



北京大學  
PEKING UNIVERSITY

# 人工智能的硬件基石

## 从物理器件到计算架构

第十一讲：AI芯片设计-II

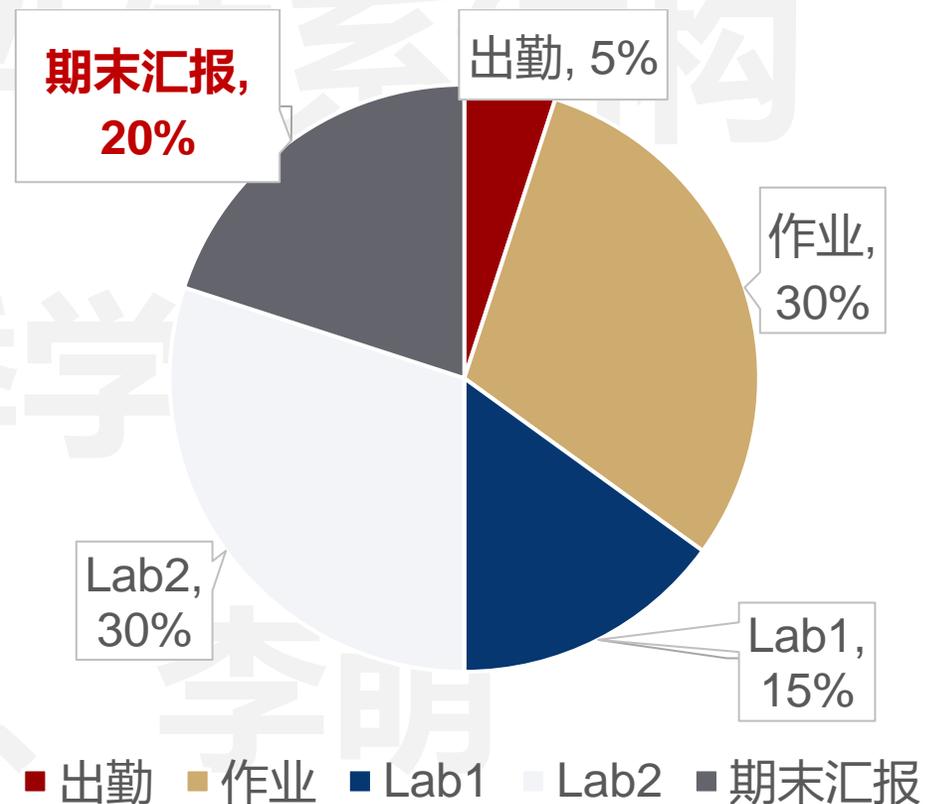
主讲：陶耀宇

2025年春季

## 注意事项

### • 课程作业情况

- **第4次作业时间：4月28号 – 5月12号**
  - 6次作业可以使用总计6个Late day
  - Late Day耗尽后，每晚交1天扣除20%当次作业分数
- **第2次lab时间：4月17日 – 6月15日**
  - **基础任务（多个test programs通过率）**
  - **Bonus（占比Lab2的50%，简单AI卷积加速核设计，开放式、无强制要求，Handout四月内放出）**
- **期末汇报**
  - 请自行挑选并阅读近5年内的芯片/架构等相关论文
  - 请同学们准备10页左右的PPT，**每位同学12分钟演讲 + 3分钟答疑**
  - **每位同学均为汇报同学打分，互评成绩占期末汇报成绩50%，请于5月19日前将汇报题目发给助教**



- 课程作业情况

- **期末汇报（每位同学12分钟汇报+3分钟答疑）**

- **AI相关的神经网络加速芯片**

- 传统AI加速器：Fused-layer cnn accelerators、Eyeries, Google TPU等
- 新兴AI加速器（大模型Transformer、Neural ODE、MANN/DNC、PINN等）

- **GPGPU等高并行并行芯片**

- 流式多处理器（Multithreaded Streaming Multiprocessors, CUDA的来源）
- FPGA架构等（可编程逻辑块、可编程路由等）

