Privacy-Preserving Joins in Untrusted Environments

Kajetan Maliszewski PhD Defense 4. April 2025



Three states of data matter



Three states of customer's sensitive data in public cloud



Trusted Execution Environments (TEEs)



Advantages over other PETs:

Performance

Enhanced functionality

Ease of deployment

Adoption of TEEs has been a bumpy road



Disadvantages over plain CPU:

Challenging hardware architecture

Performance bottlenecks

Security considerations

Joins in TEEs can significantly underperform



Thesis Goal

Analyze and improve the privacy-performance trade-off of TEEs in batch and stream join workloads

Contributions

Problem







Efficient Query Processing in TEEs



Honest-but-curious

adversary

Privacy-Preserving Stream Processing



Solution



Agenda

- 1. Motivation
- 2. Understanding equi-joins in TEEs
- 3. Relational Join for TEEs
- 4. Privacy-preserving stream joins
- 5. Conclusions









We benchmarked four families of join algorithms

		Join Algorithm	
Family	hash-based	CHT	Concise Hash Table
		PHT	Parallel Hash Table
	sort-merge	PSM	Parallel Sort-Merge
		MWAY	Multi-Way Sort-Merge
	radix-based	RHT	Radix Hash Table
		RHO	Radix Hash Optimized
		RSM	Radix Sort-Merge
	nested-based	INL	Index Nested Loop



TeeBench is a fair referee for enclave benchmarks



Plug & Play experience

We used TeeBench in all subsequent works



How to do multi-threading?





Findings Summary

We call for TEE-native solutions



We provide 7 lessons learned

None of the joins is a good choice

TEE-native approaches offer promising results

Agenda

- 1. Motivation
- 2. Understanding equi-joins in TEEs
- 3. Relational Join for TEEs
- 4. Privacy-preserving stream joins
- 5. Conclusions



Desiderata for TEE-native processing





Low memory consumption

Wait-free algorithms



Cracking-Like Philosophy



- Partition the data
- Perform sequential scans
- Consume little memory
- TBD: Barrier-free execution



b1x

....

num

Two new primitives perform the Cracking-Like Join



Cracking-Like Join algorithm

- 1: **procedure** CRKJOIN(Relation R, Relation S, int bits):
- 2: $rootR \leftarrow INITCRACKINGTREE(R)$
- 3: $rootS \leftarrow INITCRACKINGTREE(S)$
- 4: **for each** $p \in [0..(2^{bits} 1)]$ **do**
- 5: CRACKSTAGEANDBUILD(*rootR*, *p*)
- 6: CRACKSTAGEANDPROBE(rootS, p)



Partitioned relation forms independent chunks





CrkJoin outperforms the baselines in a multi-tenant scenario





CrkJoin scales in multi-core architectures



Findings summary



Efficient query processing in TEEs is possible

Our algorithm achieves up to 1000x higher performance

CrkJoin scales to multiple cores

Agenda

- 1. Motivation
- 2. Understanding equi-joins in TEEs
- 3. Relational Join for TEEs
- 4. Privacy-preserving stream joins
- 5. Conclusions



























We protect data-dependent information (i.e., graph componentes)

We allow for data-independent leakage

Stream constraints







L0: deterministic confidentiality

L1: randomized confidentiality



Informal definition

Leak current join graph and volume patterns, while hiding the graph history. Similar to randomized encryption.



L0: deterministic confidentiality

L1: randomized confidentiality

L2: oblivious lookups



Informal definition

For every tuple leak only the volume pattern.



L0: deterministic confidentiality

L2: oblivious lookups

L3: oblivious batching



Informal definition

For every batch of tuples leak the intermediate join result cardinality.



L0: deterministic confidentiality

L1: randomized confidentiality

L2: oblivious lookups

L3: oblivious batching

L4: full obliviousness



Informal definition

Leak no information about the data or access patterns.



We propose two families of stream joins secure against leakage functions





How to avoid oblivious sorting?





3

R



How to obliviously append to a sorted collection?



OblivAppend OSort micro-batch Append dummies

3 OMerge 4 Trim dummies



- 1. window is sorted
- 2. batch is unsorted
- 3. n >> m
- 4. OSort is O(nlog²n)

If n >> m: **O(nlogn)**



OblivAppend as a FK-join building block



OblivAppend

OSort micro-batch
Append dummies
OMerge
Trim dummies

FK join 5 Emit join matches

Challenges: window correctness, result uniqueness, worst-case padding

FK restrictions allows us to split the join:

 $R\Join S=((R\cup \mathcal{W}_R)\bowtie S)\cup (R\bowtie \mathcal{W}_S)$





Challenges: window correctness, result uniqueness, worst-case padding

К

S

FK restrictions allows us to split the join:

 $R \bowtie S = (\underbrace{(R \cup \mathcal{W}_R)}_{1} \bowtie S) \cup (R \bowtie \mathcal{W}_S)$





Challenges: window correctness, result uniqueness, worst-case padding

FK restrictions allows us to split the join:

 $R \bowtie S = ((R \cup \mathcal{W}_R) \bowtie S) \cup (R \bowtie \mathcal{W}_S)$



Challenges: window correctness, result uniqueness, worst-case padding

FK restrictions allows us to split the join: $W_R U$ $R\Join S=((R\cup \mathcal{W}_R)\Join S)\cup (R\bowtie \mathcal{W}_S)$ R 3 К S $|W_{S}|$ Ws



Challenges: window correctness, result uniqueness, worst-case padding

К

S

FK restrictions allows us to split the join:

 $R\Join S=((R\cup \mathcal{W}_R)\bowtie S)\cup (R\bowtie \mathcal{W}_S)$

4 Housekeeping





Challenges: window correctness, result uniqueness, worst-case padding

К

S

FK restrictions allows us to split the join:

 $R\Join S=((R\cup \mathcal{W}_R)\bowtie S)\cup (R\bowtie \mathcal{W}_S)$

Security: all operations are oblivious **Performance**: O(nlogn)



 $W_R U$

R



Privacy cost in oblivious stream joins





Privacy without obliviousness has negligible overhead

Privacy cost in oblivious stream joins





Privacy without obliviousness has negligible overhead

OCA joins perform well, while maintaining low latency

Privacy cost in oblivious stream joins





Privacy without obliviousness has negligible overhead

OCA joins perform well, while maintaining low latency

Oblivious index-based solutions significantly underperform

Contributions summary



Oblivious stream join problem definition

OSJ taxonomy

Two families of oblivious stream joins

Agenda

- 1. Motivation
- 2. Understanding equi-joins in TEEs
- 3. Relational Join for TEEs
- 4. Privacy-preserving stream joins
- 5. Conclusions



Conclusions



- TEEs offer a more private future
- Performance and security issues to solve on the way
- Solutions to three core aspects of management of sensitive data:
 - Adoption of TEEs
 - Efficient Query Processing in TEEs
 - Privacy-Preserving Stream Processing
- Future work
- Our work makes privacy-preserving data processing more accessible

