#### BumbleBee: Secure Two-party Inference Framework for Large Transformers

Wen-jie Lu





## Background

Risk of privacy leak during the AI services



#### Event 1

#### Samsung Bans ChatGPT After Engineers Use it to Fix Proprietary Code

Employees who continue to use ChatGPT will face 'disciplinary action up to and including termination of employment,' according to a staff memo.

Event 2

## Yep, human workers are listening to recordings from Google Assistant, too

## Background

Risk of privacy leak during the AI services



#### Event 1

#### Samsung Bans ChatGPT After Engineers Use it to Fix Proprietary Code

Employees who continue to use ChatGPT will face 'disciplinary action up to and including termination of employment,' according to a staff memo.

Event 2

## Yep, human workers are listening to recordings from Google Assistant, too

## Secure Two-party Computation (2PC)

- Two parties with private inputs x and y
- Compute a joint function of their inputs while preserving



(ra, rb, rc) => rc = ra \* rb or rc = ra ∧ rb

## Variants: Distributing Correlated Randomness









### **Related Work**

	Secret Sharing	Variant	Underlying 2PC lib	Transformers	Fine tune?
MPCFormer (ICRL23)	Additive	TFP	CrypTen	BERT (manually constructed)	Demanded
SHAFT (NDSS25)	Additive	TFP	CrypTen	Models written in PyTorch, e.g., BERT, GPT2, ViT	Prefer
SIGMA (PoPETs24)	Functional	TTP	EzPC (Easy Secure Multiparty Computation)	Models written in ONNX (subsets)	Not necessary
Iron (NeurIPS22)	Additive	2PC	EzPC	BERT (manually constructed)	Prefer
BOLT (Oakland24)	Additive	2PC	EzPC	BERT (manually constructed)	Demanded
BumbleBee	Additive	2PC	SecretFlow	Models written in JAX	Not necessary

### **Related Work**

	Secret Sharing	Variant	Underlying 2PC lib	Transformers	Fine tune?
MPCFormer (ICRL23)	Additive	TFP	CrypTen	BERT (manually constructed)	Demanded
SHAFT (NDSS25)	Additive	TFP	CrypTen	Models written in PyTorch, e.g., BERT, GPT2, ViT	Prefer
SIGMA (PoPETs24)	Functional	TTP	EzPC (Easy Secure Multiparty Computation)	Models written in ONNX (subsets)	Not necessary
Iron (NeurIPS22)	Additive	2PC	EzPC	BERT (manually constructed)	Prefer
BOLT (Oakland24)	Additive	2PC	EzPC	BERT (manually constructed)	Demanded
BumbleBee	Additive	2PC	SecretFlow	Models written in JAX	Not necessary

### **Related Work**

	Secret Sharing	Variant	Underlying 2PC lib	Transformers	Fine tune?
MPCFormer (ICRL23)	Additive	TFP	CrypTen	BERT (manually constructed)	Demanded
SHAFT (NDSS25)	Additive	TFP	CrypTen	Models written in PyTorch, e.g., BERT, GPT2, ViT	Prefer
SIGMA (PoPETs24)	Functional	TTP	EzPC (Easy Secure Multiparty Computation)	Models written in ONNX (subsets)	Not necessary
Iron (NeurIPS22)	Additive	2PC	EzPC	BERT (manually constructed)	Prefer
BOLT (Oakland24)	Additive	2PC	EzPC	BERT (manually constructed)	Demanded
BumbleBee	Additive	2PC	SecretFlow	Models written in JAX	Not necessary

## **Observations & Objectives**

- O1: Tremendous Communication
- ViT: 260 GB per image [1]
- Bert-base: 80 GB per 128tok
   O2: \$Computation < \$Communication</li>
- 0.09 cent per GB
- 0.025 cent per vCPU
- Minimize the Communication
- Not considering the offline/online split
- Without increasing the computation by too much





## Difficulties

- Evaluating nonlinear activations (e.g., GeLU and Softmax) in 2PC.
  - Many attempts turns to 2PC friendly alternatives. 😕 Accuracy.
- ML-Crypto Co-engineering.
  - MLer writes 2PC codes vs. Cryptographer writes ML codes.
- Large-scale matrix multiplication in 2PC.

