

2023.8.29 VLDB'23

OLIVE: Oblivious Federated Learning on TEE against the risk of Sparsification



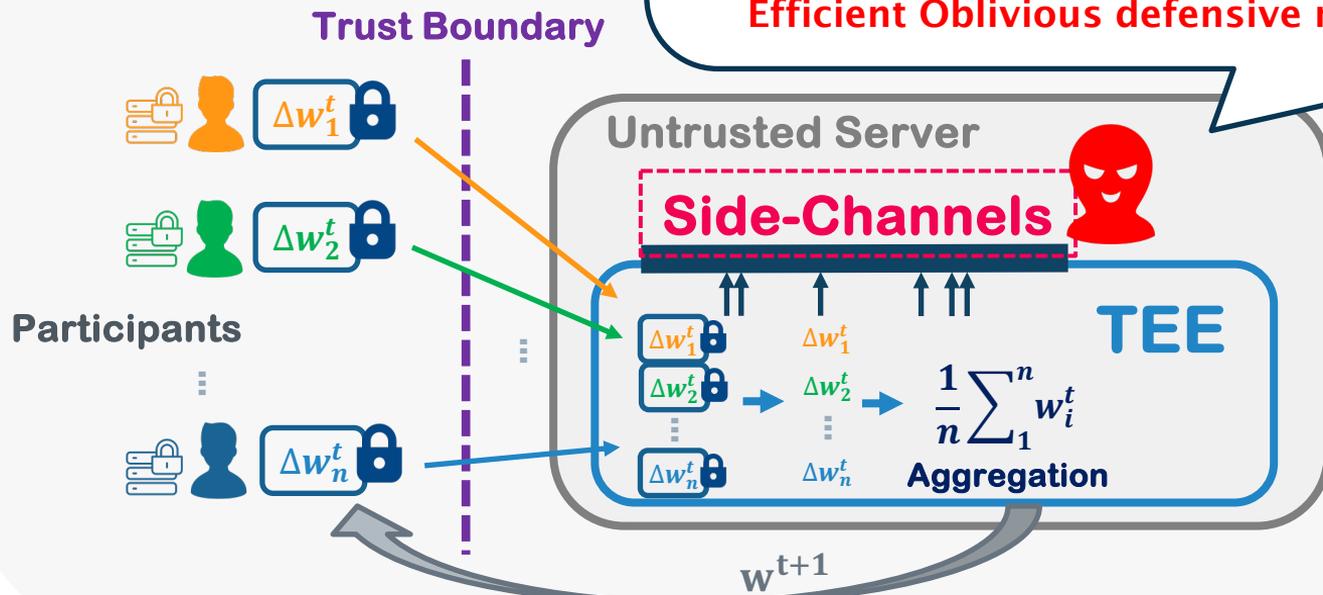
Fumiyuki Kato¹, Yang Cao², Masatoshi Yoshikawa³

¹*Kyoto University*, ²*Hokkaido University*, ³*Osaka Seikei University*,

Summary

- Motivation: FL with server-side TEE (Figure below)
- Problem: Privacy risk of Side-channels in FL with TEE
- Findings: Sparsified parameter can cause privacy leak
- Solution: Efficient Oblivious defensive mechanism

Federated Learning

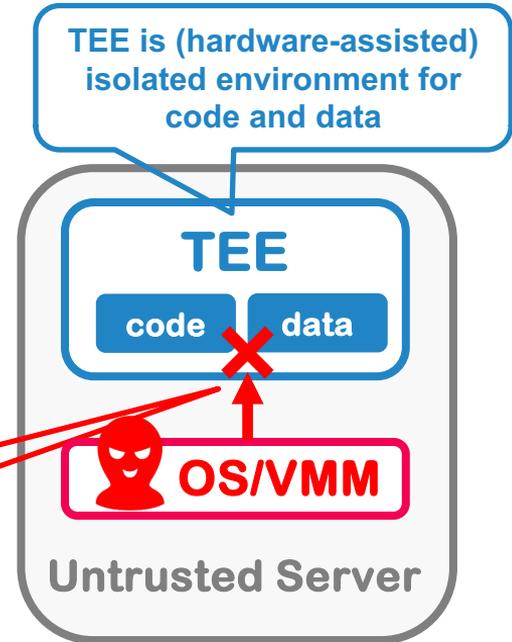


Background

Trusted Execution Environment (TEE)

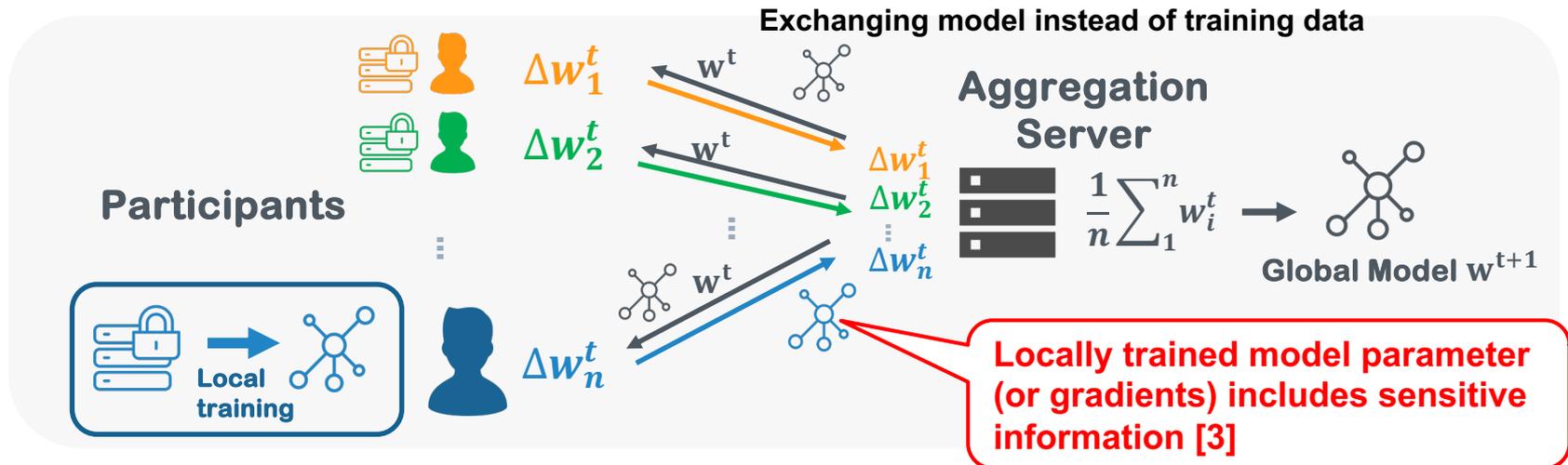
- TEE enables secure computation on remote machine
- Intel SGX – one of TEE implementations
 - 1. Memory encryption
 - Can hide code and data against privileged software (OS/VMM)
 - 2. Remote Attestation
 - Can verify the integrity of the code and data externally
- **Memory access pattern leakage via side-channels**
 - Cache-based (*Prime+Probe*) [1]
 - Page-based [2]

