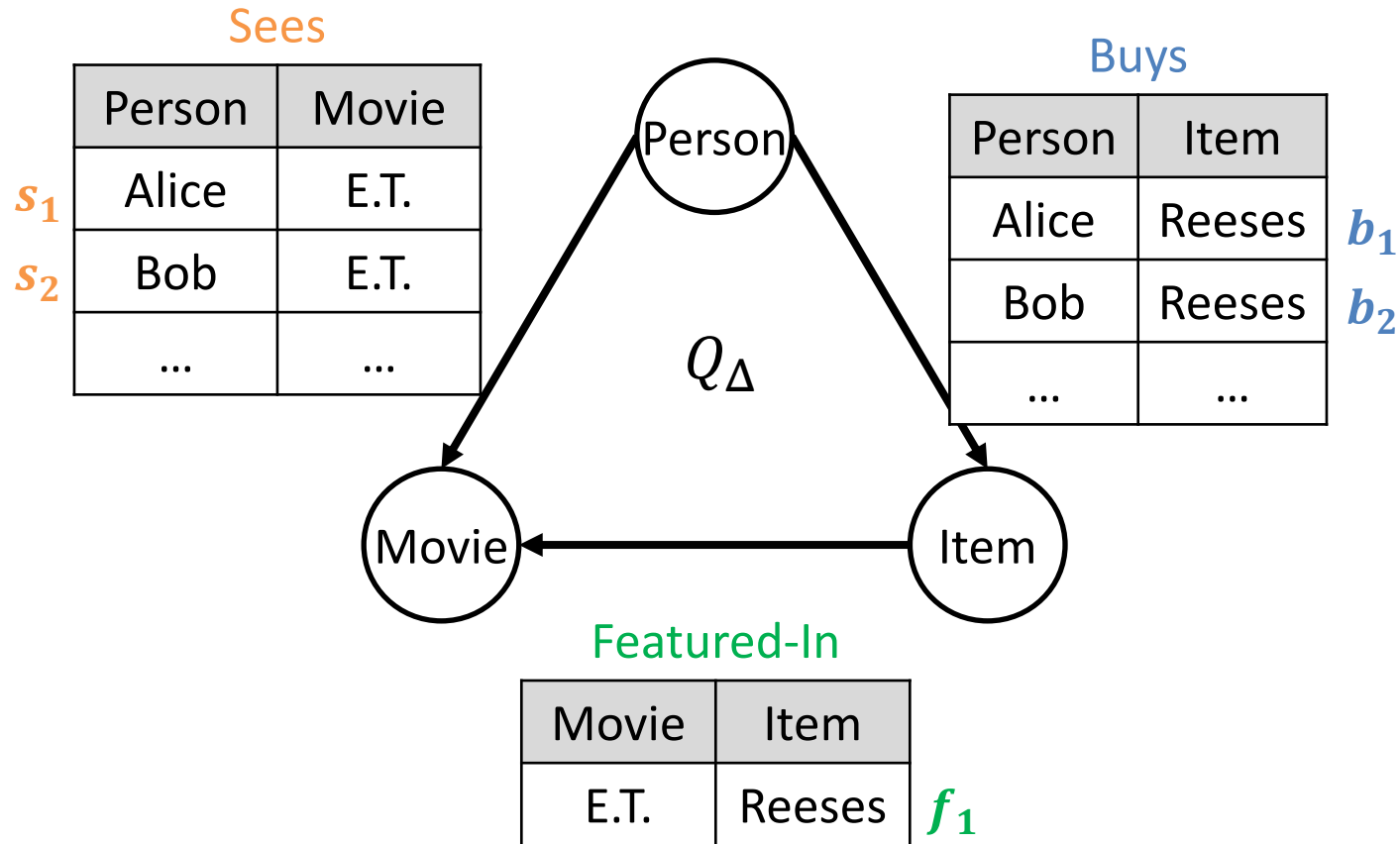
Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)



# Resilience: Example

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

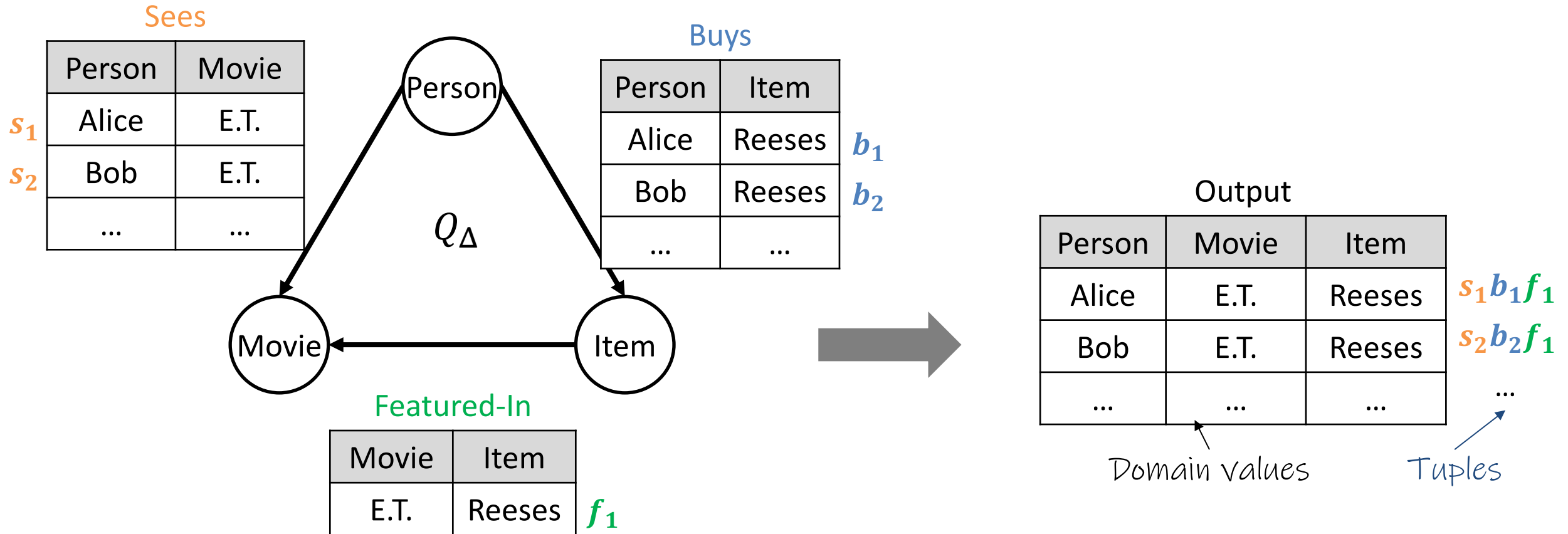
Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)



# 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)

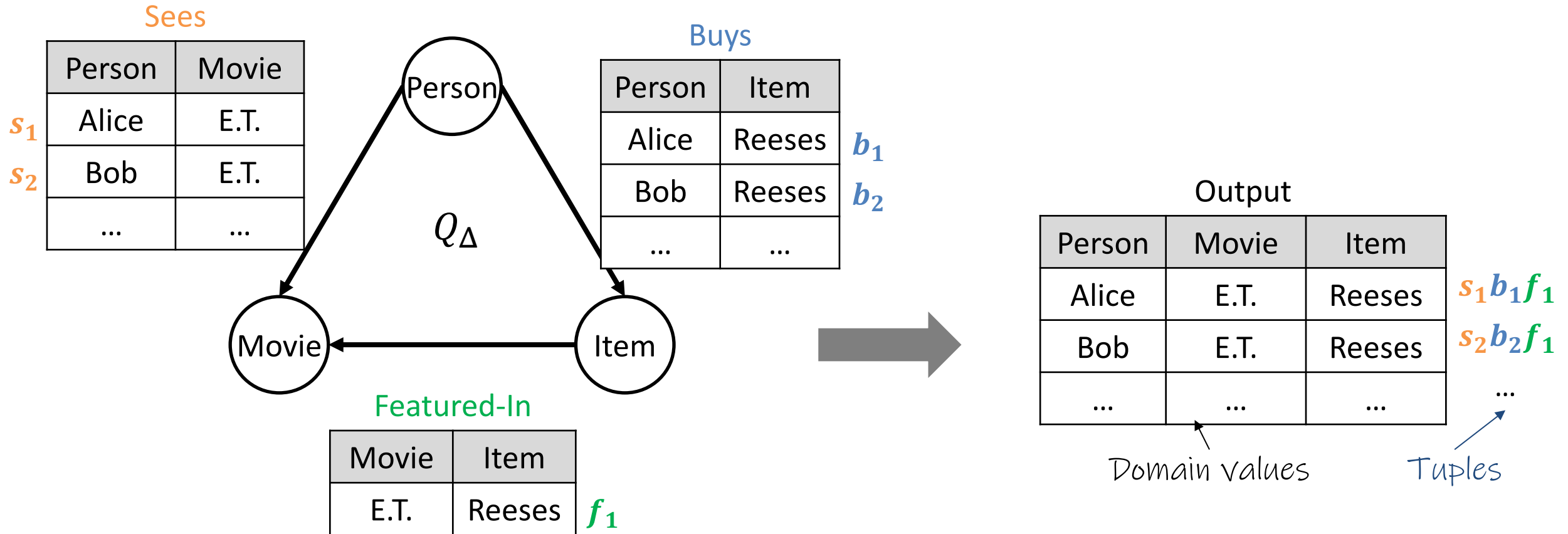




# Resilience: Example

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)

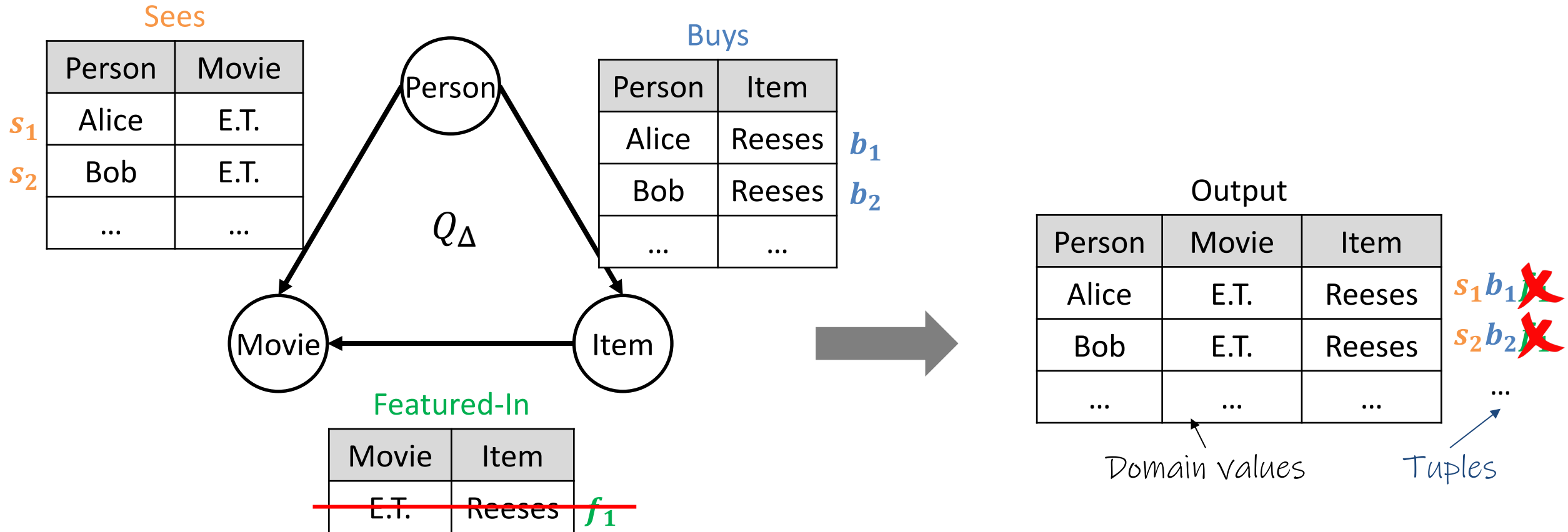


Recall: Resilience = What minimal change would it take to delete the output?"

