

# Bit Fusion: Bit-Level Dynamically Composable Architecture for Accelerating Deep Neural Networks

比特融合：对加速**DNN**的比特级动态可组合架构

ISCA'2018

---

汇报人：王铎【中国科学院大学】

*May 22, 2019*

## 1 介绍

---

- 背景
- 思想

## 2 解决方案

---

## 3 实验

---

## 背景

- 量化方法，在不损失精度下减少操作数位宽
- 对不同的DNNs的比特宽变化不同,甚至需要单独调整每一层

## 问题

- 固定位宽度加速器要么提供有限的好处,以适应最坏的情况下的比特宽要求,或不可避免地导致最终精度的下降

## 1 介绍

---

- 背景
- 思想

## 2 解决方案

---

## 3 实验

---

## Bit Fusion加速器

- 探索动态位级融合/分解新维度
- 构成了一组位级处理元素BitBrick的数组 Fusion-PE
- 动态地融合到单个DNN层的比特宽
- 最优粒度下使计算和通信最小化

## 1 介绍

---

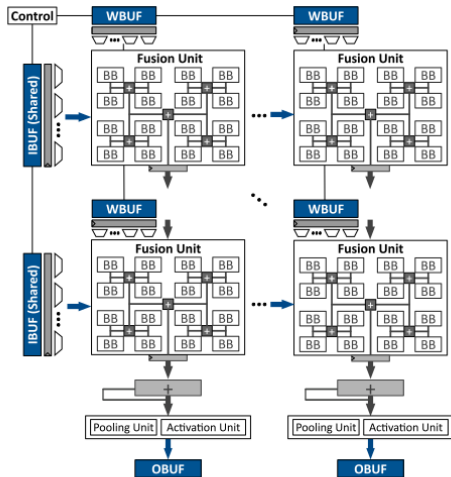
## 2 解决方案

---

- 动态融合
- ISA

## 3 实验

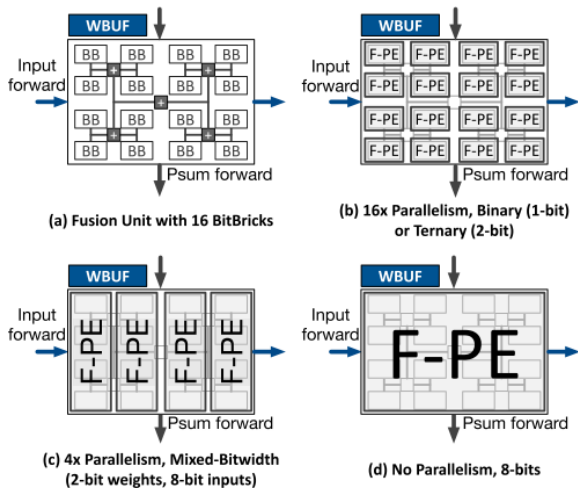
---



## BitBrick (BB)

- 均可执行单个二值(0、+1)和三值(-1、0、+1)的乘加操作

Figure 2: Bit Fusion systolic architecture comprising a collection of BitBricks (BBs) that can fuse to form FPEs.

