

# Review on

**[MiniONN, 2017], [Gazelle, 2018] and [MPCViT, 2022]**

**Xue Yufei, Southeast University**

**03/31/2023**

# Content

---

## 1. MiniONN

- ✓ Protocol for matrix-vector multiplication
- ✓ Protocol for linear/non-linear operation (Sigmoid)

## 2. GAZELLE

- ✓ FHE
- ✓ SIMDScMult/Conv

## 3. MPCViT

- ✓ Idea

# MiniONN: Multiplication

**Input:**

$S: \mathbf{w} \in \mathbb{Z}_N^n$

$C: \mathbf{r} \in \mathbb{Z}_N^n$

**Output:**

$S$ : a random number  $u \in \mathbb{Z}_N$ ;

$C$ :  $v \in \mathbb{Z}_N$ , s.t.,  $u + v \pmod N = \mathbf{w} \cdot \mathbf{r}$ .

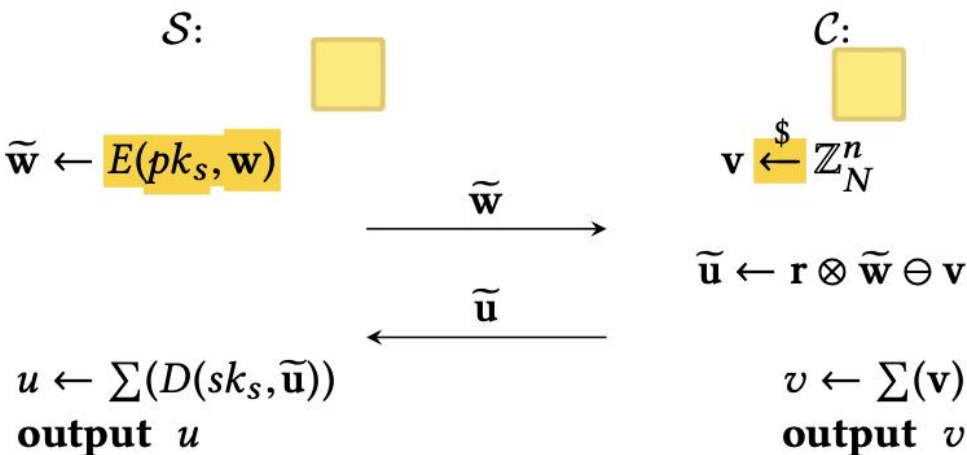


Figure 3: Dot-product triplet generation.

**Input:**

- $S$ : a vector  $\mathbf{w} \in \mathbb{Z}_N^n$ ;

- $C$ : a random vector  $\mathbf{r} \in \mathbb{Z}_N^n$ .

**Output:**

- $S$ : a random number  $u \in \mathbb{Z}_N$ ;

- $C$ :  $v \in \mathbb{Z}_N$ , s.t.,  $u + v \pmod N = \mathbf{w} \cdot \mathbf{r}$ .

Figure 1: Ideal functionality  $\mathcal{F}_{\text{triplet}}$ : generate a dot-product triplet.

$$u \rightarrow S$$

$$\mathbf{r}, v \rightarrow C$$

# MiniONN: Linear operation

**Input:**

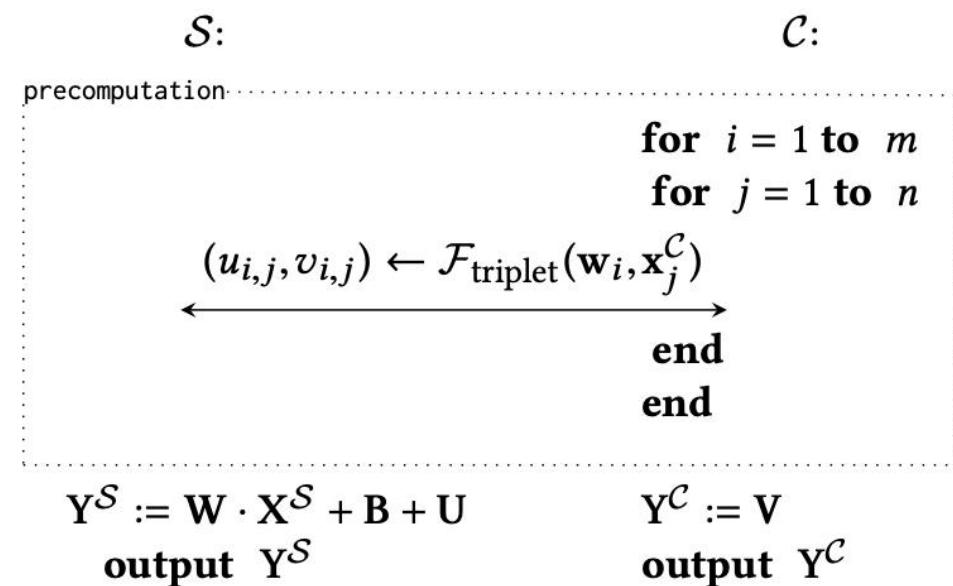
$\mathcal{S}$ :  $\mathbf{W} \in \mathbb{Z}_N^{m \times l}$ ,  $\mathbf{X}^{\mathcal{S}} \in \mathbb{Z}_N^{l \times n}$ ,  $\mathbf{B} \in \mathbb{Z}_N^{m \times n}$

