

Don't Be a Tattle-Tale: Preventing Leakages through Data Dependencies on Access Control Protected Data

Primal Pappachan^{†‡}, **Shufan Zhang**^{*}, Xi He^{*}, Sharad Mehrotra[†]

[†] University of California Irvine [‡] Pennsylvania State University ^{*} University of Waterloo

1. Inference Problem

Eid	EName	Zip	State	Role	WorkHrs	SalPerHr	Wid	Eid	DeptName	Salary	
t ₁	34	Tina	45678	WA	Student	20	40	1	34	CS	800
t ₂	56	Bobby	54321	CA	Faculty	40	200	2	56	EE	8000
t ₃	78	Dale	53567	CA	Faculty	40	180	3	78	CS	7200
t ₄	12	Khan	54321	CA	Staff	30	70	4	12	BIO	2100

Zip → State
Role → SalPerHr

RQ1. How to design an appropriate security notion to capture such leakage model?

RQ2. How to identify "minimal" number of cells to hide to prevent such leakages?

RQ3. How to design an algorithm or system to "efficiently" prevent such leakages?

2. Background

Tuple Representation

$$\forall A_i \text{ Dom}(A_i) = \{1, 2, 3\}$$

A ₁	A ₂	A ₃
c ₁ 1	c ₂ 2	c ₃ 3
c ₄ 1	c ₅ 2	c ₆ 3

Database as a collection of cells

Cell Representation

c ₁	c ₂	c ₃	c ₄	c ₅	c ₆
1	2	2	1	2	3
A ₁	A ₂	A ₃	A ₁	A ₂	A ₃

Access Control Policies:

Querier Views:

V(C)

c ₁	c ₂	c ₃	c ₄	c ₅	c ₆
1	*	2	1	2	*
A ₁	A ₂	A ₃	A ₁	A ₂	A ₃

V₀(C)

c ₁	c ₂	c ₃	c ₄	c ₅	c ₆
*	*	*	*	*	*
A ₁	A ₂	A ₃	A ₁	A ₂	A ₃

- Policy maps a Cell c_i to either sensitive or non-sensitive.
- When a cell is sensitive, its value is replaced with * (null)
- Base view where all cells are nulls

4. Tattle-Tale Condition

$TTC(\delta, V, c_i) =$ True, when all the other predicates (except $Pred(c^*)$) evaluate as True
False, otherwise

Full deniability is achieved for a sensitive cell in a view if for all relevant dependency instantiations, Tattle-Tale Condition evaluates False.

Preventing Leakages:

Hide one other cell in δ to hide the truth value of another predicate

c ₁	c ₂	c ₃	c ₄	c ₅	c ₆	c ₁₂
*	2	2	1	2	*

$$\delta_1: \neg((c_1 = c_4) \wedge (c_2 = c_5) \wedge (c_3 < c_6))$$

$$\neg((c_1 = 1) \wedge (c_2 = 2) \wedge (c_3 < *))$$

Cueset for c₁: {c₂, c₅, c₃, c₆} Unknown True Unknown

2. Related Work

- Design time prevention → Different inference channels, poor data availability [TKDE'96, CSF'98]
- Query time prevention → Weak security (fully reconstruct the sensitive cell) [TKDE'00]
- Perfect secrecy or randomized algorithm → Poor data utility [SIGMOD'04, ICDE'20]

3. Security Model

Inference channel:

- Denial Constraints (DCs)
- Function-based Constraints

$$\delta_1: \neg(Pred_1 \wedge Pred_2 \wedge \dots \wedge Pred_n)$$

$$Pred_i = (c_i \text{ op constant}) \text{ or } (c_i \text{ op } c_j)$$

where op: {=, ≠, ≤, ≥, >, <}

DC instantiation:

c ₁	c ₂	c ₃	c ₄	c ₅	c ₆
1	2	2	1	2	3
A ₁	A ₂	A ₃	A ₁	A ₂	A ₃

$$A_1 \rightarrow A_3:$$

$$\delta_1: \neg((c_1 = c_4) \wedge (c_3 \neq c_6))$$

Inference function:

$$I(c_i | V, \delta_1)$$

V(C)

