# Privacy Preserving Machine Learning

Emil Hørning & Alexander Ramos

Supervisor: Bernardo Machado David

# TABLE OF CONTENTS

# 01

# Privacy Preserving Machine Learning

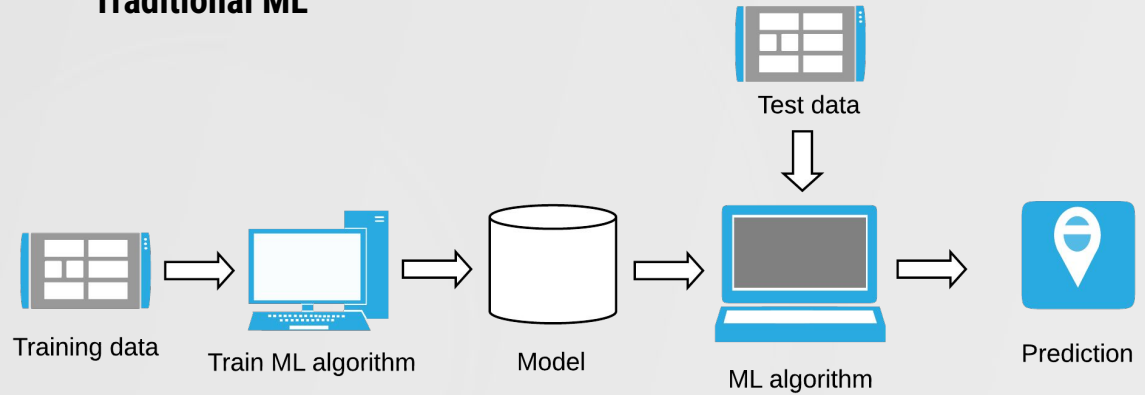A brief introduction

# Privacy Preserving Machine Learning

- Data with X number of rows (training examples), and Y number of columns (features)

- Label controls the training

- ML algorithm learns a hidden function which describes the data

| Input | Blood pressure | Fat% | … | Label |
|---|---|---|---|---|
| Alex | 132 | 20 | … | Unhealthy |
| Emil | 92 | 8 | … | Healthy |
| … | … | … | … | … |

