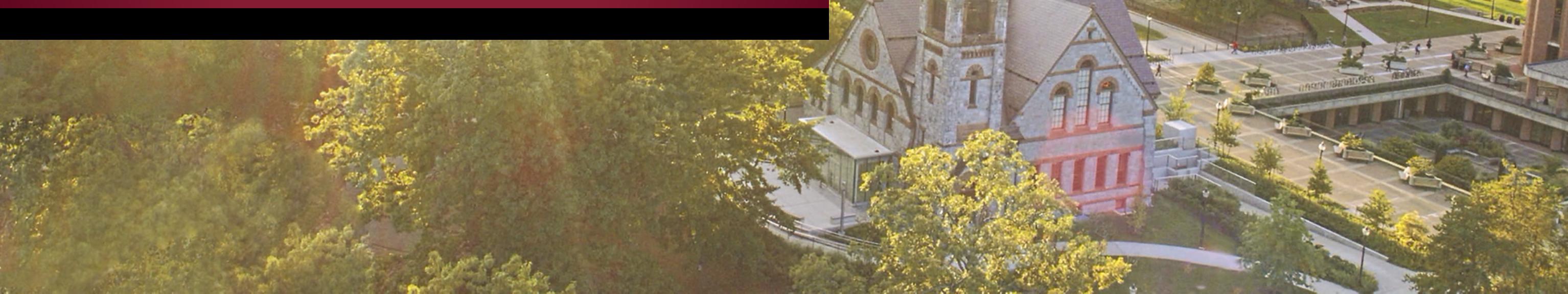


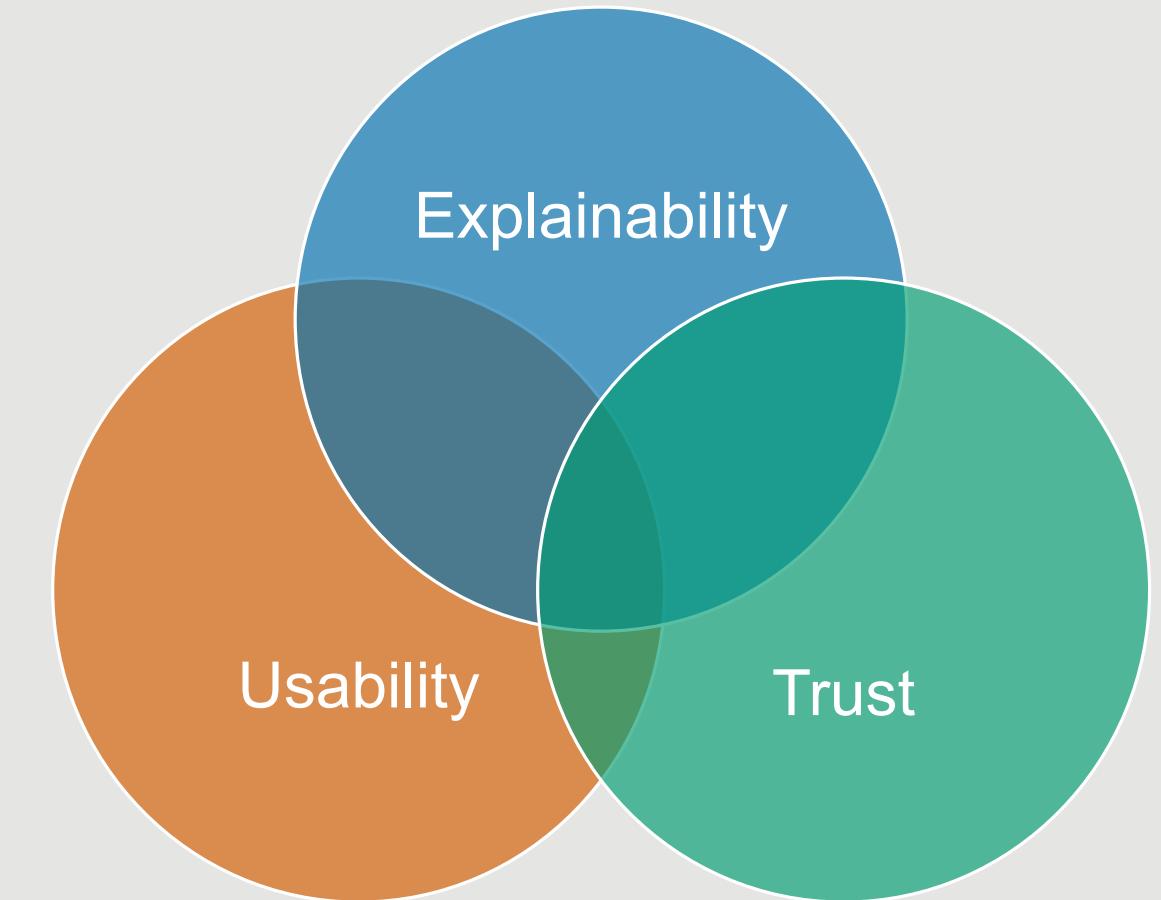
Enhancing Usability and Explainability of Data Systems

Anna Fariha

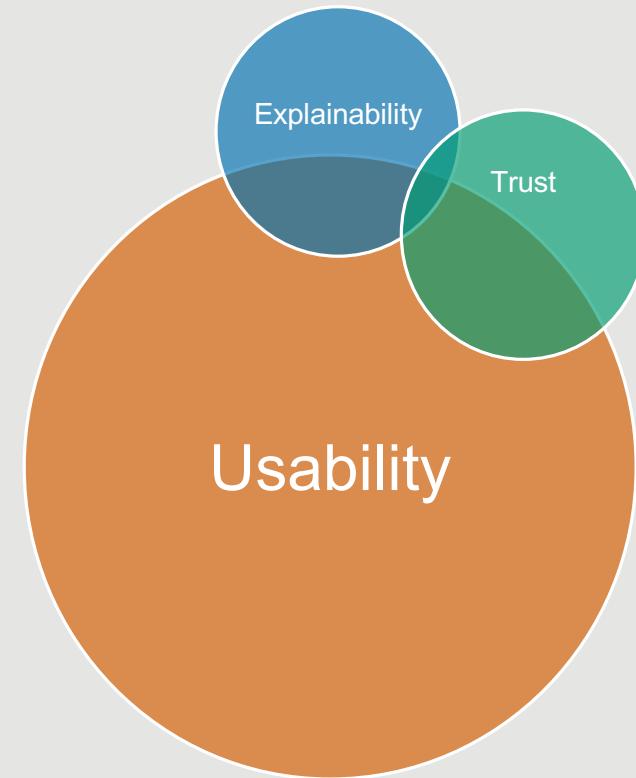
Advisor: Alexandra Meliou



Democratization of data systems



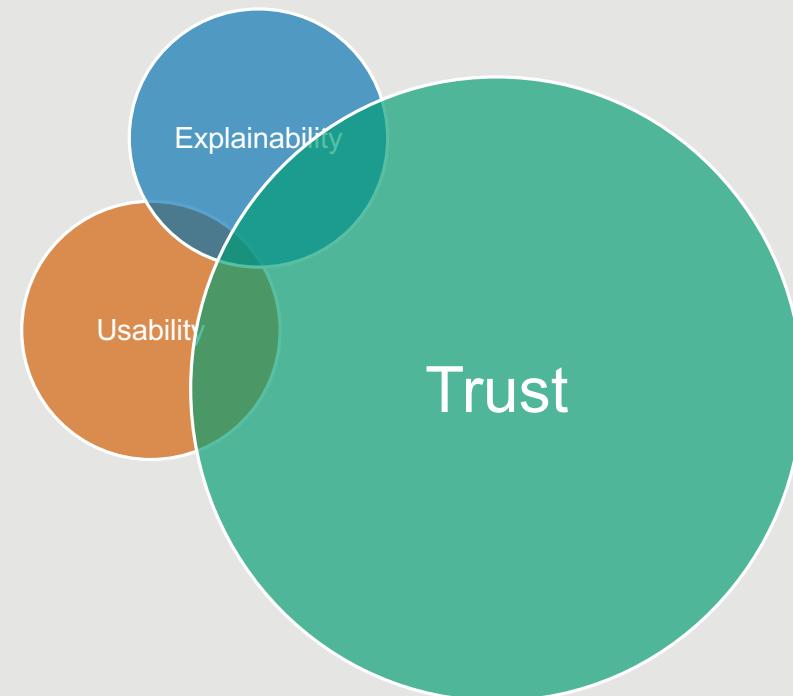
Usability



Makes data systems **accessible** to non-expert users.

- Applications
 - Data access
 - Querying relational databases
 - Data integration
 - Data transformation
 - Data visualization
 - Data summarization
 - Text document summarization

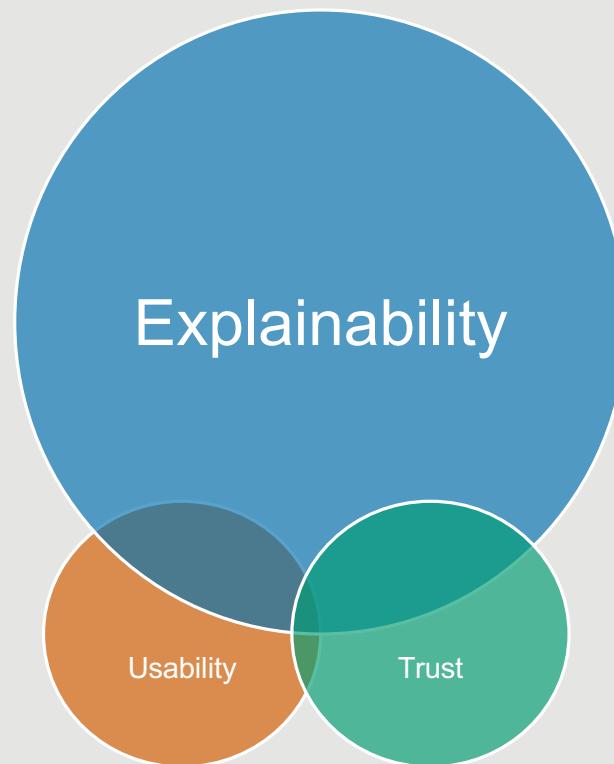
Trust



Enhances people's **confidence** towards data systems.

- Applications
 - Artificial intelligence and machine learning
 - Model predictions
 - Novel interaction mechanisms
 - Programming by example

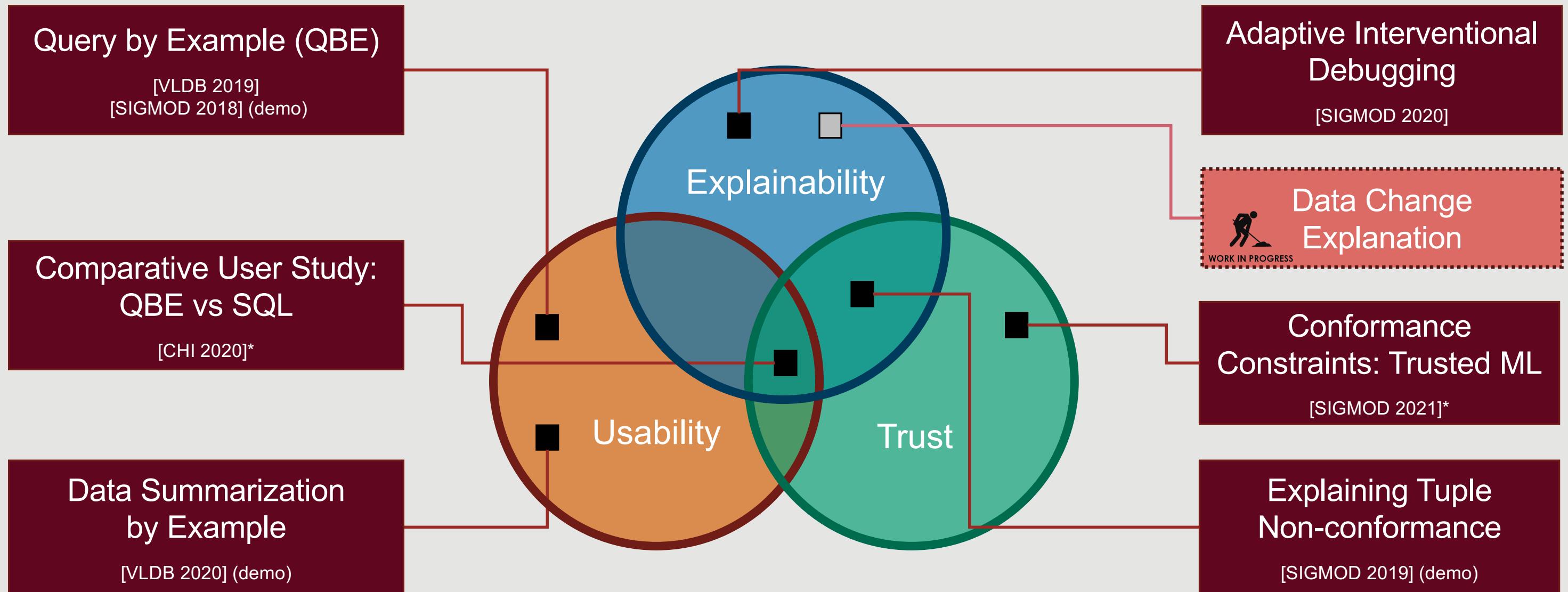
Explainability



Increases transparency of data systems.

- **Applications**
 - Machine learning
 - Model predictions
 - Distributed systems
 - Concurrent applications
 - Data evolution
 - Why/how two databases differ?
 - Fairness in algorithms/software

Dissertation outline



Part 1: Usability of Data Systems



Usability

Are data systems accessible to non-experts?



Who are our most
valuable customers?

How to express complex task specifications?



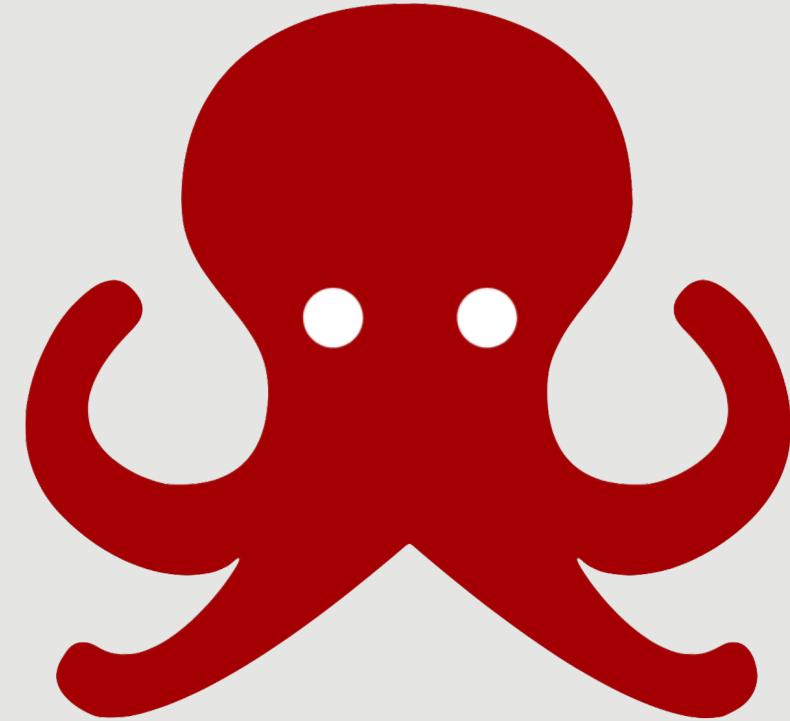
Programming by example (PBE)



- A step towards democratization of computational power.
- Enhances usability for both non-experts and experts.



Querying relational databases by example



SQuID

Semantic similarity-aware Query Intent Discovery

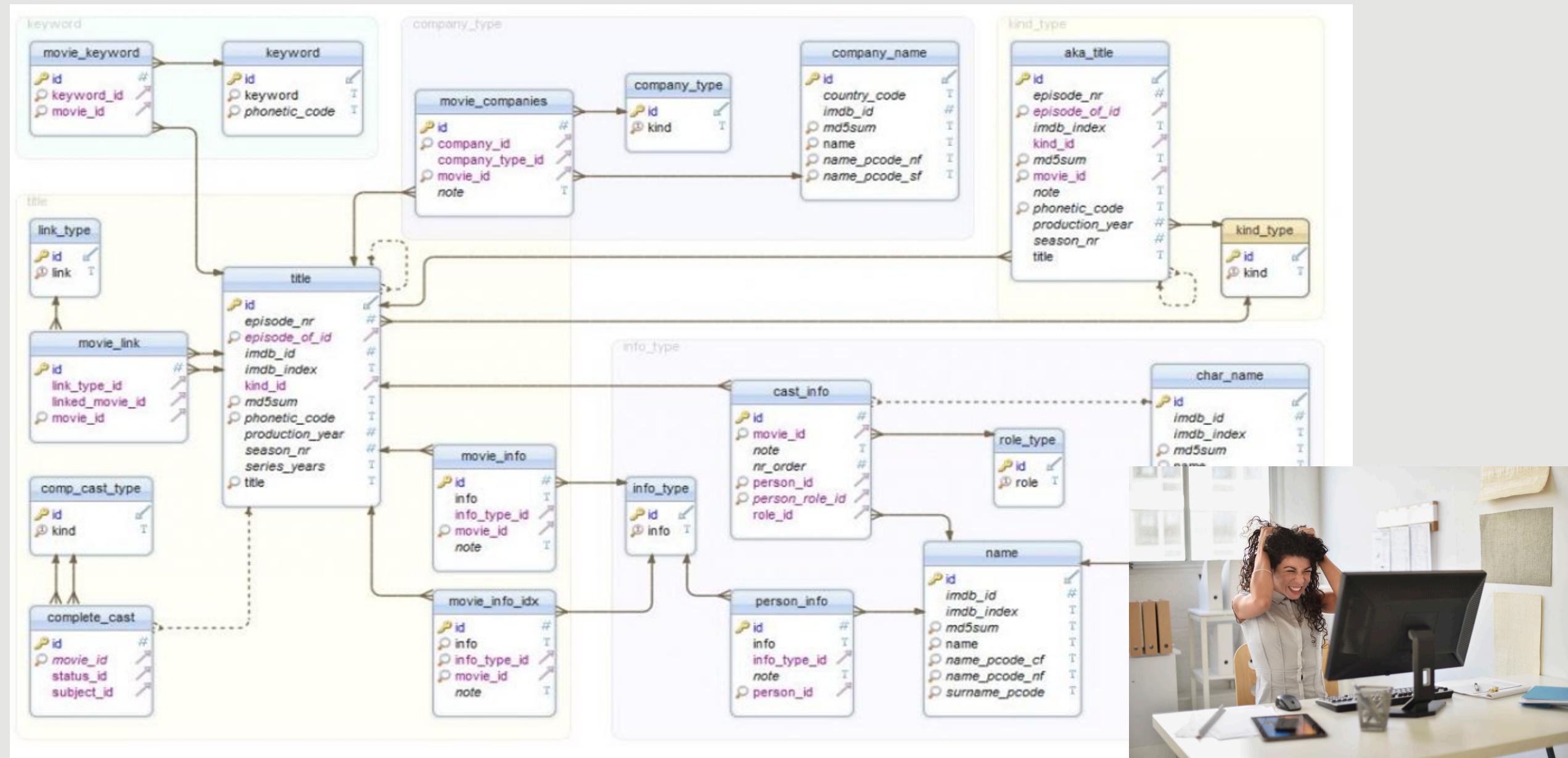
- Alice wants to find all Funny Actors from the IMDb database.



The page shows the IMDb homepage with the following elements:

- Banner:** IN THEATERS SEPTEMBER 13, The GOLDFINCH, THE STORY OF A STOLEN LIFE, WATCH TRAILER.
- Search Bar:** Find Movies, TV shows, Celebrities and more... (with a magnifying glass icon).
- Navigation:** All, Watch Now For Free, Movies, TV & Showtimes, Celebs, Events & Photos, News & Community, Watchlist, Sign in.
- Movie Trailers:** It's a wonderful life, Henry Golding, GHOSTBUSTERS 2020, WHAT WE KNOW SO FAR, Last Christmas Official Trailer, Who Will Answer the Call?, 'Ghostbusters 2020' So Far, Marvel Paves the Path to 'Shang-Chi' Trace the MCU Clues.
- Upcoming Releases:** IN THEATERS SEPTEMBER 13, The GOLDFINCH, WATCH TRAILER.
- Section:** Opening This Week.

Challenge 1: understanding the schema



Challenge 2: SQL expertise

```
SELECT person.name
FROM person, castinfo, movietogenre, genre
WHERE person.id = castinfo.person_id
    AND castinfo.movie_id = movietogenre.movie_id
    AND movietogenre.genre_id = genre.id
    AND genre.name = 'Comedy'
GROUP BY person.id
HAVING count(*) >= 40
```



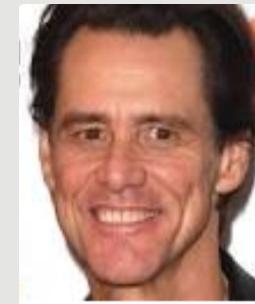
Query by example (QBE)



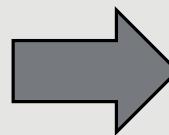
Eddie Murphy



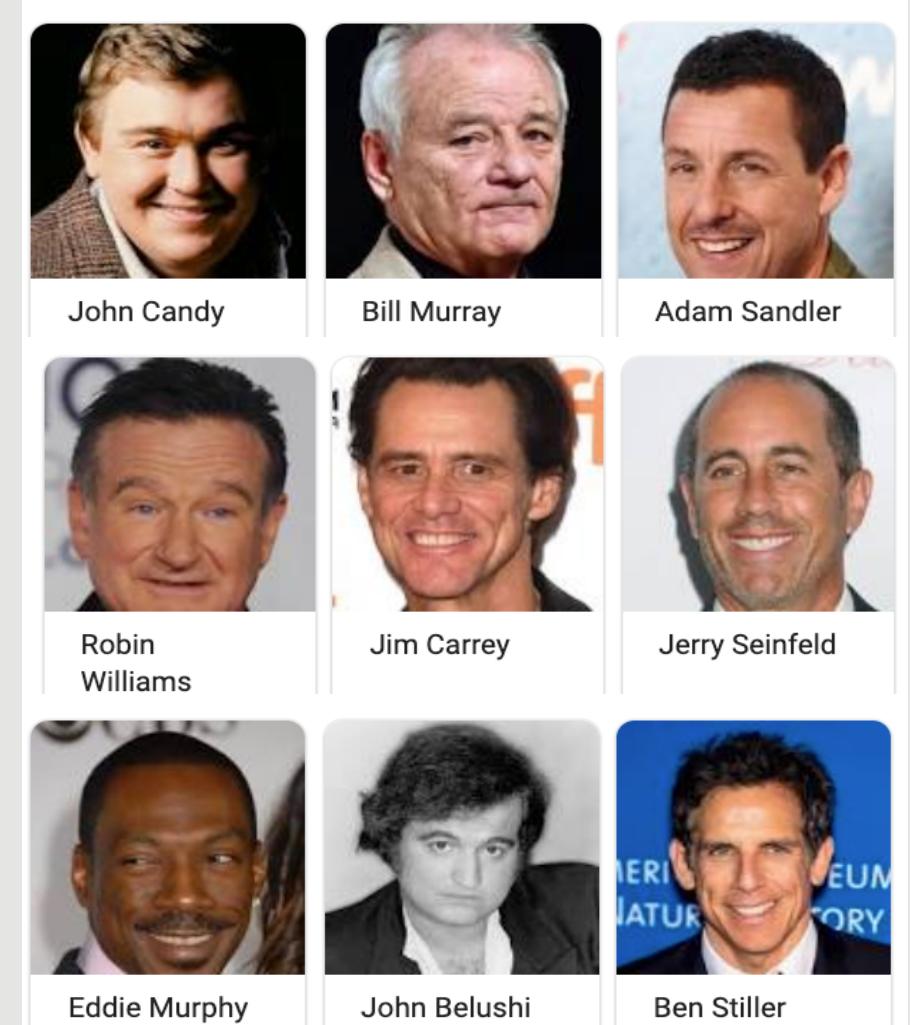
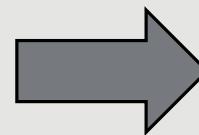
Robin Williams



Jim Carrey



QBE



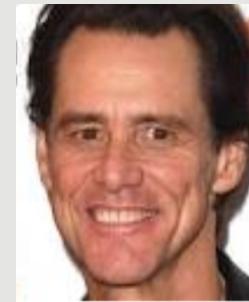
Expectation vs reality



Eddie Murphy



Robin Williams

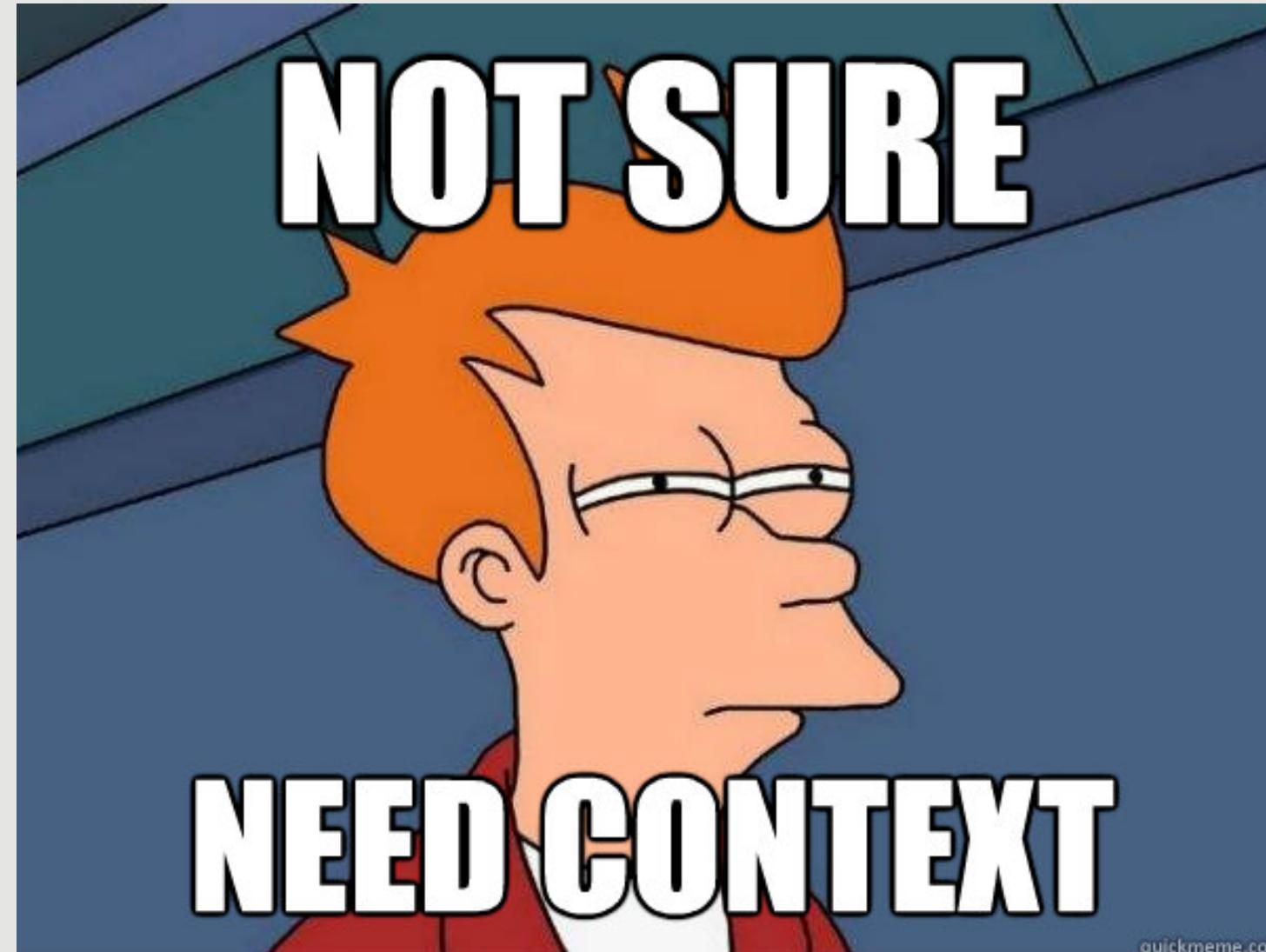


Jim Carrey

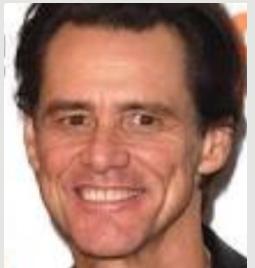


All actors

Humans use context



Discovering semantic similarity



Jim Carrey



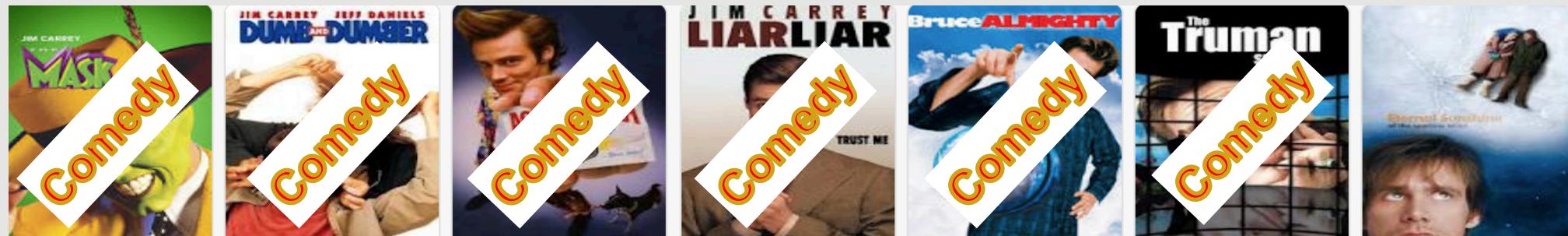
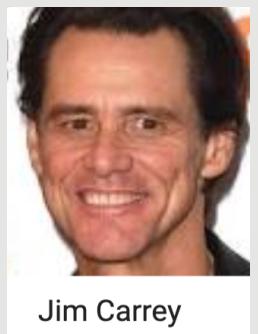
Robin
Williams



Eddie Murphy

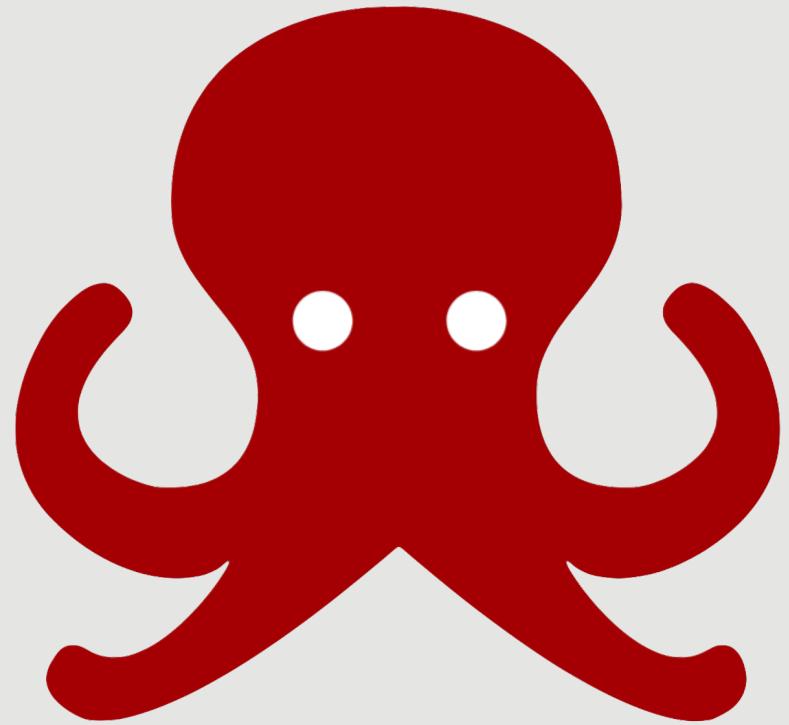
There is no “funny” attribute in the data

Discovering semantic similarity



SQuID

**Semantic Similarity-aware
Query Intent Discovery**



SQuID Outline

Modeling
Semantic
Context

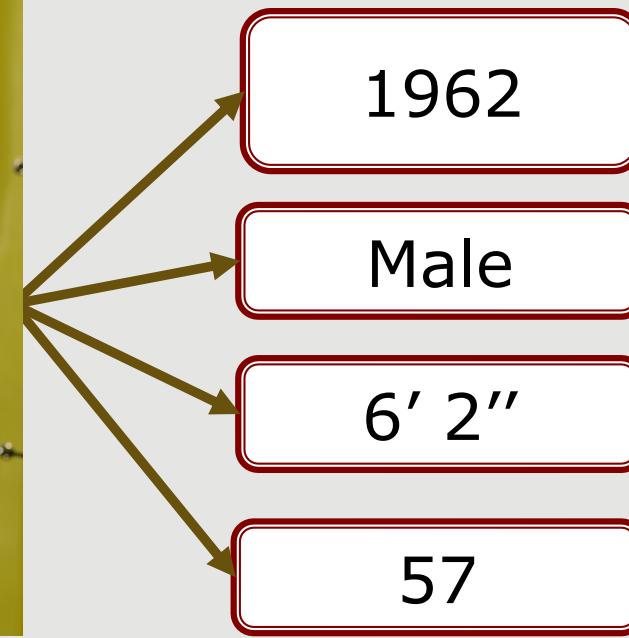
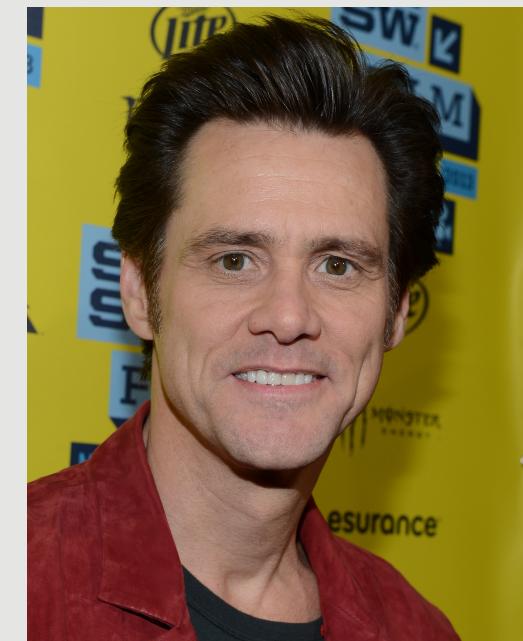
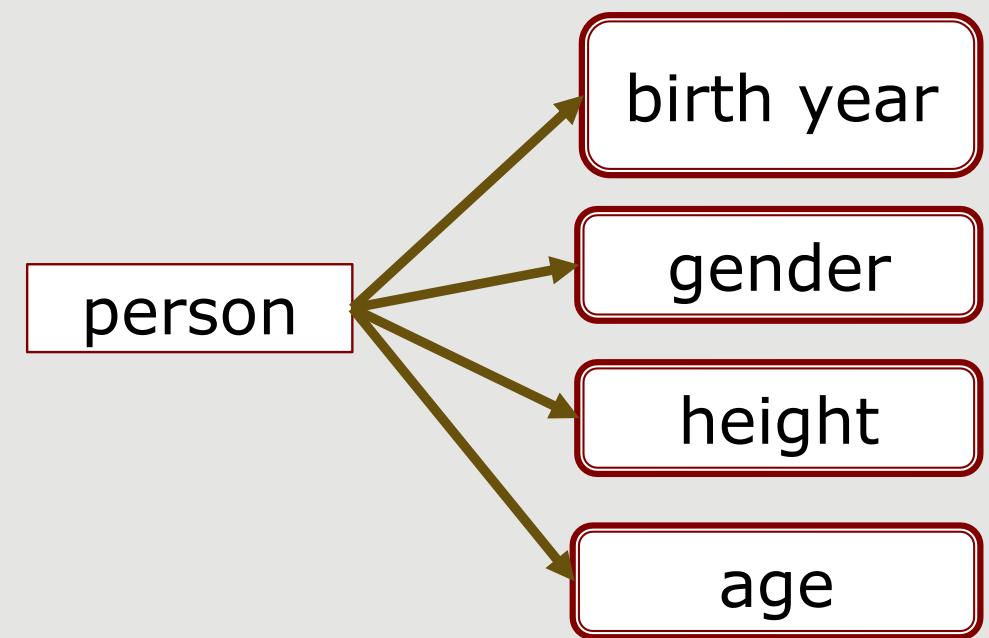
Query Intent
Discovery

Real-time
Performance

Evaluation

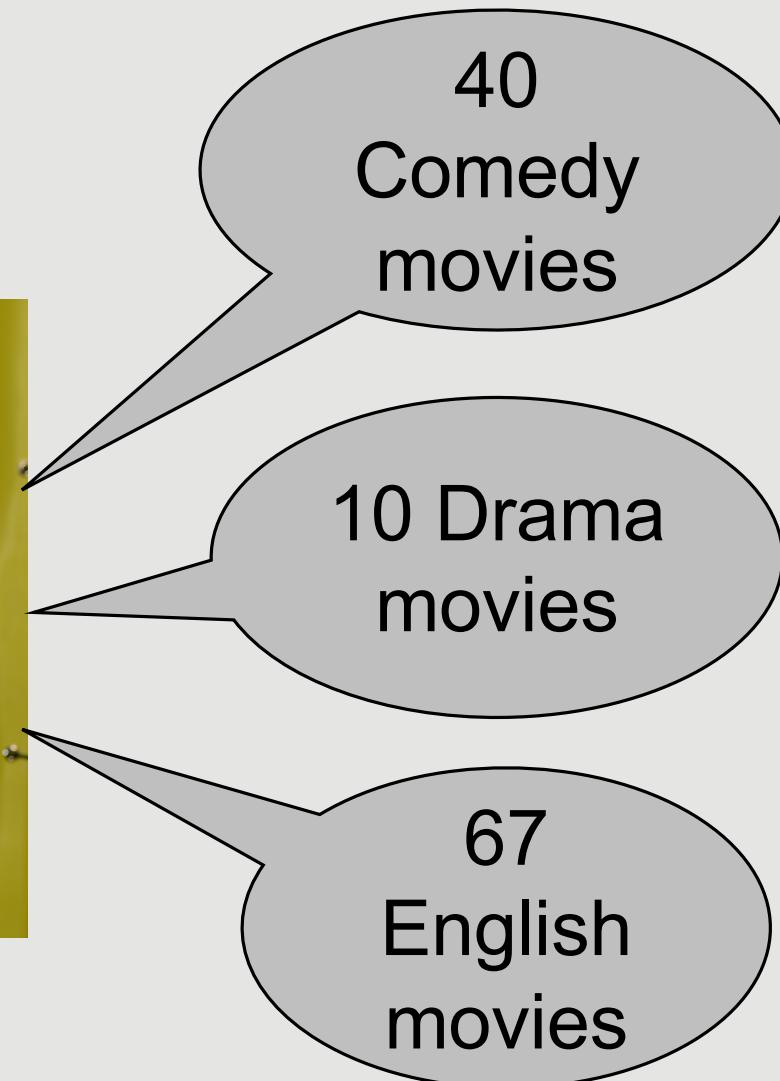
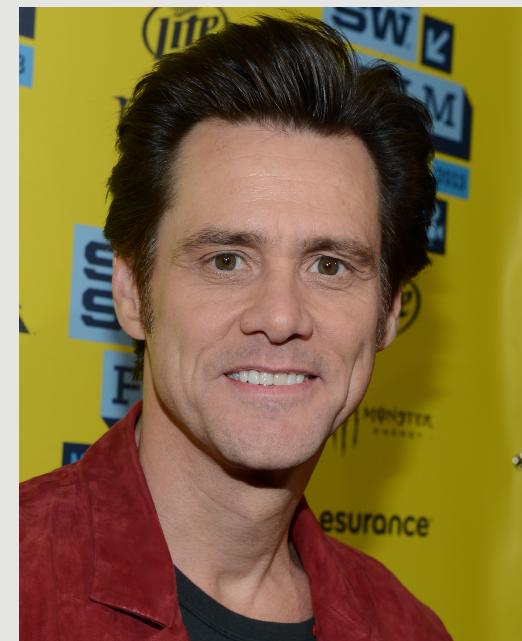
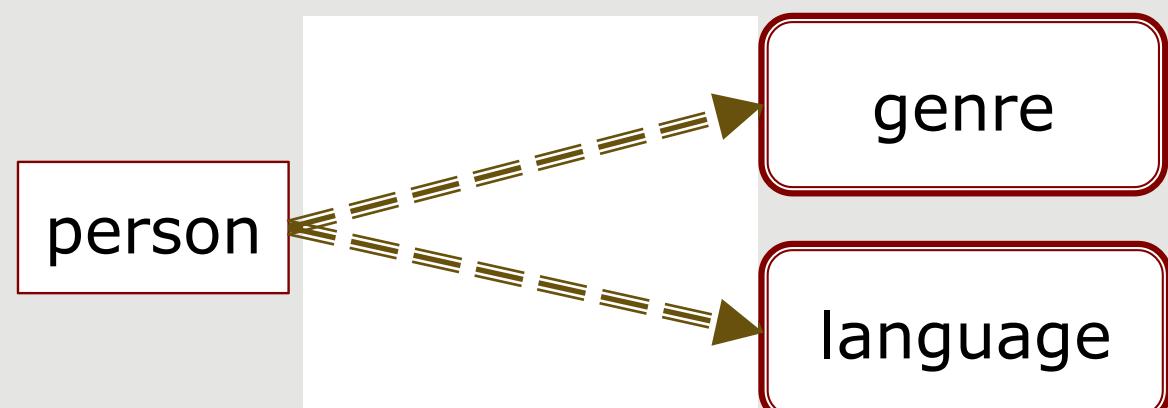
Semantic context: basic

- Directly affiliated with an entity.