# 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)

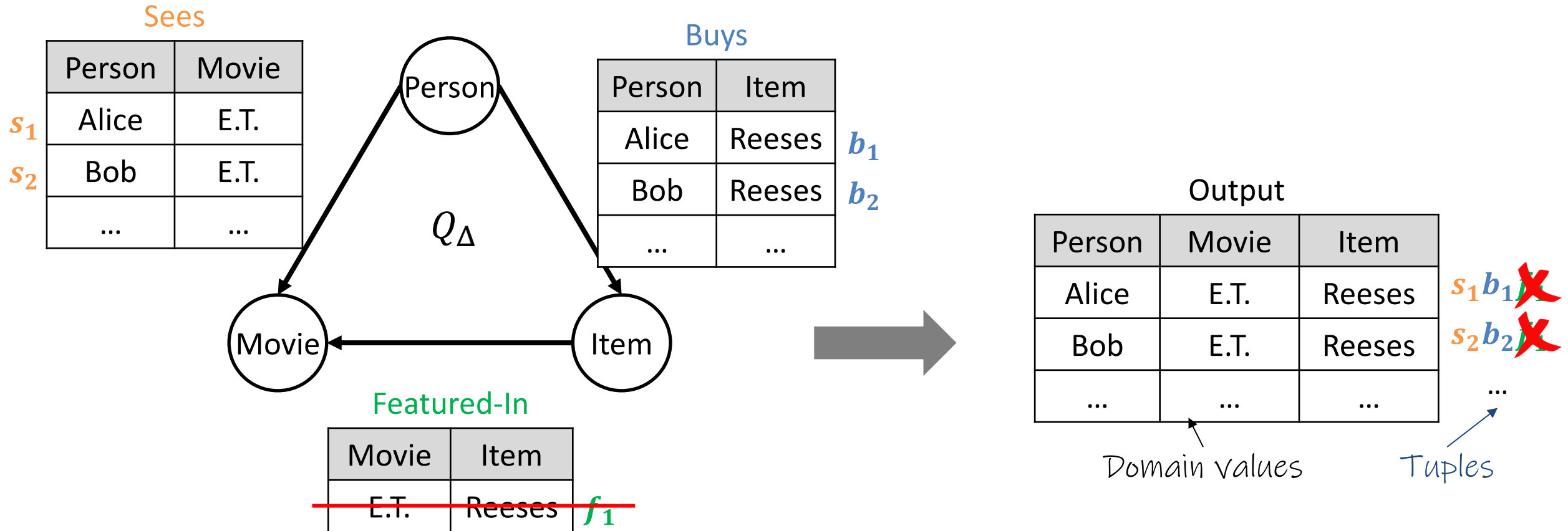


Recall: Resilience = What minimal change would it take to delete the output?"

# Resilience: Example

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

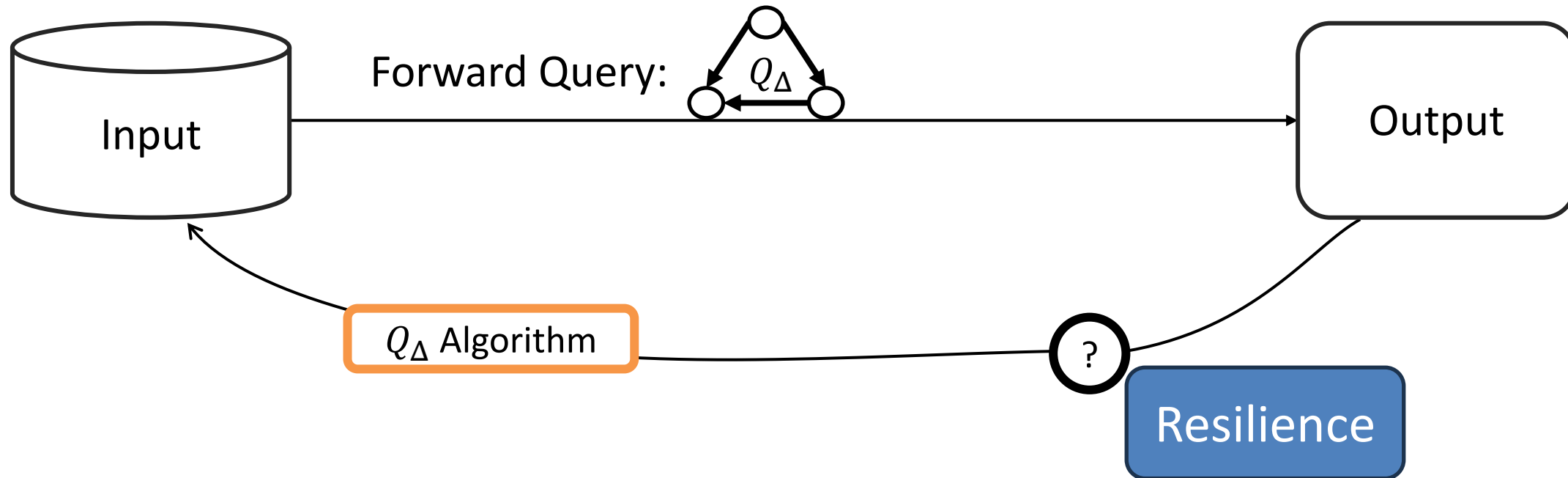
Q(person, movie, item):- **Sees**(person, movie), **Buys**(person, item), **Featured-In**(item, movie)



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

# Resilience: Example

$Q(\text{person}, \text{movie}, \text{item})$ :- Sees(person, movie), Buys(person, item), Featured-In(item, movie)

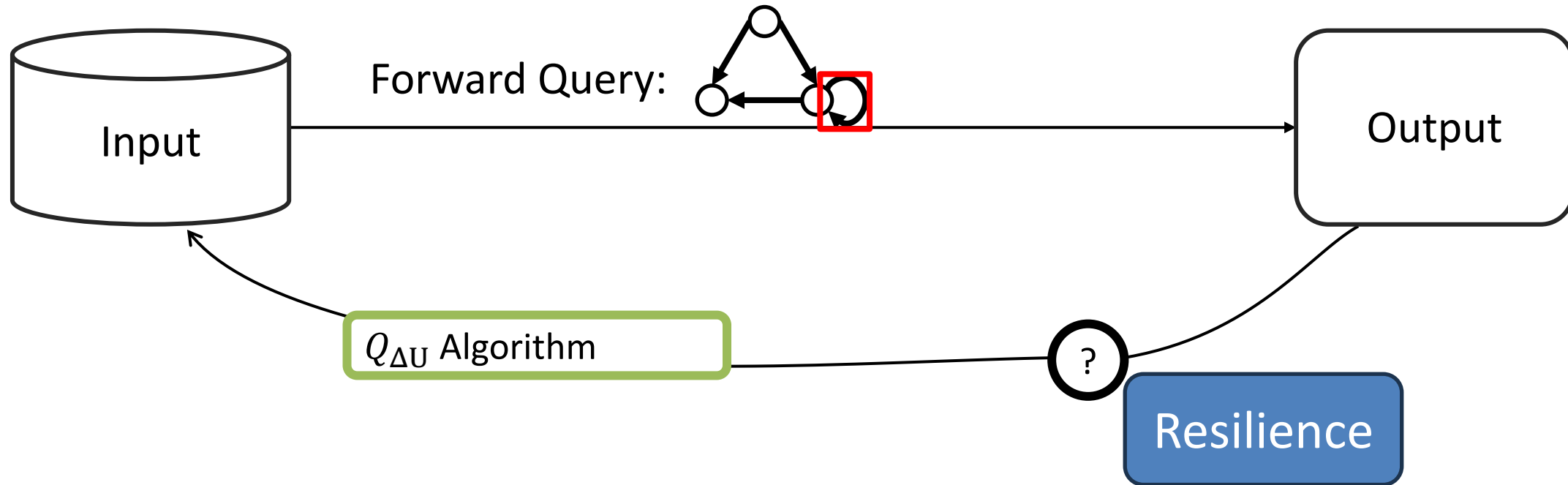


This is a NP-Hard Problem

- Optimal algorithm is exponential!  
(Unless  $P=NP$ )

# Resilience: Example

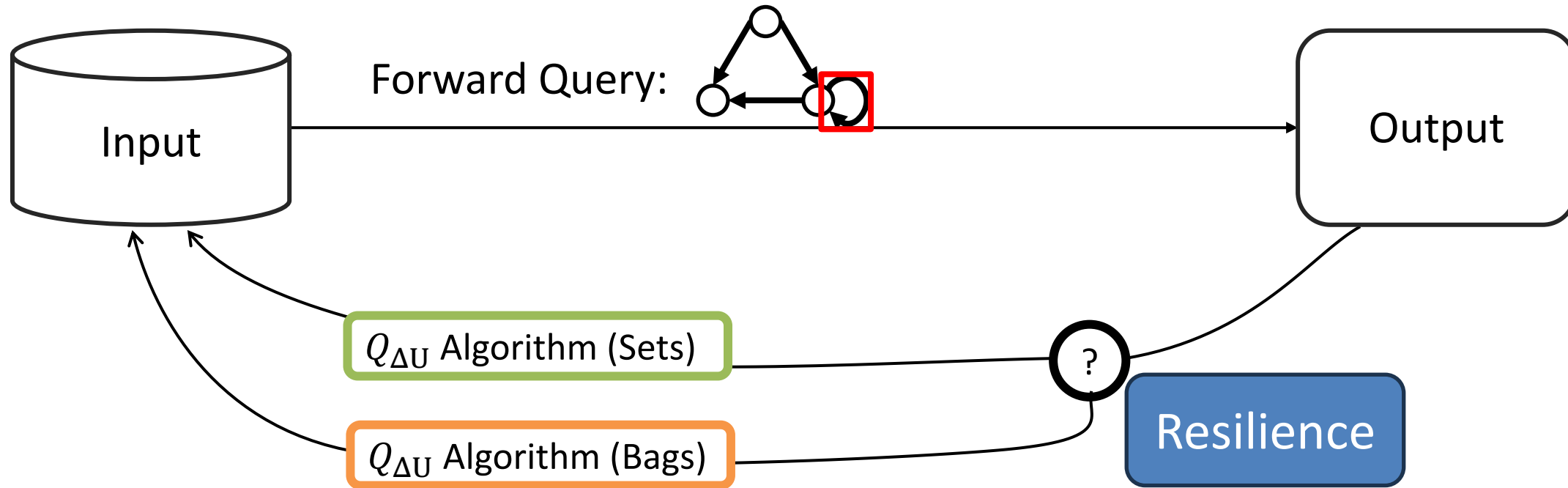
$Q(\text{person}, \text{movie}, \text{item})$ :- Sees(person, movie), Buys(person, item), Featured-In(item, movie), Rare(Item)