$\mathcal{C}$ :  $\mathbf{X}^{\mathcal{C}} \in \mathbb{Z}_N^{l \times n}$

**Output:**

$\mathcal{S}$ : A random matrix  $\mathbf{Y}^{\mathcal{S}}$

$\mathcal{C}$ :  $\mathbf{Y}^{\mathcal{C}}$  s.t.,  $\mathbf{Y}^{\mathcal{C}} + \mathbf{Y}^{\mathcal{S}} = \mathbf{W} \cdot (\mathbf{X}^{\mathcal{C}} + \mathbf{X}^{\mathcal{S}}) + \mathbf{B}$



预生成三元组, 节省开销



$$\mathbf{Y}^{\mathcal{C}} + \mathbf{Y}^{\mathcal{S}} = \mathbf{Y}$$

Figure 4: Oblivious linear transformation.

# MiniONN: Activation functions

## ReLU: garbled circuit

**Input:**

- $\mathcal{S}: y^{\mathcal{S}} \in \mathbb{Z}_N$ ;
- $\mathcal{C}: y^{\mathcal{C}}, r \in \mathbb{Z}_N$ .

**Output:**

- $\mathcal{S}: x^{\mathcal{S}} := \text{compare}(y, 0) \cdot y - r \pmod{N}$  where  $y = y^{\mathcal{S}} + y^{\mathcal{C}} \pmod{N}$ ;
- $\mathcal{C}: x^{\mathcal{C}} := r$ .

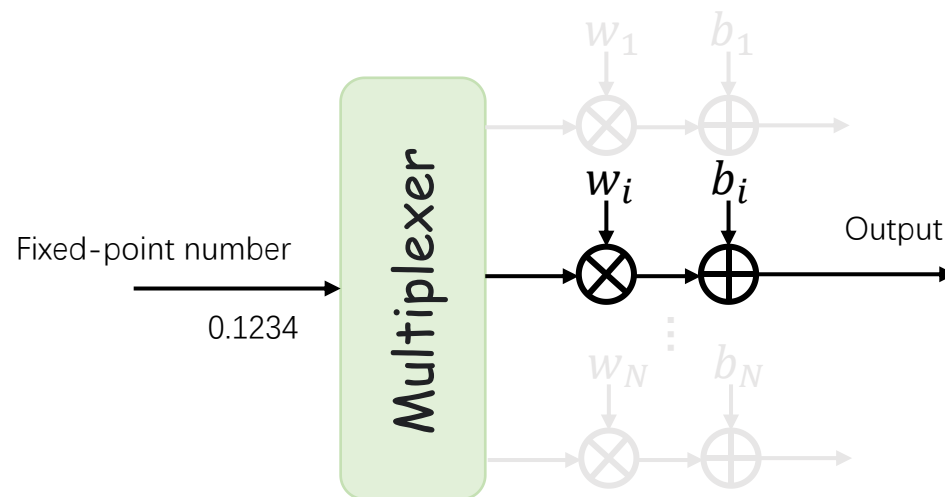
Figure 5: The ideal functionality  $\mathcal{F}_{\text{ReLU}}$ .