# Privacy Preserving Machine Learning

**Traditional ML**

- Data holders outsource ML

- Data may be sensitive



Test data

Training data → Train ML algorithm → Model → ML algorithm → Prediction

# Privacy Preserving Machine Learning

**Private machine learning**

- Data holders outsource ML

- Data may be sensitive

Test data

Training data
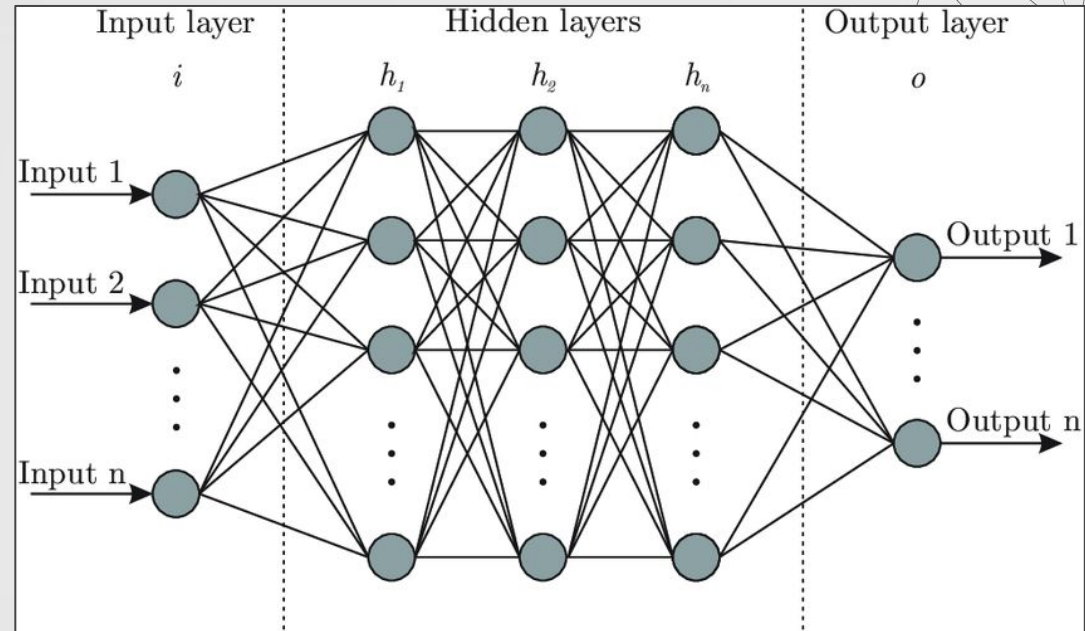
Train ML algorithm

Model

ML algorithm

Prediction

6

# 02

# Background theory

(Deep) neural networks
Multi party computation/secret shares
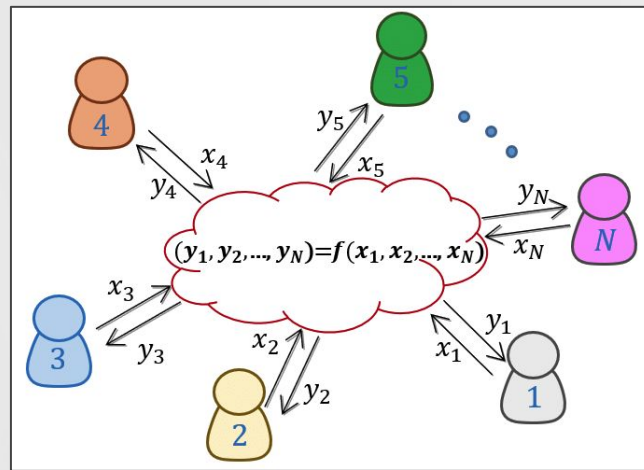Types of security/adversarial models

# 02: (Deep) Neural Networks

- A neural network estimates the parameters of "the hidden" function which describes a dataset
- Consists of **input, hidden,** and **output** layers
- Weights (parameters) are the edges
- **Forward propagation**
  - Linear combination of input values and weights $\Longrightarrow$ activation function
  - Matrix multiplication
- **Backpropagation**
  - Adjust weights according to output and true label

# 02: Secure Multi-Party Computation

- N parties want to compute some function without revealing content of own input

- Functions are implemented as protocols that are private and secure

- Such MPC protocols can be used to do private machine learning



$$(y_1, y_2, ..., y_N) = f(x_1, x_2, ..., x_N)$$

# 02: Types of security

## Input privacy

The input remains private and no party learns the secret of other parties
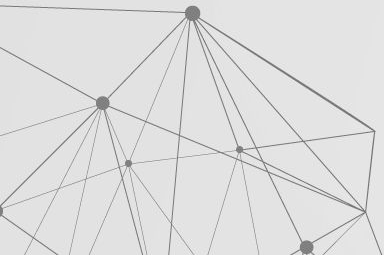
## Correctness

Any number of colluding dishonest parties participating in the protocol, should not be able to, by deviating from the protocol, force or trick an honest party to output something incorrect

### Robust
Protocol achieves this in any scenario

### Abort
Honest party detects deviation and aborts

# 02: Adversarial model

## Passive adversary

- Adversary that corrupts 1 party and observes the protocol
- Tries to learn secrets
- Follows the protocol dutifully