Minor modification makes query PTIME

# Resilience: Example

$Q(\text{person, movie, item})$ :- Sees(person, movie), Buys(person, item), Featured-In(item, movie), Rare(Item)



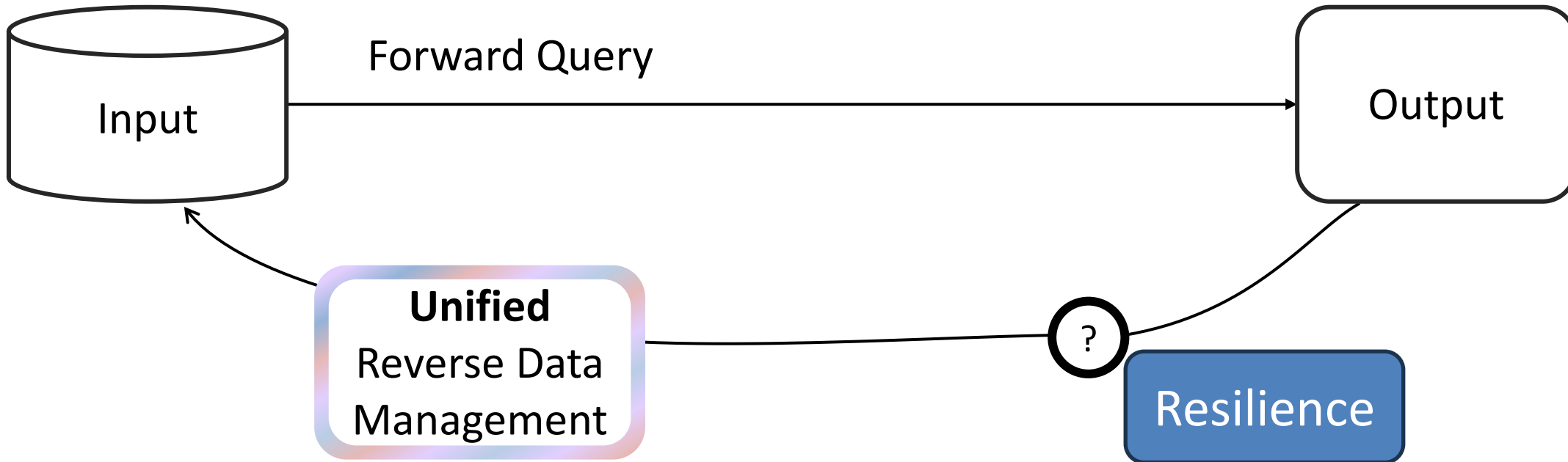
Minor modification makes query PTIME

- But only under set semantics!

If Input uses Bag semantics:

- Optimal algorithm is exponential! (Unless  $P=NP$ )

# Resilience: Example



Unified approach is *automatically* optimal

- For all known queries
- For bag / set semantics

- Unified Reverse Data Management Framework
- **Insight #1: How to build Unified Algorithms**
- Insight #2: How to build **Automatic Hardness Provers**
- What else is in the paper?
- Takeaways + Open Questions



# What is a Unified Algorithm?

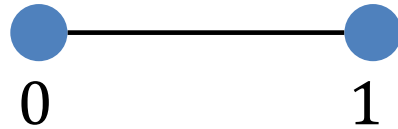
1. Unified Across Complexity (Easy or Hard)
  - Algorithm that can solve PTIME and NP-C problems
  - **Guaranteed** exact PTIME termination for all known tractable cases
2. Unified Across Settings
  - Set + **Bag** Semantics
  - All Conjunctive Queries (including **all self-joins**)
  - Can take advantage of **unspecified** Functional Dependencies

**Orange** keywords = contrast from previous work

# Detour: Integer Linear Programs

## Integer Linear Program (ILP)

$$\begin{aligned} \min w^T x \\ \text{s.t. } Ax \leq b \\ x \in \{0, 1\} \end{aligned}$$



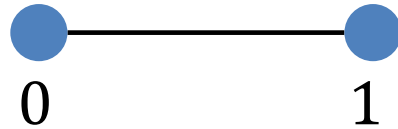
NP-Hard

Highly Optimized in Practice!

# Detour: Integer Linear Programs

## Integer Linear Program (ILP)

$$\begin{aligned} \min w^T x \\ \text{s.t. } Ax \leq b \\ x \in \{0, 1\} \end{aligned}$$



NP-Hard

Highly Optimized in Practice!

## Linear Program (LP)

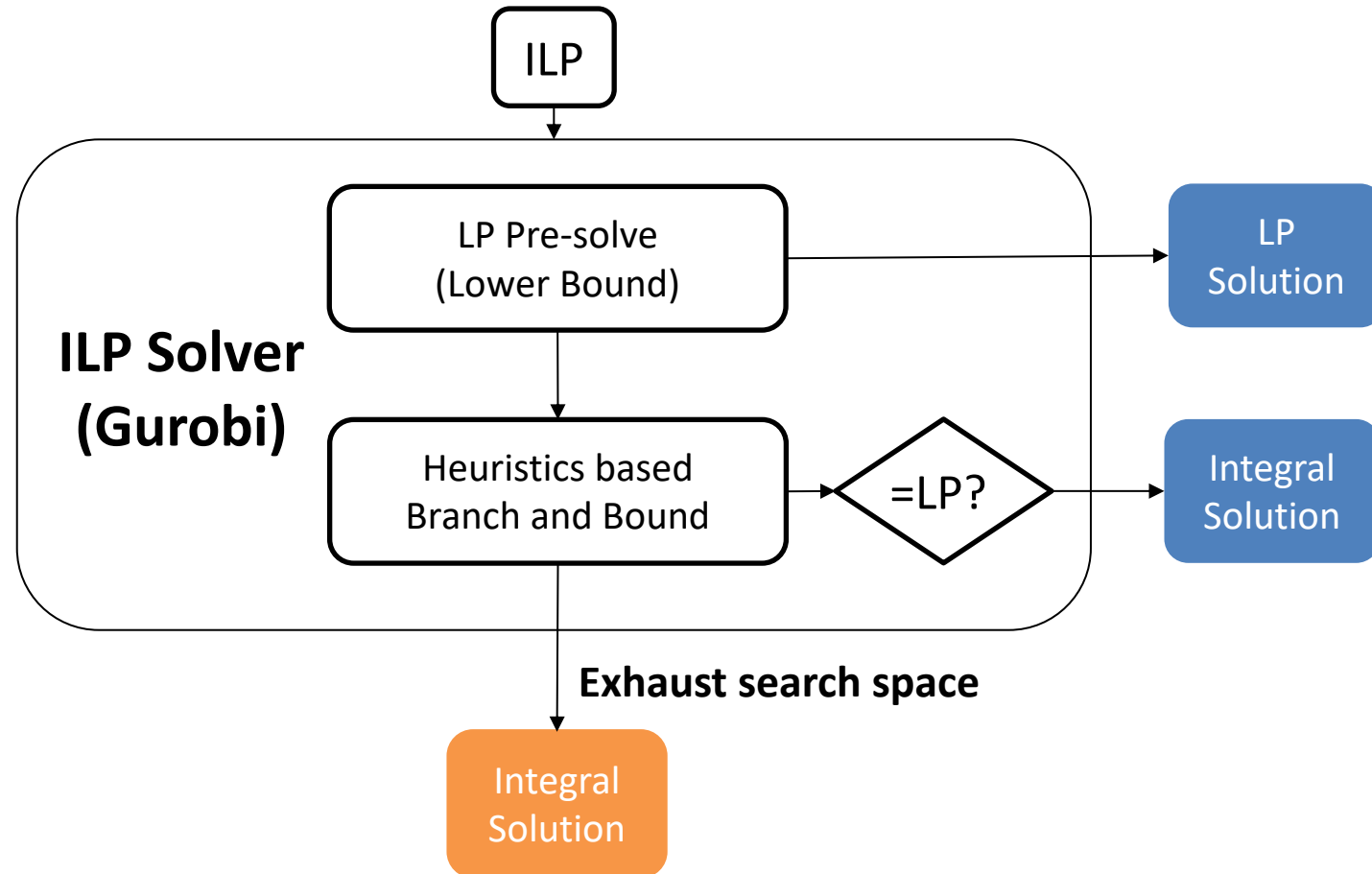
$$\begin{aligned} \min w^T x \\ \text{s.t. } Ax \leq b \\ x \in [0, 1] \end{aligned}$$



PTIME

Natural Lower Bound

# Detour: Integer Linear Programs



Insight: If **LP=ILP**, then a solution can be recovered efficiently!

# Resilience as an Integer Linear Program (ILP)

Weights if applied can go in the objective

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

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

$Q_\Delta$  Example:

Output

Person	Movie	Item
Alice	E.T.	Reeses
Bob	E.T.	Reeses
...	...	...

$s_1 b_1 f_1$

$s_2 b_2 f_1$

...



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

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

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

# Resilience as an Integer Linear Program (ILP)

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

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

$Q_\Delta$  Example:

Output

Person	Movie	Item	
Alice	E.T.	Reeses	$s_1 b_1 f_1$
Bob	E.T.	Reeses	$s_2 b_2 f_1$
...	...	...	...



Theorem.

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



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

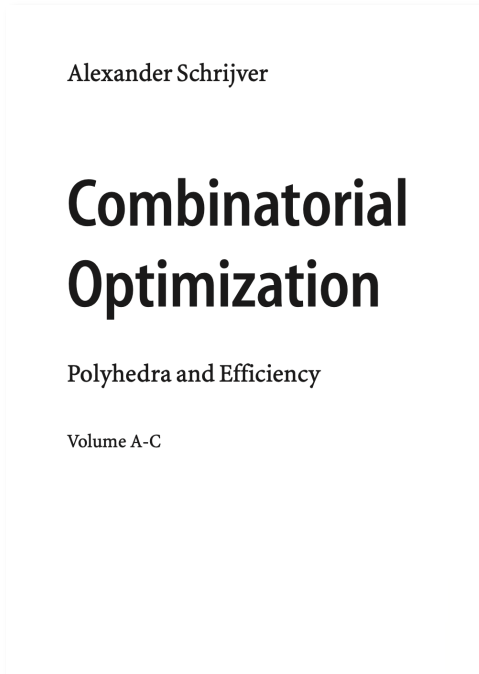
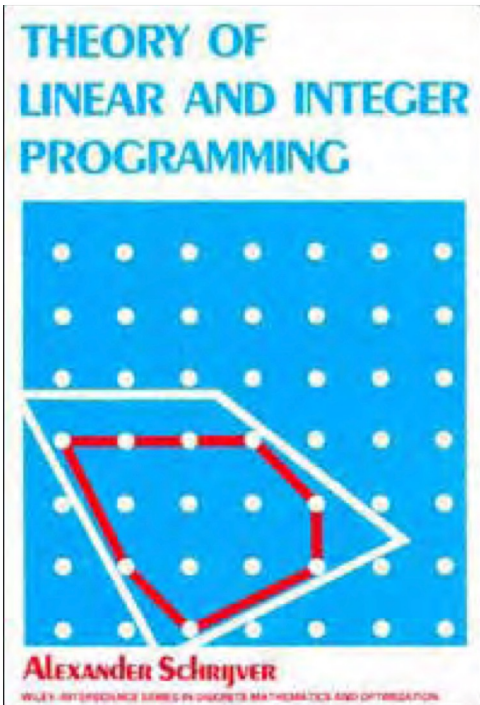
# When does LP = ILP?