#### Contributions

- Lessons for evaluating nonlinear function in 2PC.
- End-2-end code baes for private transformer inference.
  - 100% reusing JAX codes from HuggingFaces, no model finetuning.
- Efficient point-wise mul (aka OLE) and Matmul in  $\mathbb{Z}_{2^k}$  from HE
  - For example: 128x768 matrix in 0.6s and 5MB (10% of the prev)
  - A concurrent work by Jiaxing He et al. in CCS24 (Rhombus)

https://github.com/AntCPLab/OpenBumbleBee

# Three lessons for the private transformers inference

## Lesson 1: Reuse the ML codes via 2PC compilers

#### Or: Let the MLers write the ML codes

- Large transformer models could be "hard" to write from scratch.
  - Iron, MPCFormer, BOLT only run on one model by hand crafting
  - Handling the existing model weight files is also tedious
- A better way, plugin a 2PC backend to the current ML codes (eg PyTorch, TensorFlow)
  - CrypTen (SHAFT): Python-level overriding the PyTorch API
  - EzPC (SIGMA): Uses SeeDot [1] to compile TensorFlow/ONNX to a C++ backend
  - SecretFlow (BumbleBee): Multi-level IR (MLIR). Compile JAX/PyTorch to a C++ backend

## Lesson 1: Reuse the ML codes via 2PC compilers

#### Or: Let the MLers write the ML codes

- Large transformer models could be "hard" to write from scratch.
  - Iron, MPCFormer, BOLT only run on one model by hand crafting
  - Handling the existing model weight files is also tedious
- A better way, plugin a 2PC backend to the current ML codes (eg PyTorch, TensorFlow)
  - CrypTen (SHAFT): Python-level overriding the PyTorch API
  - EzPC (SIGMA): Uses SeeDot [1] to compile TensorFlow/ONNX to a C++ backend
  - SecretFlow (BumbleBee): Multi-level IR (MLIR). Compile JAX/PyTorch to a C++ backend

[1] https://github.com/mpc-msri/EzPC/tree/master/Athos#compiling-a-tensorflow-model

Prefer a middle-level IR than the eager execution

#### JAX's snippet for Softmax

- 1. **def** \_softmax(x, axis = -1):
- 2. x\_max = jnp.max(x, axis,
- $\hookrightarrow$  keepdims=**True**)
- 3. unnormalized = jnp.exp(x x\_max)
- 4. sum\_exp = jnp.sum(unnormalized, axis,
- $\leftrightarrow$  keepdims=**True**)
- 5. result = unnormalized / sum\_exp
- 6. return result

#### SHAFT's snippet for Softmax

```
if method == "reciprocal":
```

maximum\_value = self.max(dim, keepdim=True)[0]

```
logits = self - maximum_value
```

```
numerator = logits.exp()
```

```
with cfg.temp_override({"functions.reciprocal_all_pos": True}):
```

```
inv_denominator = numerator.sum(dim, keepdim=True).reciprocal()
```

#### return numerator \* inv\_denominator

[1] https://github.com/jax-ml/jax/blob/ed952c8e651bf8318687e95ac545358a884b7bf3/jax/\_src/nn/functions.py#L601
 [2] https://github.com/andeskyl/SHAFT/blob/8ade8e3858611983dee1e77f6da6d346e7e08831/crypten/common/functions/approximations.py#L604

softmax
$$(\vec{x})[j] = \frac{\exp(x_j)}{\sum_i \exp(x_i)}$$

#### JAX's snippet for Softmax

- 1. **def** \_softmax(x, axis = -1):
- 2. x\_max = jnp.max(x, axis,
- $\hookrightarrow$  keepdims=**True**)
- 3. unnormalized =  $jnp.exp(x x_max)$
- 4. sum\_exp = jnp.sum(unnormalized, axis,
- $\leftrightarrow$  keepdims=**True**)
- 5. result = unnormalized / sum\_exp
- 6. return result

#### SHAFT's snippet for Softmax

```
if method == "reciprocal":
```

maximum\_value = self.max(dim, keepdim=True)[0]

```
logits = self - maximum_value
```

```
numerator = logits.exp()
```

```
with cfg.temp_override({"functions.reciprocal_all_pos": True}):
```

inv\_denominator = numerator.sum(dim, keepdim=True).reciprocal()

```
return numerator * inv_denominator
```

[1] https://github.com/jax-ml/jax/blob/ed952c8e651bf8318687e95ac545358a884b7bf3/jax/\_src/nn/functions.py#L601
 [2] https://github.com/andeskyl/SHAFT/blob/8ade8e3858611983dee1e77f6da6d346e7e08831/crypten/common/functions/approximations.py#L604

softmax
$$(\vec{x})[j] = \frac{\exp(x_j)}{\sum_i \exp(x_i)}$$

- "keepdims=True" aims to align the operands shape
- However, it increases the costs of division by 10x 100x times!
  - x.div(y.broadcast(x.shape)):
     O(n) div and O(n) mul ☺
  - x.mul(y.recip.broadcast(x.shape)):
     O(1) recip and O(n) mul ☺