---

## Active adversary

- Adversary that corrupts 1 party and deviates from the protocol
- Tampers with inputs and outputs to alter result of protocol
- Tampers with inputs and outputs to learn secrets
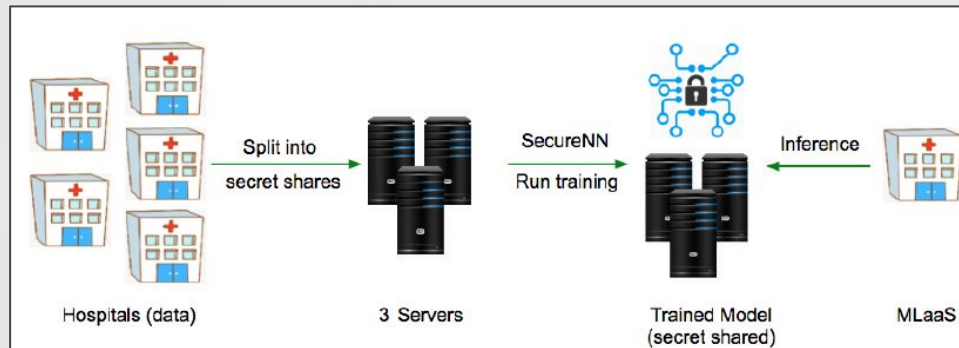- May halt the protocol to prevent completion

# 03

# SecureNN

SecureNNs PPML approach
Security guarantees
Protocol structure
End-to-end protocol
Communication bottlenecks

# 03: SecureNNs PPML approach

- MPC protocols (3 parties)

- Input $\implies$ secret shares $\implies$ send to the 3 parties

- Run interactive protocol to train a neural network

- Trained model retained as secret shares
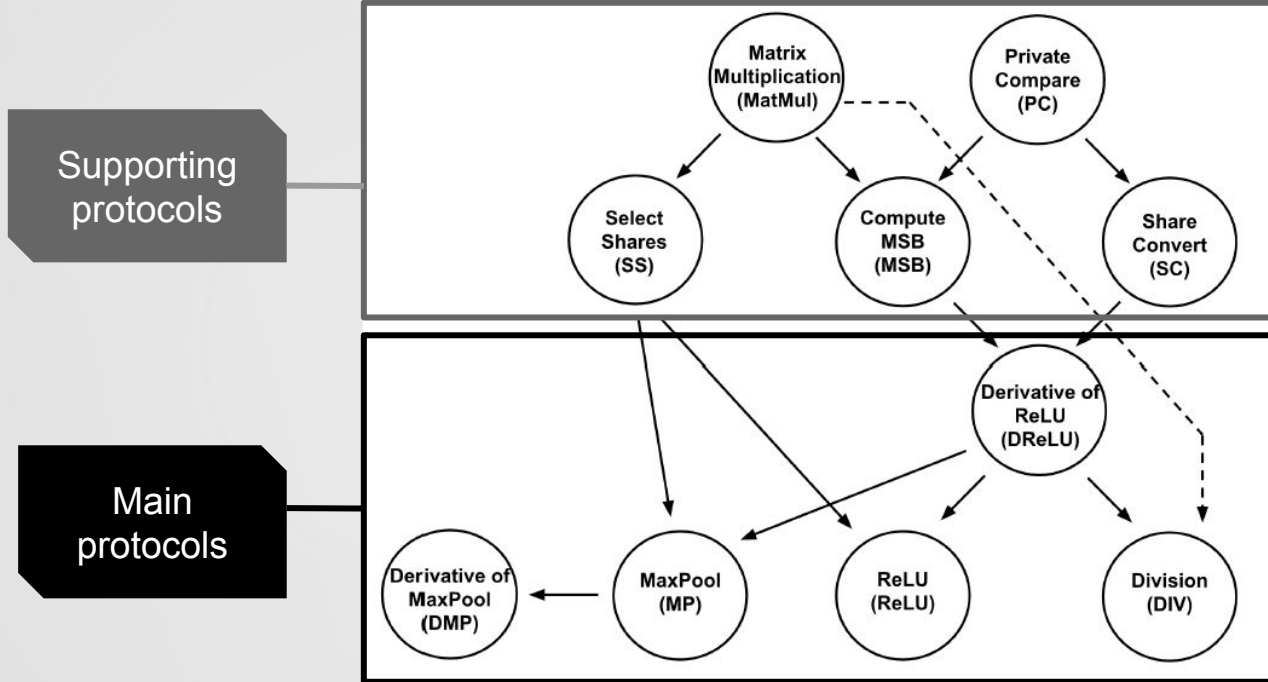
- Private prediction