- **安全与通信领域处理器芯片**

- 各类密钥编码（AES、RSA）、视频编码（MPEG等）、通信编码（LDPC、Polar等）

- **传统CPU芯片**

- 优化Branch Predictor、Load-Store、缓存预读取、众核缓存一致性等

- **后摩尔时代的新兴芯片**

- 存算一体/感存算一体、量子计算、生物信息处理、高维NoC、区块链
- 基于后摩尔非CMOS器件的架构（模拟计算架构、动力学计算架构等）

# 目录

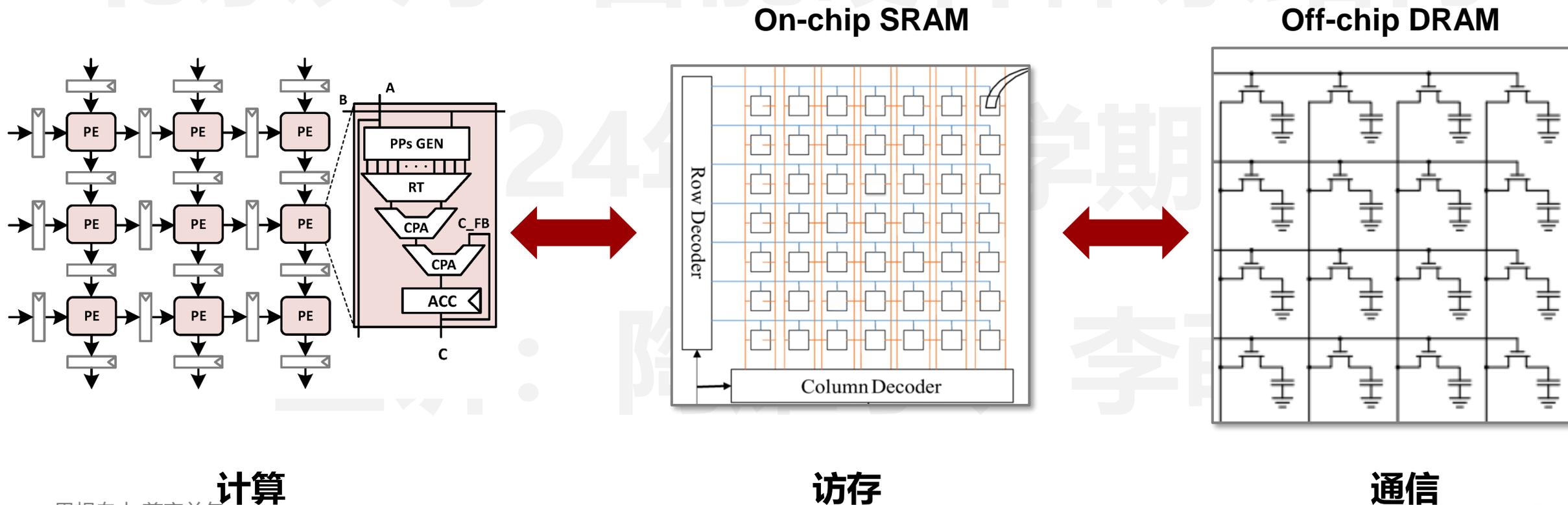
CONTENTS



01. 典型AI芯片架构
02. AI芯片软硬件协同设计
03. AI大模型芯片设计
04. 存算一体AI芯片架构

# AI加速器架构基础-Recap

- 我们介绍了AI加速器的基础架构，并着重介绍了AI加速器中的计算单元、片上访存以及片间/片外通信
- 我们将介绍有代表性的**AI加速器架构**



# AI加速器架构发展——DianNao

- ASPLOS 2014最佳论文
- DianNao对于神经网络加速器设计进行了系统讨论
  - 如何设计计算阵列?
  - 如何设计片上存储?
  - 如何提升数据复用, 降低访存开销?
  - 如何灵活支持不同操作?
- DianNao面向MLP和CNN, 相关参数如下

Performance	452 GOPS
Area	3.02 mm <sup>2</sup>
Power	485 mW
Technology Node	65 nm
Baseline	117.87X speedup over SIMD processor



## DianNao: A Small-Footprint High-Throughput Accelerator for Ubiquitous Machine-Learning

Tianshi Chen  
SKLCA, ICT, China

Zidong Du  
SKLCA, ICT, China

Ninghui Sun  
SKLCA, ICT, China

Jia Wang  
SKLCA, ICT, China

Chengyong Wu  
SKLCA, ICT, China

Yunji Chen  
SKLCA, ICT, China

Olivier Temam  
Inria, France

### Abstract

Machine-Learning tasks are becoming pervasive in a broad range of domains, and in a broad range of systems (from embedded systems to data centers). At the same time, a small set of machine-learning algorithms (especially Convolutional and Deep Neural Networks, i.e., CNNs and DNNs) are proving to be state-of-the-art across many applications. As architectures evolve towards heterogeneous multi-cores composed of a mix of cores and accelerators, a machine-learning accelerator can achieve the rare combination of efficiency (due to the small number of target algorithms) and broad application scope.

Until now, most machine-learning accelerator designs have focused on efficiently implementing the computational part of the algorithms. However, recent state-of-the-art CNNs and DNNs are characterized by their large size. In this study, we design an accelerator for large-scale CNNs and DNNs, with a special emphasis on the impact of memory on accelerator design, performance and energy.

We show that it is possible to design an accelerator with a high throughput, capable of performing 452 GOPs (key NN operations such as synaptic weight multiplications and

neurons outputs additions) in a small footprint of 3.02 mm<sup>2</sup> and 485 mW; compared to a 128-bit 2GHz SIMD processor, the accelerator is 117.87x faster, and it can reduce the total energy by 21.08x. The accelerator characteristics are obtained after layout at 65nm. Such a high throughput in a small footprint can open up the usage of state-of-the-art machine-learning algorithms in a broad set of systems and for a broad set of applications.