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 criteria!



83	<del>Balanced and unimodular hypergraphs</del>	1439
83.1	<del>Balanced hypergraphs</del>	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	<del>Unimodular hypergraphs</del>	1448
83.3a	<del>Further notes</del>	1450

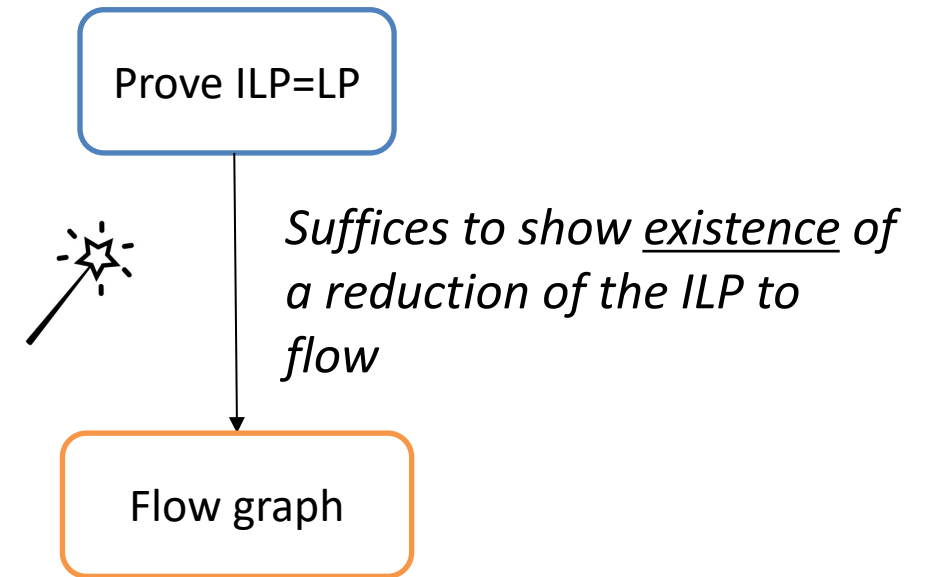
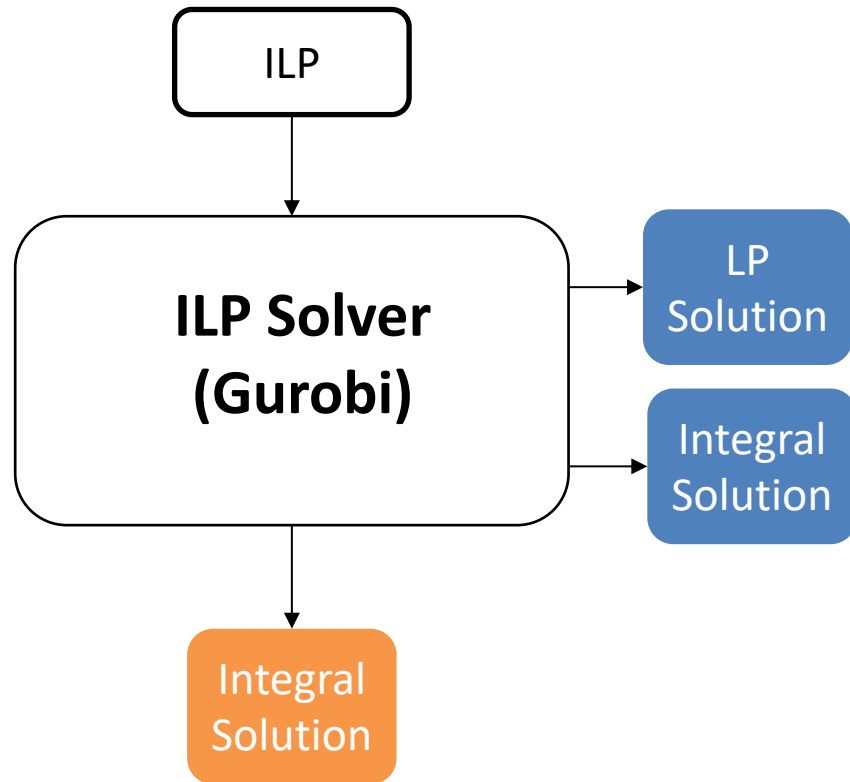


Show correspondence to a flow graph as a criterion for when LP=ILP

# Unified Approach for Resilience

Typical Goal: Find algorithm for easy cases

Our Goal: *Just* prove case is easy



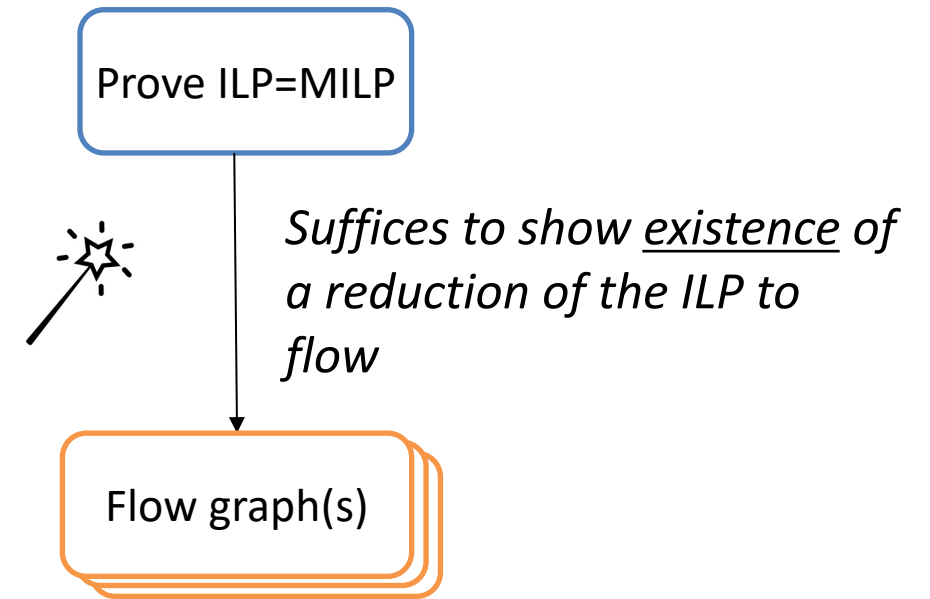
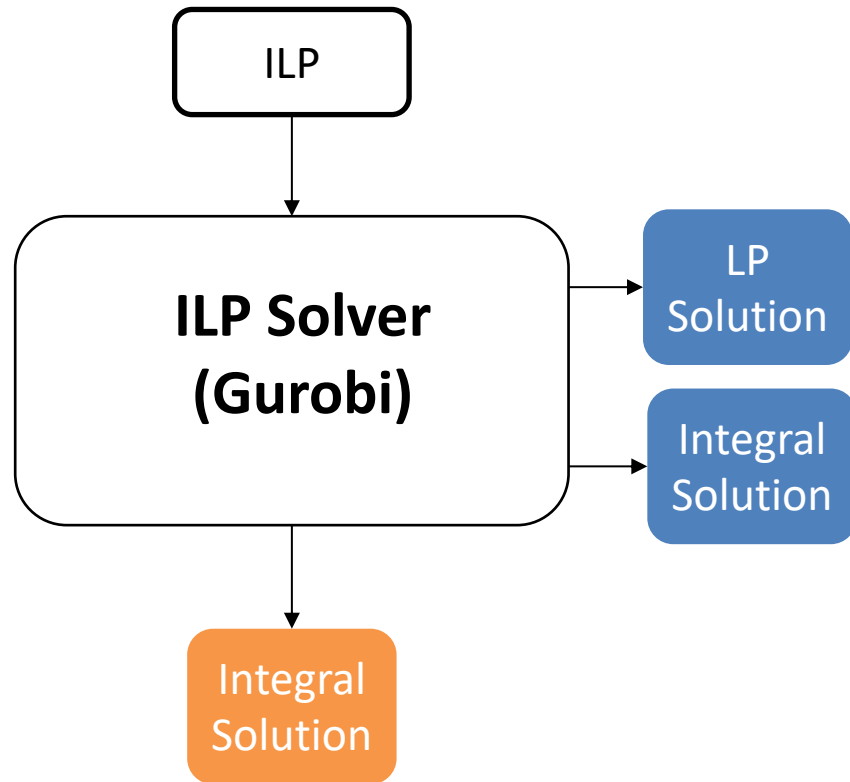
Use classical Max-Flow Min-Cut Theorem to show LP = ILP



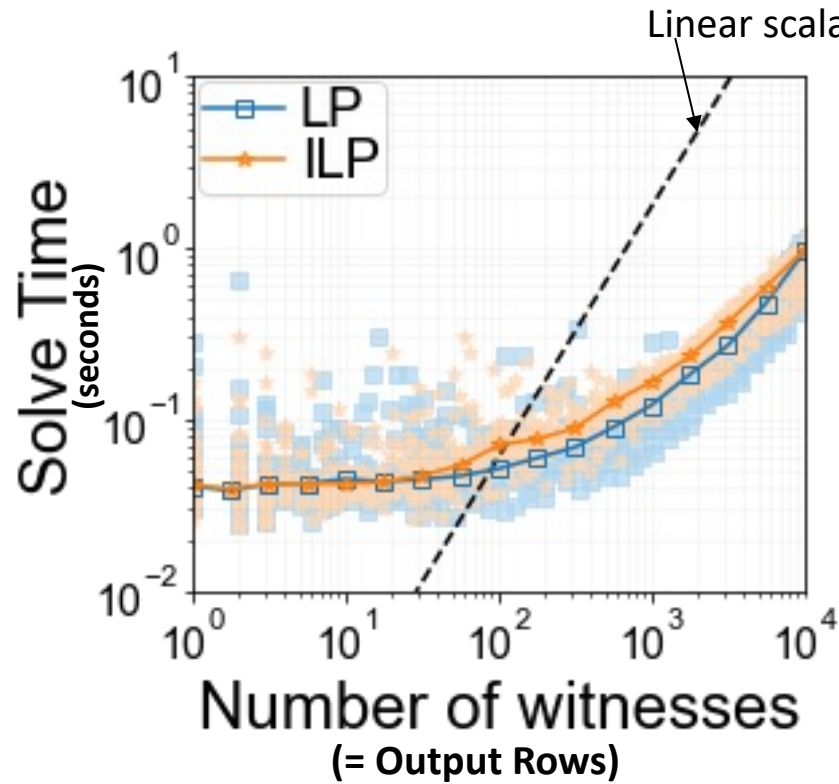
# Unified Approach for Causal Responsibility