## Sigmoid:

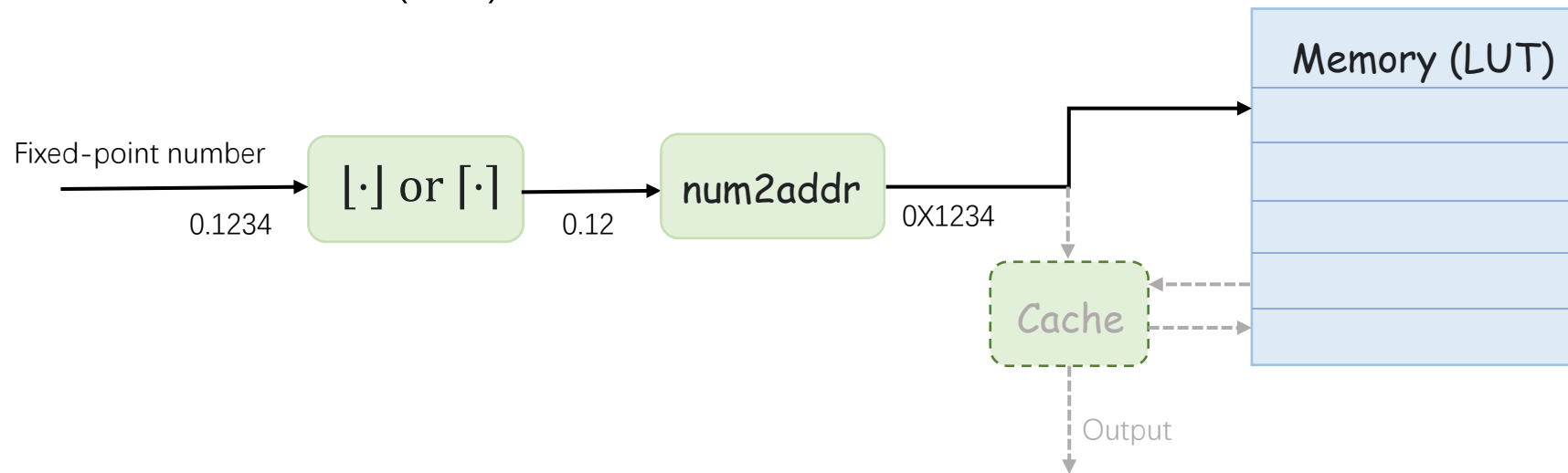
1. 多项式拟合 例:  $x^3 + x^2 + x = x \times (x^2 + x + 1) \rightarrow x \times x \times (x + 1)$
2. 分段计算, 会用到 $\text{compare}$ 单元