# 03: Security guarantees

## Passive adversary

- SecureNN provides full input privacy against a single passive corrupted party
- Proved with simulation proof
- Correctness of the protocol follows trivially

## Active adversary

- SecureNN provides full input privacy against an active adversary corrupting 1 party
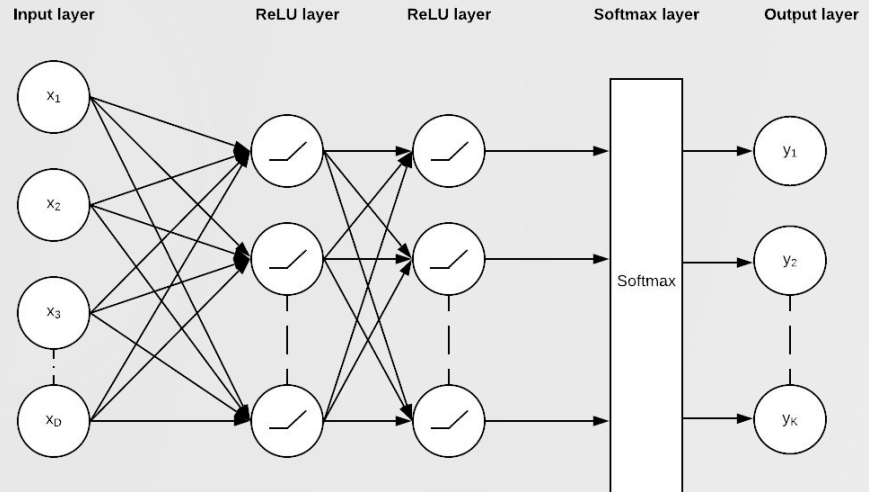- SecureNN can't provide correctness for an active adversary

# 03: Protocol structure



Supporting protocols

Main protocols

# 03: End-to-end protocol

- Construct a NN using a sequence of protocol calls

- E.g. a 3 layer network with a softmax output layer

1. Call MatMul
2. Call ReLU
3. Call MatMul
4. Call ReLu
5. Call DIV

$$ASM(u_i) = \frac{ReLU(u_i)}{\sum ReLU(u_i)}$$



Input layer   ReLU layer   ReLU layer   Softmax layer   Output layer

# 03: Communication bottlenecks

- Function **DReLU** uses 9 rounds of communication

- Uses trick to obtain MSB by converting to odd ring

**Algorithm 6** ReLU', $\Pi_{\text{DReLU}}(\{P_0, P_1\}, P_2)$:

**Input:** $P_0, P_1$ hold $\langle a \rangle_0^L$ and $\langle a \rangle_1^L$, respectively.

**Output:** $P_0, P_1$ get $\langle \text{ReLU}'(a) \rangle_0^L$ and $\langle \text{ReLU}'(a) \rangle_1^L$.

**Common Randomness:** $P_0, P_1$ hold random shares of $0$ over $\mathbb{Z}_L$, denoted by $u_0$ and $u_1$ resp.

1: For $j \in \{0, 1\}$, parties $P_j$ computes $\langle c \rangle_j^L = 2\langle a \rangle_j^L$.

2: $P_0, P_1, P_2$ run $\Pi_{\text{SC}}(\{P_0, P_1\}, P_2)$ with $P_0, P_1$ having inputs $\langle c \rangle_j^L$ & $\langle c \rangle_1^L$ & $P_0, P_1$ learn $\langle y \rangle_0^{L-1}$ & $\langle y \rangle_1^{L-1}$, resp.

3: $P_0, P_1, P_2$ run $\Pi_{\text{MSB}}(\{P_0, P_1\}, P_2)$ with $P_j, j \in \{0, 1\}$ having input $\langle y \rangle_j^{L-1}$ & $P_0, P_1$ learn $\langle \alpha \rangle_0^L$ & $\langle \alpha \rangle_1^L$, resp.