Semantic context: derived

- **Aggregate over a basic property of an associated entity.**
 - number of comedy movies an actor appeared in.



Filters

- Encode semantic context.

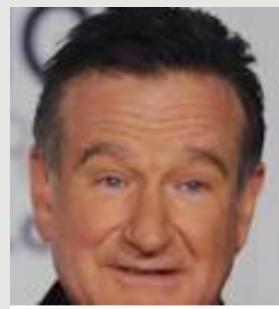
```
SELECT
    person
FROM
    people
WHERE
    color = orange
```



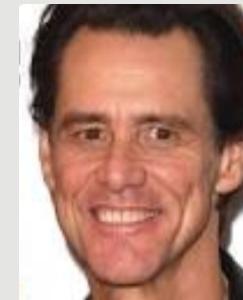
Intended or co-incidental?



Eddie Murphy



Robin Williams



Jim Carrey

- Male
- Born in North America
- Appeared in 80+ Hollywood movies
- **Appeared in 40+ comedy movies**
- Appeared in 20+ drama movies
- Height above 5 feet
- Born after 1940
- ...

Abduction

- Most likely **explanation** of an **observation**.
- Most likely **query** given the **examples**.

Maximum likelihood estimation is abduction!

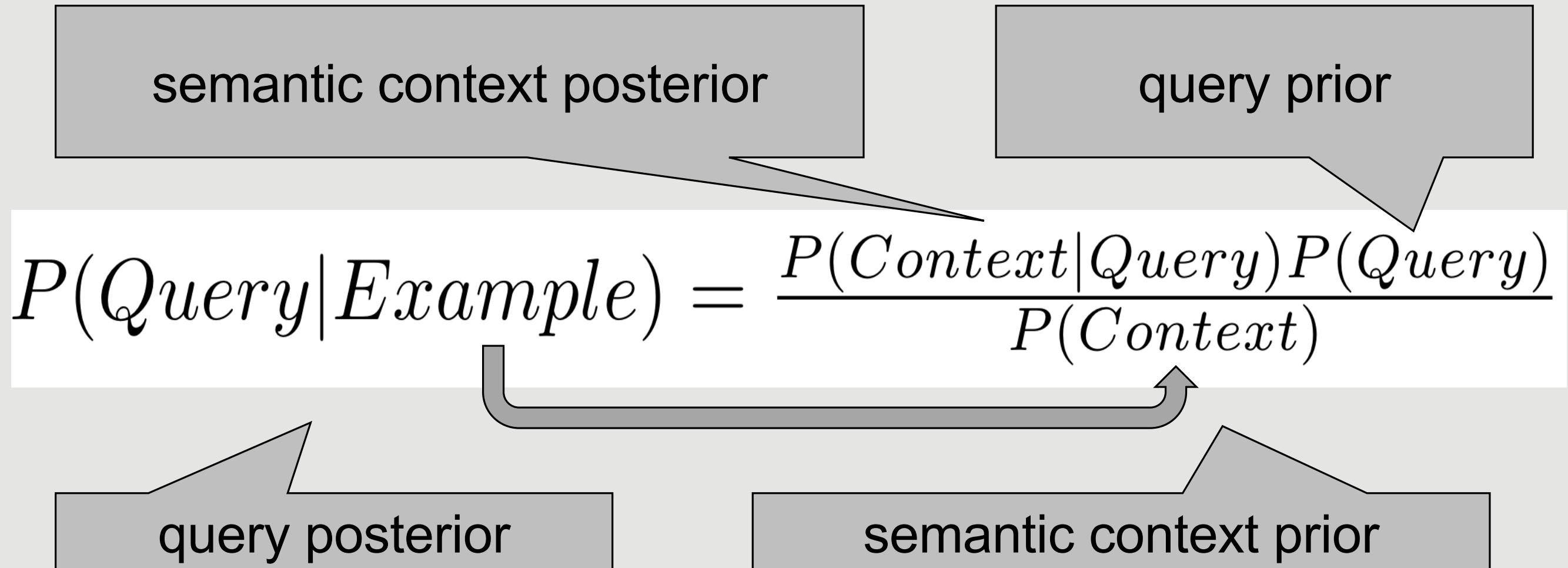
Problem definition

Query intent discovery: given a *Database* and *Example*, find *Query* such that:

$$\textit{Example} \subseteq \textit{Query}(\textit{Database})$$

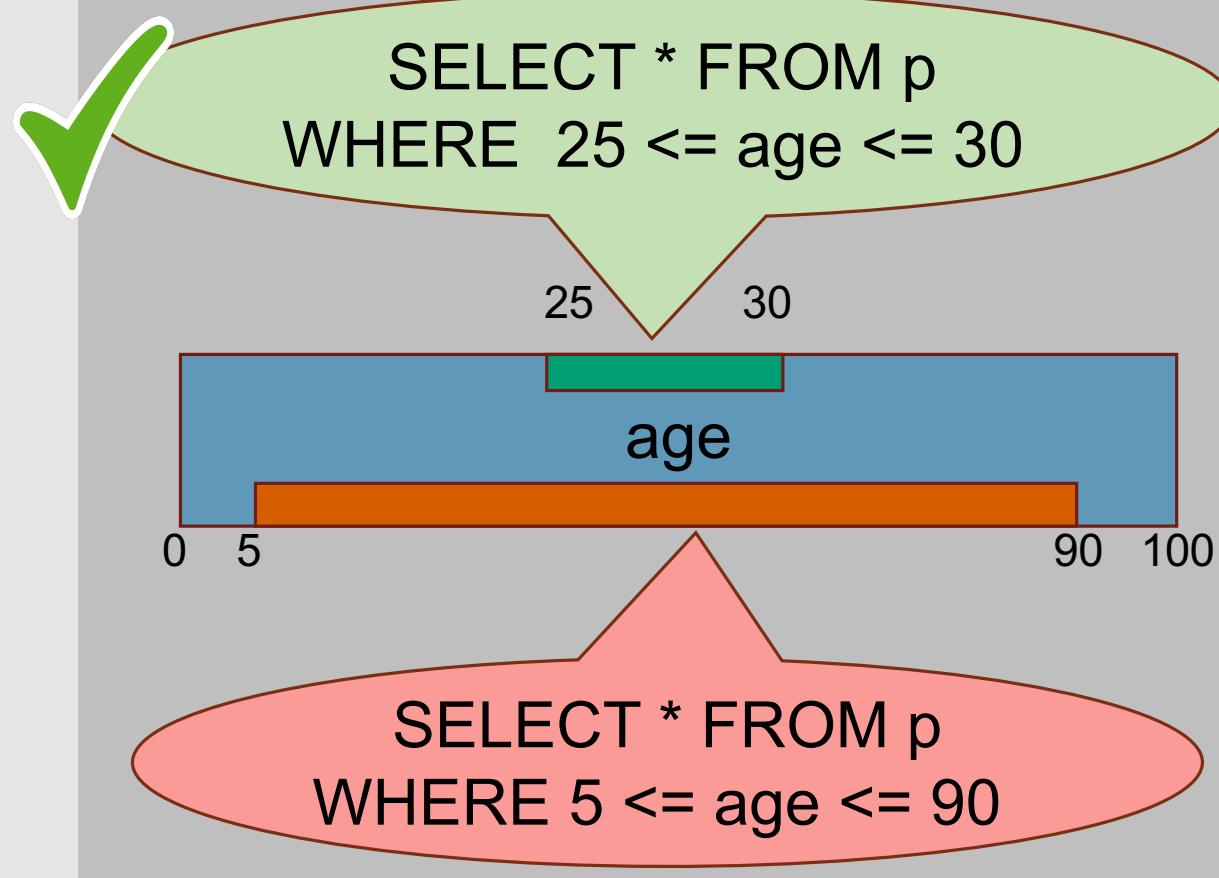
$$\textit{Query} = \arg \max_q P(q|\textit{Example})$$

Probabilistic abduction model

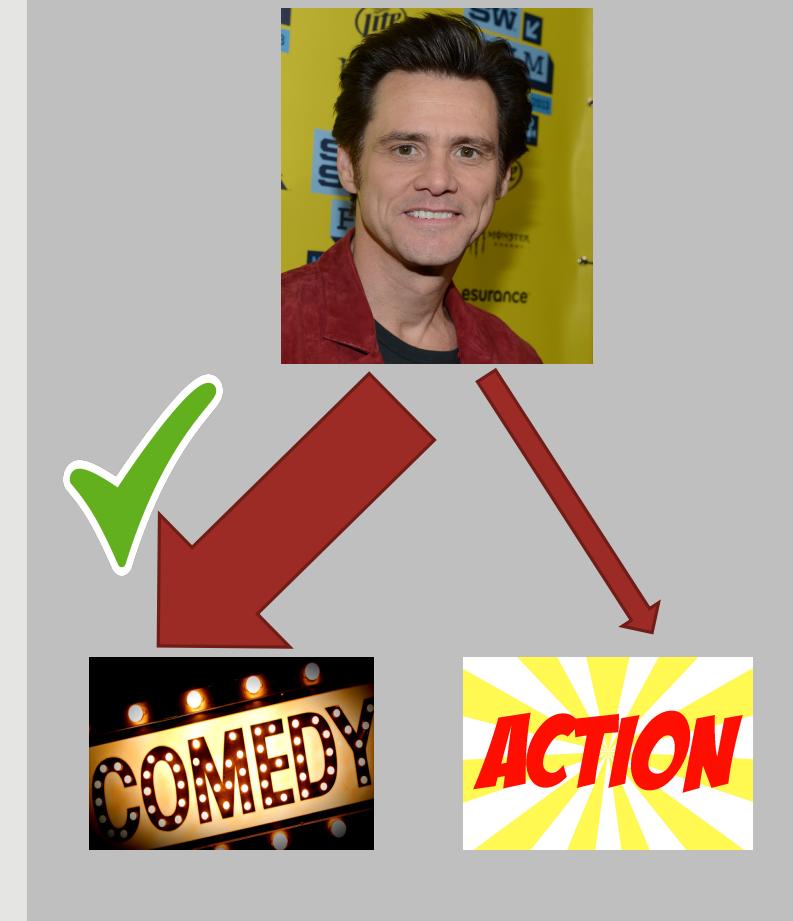


$$\frac{P(\text{Context}|\text{Query})P(\text{Query})}{P(\text{Context})}$$

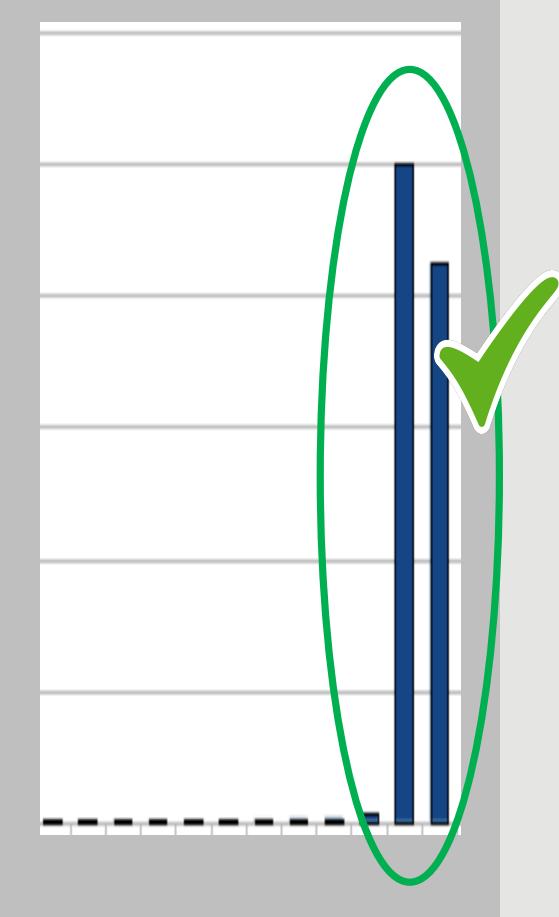
Domain selectivity



Association strength



Outlier



$$\frac{P(\text{Context}|\text{Query})P(\text{Query})}{P(\text{Context})}$$

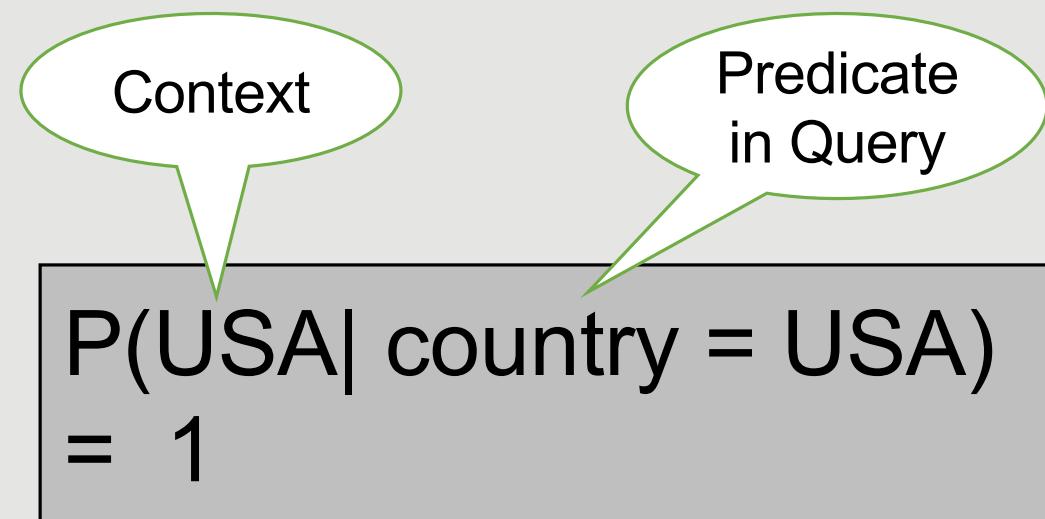
- Data selectivity

USA: 80%

	country
...	USA
...	CAN
...	USA
...	CAN
...	USA
...	USA

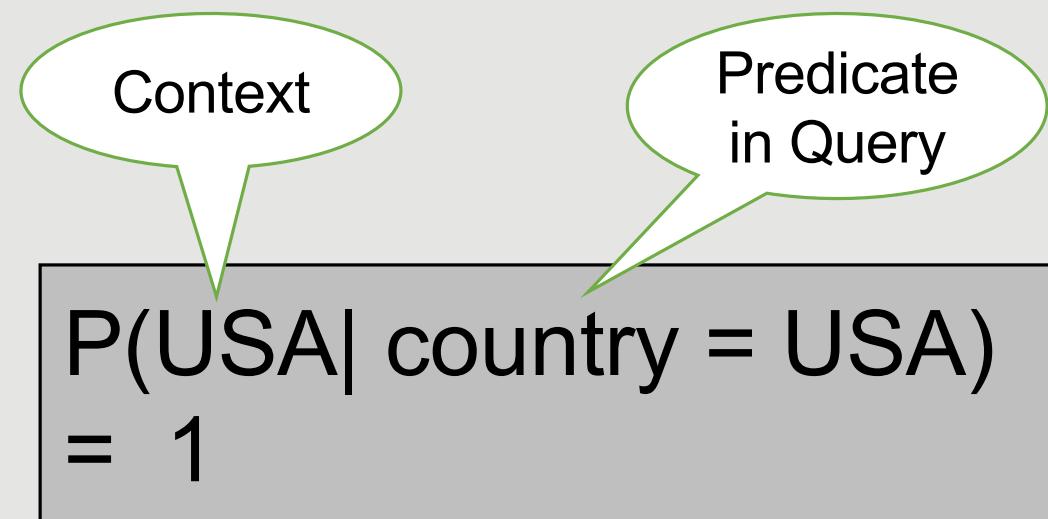
CAN: 20%

$$\frac{P(\text{Context}|\text{Query})P(\text{Query})}{P(\text{Context})}$$



	country
...	USA
...	CAN
...	USA
...	CAN
...	USA
...	USA

$$\frac{P(\text{Context}|\text{Query})P(\text{Query})}{P(\text{Context})}$$

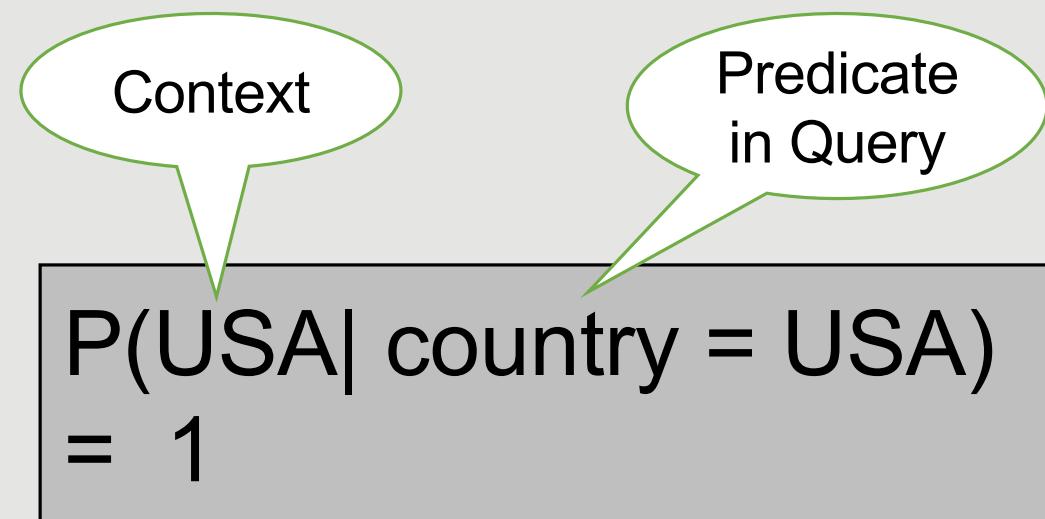


$P(\text{USA} | \text{No Filter}) = 0.8$

	country
...	USA
...	CAN
...	USA
...	CAN
...	USA
...	USA

	country
...	USA

$$\frac{P(\text{Context}|\text{Query})P(\text{Query})}{P(\text{Context})}$$

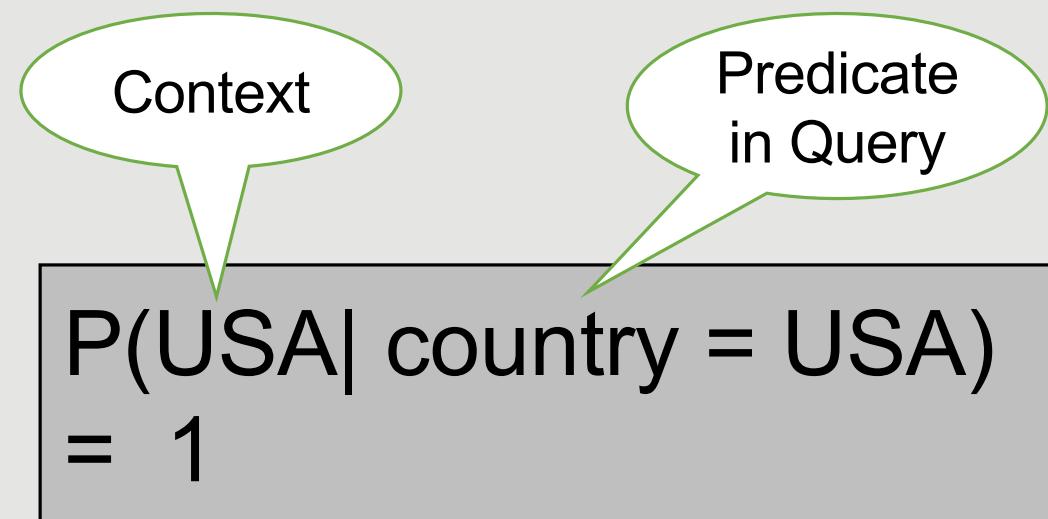


$$\begin{aligned} P(\text{USA} | \text{No Filter}) &= 0.8 * 0.8 \\ &= 0.64 \end{aligned}$$

	country
...	USA
...	CAN
...	USA
...	CAN
...	USA
...	USA

	country
...	USA
...	USA

$$\frac{P(\text{Context}|\text{Query})P(\text{Query})}{P(\text{Context})}$$

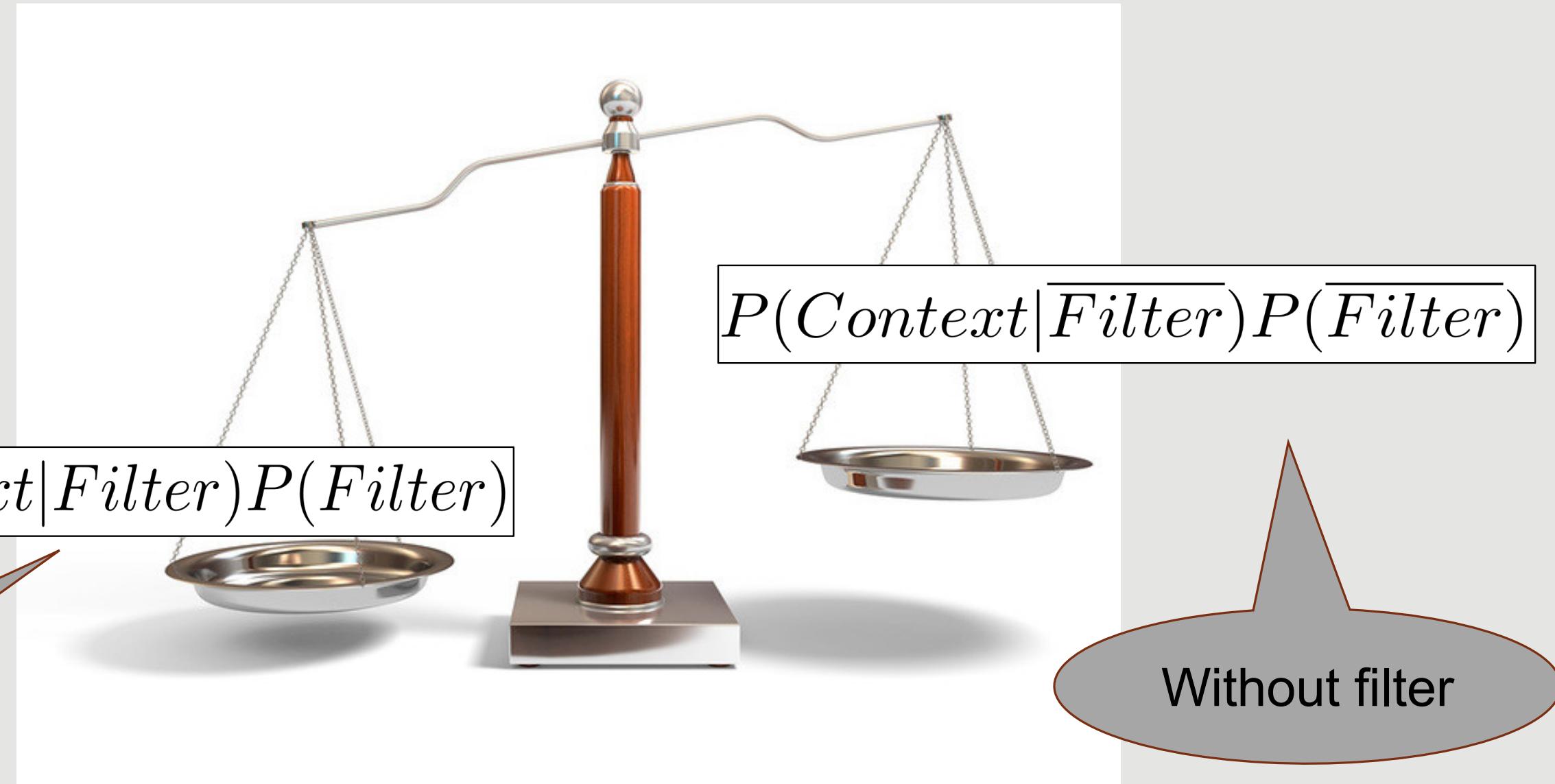


$P(\text{USA} | \text{No Filter})$
= $0.8 * 0.8 * 0.8$
= 0.51

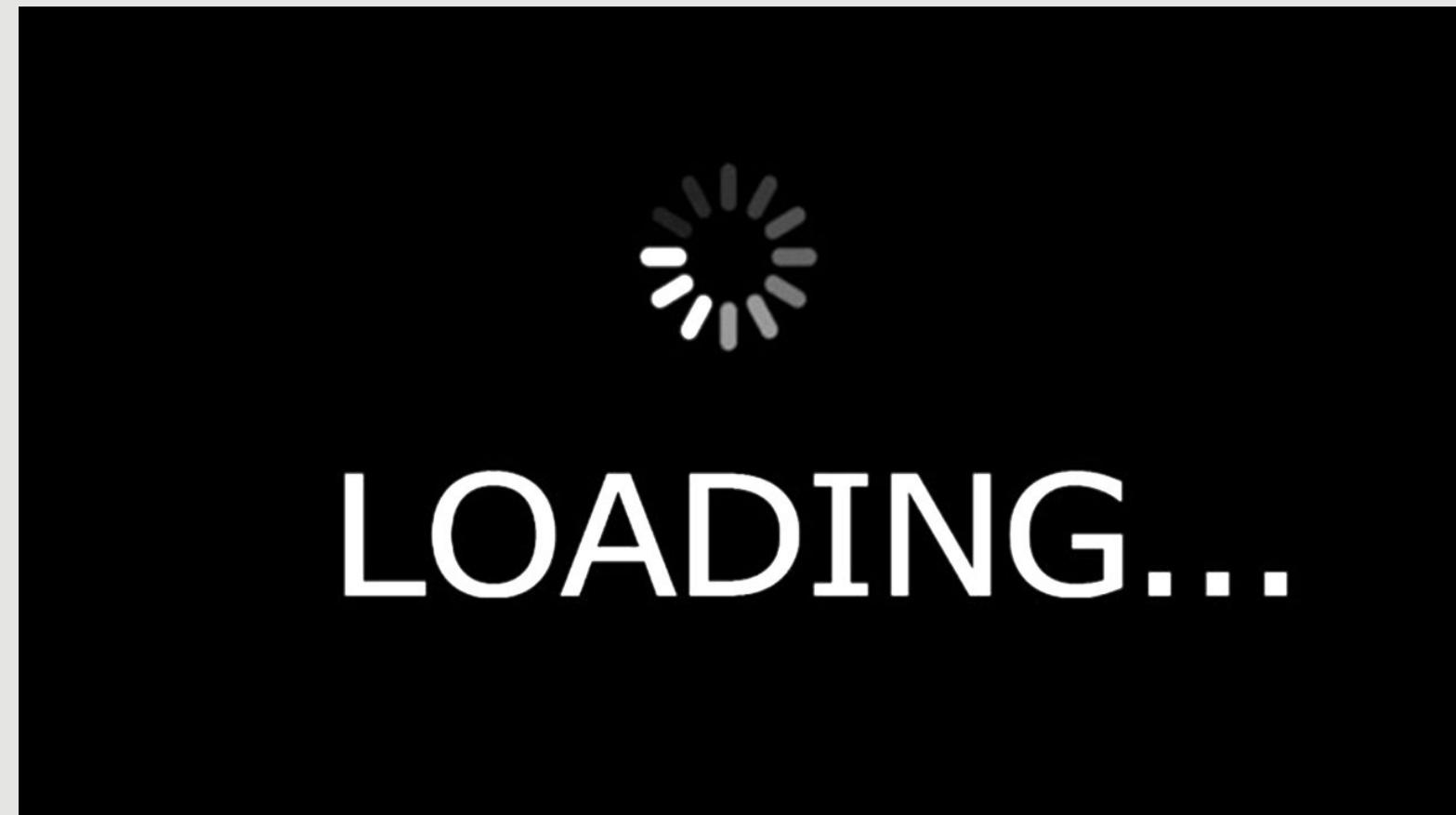
	country
...	USA
...	CAN
...	USA
...	CAN
...	USA
...	USA

	country
...	USA
...	USA
...	USA

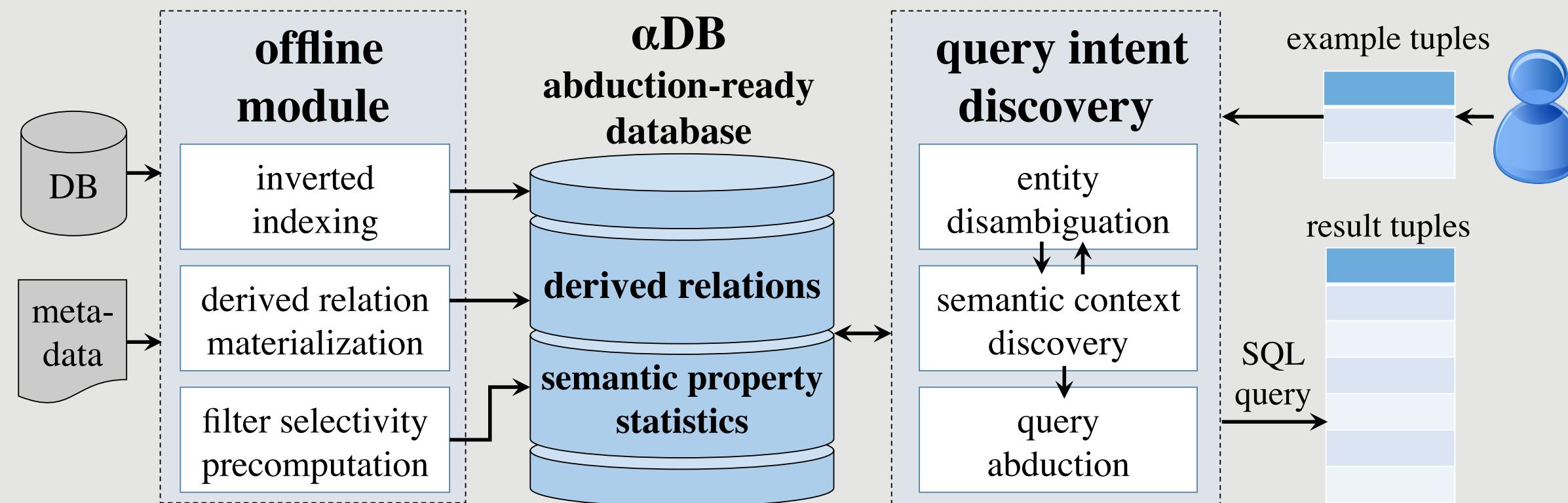
SQuID algorithm: to pick or drop filters?



Real-time performance



Abduction-ready database



Evaluation

1. How efficient is SQuID for large datasets and many examples?
2. Does SQuID infer the right query?
3. Can alternative techniques be effective in intent discovery?
 - Query Reverse Engineering (TALOS, 2014)
 - Positive and Unlabeled Learning (Elkan et al., 2008)
 - Query run-time comparison
 - Case studies

Datasets



633 MB

15 relations

- person: 6M rows
- movies: 1M rows
- castinfo: 14M rows

16 benchmark
queries



dblp

computer science bibliography

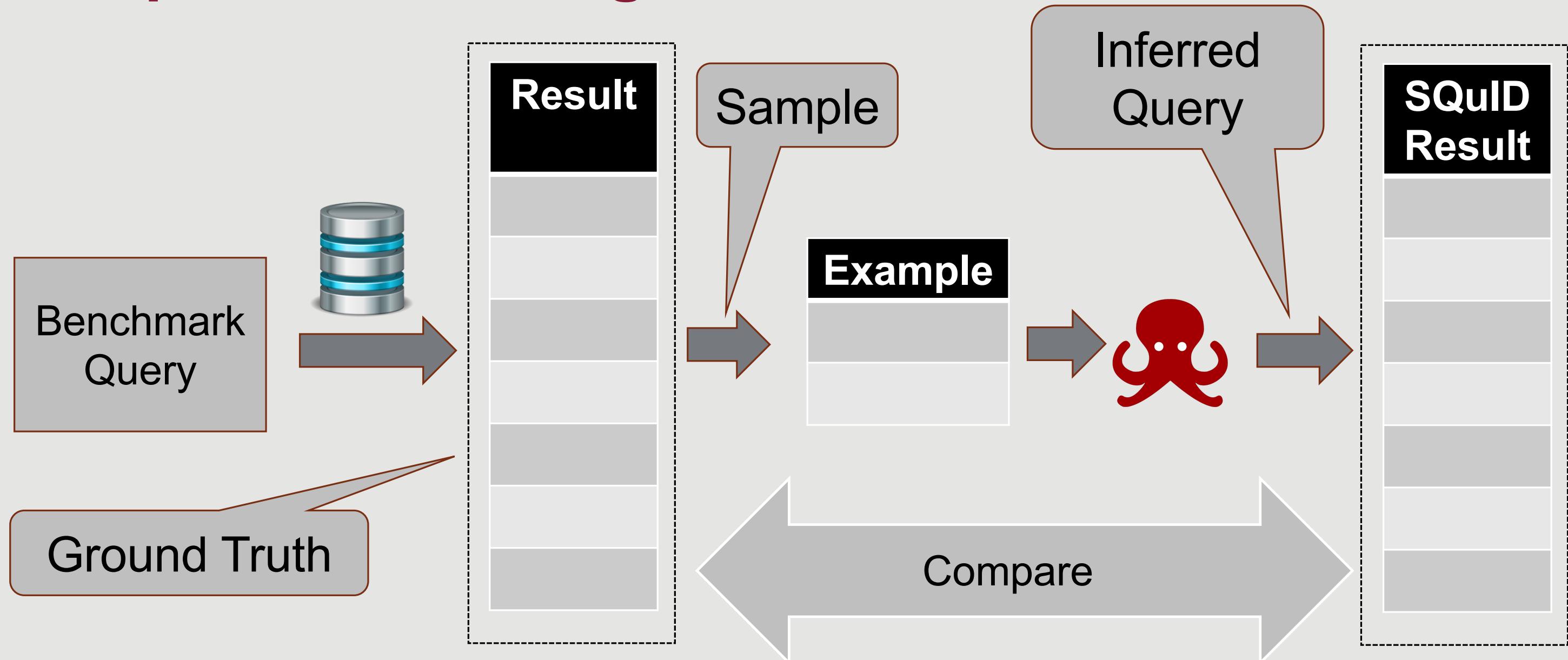


Machine Learning Repository

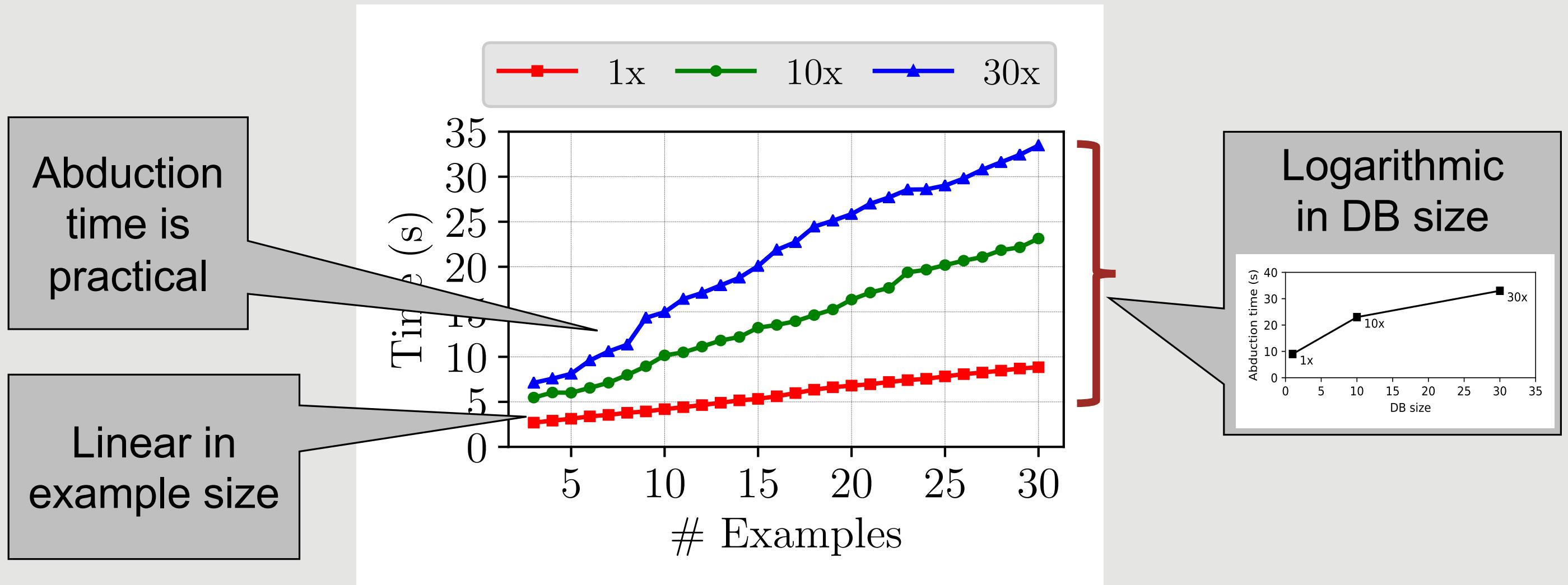
5 benchmark
queries

20 benchmark
queries

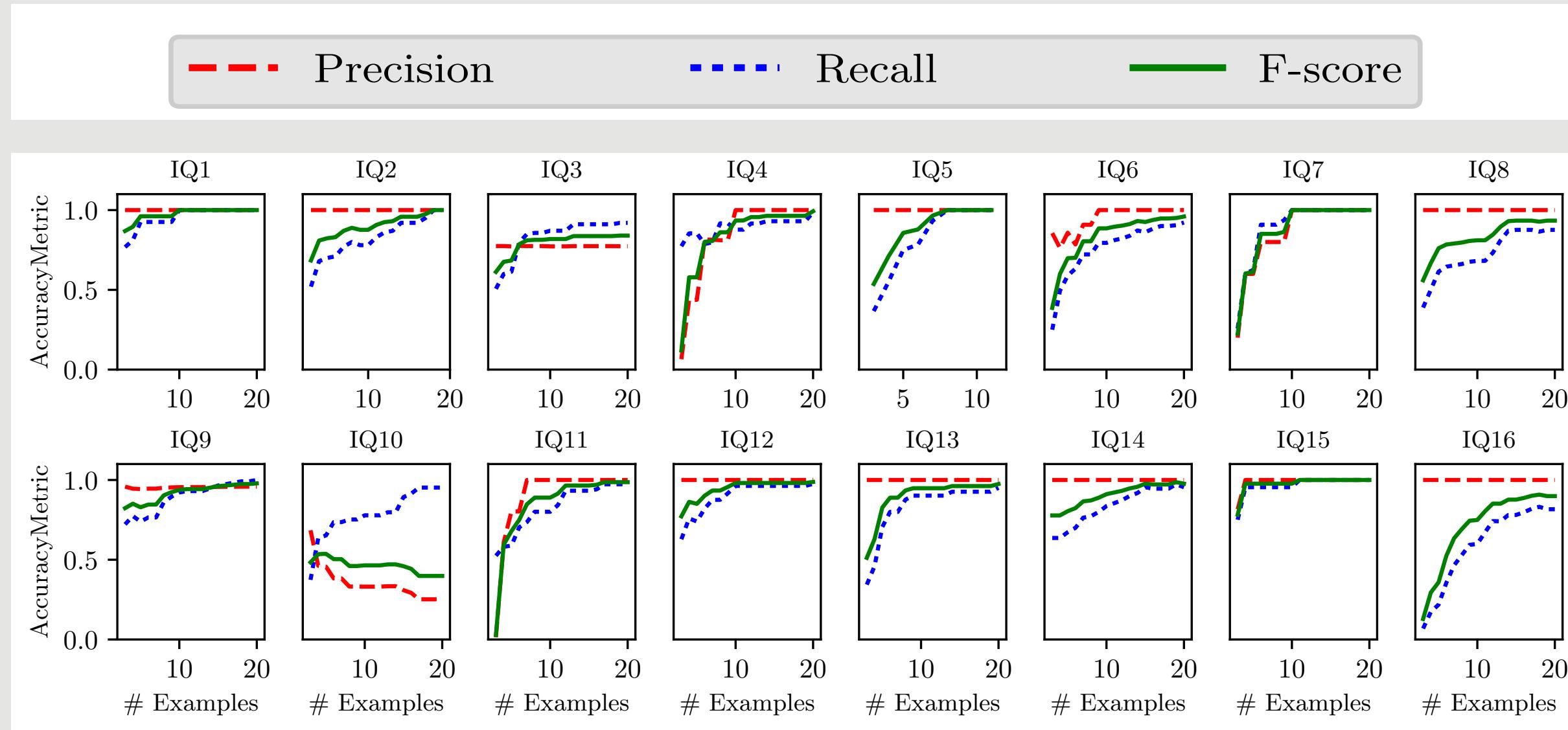
Experiment settings



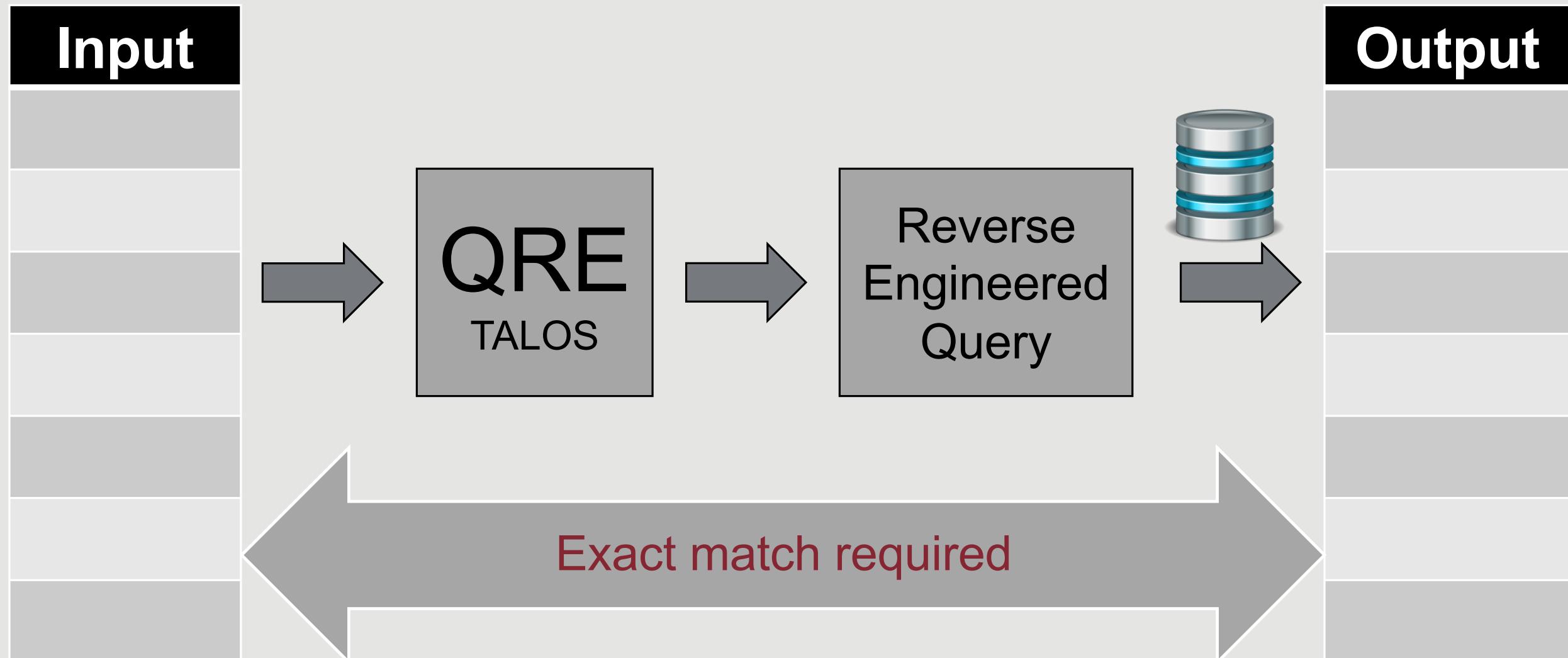
How does SQuID perform with large datasets or many examples?



