# Pencil: Private and Extensible Collaborative Machine Learning without the Non-Colluding Assumption

Xuanqi Liu, Zhuotao Liu, Qi Li, Ke Xu, Mingwei Xu
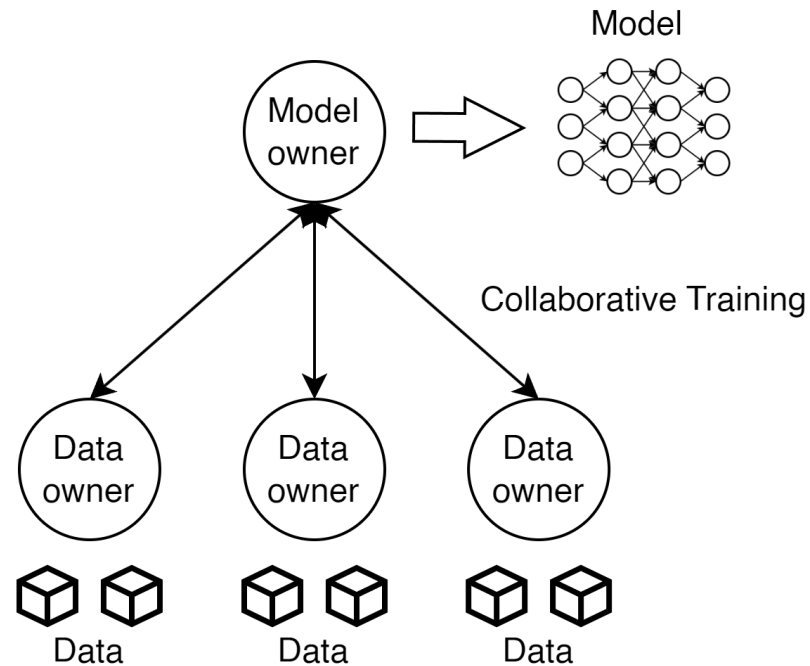
*Tsinghua University*

# Overview

- Key problem: Machine learning when model and data are separated

# Overview

- Key problem: Machine learning when model and data are separated
  - A Model Owner (MO) wishes to use the data of multiple Data Owners (DOes) to train a model.
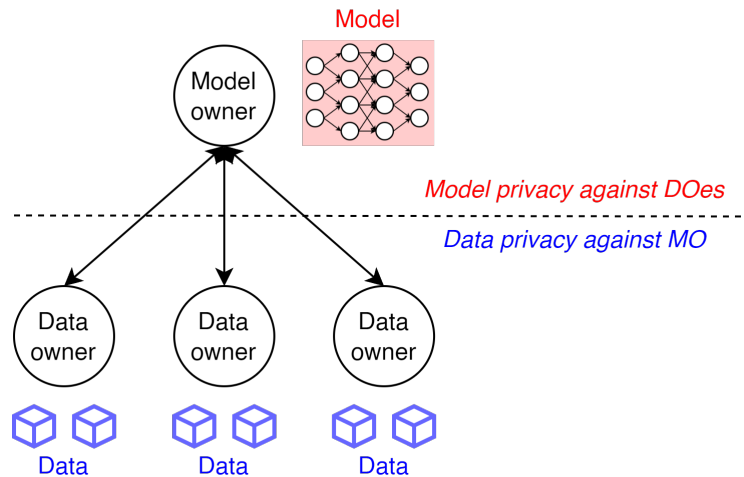
# Overview

- Requirements
    1) **Data privacy** – Fundamentally, raw data of DOes should not be leaked.

# Overview

- Requirements

  1) **Data privacy** – Fundamentally, raw data of DOes should not be leaked.

  2) **Model privacy** – Since MO may wish to finetune an existing model, the model parameters should not be leaked to other participants. Furthermore, the model should be able to be deployed independently by MO after training.

# Overview

- Requirements

  1) **Data privacy** – Fundamentally, raw data of DOes should not be leaked.

  2) **Model privacy** – Since MO may wish to finetune an existing model, the model parameters should not be leaked to other participants. Furthermore, the model should be able to be deployed independently by MO after training.

  3) **Extensibility** – The framework should scale to more DOes without significant increased cost
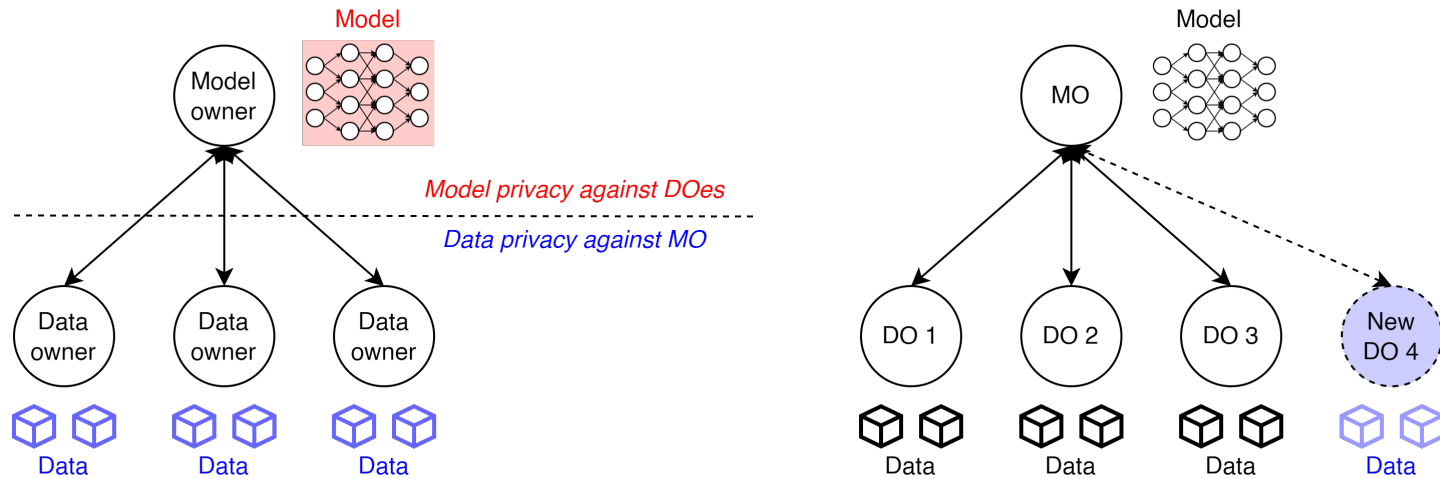
# Overview