4: For $j \in \{0, 1\}$, $P_j$ outputs $\langle \gamma \rangle_j^L = j - \langle \alpha \rangle_j^L + u_j$.

| Protocol | Rounds | Communication |
|---|---|---|
| MatMul$_{m,n,v}$ | 2 | $2(2mn + 2nv + mv)\ell$ |
| MatMul$_{m,n,v}$ (with PRF) | 2 | $(2mn + 2nv + mv)\ell$ |
| SelectShare | 2 | $5\ell$ |
| PrivateCompare | 1 | $2\ell \log p$ |
| ShareConvert | 4 | $4\ell \log p + 6\ell$ |
| Compute MSB | 5 | $4\ell \log p + 13\ell$ |

17

# 04

## Improving SecureNN

Initial attempt
Alternative approach (The Genome Paper)
Bit-decomposition
Optimised bit-decomposition (ComposeNet)

# 04: Initial attempt

- Remove 'Share Convert' and assume shares to be in $Z_{L-1}$ already

- First month was spent trying to compile/run SecureNN code

- We were in contact with main author Sameer Wagh, but it wasn't very helpful

- Later we realized our logic was wrong about removing 'Share Convert' and started looking for other ways to optimize ComputeMSB

# 04: Alternative PPML approach



High Performance Logistic Regression for Privacy-Preserving Genome Analysis

- The Genome Paper* released during our thesis
- Presents 2-PC method for performing logistic regression
- Same method can be used for private neural networks
- Won first prize in Track 4 of iDash 2019
- Secure bit-decomposition to obtain shares of MSB

*"High performance logistic regression for privacy preserving genome analysis" by De Cock et. al.

# 04: Bit-decomposition

- Protocols exists that produces XOR shares of all bits of a shared value.

- Utilize this to get shares of MSB.

- Genome Paper presents 3 Bit Decomposition protocols for $\lambda$-bit integers
  - Linear Adder circuit (Linear, $\lambda$ rounds)
  - Speculative optimized ($2 + \log_2(\lambda)$ rounds)
  - Hyper Optimized ComposeNET ($1 + \log_2(\lambda)$ rounds) (BitDecompOPT)

- We implemented ComputeMSB and BitDecompOPT to compare.

# 04: BitDecompOPT

- Based on the speculative method
  - Logarithmic rounds

- **Observation**: The carry bit depends on 2 signals

- Computing all matrix compositions $M_1$ to $M_i$ yields the carry vector in upper right hand entry of resulting matrix

$$s_i = a_i \oplus b_i \oplus c_{i-1}$$

$$c_i = a_i b_i \oplus a_i c_{i-1} \oplus b_i c_{i-1}$$

**Generate**

$$g_i = a_i b_i$$

**Propagate**

$$p_i = a_i \oplus b_i$$

$$s_i = p_i \oplus c_{i-1}$$

$$c_i = g_i + p_i c_{i-1}$$

$$\begin{bmatrix} c_i \\ 1 \end{bmatrix} = \begin{bmatrix} p_i & g_i \\ 0 & 1 \end{bmatrix} \begin{bmatrix} c_{i-1} \\ 1 \end{bmatrix} = M_i \begin{bmatrix} c_{i-1} \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} c_i \\ 1 \end{bmatrix} = \begin{bmatrix} p_i & g_i \\ 0 & 1 \end{bmatrix} \begin{bmatrix} p_{i-1} & g_{i-1} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} c_{i-2} \\ 1 \end{bmatrix} = M_i M_{i-1} \begin{bmatrix} c_{i-2} \\ 1 \end{bmatrix}$$