Access patterns can be visible regardless of memory encryption



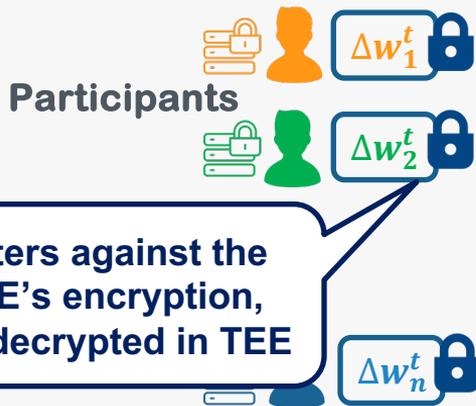
Federated Learning (FL)

- Collaborative ML scheme with
 - many participants
 - a central aggregation server
- **Problem: Locally trained model is sufficient to leak sensitive information**



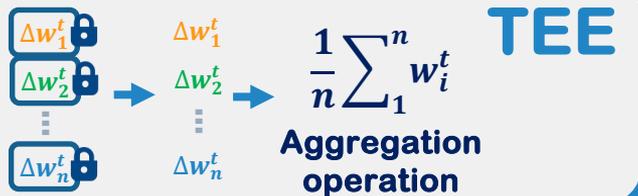
FL with server-side TEE

FL with TEE



Hide parameters against the server by TEE's encryption, which is only decrypted in TEE

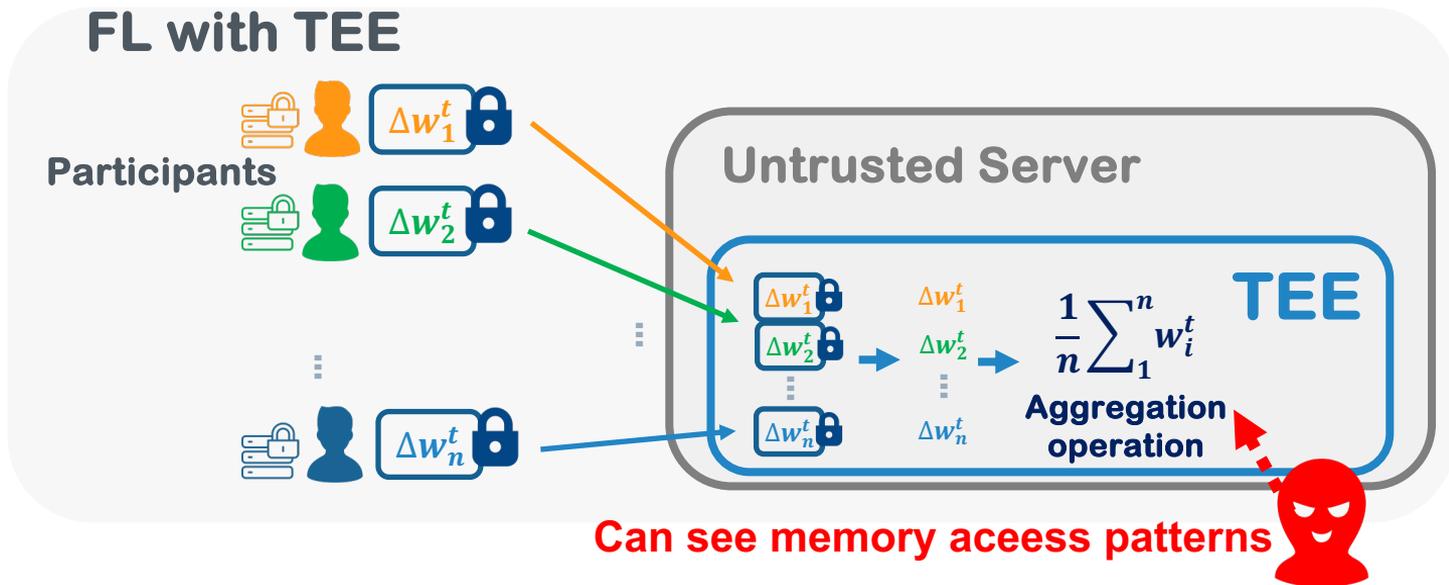
Untrusted Server



The combination of CDP improves the utility of DP-FL. (vs Shuffle DP)

	Trust model	Utility
CDP-FL [4, 28, 50, 84]	Trusted server	Good
LDP-FL [45, 74, 81, 92]	Untrusted server	Limited
Shuffle DP-FL [29, 44]	Untrusted server + Shuffler	Shuffle DP-FL \leq CDP-FL
OLIVE (Ours)	Untrusted Server with TEE	OLIVE = CDP-FL

Problem of FL with server-side TEE



Problem: The impact of side-channels of TEE is unknown

- What is the specific privacy risks?
- What is practical protection against the attacks?

Contributions

- **In FL with server-side TEE, we study both of Attack and Defense in terms of memory access pattern leaks**
- **Attack**
 - We show that the sparsified parameters often used in FL can leak sensitive information via memory access patterns
 - We demonstrate that privacy attacks are possible using information obtained from memory access patterns
- **Defense**
 - We design an efficient oblivious FL aggregation algorithm
 - We evaluate the proposed defensive mechanism on real-world scales