- Requirements

    1) **Data privacy** – Fundamentally, raw data of DOes should not be leaked.

    2) **Model privacy** – Since MO may wish to finetune an existing model, the model parameters should not be leaked to other participants. Furthermore, the model should be able to be deployed independently by MO after training.

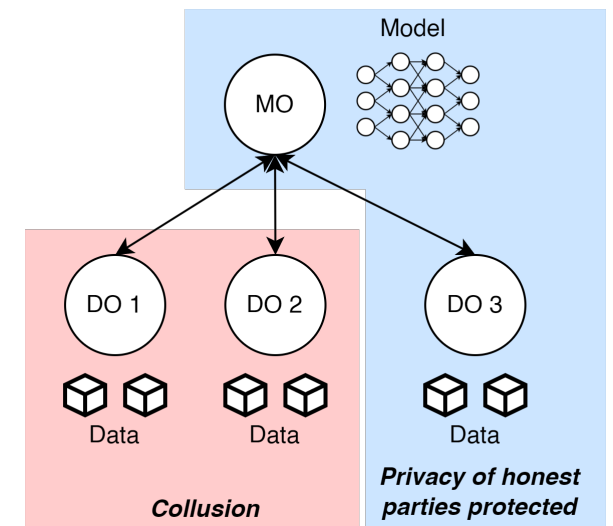    3) **Extensibility** – The framework should scale to more DOes without significant increased cost

    4) **Against collusion** – Colluding parties should not have advantage to break privacy of other honest parties

# Problems of prior works

- Federated Learning
  - In every iteration, Server (MO) distributes the model to Clients (DOes).
  - Clients train with local data and upload the updates for aggregation.

# Problems of prior works

- Federated Learning
  - In every iteration, Server (MO) distributes the model to Clients (DOes).
  - Clients train with local data and upload the updates for aggregation.



Pros
- Raw data never leave the client;
- Extensible (arbitrary joining and leaving)

# Problems of prior works

- Federated Learning
  - In every iteration, Server (MO) distributes the model to Clients (DOes).
  - Clients train with local data and upload the updates for aggregation.



Pros
- Raw data never leave the client;
- Extensible (arbitrary joining and leaving)

Cons
- The model privacy isn't protected at all!

# Problems of prior works

- Secure Multiparty Computation (MPC)
  - MO and DOes participate in n-party computation as servers.



Participants' roles are symmetric in MPC.
All model and data are secret-shared.

# Problems of prior works

- Secure Multiparty Computation (MPC)
  - MO and DOes participate in n-party computation as servers.



Participants' roles are symmetric in MPC.
All model and data are secret-shared.

Pros
- Model (and data) privacy guaranteed.

# Problems of prior works

- Secure Multiparty Computation (MPC)
  - MO and DOes participate in n-party computation as servers.



Participants' roles are symmetric in MPC.
All model and data are secret-shared.

Pros
- Model (and data) privacy guaranteed.

Cons
- Not extensible: introducing new DOes requires new protocol design.

# Problems of prior works

- Secure Multiparty Computation (MPC)
  - MO and DOes participate in n-party computation as servers.



Participants' roles are symmetric in MPC.
All model and data are secret-shared.

Pros
- Model (and data) privacy guaranteed.

Cons
- Not extensible: introducing new DOes requires new protocol design.
- Not secure against collusion
- Huge communication overhead.

# Problems of prior works

- Pure Homomorphic Encryption (HE) Approaches
  - DO uploads encrypted data for training; MO obtains an encrypted model.

# Problems of prior works

- Pure Homomorphic Encryption (HE) Approaches
  - DO uploads encrypted data for training; MO obtains an encrypted model.



The encrypted model could only be used by the DO assisting training.

# Problems of prior works

- Pure Homomorphic Encryption (HE) Approaches
  - DO uploads encrypted data for training; MO obtains an encrypted model.

MO ⟹ Encrypted model

↑ Encrypted data

DO

The encrypted model could only be used by the DO assisting training.

Pros
- Low communication overhead.

# Problems of prior works

- Pure Homomorphic Encryption (HE) Approaches
  - DO uploads encrypted data for training; MO obtains an encrypted model.



The encrypted model could only be used by the DO assisting training.

Pros
- Low communication overhead.