# 04: BitDecompOPT

- Matrix composition network: **ComposeNet**
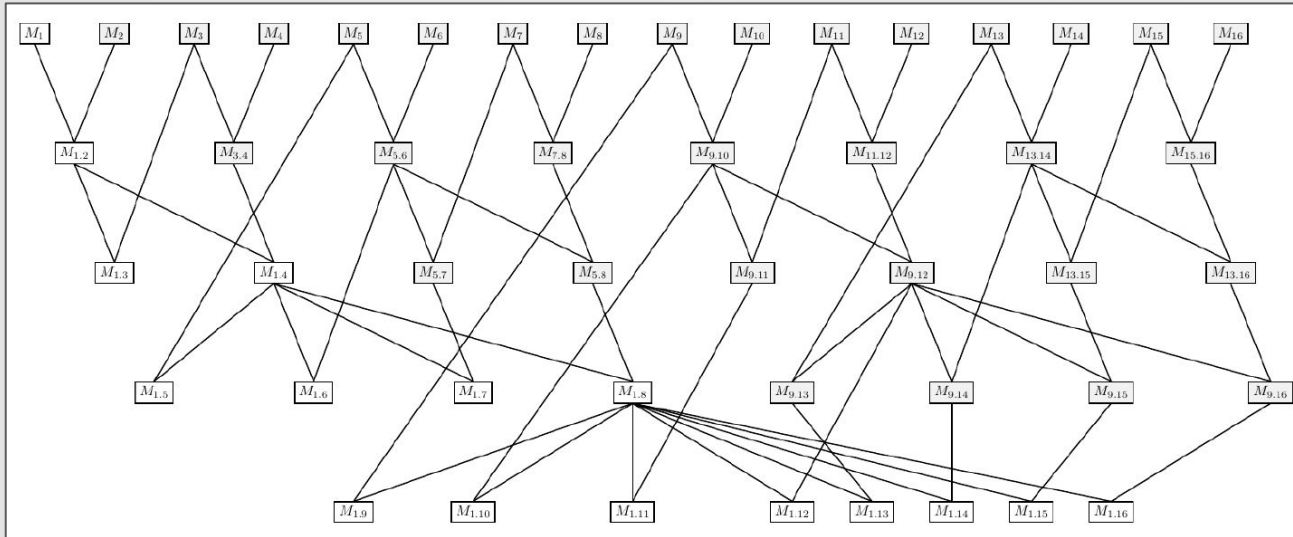
- Logarithmic depth network

# 05
# Experiments

Our implementation
Design of experiments

# 05: Our implementation

**Architectural implementations**

- Secret share logic
- Data structures for handling integers in $Z_L$
- Networking infrastructure

**Protocol implementations**

- Multiplication 2-party: Single-MatMul, Single-Mult, List-MatMul, List-Mult
- Multiplication 3-party: Single-Mult
- SC (Share Convert)
- PC (Private Compare)
- ComputeMSB
- ComposeNET
- BitDecomp
- BitDecompOPT

# 05: Design of experiments

- Bitlength L = $2^{64}$, and prime $p$ = 67
- Each protocol was run 1000 times

**Distributed experiments**



**Local experiments**

Windows 10 (64 bit)
4 cores, 16 GB RAM
Local ping: 0.19 ms

London

A
8.05 ms

B
12.23 ms

Digitalocean.com
3x Ubuntu 18.04 lts (64 bit)
1 core, 1 GB RAM

Amsterdam

Frankfurt

C
18.21 ms

# 06
## Results

# 06: Results

## Local vs distributed

| | Local time (s) | | Dist. time (s) | | Comm (bytes) | |
|---|---|---|---|---|---|---|
| Input size → | 1 | 1000 | 1 | 1000 | 1 | 1000 |
| BitDecompOPT A | - | - | 0.0686 | 68.58 | 1674 | 1674000 |
| BitDecompOPT B | - | - | 0.0726 | 72.58 | 1674 | 1674000 |
| BitDecompOPT C | - | - | 0.0825 | 82.47 | 1674 | 1674000 |
| BitDecompOPT Avg | 0.0304 | 30.48 | 0.0745 | 74.54 | 1674 | 1674000 |
| SC + ComputeMSB | 0.0221 | 22.04 | 0.1015 | 101.52 | 1987.493 | 1987493 |