Attack analysis

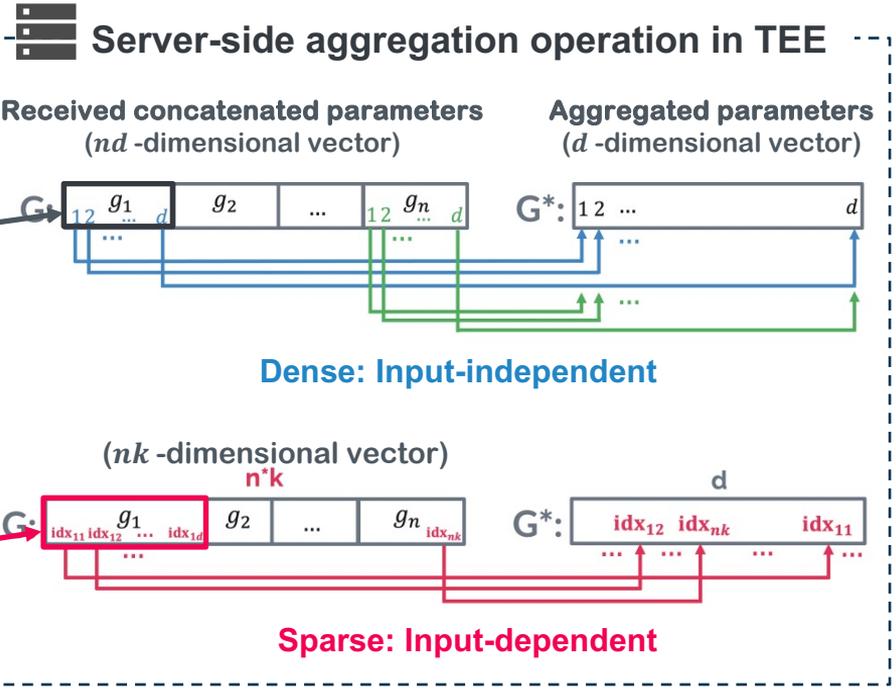
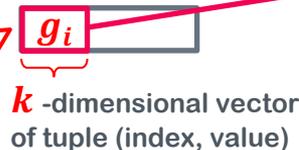
Memory Access Pattern Analysis on Aggregation Operation of FL

- If parameter is *Dense*, the memory access of the aggregation operation is independent from inputs



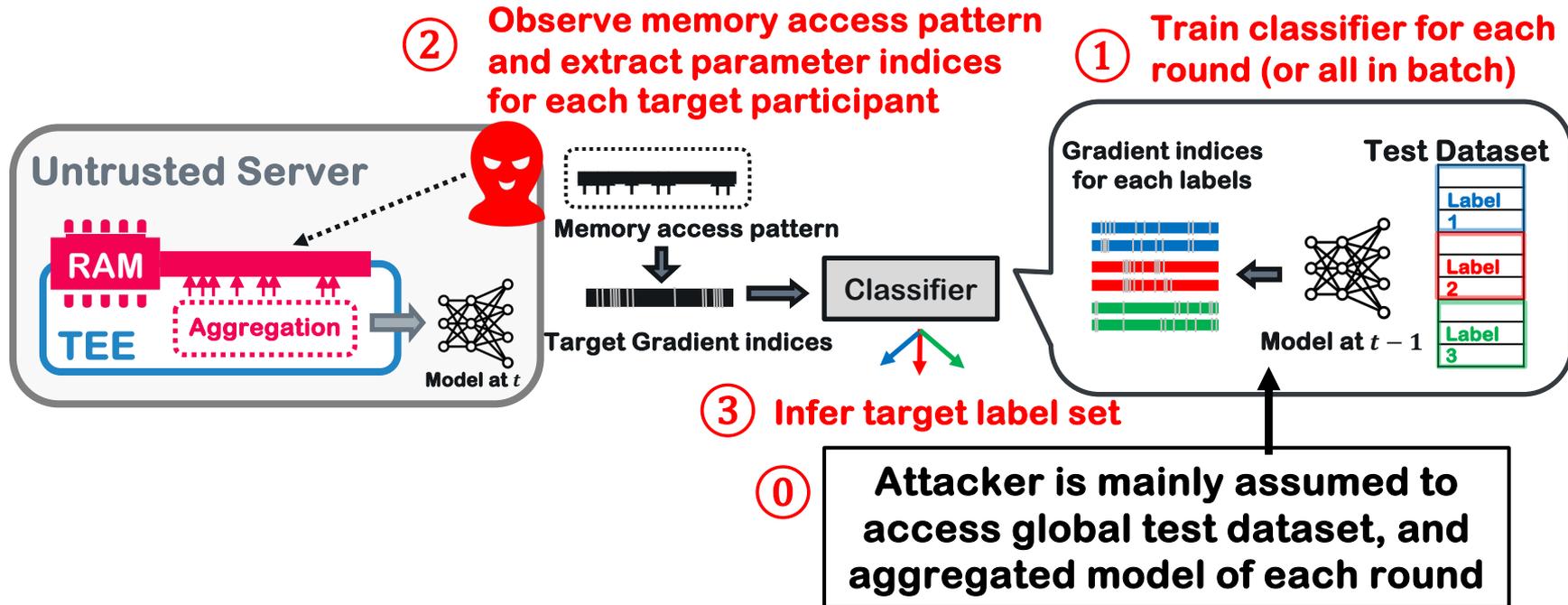
- In Sparsified setting**, memory access pattern can leak the selected parameter indices (data-dependent)

(Data-dependent) Top- k sparsification is often used in FL for better Comms-cost and better utility [4]



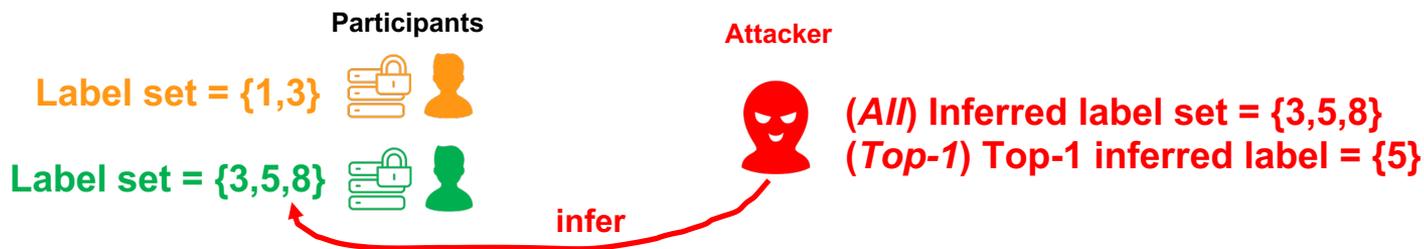
Overview of Attack Design

- To show the leaked information can leak private information
- The goal is to infer the sensitive label set of the target participant