SQuID works with few examples

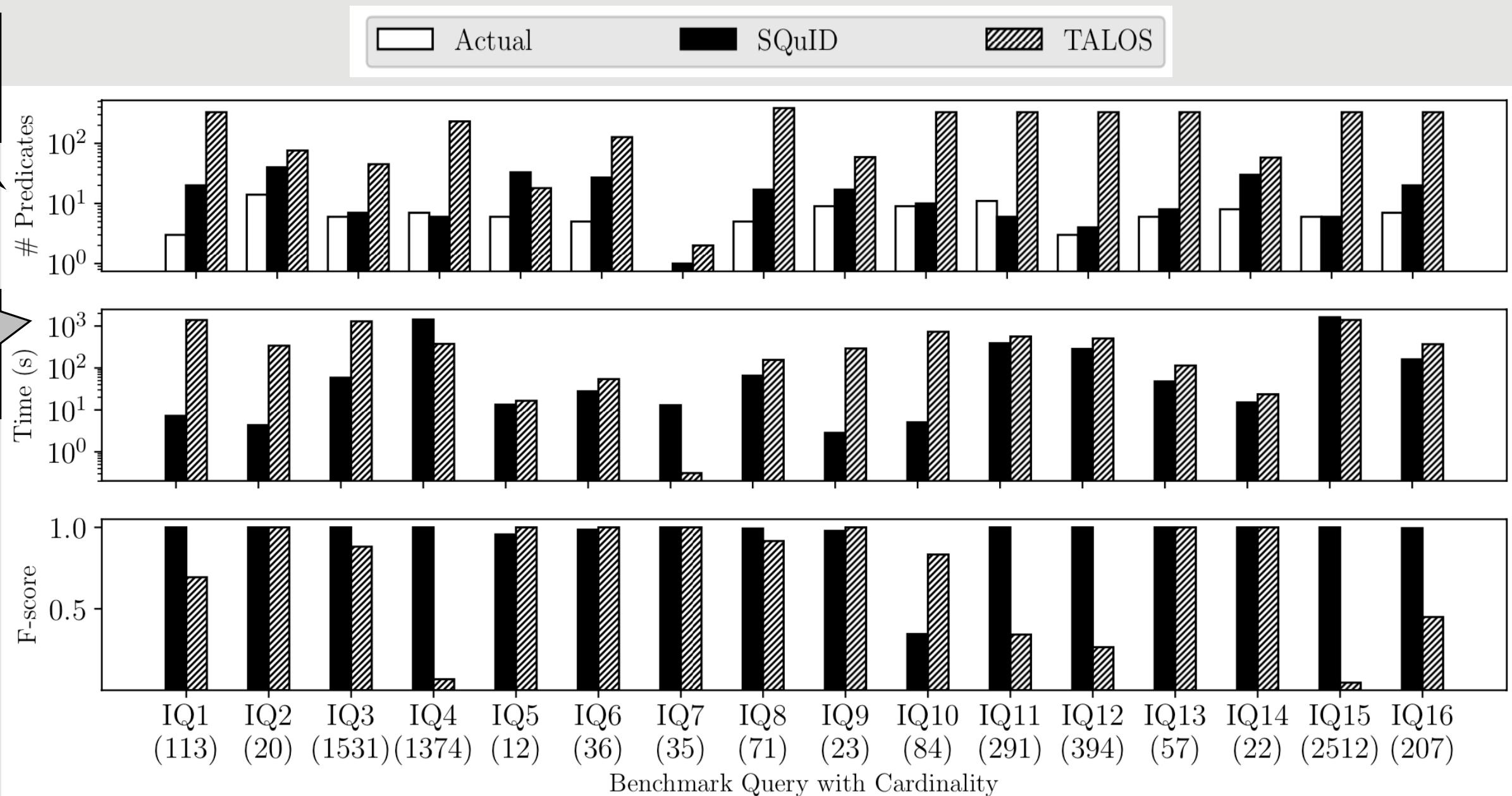


Query reverse engineering (QRE)

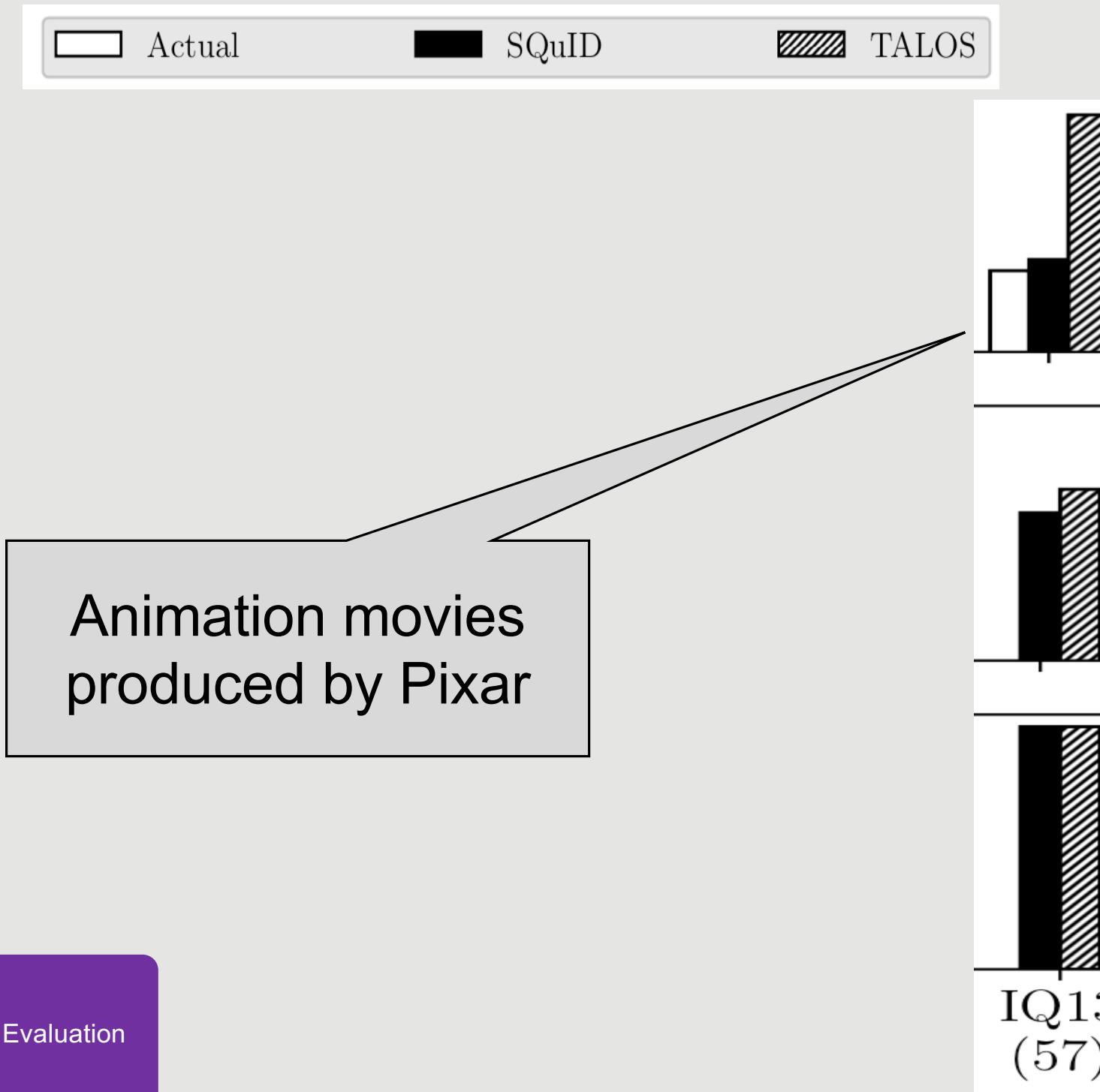


SQuID outperforms QRE

Log scale



SQuID outperforms QRE



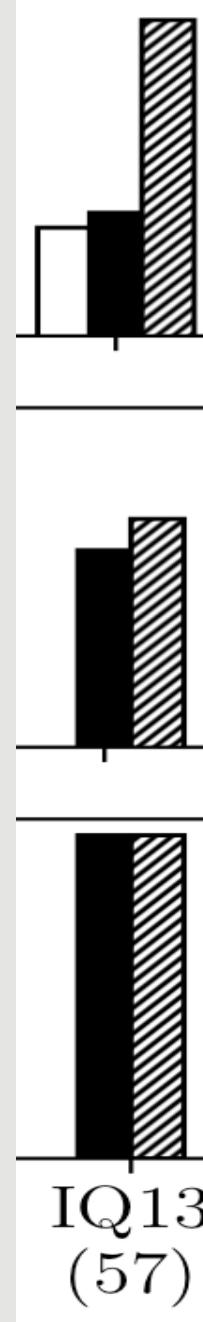
Animation movies produced by Pixar

SQuID outperforms QRE



Original Query

```
SELECT  
    DISTINCT movie.title  
FROM  
    movie, production, company, movietogenre, genre  
WHERE  
    movie.id = production.movie_id AND  
    production.company_id = company.id AND  
    company.name LIKE '%Pixar%' AND  
    movie.id = movietogenre.movie_id AND  
    movietogenre.genre_id = genre.id AND  
    genre.name = 'Animation';
```



SQuID outperforms QRE



Original Query	SQuID Query
<pre>SELECT DISTINCT movie.title FROM movie, production, company, movietogenre, genre WHERE movie.id = production.movie_id AND production.company_id = company.id AND movie.id = movietogenre.movie_id AND movietogenre.genre_id = genre.id AND movie.country = USA AND genre.name = 'Animation' AND company.name = 'Pixar'</pre>	<pre>SELECT DISTINCT movie.title FROM movie, production, company, movietogenre, genre WHERE movie.production_year >= 1984 AND movie.production_year <= 2021 AND movie.country = USA AND genre.name = 'Animation' AND company.name = 'Pixar' AND movie.id = movietogenre.movie_id AND genre.id = movietogenre.genre_id AND movie.id = movietoproduction.movie_id AND company.id = movietoproduction.company_id</pre>



SQuID outperforms QRE

UMassAmherst

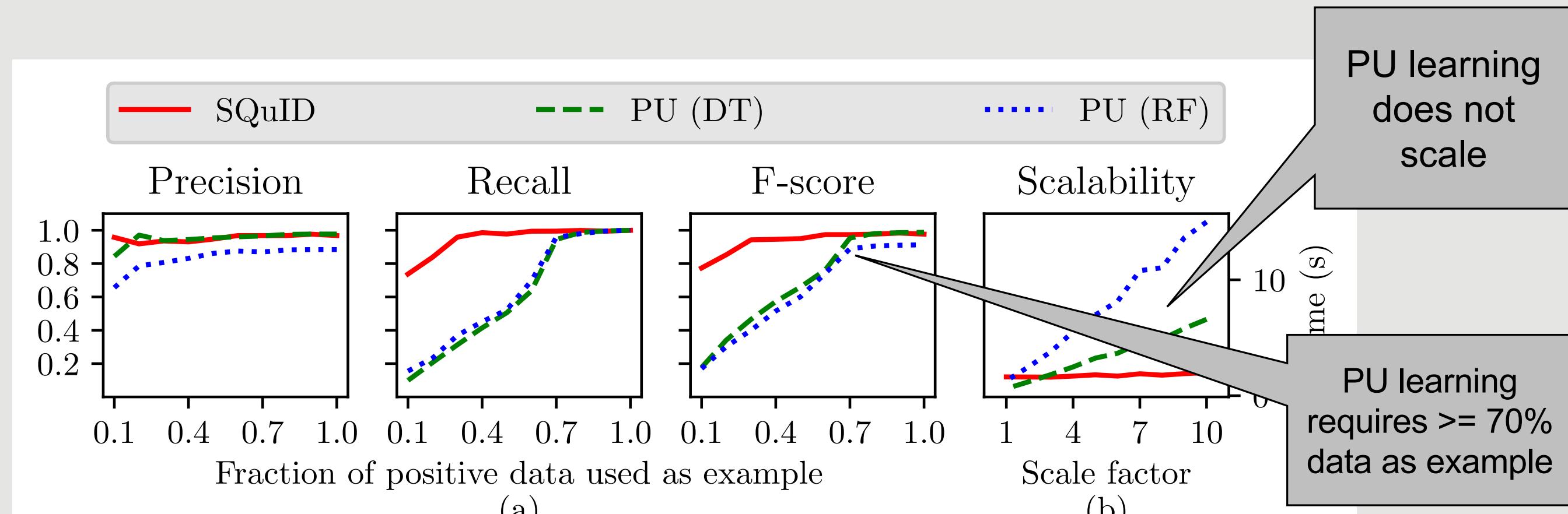
College of Information
& Computer Sciences



Original Query	SQuID Query	TALOS Query
<pre>SELECT DISTINCT FROM movie, production_company, genre_name WHERE movie.id = production_company.movie_id AND movie.id = genre_name.movie_id</pre>	<pre>SELECT DISTINCT FROM movie, production_company, genre_name WHERE movie.id = production_company.movie_id AND movie.id = genre_name.movie_id</pre>	<pre>SELECT distinct title FROM movie WHERE movie.production_company_id = movie.id AND movie.genre_name_id = movie.id</pre>

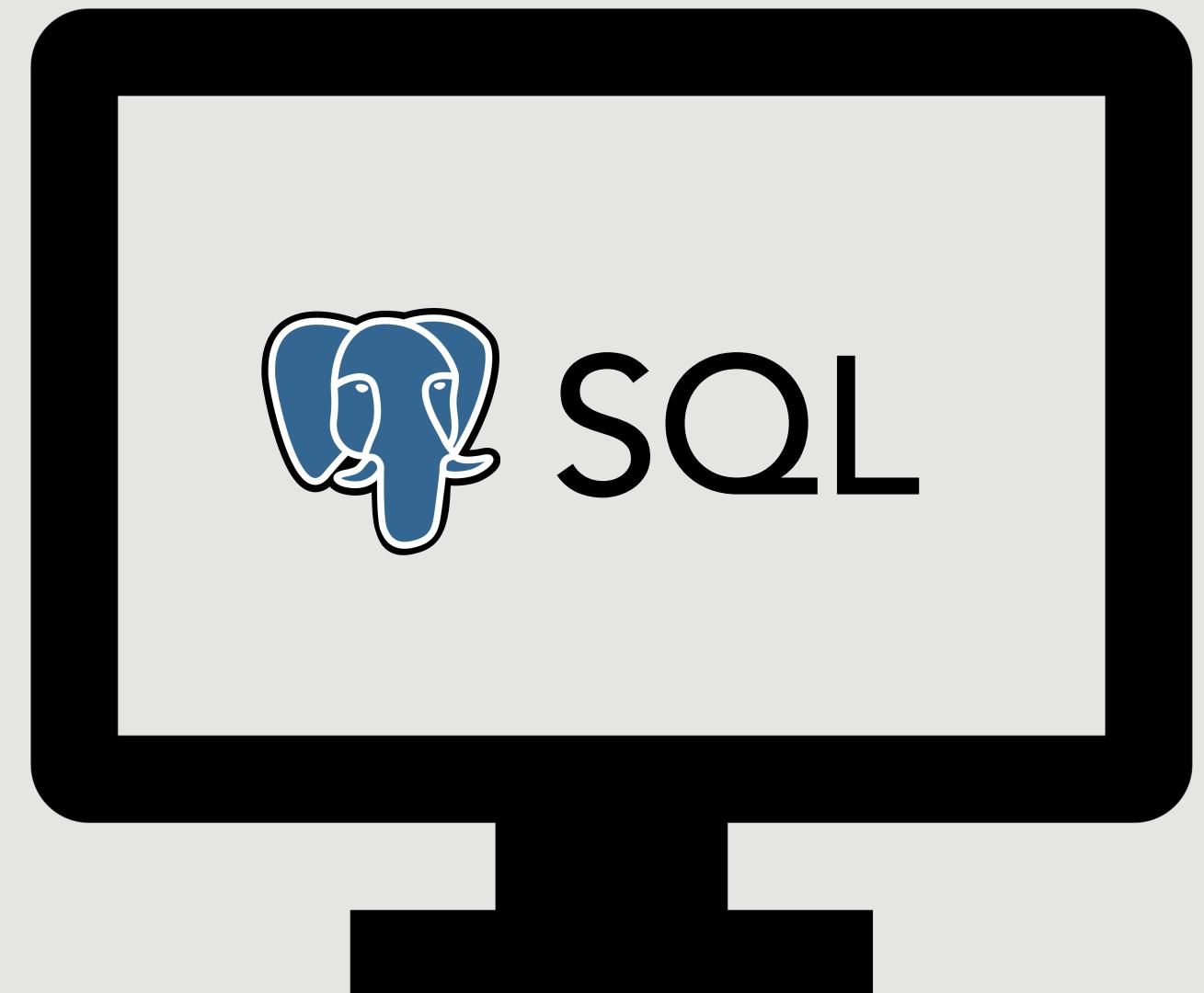
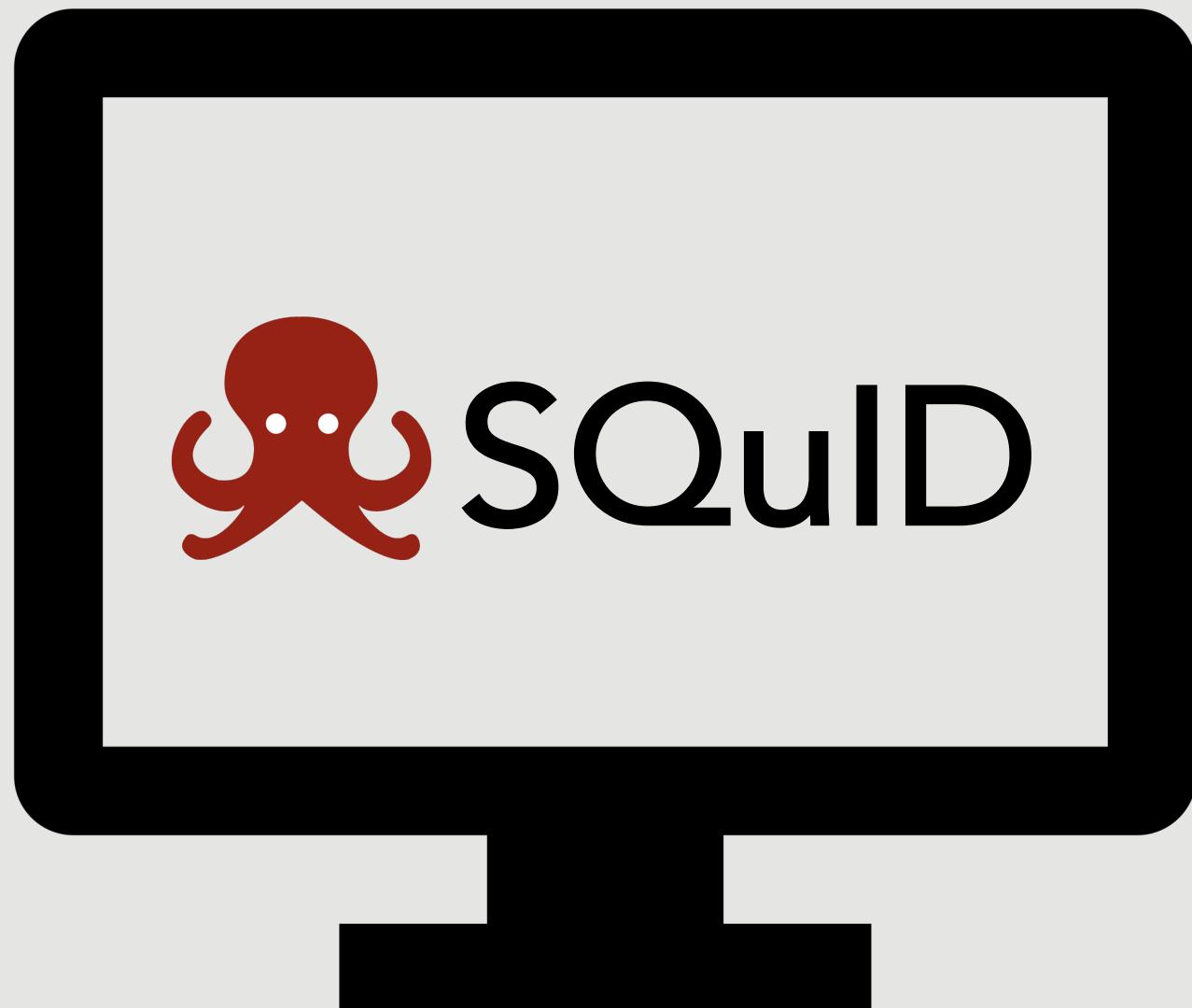
Query Reverse
Engineering
overfits

SQuID outperforms machine learning

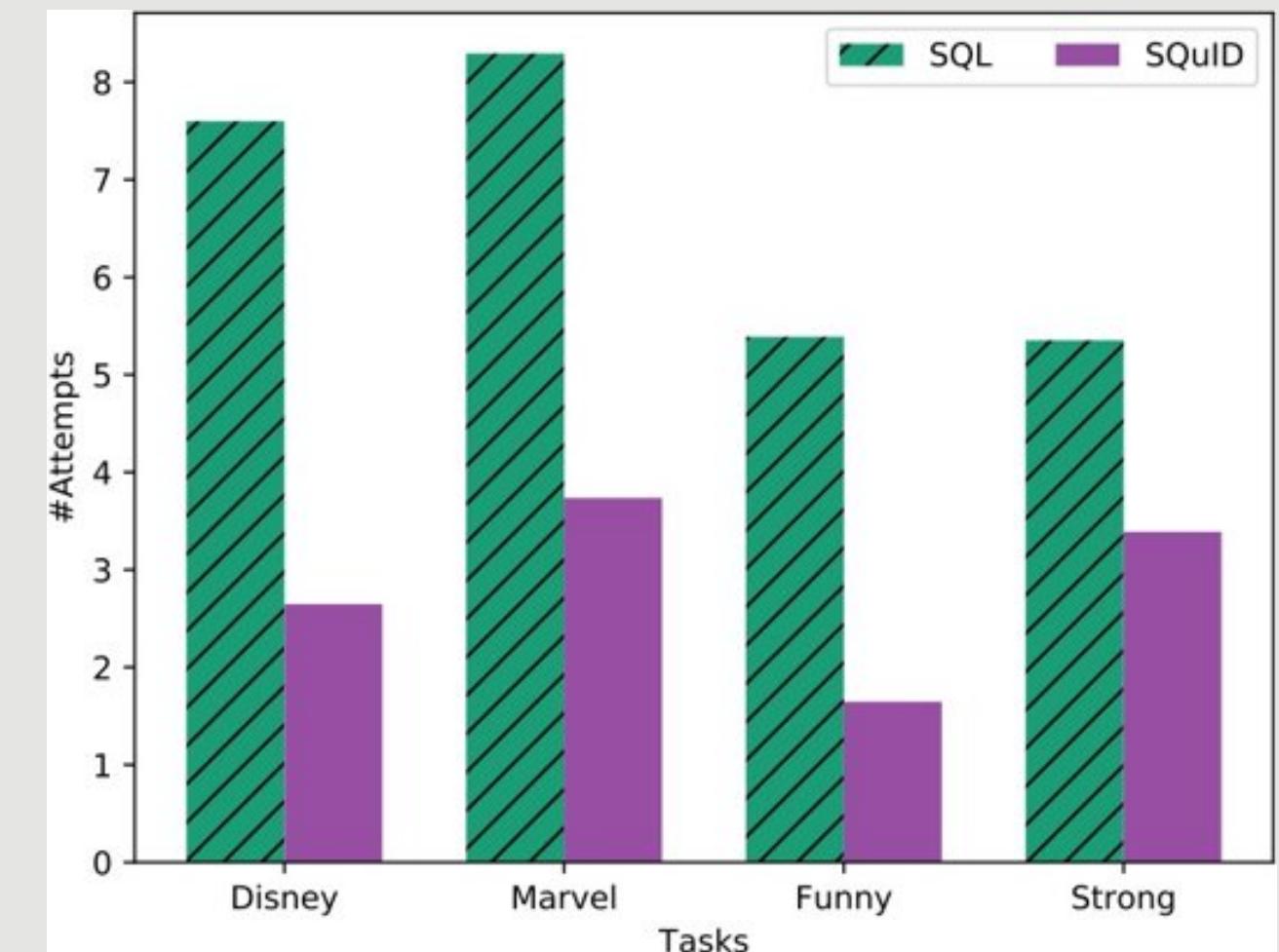
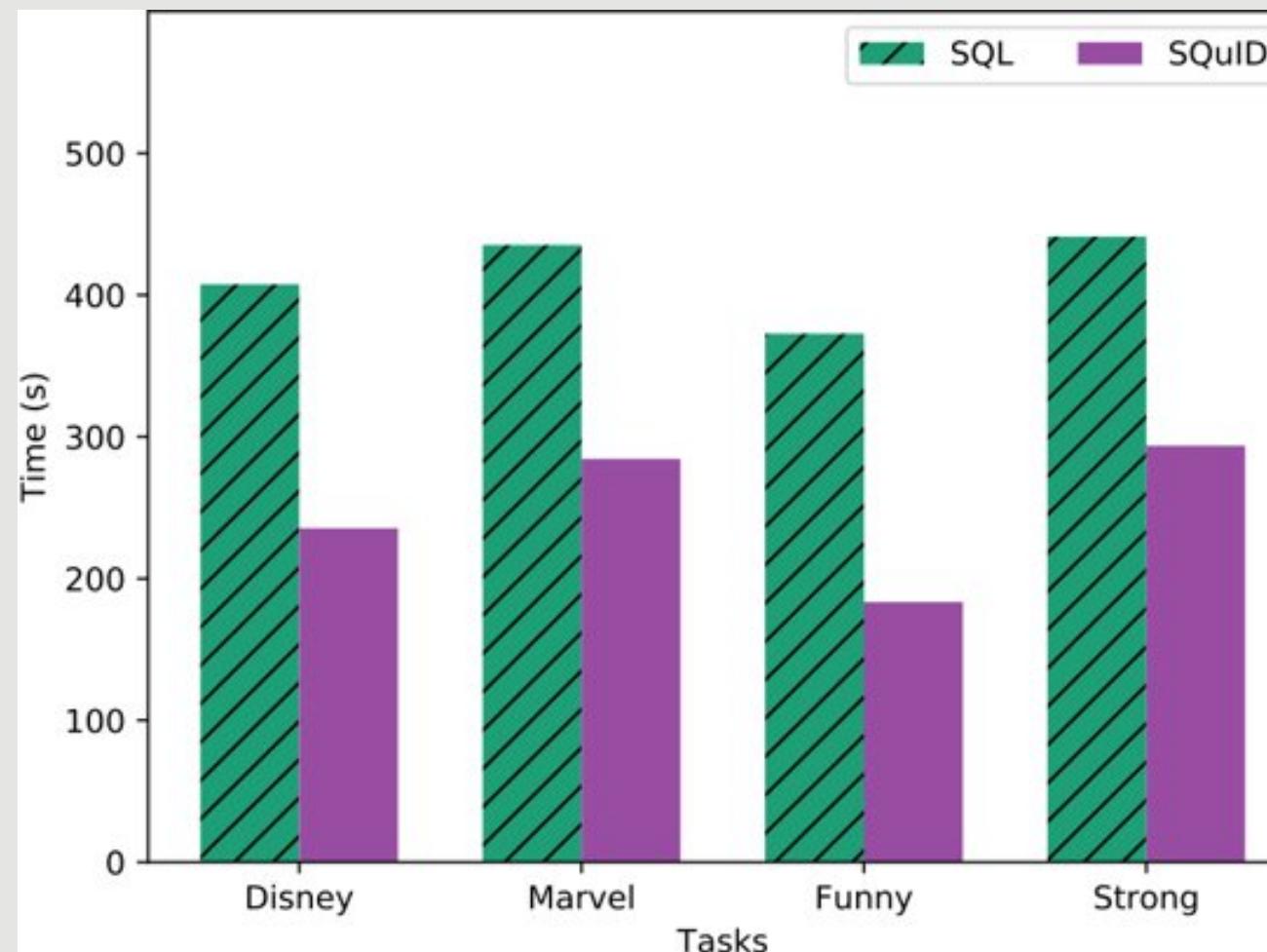


Generic machine learning cannot model RDBMS specific assumptions

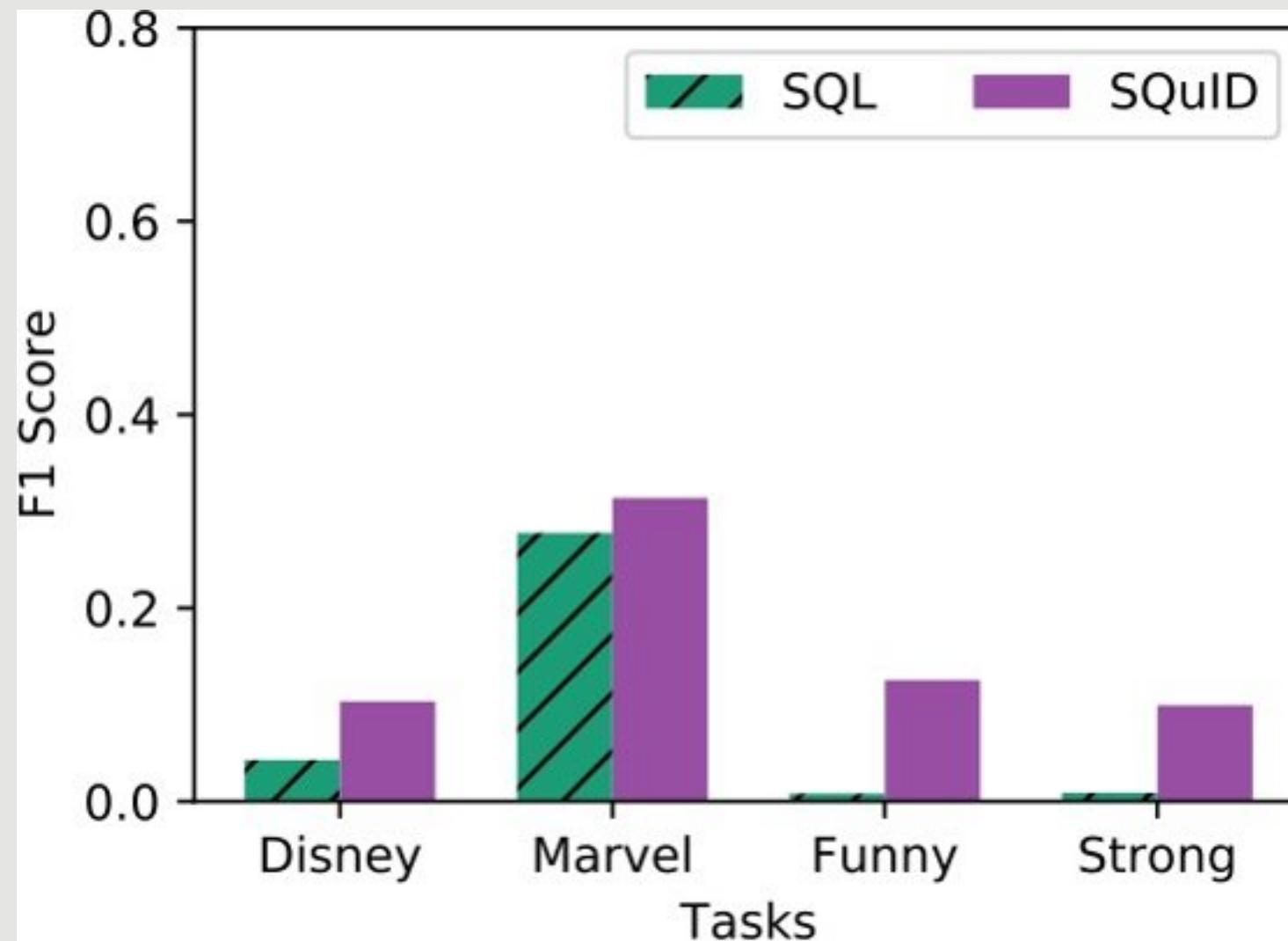
Comparative user studies: QBE vs SQL



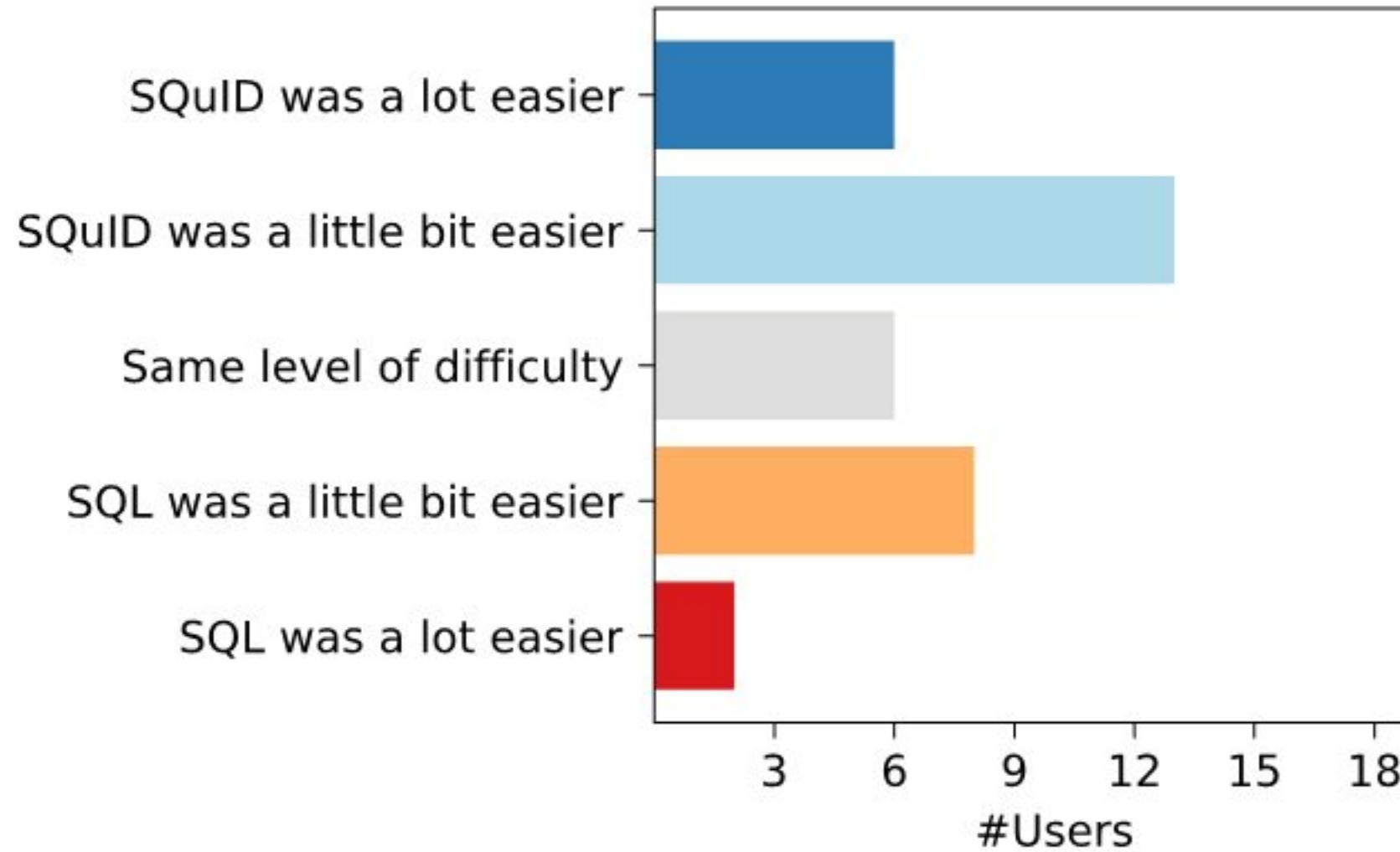
SQuID increased user efficiency



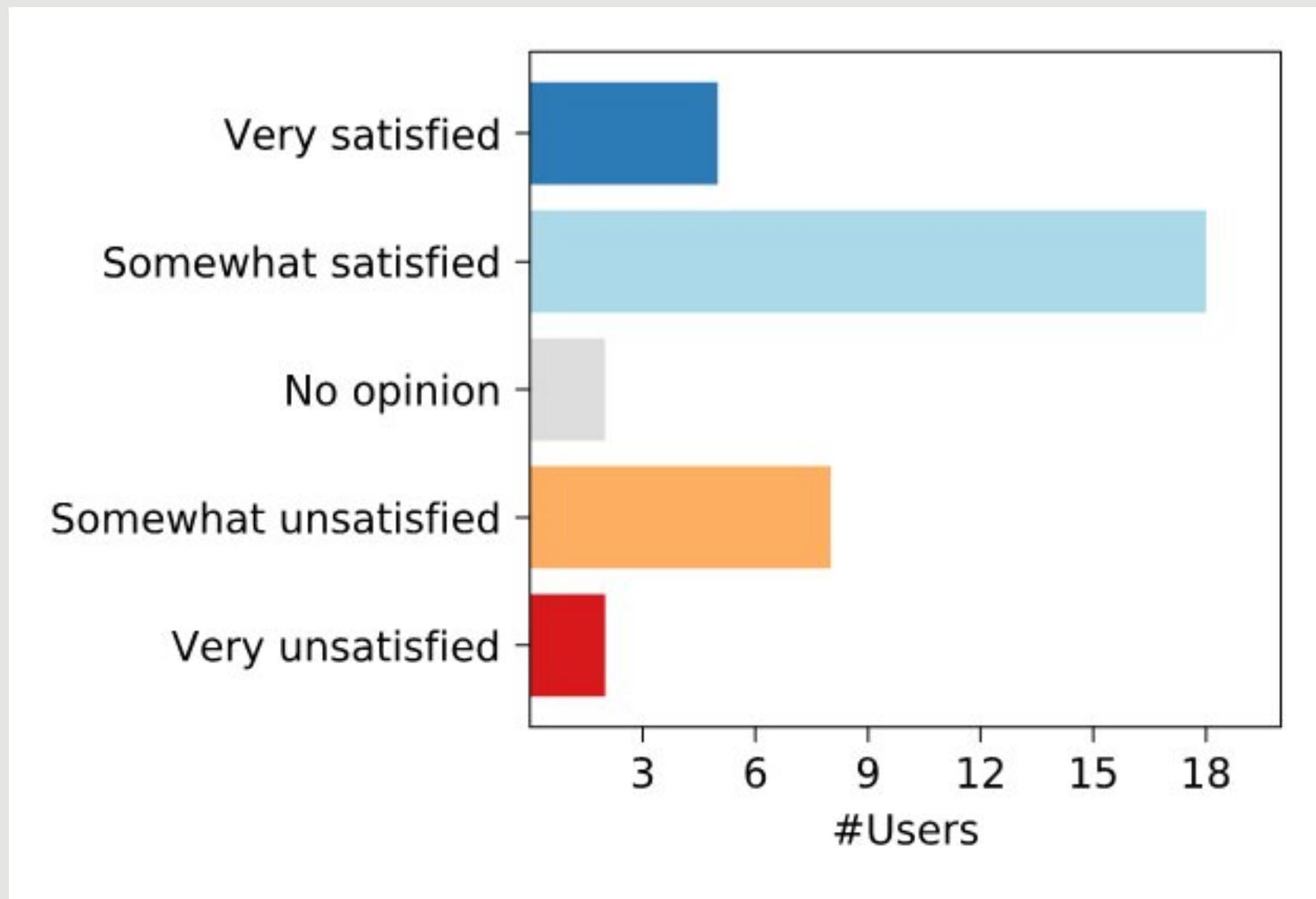
Overall, SQuID generated more accurate results