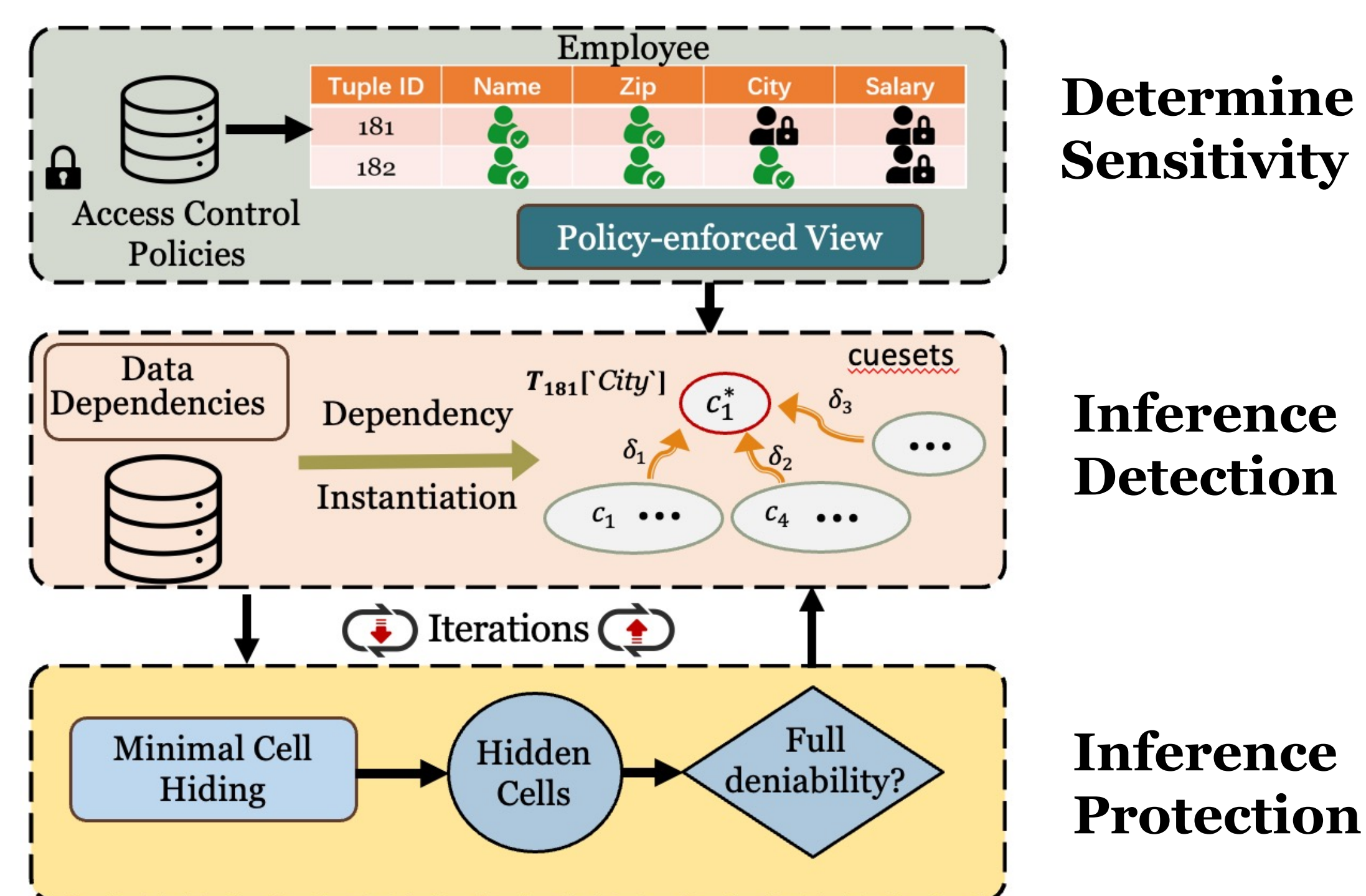
c ₁	c ₂	c ₃	c ₄	c ₅	c ₆
1	*	2	1	2	*
A ₁	A ₂	A ₃	A ₁	A ₂	A ₃

Full deniability:

$$\text{For } S_\Delta = \{\delta_i\}, \forall c_i \in C^S, I(c_i | V, S_\Delta) = I(c_i | V_0, S_\Delta)$$

Adversary can learn nothing about the sensitive cells beyond what is given in the base view.