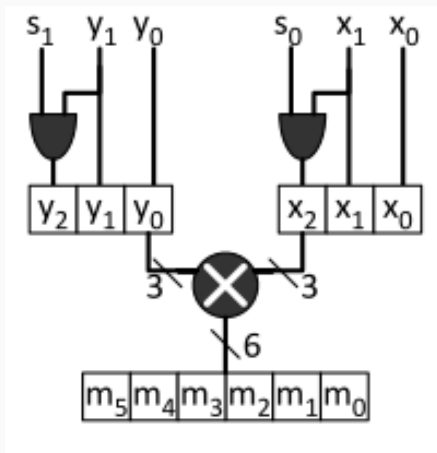


**Figure 3: Dynamic composition of BitBricks (BBs) to construct Fused Processing Engines (FPEs).**

## Fusion Unit

- 一个具有16个BitBrick
- 可以提  
供1、2、4、8和16个  
有不同操作数位宽  
的Fused-PE





**Figure 5: A single BitBrick.**

## BitBrick

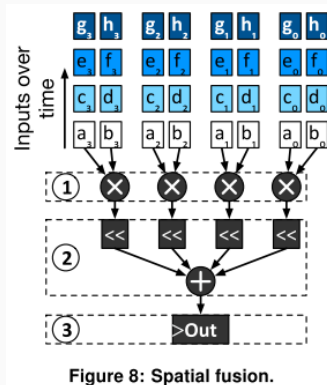
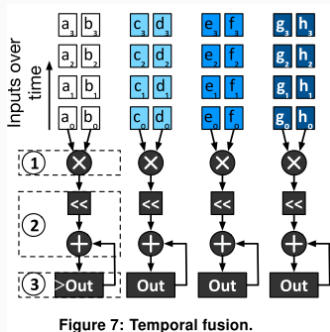
- 两个2-bit操作数 $x_1, x_0$  和 $y_1, y_0$ , 两个符号位 $s_1, s_0$
- 符号位定义了2-bit操作数 $x_1, x_0$ 带符号（-2到1）或者无符号（0到3）

## Bitwidth Mode

## Operation

2n-bit × 2n-bit	$2^{2n} \times (A_{2n})_{hi} \otimes (B_{2n})_{hi} + 2^n \times (A_{2n})_{hi} \otimes (B_{2n})_{lo}$ $+ 2^n \times (A_{2n})_{lo} \otimes (B_{2n})_{hi} + 2^0 \times (A_{2n})_{lo} \otimes (B_{2n})_{lo}$
2n-bit × n-bit	$2^n \times (A_{2n})_{hi} \otimes B_n + 2^0 \times (A_{2n})_{lo} \otimes B_n$ $+ 2^n \times (C_{2n})_{hi} \otimes D_n + 2^0 \times (C_{2n})_{lo} \otimes D_n$
n-bit × 2n-bit	$2^n \times A_n \otimes (B_{2n})_{hi} + 2^n \times C_n \otimes (D_{2n})_{hi}$ $+ 2^0 \times A_n \otimes (B_{2n})_{lo} + 2^0 \times C_n \otimes (D_{2n})_{lo}$
n-bit × n-bit	$2^0 \times A_n \otimes B_n + 2^0 \times C_n \otimes D_n$ $+ 2^0 \times E_n \otimes F_n + 2^0 \times G_n \otimes H_n$

**Figure 6: Supported modes for a fuse\_multiply<sub>n→2n</sub>**



空间融合比时间融合需要更多面积但是性能较好

混合模式：空间融合组合16个BitBrick以实现 $\text{fuse\_multiply}_{2 \rightarrow 8}$ 算子，再利用时间融合组合每四个周期发射 $\text{fuse\_multiply}_{2 \rightarrow 16}$

## 1 介绍

---

## 2 解决方案

---

- 动态融合
- **ISA**

## 3 实验

---

Table 1: Bit Fusion Instruction Set.

OpCode	Operand Specification		Loop Identifier	Immediate
5-bits	6-bits		5-bits	16-bits
<i>setup</i>	op0.bitwidth	op1.bitwidth	X	X
<i>ld-mem</i>	scratchpad-type	mem.bitwidth	loop-id	num-words
<i>st-mem</i>				
<i>rd-buf</i>		X		X
<i>wr-buf</i>				
<i>gen-addr</i>		ld/st		stride
<i>compute</i>	fn			X
<i>loop</i>	X	loop-level		num-iterations
<i>block-end</i>	Address of next instruction			

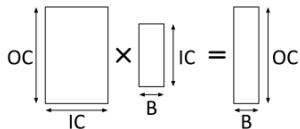


Figure 10: A single Fully-Connected Layer. The  $\times$  symbol represents matrix-matrix multiplication.