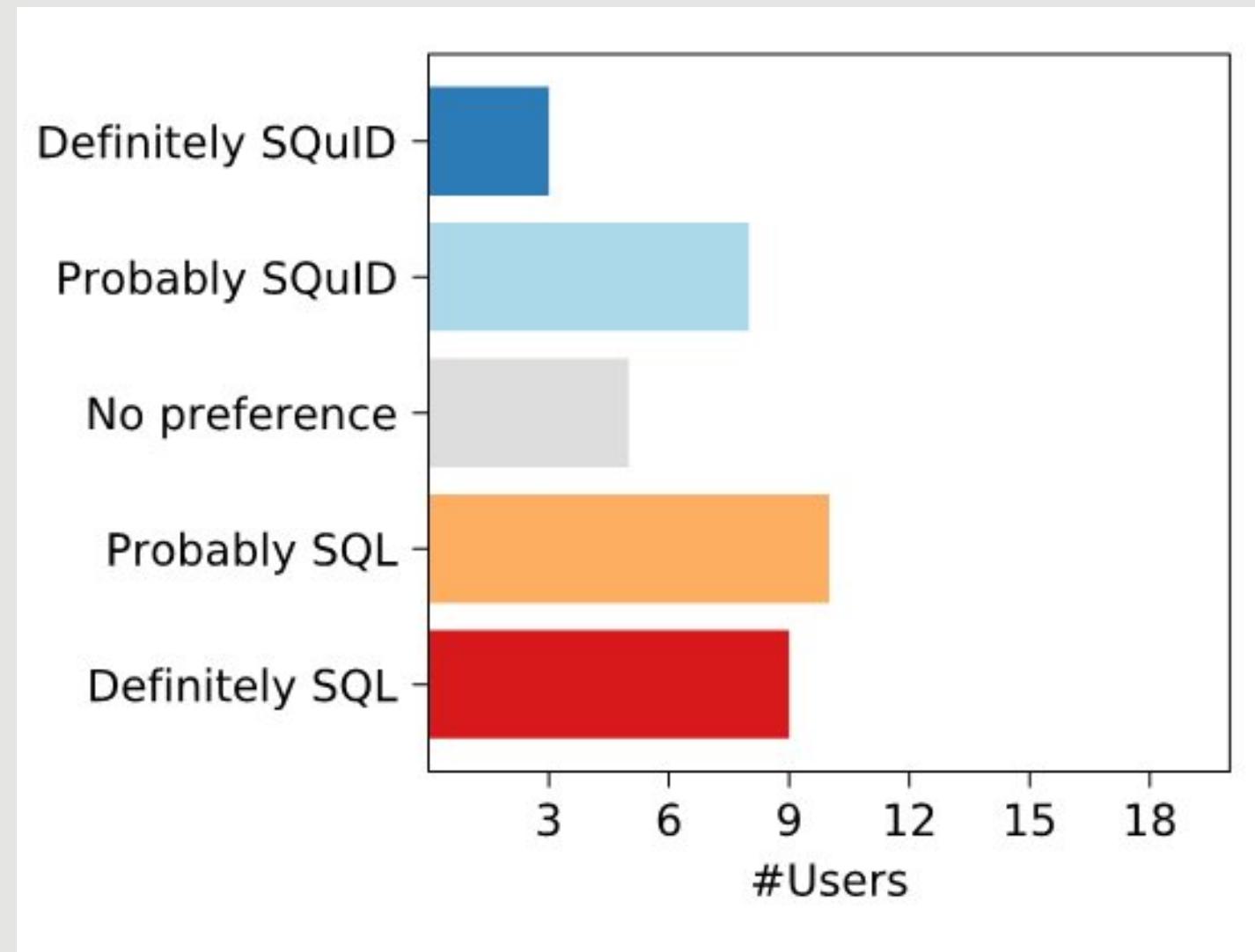
SQuID was easier to use



Participants were satisfied with SQuID results



SQuID or SQL?



Anecdotal comments



*“Even if I forget about syntax . . . figuring out **how to go about writing the pseudo-code** query for funny actors [is difficult]”*

*“Vague tasks are generally a lot more open to interpretation. Coding up a query that meets someone’s **vague specifications** [is] hard . . . It was **very hard to nail down what the correct definition of funny is.**”*

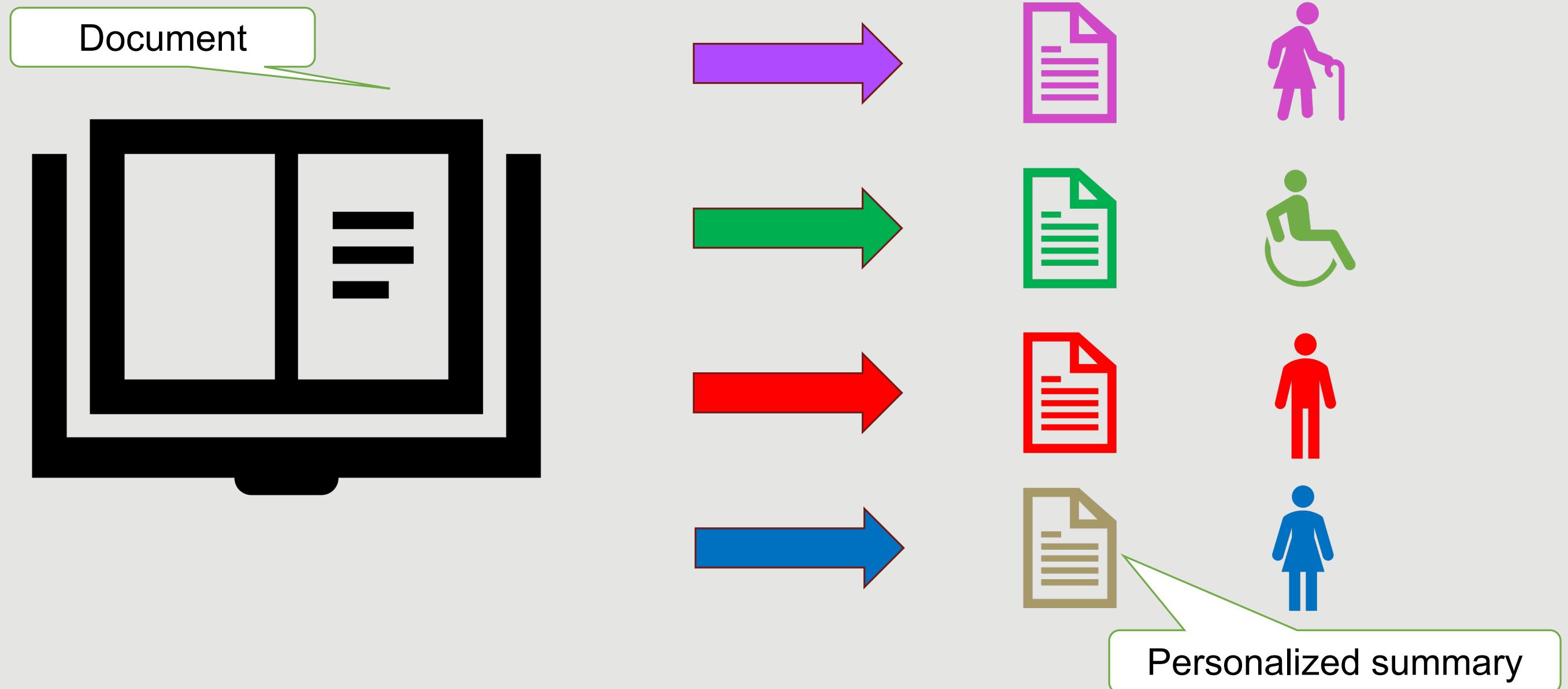


Personalized text document summarization

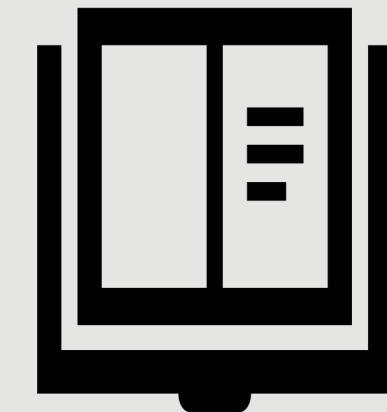
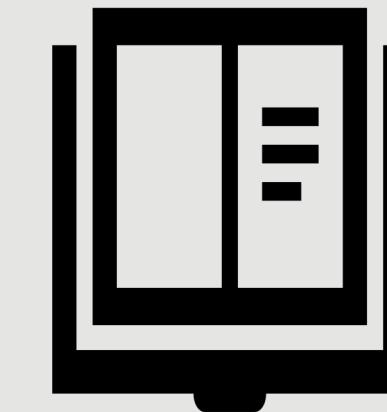
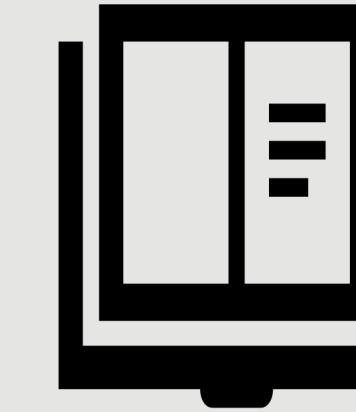
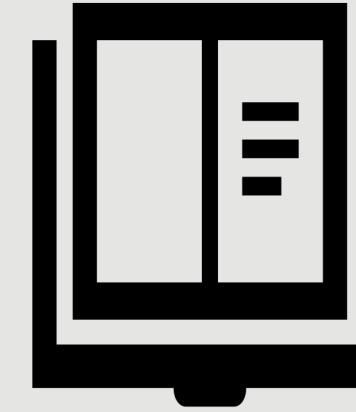


**SuDocu: Summarizing
Documents by Example**

Personalized summarization



Summarization by example



Example Summaries

Automatic Summaries

SuDocu interface

Summary Input

1 Utah

Sentences (120):

In 1957, Utah created the [Utah State Parks](#) Commission with four parks. Today, Utah State Parks manages 43 parks and several undeveloped areas totaling over 95,000 acres of land and more than 1,000,000 acres of water. Utah's state parks are scattered throughout Utah, from [Bear Lake State Park](#) at the Utah/Idaho border to [Edge of the Cedars State Park](#) Museum deep in the [Four Corners](#) region and everywhere in between. Utah State Parks is also home to the state's [off highway vehicle](#) office, state boating office and the trails program. [33]

Submit Summary

Generated Summaries

4 Massachusetts

It borders on the Atlantic Ocean to the east, the states of Connecticut and Rhode Island to the south, New Hampshire and Vermont to the north, and New York to the west. The large coastal plain of the Atlantic Ocean in the eastern section of the state contains Greater Boston, along with most of the state's population, as well as the distinctive Cape Cod peninsula. Along the western border of Western Massachusetts lies the highest elevated part of the state, the Berkshires. Most of Massachusetts has a humid continental, with cold winters and warm summers. The climate of Boston is quite representative for the commonwealth, characterized by summer highs of around 81 °F (27 °C) and winter highs of 35 °F (2 °C), and is quite wet. Frosts are frequent all winter, even in coastal areas due to prevailing inland winds.

Example Summaries

Utah

The state of Utah relies heavily on income from tourists and travelers visiting the state's parks and ski resorts. Today, Utah State Parks manages 43 parks and several undeveloped areas totaling over 95,000 acres of land and more than 1,000,000 acres of water. With five national parks (Arches, Bryce Canyon, Canyonlands, Capitol Reef, and Zion), Utah has the third most national parks of any state after Alaska and California. Temperatures dropping below 0 °F (-18 °C) should be expected on occasion in most areas of the state most years.

Arizona

Arizona is well known for its desert Basin and Range region in the state's southern portions, which is rich in a landscape of xerophyte plants such as the cactus. The canyon is one of the Seven Natural Wonders of the World and is largely contained in the Grand Canyon National Park—one of the first national parks in the United States. Extremely cold temperatures are not unknown; cold air systems from the northern states and Canada occasionally push into the state, bringing temperatures below 0 °F (-18 °C) to the state's northern parts.

Montana

The Rocky Mountain Front is a significant feature in the state's north-central portion, and isolated island ranges that interrupt the prairie landscape common in the central and eastern parts of the state. It contains the state's highest point, Granite Peak, 12,799 feet high. Farther east, areas such as Makoshika State Park near Glendive and Medicine Rocks State Park near Ekalaka contain some of the most scenic badlands regions in the state. The coldest temperature on record for Montana is also the coldest temperature for the contiguous United States. On January 20, 1954, -70 °F or -56.7 °C was recorded at a gold mining camp near Rogers Pass. Temperatures vary greatly on cold nights.

Summarize

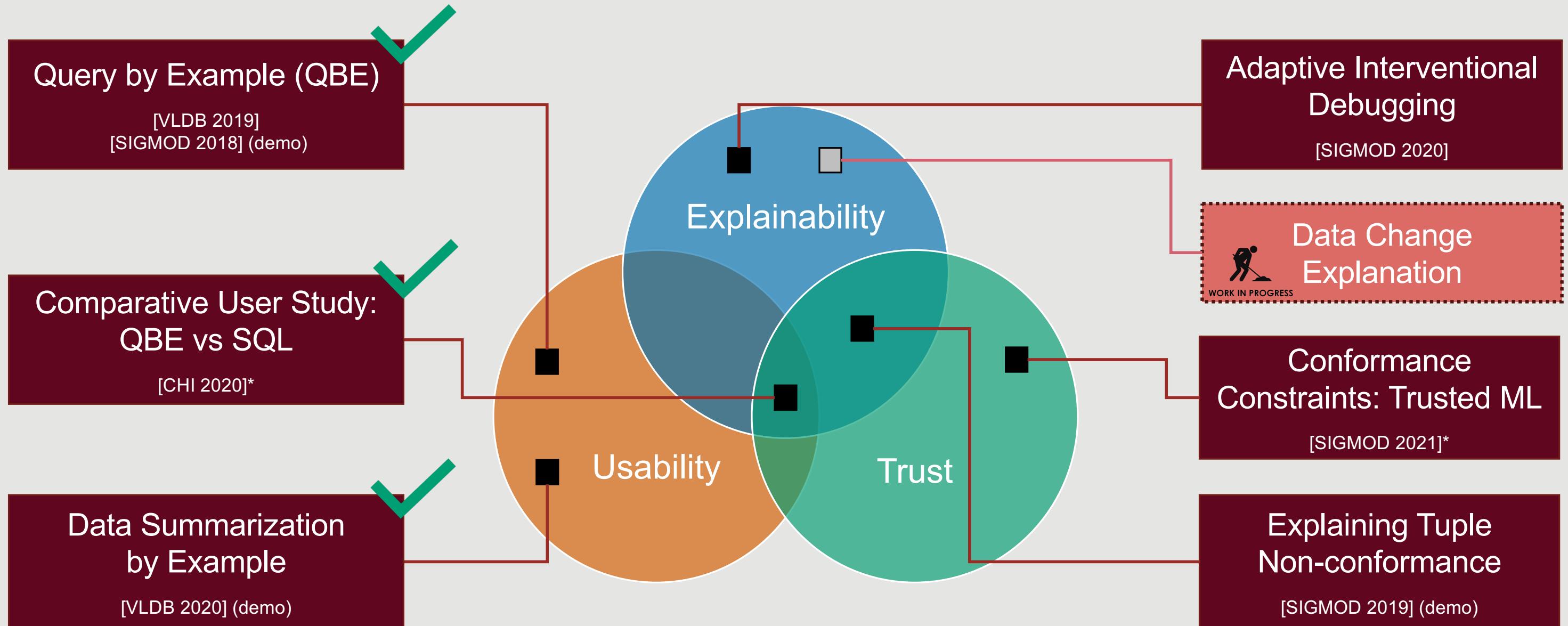
Explanation (PaQL)

```

SELECT PACKAGE(*)
FROM state_sentences
WHERE state = 'Massachusetts'
SUCH THAT
  SUM(topic_1) BETWEEN 0.06 AND 0.45 AND
  SUM(topic_2) BETWEEN 0.24 AND 0.79 AND
  SUM(topic_3) BETWEEN 0.41 AND 0.84 AND
  SUM(topic_4) BETWEEN 0.83 AND 1.85 AND
  SUM(topic_5) BETWEEN 0.95 AND 1.29 AND
  SUM(topic_6) BETWEEN 2.64 AND 3.20 AND
  SUM(topic_7) BETWEEN 2.14 AND 4.72 AND
  SUM(topic_8) BETWEEN 0.07 AND 0.43 AND
  SUM(topic_9) BETWEEN 0.07 AND 0.41 AND
  SUM(topic_10) BETWEEN 0.58 AND 0.84
MAXIMIZE
  SUM(m_score)
topic_6: climate, temperature, summer, winter, ...
  
```



Dissertation outline



Part 2: Trust in Data Systems



Trust



To trust or not to trust?



Daily **Mail**
.com

Science & Tech

IBM's Watson AI suggested 'often inaccurate' and 'unsafe' treatment recommendations for cancer patients, internal documents show

To trust or not to trust?

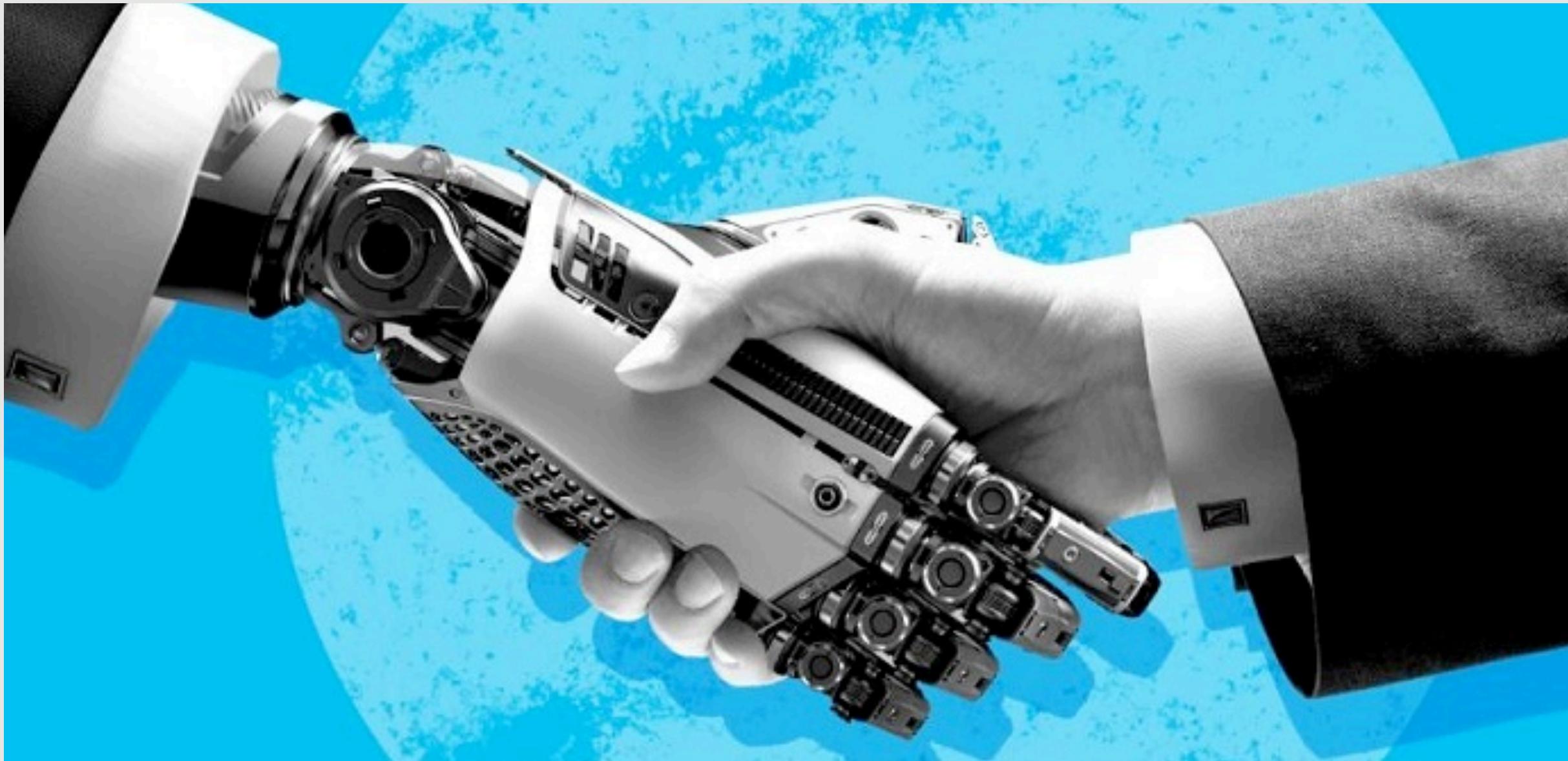


The New York Times

Self-Driving Uber Car Kills Pedestrian in Arizona, Where Robots Roam

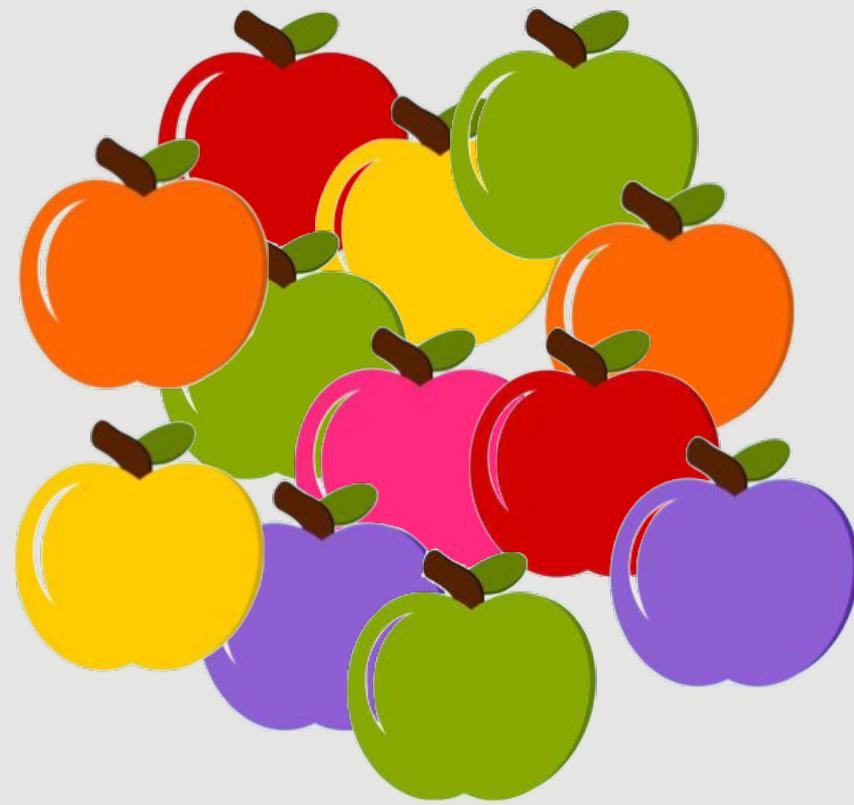


Conformance constraints: trusted machine learning

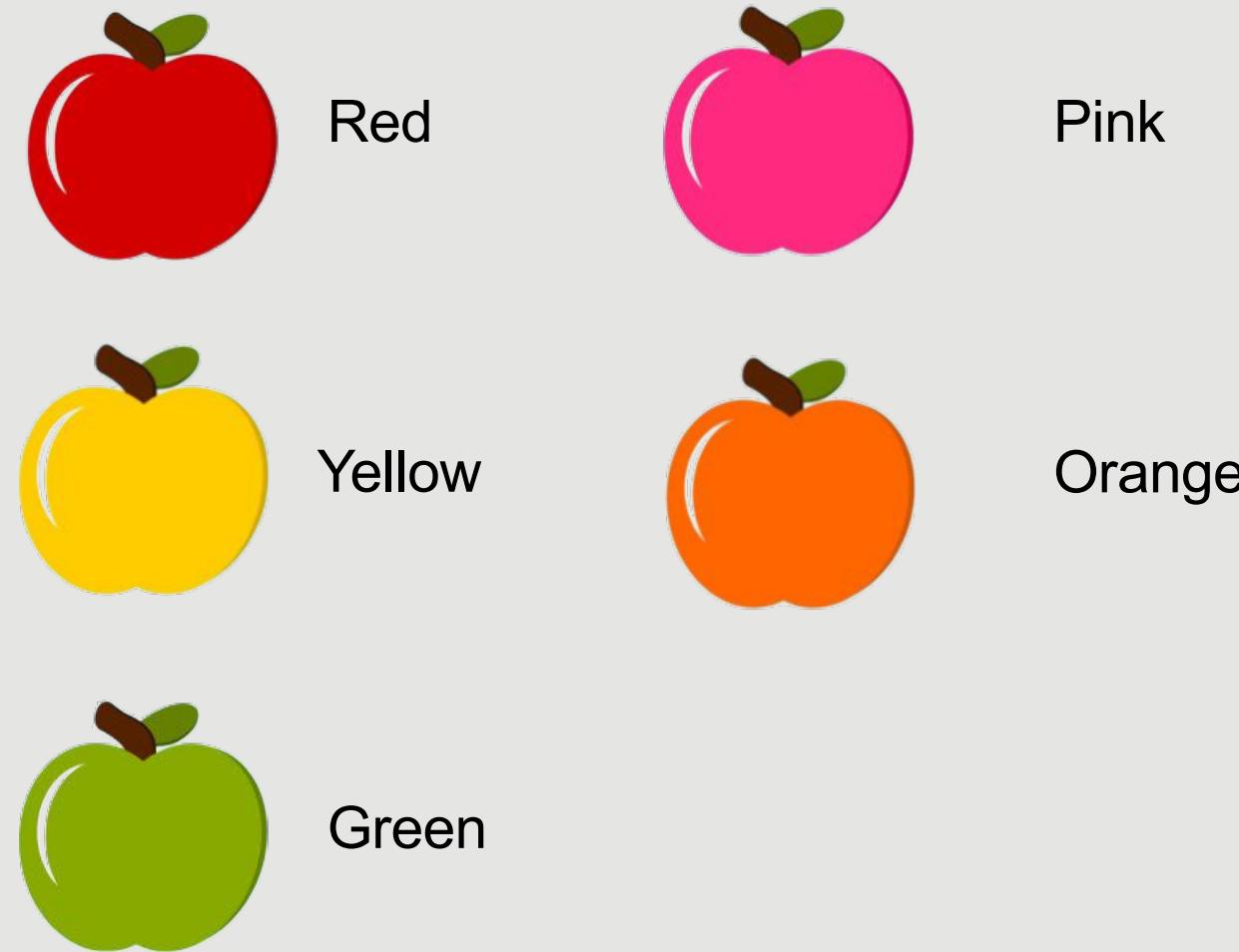


Trusting ML predictions

Training data

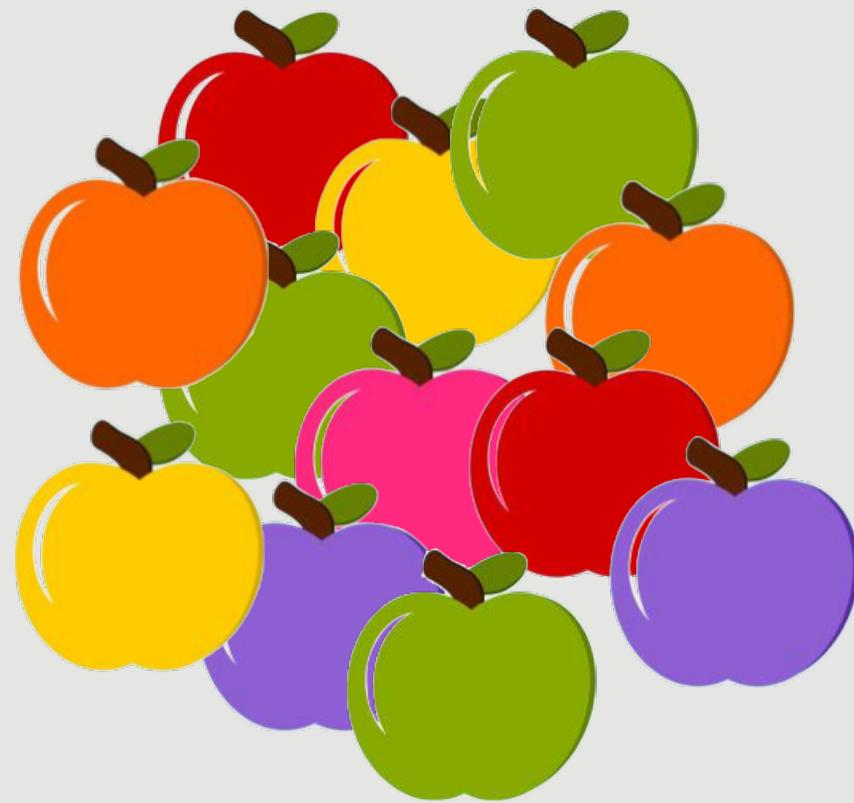


New data

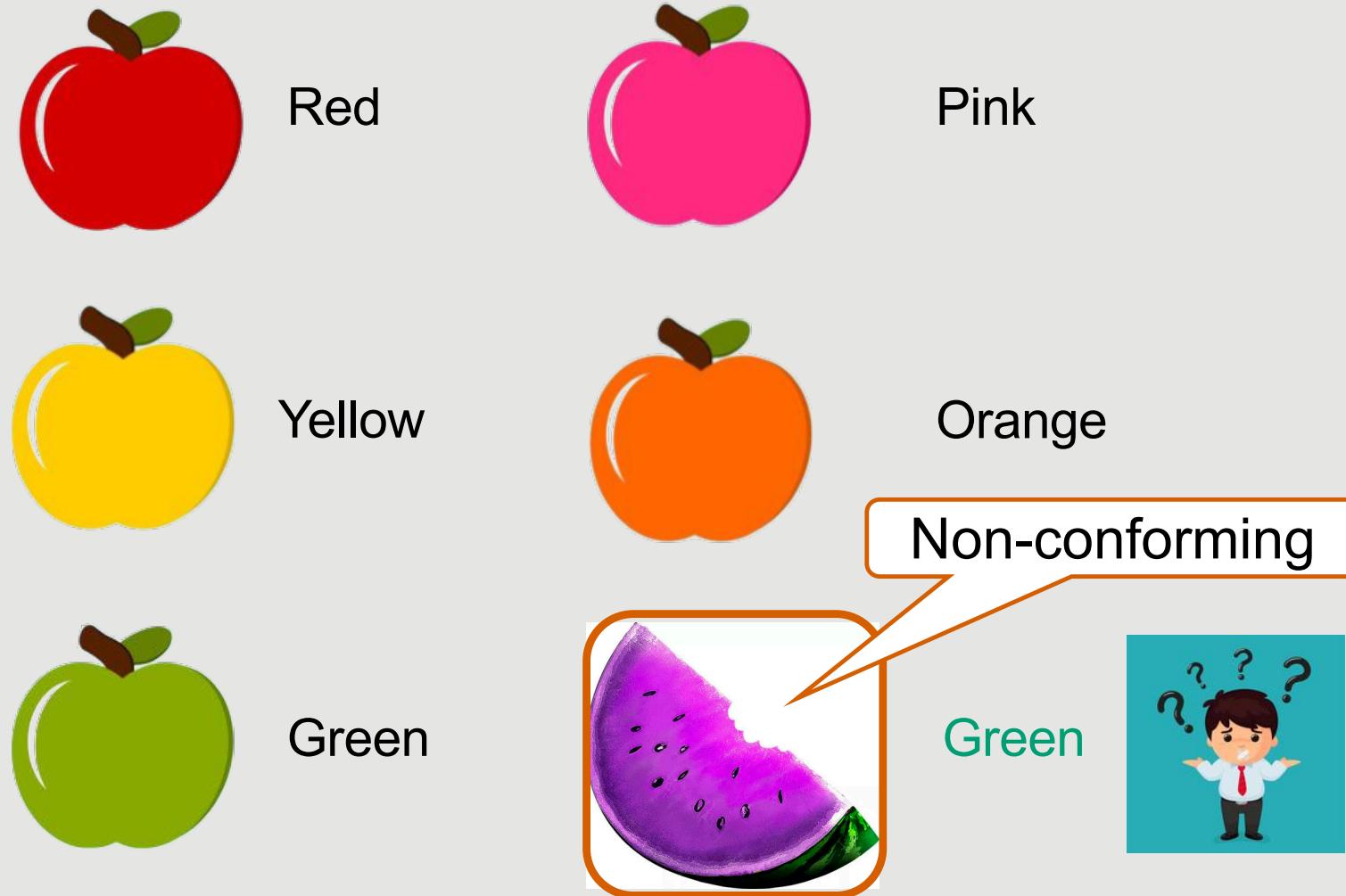


Trusting ML predictions

Training data



New data



Non-conformance = untrustworthy prediction

Detection

Is it
non-conforming?



A real-world example: airlines dataset

Regression task: predict arrival delay

...	dep_date	dep_time	arr_time	duration (minute)
...	May 2	14:30	18:20	230
...	July 22	09:05	12:15	195
...	June 6	10:20	20:00	582
...	May 19	11:10	13:05	117
...	April 7	22:30	06:10	458

DAYTIME flights

OVERNIGHT flight

A real-world example: airlines dataset

- Trained with DAYTIME flights only
- Constraints observed in DAYTIME flights
 - “departure time is earlier than arrival time”
 - “their difference is very close to flight duration”
- OVERNIGHT flights
 - violate DAYTIME flights’ constraints
 - incur high regression error

Constraint violation correlates with high regression error

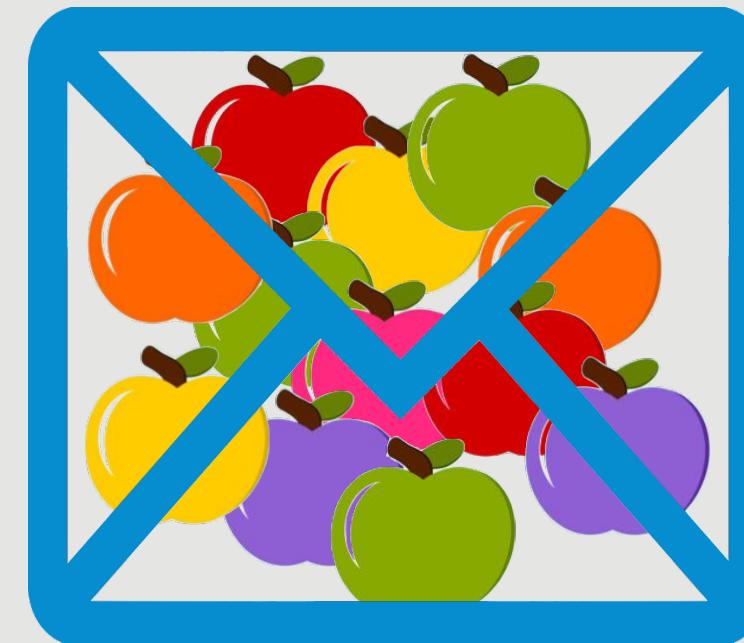
Conformance constraints (CCs)

ML pipelines **drop** low-variance dimensions to achieve dimensionality reduction.

ML models **assume** that training data's constraints/properties will continue to hold during serving.