Cons
- Not extensible: only one DO!
- Heavy computation

# Our solution

- 2-party training: HE + MPC
  - Data and model privacy guaranteed.

**Single training step**



MO ←→ DO

Efficient HE/MPC training protocol

# Our solution

- 2-party training: HE + MPC
  - Data and model privacy guaranteed.
  - Model updates are given only to MO.

**Single training step**



Efficient HE/MPC
training protocol

MO ⟷ DO

# Our solution

- 2-party training: HE + MPC
  - Data and model privacy guaranteed.
  - Model updates are given only to MO.
    - i.e. the model is not shared nor encrypted w.r.t specific DO(es)

**Single training step**

Efficient HE/MPC
training protocol

MO ⟷ DO

# Our solution

- 2-party training: HE + MPC
  - Data and model privacy guaranteed.
  - Model updates are given only to MO.
    - i.e. the model is not shared nor encrypted w.r.t specific DO(es)
- Extensibility and collusion defence
  - MO trains with a different DO in each step.

**Single training step**

MO ⟷ DO

Efficient HE/MPC training protocol

**Multiple steps across DOes**

MO

Step 1    Step $i$    Step $n$

DO 1    ......    DO $i$    ......    DO $n$

# Our solution

- 2-party training: HE + MPC
  - Data and model privacy guaranteed.
  - Model updates are given only to MO.
    - i.e. the model is not shared nor encrypted w.r.t specific DO(es)

- Extensibility and collusion defence
  - MO trains with a different DO in each step.
  - Since no privacy leaks in 2-party, collusion could not break the privacy of any party.

**Single training step**



Efficient HE/MPC training protocol

MO ←→ DO

**Multiple steps across DOes**



MO

Step 1    Step i    Step n

DO 1    ......    DO i    ......    DO n

# Pencil training overview

- Secret shares throughout FP and BP

# Pencil training overview

- Secret shares throughout FP and BP
  - $l$-layer sequential model
  - $X_i = f_i(X_{i-1})$

# Pencil training overview

- Secret shares throughout FP and BP
  - $l$-layer sequential model
  - $X_i = f_i(X_{i-1})$
  - In FP, the two parties keep $\langle X_i \rangle$ in secret shares
    $$\langle \mathbf{X}_i \rangle_0 + \langle \mathbf{X}_i \rangle_1 = f_i(\langle \mathbf{X}_{i-1} \rangle_0 + \langle \mathbf{X}_{i-1} \rangle_1).$$

# Pencil training overview

- Secret shares throughout FP and BP
  - $l$-layer sequential model
  - $X_i = f_i(X_{i-1})$
  - In FP, the two parties keep $\langle X_i \rangle$ in secret shares

  $$\langle \mathbf{X}_i \rangle_0 + \langle \mathbf{X}_i \rangle_1 = f_i(\langle \mathbf{X}_{i-1} \rangle_0 + \langle \mathbf{X}_{i-1} \rangle_1).$$

  - In BP, similarly $\langle \nabla_{X_i} \rangle$ is shared;

  $$\langle \nabla_{\mathbf{X}_{i-1}} \rangle_0 + \langle \nabla_{\mathbf{X}_{i-1}} \rangle_1 = (\langle \nabla_{\mathbf{x}_i} \rangle_0 + \langle \nabla_{\mathbf{x}_i} \rangle_1) \odot_x \frac{\partial f_i(\mathbf{X}_{i-1})}{\partial \mathbf{X}_{i-1}}$$

# Pencil training overview

- Secret shares throughout FP and BP
  - $l$-layer sequential model
  - $X_i = f_i(X_{i-1})$
  - In FP, the two parties keep $\langle X_i \rangle$ in secret shares

$$\langle \mathbf{X}_i \rangle_0 + \langle \mathbf{X}_i \rangle_1 = f_i(\langle \mathbf{X}_{i-1} \rangle_0 + \langle \mathbf{X}_{i-1} \rangle_1).$$

  - In BP, similarly $\langle \nabla_{X_i} \rangle$ is shared;

$$\langle \nabla_{\mathbf{X}_{i-1}} \rangle_0 + \langle \nabla_{\mathbf{X}_{i-1}} \rangle_1 = (\langle \nabla_{\mathbf{X}_i} \rangle_0 + \langle \nabla_{\mathbf{X}_i} \rangle_1) \odot_x \frac{\partial f_i(\mathbf{X}_{i-1})}{\partial \mathbf{X}_{i-1}}$$

  - ... but the weight gradients are given to the MO.

$$\nabla_{\mathbf{b}_i} = \nabla_{\mathbf{X}_i} \odot_b \frac{\partial f_i(\mathbf{X}_{i-1}; \mathbf{W}_i, \mathbf{b}_i)}{\partial \mathbf{b}_i}$$

$$\nabla_{\mathbf{w}_i} = \nabla_{\mathbf{X}_i} \odot_W \frac{\partial f_i(\mathbf{X}_{i-1}; \mathbf{W}_i, \mathbf{b}_i)}{\partial \mathbf{W}_i}$$



MO · DO

Forward propagation

$X_{i-1}^S$ · $X_{i-1}^C$

$$\mathbf{X}_i = f_i(\mathbf{X})$$

$X_i^S$ · $X_i^C$

Backward propagation

$X_{i-1}^S, \nabla_{X_i}^S$ · $X_{i-1}^C, \nabla_{X_i}^C$

$$\nabla_{\mathbf{X}_{i-1}} = \nabla_{\mathbf{X}_i} \odot_x \frac{\partial f_i(\mathbf{X}_{i-1})}{\partial \mathbf{X}_{i-1}}$$