## Raw local computation

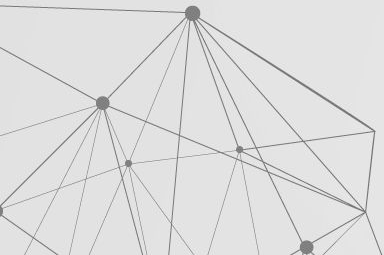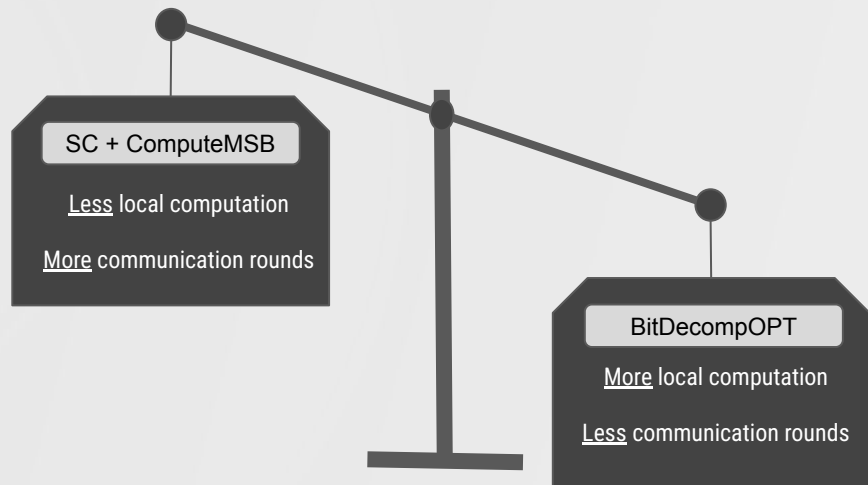| | time (s) | |
|---|---|---|
| Input size → | 1 | 1000 |
| SC + ComputeMSB | 0.000147 | 0.147 |
| BitDecompOPT | 0.0111 | 11.1 |

# 07

# Discussion

ComputeMSB vs. BitDecompOPT - local setting
Theoretical vs. actual communication
Combining approaches
Utilizing full bit-decomposition

# 07: ComputeMSB vs BitDecompOPT – Local setting

- SecureNN's ComputeMSB runs faster than BitDecompOPT in a local setting

- **Tradeoff**: Local computation and communication rounds

- BitDecompOPT introduces significant computational overhead.

SC + ComputeMSB

Less local computation

More communication rounds

BitDecompOPT

More local computation

Less communication rounds

# 07: Theoretical vs actual Communication

- Theoretical communication is based on the assumption that the implementation executes with no overhead.

- Requires intricate knowledge of low level programming and networking, which was out of scope for this project.

- Our results are achieved from an implementation that is far from the theoretical data transfer.

Theoretical Vs. Actual
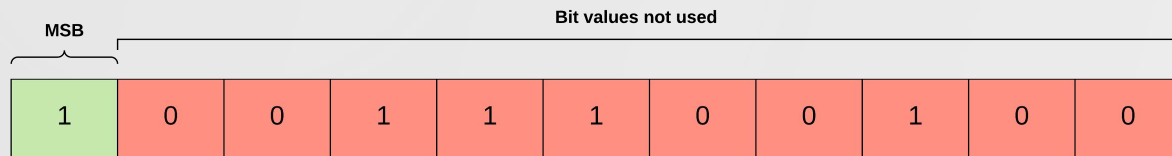
# 07: Combining approaches

- Implement BitDecomOPT in SecureNN to reduce DReLU rounds 9 $\Longrightarrow$ 7

- Both SecureNN and The Genome Paper claim same security.

- Problem with 3 party $\Longrightarrow$ 2 party

- SecureNN already uses a trusted initializer for common randomness, use it for beaver triplets to properly implement BitDecompOPT

# 07: Utilizing full bit-decomposition

- Bit-decomposition faster than computing MSB directly

- **But** needs to compute the full bit-decomposition of a value to obtain MSB

- Computes ƛ-1 bit values which are not used for anything else

- Further optimize by utilizing values computed during bit-decomposition

  - Other activation functions?

# THANKS