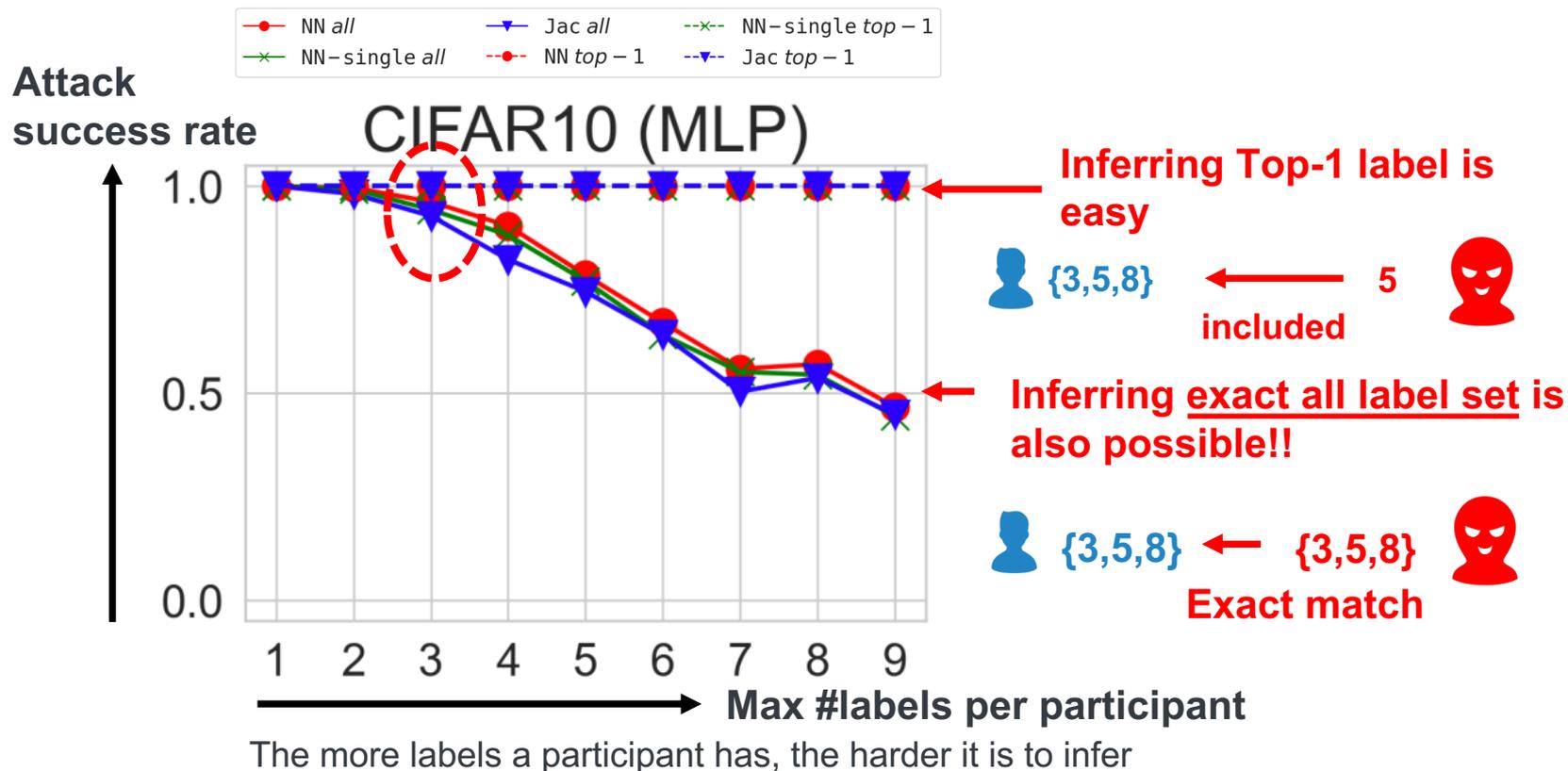


Empirical Evaluation: Setup

- Dataset
 - MNIST and CIFAR100 (, and Purchase100 (tabular dataset) in the paper)
- FL setting
 - Sparsification: with Top-10% sparsification
 - The maximum #Labels of each participant is controlled (#Participants: 1000)
- Evaluation Metrics
 - *All*: The ratio that predicted labels exactly match target label set
 - *Top-1*: The ratio that top-1 predict-scored label is included in target label set
 - Weakest privacy leak

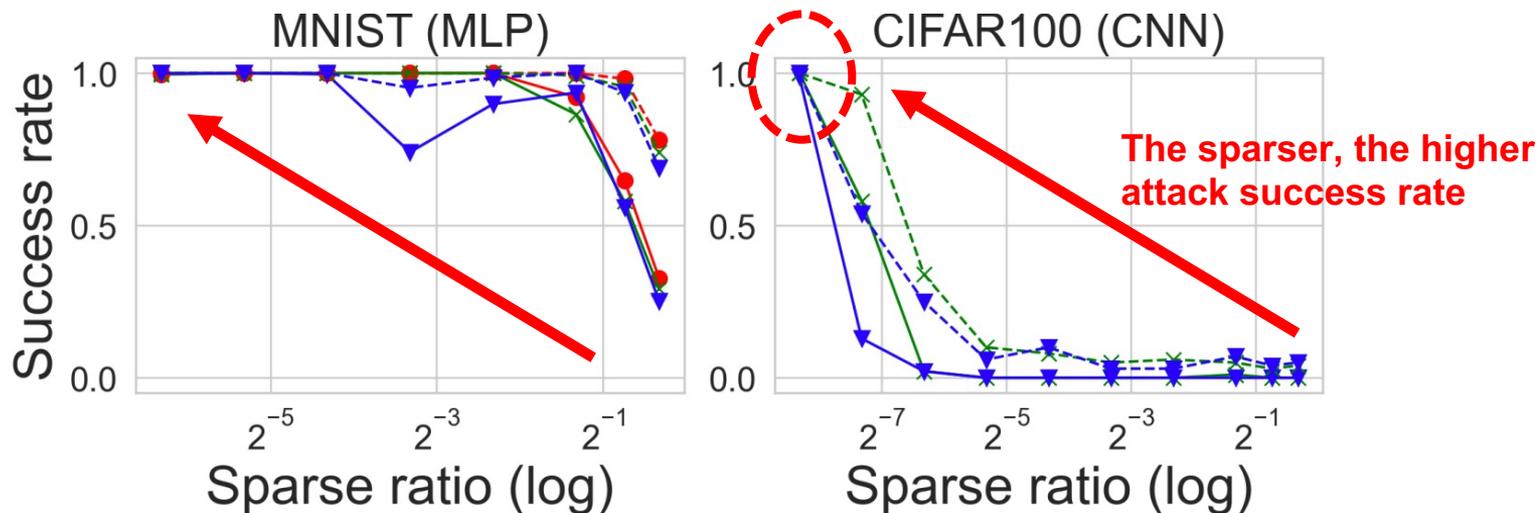


Result: Attack is overall successful



The sparser, the easier the attack

When sparse ratio=0.3%, attacks are almost 100% successful even we have 100 class labels.
(Common sparse ratio in FL is, for example, 0.1% [5])