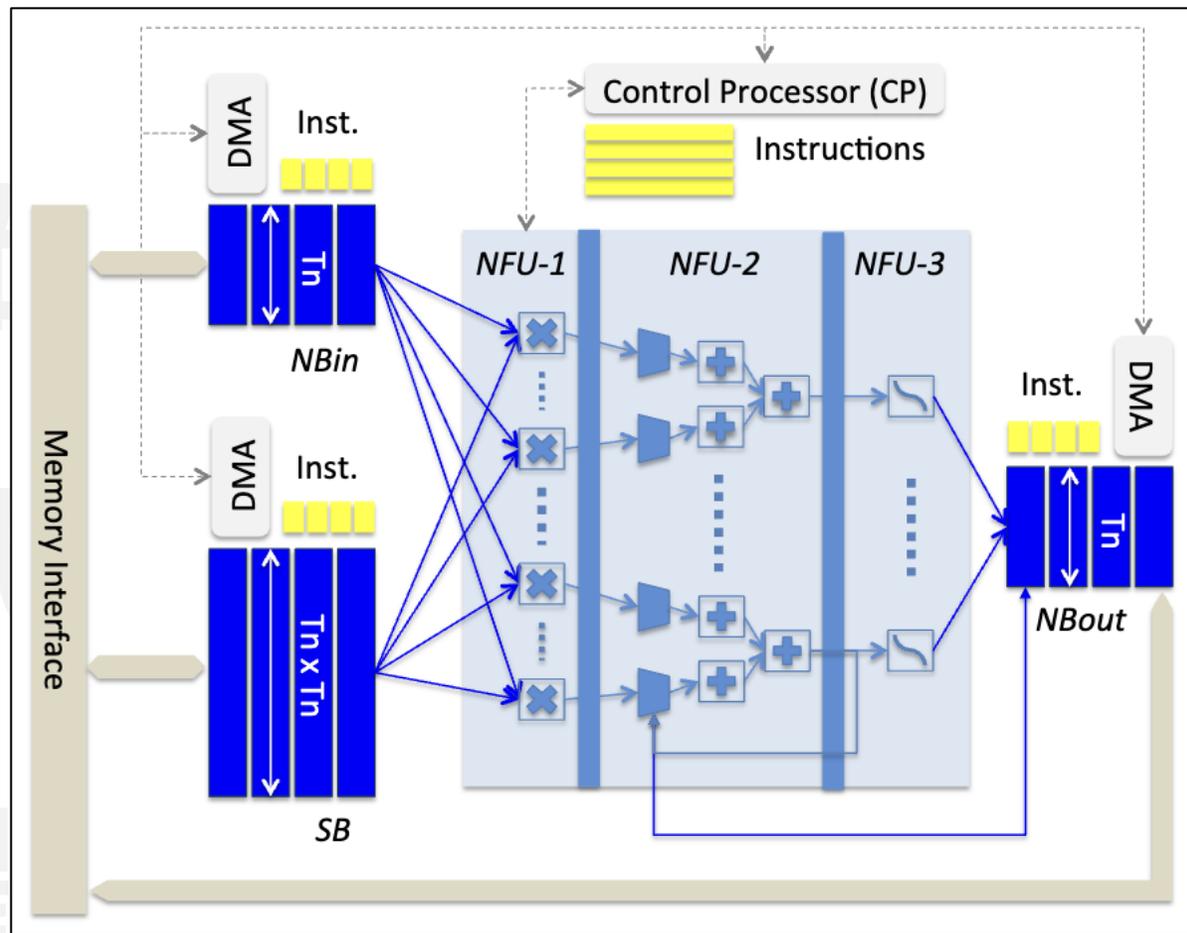
### 1. Introduction

As architectures evolve towards heterogeneous multi-cores composed of a mix of cores and accelerators, designing accelerators which realize the best possible tradeoff between flexibility and efficiency is becoming a prominent issue.

The first question is for which category of applications one should primarily design accelerators? Together with the architecture trend towards accelerators, a second simultaneous and significant trend in high-performance and embedded applications is developing: many of the emerging high-performance and embedded applications, from image/video/audio recognition to automatic translation, business analytics, and all forms of robotics rely on *machine-learning techniques*. This trend even starts to percolate in our community where it turns out that about half of the benchmarks of PARSEC [2], a suite partly introduced to highlight the emergence of new types of applications, can be implemented using machine-learning algorithms [4]. This trend in application comes together with a third and equally remarkable trend in machine-learning where a small number of techniques, based on neural networks (especially Convolutional Neural Networks [27] and Deep Neural Networks

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
ASPLOS '14, March 1-5, 2014, Salt Lake City, Utah, USA.  
Copyright © 2014 ACM 978-1-4503-2305-5/14/03...\$15.00.  
<http://dx.doi.org/10.1145/http://dx.doi.org/10.1145/2541940.2541967>

- 如何进行神经网络加速器设计?
  - 直接实例化完整神经网络 vs 通用神经网络计算架构实现多算子支持
- 如何进行计算阵列设计?
- 如何设计访存?
  - High associative cache vs 便签存储器



# AI加速器架构发展——DianNao

```
for (int nnn = 0; nnn < Nn; nnn += Tnn) { // tiling for output neurons;
  for (int iii = 0; iii < Ni; iii += Tii) { // tiling for input neurons;
    for (int nn = nnn; nn < nnn + Tnn; nn += Tn) {
      for (int n = nn; n < nn + Tn; n++)
        sum[n] = 0;
      for (int ii = iii; ii < iii + Tii; ii += Ti)
        // — Original code —
        for (int n = nn; n < nn + Tn; n++)
          for (int i = ii; i < ii + Ti; i++)
            sum[n] += synapse[n][i] * neuron[i];
      for (int n = nn; n < nn + Tn; n++)
        neuron[n] = sigmoid(sum[n]);
    }
  }
}
```

## 全连接层

```
for (int yy = 0; yy < Nyin; yy += Ty) {
  for (int xx = 0; xx < Nxin; xx += Tx) {
    for (int nnn = 0; nnn < Nn; nnn += Tnn) {
      // — Original code — (excluding nn, ii loops)
      int yout = 0;
      for (int y = yy; y < yy + Ty; y += sy) { // tiling for y;
        int xout = 0;
        for (int x = xx; x < xx + Tx; x += sx) { // tiling for x;
          for (int nn = nnn; nn < nnn + Tnn; nn += Tn) {
            for (int n = nn; n < nn + Tn; n++)
              sum[n] = 0;
            // sliding window;
            for (int ky = 0; ky < Ky; ky++)
              for (int kx = 0; kx < Kx; kx++)
                for (int ii = 0; ii < Ni; ii += Ti)
                  for (int n = nn; n < nn + Tn; n++)
                    for (int i = ii; i < ii + Ti; i++)
                      // version with shared kernels
                      sum[n] += synapse[ky][kx][n][i]
                        * neuron[ky + y][kx + x][i];
                      // version with private kernels
                      sum[n] += synapse[yout][xout][ky][kx][n][i]
                        * neuron[ky + y][kx + x][i];
                  for (int n = nn; n < nn + Tn; n++)
                    neuron[yout][xout][n] = non_linear_transform(sum[n]);
          } xout++; } yout++;
        }
      }
    }
  }
}
```

## 卷积层

- 如何处理池化层?