Conformance constraints

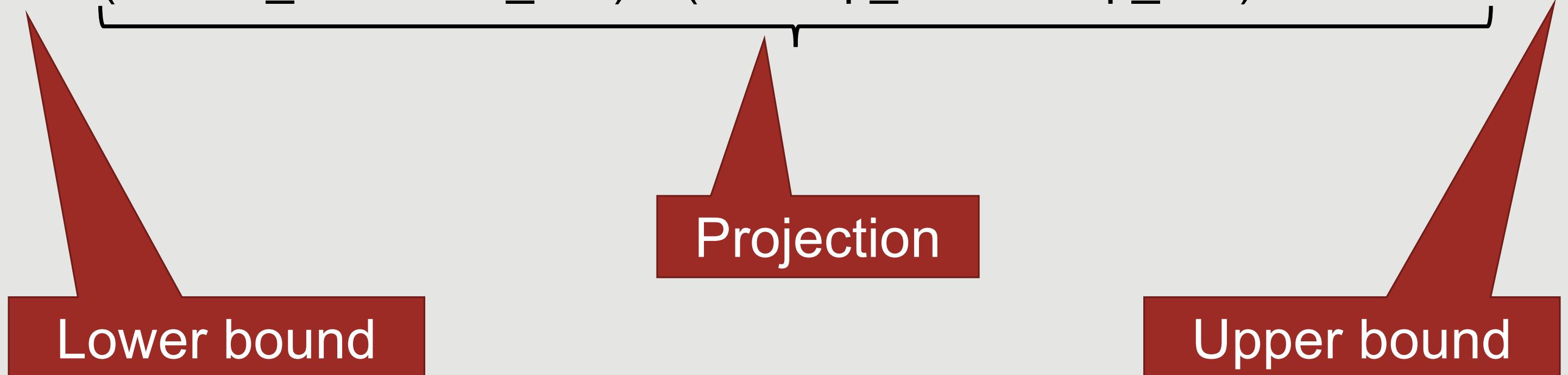
- constraints that the data satisfies
- capture the invariants of the data



Conformance constraints

- Encode linear arithmetic relationship over multiple attributes.

$$-\epsilon \leq (60 \cdot \text{arr_hour} + \text{arr_min}) - (60 \cdot \text{dep_hour} + \text{dep_min}) - \text{duration} \leq \epsilon$$



Projection

Lower bound

Upper bound

Conformance constraints: example

Height	Weight	BMI
6 feet	142 lbs	19.3
5 feet	170 lbs	33.2
5 feet	130 lbs	25.4

$10 \leq \text{BMI} \leq 40$

$-40 \leq (28 \times \text{Height} - \text{Weight}) \leq 30$

Violation of conformance constraint

$10 \leq \text{BMI} \leq 40$

Height	Weight	BMI
6 feet	142 lbs	19.3
5 feet	170 lbs	33.2
5 feet	130 lbs	25.4
6 feet	170 lbs	231

Degree of violation

$10 \leq \text{BMI} \leq 40$

Height	Weight	BMI
6 feet	142 lbs	19.3
5 feet	170 lbs	33.2
5 feet	130 lbs	25.4
6 feet	170 lbs	231
6 feet	170 lbs	20000

Projection

$$-\epsilon \leq (60 \cdot \text{arr_hour} + \text{arr_min}) - (60 \cdot \text{dep_hour} + \text{dep_min}) - \text{duration} \leq \epsilon$$

Lower bound

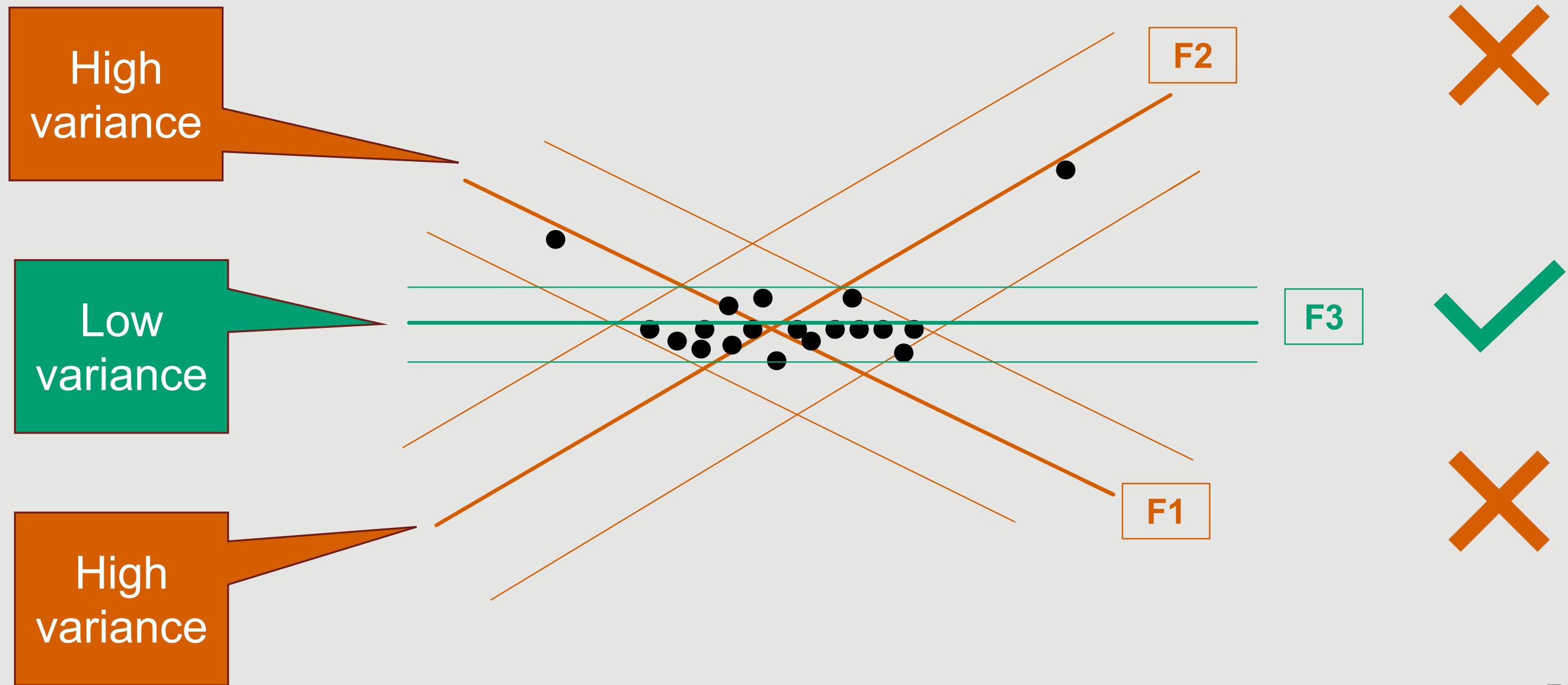
Projection

Upper bound

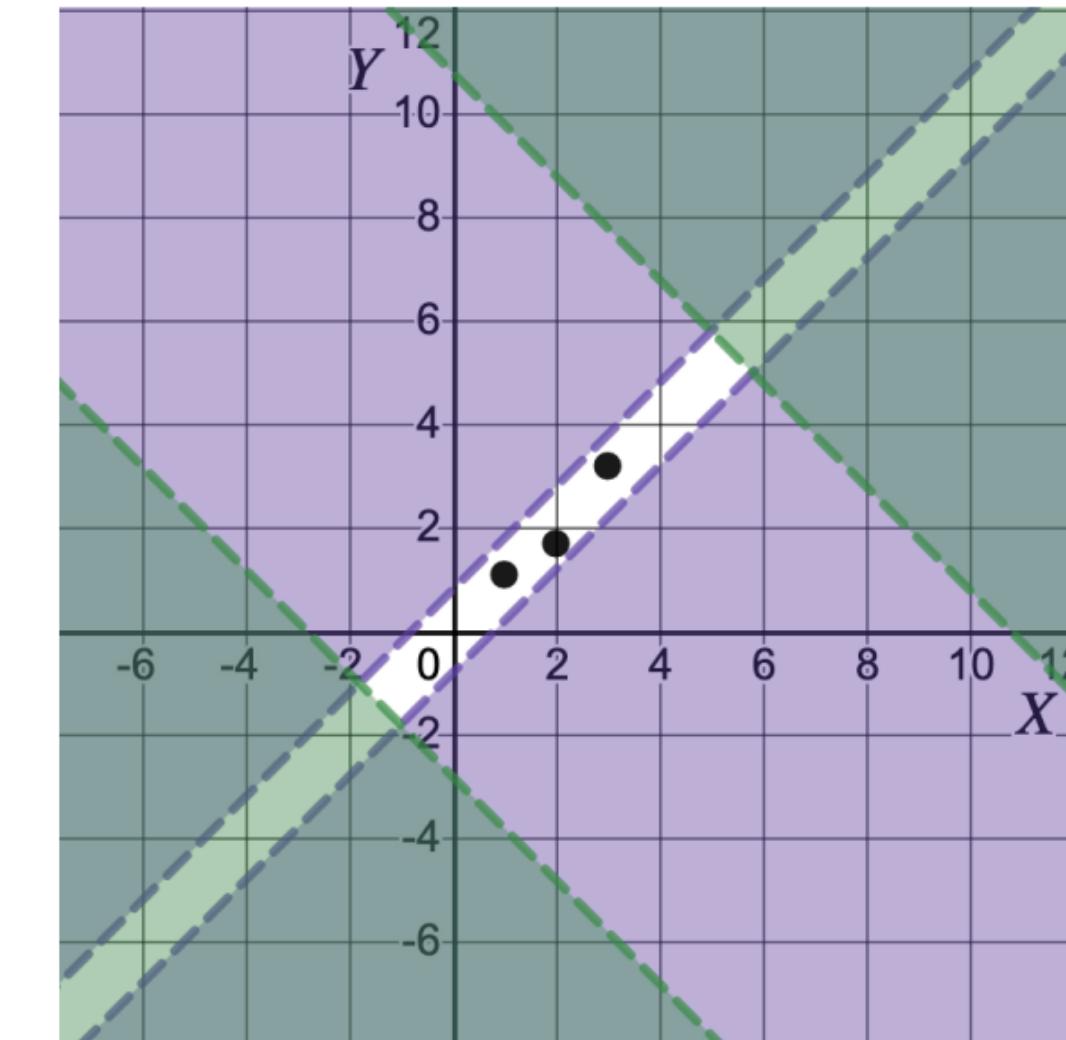
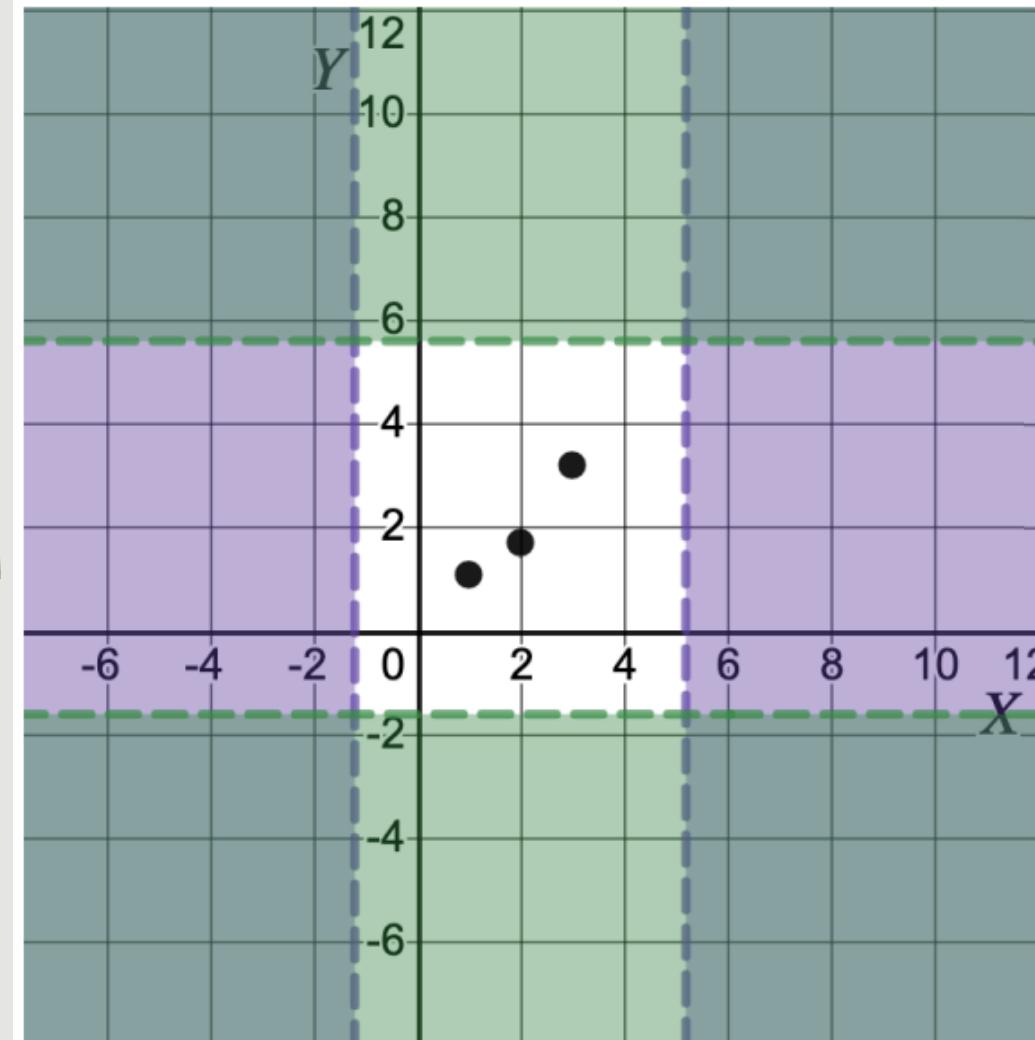
What are “good” projections?

- Infinitely many projections possible
 - Pick the **low-variance** projections.
 - Because?
 - They more useful in detecting **trends** in the data.
- Do we pick all low-variance projections?
 - Pick a set of projections with **low pair-wise correlations**.
 - Because?
 - They **complement** each other.

Low-variance projections

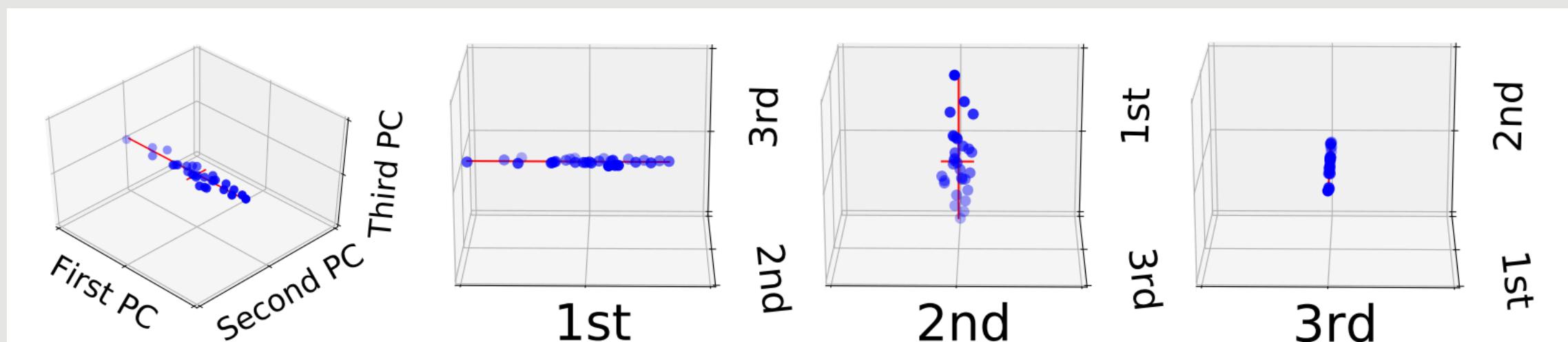


Projections with small mutual correlation



Discovering projections: PCA

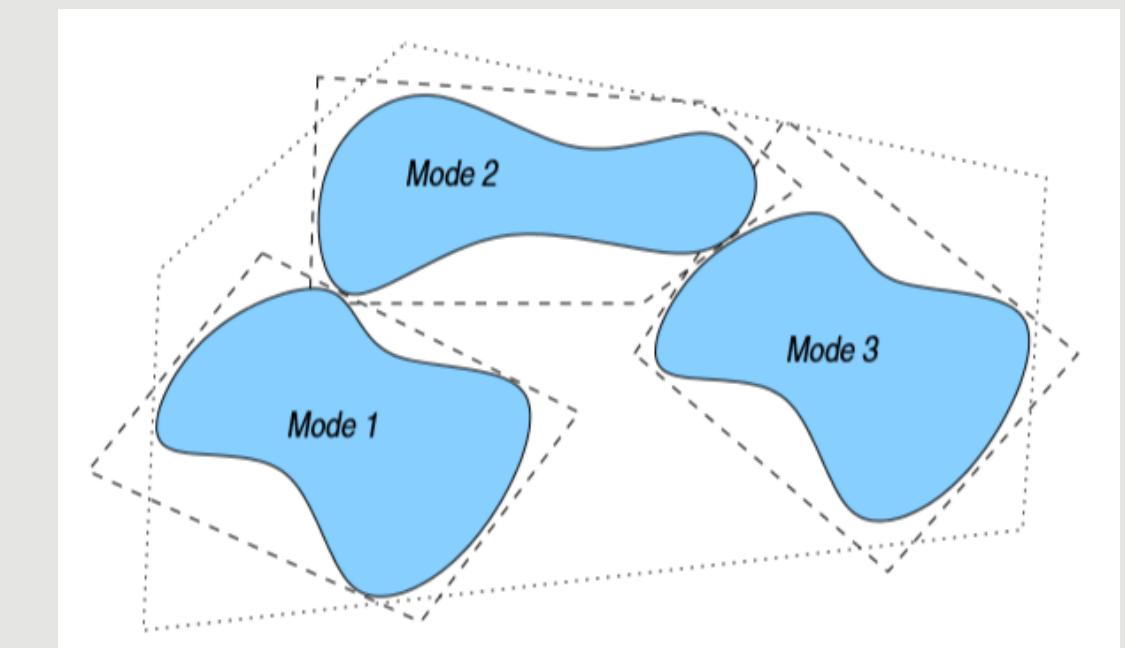
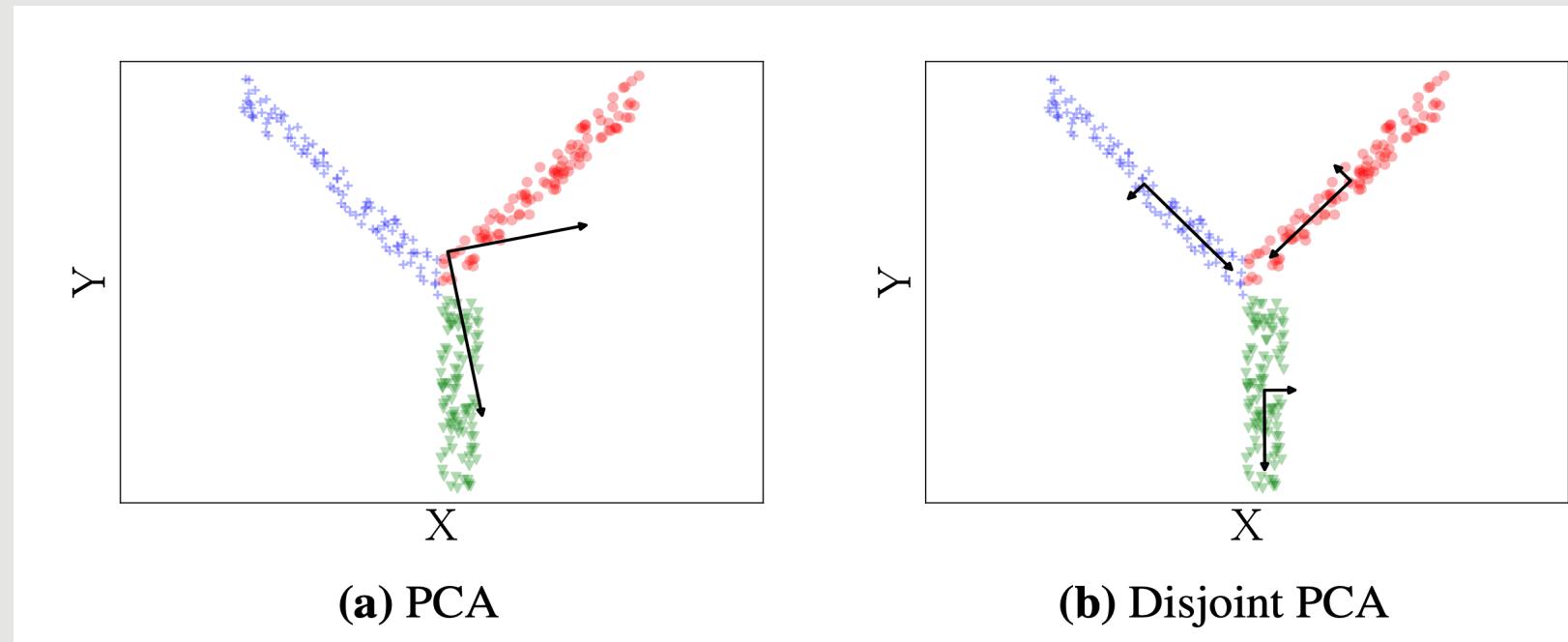
- Principal Component Analysis (PCA)
 - Produces projections with small mutual correlations
 - Intuition: principal components are **orthogonal** to each other
 - Computing violation
 - Weigh CCs with **low** variance projections **more**
 - Weigh CCs with **high** variance projections **less**



Disjunctive conformance constraints

- Divide the dataset into **disjoint** partitions.
- Learn CCs for each partition.
- Compute **disjunctive** CCs.

$$\begin{aligned}\psi_2 : M = \text{"May"} \triangleright -2 \leq AT - DT - DUR \leq 0 \\ \vee \quad M = \text{"June"} \triangleright \quad 0 \leq AT - DT - DUR \leq 5 \\ \vee \quad M = \text{"July"} \triangleright -5 \leq AT - DT - DUR \leq 0\end{aligned}$$



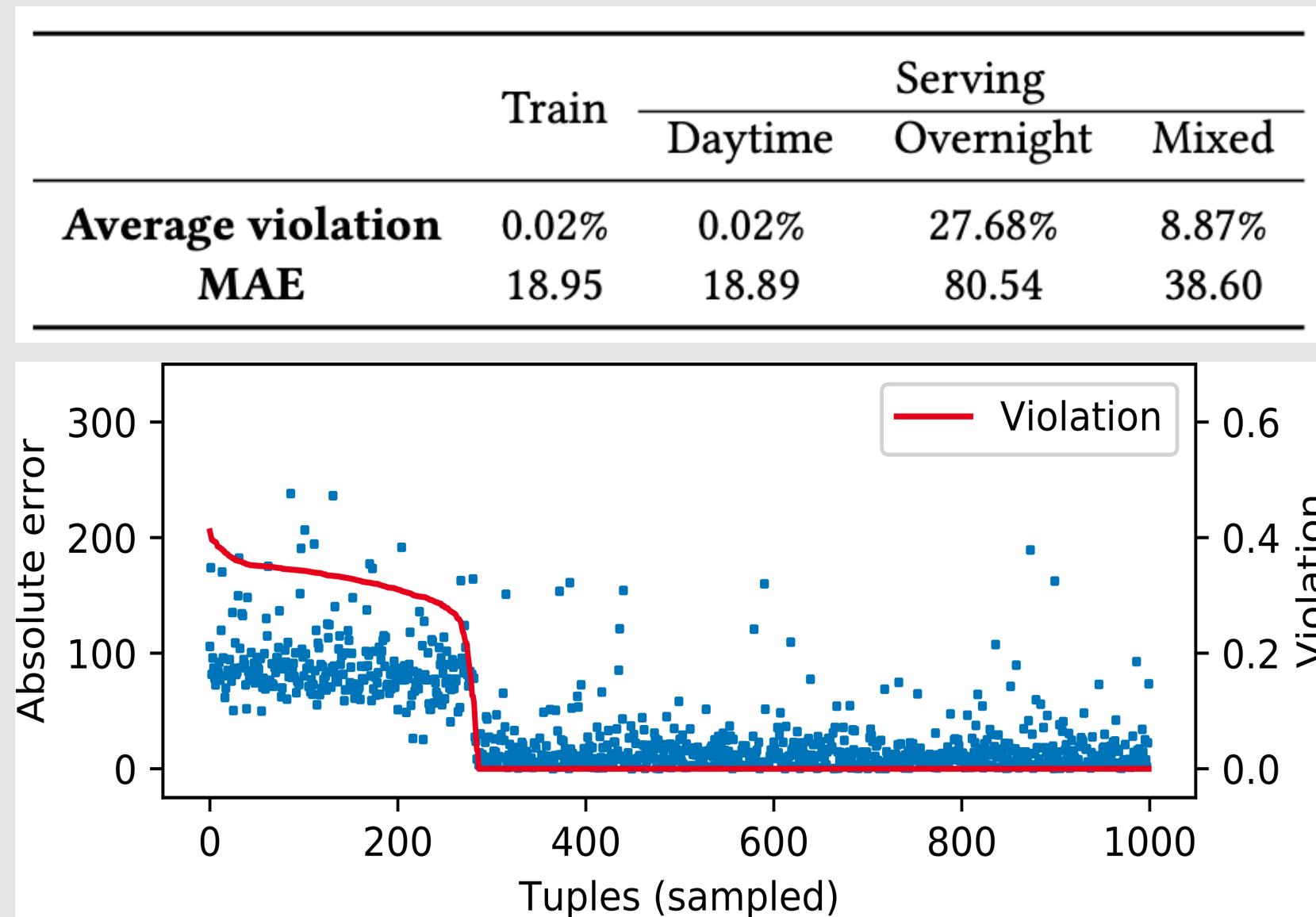
Complexity analysis

- Runtime
 - Linear in number of tuples in the dataset
 - Cubic in number of attributes
 - Highly parallelizable
- Memory
 - Quadratic in number of attributes

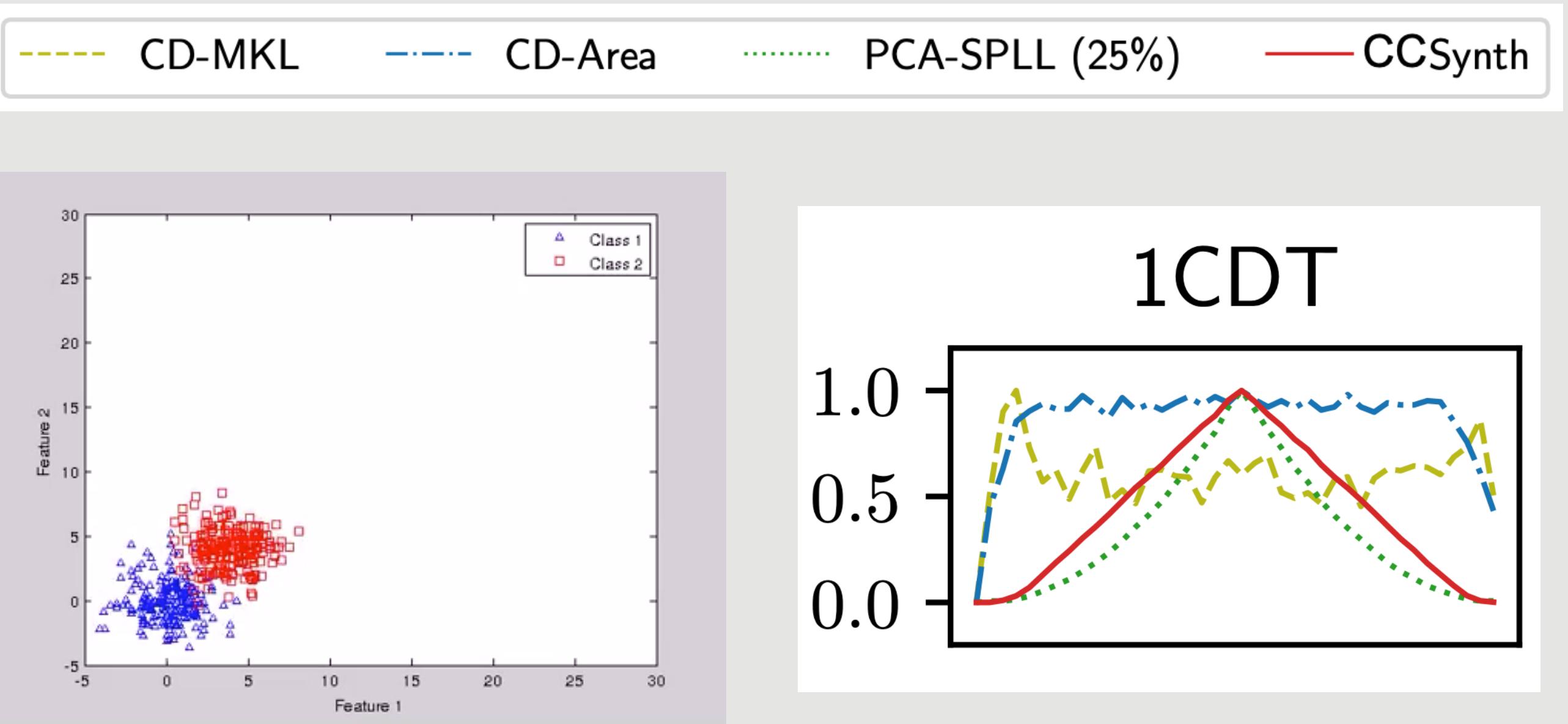
Experimental results: two applications

- Trusted Machine Learning
 - Is there a relationship between CC violation and the ML model's prediction accuracy?
- Data-drift
 - Can CCs be used to quantify data drift?

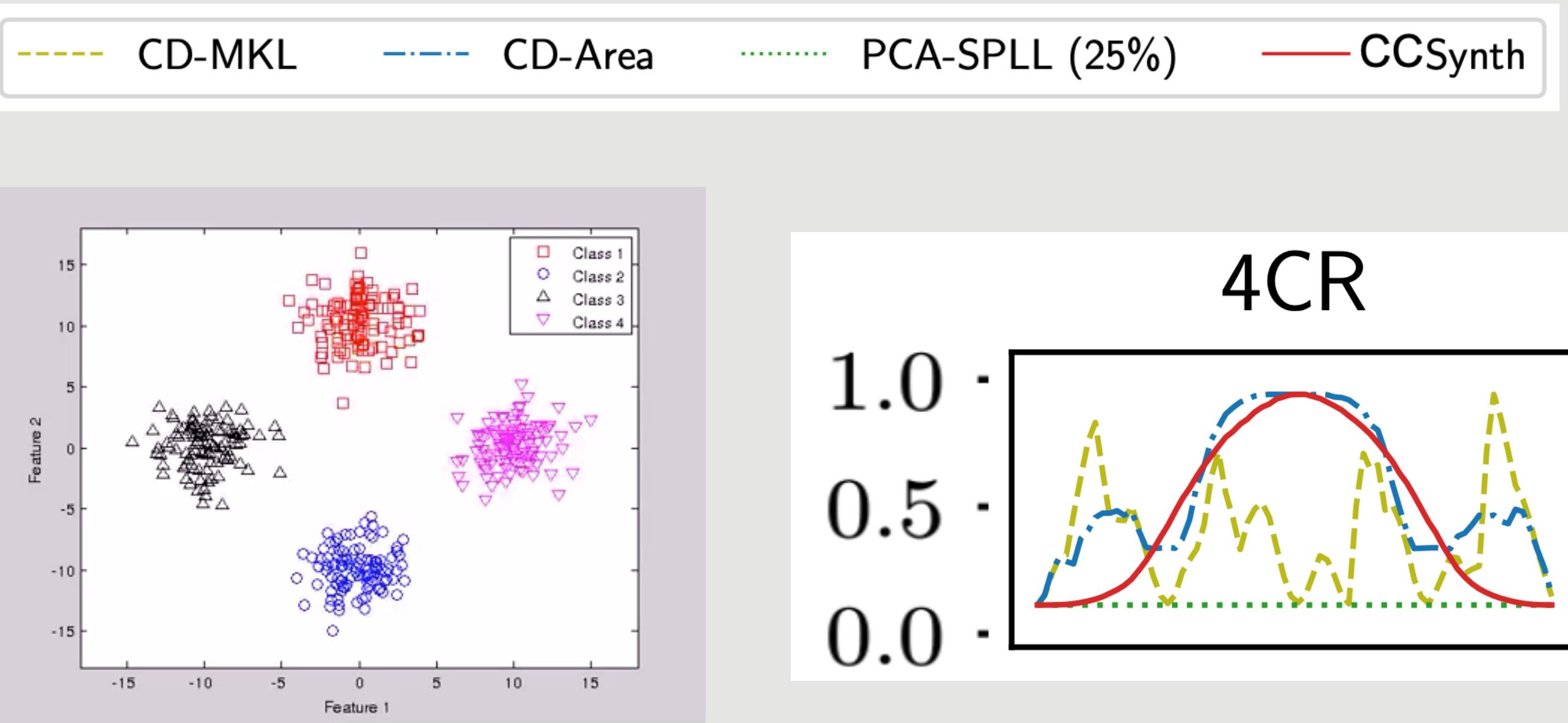
Trusted machine learning: airlines dataset



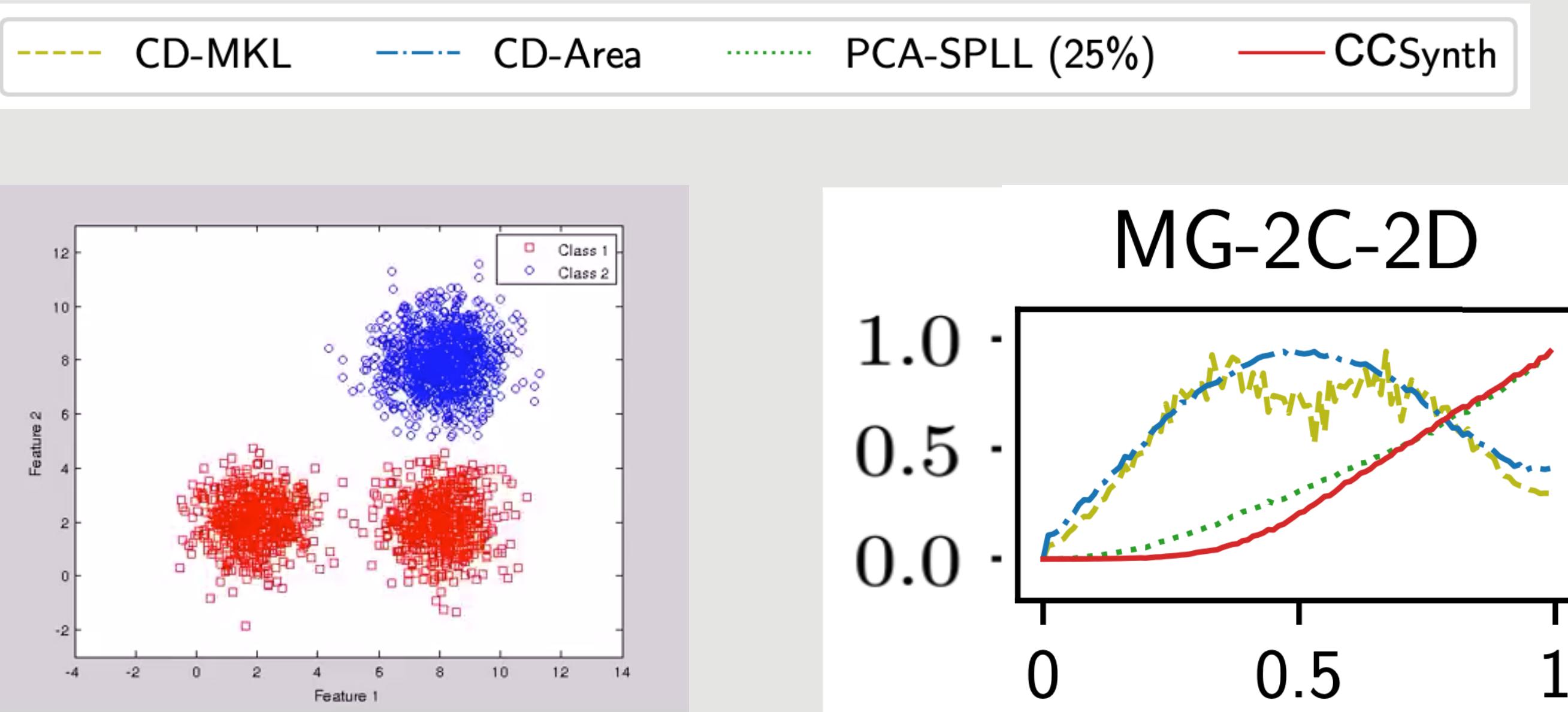
Data drift: EVL benchmark (1/4)



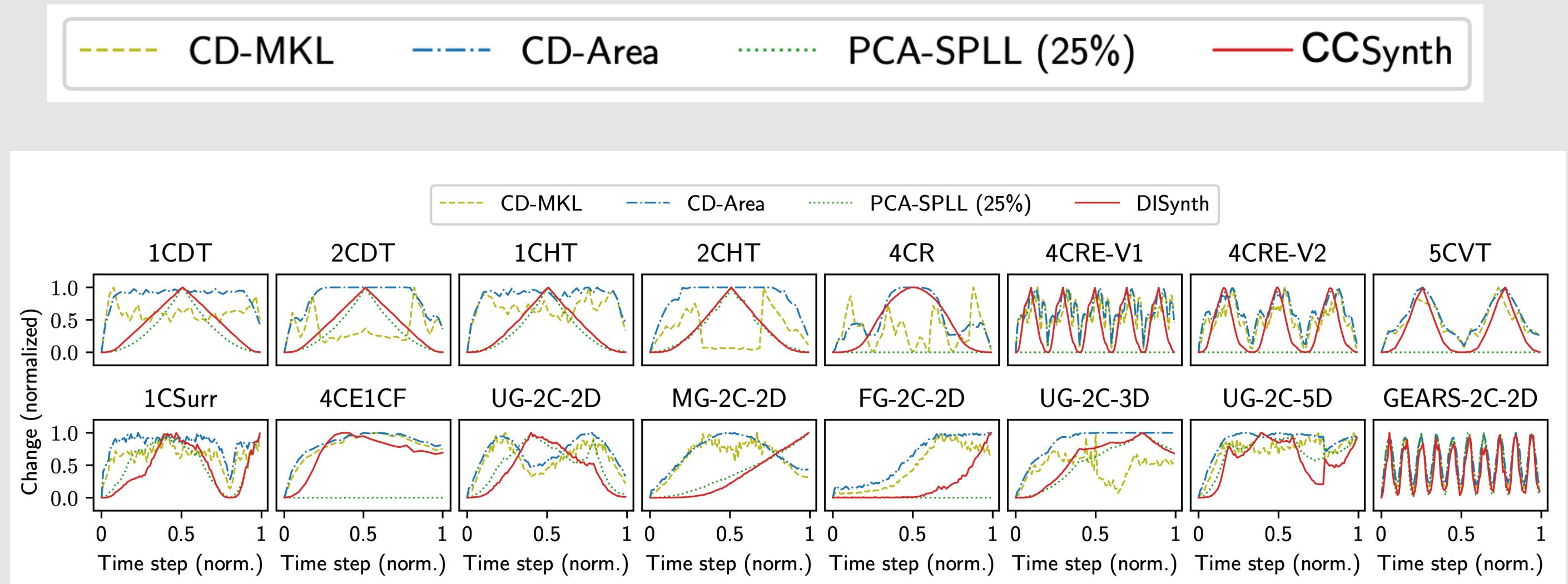
Data drift: EVL benchmark (2/4)



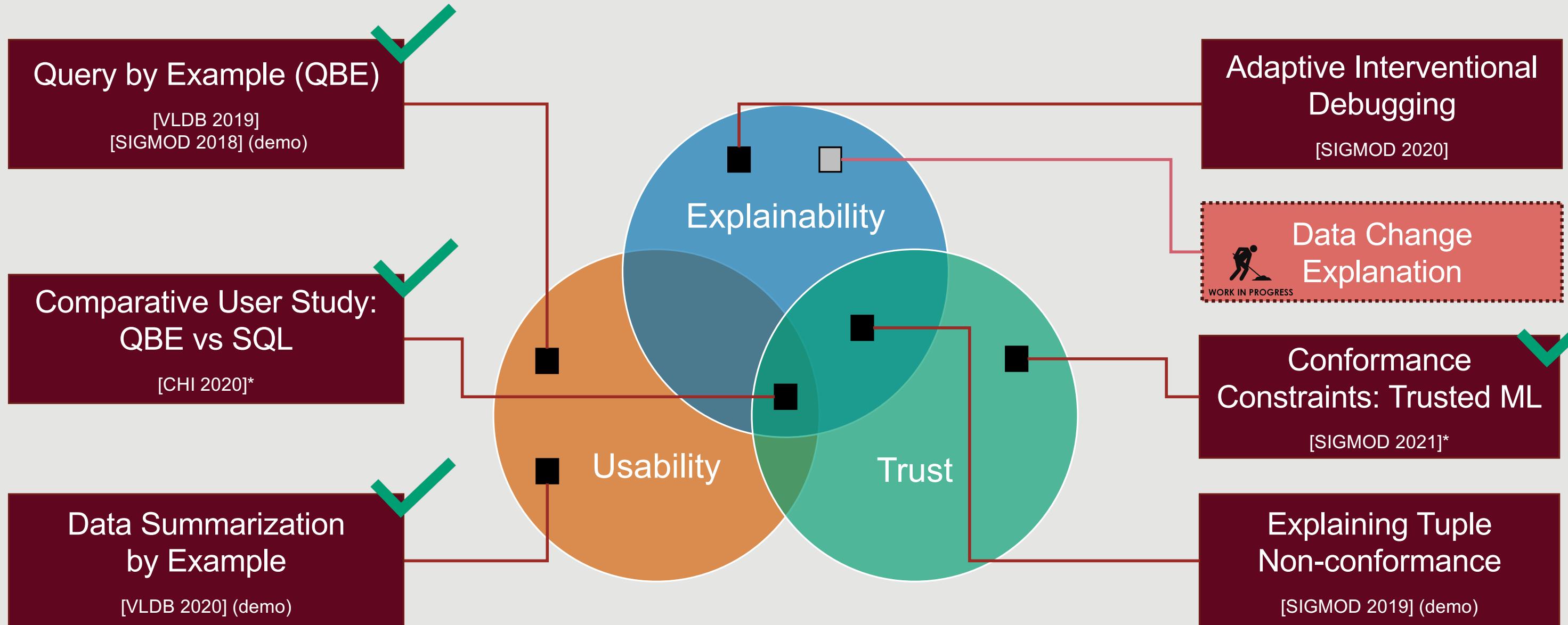
Data drift: EVL benchmark (3/4)