Compute     Memory  
 Loop     Buffer

```

loop: oc -> (OC)
  loop: ic -> (IC)
    ld-mem, IBUF, 1
    ld-mem, WBUF, 1
    ld-mem, OBUF, 1
    rd-buf, IBUF -> in
    rd-buf, WBUF -> wgt
    rd-buf, OBUF -> out
    out += in * wgt
    wr-buf: out -> OBUF
    st-mem: OBUF, 1
    
```

(a) Initial code

```

loop: toc -> (1, #tileoc)
  ld-mem, OBUF, tileoc
  loop: tic -> (1, #tileic)
    ld-mem, IBUF, tileic
    ld-mem, WBUF, tilewgt
    loop: oc -> (1, tileoc)
      rd-buf, OBUF -> out
      loop: ic -> (1, tileic)
        rd-buf, IBUF -> in
        rd-buf, WBUF -> wgt
        out += in * wgt
      wr-buf: out -> OBUF
    st-mem, OBUF, tileoc
    
```

(b) Optimized code

Figure 11: (a) The equivalent code for the Fully Connected Layer. (b) Optimized code using loop tiling and ordering. setup and gen-addr instructions omitted for clarity.

循环排序: 优化外层循环和存储指令顺序以降低片外访问, 在输入稳定、输出稳定和权重稳定之间切换。

循环切片: 使得循环操作所需的数据适合于片上 scratchpad。

## 1 介绍

---

## 2 解决方案

---

## 3 实验

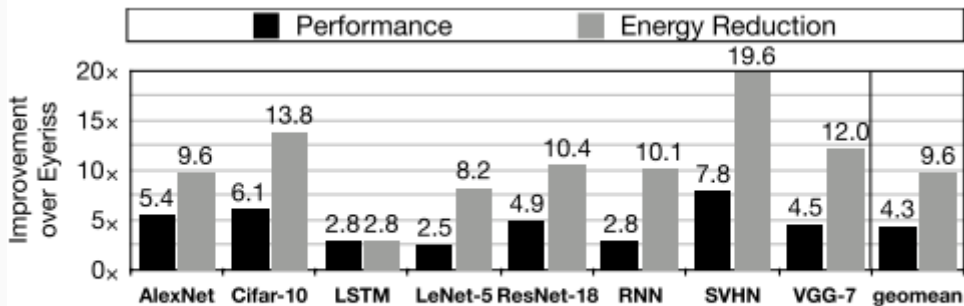
---

- 实验

**Table 3: The evaluated ASIC and GPU platforms. \*Stripes entries are per-tile.**

	ASIC				GPU	
Chip	BitFusion	Eyeriss	Stripes*	Chip	Titan X	Tegra X2
Cores (1.1 mm <sup>2</sup> )	512 FU	168 Pes	4096 SIPs	Cores	3,584	256
On-chip Memory	112 KB	108 KB	2 MB eDRAM 16 KB SRAM	Memory	12 GB	8 GB
Chip Area (mm <sup>2</sup> )	5.87	5.87	3.62	TDP	250 W	7.5 W
Frequency	500 MHz	500 MHz	980 MHz	Frequency	1,531 MHz	875 MHz
Technology	45	45	45	Technology	16 nm	16 nm





**Figure 12: Bit Fusion performance and energy improvements over Eyeriss.**

## Eyeriss

4.3 $\times$ 加速和9.6 $\times$ 能耗节省

## Stripes<sup>a</sup>

2.4 $\times$ 加速和4.1 $\times$ 能耗节省

## GPU

比Tegra X2有4.3 $\times$ 加速

---

<sup>a</sup>P. Judd, J. Albericio, T. Hetherington, T. M. Aamodt, and A. Moshovos, “Stripes: Bit-serial deep neural network computing,” in MICRO, 2016.

## 动态位级融合和分解

- 将比特级可组合的PE与DNN层的变化位宽相匹配
- 最优粒度下使计算和通信最小化

## 对位级组合能力的微结构设计

- 一种二维的BitBricks阵列
- 构造了一个可融合的处理引擎
- 在不同的位宽下执行DNN计算

## 硬件-软件抽象以比特灵活的加速 **Fusion-ISA**

# Thanks

---

**Bit Fusion: Bit-Level Dynamically Composable  
Architecture for Accelerating Deep Neural  
Networks**

比特融合：对加速**DNN**的比特级动态可组合架构