For linear layers:

$$\nabla_{\mathbf{b}_i} = \nabla_{\mathbf{X}_i} \odot_b \frac{\partial f_i(\mathbf{X}_{i-1})}{\partial \mathbf{b}_i}$$

$$\nabla_{\mathbf{W}_i} = \nabla_{\mathbf{X}_i} \odot_W \frac{\partial f_i(\mathbf{X}_{i-1})}{\partial \mathbf{W}_i}$$

$\nabla_{X_{i-1}}^S$ · $\nabla_{X_{i-1}}^C$

$\nabla_{b_i}, \nabla_{w_i}$

# Training: Linear layers = HE

- Forward propagation
  - 2-round protocol with HE



**MO**

Inputs
$\mathbf{W}, \mathbf{b}, \langle \mathbf{X} \rangle_0$

Samples $\mathbf{s}$

Outputs
$\langle \mathbf{Y} \rangle_0$
$= \mathbf{s} + \mathbf{b}$

$[\![\langle \mathbf{X} \rangle_1]\!]$

$[\![\langle \mathbf{Y}_1 \rangle]\!]$
$= \mathbf{W} \circ (\langle \mathbf{X} \rangle_0 + [\![\langle \mathbf{X} \rangle_1]\!]) - \mathbf{s}$

**DO**

Inputs
$\langle \mathbf{X} \rangle_1$

Outputs
$\langle \mathbf{Y} \rangle_1$

# Training: Linear layers = HE

- Forward propagation
  - 2-round protocol with HE
  - As a general solution, this algorithm does not specify how $\mathbf{W} \circ [\![\mathbf{X}]\!]$ is evaluated.
  - Our implementation uses batched polynomial encoding, but other methods (e.g. Gazelle's encoding) could be used.

**MO**

Inputs
$\mathbf{W}, \mathbf{b}, \langle\mathbf{X}\rangle_0$

$[\![\langle\mathbf{X}\rangle_1]\!]$

⟵

Samples $\mathbf{s}$

$[\![\langle\mathbf{Y}_1\rangle]\!]$
$= \mathbf{W} \circ (\langle\mathbf{X}\rangle_0 + [\![\langle\mathbf{X}\rangle_1]\!]) - \mathbf{s}$

⟶

Outputs
$\langle\mathbf{Y}\rangle_0$
$= \mathbf{s} + \mathbf{b}$

**DO**

Inputs
$\langle\mathbf{X}\rangle_1$

Outputs
$\langle\mathbf{Y}\rangle_1$

# Training: Linear layers = HE

- Backpropagation
  - Gradient of x, shared: $\nabla_x = W \odot_W \nabla_y$
    - E.g. For FC, $\langle Y \rangle = W \langle X \rangle, \langle \nabla X \rangle = \langle \nabla Y \rangle \cdot W^T$

# Training: Linear layers = HE

- Backpropagation
  - Gradient of x, shared: $\nabla_x = W \odot_W \nabla_y$
    - E.g. For FC, $\langle Y \rangle = W \langle X \rangle, \langle \nabla X \rangle = \langle \nabla Y \rangle \cdot W^T$
  - Gradient of b, given to MO: $\nabla_b = \nabla_y \odot_b \dfrac{\partial f}{\partial b}.$
    - For FC and Conv, simply sum up $\nabla y$ and reveal to MO.

# Training: Linear layers = HE

- Backpropagation
  - Gradient of x, shared: $\nabla_x = W \odot_W \nabla_y$
    - E.g. For FC, $\langle Y \rangle = W \langle X \rangle$, $\langle \nabla X \rangle = \langle \nabla Y \rangle \cdot W^T$
  - Gradient of b, given to MO: $\nabla_b = \nabla_y \odot_b \frac{\partial f}{\partial b}$.
    - For FC and Conv, simply sum up $\nabla y$ and reveal to MO.
  - Gradient of W, given to MO

# Training: Linear layers = HE

- For gradient of W, it is the product of two secret-shared values

$$\nabla_{\mathbf{W}} = \nabla_{\mathbf{Y}} \odot \frac{\partial f(\mathbf{X}; \mathbf{W}, \mathbf{b})}{\partial \mathbf{W}} = \nabla_{\mathbf{Y}} \odot \mathbf{X}$$

# Training: Linear layers = HE

- For gradient of W, it is the product of two secret-shared values

$$\nabla_{\mathbf{W}} = \nabla_{\mathbf{Y}} \odot \frac{\partial f(\mathbf{X}; \mathbf{W}, \mathbf{b})}{\partial \mathbf{W}} = \nabla_{\mathbf{Y}} \odot \mathbf{X}$$

- Solution: HE for cross terms

(no need for cipher-cipher multiplication)

$$[\![\nabla_{\mathbf{W}}]\!] = (\langle\nabla_{\mathbf{Y}}\rangle_0 + [\![\langle\nabla_{\mathbf{Y}}\rangle_1]\!]) \odot (\langle\mathbf{X}\rangle_0 + [\![\langle\mathbf{X}\rangle_1]\!]).$$

# Training: Linear layers = HE

- For gradient of W, it is the product of two secret-shared values

$$\nabla_{\mathbf{W}} = \nabla_{\mathbf{Y}} \odot \frac{\partial f(\mathbf{X}; \mathbf{W}, \mathbf{b})}{\partial \mathbf{W}} = \nabla_{\mathbf{Y}} \odot \mathbf{X}$$