# MiniONN: Sigmoid

## 1. 使用分段线性拟合

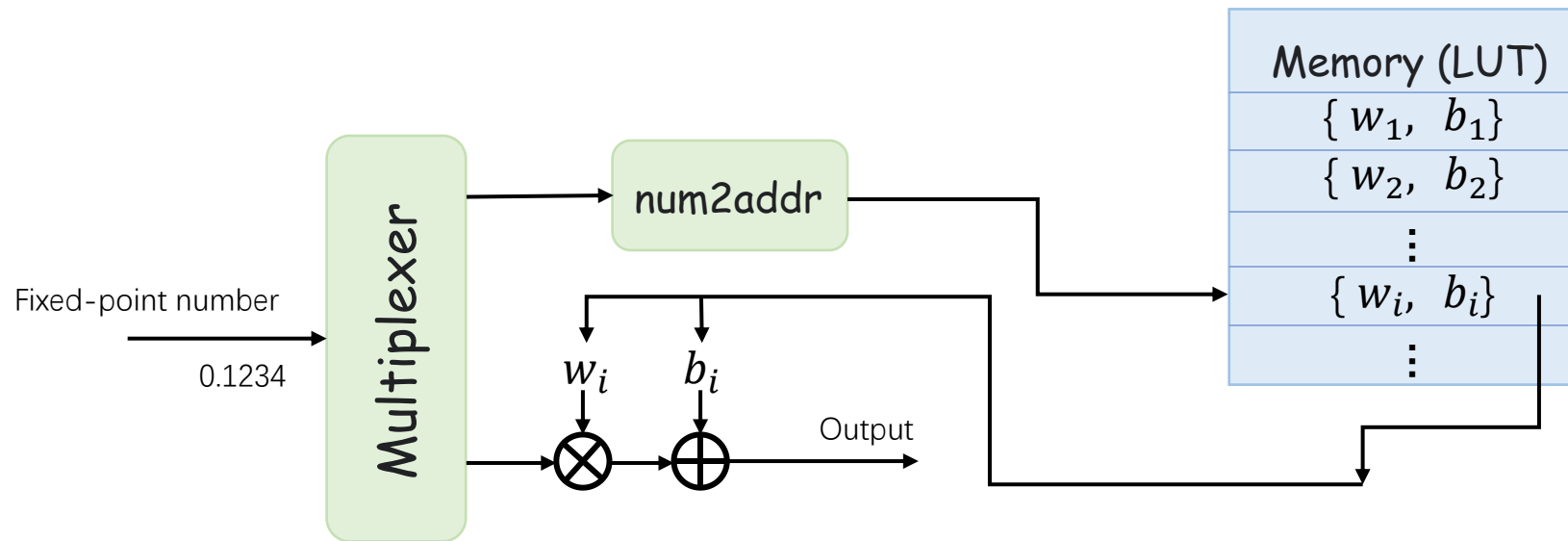


## 2. 使用LUT计算非线性函数 ( $A^{-1}$ )



# MiniONN: Sigmoid

## 3. 复用计算单元



# GAZELLE: FHE

---

加法同态:  $\mathbf{En}(x + y) = \mathbf{En}(x) \oplus \mathbf{En}(y)$

乘法同态:  $\mathbf{En}(x \times y) = \mathbf{En}(x) \otimes \mathbf{En}(y)$

**v.s.**

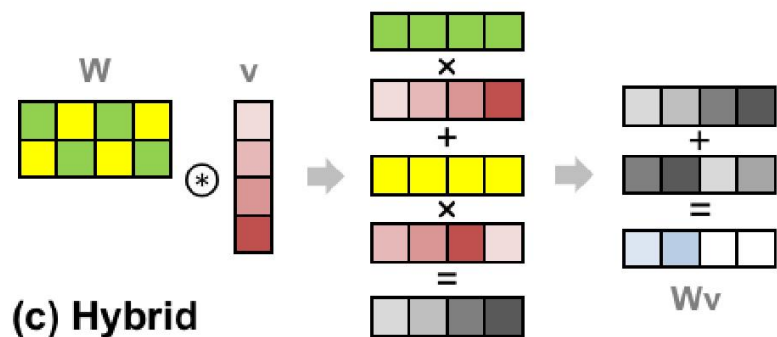
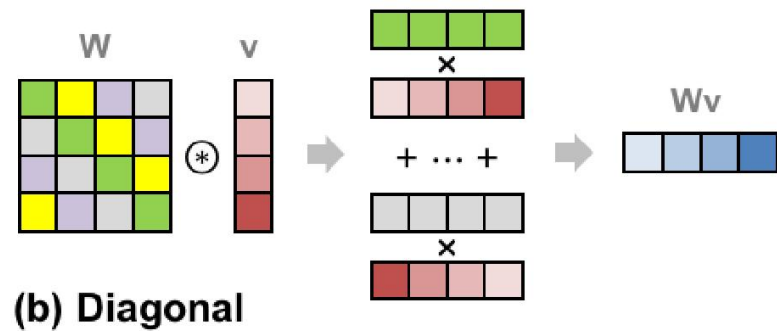
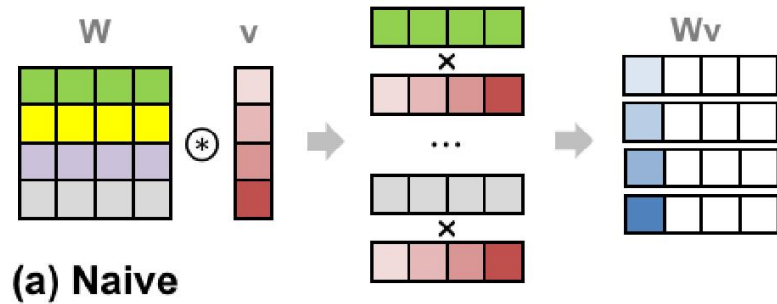
$\mathbf{FFT}(x + y) = \mathbf{FFT}(x) + \mathbf{FFT}(y)$

$\mathbf{FFT}(x \times y) = \mathbf{FFT}(x) * \mathbf{FFT}(y)$

线性系统的线性性质



# GAZELLE: SIMDScMult



- How did the authors think about it?
- Is there an optimization?  
(To note)

# GAZELLE: SIMDScMult

---

# GAZELLE: SIMDScMult

---

# GAZELLE: SIMDScMult

---

# GAZELLE: SIMDScMult

---

# GAZELLE: Conv

## Optimization for padding

$$\begin{array}{|c|c|c|} \hline 5 & 6 & 7 \\ \hline 8 & 0 & 1 \\ \hline 2 & 3 & 4 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline 6 & 7 & 8 \\ \hline 0 & 1 & 2 \\ \hline 3 & 4 & 5 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline 7 & 8 & 0 \\ \hline 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \quad \longrightarrow \quad \boxed{\text{red}} \text{ 的部分保留, 其余部分置零}$$