#### JAX's snippet for Softmax

- 1. **def** \_softmax(x, axis = -1):
- 2. x\_max = jnp.max(x, axis,
- $\hookrightarrow$  keepdims=**True**)
- 3. unnormalized =  $jnp.exp(x x_max)$
- 4. sum\_exp = jnp.sum(unnormalized, axis,
- $\leftrightarrow$  keepdims=**True**)
- 5. result = unnormalized / sum\_exp
- 6. return result

#### SHAFT's snippet for Softmax

```
if method == "reciprocal":
```

maximum\_value = self.max(dim, keepdim=True)[0]

```
logits = self - maximum_value
```

```
numerator = logits.exp()
```

```
with cfg.temp_override({"functions.reciprocal_all_pos": True}):
```

inv\_denominator = numerator.sum(dim, keepdim=True).reciprocal()

```
return numerator * inv_denominator
```

[1] https://github.com/jax-ml/jax/blob/ed952c8e651bf8318687e95ac545358a884b7bf3/jax/\_src/nn/functions.py#L601
 [2] https://github.com/andeskyl/SHAFT/blob/8ade8e3858611983dee1e77f6da6d346e7e08831/crypten/common/functions/approximations.py#L604

softmax
$$(\vec{x})[j] = \frac{\exp(x_j)}{\sum_i \exp(x_i)}$$

- "keepdims=True" aims to align the operands shape
- However, it increases the costs of division by 10x – 100x times!
  - x.div(y.broadcast(x.shape)):
     O(n) div and O(n) mul ☺
  - x.mul(y.recip.broadcast(x.shape)):
     O(1) recip and O(n) mul ☺



softmax $(\vec{x})[j] = \frac{\exp(x_j)}{\sum_i \exp(x_i)}$ 

- "keepdims=True" aims to align the operands shape
- However, it increases the costs of division by 10x – 100x times!
  - x.div(y.broadcast(x.shape)):
     O(n) div and O(n) mul ☺
  - x.mul(y.recip.broadcast(x.shape)):
     O(1) recip and O(n) mul <sup>(C)</sup>

#### A sad story ...

In BumbleBee, the division (reciprocal) in softmax takes less than 1% of the total costs.



SHAFT [3]

Apply a numerical method to avoid the division which can be handled more easily actually.

Wang et al. Characterization of MPC-based Private Inference for Transformer-based Models
 Li et al. MPCFORMER: FAST, PERFORMANT AND PRIVATE TRANSFORMER INFERENCE WITH MPC
 Andes Y. L. Kei et al. SHAFT: Secure, Handy, Accurate and Fast Transformer Inference

# Lesson 2: Piecewise rather than numerical methods

All about is the numerical stability

 $Gelu(x) = 0.5x(1 + \tanh(\sqrt{2/\pi}(x + 0.044715x^3)))$ 



$$Gelu(x) = 0.5x(1 + \tanh(\sqrt{2/\pi}(x + 0.044715x^3)))$$



- Need clipping on [-3, 3]
- Heavy division

[LHZWH] Squirrel: A Scalable Secure Two-Party Computation Framework for Training Gradient Boosting Decision Tree

Attempt 3

$$\mathsf{Gelu}(x) = 0.5x(1 + \tanh(\sqrt{2/\pi}(x + 0.044715x^3)))$$

#### Attempt 1

$$\tanh(x) \approx \frac{x + x^3/9 + x^5/945}{1 + 4x^2/9 + x^4/63}$$

- Need clipping on [-3, 3]
- Heavy division

tanh(x) 
$$\approx \sum_{j=0}^{7} F_j(x)$$

- 8 degree Fourier series approximation
- Need clipping on [-4, -4]
- 18 MULs

#### Attempt 3

$$\operatorname{Relu}(x) - \operatorname{Gelu}(x) \approx \sum_{j=0}^{7} G_j(x)$$

- 8 degree Fourier series approximation
- Need clipping on [-4, -4]
- 14 MULs + max(0, x)
- > 1450bytes per Gelu

 $Gelu(x) = 0.5x(1 + \tanh(\sqrt{2/\pi}(x + 0.044715x^3)))$ 



1. One-vs-Many Comparisons

2. Smoothness

$$\begin{cases} -\epsilon & x \le -4 \\ P(x) = \sum_{i=0}^{i=3} a_i x^i & -4 < x \le -1.85 \\ Q(x) = \sum_{i=0}^{i=6} b_i x^i & -1.85 < x \le 3 \\ x - \epsilon & x > 3 \end{cases}$$