- Solution: HE for cross terms
(no need for cipher-cipher multiplication)

$$[\![\nabla_{\mathbf{W}}]\!] = (\langle \nabla_{\mathbf{Y}} \rangle_0 + [\![\langle \nabla_{\mathbf{Y}} \rangle_1]\!]) \odot (\langle \mathbf{X} \rangle_0 + [\![\langle \mathbf{X} \rangle_1]\!]).$$

**MO**

Inputs
$\langle \mathbf{X} \rangle_0, \langle \nabla_{\mathbf{Y}} \rangle_0$

$\longleftarrow \quad [\![\langle \mathbf{X} \rangle_1]\!], [\![\langle \nabla_{\mathbf{Y}} \rangle_1]\!]$

Samples
random $\mathbf{s}$

(Masked cross term)
$[\![\tilde{\nabla}_{\mathbf{W}}^{\mathsf{cross}}]\!] = \langle \nabla_{\mathbf{Y}} \rangle_0 \odot [\![\langle \mathbf{X} \rangle_1]\!]$
$+ [\![\langle \nabla_{\mathbf{Y}} \rangle_1]\!] \odot \langle \mathbf{X} \rangle_0 - \mathbf{s}$

$\longrightarrow$

$\tilde{\nabla}_{\mathbf{W}} = \tilde{\nabla}_{\mathbf{W}}^{\mathsf{cross}} + \mathbf{e}$
$+ \langle \nabla_{\mathbf{Y}} \rangle_1 \odot \langle \mathbf{X} \rangle_1$

Outputs
$\nabla_{\mathbf{W}} = \tilde{\nabla}_{\mathbf{W}} +$
$\langle \nabla_{\mathbf{Y}} \rangle_0 \odot \langle \mathbf{X} \rangle_0$

$\longleftarrow$

**DO**

Inputs
$\langle \mathbf{X} \rangle_1, \langle \nabla_{\mathbf{Y}} \rangle_1$

Samples
DP noise $\mathbf{e}$

# Training: Non-linear layers = MPC

- Build with Two-party MPC

# Training: Non-linear layers = MPC

- Build with Two-party MPC
- Example: $\mathrm{ReLU}(x) = \mathrm{DReLU}(x) \cdot x$

# Training: Non-linear layers = MPC

- Build with Two-party MPC
- Example: $\text{ReLU}(x) = \text{DReLU}(x) \cdot x$
  - DReLU => secure comparison protocol
  - Boolean-arithmetic multiplication => OT-based multiplexing

# Optimzing

- Substantial part of computation lies in HE linear operation evaluation

# Optimzing

- Substantial part of computation lies in HE linear operation evaluation
  - $W \circ [\![X]\!]$ in FP, results shared
  - $W \odot_x [\![\nabla_Y]\!]$ in BP, results shared
  - $\langle \nabla_Y \rangle_0 \odot [\![\langle X \rangle_1]\!] + [\![\langle \nabla_Y \rangle_1]\!] \odot \langle X \rangle_0$ in BP, results shared

# Optimzing

- Substantial part of computation lies in HE linear operation evaluation
  - $W \circ [\![ X ]\!]$ in FP, results shared
  - $W \odot_x [\![ \nabla_Y ]\!]$ in BP, results shared
  - $\langle \nabla_Y \rangle_0 \odot [\![ \langle X \rangle_1 ]\!] + [\![ \langle \nabla_Y \rangle_1 ]\!] \odot \langle X \rangle_0$ in BP, results shared
- **Generalization**: Is there a way to accelerate online evaluation of general operator $u \circ v$, each party holding one operand?

# Optimzing

- Substantial part of computation lies in HE linear operation evaluation
  - $W \circ [\![X]\!]$ in FP, results shared
  - $W \odot_x [\![\nabla_Y]\!]$ in BP, results shared
  - $\langle \nabla_Y \rangle_0 \odot [\![\langle X \rangle_1]\!] + [\![\langle \nabla_Y \rangle_1]\!] \odot \langle X \rangle_0$ in BP, results shared
- **Generalization**: Is there a way to accelerate online evaluation of general operator $u \circ v$, each party holding one operand?
- First, let's consider a fixed $u$ and variable $v$'s.

# Traditional preprocessing: Beaver triples



**Alice**
Inputs fixed $\mathbf{u}$

Obtains $\langle \mathbf{w} \rangle_0$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_0 = \langle \mathbf{w} \rangle_0$
$+ \mathbf{u} \circ (\mathbf{v} - \mathbf{v}')$

**Bob**
Samples
random $\mathbf{v}'$

Obtains $\langle \mathbf{w} \rangle_1$

Inputs $\mathbf{v}$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_1$
$= \langle \mathbf{w} \rangle_1$

Preprocess (heavy)

HE/MPC protocol

Return shares of $\mathbf{w} = \mathbf{u} \circ \mathbf{v}'$

Online (lightweight)

Masked input $\mathbf{v}' - \mathbf{v}$

# Traditional preprocessing: Beaver triples



- Preprocessed shares cannot be reused, or info is leaked.