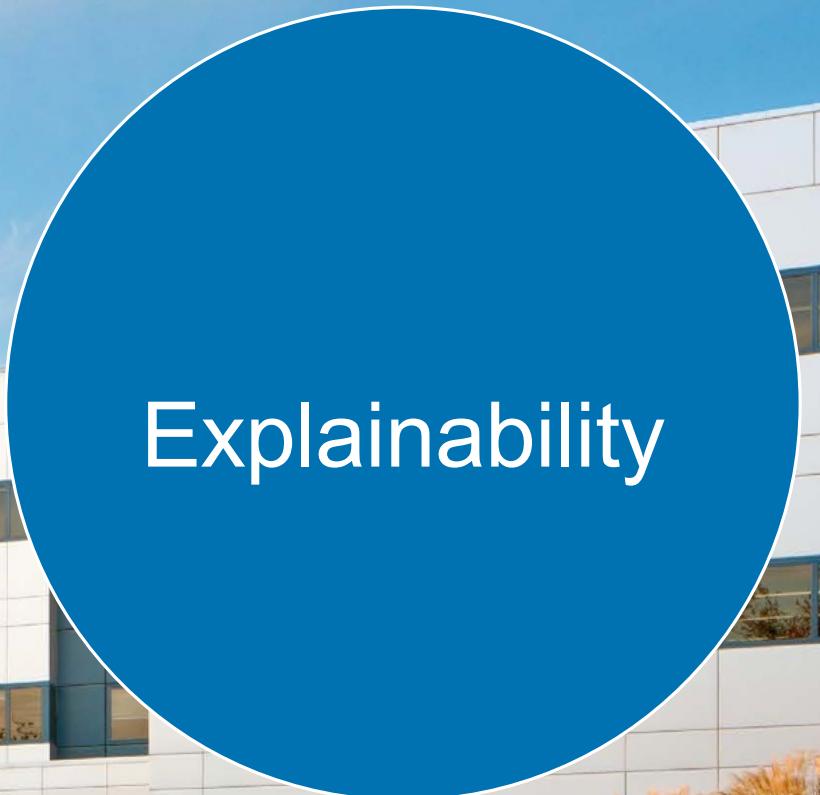
Data drift: EVL benchmark (4/4)



Dissertation outline

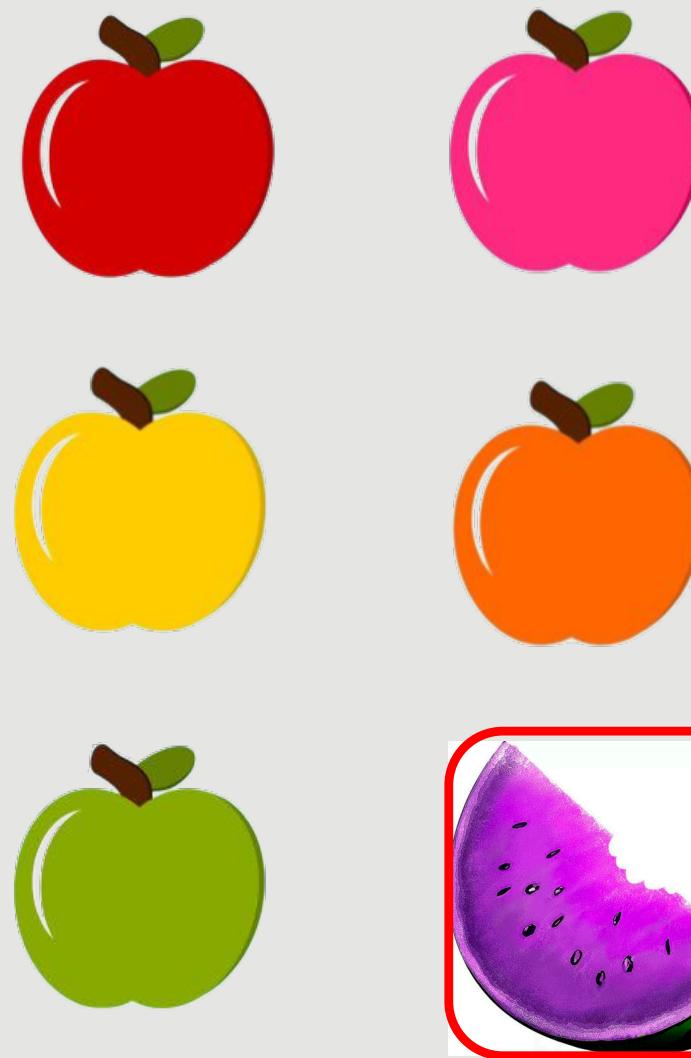
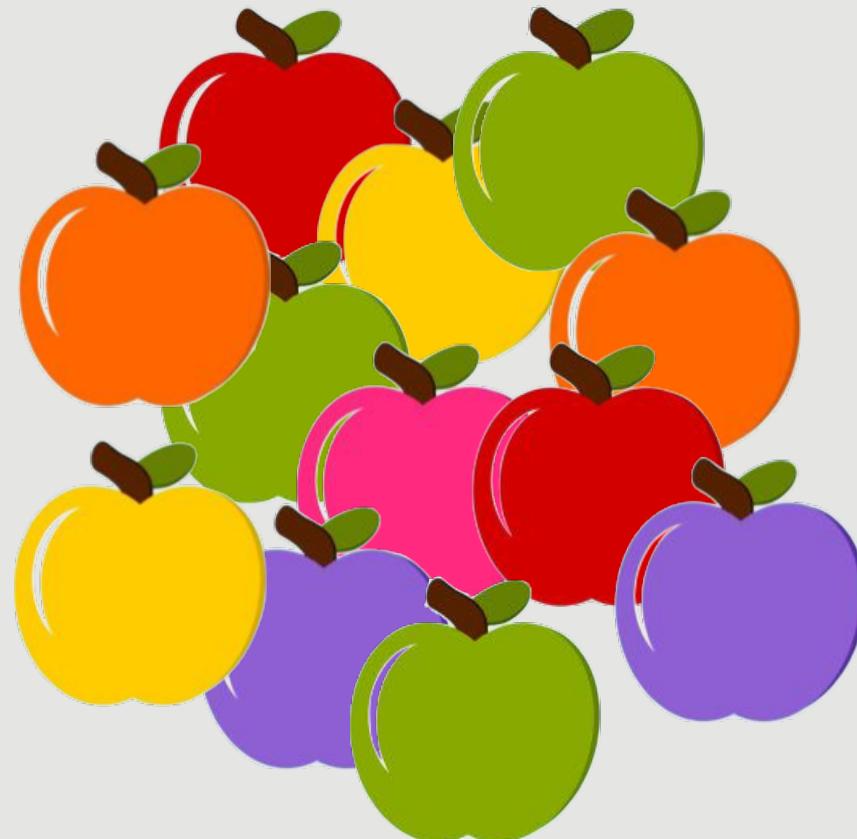


Part 3: Explanation Frameworks

A large, solid blue circle is positioned on the right side of the slide, partially overlapping the building image. It contains the word "Explainability" in a white, sans-serif font.

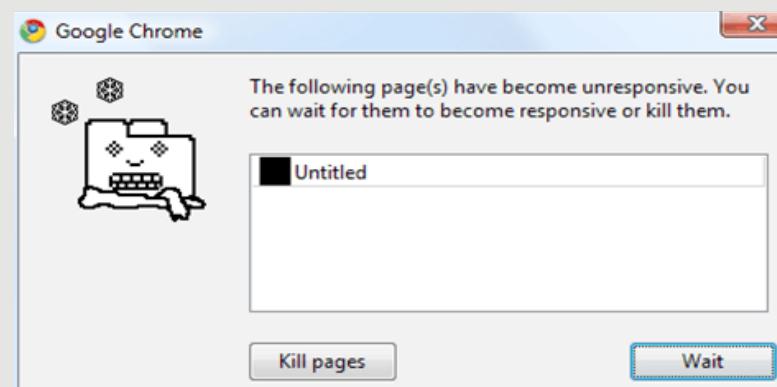
Explainability

Why ML models fail for certain tuples?



**How is this
different?**

Why do systems (sometimes) behave unexpectedly?



Why did the
system crash?



ExTuNe

Explaining Tuple Non-conformance

Conformance constraints

Detection

Is it
non-conforming?

ExTuNe

Explanation



Why is it
non-conforming?

Tuple-level explanation



Tuple-level explanation

Height	Weight	BMI
6 feet	142 lbs	19.3
5 feet	170 lbs	33.2
5 feet	130 lbs	25.4
6 feet	170 lbs	231

Tuple-level explanation

Height	Weight	BMI
6 feet	142 lbs	19.3
5 feet	170 lbs	33.2
5 feet	130 lbs	25.4
6 feet	170 lbs	231

Intervention reveals causality

Height	Weight	BMI
6 feet	142 lbs	19.3
5 feet	170 lbs	33.2
5 feet	130 lbs	25.4
6 feet	170 lbs	25.9

Mean = 25.9

change

Intervention reveals causality

Height	Weight	BMI
6 feet	142 lbs	19.3
5 feet	170 lbs	33.2
5 feet	130 lbs	25.4
6 feet	170 lbs	25.9

Blame!

Intervention reveals causality

Height	Weight	BMI
6 feet	142 lbs	19.3
5 feet	170 lbs	33.2
5 feet	130 lbs	25.4
16 feet	70 lbs	25.9

Blame!

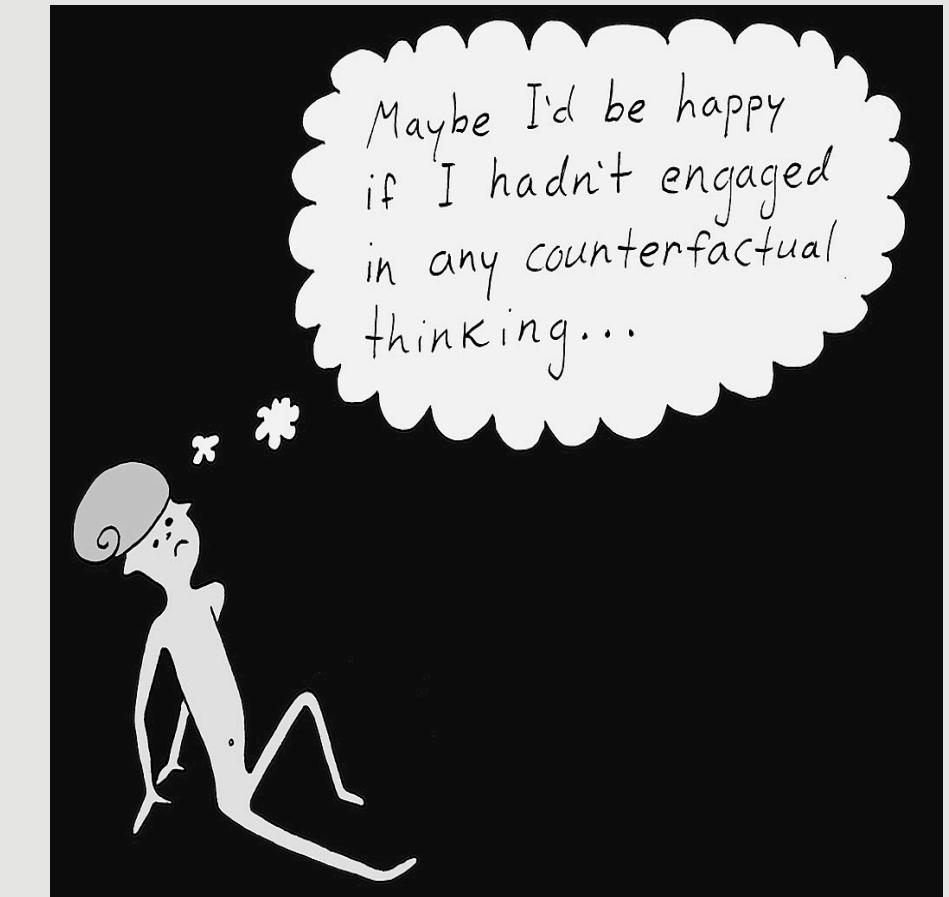
Blame!

ExTuNe principles: actual causality

Actual Causality

“When **K** other events are removed, then **C** is a **counterfactual** cause of **E**”

- **C** is an actual cause of **E**
- **C**’s responsibility is $1/(K + 1)$

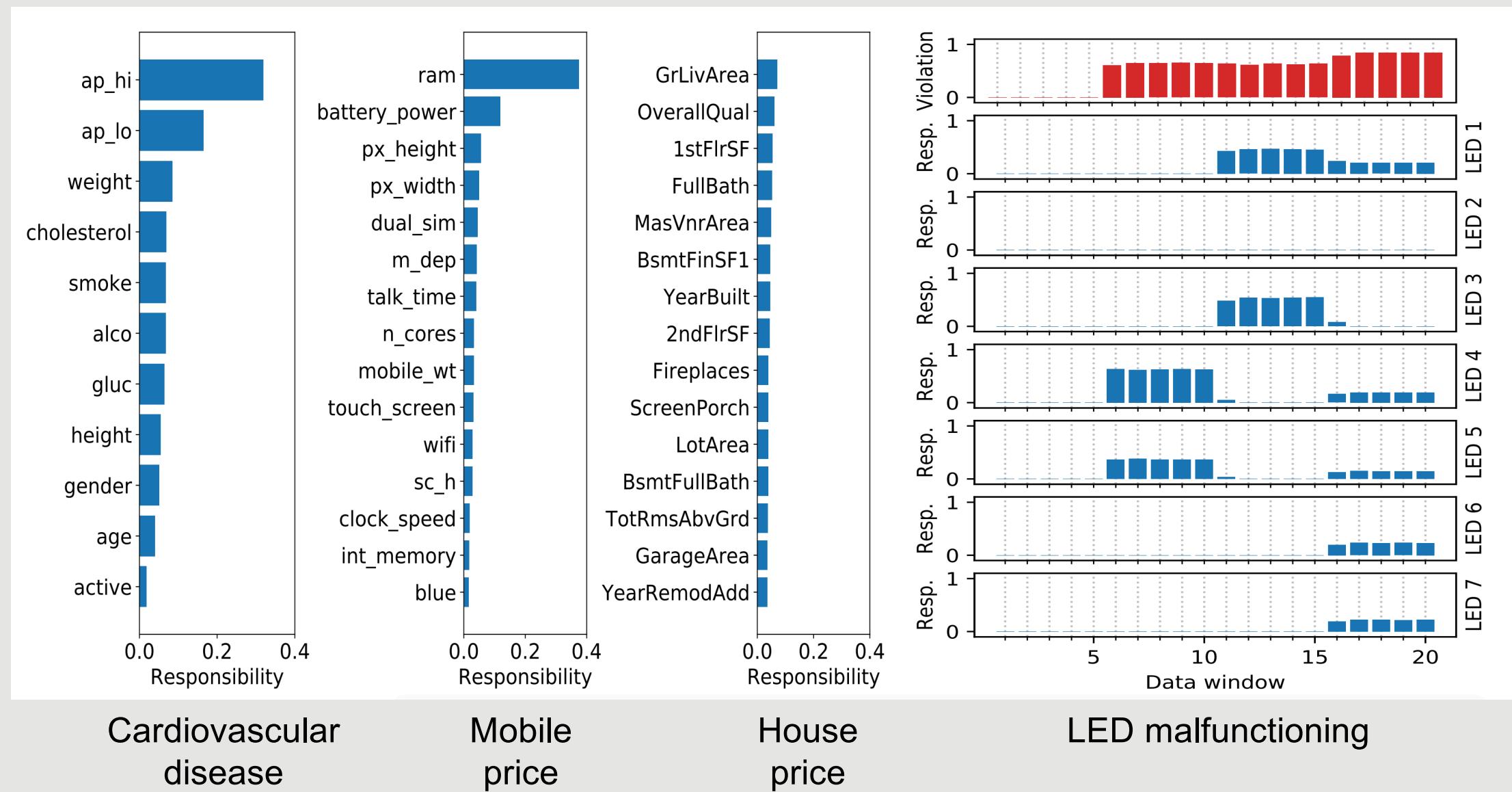


Maybe I'd be happy
if I hadn't engaged
in any counterfactual
thinking...

ExTuNe interface



ExTuNe evaluation: case studies



Anomaly in COVID dataset

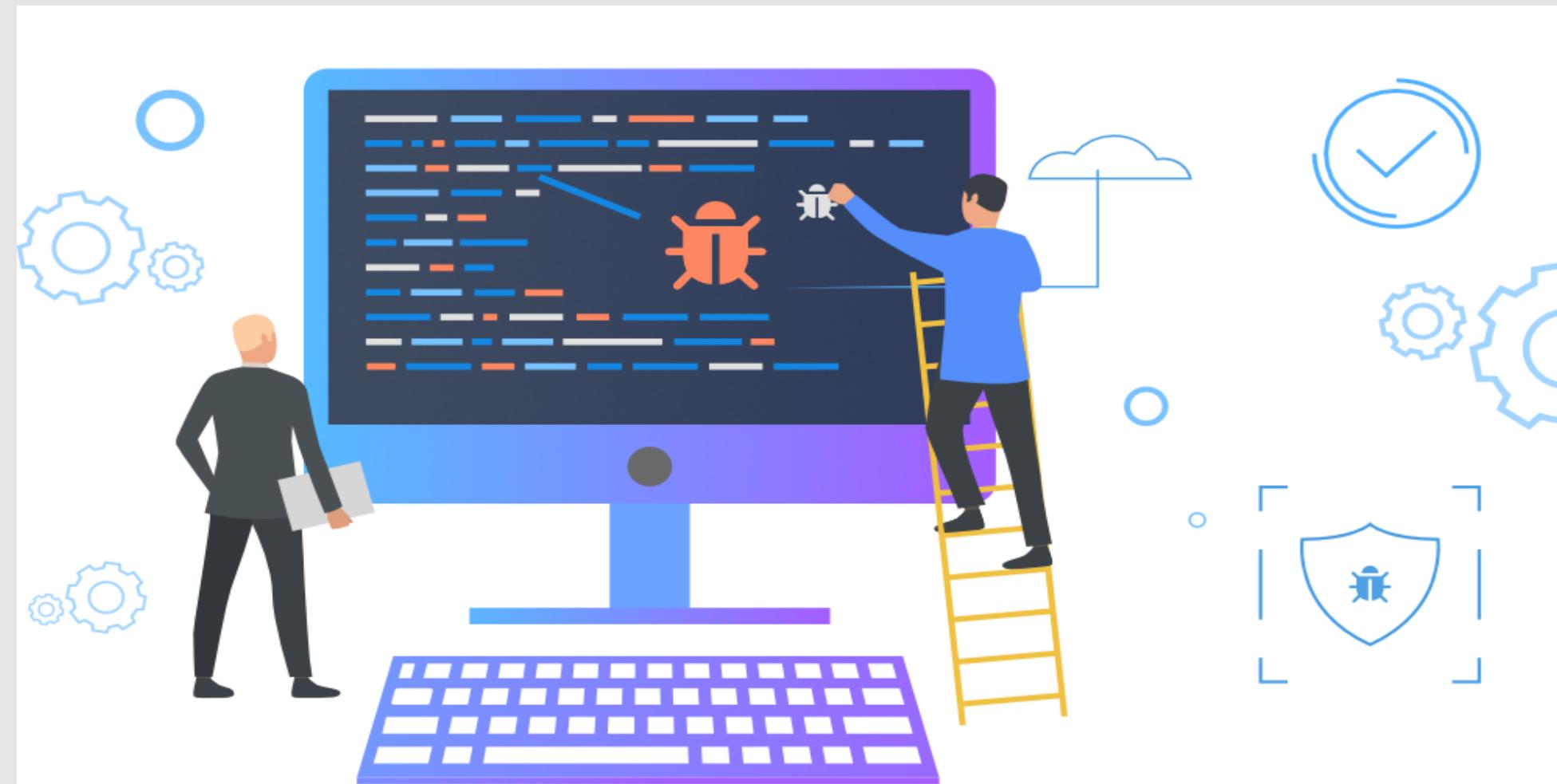
Conformance constraint: #positive + #negative = #total

Violation		date	state	positive	negative	pending	hospitalized	death	total	population	hospital_beds
201	0.350000	20200321	NY	10356	35081	0	1603	44	45437	19453561	52524
36	0.350000	20200324	NY	25665	65605	0	3234	210	91270	19453561	52524
256	0.290000	20200320	NY	7102	25325	0	0	35	32427	19453561	52524
4	0.240000	20200324	CA	2102	13452	12100	0	40	27654	39512223	71122
59	0.240000	20200323	CA	1733	12567	12100	0	27	26400	39512223	71122
311	0.220000	20200319	NY	4152	18132	0	0	12	22284	19453561	52524
91	0.200000	20200323	NY	20875	57414	0	2635	114	78289	19453561	52524
88	0.190000	20200323	NJ	2844	359	94	0	27	3297	8882190	21317
471	0.130000	20200316	NJ	178	120	20	0	2	218	8882190	21317
389	0.130000	20200317	CA	483	7981	0	0	11	8407	39512223	71122
51	0.130000	20200324	WA	2221	31712	0	0	110	33933	7614893	12945

$$178 + 120 \\ = 298$$

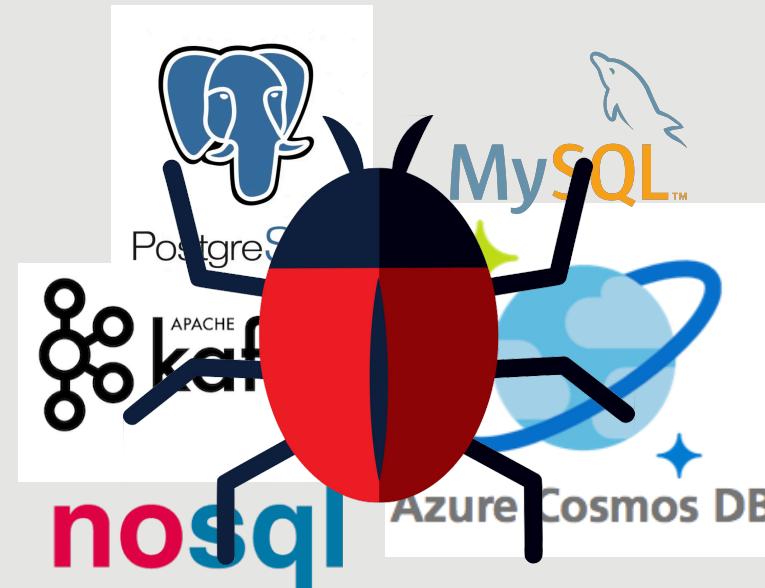
$$178 + 120 \\ \neq 218$$

Explaining data systems' failure

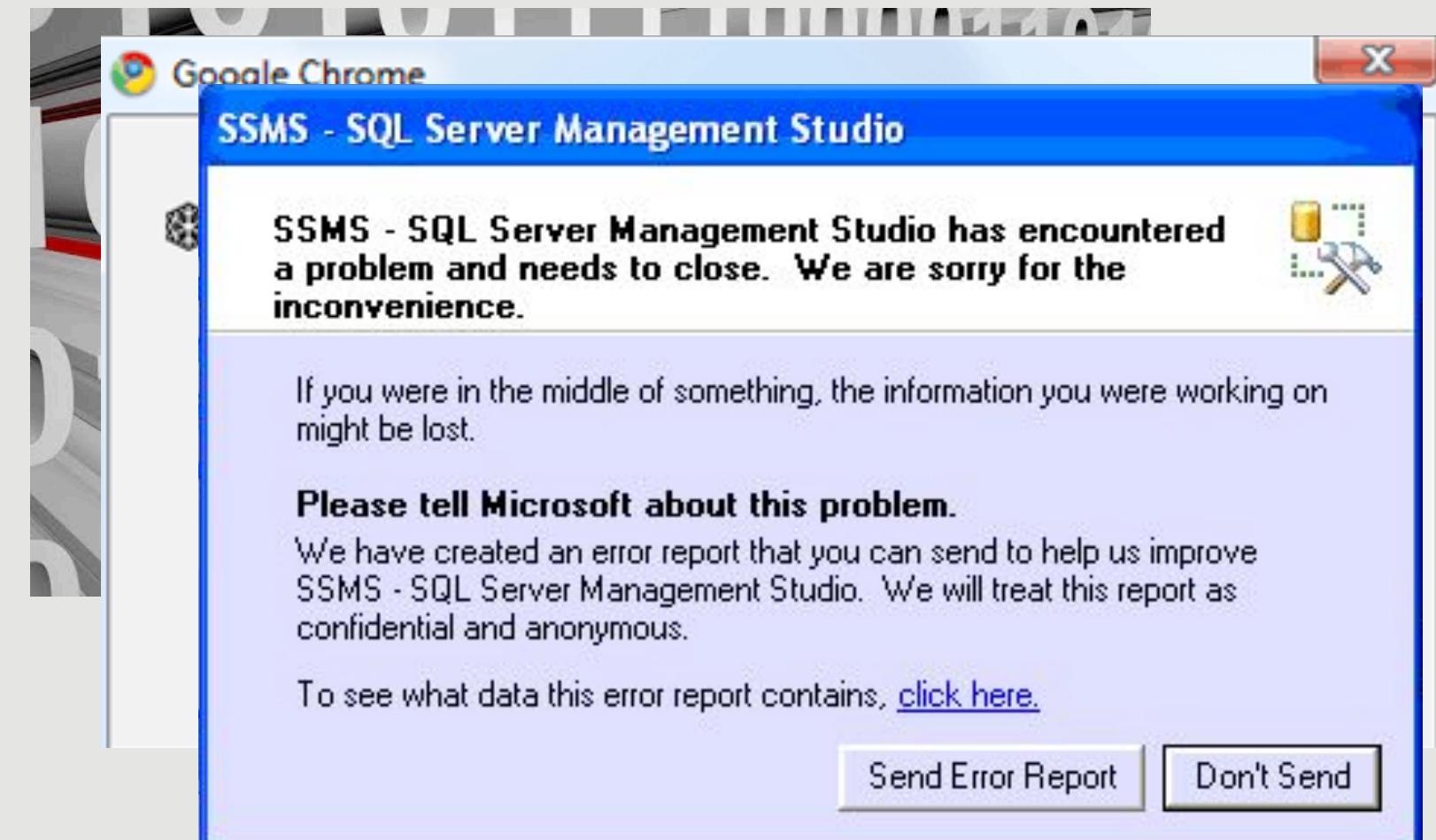


AID: Causality-guided Adaptive Interventional Debugging

DBMS are complex and contain bugs

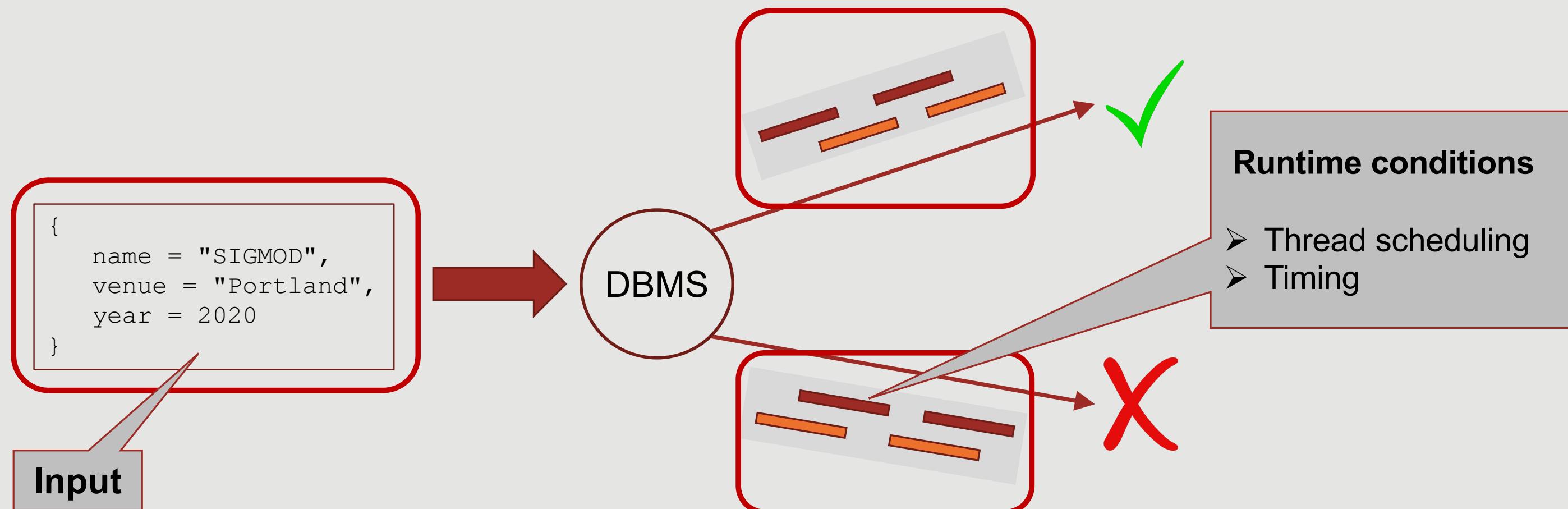


- concurrent
- parallel
- asynchronous

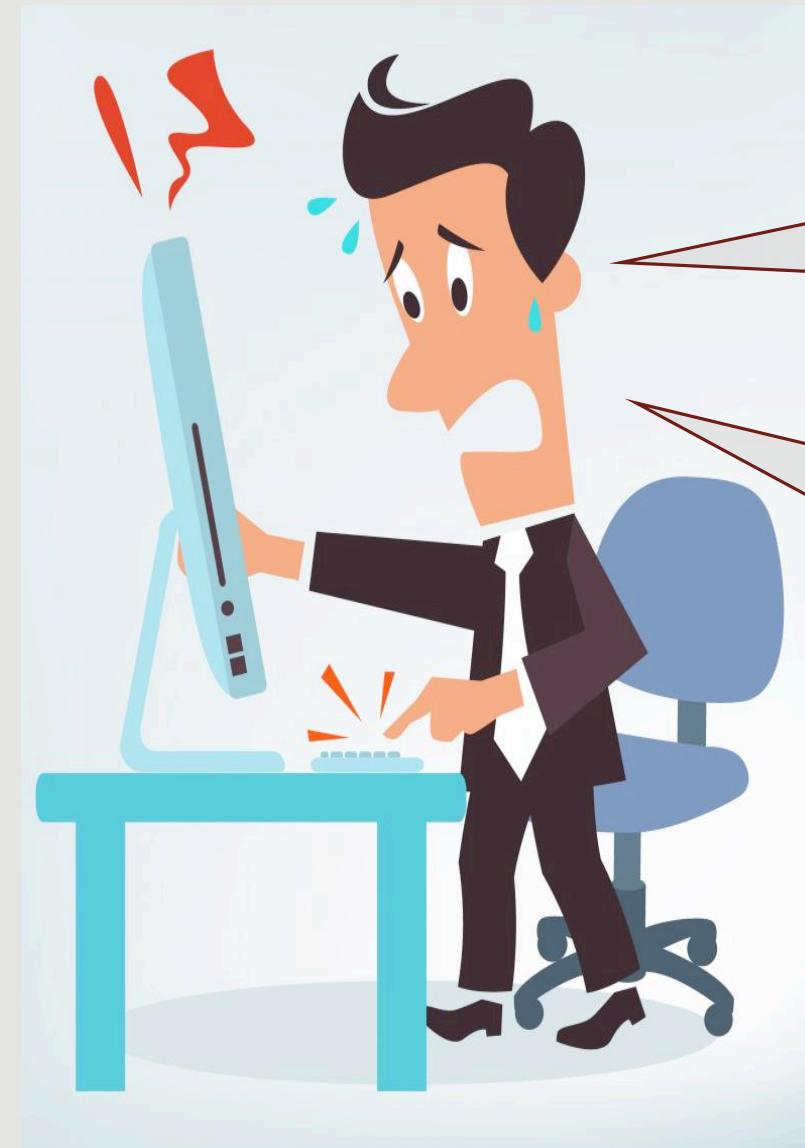


Intermittent failure

sometimes succeeds sometimes fails



Motivation and goal



Investigate root causes
of intermittent failure

Npgsql intermittent failure

[ADO.NET data provider for PostgreSQL]

The screenshot shows a GitHub issue page for the 'npgsql / npgsql' repository. The issue is titled 'Race condition in PoolManager.TryGetValue #2485'. It is marked as 'Closed' by 'thetranman' on May 29, 2019, with 3 comments. The issue details section includes a 'Steps to reproduce' section with instructions to fill in connection string values for a VolatileTest, and a 'The issue' section mentioning related issue #2146 and production code issues with NpgsqlConnection. The right sidebar displays project metadata: Assignees (thetranman), Labels (bug), Projects (None yet), Milestone (4.0.8), and Linked pull requests.

npgsql / npgsql

Used by 10.7k Watch 175 Star 2k Fork 628

Code Issues 167 Pull requests 33 Actions Security 0 Insights New issue

Race condition in PoolManager.TryGetValue #2485

Closed thetranman opened this issue on May 29, 2019 · 3 comments

thetranman commented on May 29, 2019

Contributor

...

Steps to reproduce

I've created a test that can reproduce the issue. All you have to do is fill in the values for the connection string. The test is VolatileTest as seen here:
<https://github.com/thetranman/npgsql/pull/1/files>

The issue

Could be related to: #2146

In our production code, we are running into issues when trying to create a new Postgres connection (Specifically when we call: var connection = new NpgsqlConnection(ConnectionString);).

This can intermittently occur when we are trying to start our service on a server which can contain

Assignees

thetranman

Labels

bug

Projects

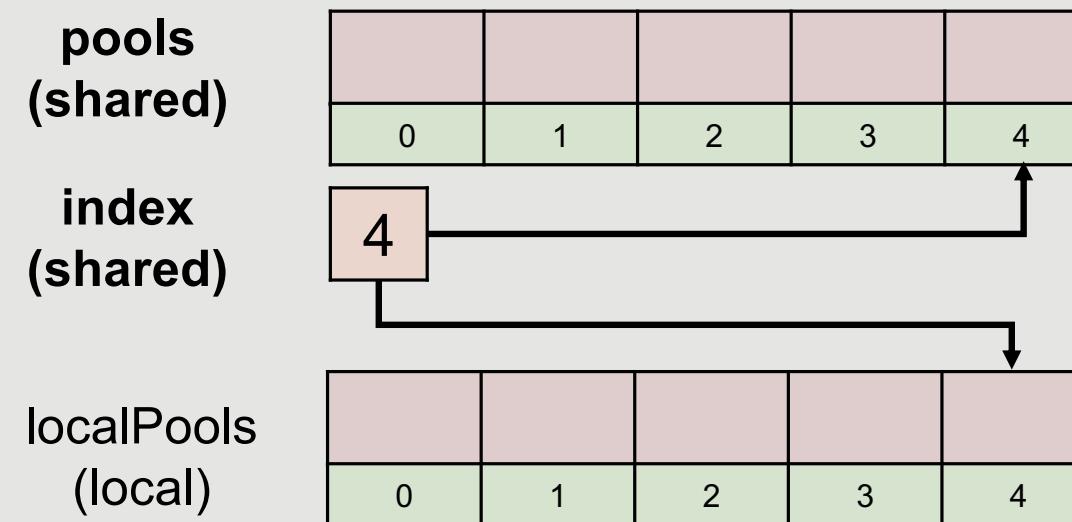
None yet

Milestone

4.0.8

Linked pull requests

Npgsql intermittent failure



Thread 1

Find(key):

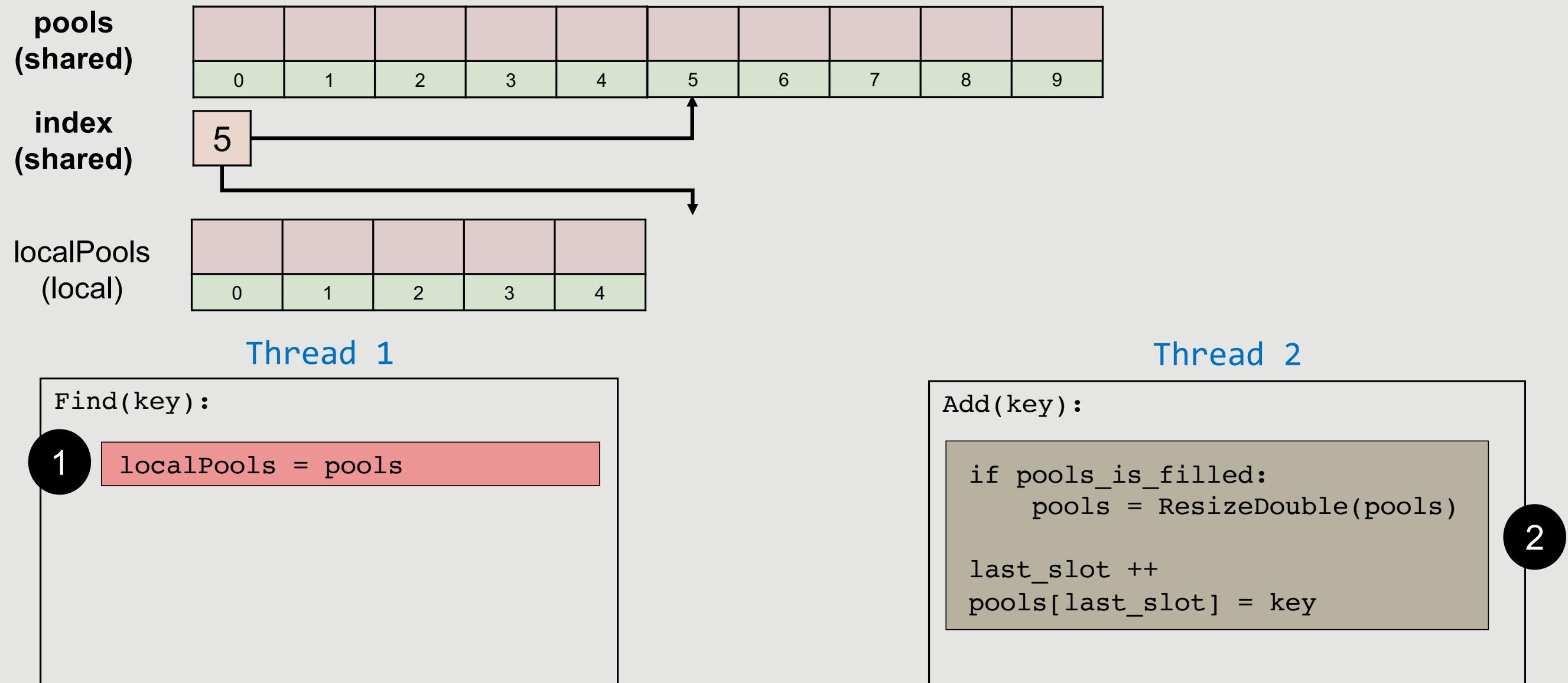
1

localPools = pools

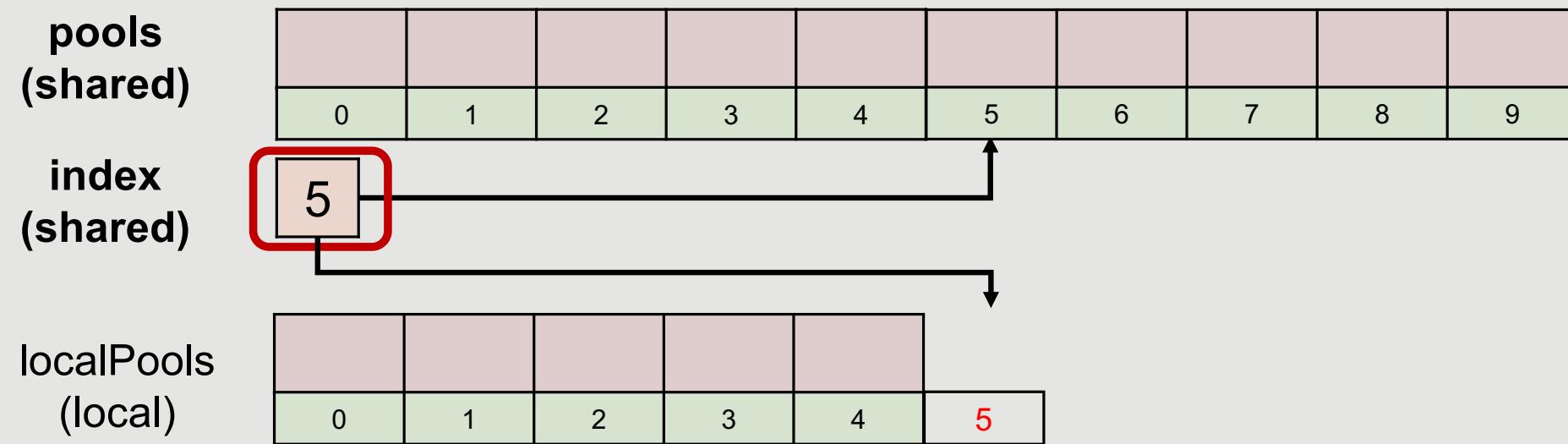
Thread 2

Add(key):