- The one-vs-many comparisons has a smaller amortized costs in 2PC.
- 3 comparisons is same as Attempt 3
- Low-degree polys share terms. 6 MULs vs 14 MULs

1. One-vs-Many Comparisons

#### 2. Smoothness

$$\begin{cases} -\epsilon & x \le -4 \\ P(x) = \sum_{i=0}^{i=3} a_i x^i & -4 < x \le -1.85 \\ Q(x) = \sum_{i=0}^{i=6} b_i x^i & -1.85 < x \le 3 \\ x - \epsilon & x > 3 \end{cases}$$

.

- Smoothness allows an error-tune comparison (eg. ignoring some low-end bits)
- 718 bytes per gelu (50% off to Attempt 3)



# Lesson 3: Simplicity is good

All about is the numerical stability again

softmax
$$(\vec{x})[j] = \frac{\exp(x_j)}{\sum_i \exp(x_i)} = \frac{\exp(x_j - \hat{x})}{\sum_i \exp(x_i - \hat{x})} \quad \hat{x} = \max(\vec{x})$$

- Exp(x) on negative inputs are bounded.
- Good enough approximation by low-degree Taylor with one clipping

softmax
$$(\vec{x})[j] = \frac{\exp(x_j)}{\sum_i \exp(x_i)} = \frac{\exp(x_j - \hat{x})}{\sum_i \exp(x_i - \hat{x})}$$

- Exp(x) on negative inputs are bounded.
- Good enough approximation by low-degree with one clipping



#### "The reviewer" C:

Also, recent work at ACSAC 2023 [c] demonstrates improved secure softmax computation. Integrating recent developments could provide a better review of related work and may also improve the practicality and efficiency of the proposed protocol.

#### "The reviewer" C:

Also, recent work at ACSAC 2023 [c] demonstrates improved secure softmax computation. Integrating recent developments could provide a better review of related work and may also improve the practicality and efficiency of the proposed protocol.

```
# self \in [-iter_num, iter_num]
def ACSAC_softmax(self, iter_num):
    dim = self.shape[-1]
    x = self / iter_num
    g = np.ones(self.shape) / dim
    for _ in range(iter_num):
        z = x * g
        g = g + (z - sum(z) * g)
    return g
```

Pros:

- Division-free
- Maximum-free

Cons:

• #Iteration grows linear with the approximation range (which could be huge in LLM)

#### "The reviewer" C:

Also, recent work at ACSAC 2023 [c] demonstrates improved secure softmax computation. Integrating recent developments could provide a better review of related work and may also improve the practicality and efficiency of the proposed protocol.

```
# self \in [-iter_num, iter_num]
def ACSAC_softmax(self, iter_num):
    dim = self.shape[-1]
    x = self / iter_num
    g = np.ones(self.shape) / dim
    for _ in range(iter_num):
        z = x * g
        g = g + (z - sum(z) * g)
    return g
```

Pros:

- Division-free
- Maximum-free

Cons:

 #Iteration grows linear with the approximation range (which could be huge in LLM)

#### "The reviewer" C:

Also, recent work at ACSAC 2023 [c] demonstrates improved secure softmax computation. Integrating recent developments could provide a better review of related work and may also improve the practicality and efficiency of the proposed protocol.

```
# self \in [-iter_num, iter_num]
def ACSAC_softmax(self, iter_num):
    dim = self.shape[-1]
    x = self / iter_num
    g = np.ones(self.shape) / dim
    for _ in range(iter_num):
        z = x * g
        g = g + (z - sum(z) * g)
    return g
```

What if applying range clipping [1]?

- O(2n) comparisons
- However, the maximum only needs O(n) comparisons

# Thanks!

- <u>1678.pdf (iacr.org)</u>
- AntCPLab/OpenBumbleBee (github.com)

## Anecdote: Why Naming BumbleBee

#### **Cheetah: Lean and Fast Secure Two-Party Deep Neural Network Inference**

Zhicong Huang Alibaba Group Wen-jie Lu Alibaba Group Cheng Hong Alibaba Group

Jiansheng Ding Alibaba Group





## Anecdote: Why Naming BumbleBee



Cheetah + Transfomer => Cheetor ("Maximize!")

The very first computer-aid 3D animation in 1996.

## Anecdote: Why Naming BumbleBee



Team Leader: Cheetor is so unknown! Let be Bumblebee.