### Alice
Inputs fixed $\mathbf{u}$

Preprocess (heavy)

HE/MPC protocol

Obtains $\langle \mathbf{w} \rangle_0$    Return shares of $\mathbf{w} = \mathbf{u} \circ \mathbf{v}'$

### Bob

Samples random $\mathbf{v}'$

Obtains $\langle \mathbf{w} \rangle_1$

Online (lightweight)    Inputs $\mathbf{v}$

Masked input $\mathbf{v}' - \mathbf{v}$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_0 = \langle \mathbf{w} \rangle_0$
$+ \mathbf{u} \circ (\mathbf{v} - \mathbf{v}')$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_1$
$= \langle \mathbf{w} \rangle_1$

# Traditional preprocessing: Beaver triples



- Preprocessed shares cannot be reused, or info is leaked.
  - If, for $v_0$ and $v_1$, the same sharing of $w = u' \circ v'$ was used, Bob would send $v' - v_0$ and $v' - v_1$ to Alice, so Alice would obtain the difference $v_0 - v_1$ (a direct linear combination of 2 input values)

# Traditional preprocessing: Beaver triples



- Preprocessed shares cannot be reused, or info is leaked.
  - If, for $v_0$ and $v_1$, the same sharing of $w = u' \circ v'$ was used, Bob would send $v' - v_0$ and $v' - v_1$ to Alice, so Alice would obtain the difference $v_0 - v_1$ (a direct linear combination of 2 input values)
- Total communication is not reduced, while total computation is even increased.

# Novel approach: Multiple masks

- Solution: use multiple masks to increase the revealed linear combination complexity.

# Novel approach: Multiple masks

- Solution: use multiple masks to increase the revealed linear combination complexity.
  - Bob samples $m$ random $v_i'$'s to conduct preprocessing.



**Alice**

Inputs fixed $\mathbf{u}$

Preprocess (heavy)

HE/MPC protocol

Obtains $\langle \mathbf{w}_i \rangle_0$    Return shares of $\mathbf{w}_i = \mathbf{u} \circ \mathbf{v}_i'$

Online (lightweight)

$k_i, \tilde{\mathbf{v}} = \mathbf{v} - \sum_i k_i \mathbf{v}_i'$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_0 = \mathbf{u} \circ \tilde{\mathbf{v}}$
$+ \sum_i k_i \langle \mathbf{w}_i \rangle_0$

**Bob**

For $i \in [m]$ samples random $\mathbf{v}_i'$

Obtains $\langle \mathbf{w}_i \rangle_1$

Inputs $\mathbf{v}$
Samples $k_i$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_1$
$= \sum_i k_i \langle \mathbf{w}_i \rangle_1$

# Novel approach: Multiple masks

- Solution: use multiple masks to increase the revealed linear combination complexity.
  - Bob samples $m$ random $v_i'$'s to conduct preprocessing.
  - Reuse the shares $\langle w_i \rangle$ for multiple online executions



**Alice**
Inputs fixed $\mathbf{u}$

Preprocess (heavy)

For $i \in [m]$ samples random $\mathbf{v}_i'$

HE/MPC protocol

Obtains $\langle \mathbf{w}_i \rangle_0$    Return shares of $\mathbf{w}_i = \mathbf{u} \circ \mathbf{v}_i'$    Obtains $\langle \mathbf{w}_i \rangle_1$

**Bob**

Online (lightweight)

Inputs $\mathbf{v}$
Samples $k_i$

$k_i, \tilde{\mathbf{v}} = \mathbf{v} - \sum_i k_i \mathbf{v}_i'$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_0 = \mathbf{u} \circ \tilde{\mathbf{v}} + \sum_i k_i \langle \mathbf{w}_i \rangle_0$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_1 = \sum_i k_i \langle \mathbf{w}_i \rangle_1$

# Multiple masks

- Extending to variable $u$'s
  - Similarly, Alice samples multiple masks $u_i'$

# Multiple masks

- Extending to variable $u$'s
  - Similarly, Alice samples multiple masks $u_i'$

| Alice | | Bob |
|---|---|---|
| | | |

**Alice**

For $i \in [m]$ samples random $\mathbf{u}_i'$

Obtains $\langle \mathbf{w}_{ij} \rangle_0$

Inputs $\mathbf{u}$
Samples $k_i$

Computes $\langle \mathbf{t}_j \rangle_0$
$= \langle \mathbf{u} \circ \mathbf{v}_j' \rangle_0$
$= \sum_i k_i \langle \mathbf{w}_{ij} \rangle_0$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_0 = \mathbf{u} \circ \tilde{\mathbf{v}}$
$+ \sum_j l_j \langle \mathbf{t}_j \rangle_0$

**Preprocess (heavy)**

HE/MPC protocol

Return shares of $\mathbf{w}_{ij} = \mathbf{u}_i' \circ \mathbf{v}_j'$

**Online (lightweight)**

$k_i, \tilde{\mathbf{u}} = \mathbf{u} - \sum_i k_i \mathbf{u}_i'$

$l_j, \tilde{\mathbf{v}} = \mathbf{v} - \sum_j l_j \mathbf{v}_j'$

**Bob**

For $j \in [m]$ samples random $\mathbf{v}_j'$

Obtains $\langle \mathbf{w}_{ij} \rangle_1$

Inputs $\mathbf{v}$
Samples $l_j$

Computes $\langle \mathbf{t}_j \rangle_1 =$
$\langle \mathbf{u} \circ \mathbf{v}_j' \rangle_1 = \tilde{\mathbf{u}} \circ \mathbf{v}_j'$
$+ \sum_i k_i \langle \mathbf{w}_{ij} \rangle_1$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_1$
$= \sum_j l_j \langle \mathbf{t}_j \rangle_1$