- 池化层和卷积层具有不同的访存需求
- 通过局部转置的方式，在DMA读入输入激活值的时候，采用interleaving的方式，将输入激活值存储在便签存储器中

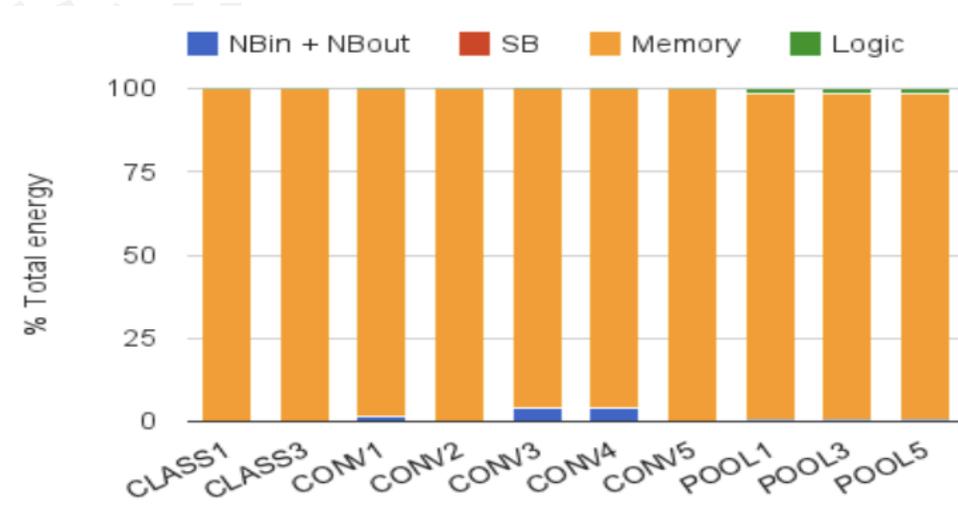
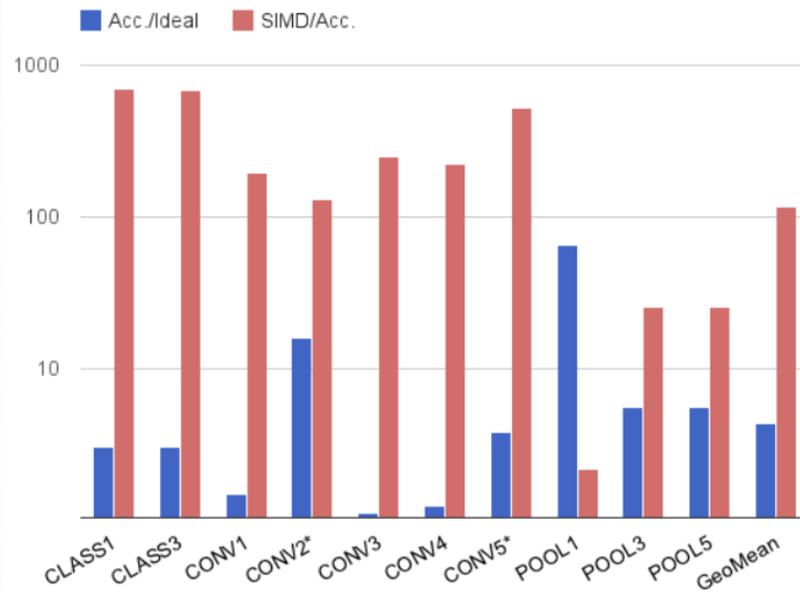
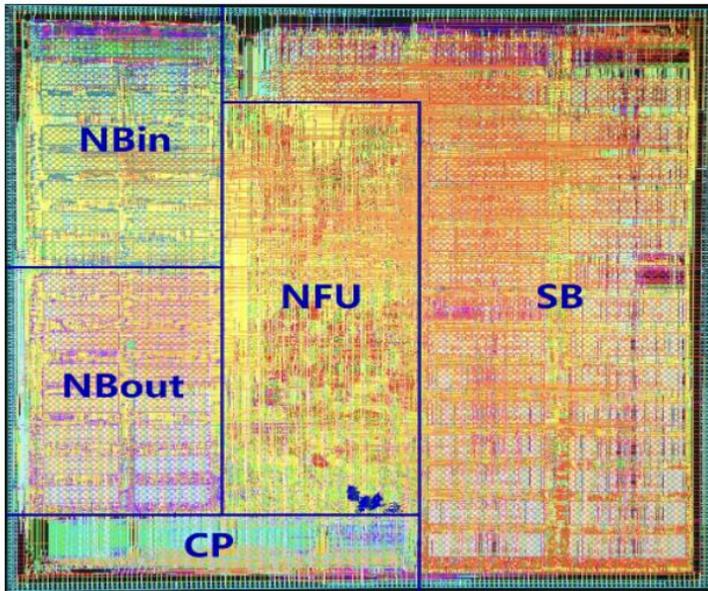
- 这种设计有没有什么问题?

```
for (int yy = 0; yy < Nyin; yy += Ty) {
  for (int xx = 0; xx < Nxin; xx += Tx) {
    for (int iii = 0; iii < Ni; iii += Tii)
      // — Original code — (excluding ii loop)
      int yout = 0;
      for (int y = yy; y < yy + Ty; y += sy) {
        int xout = 0;
        for (int x = xx; x < xx + Tx; x += sx) {
          for (int ii = iii; ii < iii + Tii; ii += Ti)
            for (int i = ii; i < ii + Ti; i++)
              value[i] = 0;
          for (int ky = 0; ky < Ky; ky++)
            for (int kx = 0; kx < Kx; kx++)
              for (int i = ii; i < ii + Ti; i++)
                // version with average pooling;
                value[i] += neuron[ky + y][kx + x][i];
                // version with max pooling;
                value[i] = max(value[i], neuron[ky + y][kx + x][i]);
              } } } }
          // for average pooling;
          neuron[xout][yout][i] = value[i] / (Kx * Ky);
          xout++; } yout++;
        } } } }
```

# AI加速器架构发展——DianNao

- 相较于CPU的SIMD计算，实现了显著的计算速度提升
- 但是，相较于理论加速比，仍然存在一定的差距，即硬件仍然面临欠利用问题
- 主要功耗来源于DRAM的访问