5. System Overview



Inference Detection:

- Done naively could result in quadratic blow up
- Optimization: Instantiate only the set of cells the lead to $Pred(c^*)$ evaluating as False

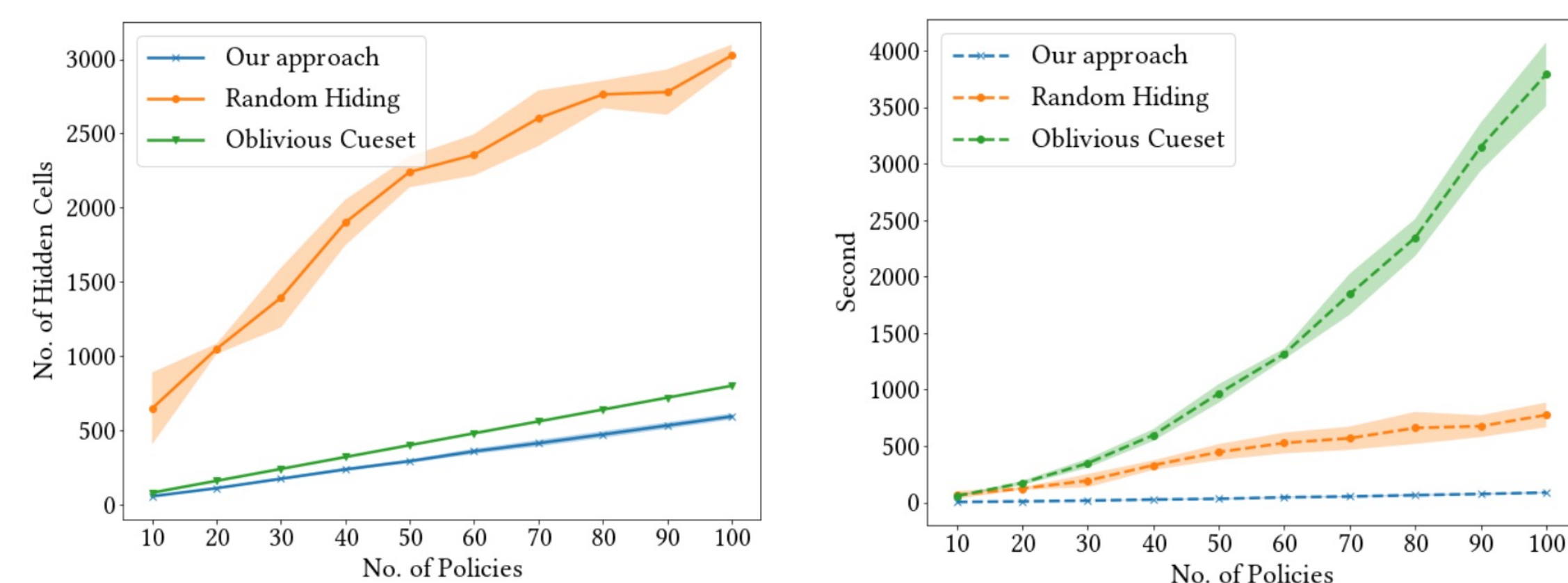
Inference Protection:

- Perform Minimum Subset Cover (MSC) on the cuesets to find a minimum cover of the subsets
- Repeat until all cuesets corresponding to sensitive cells and hiddenCells are protected

6. Evaluation

Dataset: Tax [ICDE'07]