# Multiple masks

- Extending to variable $u$'s
  - Similarly, Alice samples multiple masks $u_i'$

- If the online phase is executed $T$ times:
  - Traditional: $T$ times HE/MPC evaluation of $\circ$
  - Ours: $m^2$ times of HE/MPC evaluation, regardless of $T$

# Multiple masks

- Extending to variable $u$'s
  - Similarly, Alice samples multiple masks $u_i'$

- If the online phase is executed $T$ times:
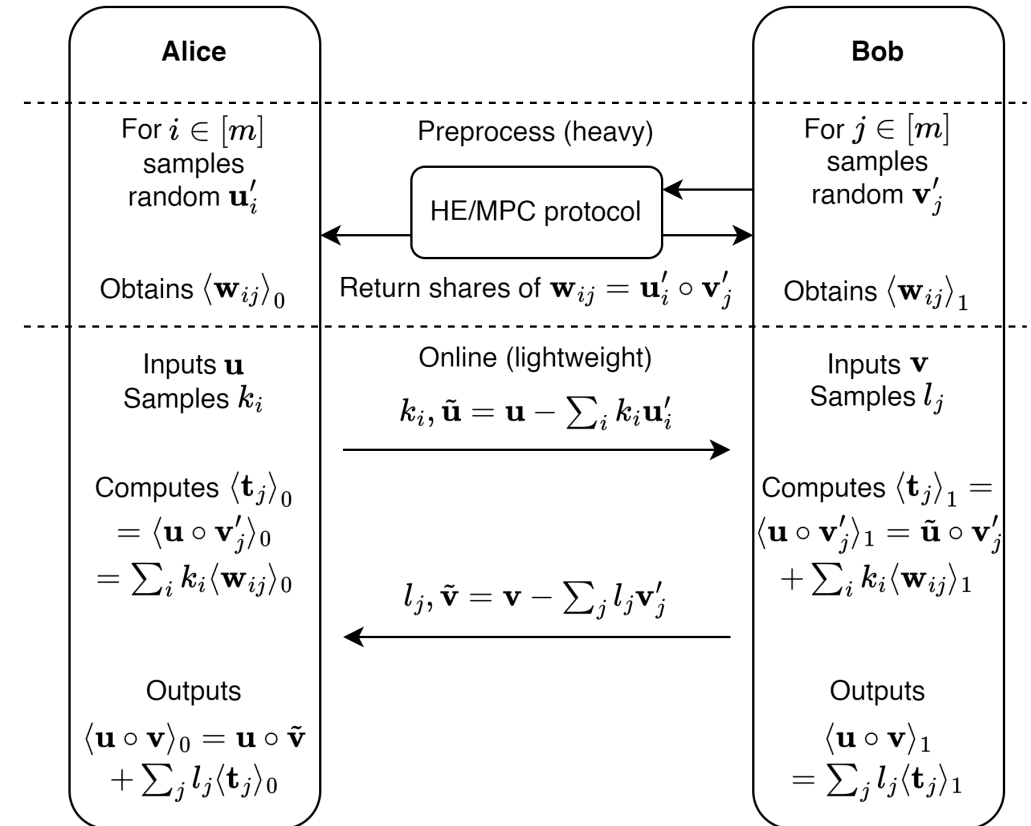  - Traditional: $T$ times HE/MPC evaluation of $\circ$
  - Ours: $m^2$ times of HE/MPC evaluation, regardless of $T$

- Security analysis shows
  - To eliminate masks, an attacker would require at least $m+1$ equations
  - Complexity of breaking one $u$ or $v$ is $O(2^{fm})$, $f$ being the fixed-point precision



Alice

For $i \in [m]$ samples random $\mathbf{u}_i'$

Preprocess (heavy)

HE/MPC protocol

Bob

For $j \in [m]$ samples random $\mathbf{v}_j'$

Obtains $\langle \mathbf{w}_{ij} \rangle_0$

Return shares of $\mathbf{w}_{ij} = \mathbf{u}_i' \circ \mathbf{v}_j'$

Obtains $\langle \mathbf{w}_{ij} \rangle_1$

Inputs $\mathbf{u}$
Samples $k_i$

Online (lightweight)

$k_i, \tilde{\mathbf{u}} = \mathbf{u} - \sum_i k_i \mathbf{u}_i'$

Inputs $\mathbf{v}$
Samples $l_j$

Computes $\langle \mathbf{t}_j \rangle_0$
$= \langle \mathbf{u} \circ \mathbf{v}_j' \rangle_0$
$= \sum_i k_i \langle \mathbf{w}_{ij} \rangle_0$

$l_j, \tilde{\mathbf{v}} = \mathbf{v} - \sum_j l_j \mathbf{v}_j'$

Computes $\langle \mathbf{t}_j \rangle_1 =$
$\langle \mathbf{u} \circ \mathbf{v}_j' \rangle_1 = \tilde{\mathbf{u}} \circ \mathbf{v}_j'$
$+ \sum_i k_i \langle \mathbf{w}_{ij} \rangle_1$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_0 = \mathbf{u} \circ \tilde{\mathbf{v}}$
$+ \sum_j l_j \langle \mathbf{t}_j \rangle_0$

Outputs
$\langle \mathbf{u} \circ \mathbf{v} \rangle_1$
$= \sum_j l_j \langle \mathbf{t}_j \rangle_1$

# Evaluation: Training costs

Since the preprocessing technique introduces different security properties, we denote Pencil with or without it as Pencil and Pencil+ respectively.