Typical Goal: Find algorithm for easy cases

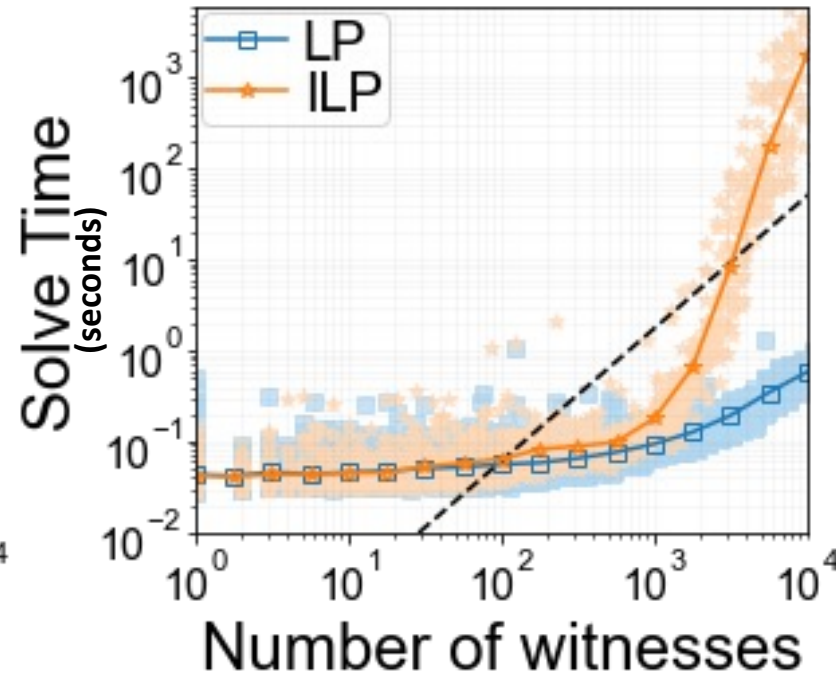
Our Goal: *Just* prove case is easy



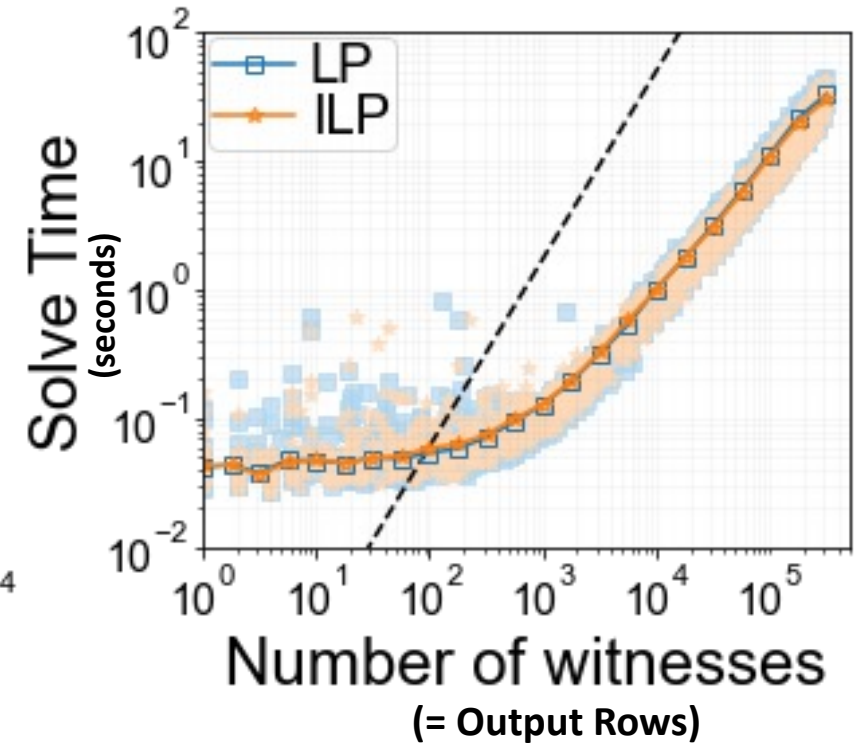
# Experiments: Unified Resilience ILP



Q1: SJ-Confluence  
Known **PTIME**



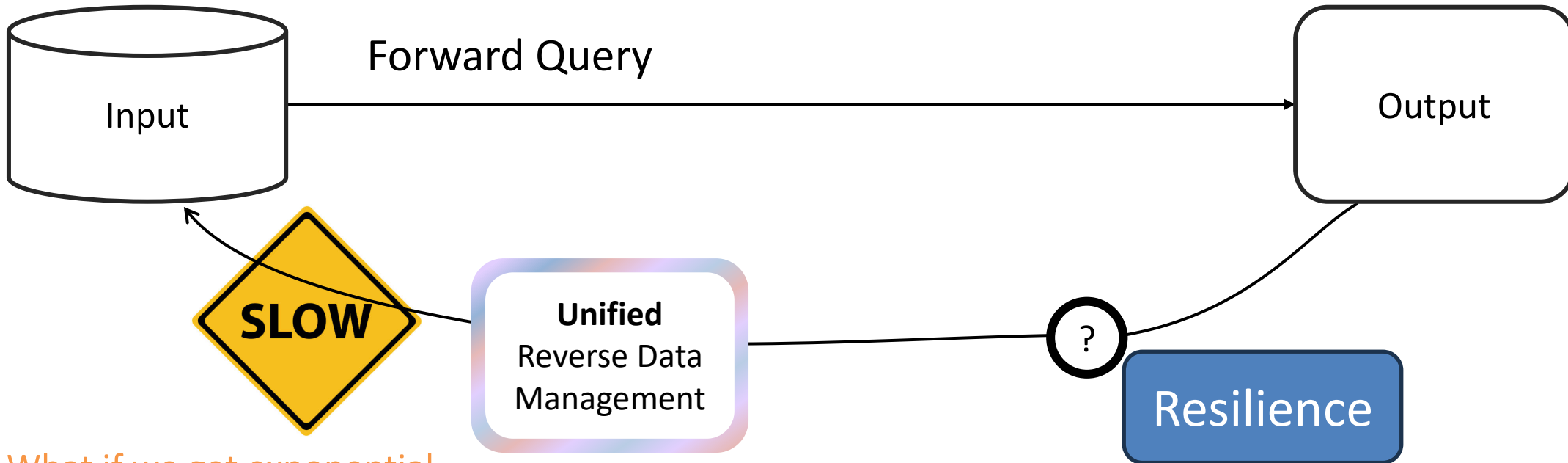
Q2: SJ-Chain  
Known **NP-Complete**



Q3: SJ-z6  
Previously **Open**  
We show **NP-Complete**

- Unified Reverse Data Management Framework
- Insight #1: How to build **Unified Algorithms**
- **Insight #2: How to build Automatic Hardness Provers**
- What else is in the paper?
- Takeaways + Open Questions

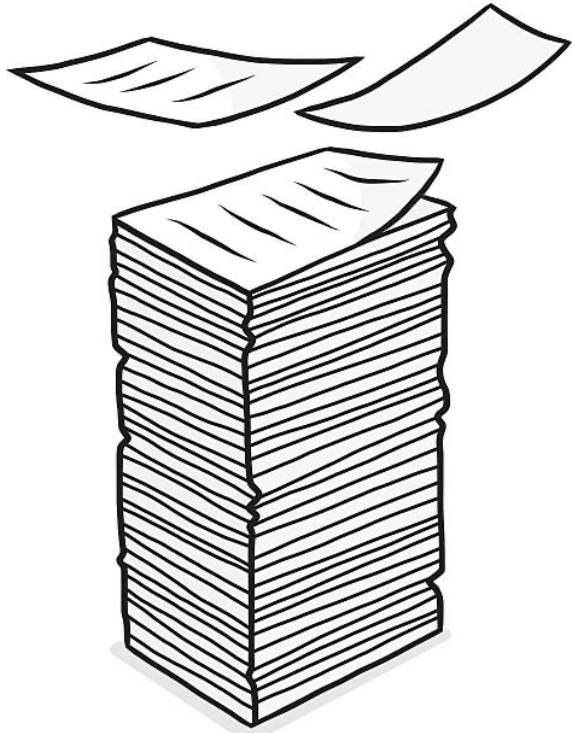
# Unified Approach to prove Hardness



What if we get exponential complexity?

How to prove that exponential complexity is optimal (best we can do)?

# Unified Approach to prove Hardness

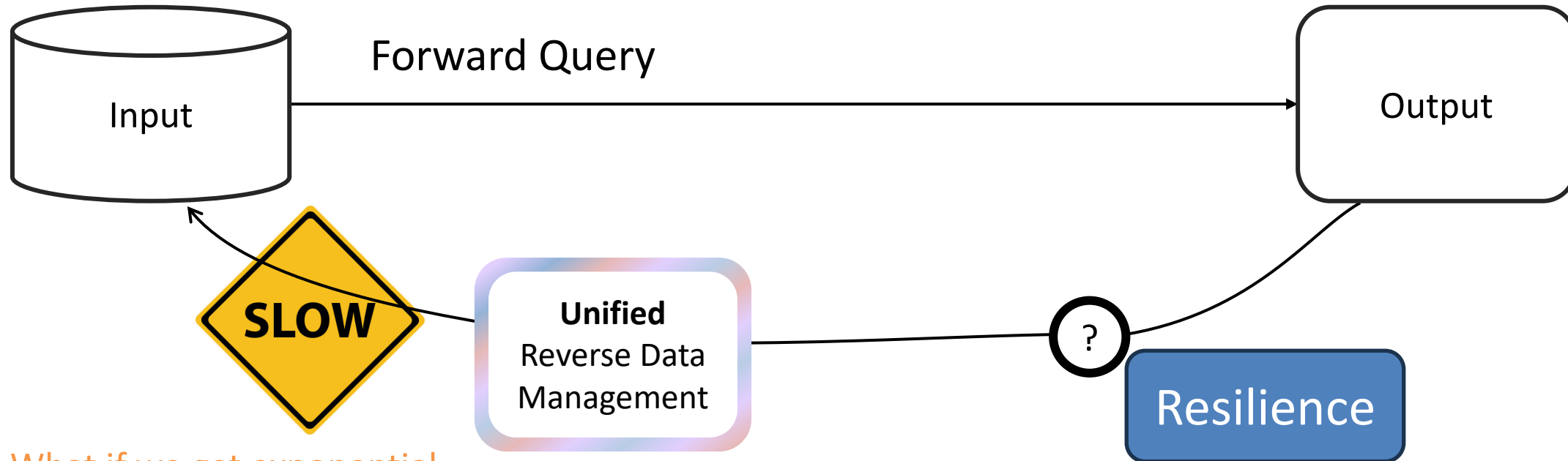


*..... Hence, If you can solve Problem A, you can solve a NP-hard problems like Vertex Cover*