- 10k tuples, 14 attributes, with 11 dependencies

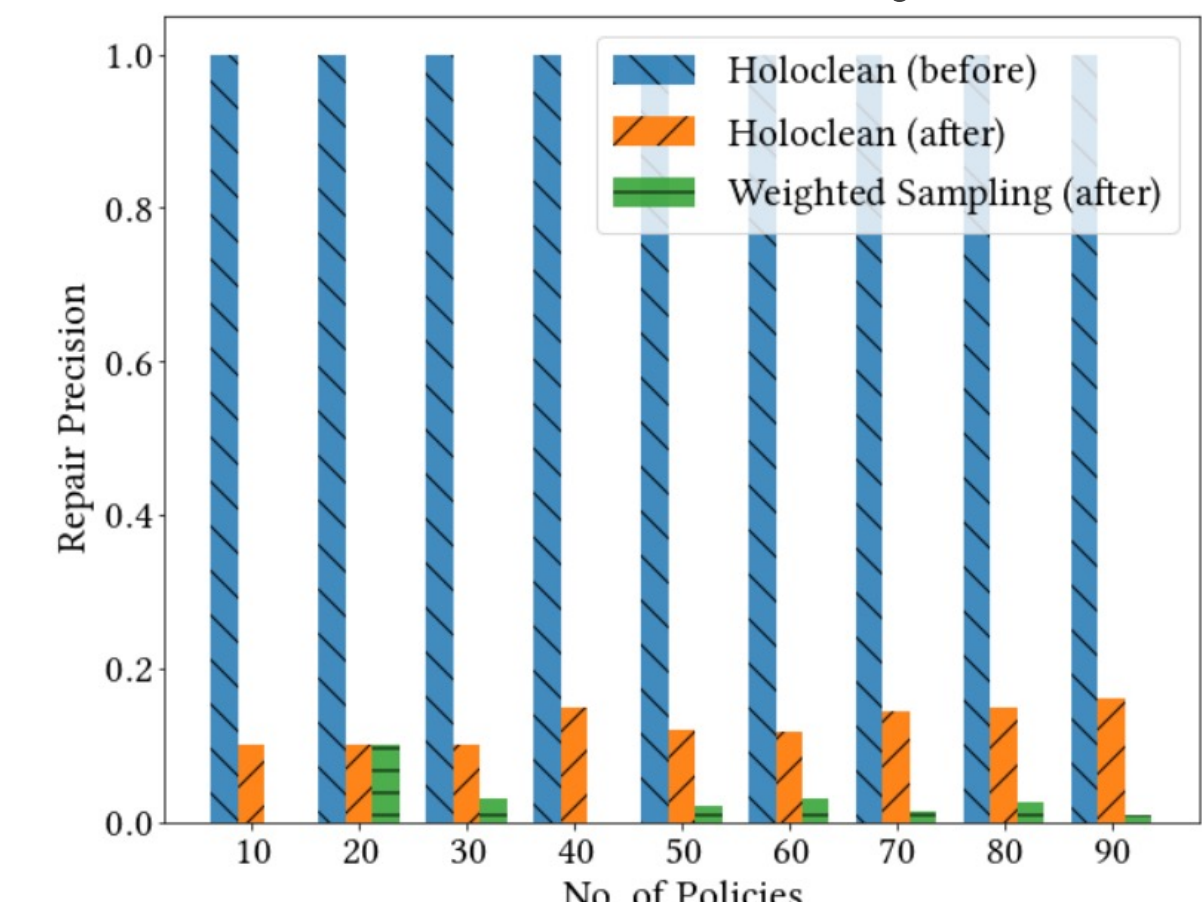


Data Utility:

- With increasing number of sensitive cells, the number of hidden cells increases linearly
- Our approach always hides less cells than the baselines

Performance:

- The overhead of our approach is small, compared to the baselines
- When increasing the number of sensitive cells, the performance overhead scales linearly for our approach



Holoclean [VLDB'17] as an Adversary:

- Holoclean can reconstruct 100% protected cells before applying our approach
- After applying our approach, it only recovers 10%-15% cells, with marginal improvement over random guess

7. Takeaways

Leakage attack based on two types of data dependencies

- Denial Constraints
- Function-based Constraints

Strong security model

- Tattle-Tale Condition
- Full Deniability
- Relaxing assumptions in the model

End-to-end system

- Utility, efficiency, scalability and convergence
- Optimizations to improve performance



Check out the project repo!



UNIVERSITY OF WATERLOO



UCI



PennState