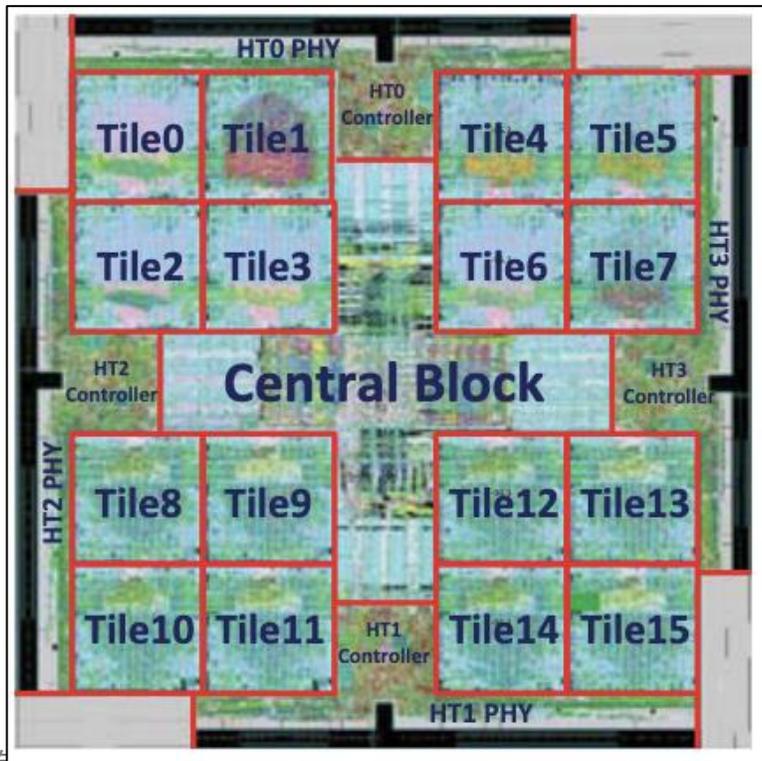
北京大学-智能硬件体系结构



# AI加速器架构发展——DaDianNao

- MICRO 2014最佳论文
- Multi-chip神经网络加速器：机器学习超级计算机

北京大学-智能硬件



## DaDianNao: A Machine-Learning Supercomputer

Yunji Chen<sup>1</sup>, Tao Luo<sup>1,3</sup>, Shaoli Liu<sup>1</sup>, Shijin Zhang<sup>1</sup>, Liqiang He<sup>2,4</sup>, Jia Wang<sup>1</sup>, Ling Li<sup>1</sup>,  
Tianshi Chen<sup>1</sup>, Zhiwei Xu<sup>1</sup>, Ninghui Sun<sup>1</sup>, Olivier Temam<sup>2</sup>

<sup>1</sup> SKL of Computer Architecture, ICT, CAS, China

<sup>2</sup> Inria, Scalay, France

<sup>3</sup> University of CAS, China

<sup>4</sup> Inner Mongolia University, China

**Abstract**—Many companies are deploying services, either for consumers or industry, which are largely based on machine-learning algorithms for sophisticated processing of large amounts of data. The state-of-the-art and most popular such machine-learning algorithms are Convolutional and Deep Neural Networks (CNNs and DNNs), which are known to be both computationally and memory intensive. A number of neural network accelerators have been recently proposed which can offer high computational capacity/area ratio, but which remain hampered by memory accesses.

However, unlike the memory wall faced by processors on general-purpose workloads, the CNNs and DNNs memory footprint, while large, is not beyond the capability of the on-chip storage of a multi-chip system. This property, combined with the CNN/DNN algorithmic characteristics, can lead to high internal bandwidth and low external communications, which can in turn enable high-degree parallelism at a reasonable area cost. In this article, we introduce a custom multi-chip machine-learning architecture along those lines. We show that, on a subset of the largest known neural network layers, it is possible to achieve a speedup of 450.65x over a GPU, and reduce the energy by 150.31x on average for a 64-chip system. We implement the node down to the place and route at 28nm, containing a combination of custom storage and computational units, with industry-grade interconnects.

### I. INTRODUCTION

Machine-Learning algorithms have become ubiquitous in a very broad range of applications and cloud services; examples include speech recognition, e.g., Siri or Google Now, click-through prediction for placing ads [27], face identification in Apple iPhoto or Google Picasa, robotics [20], pharmaceutical research [9], and so on. It is probably not exaggerated to say that machine-learning applications are in the process of displacing scientific computing as the major driver for high-performance computing. Early symptoms of this transformation are Intel calling for a refocus on Recognition, Mining and Synthesis applications in 2005 [14] (which later led to the PARSEC benchmark suite [3]), with Recognition and Mining largely corresponding to machine-learning tasks, or IBM developing the Watson supercomputer, illustrated with the Jeopardy game in 2011 [19].

Remarkably enough, at the same time this profound shift in applications is occurring, two simultaneous, albeit apparently unrelated, transformations are occurring in the machine-learning and in the hardware domains. Our community is well aware of the trend towards heterogeneous computing where architecture specialization is seen as a promising path to achieve high performance at low energy [21], provided we can find ways to reconcile architecture specialization and flexibility. At the same time, the machine-learning domain has profoundly evolved since 2006, where a category of algorithms, called Deep Learning (Convolutional and Deep Neural Networks), has emerged as state-of-the-art across a broad range of applications [33], [28], [32], [34]. In other words, at the time where architects need to find a good tradeoff between flexibility and efficiency, it turns out that just one category of algorithms can be used to implement a broad range of applications. In other words, there is a fairly unique opportunity to design highly specialized, and thus highly efficient, hardware which will benefit many of these emerging high-performance applications.