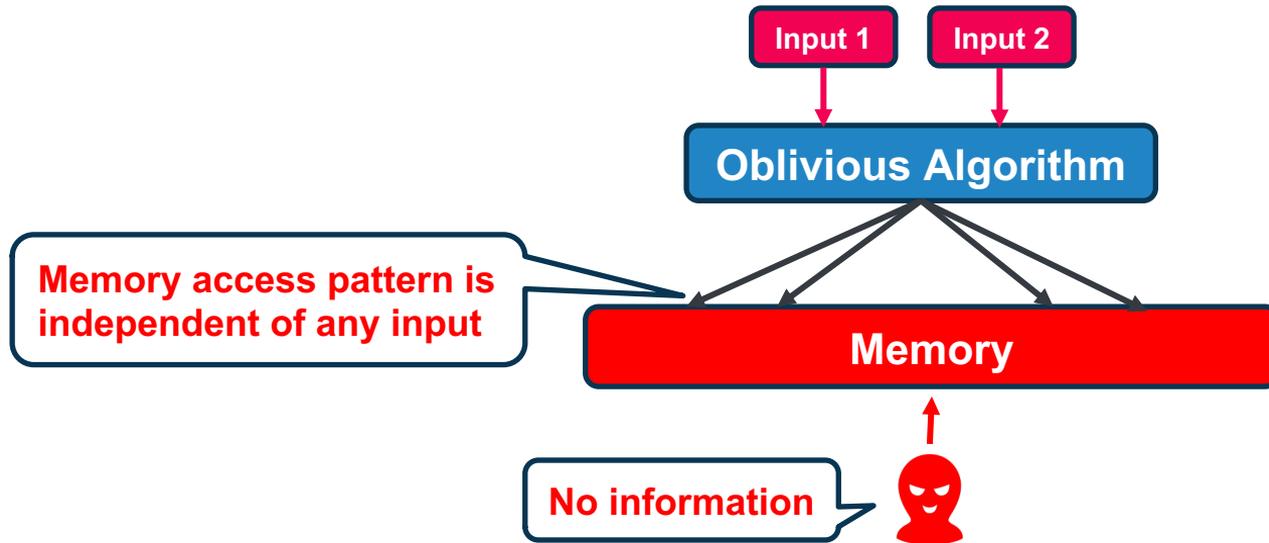


The number of labels of participants is fixed at 2.

Defensive mechanism

Oblivious Algorithm

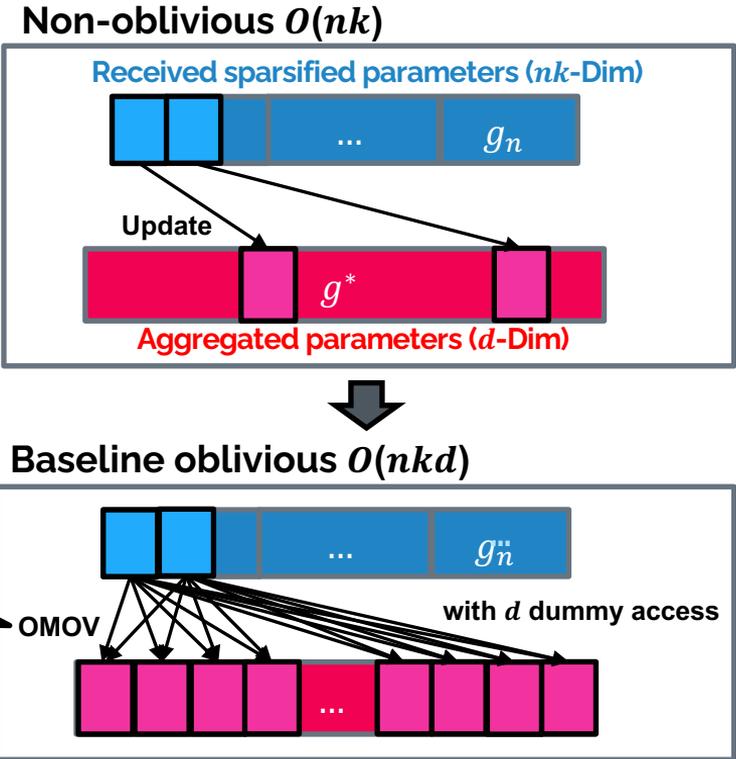
- **Oblivious algorithm** is an algorithm whose memory access pattern is independent of the input values
 - **No problem if memory access pattern leaks from side-channels**



Oblivious Algorithm: Baseline

- Non-oblivious method
 - $O(nk)$
- Baseline method (oblivious)
 - Full memory access approach
 - $O(nkd)$
 - n : #Participants
 - k : Dimensions of the sparsified model
 - d : Dimensions of the dense model

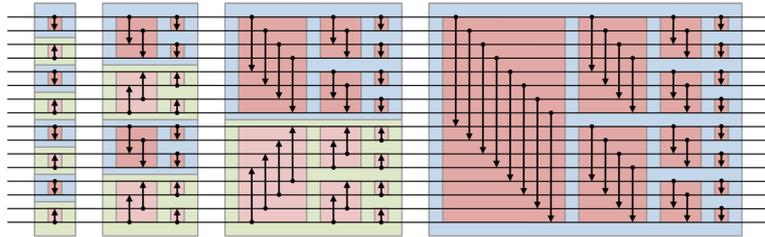
Using CMOV (of x86 instruction)-based oblivious primitive (**OMOV**) to ensure the program's execution path, including conditional branches, oblivious



But, can we make more efficient oblivious algorithm for this purpose?

Oblivious Algorithm: Advanced

- Advanced method
 - $O((nk + d) \log^2(nk + d))$
 - Using *oblivious sort*
 - Bitonic sort causes fixed memory access pattern



Sorting network of Bitonic sort
(https://en.wikipedia.org/wiki/Bitonic_sorter)

Received sparsified parameters (nk -Dim)