I computed a hardness certificate you can easily verify

# Unified Approach to prove Hardness

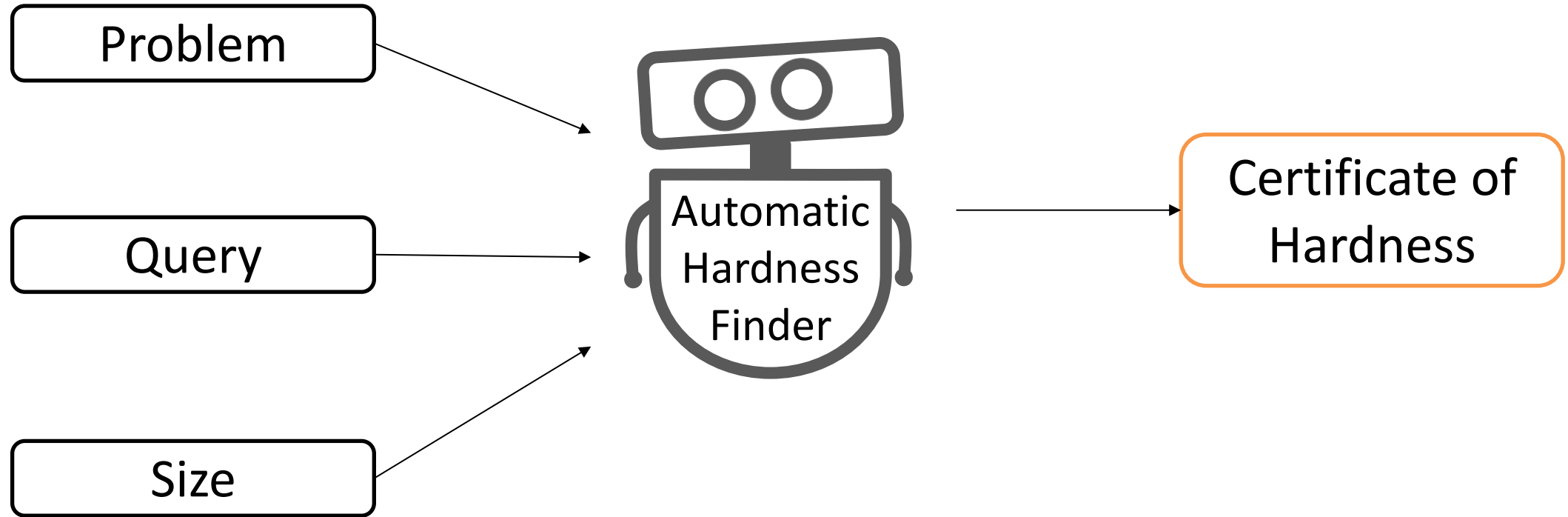


What if we get exponential complexity?

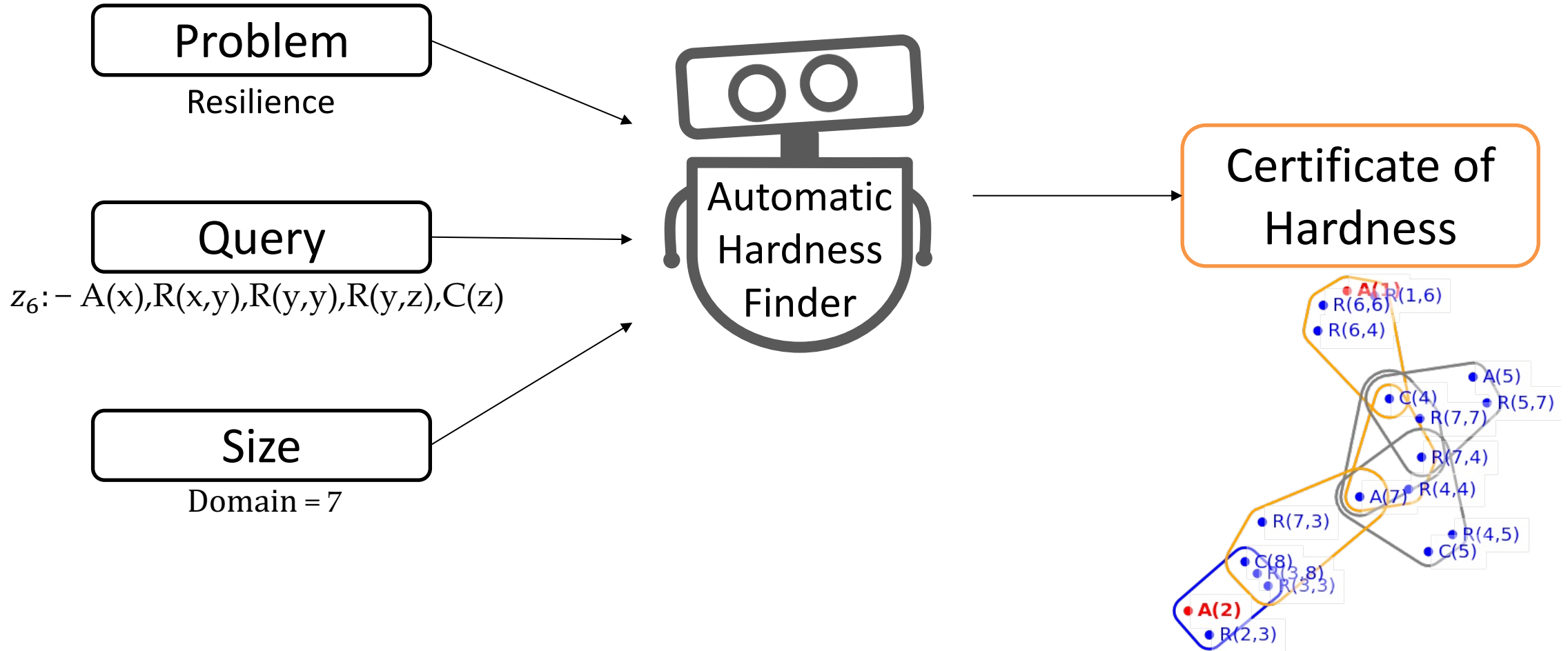
*automatically*

How to <sup>↑</sup>prove that exponential complexity is optimal (best we can do)?

# Unified Approach to prove Hardness



# Unified Approach to prove Hardness

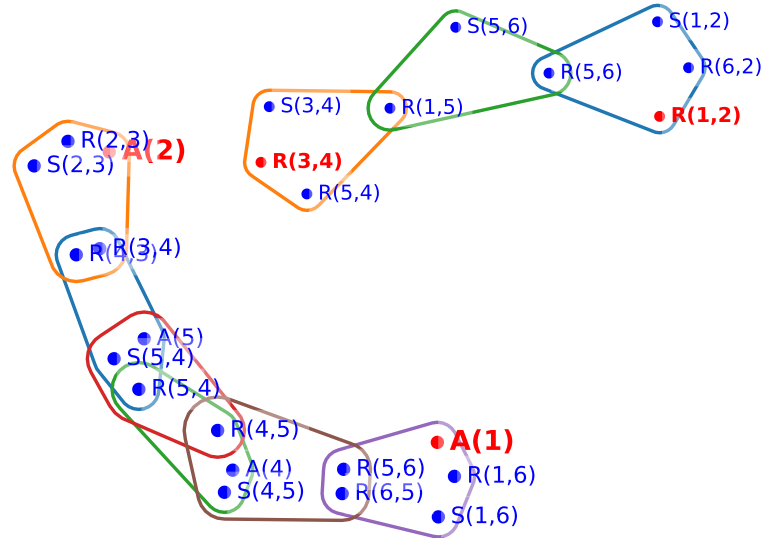




# 5 New Hardness Gadgets

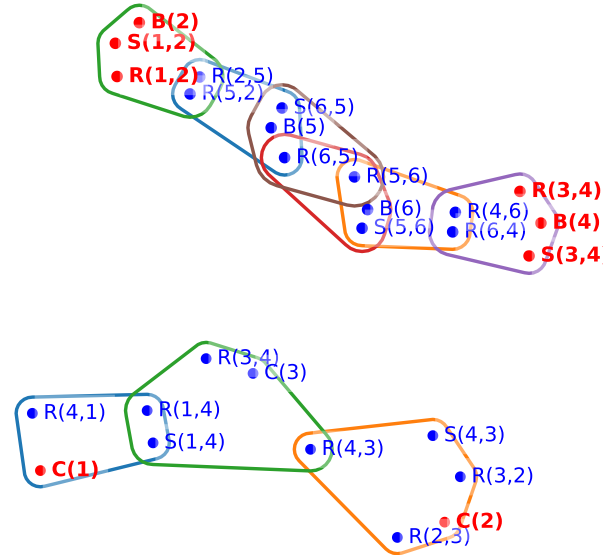
- Using the Automatic Hardness Finder, we proved 5 queries hard (out of 7 previously **open** from Freire+20)

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

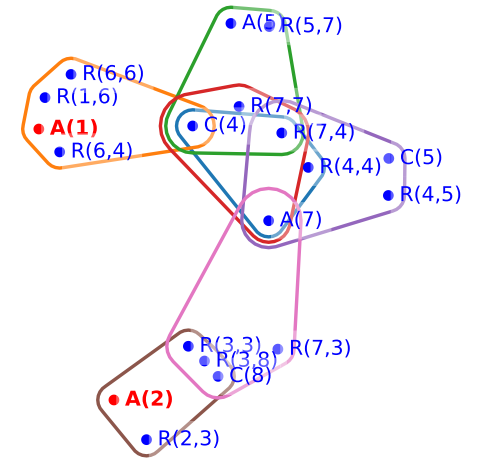


$$q_{3perm-R}^{AS_{xy}} :- A(x), S(x,y), R(x,y), R(y,z), R(z,y)$$

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



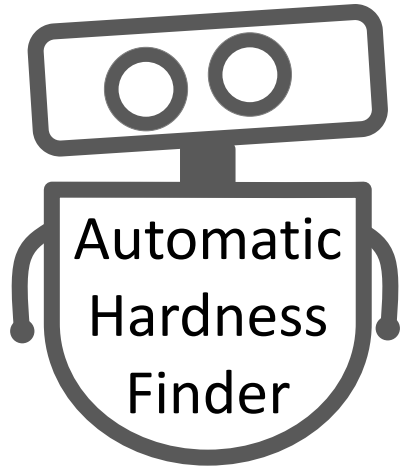
$$q_{3perm-R}^{S_{xy}C} :- S(x,y), R(x,y), R(y,z), R(z,y), C(z)$$