Npgsql intermittent failure



Npgsql intermittent failure



Thread 1

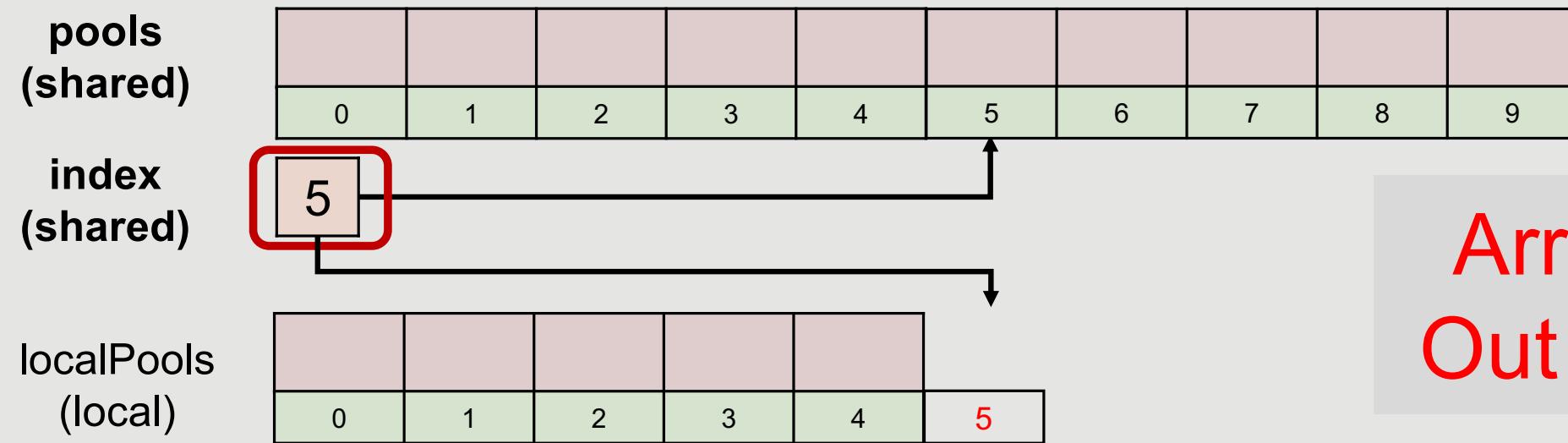
```
Find(key):  
1   localPools = pools  
  
for i in range(0,last_slot+1):  
    if (localPools[i] == key)  
        return i  
  
return null
```

Thread 2

```
Add(key):  
if pools_is_filled:  
    pools = ResizeDouble(pools)  
  
last_slot ++  
pools[last_slot] = key
```

2

Npgsql intermittent failure



Array Index
Out of Bound

Thread 1

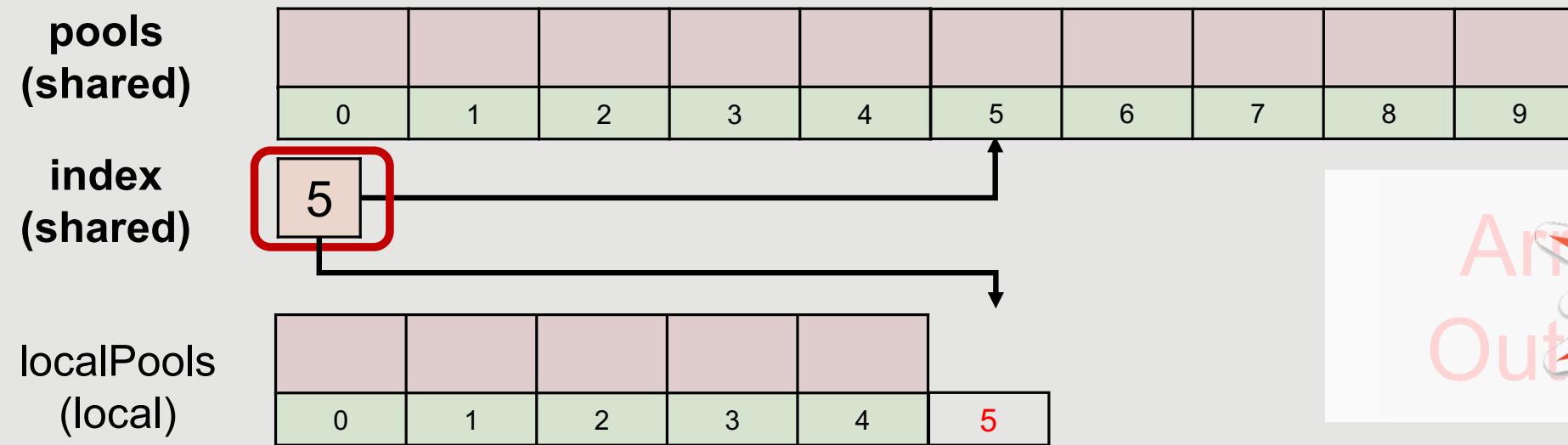
```
Find(key):  
1   localPools = pools  
  
for i in range(0,last_slot+1):  
    if (localPools[i] == key)  
        return i  
  
return null
```

Thread 2

```
Add(key):  
  
if pools_is_filled:  
    pools = ResizeDouble(pools)  
  
last_slot ++  
pools[last_slot] = key
```

2

Npgsql intermittent failure



Thread 1

```
Find(key):  
1   localPools = pools  
  
for i in range(0,last_slot+1):  
    if (localPools[i] == key)  
        return i  
  
return null
```

Thread 2

```
Add(key):  
  
if pools_is_filled:  
    pools = ResizeDouble(pools)  
  
last_slot ++  
pools[last_slot] = key
```

2

Npgsql intermittent failure

1

```
localPools = pools
```

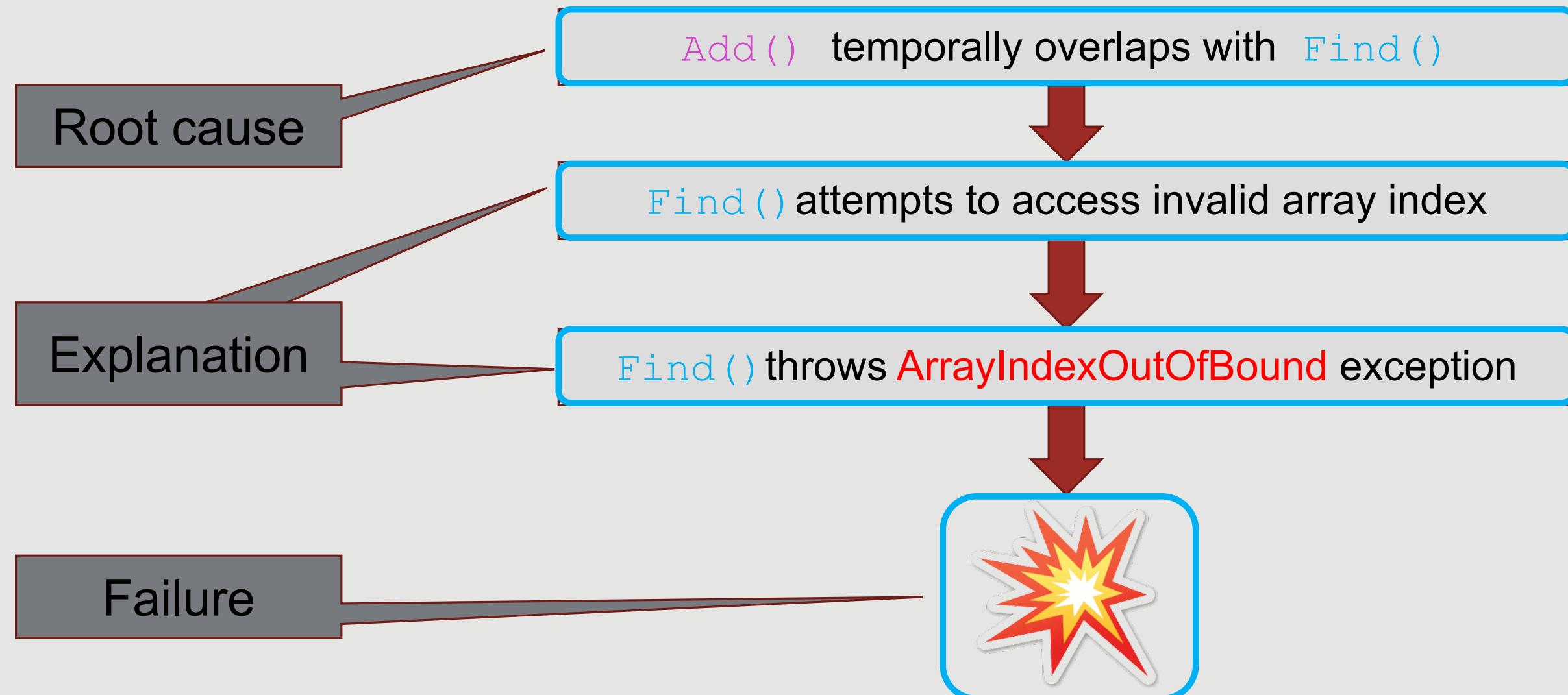
2

```
if pools_is_filled:  
    pools = ResizeDouble(pools)  
  
last_slot ++  
pools[last_slot] = key
```

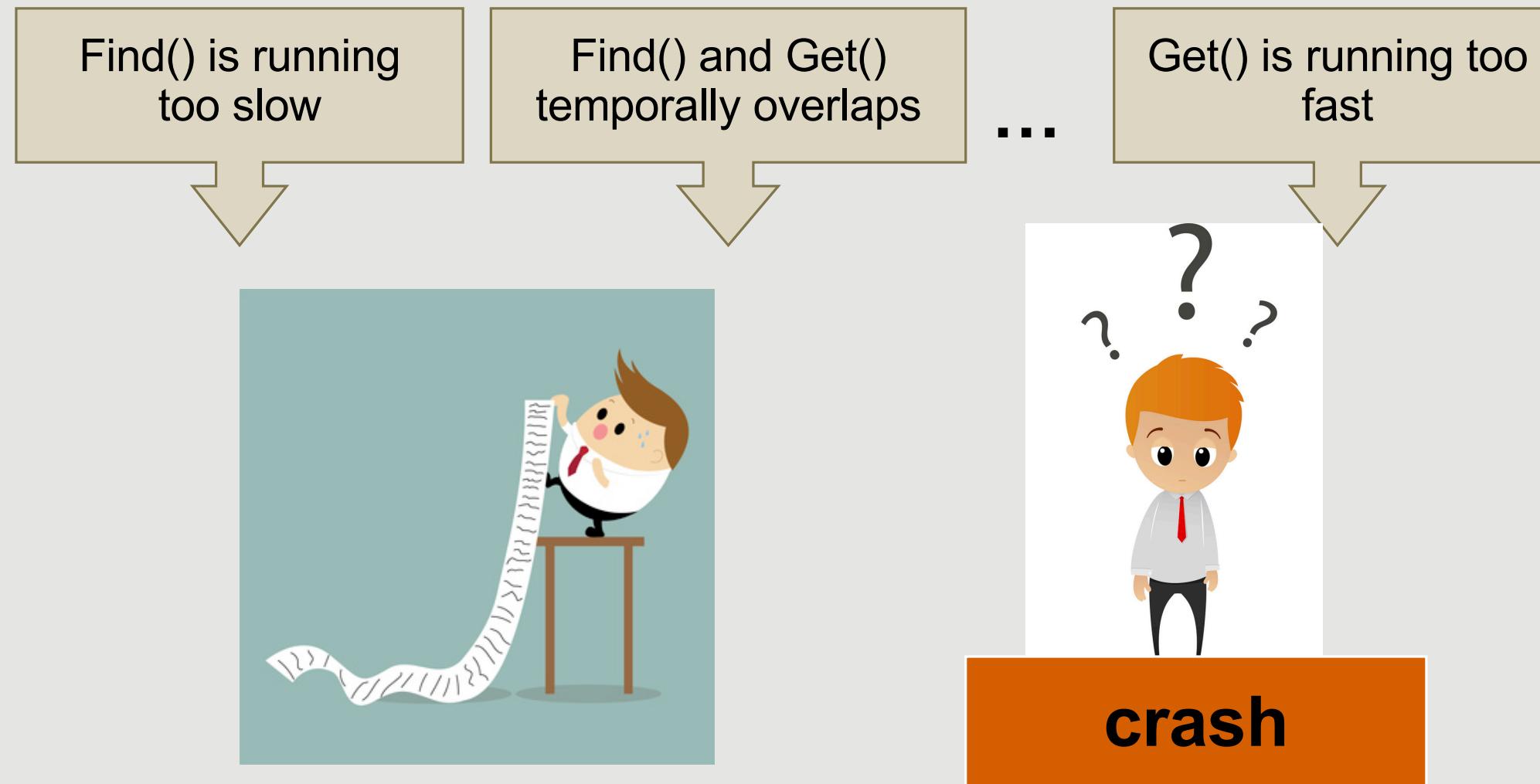
3

```
for i in range(0,last_slot+1):  
    if (localPools[i] == key)  
        return i  
  
return null
```

Investigating Npgsql crash



Limitations of statistical debugging



Our goals

Root-cause identification

Find() and Get()
temporally overlaps

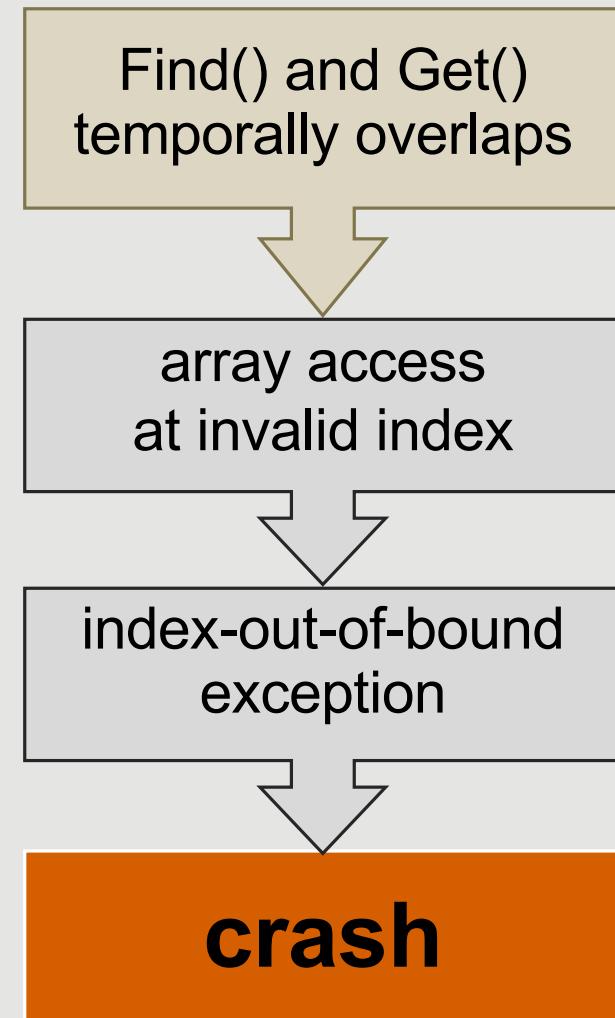


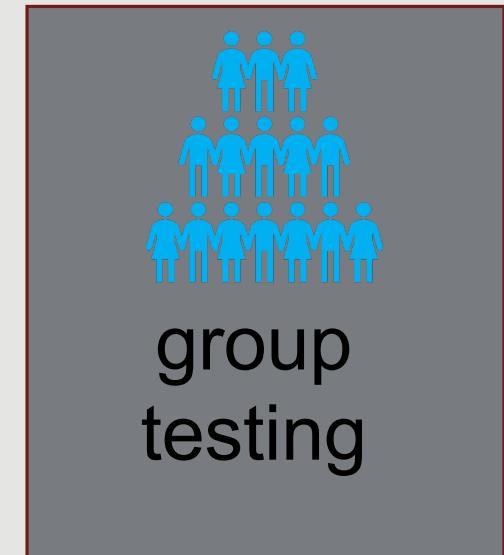
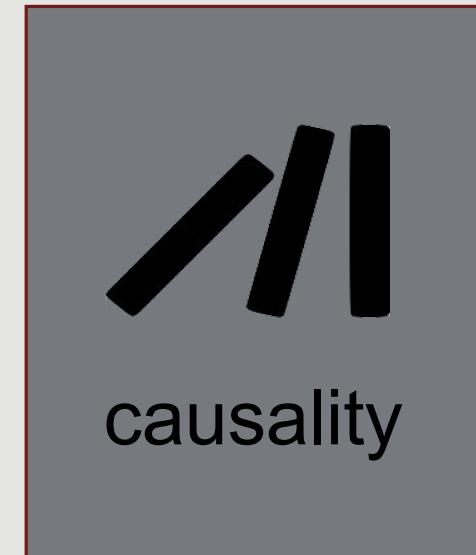
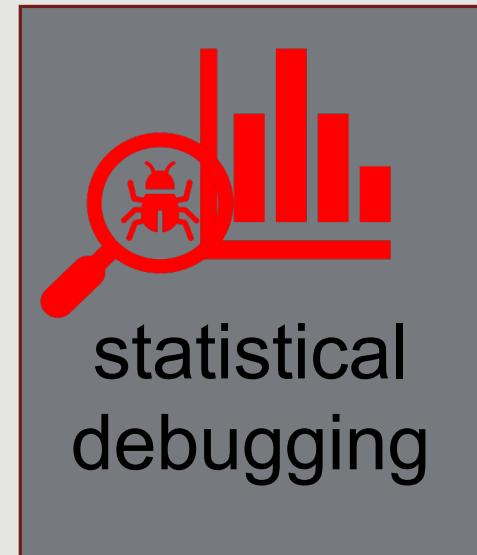
crash

Our goals

Root-cause identification

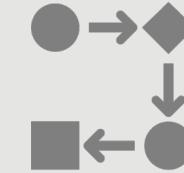
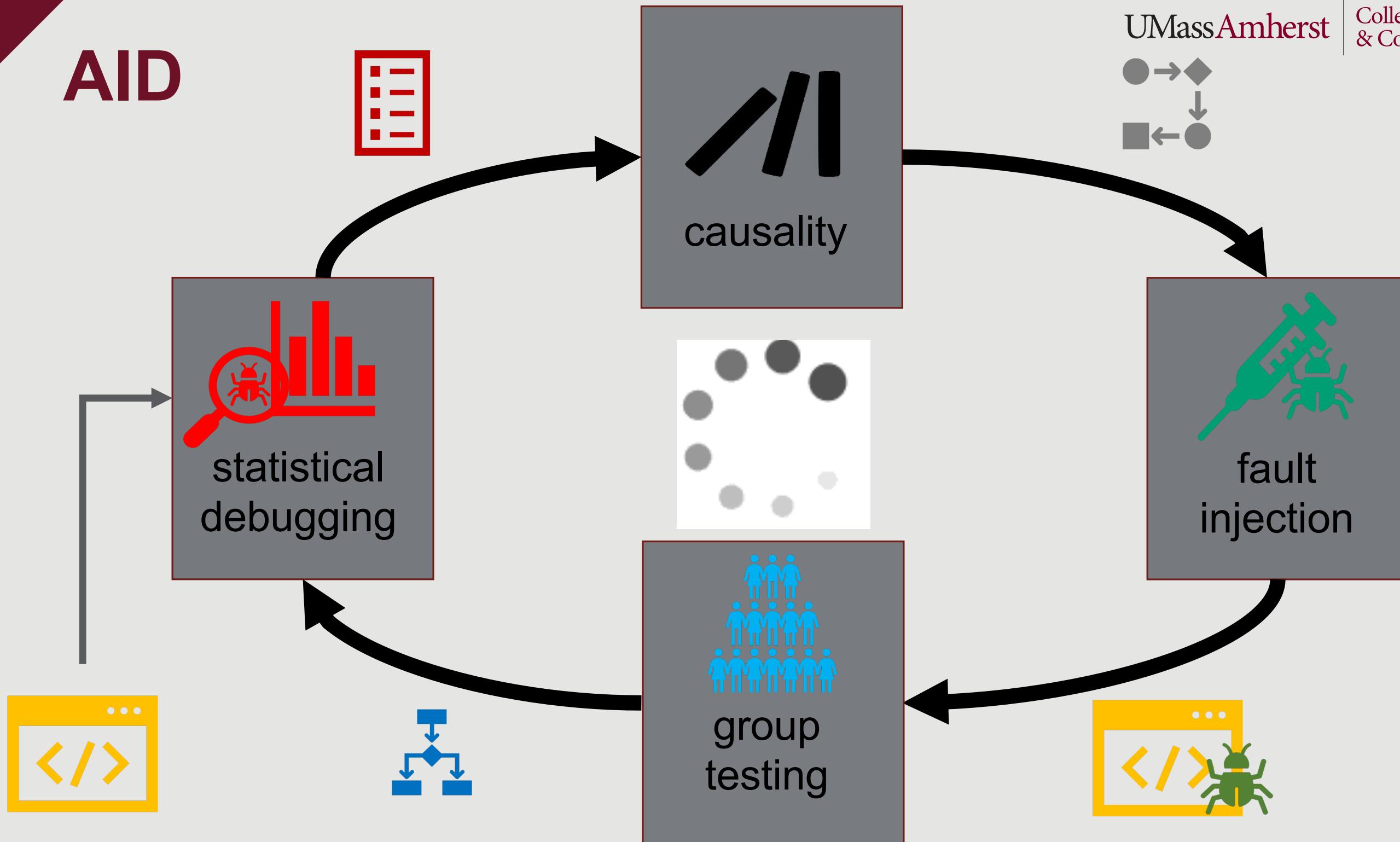
Explanation





AID: Adaptive Interventional Debugging

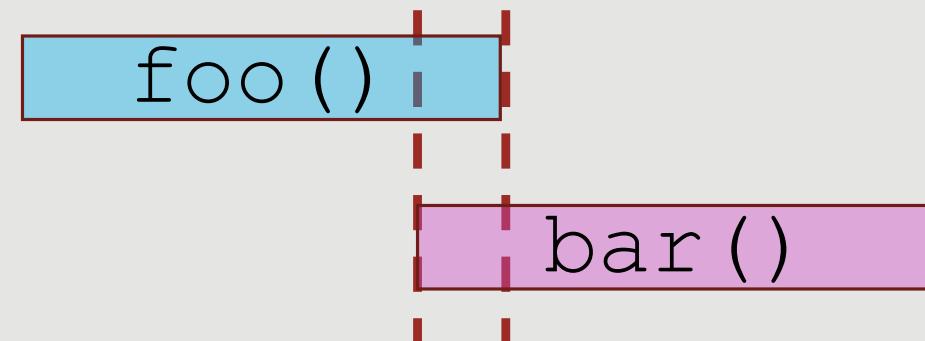
AID



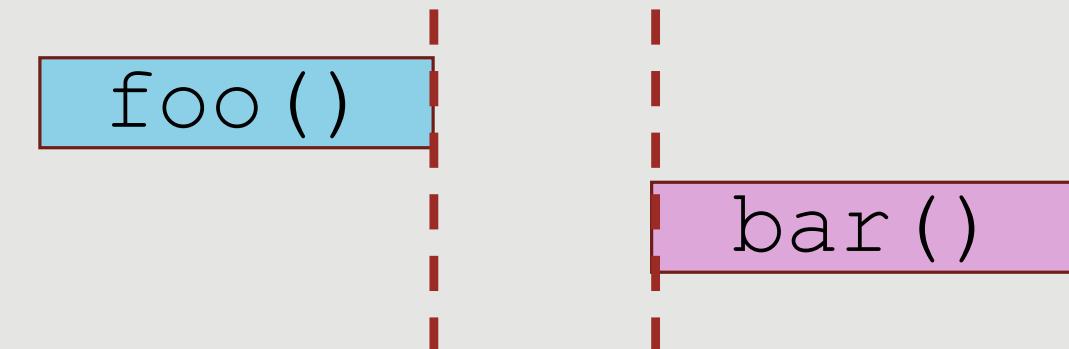
Finding candidate predicates

- Step 1: Program instrumentation finds all predicates
- Step 2: Statistical debugging finds correlated predicates

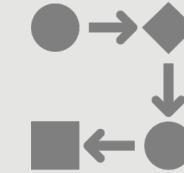
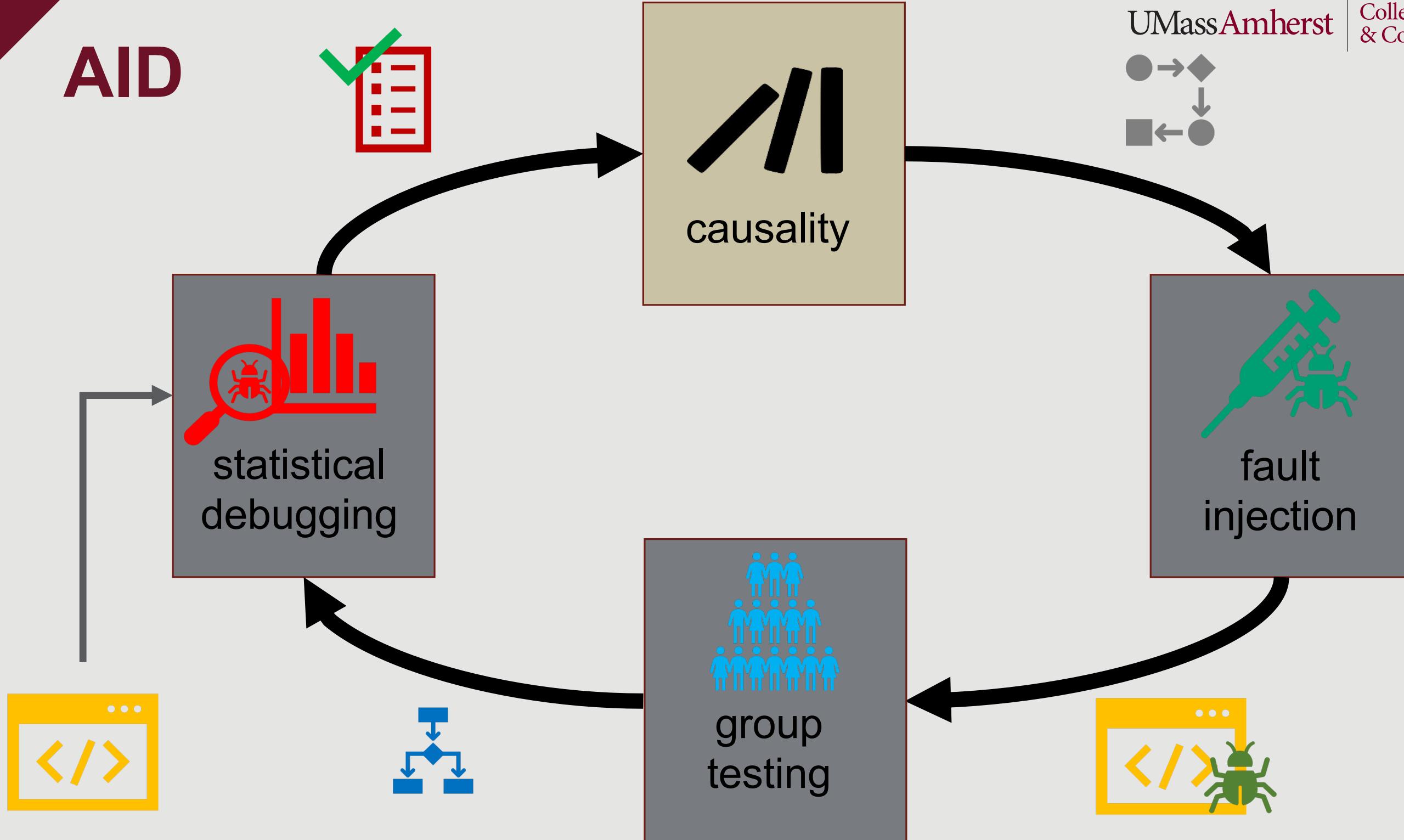
Always appear in
failed executions



Never appear in
successful executions

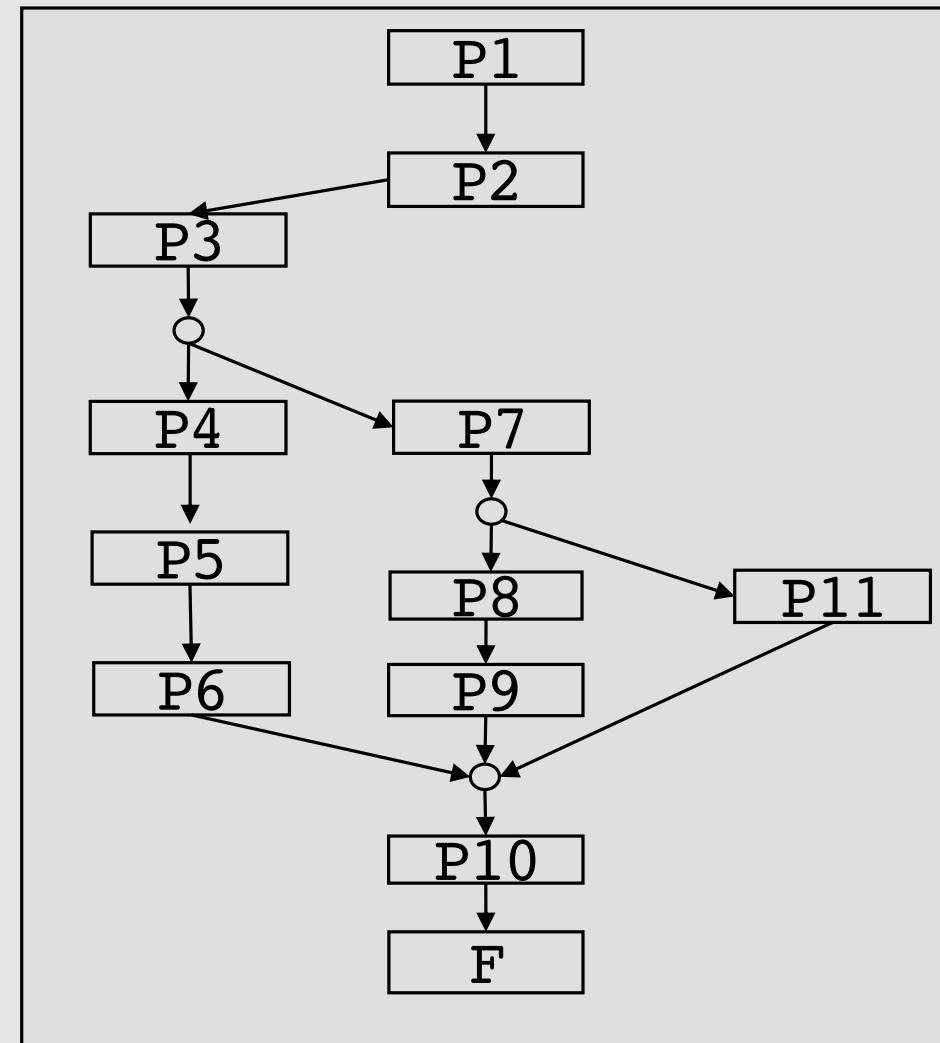


AID

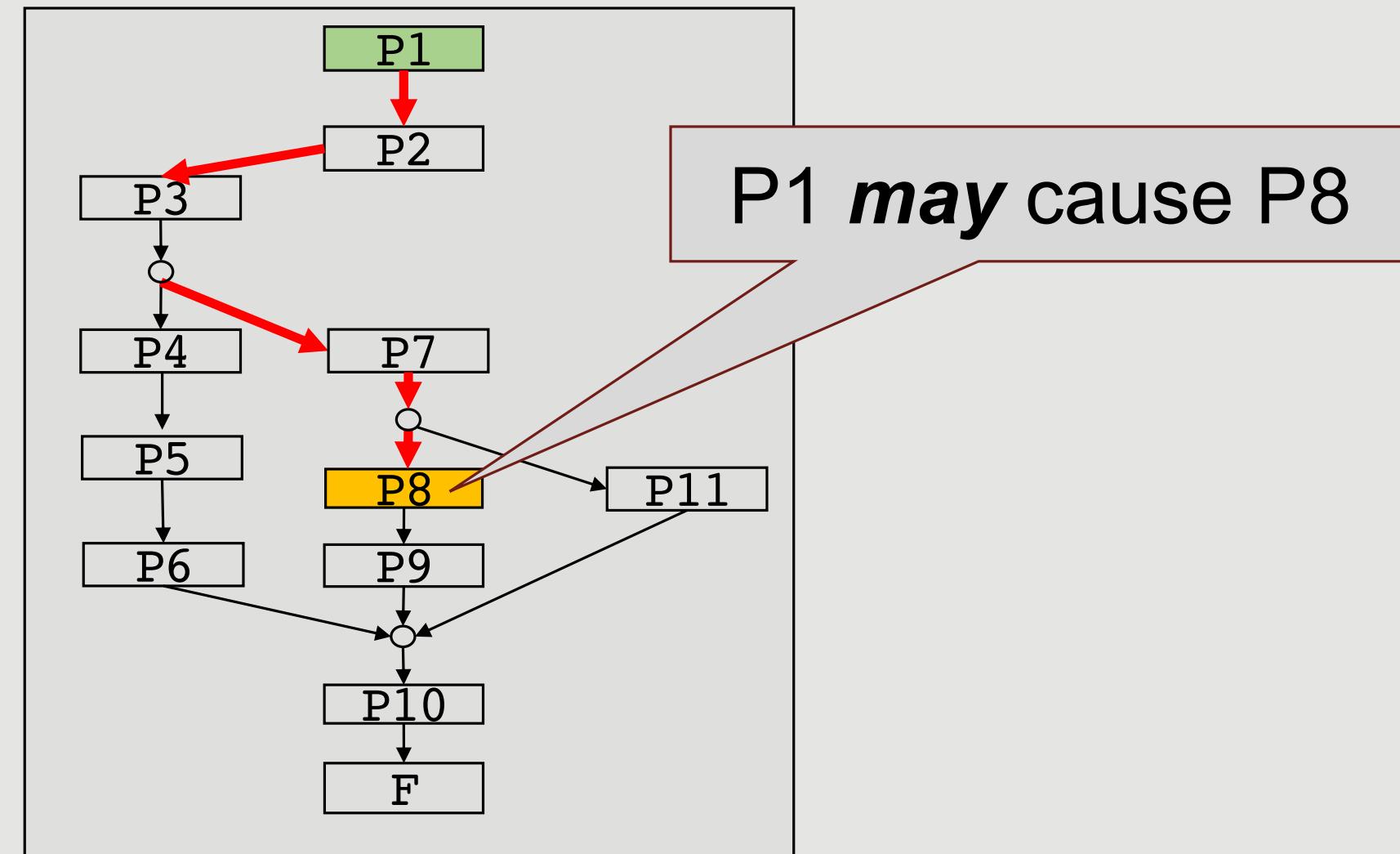


Cause must temporally precede effect

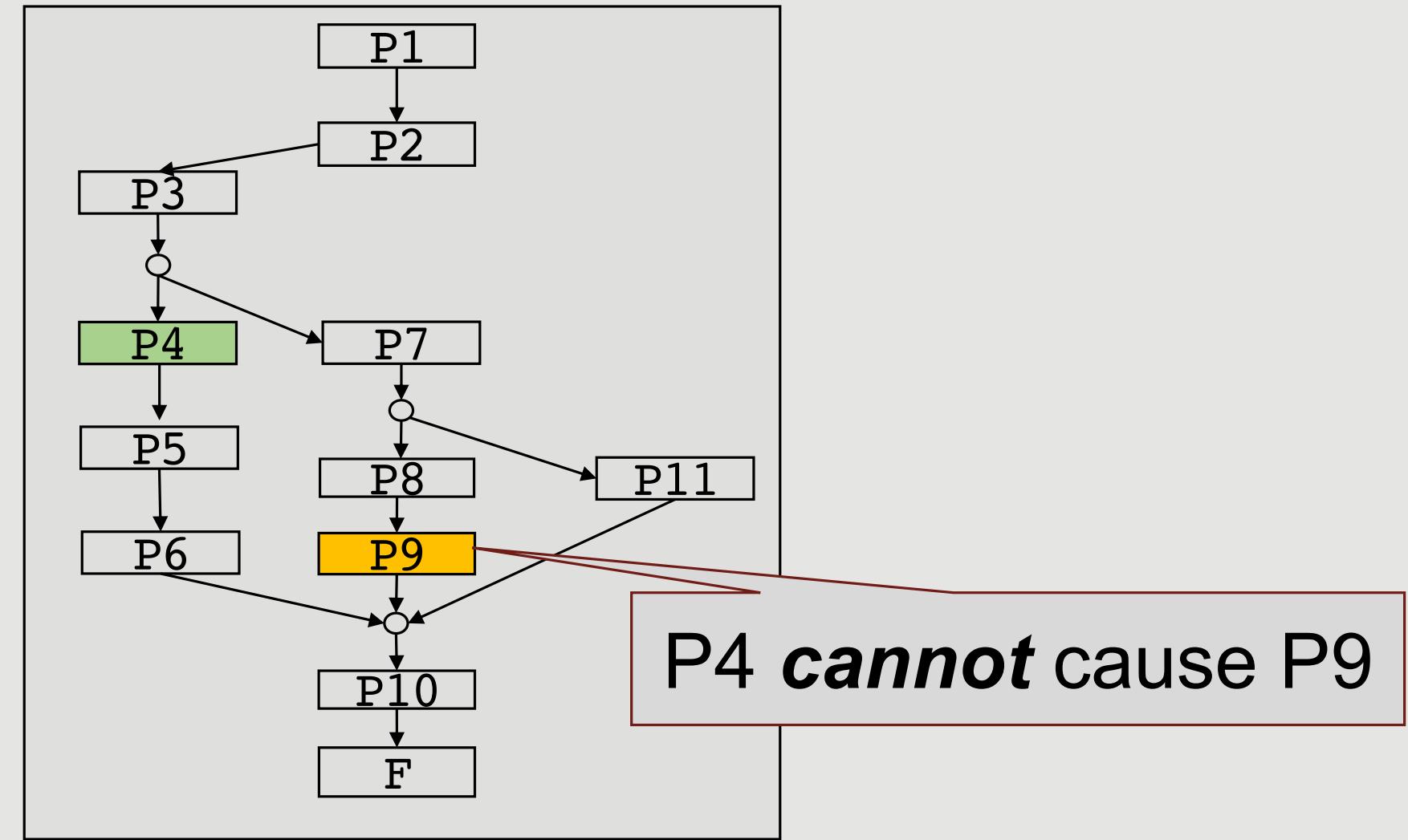
Temporal precedence graph



Approximating causality



Approximating causality



Counterfactual causality

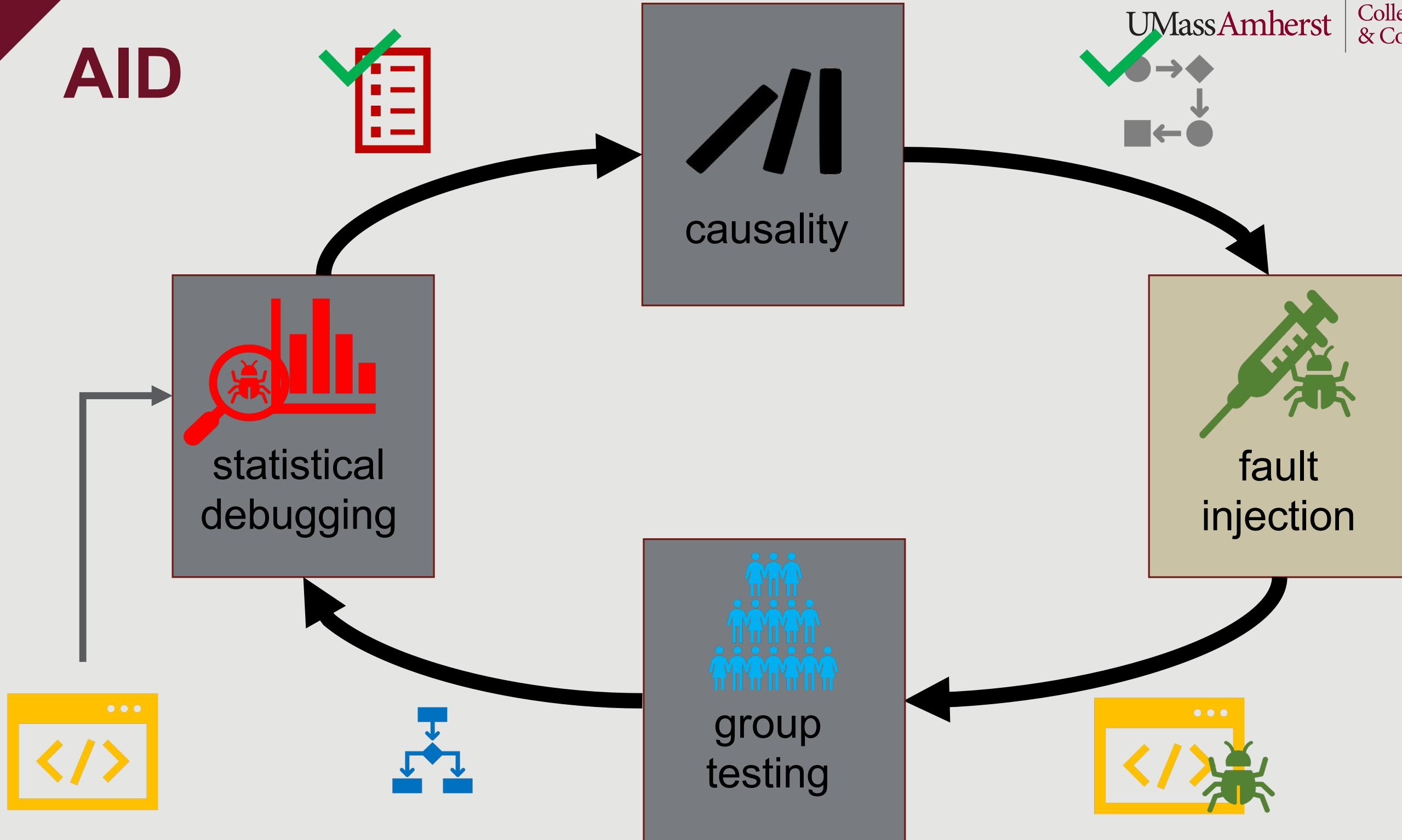
C is a *counterfactual cause* of E
If C had not occurred
E would not have occurred