Advanced method

Next page 

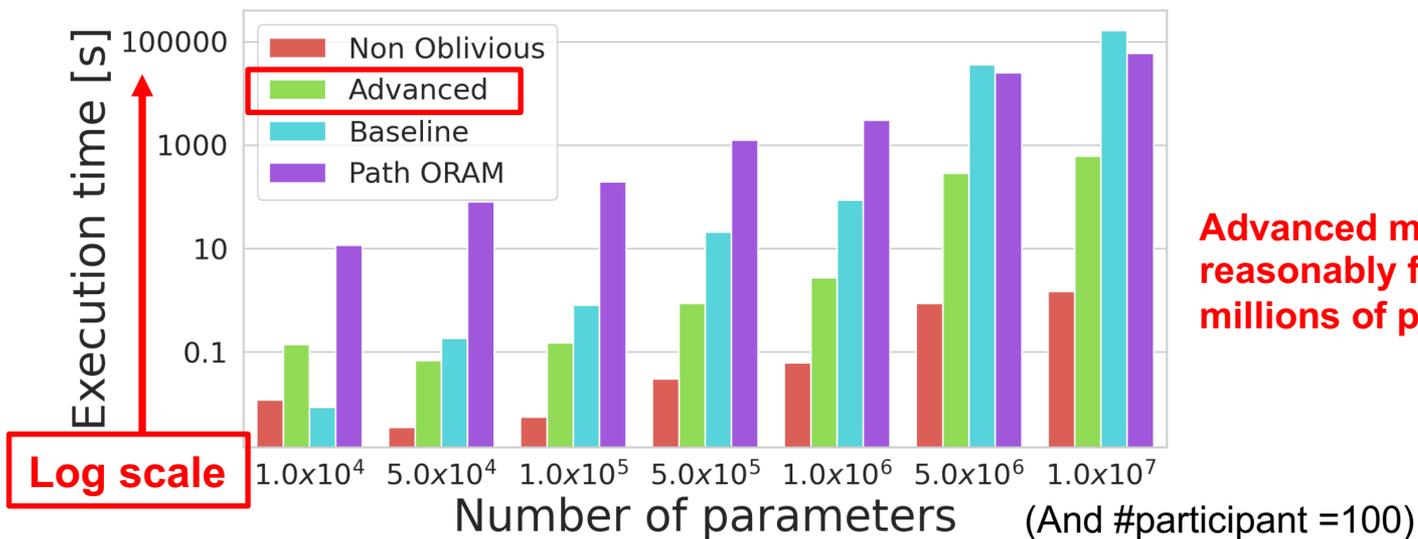


Aggregated parameters (d -Dim)

Evaluation

**Most efficient ORAM implementation
But, with some modification to
implement on SGX's security model [6]**

- Non Oblivious, Path ORAM (ZeroTrace [6]), Baseline, Advanced
- with Intel SGX (PRM: 128 MB)



**Advanced method is
reasonably fast with over
millions of parameters**

[6] Sajin Sasy, Sergey Gorbunov, and Christopher Fletcher. "ZeroTrace: Oblivious memory primitives from Intel SGX." Symposium on Network and Distributed System Security (NDSS). 2018.

Conclusions

- **In FL with server-side TEE, we studied both of Attack and Defense**
- **Attack**
 - We show that the sparsified parameters often used in FL can leak sensitive information via memory access patterns
 - We demonstrated that privacy attacks are possible using information obtained from memory access patterns
- **Defense**
 - We designed an efficient oblivious FL aggregation algorithm
 - We evaluated the proposed defensive mechanism on real-world scales

Appendix

Assumptions

- Attacker can access
 - Test dataset
 - Aggregated global model in each round
 - Observe access patterns (from side-channels)
- Top-k sparsification is used

Attack: Cache line-level

- Even if the granularity of the observation becomes a cachet line, the results don't change much

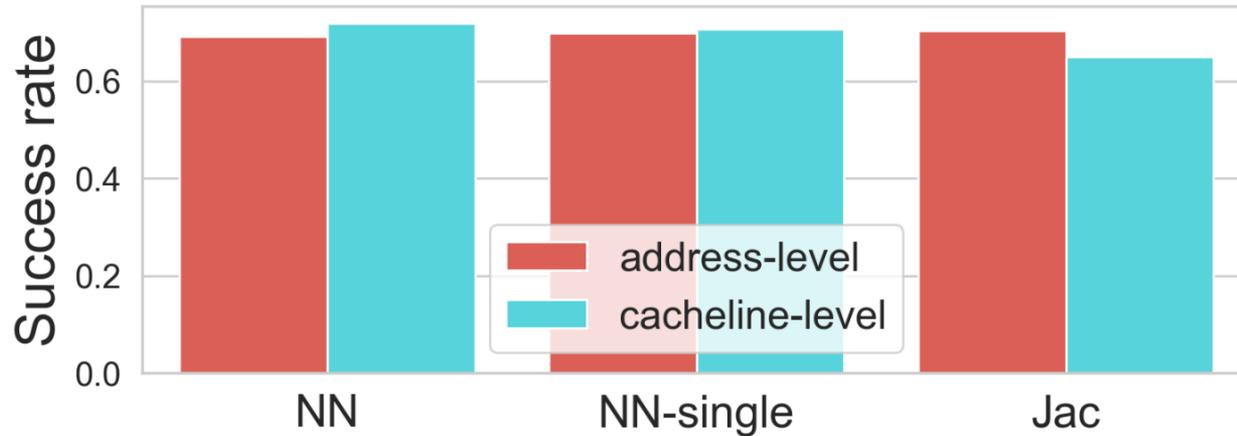
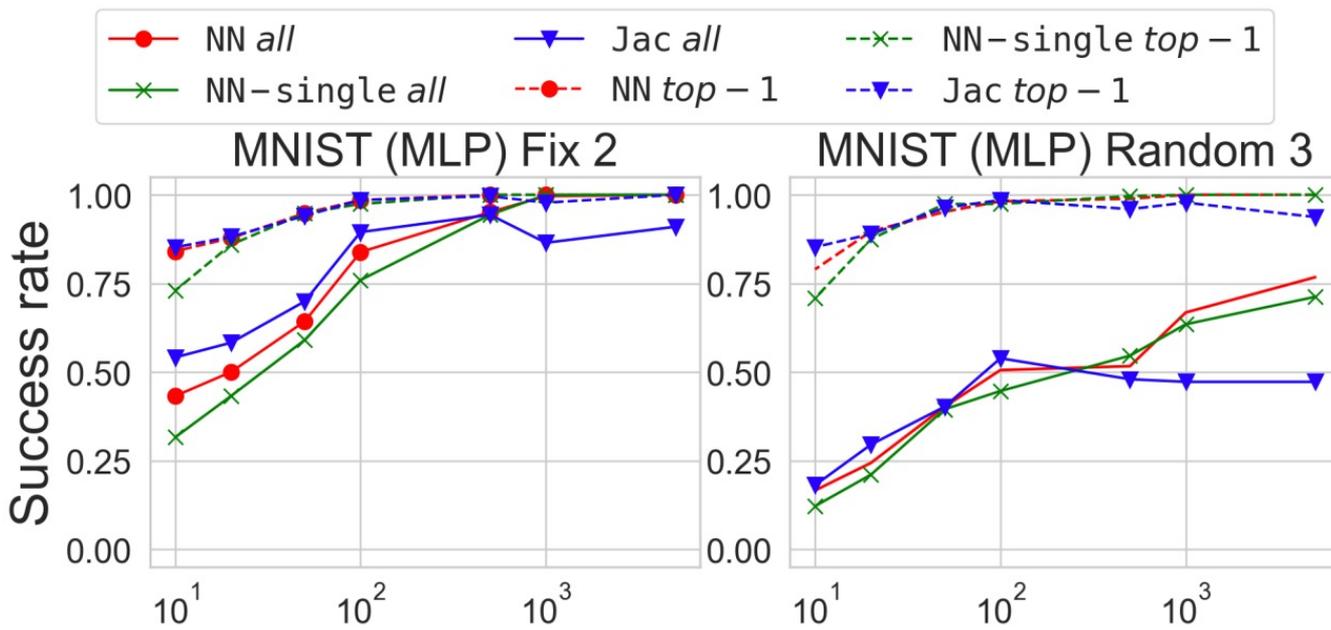


Figure 7: Cacheline-level leakage on CNN of CIFAR10: Attacks are possible with at least slightly less accuracy.