# Evaluation: Training costs

Since the preprocessing technique introduces different security properties, we denote Pencil with or without it as Pencil and Pencil+ respectively.

| Scenario | Task | Model | Pencil | | | Pencil+ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | Online | | | Preprocessing | | Online | | |
| | | | $TP_{LAN}$ | $TP_{WAN}$ | C | $T_{prep}$ | $C_{prep}$ | $TP_{LAN}$ | $TP_{WAN}$ | C |
| Train from scratch | MNIST | NN1 | $9.73 \times 10^4$ | $5.12 \times 10^4$ | 1.66 | 0.02 | 3.35 | $26.52 \times 10^4$ | $19.87 \times 10^4$ | 0.23 |
| | MNIST | NN2 | $7.70 \times 10^4$ | $4.43 \times 10^4$ | 1.71 | 0.02 | 4.13 | $13.72 \times 10^4$ | $10.75 \times 10^4$ | 0.36 |
| | CIFAR10 | NN3 | $2.58 \times 10^4$ | $1.62 \times 10^4$ | 4.11 | 0.05 | 10.26 | $2.90 \times 10^4$ | $1.98 \times 10^4$ | 2.86 |
| | CIFAR10 | NN4 | $0.18 \times 10^4$ | $0.12 \times 10^4$ | 44.89 | 0.70 | 83.12 | $0.22 \times 10^4$ | $0.15 \times 10^4$ | 34.90 |
| Transfer learning | CIFAR10 | NN5 | $0.52 \times 10^4$ | $0.39 \times 10^4$ | 11.33 | 0.91 | 46.00 | $1.55 \times 10^4$ | $1.24 \times 10^4$ | 2.90 |
| | CIFAR10 | NN6 | $1.83 \times 10^4$ | $1.17 \times 10^4$ | 5.48 | 0.30 | 15.96 | $8.05 \times 10^4$ | $5.89 \times 10^4$ | 0.82 |

TABLE III: Training costs for different ML tasks. For the online phase, TP stands for the throughput (images/hour) of the training system, and subscript LAN, WAN indicate the network settings; C stands for the online communication (MB) per image. For Pencil+, we also report the time ($T_{prep}$, hours) and communication ($C_{prep}$, GB) of preprocessing. Note that the preprocessing overhead is one-time overhead.

With Pencil+ and transfer learning, a model for CIFAR10 classification could be trained within 6.5 hours (10 epochs)

# Evaluation: Training costs

| | Throughput ($10^4$ img/h) | | | | Comm. (MB/img) | | |
| Model | [2] | [12] | P | P$^+$ | [12] | P | P$^+$ |
|---|---|---|---|---|---|---|---|
| $2 \times 128$FC | 0.7 | 0.11 | 9.7 | 29.3 | 552 | 1.7 | 0.2 |
| $3 \times 128$FC | 0.6 | 0.10 | 8.1 | 18.9 | 658 | 2.2 | 0.3 |
| $2 \times 512$FC | 0.2 | 0.03 | 2.6 | 13.2 | 3470 | 5.2 | 0.8 |

TABLE VII: Performance comparison with QUOTIENT [2] and Semi2k [12] in the 2 party setting. The models are represented as $n \times m$FC, as used by [2]. P represents Pencil and P$^+$ represents Pencil$^+$.

- Comparison with previous 2PC works shows improvements of up to 2 orders of magnitude.

# Evaluation: Training costs

| | Throughput ($10^4$ img/h) | | | | Comm. (MB/img) | | |
|---|---|---|---|---|---|---|---|
| Model | [2] | [12] | P | P+ | [12] | P | P+ |
| $2 \times 128$FC | 0.7 | 0.11 | 9.7 | 29.3 | 552 | 1.7 | 0.2 |
| $3 \times 128$FC | 0.6 | 0.10 | 8.1 | 18.9 | 658 | 2.2 | 0.3 |
| $2 \times 512$FC | 0.2 | 0.03 | 2.6 | 13.2 | 3470 | 5.2 | 0.8 |

TABLE VII: Performance comparison with QUOTIENT [2] and Semi2k [12] in the 2 party setting. The models are represented as $n \times m$FC, as used by [2]. P represents Pencil and P+ represents Pencil+.

| | Throughput ($10^3$ img/h) | | | Comm. (per img) | | |
|---|---|---|---|---|---|---|
| Model | [12] | Pencil | Pencil+ | [12] | Pencil | Pencil+ |
| 2 parties | 1.11 | 97 | 293 | 0.55GB | 1.7MB | 0.2MB |
| 3 parties | 0.61 | 97 | 293 | 2.58GB | 1.7MB | 0.2MB |
| 4 parties | 0.41 | 97 | 293 | 6.06GB | 1.7MB | 0.2MB |
| 5 parties | 0.07 | 97 | 293 | 57.69GB | 1.7MB | 0.2MB |

TABLE VIII: Performance comparison with Semi2k [12] in multiple party setting.

- Comparison with previous 2PC works shows improvements of up to 2 orders of magnitude.

- Unlike previous general n-PC frameworks, extending to multiple DOes does not introduce extra overhead for Pencil.

# Thank you for listening!