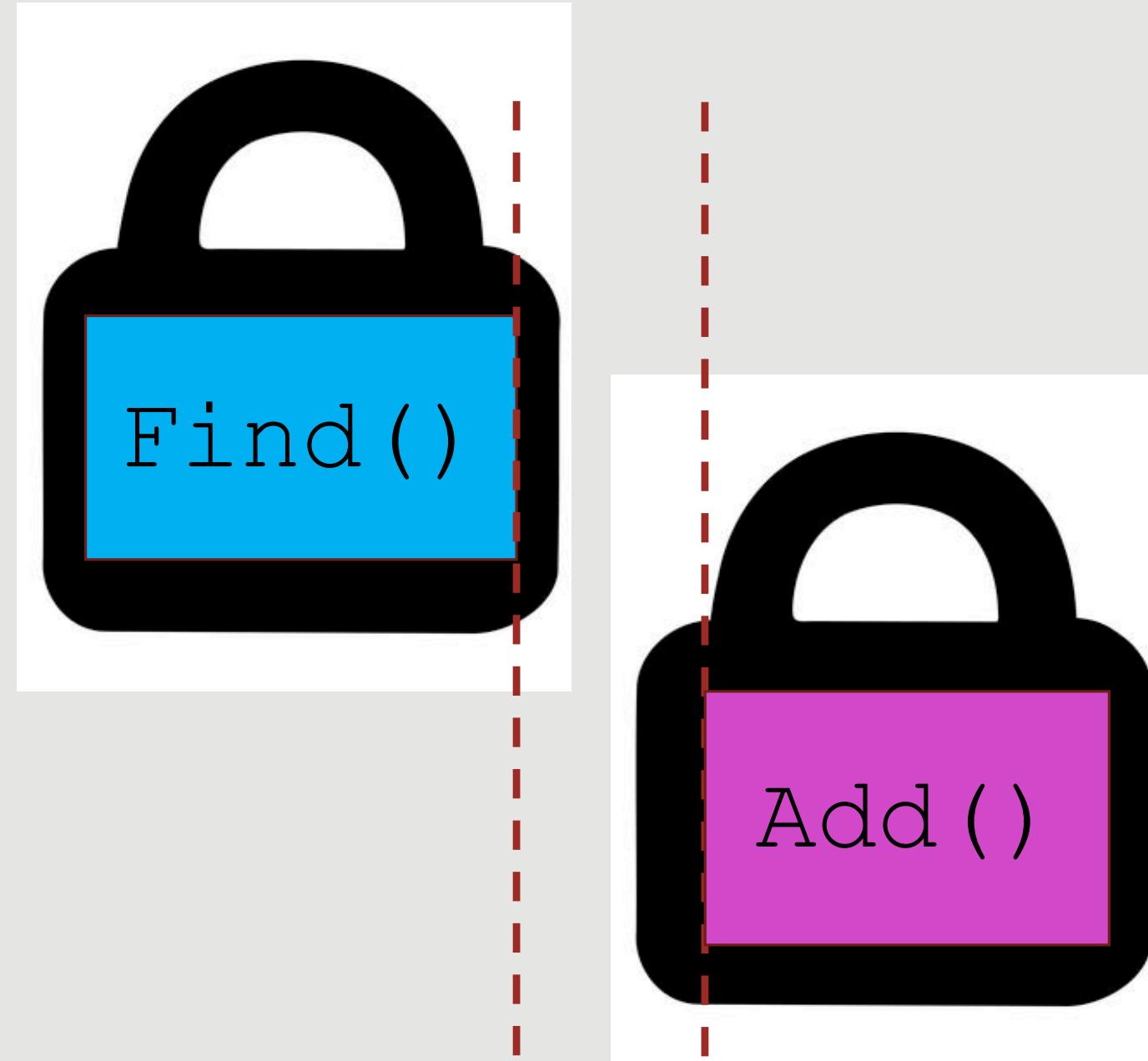
Intervention



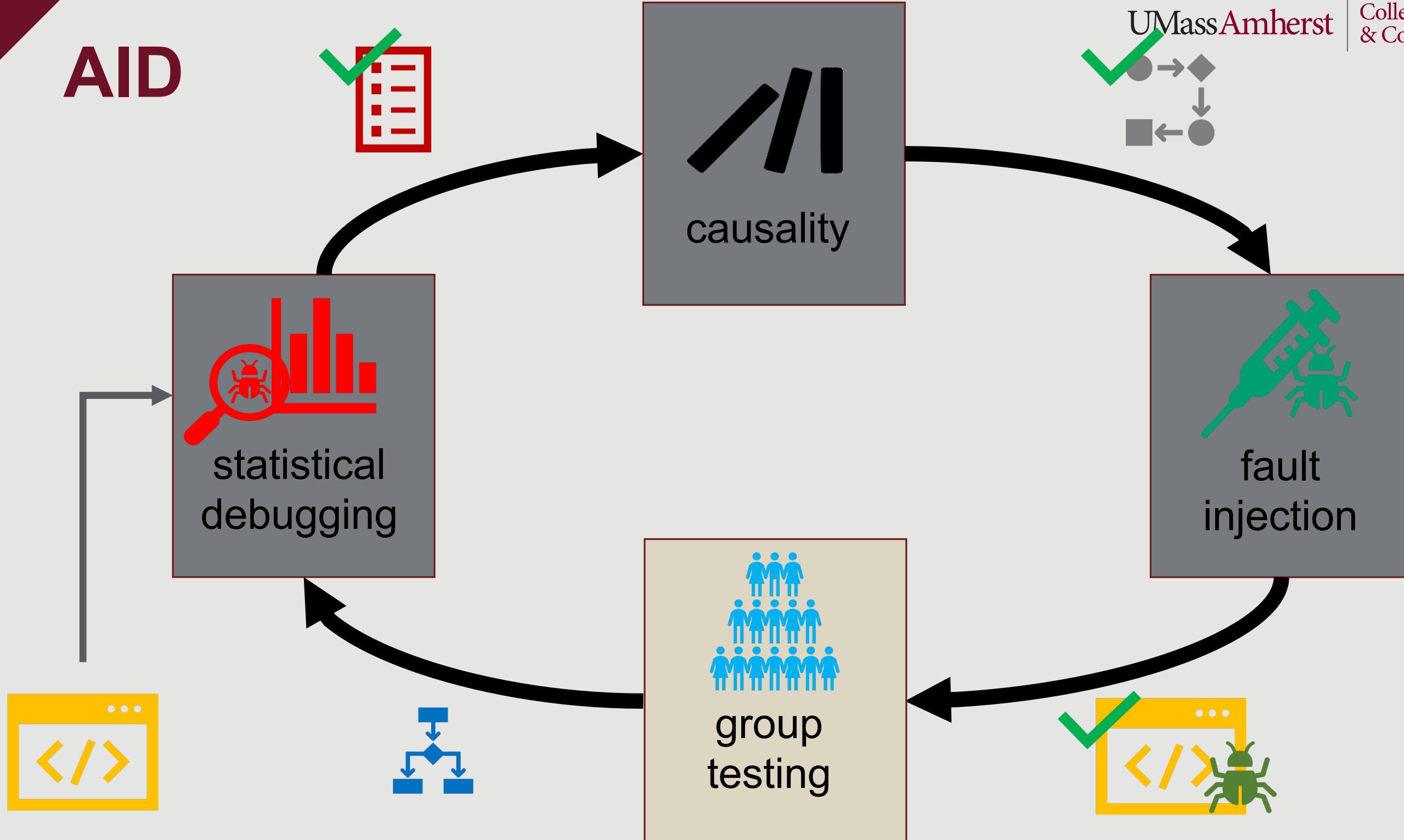
AID



Fault injection



AID



Group testing

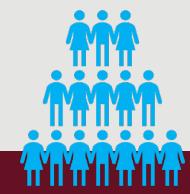
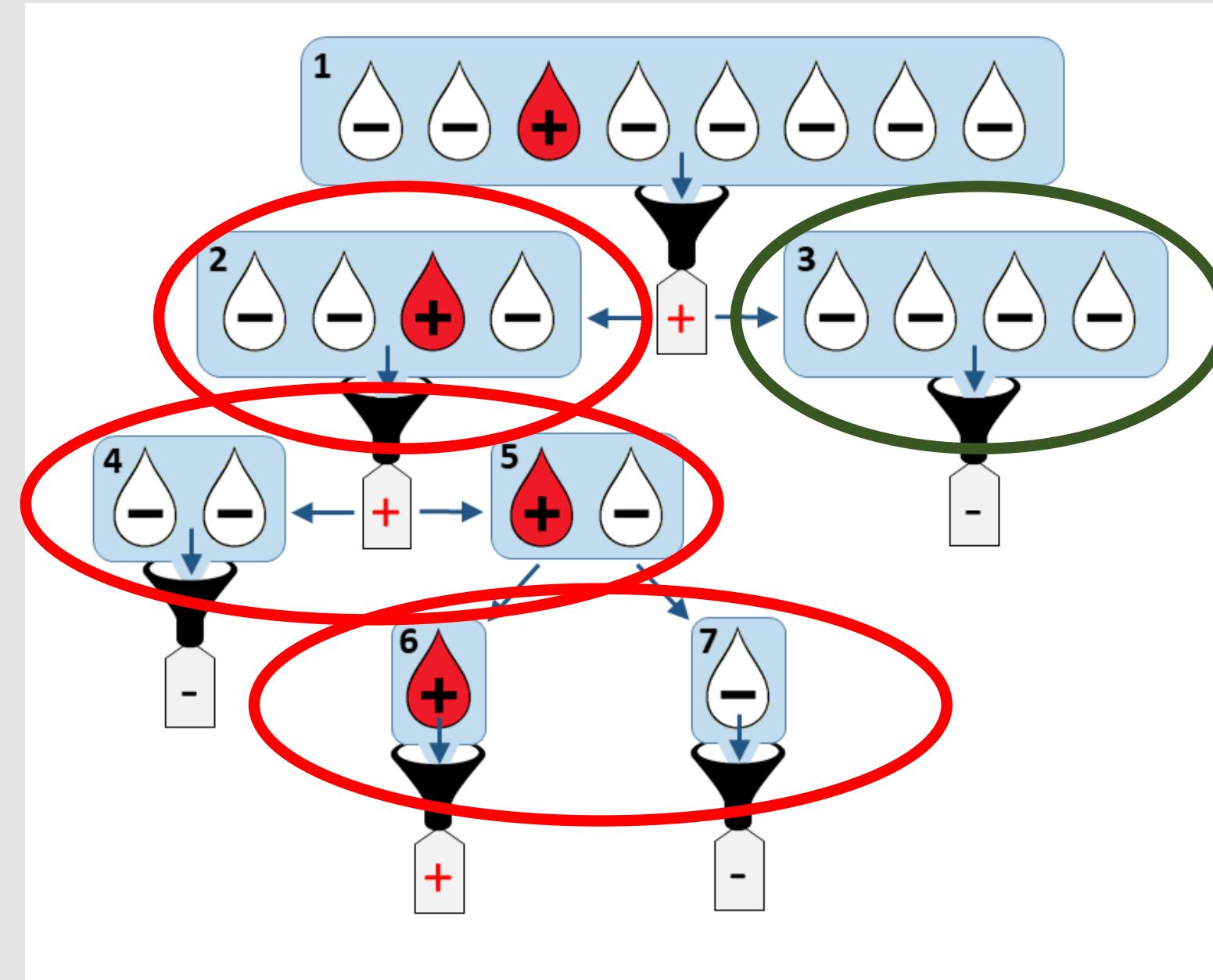
PUBLIC HEALTH

Coronavirus Test Shortages Trigger a New Strategy: Group Screening

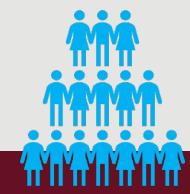
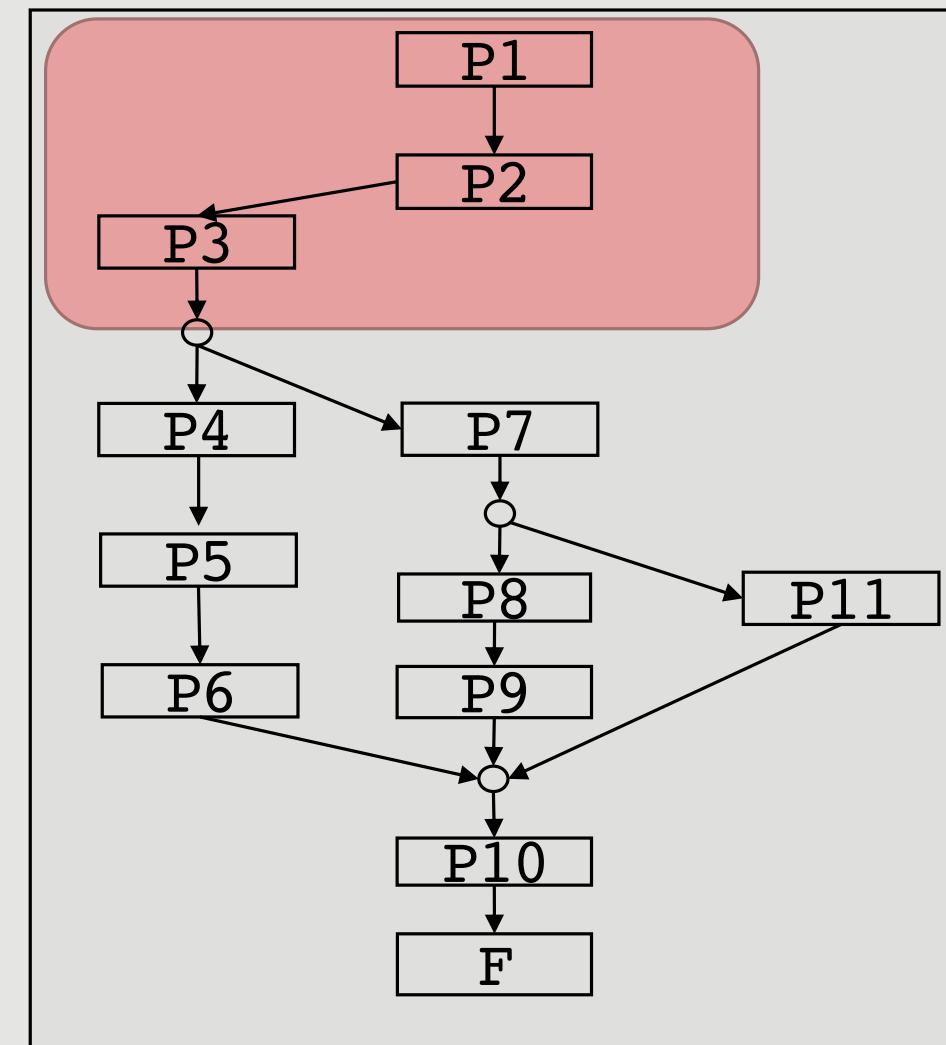
Pooling diagnostic samples, and using a little math, lets more people get tested with fewer assays



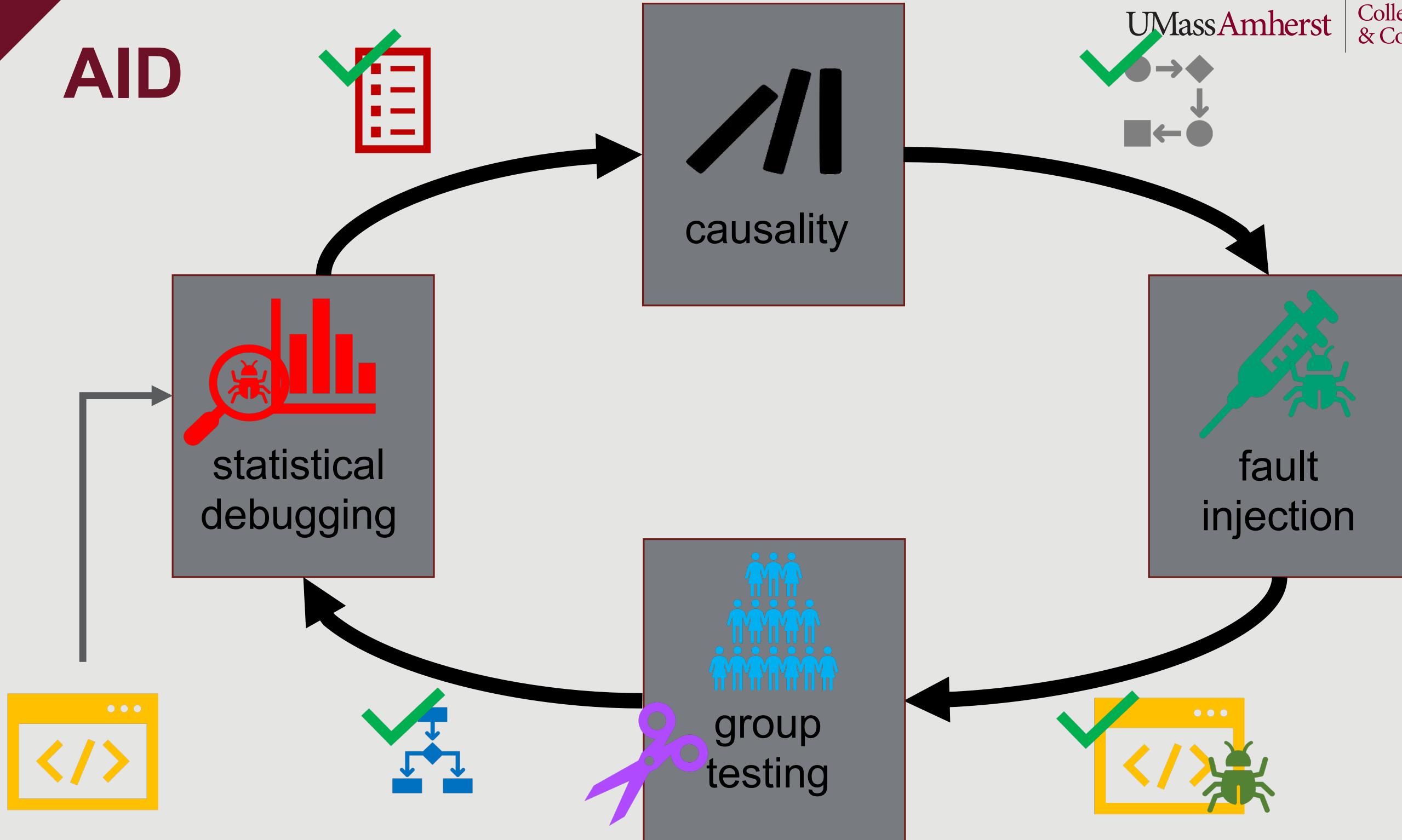
Adaptive group testing



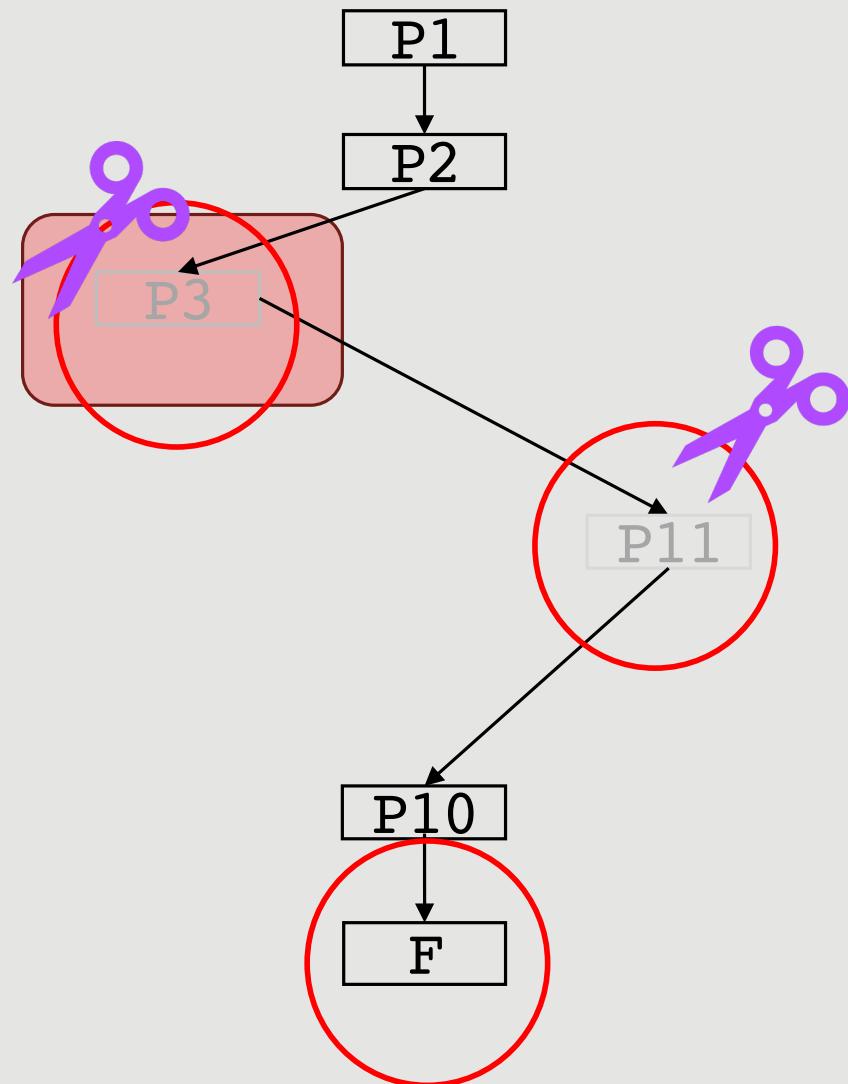
AID applies group intervention



AID



AID pruning





EVALUATION



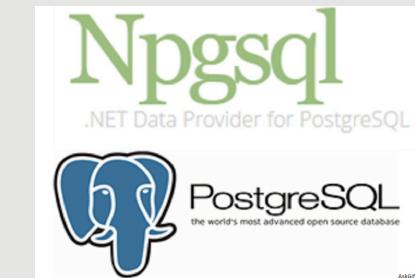
Six real-world bugs



Data race ✓



Use-after-free ✓



Timing-bug ✓

Network



Microsoft

Random number
collision ✓

BuildAndTest



Microsoft

Order violation ✓

HealthTelemetry

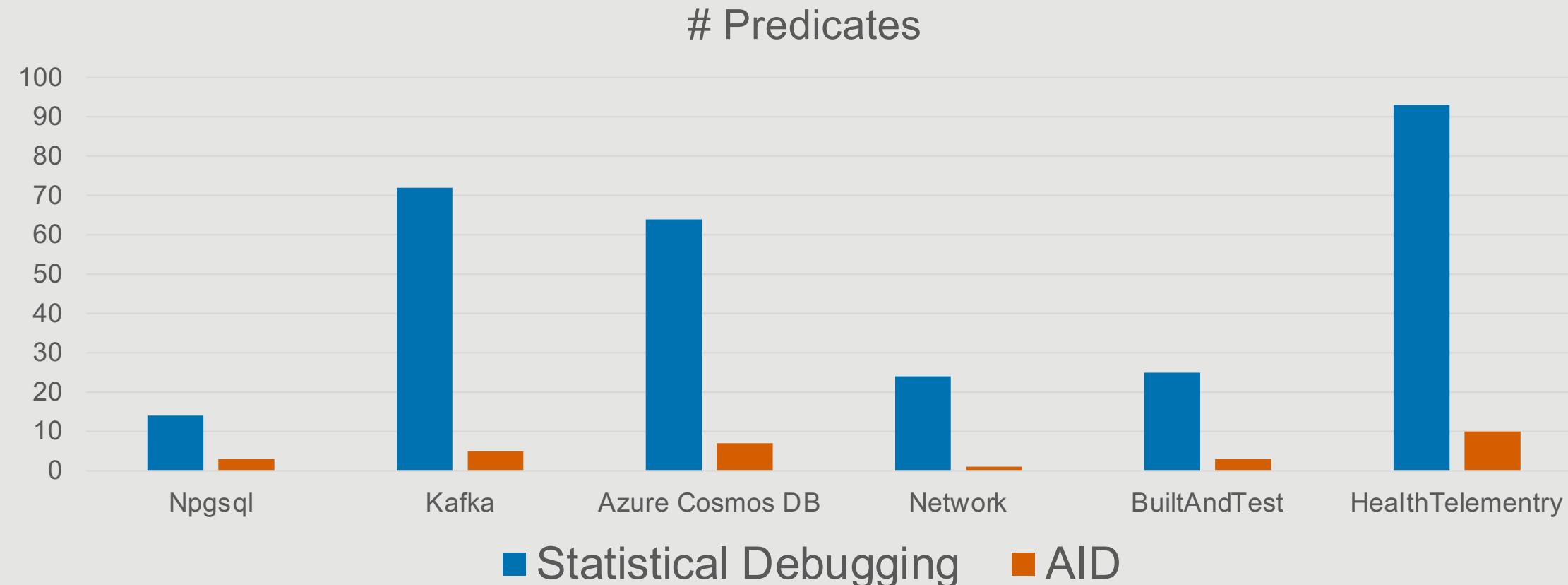


Microsoft

Race condition ✓

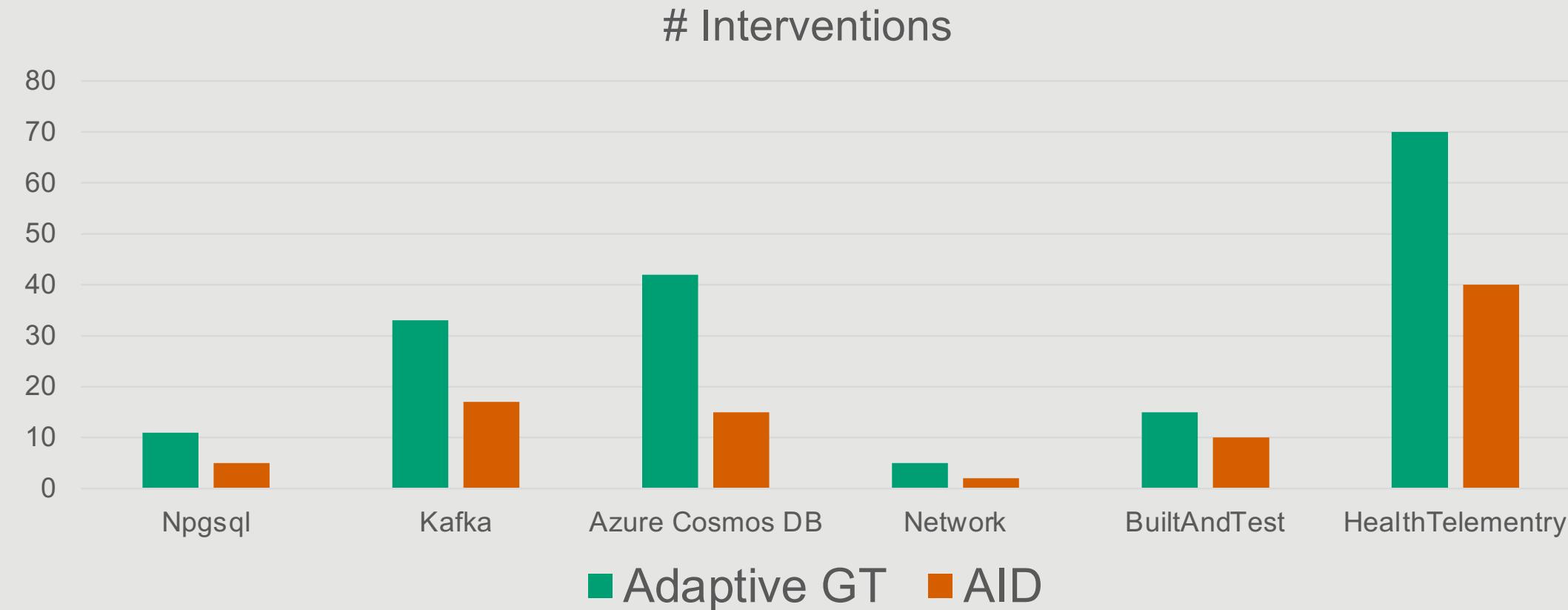
Statistical debugging vs AID

AID produces no false positives



Adaptive group testing vs AID

AID's pruning reduces #Interventions



Theoretical analyses

CPD: Causal Path Discovery

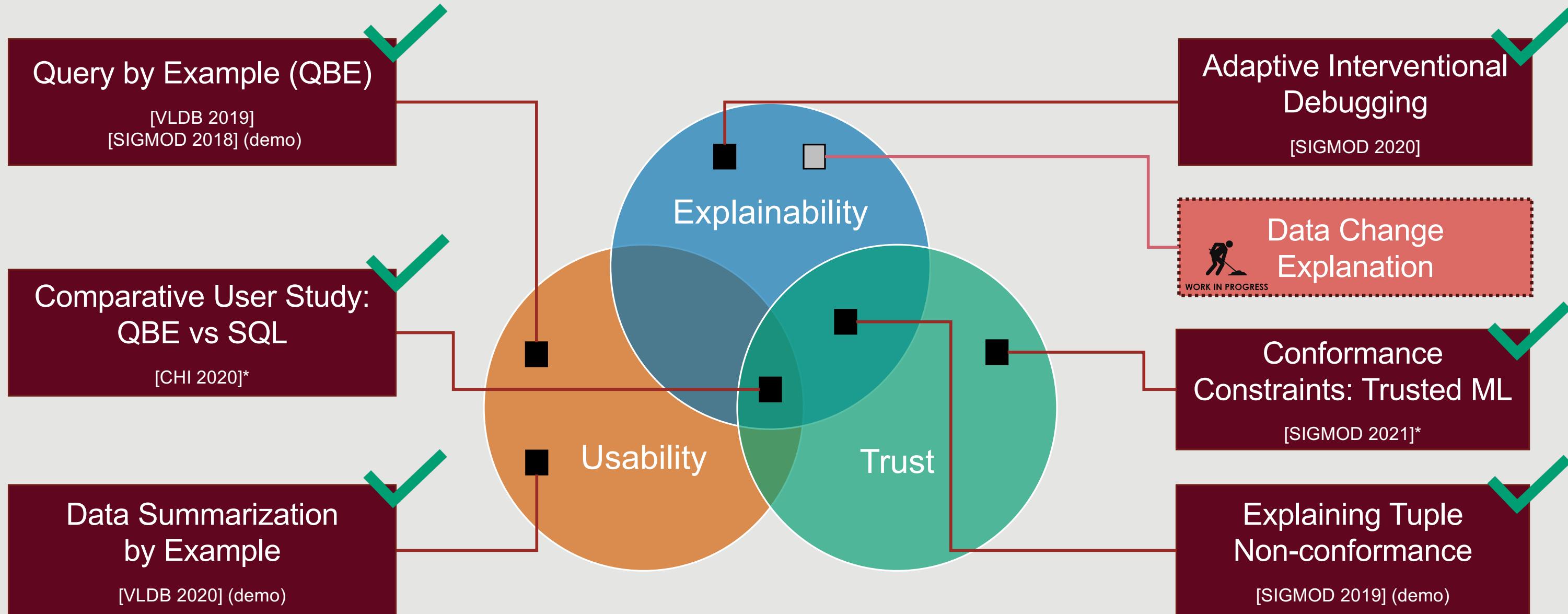
GT: Group Testing

AID: Adaptive Interventional Debugging

TAGT: Traditional Adaptive Group Testing

	Search space	#Interventions	
		Lower bound	Upper bound (AID/TAGT)
CPD	$(B(2^n - 1) + 1)^J$	$\frac{JBn}{JBn + DS_1} \log \binom{JBn}{D}$	$J \log B + D \log(Jn) - \frac{D(D-1)S_2}{2Jn}$
GT	2^{JBn}	$\log \binom{JBn}{D}$	$D \log B + D \log(Jn) - \frac{D(D-1)}{2JBn}$

Dissertation outline



Part 4: Proposed Contributions & Tentative Timeline

Data Change
Explanation

How did my data change over last couple years?

15:50:23,,0.5,69,,11425,,, "271504218477G",32,,,...,11,,,...
11 - 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,41,69,32,2018-10-10
15:50:23,,0.5,69,,11428,,, "271507491568G",32,,,...,11,,,...
12 - 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,41,69,32,2018-10-10
15:50:23,,0.5,69,,11484,,, "271508481857G",32,,,...,11,,,...
13 - 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,77,69,53,2018-10-11
11:35:05,,0.05,134,,11447,,, "874231098887G",53,,,...,11,,,...
14 - 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,77,69,53,2018-10-11
11:35:05,,0.05,134,,11448,,, "874231135374G",53,,,...,11,,,...
15 - 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,77,69,53,2018-10-11
11:35:05,,0.05,134,,11479,,, "874231461234G",53,,,...,11,,,...
16 - 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,87,69,59,2018-10-11
13:43:24,,0.05,34,,11487,,, "874231676529G",59,,,...,11,,,...
17 - 73,"FIVE, PERSON",,FIVE,PERSON,,,...,Y,G,23,73,19,2018-10-08
22:25:59,,0.75,73,,14,,, "271508486757G",19,,,...,11,,,...
18 - 73,"FIVE, PERSON",,FIVE,PERSON,,,...,Y,G,23,73,19,2018-10-08
22:25:59,,0.75,73,,11512,,, "874231926046G",19,,,...,11,,,...

15:50:23,,0.5,69,,11425,,, "271504218477",32,,,...,11,,,...
11 + 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,41,69,32,2018-10-10
15:50:23,,0.5,69,,11428,,, "271507491568",32,,,...,11,,,...
12 + 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,41,69,32,2018-10-10
15:50:23,,0.5,69,,11484,,, "271508481857",32,,,...,11,,,...
13 + 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,77,69,53,2018-10-11
11:35:05,,0.05,134,,11447,,, "874231098887",53,,,...,11,,,...
14 + 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,77,69,53,2018-10-11
11:35:05,,0.05,134,,11448,,, "874231135374",53,,,...,11,,,...
15 + 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,77,69,53,2018-10-11
11:35:05,,0.05,134,,11479,,, "874231461234",53,,,...,11,,,...
16 + 69,"FOUR, PERSON",,FOUR,PERSON,,,...,Y,G,87,69,59,2018-10-11
13:43:24,,0.05,34,,11487,,, "874231676529",59,,,...,11,,,...
17 + 73,"FIVE, PERSON",,FIVE,PERSON,,,...,Y,G,23,73,19,2018-10-08
22:25:59,,0.75,73,,14,,, "271508486757",19,,,...,11,,,...
18 + 73,"FIVE, PERSON",,FIVE,PERSON,,,...,Y,G,23,73,19,2018-10-08
22:25:59,,0.75,73,,11512,,, "874231926046",19,,,...,11,,,...

Prior work

- Existing approaches mostly focus on **syntactic** changes.
- Fail to provide **consumable summary** of changes.

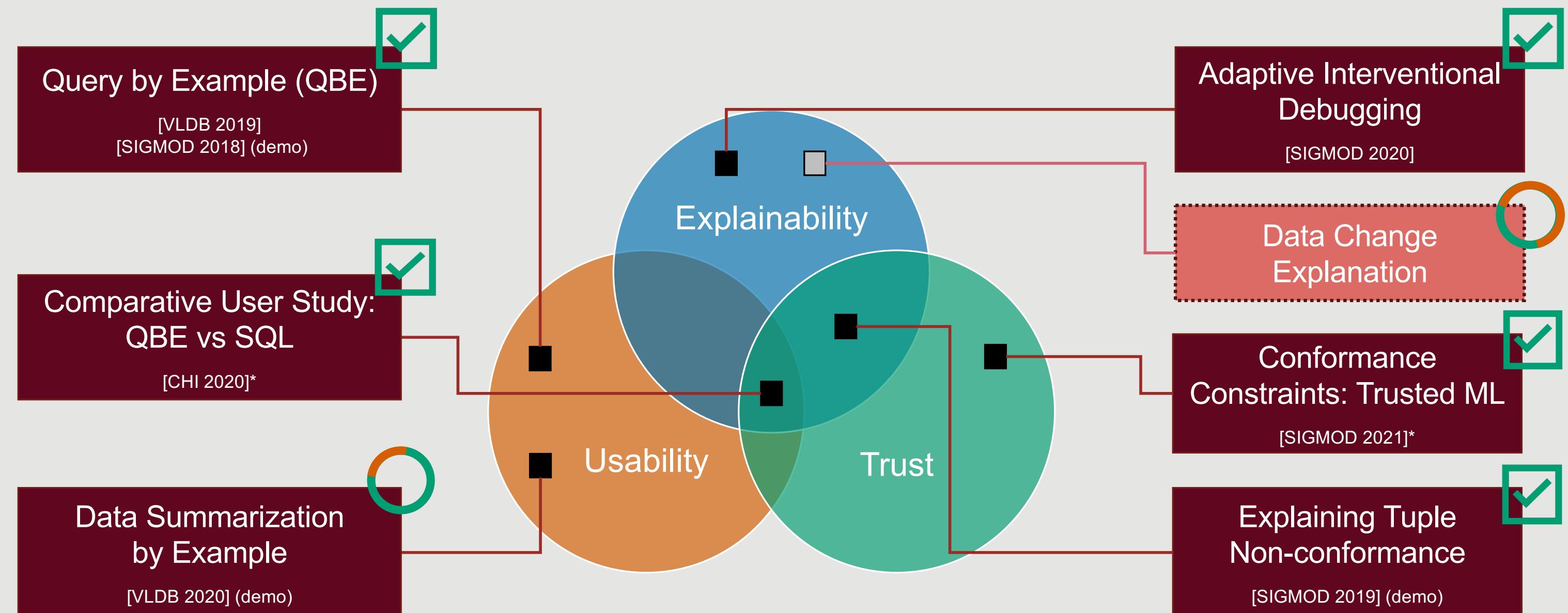
Our goal

- Provide a consumable summary of **semantic changes** that **explains** how two databases differ.
- Explains database **evolution**.
- Reveals **patterns** in data change.

Evaluating SuDocu

- Data collection
- Tuning SuDocu's learning algorithm
- Evaluation
 - Against ground-truth summaries
 - Comparison with other baselines
 - User study

Current status



Tentative timeline

- October 2020: proposal defense
- November – December 2020: evaluating SuDocu
- January 2020: submit to VLDB 2021
- January – June 2021: work on Data Change Explanation Framework
- July 2021: submit to SIGMOD 2022
- June – August 2021: work on dissertation
- August 2021: final defense

Other project affiliations

- Fair classifiers: experiment and evaluation



- Data profile debugger



- Data sampling by example



Acknowledgements



Acknowledgements

- Committee



- Mentors and collaborators



Acknowledgements

- Armand Asnani, UMass
- Lucy Cousins, UMass
- Nischal Dave, UMass
- Larkin Flodin, UMass
- Juliana Freire, NYU
- Sainyam Galhotra, UMass
- Maliha Tashfia Islam, UMass
- Eunice Jun, UW
- Beryl Larson, Wellesley College
- Genglin Liu, UMass
- Raoni Lourenço, NYU
- Raj Kumar Maity, UMass
- Kancha Masalia, UMass
- Sheshera Mysore, UMass
- Tony Ohmann, UMass
- Vincent Pun, UMass
- Sheikh Muhammad Sarwar, UMass
- Michael Satanovsky, Hopkins School
- Divesh Srivastava, AT&T
- Zoey Sun, Smith College
- Nishant Yadav, UMass



COMPUTING FOR THE COMMON GOOD

Anna Fariha

afariha@cs.umass.edu

people.cs.umass.edu/afariha