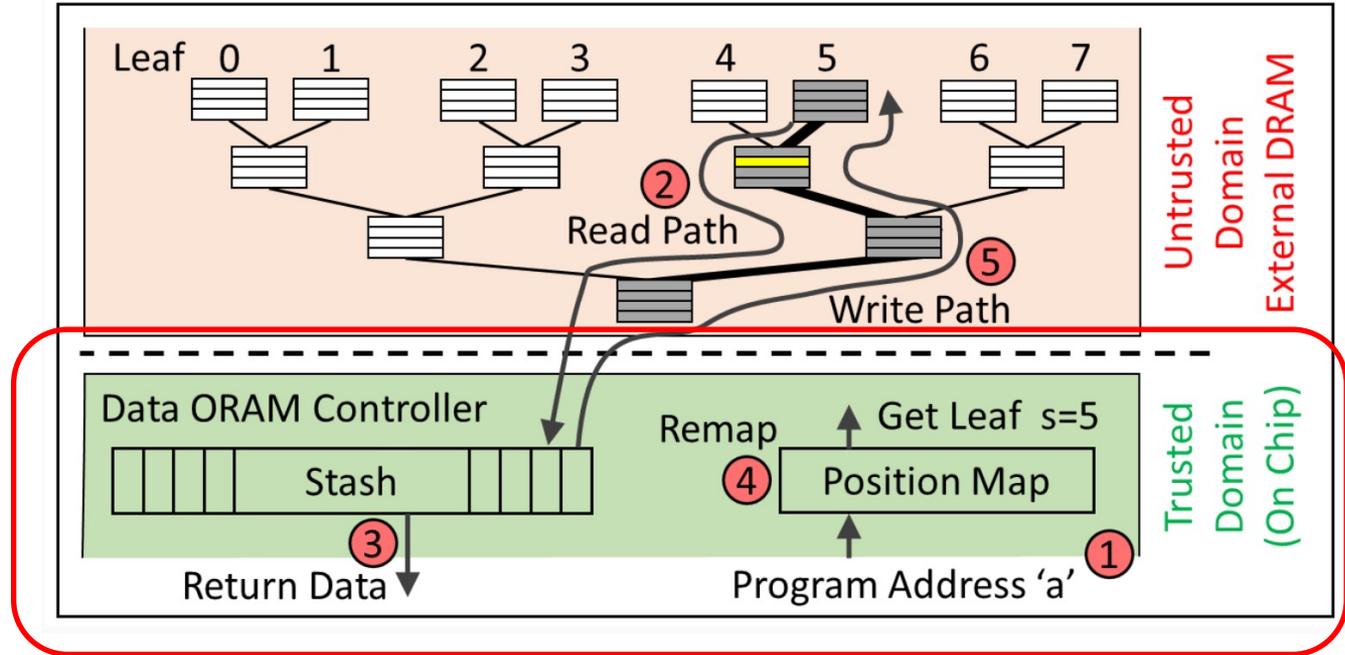
Attack: Various test dataset

- Test dataset can be very small



Path ORAM [A.1]

Access cost can be reduced to $O(\log N)$



In the trust model of SGX, Trusted Domain of ORAM also needs to be oblivious [A.2]

[Image source: <https://scl.engr.uconn.edu/research/oram.php>]

[A.1] Stefanov, E., Van Dijk, M., Shi, E., Fletcher, C., Ren, L., Yu, X. and Devadas, S., 2013, November. Path ORAM: an extremely simple oblivious RAM protocol. In Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security (pp. 299-310). ACM.

[A.2] Sajin Sasy, Sergey Gorbunov, and Christopher Fletcher. "ZeroTrace: Oblivious memory primitives from Intel SGX." Symposium on Network and Distributed System Security (NDSS). 2018. 26

Discussion: Differentially Obliviousness

- Chen et al [A.3] formalize DO algorithms

Definition 2.1 (Differentially oblivious (stateless) algorithms). *Let ϵ, δ be functions in a security parameter λ . We say that the stateless algorithm M satisfies (ϵ, δ) -differential obliviousness, iff for any neighboring inputs I and I' , for any λ , for any set S of access patterns, it holds that*

$$\Pr[\mathbf{Accesses}^M(\lambda, I) \in S] \leq e^{\epsilon(\lambda)} \cdot \Pr[\mathbf{Accesses}^M(\lambda, I') \in S] + \delta(\lambda)$$

where $\mathbf{Accesses}^M(\lambda, I)$ is a random variable denoting the ordered sequence of memory accesses the algorithm M makes upon receiving the input λ and I .

- NIPS '19 [A.4] and CCS '18 [A.5] proposed similar algorithms to guarantee DO

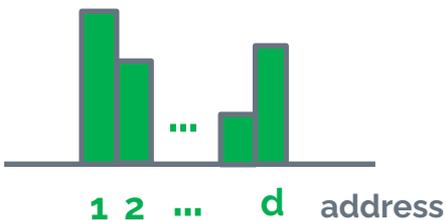
[A.3] Chan et al, Foundations of differentially oblivious algorithms. SIAM 2019.

[A.4] Joshua et al, An Algorithmic Framework For Differentially Private Data Analysis on Trusted Processors. NIPS 2019

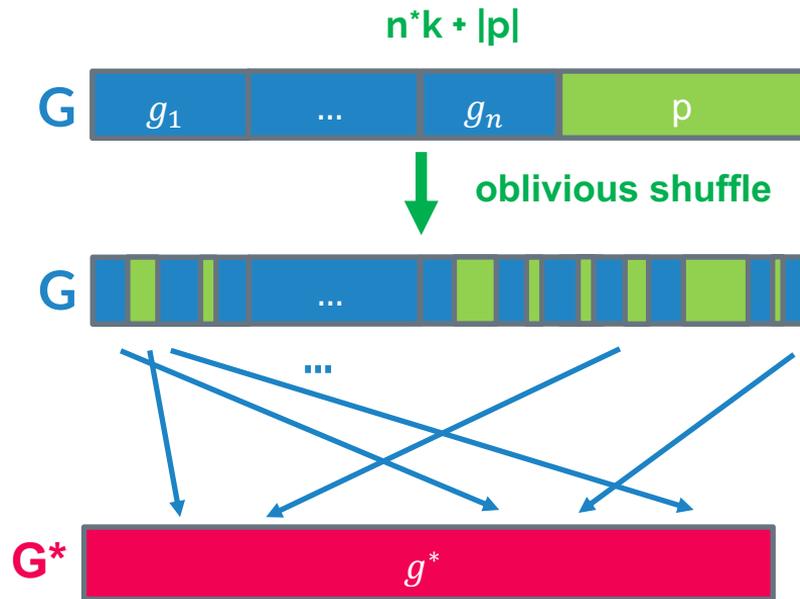
[A.5] Mazloom et al. Secure Computation with Differentially Private Access Patterns. CCS 2018