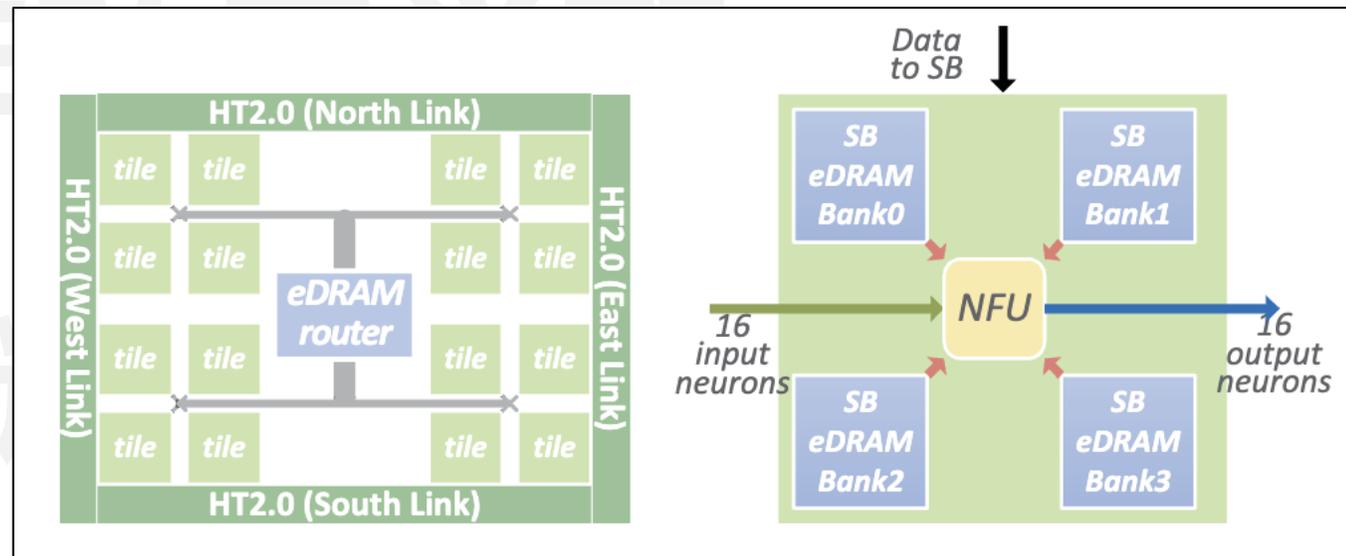
A few research groups have started to take advantage of this special context to design accelerators meant to be integrated into heterogeneous multi-cores. Temam [47] proposed a neural network accelerator for multi-layer perceptrons, though it is not a deep learning neural network, Esmailzadeh et al. [16] propose to use a hardware neural network called NPU for approximating any program function, though not specifically for machine-learning applications, Chen et al. [5] proposed an accelerator for Deep Learning (CNNs and DNNs). However, all these accelerators have significant neural network size limitations: either small neural networks of a few tens of neurons can be executed, or the neurons and synapses (i.e., weights of connections between neurons) intermediate values have to be stored in main memory. These two limitations are severe, respectively from a machine-learning or a hardware perspective.

From a machine-learning perspective, there is a significant trend towards increasingly large neural networks. The recent work of Krizhevsky et al. [32] achieved state-of-the-art accuracy on the ImageNet database [13] with “only” 60

- DaDianNao的设计特点:
  - 分块 (tile) 设计, 模型权重存储靠近计算单元, 采用4-bank eDRAM进行存储
  - 采用胖树的路由设计, 进行输入和输出激活值的广播
  - 芯片中间设置2个eDRAM bank, 进行输入和输出激活值的存储
  - 主要针对MLP和CNN设计, 同时支持训练和推理

Parameters	Settings	Parameters	Settings
Frequency	606MHz	tile eDRAM latency	~3 cycles
# of tiles	16	central eDRAM size	4MB
# of 16-bit multipliers/tile	256+32	central eDRAM latency	~10 cycles
# of 16-bit adders/tile	256+32	Link bandwidth	6.4x4GB/s
tile eDRAM size/tile	2MB	Link latency	80ns

Table III: Architecture characteristics.



- 2016年, MIT团队相继在HPCA、ISSCC、ISSC发表论文Eyeriss, 对于学术界产生了重要影响
- Eyeriss在对于不同数据流进行系统比较的基础上, 提出了新的row stationary数据流, 并流片验证

## Eyeriss: A Spatial Architecture for Energy-Efficient Dataflow for Convolutional Neural Networks

Yu-Hsin Chen<sup>1</sup>, Joel Emer<sup>1\*</sup> and Vivienne Sze<sup>2</sup>

<sup>1</sup>EECS, MIT  
Cambridge, MA 02139

<sup>2</sup>NVIDIA Research, NVIDIA  
Westford, MA 01886

\*{yhchen, jsemer, sze}@mit.edu

**Abstract**—Deep convolutional neural networks (CNNs) are widely used in modern AI systems for their superior accuracy but at the cost of high computational complexity. The complexity comes from the need to simultaneously process hundreds of filters and channels in the high-dimensional convolutions, which involve a significant amount of data movement. Although highly-parallel compute paradigms, such as SIMD/SIMT, effectively address the computation requirement to achieve high throughput, energy consumption still remains high as data movement can be more expensive than computation. Accordingly, finding a dataflow that supports parallel processing with minimal data movement cost is crucial to achieving energy-efficient CNN processing without compromising accuracy.

In this paper, we present a novel dataflow, called *row-stationary* (RS), that minimizes data movement energy consumption on a spatial architecture. This is realized by exploiting local data reuse of filter weights and feature map pixels, i.e., activations, in the high-dimensional convolutions, and minimizing data movement of partial sum accumulations. Unlike dataflows used in existing designs, which only reduce certain types of data movement, the proposed RS dataflow can adapt to different CNN shape configurations and reduces all types of data movement through maximally utilizing the processing engine (PE) local storage, direct inter-PE communication and spatial parallelism. To evaluate the energy efficiency of the different dataflows, we propose an analysis framework that compares energy cost under the same hardware area and processing parallelism constraints. Experiments using the CNN configurations of AlexNet show that the proposed RS dataflow is more energy efficient than existing dataflows in both convolutional (1.4x to 2.5x) and fully-connected layers (at least 1.3x for batch size larger than 16). The RS dataflow has also been demonstrated on a fabricated chip, which verifies our energy analysis.

### I. INTRODUCTION

The recent popularity of deep learning [1], specifically deep convolutional neural networks (CNNs), can be attributed to its ability to achieve unprecedented accuracy for tasks ranging from object recognition [2–5] and detection [6, 7] to scene understanding [8]. These state-of-the-art CNNs [2–5] are orders of magnitude larger than those used in the 1990s [9], requiring up to hundreds of megabytes for filter weight storage and 30k–600k operations per input pixel.