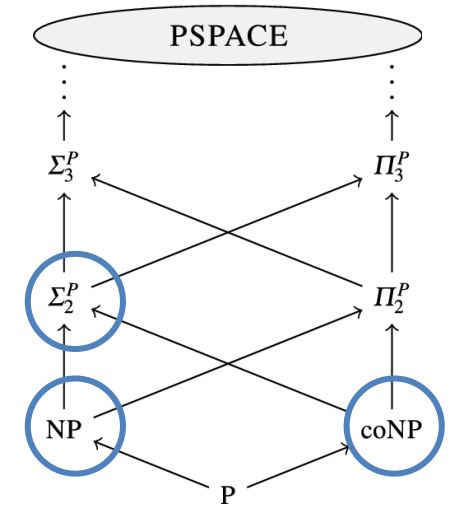
$$z_6 :- A(x), R(x,y), R(y,y), R(y,z), C(z)$$

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

# How to Build the Automatic Hardness Finder?




1. Declarative Program
2. Semantic Specification of “Hardness”
  - We show 5 properties that are **sufficient** to show hardness
  - These properties can be easily tested
3. Solve NP-Hard problems as a sub-routine
  - Expressivity used:  $\Sigma_2^P$
  - Disjunctive Logic Programs



- Unified Reverse Data Management Framework
- Insight #1: How to build **Unified Algorithms**
- Insight #2: How to build **Automatic Hardness Provers**
- **What else is in the paper?**
- Takeaways + Open Questions

# What else is in the paper?

- Unified Algorithms for Resilience and Casual Responsibility
  - Automatic Hardness Finder
  - Complexity Dichotomy for Resilience and Causal Responsibility under Bag semantics
  - More tractable cases:
    - Read-Once Instances, Functional Dependencies...
  - Approximation algorithms
  - Experimental Verification
- 
- In this talk*

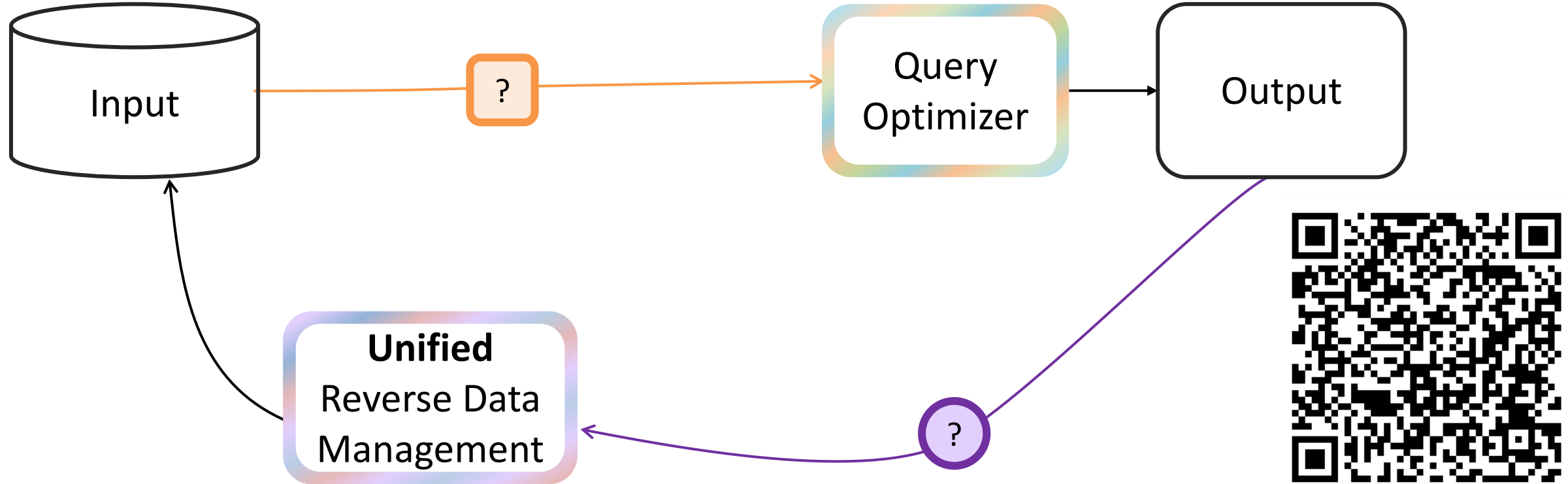
# Takeaways

- 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
  - Minimal Factorization of Provenance of CQs (**at PODS 5pm tomorrow!**)
  - Causal Responsibility
  - ..... **Claim: many more: Deletion Propagation, Algorithmic Fairness, ...**
- Uncovering more complexity results across different problems
- Build a “Reverse Query Optimizer”

# Takeaways



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, SIGMOD 2024
- Makhija, Gatterbauer. Minimally Factorizing the Provenance of Self-join Free Conjunctive Queries, PODS 2024

# Appendix

# 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”

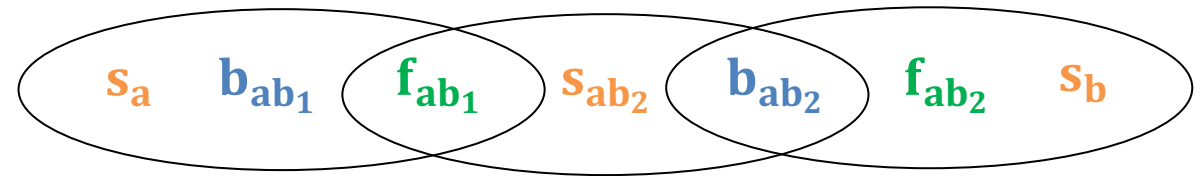
4. OR – property

5. Composability

“Key” properties:

- Semantically defined

- I will just show intuition

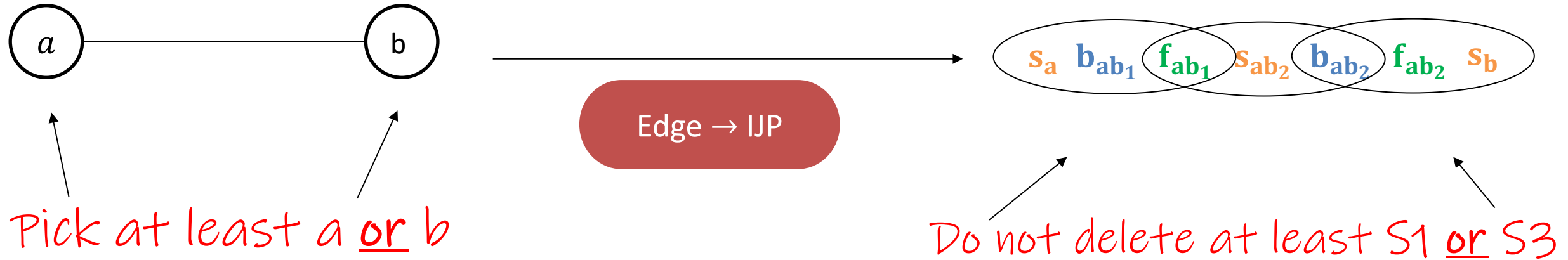


Data Hypergraph



# Unified Approach to prove Hardness

## Key Property #1: OR property



# Unified Approach to prove Hardness

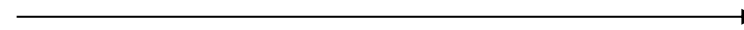
## Key Property #1: OR property



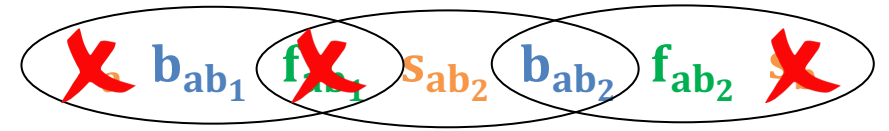
Pick at least a or b

VC = ??