Discussion: Differentially Obliviousness

- NIPS '19 [A.4]
 - p padding, **oblivious shuffle**
 - $O((nk + |p|)\log^2(nk + |p|))$
 - leaks differentially private histogram of all indices



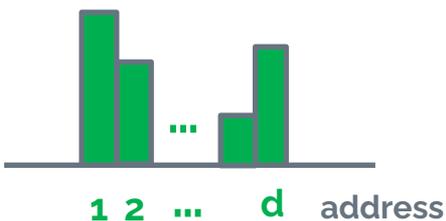
Access pattern histogram with DP



- It doesn't work due to huge padding size**
- (1) **Sensitivity (i.e., d) can be too large ($|p|$ is $O(kd)$)**
 - (2) **Can only use one-sided noise**

Discussion: Differentially Obliviousness

- NIPS '19 [A.4]
 - p padding, oblivious shuffle
 - $O((nk + |p|) \log^2(nk + |p|))$
 - leaks differentially private histogram of all indices



Access pattern histogram with DP

p is too large in FL setting

- noise vector $z \in \mathbb{R}^d$
 - $z_i \sim \text{Lap}(2k/\epsilon)$
 - d -dimensional
- Then, p is $O(kd)$
 - kd is very large
- Remember
 - $O((nk + p) \log^2(nk + p))$
- Moreover, it is necessary to allocate memory for the padded data, which is very incompatible with SGX that has poor memory

Algorithms

Algorithm 3 Baseline

Input: $g = g_1 \parallel \dots \parallel g_n$: concatenated gradients, nk length

Output: g^* : aggregated parameters, d length

```
1: initialize aggregated gradients  $g^*$ 
2: for each  $(idx, val) \in g$  do
3:   /*  $c$  is the number of weights included in one cacheline */
4:   /* offset indicates the position of  $idx$  in the cacheline */
5:   for each  $(idx^*, val^*) \in g^*$  if  $idx^* \equiv \text{offset} \pmod{c}$  do
6:      $flag \leftarrow idx^* == idx$            ▶ target index or not
7:      $val' \leftarrow \text{o\_mov}(flag, val^*, val^* + val)$ 
8:     write  $val'$  into  $idx^*$  of  $g^*$ 
9: return  $g^*$ 
```

Algorithm 4 Advanced

Input: $g = g_1 \parallel \dots \parallel g_n$: concatenated gradients, nk length

Output: g^* : aggregated parameters, d length

```
1: /* initialization: prepare zero-valued gradients for each index */
2:  $g' \leftarrow \{(1, 0), \dots, (d, 0)\}$            ▶ all values are zero
3:  $g \leftarrow g \parallel g'$                        ▶ concatenation
4: /* oblivious sort in  $O((nk + d) \log^2(nk + d))$  */
5: oblivious sort  $g$  by index
6: /* oblivious folding in  $O(nk + d)$  */
7:  $idx \leftarrow$  index of the first weight of  $g$ 
8:  $val \leftarrow$  value of the first weight of  $g$ 
9: for each  $(idx', val') \in g$  do ▶ Note: start from the second weight of  $g$ 
10:   $flag \leftarrow idx' == idx$ 
11:  /*  $M_0$  is a dummy index and very large integer */
12:   $idx_{prior}, val_{prior} \leftarrow \text{o\_mov}(flag, (idx, val), (M_0, 0))$ 
13:  write  $(idx_{prior}, val_{prior})$  into  $idx' - 1$  of  $g$ 
14:   $idx, val \leftarrow \text{o\_mov}(flag, (idx', val'), (idx, val + val'))$ 
15: /* oblivious sort in  $O((nk + d) \log^2(nk + d))$  */
16: oblivious sort  $g$  by index again
17: return take the first  $d$  values as  $g^*$ 
```

Oblivious Primitives

- **O_MOV**
 - Using CMOV (x86 instruction)
 - set the value of either “a” or “b” in the register depending on the conditional flag
 - **adversary cannot see**
 - constructs o_mov (oblivious move), o_swap (oblivious swap), o_write (oblivious write)

```
38  #[inline]
39  pub fn o_swap<T>(flag: isize, x: &T, y: &T) {
40      unsafe {
41          llvm_asm!(
42              "test %rax, %rax \n\t"
43              movq (%r8), %r10 \n\t"
44              movq (%rdx), %r9 \n\t"
45              mov %r9, %r11 \n\t"
46              cmovnz %r10, %r9 \n\t"
47              cmovnz %r11, %r10 \n\t"
48              movq %r9, (%rdx) \n\t"
49              movq %r10, (%r8) \n\t"
50              :
51              : "{rax}"(flag), "{rdx}"(x), "{r8}"(y)
52              : "rax", "rdx", "r8", "r9", "r10", "r11"
53              : "volatile"
54          );
55      }
56  }
57
```

o_swap by x86

Oblivious Primitives

- **O_WRITE**
 - n times oblivious mov (o_mov)
 - x16 faster by cache-line optimization in Baseline method

O_MOV



```
for i in 0..(num_of_weights / CACHE_LINE_NUM_OF_WEIGHT) {  
  // Only one of the values to be read needs to be written.  
  // Since it is not known from the outside which value was written  
  let cache_line_addr = dst_base.offset(CACHE_LINE_NUM_OF_WEIGHT * i + ith) as *mut f32;  
  let flag = (cache_line_addr == addr) as isize;  
  
  let ret: f32;  
  llvm_asm!(  
    "xor %rcx, %rcx \n\t"  
    "mov %r8, %rcx \n\t"  
    "test %rcx, %rcx \n\t"  
    "cmovnz %rdx, %rax \n\t"  
    : "{rax}"(ret)  
    : "{r8}"(flag), "{rax}" (*cache_line_addr), "{rdx}" (val)  
    : "rax", "rcx", "rdx", "r8"  
    : "volatile"  
  );  
  *cache_line_addr = ret;  
  last += 1;  
}
```

} o_mov