The large size of such networks poses both throughput and energy efficiency challenges to the underlying processing hardware. Convolutions account for over 90% of the CNN operations and dominates runtime [10]. Although these operations can leverage highly-parallel compute paradigms, such as SIMD/SIMT, throughput may not scale accordingly due to the accompanying bandwidth requirement, and the energy consumption remains high as data movement can be more expensive than computation [11–13]. In order to achieve energy-efficient CNN processing without compromising throughput, we need to develop dataflows that support parallel processing with minimal data movement. The differences in data movement energy cost based on where the data is stored also needs to be accounted for. For instance, fetching data from off-chip DRAMs costs orders of magnitude more energy than from on-chip storage [11, 12].

Many previous papers have proposed specialized CNN dataflows on various platforms, including GPU [14], FPGA [15–21], and ASIC [22–26]. However, due to differences in technology, hardware resources and system setup, a direct comparison between different implementations does not provide much insight into the relative energy efficiency of different dataflows. In this paper, we evaluate the energy efficiency of various CNN dataflows on a spatial architecture under the same hardware resource constraints, i.e., area, processing parallelism and technology. Based on this evaluation, we will propose a novel dataflow that maximizes energy efficiency for CNN acceleration.

To evaluate energy consumption, we categorize the data movements in a spatial architecture into several levels of hierarchy according to their energy cost, and then analyze each dataflow to assess the data movement at each level. This analysis framework provides insights into how each dataflow exploits different types of data movement using various architecture resources. It also offers a quantifiable way to examine the differences in energy efficiency between different dataflows.

Previously proposed dataflows typically optimize a certain type of data movement, such as input data reuse or partial

## ISSCC 2016 / SESSION 14 / NEXT-GENERATION PROCESSING / 14.5

### 14.5 Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks

Yu-Hsin Chen<sup>1</sup>, Tushar Krishna<sup>1</sup>, Joel Emer<sup>1\*</sup>, Vivienne Sze<sup>2</sup>

<sup>1</sup>Massachusetts Institute of Technology, Cambridge, MA, <sup>2</sup>Nvidia, Westford, MA

Deep learning using convolutional neural networks (CNN) gives state-of-the-art accuracy on many computer vision tasks (e.g. object detection, recognition, segmentation). Convolutions account for over 90% of the processing in CNNs for both inference/testing and training, and fully convolutional networks are increasingly being used. To achieve state-of-the-art accuracy requires CNNs with not only a larger number of layers, but also millions of filter weights, and varying shapes (i.e. filter sizes, number of filters, number of channels) as shown in Fig. 14.5.1. For instance, AlexNet [1] uses 2.3 million weights (4.6MB of storage) and requires 666 million MACs per 227x227 image (136MACs/pixel). VGG16 [2] uses 14.7 million weights (28.4MB of storage) and requires 15.3 billion MACs per 224x224 image (306MACs/pixel). The large number of filter weights and channels results in substantial data movement, which consumes significant energy.

Existing accelerators do not support the configurability necessary to efficiently support large CNNs with different shapes [2], and using mobile GPUs can be expensive [3]. This paper describes an accelerator that can deliver state-of-the-art accuracy with minimum energy consumption in the system (including DRAM) in real-time, by using two key methods: (1) efficient dataflow and supporting hardware (spatial array, memory hierarchy and on-chip network) that minimize data movement by exploiting data reuse and support different shapes; (2) exploit data statistics to minimize energy through zeros skipping/gating to avoid unnecessary reads and computations; and data compression to reduce off-chip memory bandwidth, which is the most expensive data movement.

Figure 14.5.2 shows the top-level architecture and memory hierarchy of the accelerator. Data movement is optimized by buffering input image data (img), filter weights (Filt) and partial sums (Psum) in a shared 108KB SRAM buffer, which facilitates temporal reuse of loaded data. Image data and filter weights are read from DRAM to the buffer and streamed into the spatial computation array allowing for overlap of memory traffic and computation. The streaming and reuse allows the system to achieve high computational efficiency even when running the memory link at a lower clock frequency than the spatial array. The spatial array computes inner products between the image and filter weights, generating partial sums that are returned from the array to the buffer and then, optionally reduced (ReLU) and compressed, to the DRAM. Run-length-based compression reduces the average image bandwidth by 2x. Configurable support for image and filter sizes that do not fit completely into the spatial array is achieved by saving partial sums in the buffer and later restoring them to the spatial array. The sizes of the spatial array and buffer determine the number of such “passes” needed to do the calculations for a specific layer. Unused PEs are clock gated.

Figure 14.5.3 shows the dataflow within the array for filter weights, image values and partial sums. If the filter height (R) equals the number of rows in the array (in our case 12), the logical dataflow would be as follows: (1) filter weights are read from the buffer into the left column of the array (one filter row per PE) and the filter weights move from left to right within the array; (2) image values are fed into the left column and bottom row of the array (one image row per PE) and image values move up diagonally; (3) partial sums for each output row move up vertically, and can be read out of the top row at the end of the computational pass. If the partial sums are used in the next pass, they are fed into the bottom row of the array from the buffer at the beginning of the next computational pass.

In order to maximize utilization of a fixed-size array for different shapes, the mapping may require either folding or replication if the shape size is larger or smaller than the array dimension, respectively. Replication results in increased throughput as compared to the purely logical dataflow described above. Cases II, III, IV, and V in Fig. 14.5.3 illustrate the replication and folding of image values for various layers of AlexNet. The same data values are used in the spatial array. Across the six example cases, which include physical mapping of filter weights, image values and partial sums onto the fixed-size spatial array, we see the logical dataflow patterns translating to myriad physical dataflow patterns that need to be

supported. Furthermore, the same data value is often needed by multiple PEs, whose physical location in the array depends on the data type (filter, image or partial sum) and layer.