Violates constraints



Edge  $\rightarrow$  IJP



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

RES = 3

Violates minimality

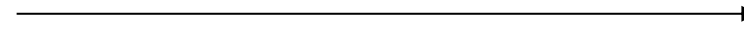
# Unified Approach to prove Hardness

## Key Property #1: OR property

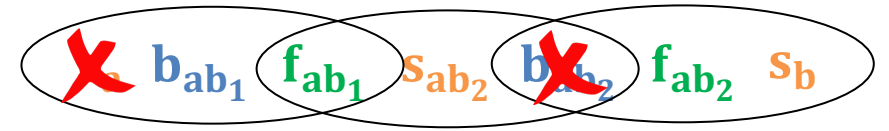


Pick at least a or b

VC = 1



Edge  $\rightarrow$  IJP



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

RES = 2

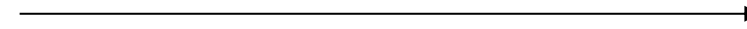
# Unified Approach to prove Hardness

## Key Property #1: OR property

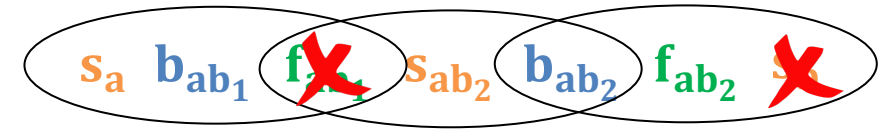


Pick at least a or b

VC = 1



Edge  $\rightarrow$  IJP



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

RES = 2

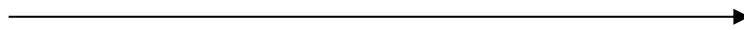
# Unified Approach to prove Hardness

## Key Property #1: OR property

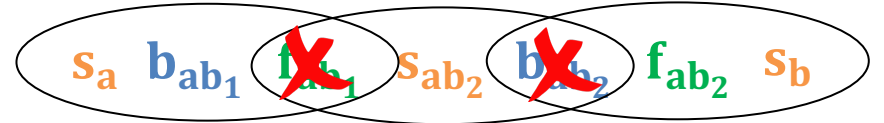


Pick at least a or b

VC = 2



Edge  $\rightarrow$  IJP

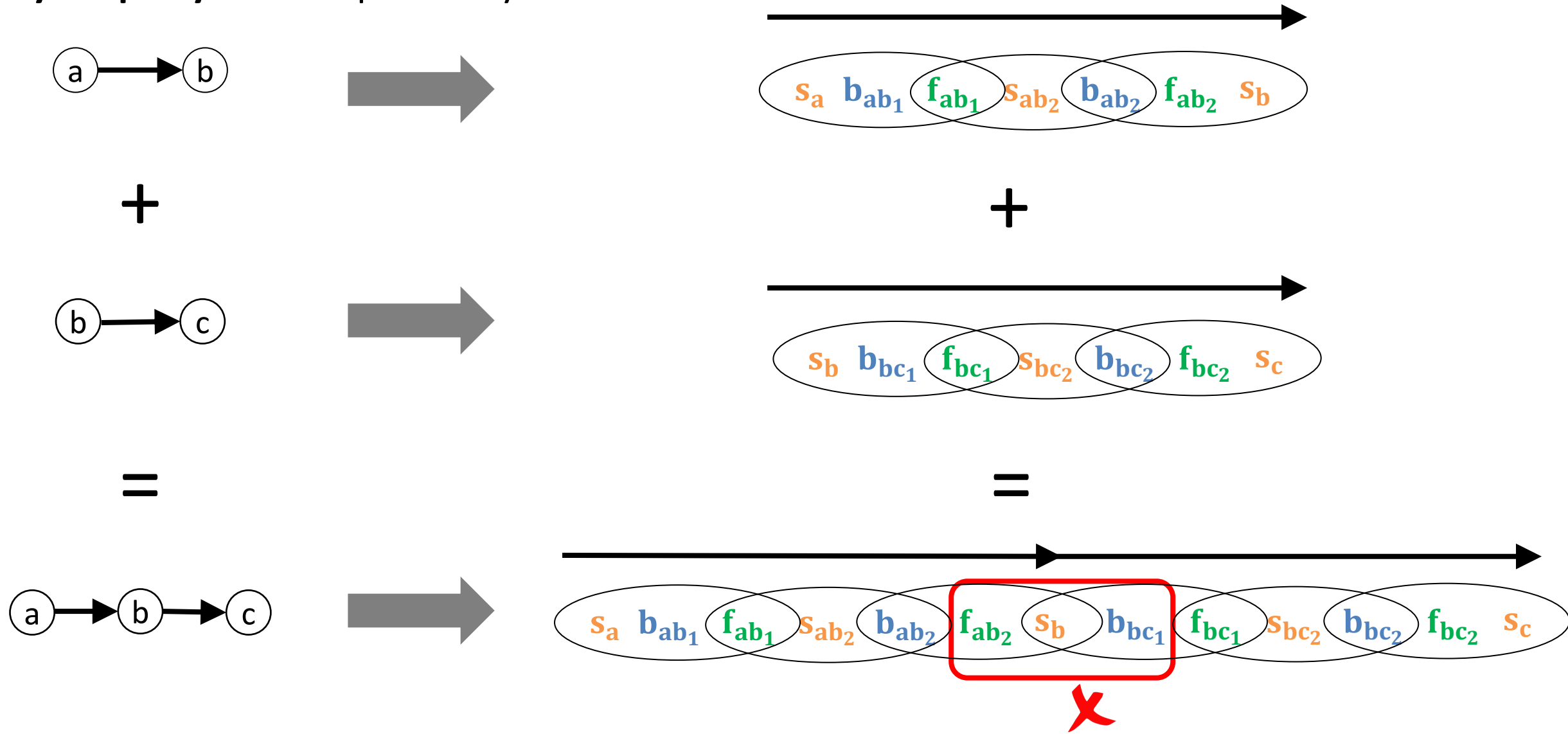


Do not delete at least  $S_1$  or  $S_3$

RES = 2

# Unified Approach to prove Hardness

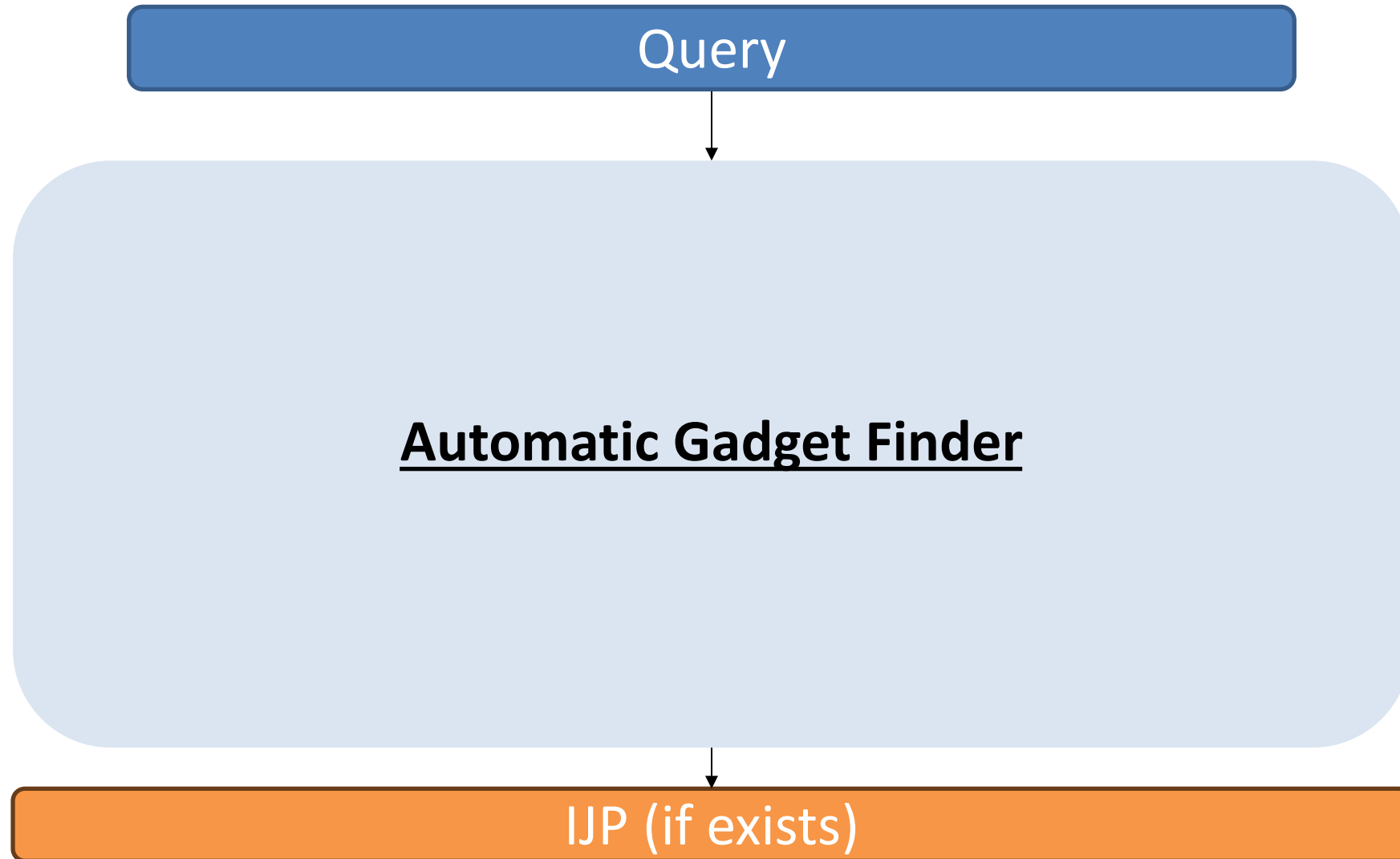
## Key Property #2: Composability of IJPs



Composing two IJPs should not lead to additional witnesses

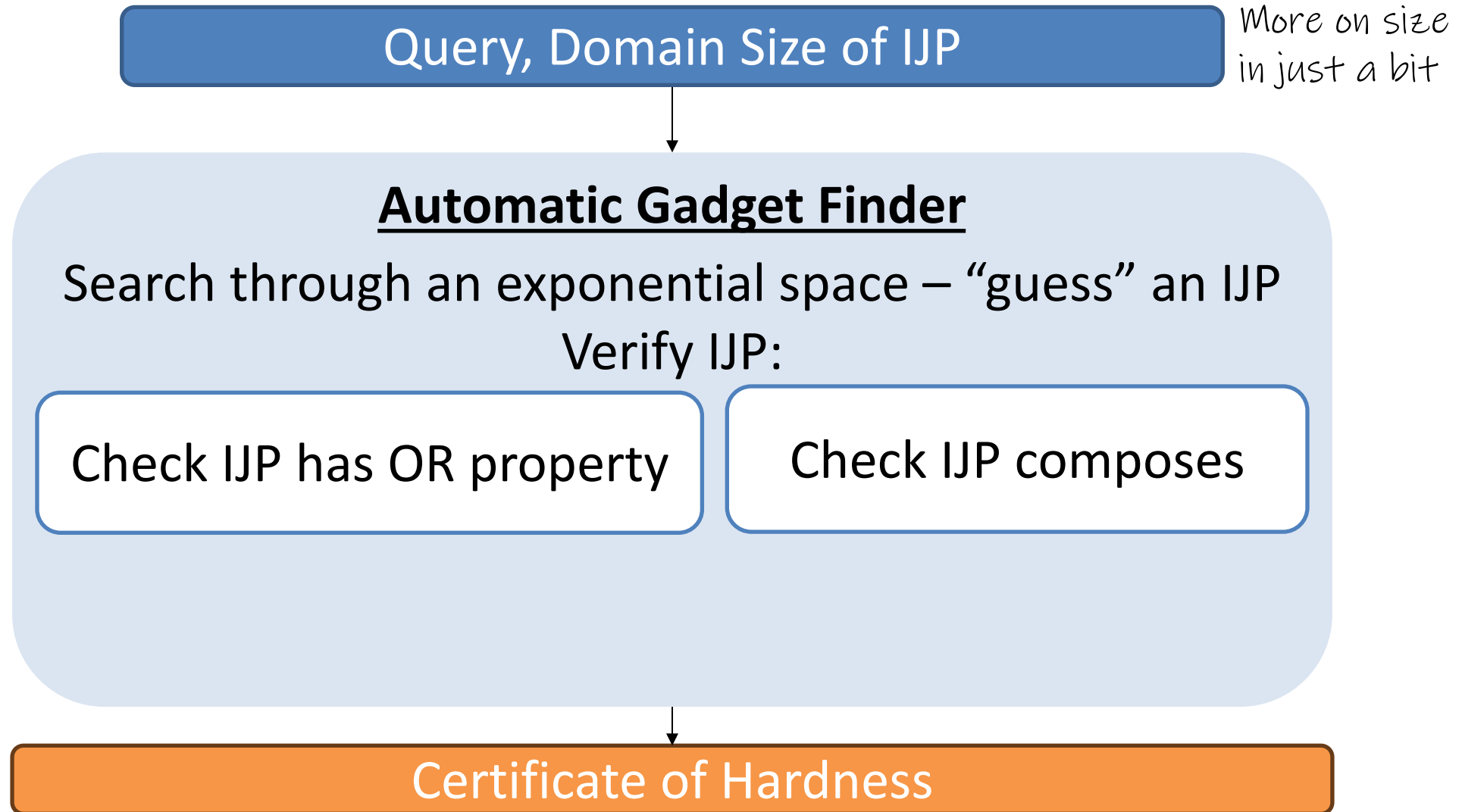
# Our Goal

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

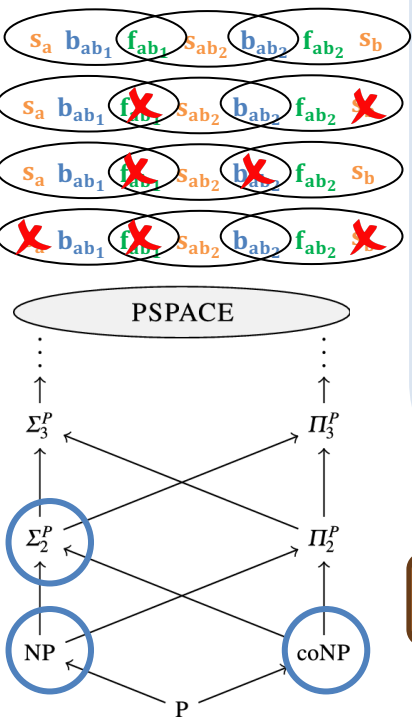
Verify IJP:

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