$\pi_{-4}(\text{in})$     $f_{(-1,-1)}$     $\pi_{-3}(\text{in})$     $f_{(0,-1)}$     $\pi_{-2}(\text{in})$     $f_{(1,-1)}$

$$\begin{array}{|c|c|c|} \hline 8 & 0 & 1 \\ \hline 2 & 3 & 4 \\ \hline 5 & 6 & 7 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline 3 & 4 & 5 \\ \hline 6 & 7 & 8 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 0 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} +$$

$\pi_{-1}(\text{in})$     $f_{(-1,0)}$     $\pi_0(\text{in})$     $f_{(0,0)}$     $\pi_1(\text{in})$     $f_{(1,0)}$

$$\begin{array}{|c|c|c|} \hline 2 & 3 & 4 \\ \hline 5 & 6 & 7 \\ \hline 8 & 0 & 1 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline 3 & 4 & 5 \\ \hline 6 & 7 & 8 \\ \hline 0 & 1 & 2 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline 4 & 5 & 6 \\ \hline 7 & 8 & 0 \\ \hline 1 & 2 & 3 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline 3 & 4 & 5 \\ \hline 6 & 7 & 8 \\ \hline \end{array}$$

$\pi_2(\text{in})$     $f_{(-1,1)}$     $\pi_3(\text{in})$     $f_{(0,0)}$     $\pi_4(\text{in})$     $f_{(1,1)}$    **out**

## Channel Packing:

Input  $c_i$ , output  $c_o$ , integer  $c_n \rightarrow$  Input  $\frac{c_i}{c_n}$ , output  $\frac{c_o}{c_n}$

计算结果一样吗? (to note)

# GAZELLE: Conv

---

# MPCViT:

General idea:

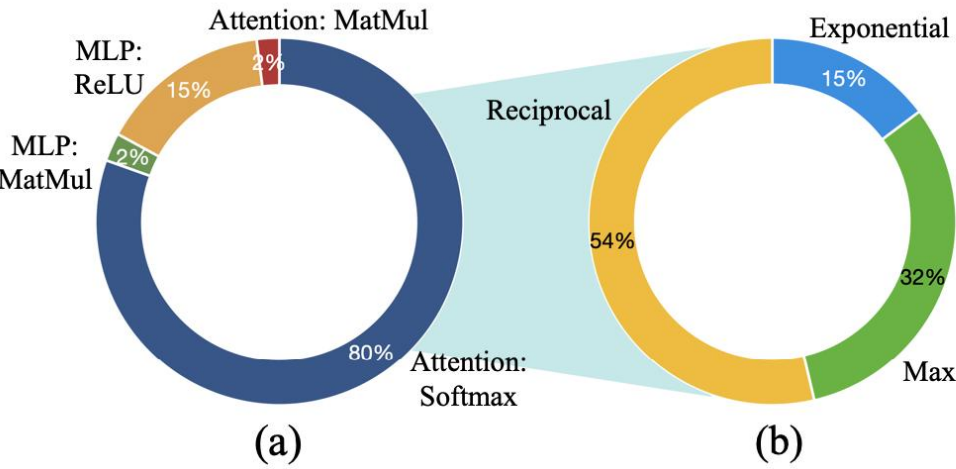
Low-latency

NLP skills → CV

MPCViT: Searching for MPC-friendly Vision Transformer  
with Heterogeneous Attention  
 $\alpha A + (1 - \alpha)B$

Step:

1. find most time-consuming operation
2. get rid of  $e^x$ , replaced by ScaleAttn+RSAttn



Type	Top-1 Acc. (%)	Latency (s)
Softmax Attention	92.69	6.82
ReLU Attention	fail	fail
ReLU6 Attention	90.50	3.02
Sparsemax Attention [24]	91.23	3.23
XNorm Attention [10]	91.24	13.25
Square Attention	91.27	0.72
2Quad Attention [23]	91.86	4.22
ScaleAttn [12]	91.52	0.66
ReLU Softmax Attention (RSAttn) [14]	92.31	5.32



# MPCViT:

---

## 1. Differentiable NAS

$$\alpha \cdot \text{ReLU}\text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V + (1 - \alpha) \cdot \frac{\text{ScaleAttn}(Q, K, V)}{\sqrt{d_k}}$$

取RSAttn的准确性，取ScaleAttn的低时延，更倾向ScaleAttn

---

## 2. How to effectively train?

Token-wise feature-based KD