Since different layers have different shapes and hence different mappings, a design-time fixed interconnect topology will not work. Every PE can potentially be a destination for a piece of data in some particular configuration, and so a Network-on-Chip (NoC) is needed to support address-based data delivery. However, traditional NoC designs with switches at every PE to buffer/forward data to one or multiple targets would result in multi-cycle delays. A full-chip broadcast to every PE could work, but would consume enormous power.

To optimize data movement, it is important to exploit spatial reuse, where a single buffer read can be used by multiple PEs (i.e. multicast). Fig. 14.5.4 shows our NoC that supports configurable data patterns, and provides an energy-efficient multicast to a variable number of PEs within a single cycle. The NoC comprises one Global Y bus, and 12 Global X buses (one per row). Each PE is configured with a (row, col) ID at the beginning of processing via a scan chain. Multicast to any subset of PEs is achieved by addressing the same ID to multiple PEs. Data from the buffer is tagged with the target PEs’ (row, col) ID, and multicast controllers at the input of each X bus and each PE deliver data only to those X buses and PEs, respectively, that match the target ID to avoid unnecessary switching. Data is sent on the buses only if all target PEs are ready (i.e. have an energy buffer) to receive. To support high bandwidth, we use separate input NoCs for filter, image, and partial sums. The partial sum NoC has a separate set of output links to the buffer to write the final partial sums. The NoC data delivery for four of the cases from Fig. 14.5.3 is shown in Fig. 14.5.4.

Each processing engine, shown in Fig. 14.5.5, is a three-stage pipeline responsible for calculating the inner product of the input image and filter weights for a single row of the filter. The sequence of partial sums for the sliding filter window is computed sequentially. The partial sums for the row are passed on a local link to the neighboring PE (see Fig. 14.5.4), where the cross-row partial sums are computed. Local scratch pads allow for temporary energy-efficient storage of input image and filter weights by recirculating values needed by different windows. A partial sum scratch pad allows for temporal reuse of partial sums being generated for different images and/or channels and filters. Data gating is achieved by recording the input image values of zero in a “zero buffer” and skipping filter reads and computation for those values resulting in a 45% power savings in the PE.

The test chip is implemented in 65nm CMOS. It operates at 200MHz core clock and 60MHz link clock, which results in a frame rate of 34.7fps on the five convolutional layers in AlexNet and a measured power of 278mW at 1V. The PE array, NoC and on-chip buffer consume 77.8%, 15.6% and 2.7% of the total power, respectively. The core and link clocks can scale up to 250MHz and 90MHz, respectively. This enables us to achieve a throughput of 44.8fps at 1.17V. Fig. 14.5.6 shows the performance at each layer, including compression ratio, power consumption, PE utilization, and memory access to highlight the reduction in DRAM bandwidth, efficiency of the reconfigurable mapping and reduced data access due to data reuse, respectively. A die photo of the chip and the range of the shapes it can support natively are shown in Fig. 14.5.7.

### Acknowledgments

This work is funded by the DARPA YFA grant N66011-14-1-4039, MIT Center for Integrated Circuits & Systems, and a gift from Intel. The authors would also like to thank Mehul Tiwari and Michael Price for their technical assistance.

### References

- [1] A. Krizhevsky, I. Sutskever, G. E. Hinton, “Imagenet Classification with Deep Convolutional Neural Networks,” *Neural Information Processing Systems*, pp. 1097–1105, 2012.
- [2] K. Simonyan, A. Zisserman, “Very Deep Convolutional Networks for Large-Scale Image Recognition,” *CoRR*, abs/1409.1556, 2014.
- [3] S. Park et al., “A 1.93TOPS/W Scalable Deep Learning/Inference Processor with Tetra-parallel MIMD Architecture for Big Data Applications,” *ISSCC Dig. Tech. Papers*, pp. 80–81, 2015.
- [4] S. Chellur et al., “cuDNN: Efficient Primitives for Deep Learning,” *CoRR*, abs/1410.0759, 2014.

This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.  
IEEE JOURNAL OF SOLID-STATE CIRCUITS

## Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks

Yu-Hsin Chen, Student Member, IEEE, Tushar Krishna, Member, IEEE, Joel S. Emer, Fellow, IEEE, and Vivienne Sze, Senior Member, IEEE

**Abstract**—Eyeriss is an accelerator for state-of-the-art deep convolutional neural networks (CNNs). It optimizes for the energy efficiency of the entire system, including the accelerator chip and off-chip DRAM, for various CNN shapes by reconfiguring the architecture. CNNs are widely used in modern AI systems but also bring challenges on throughput and energy efficiency to the underlying hardware. This is because its computation requires a large amount of data, creating significant data movement from on-chip and off-chip that is more energy-consuming than computation. Minimizing data movement energy cost for any CNN shape, therefore, is the key to high throughput and energy efficiency. Eyeriss achieves these goals by using a proposed processing dataflow, called row stationary (RS), on a spatial architecture with 168 processing elements. RS dataflow reconfigures the computation mapping of a given shape, which optimizes energy efficiency by maximally reusing data locally to reduce expensive data movement, such as DRAM accesses. Compression and data gating are also applied to further improve energy efficiency. Eyeriss processes the convolutional layers at 35 frames/s and 0.0029 DRAM access/MAC/s and accumulation (MAC) for AlexNet at 278 mW (batch size  $N = 4$ ), and 0.7 frames/s and 0.0035 DRAM access/MAC for VGG-16 at 236 mW ( $N = 3$ ).

**Index Terms**—Convolutional neural networks (CNNs), dataflow processing, deep learning, energy-efficient accelerators, spatial architecture.

### I. INTRODUCTION

DEEP learning using convolutional neural networks (CNNs) [1] has achieved unprecedented accuracy on many modern AI applications [2]–[3]. However, state-of-the-art CNNs require tens to hundreds of megabytes of parameters on billions of operations in a single inference pass, creating significant data movement from on-chip and off-chip to support the computation. Since data movement can be more energy-consuming than computation [10], [11], the

Manuscript received May 5, 2016; revised July 31, 2016; accepted September 28, 2016. This paper was approved by Associate Editor Dejan Markovic.  
Y.-H. Chen and V. Sze are with the Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139 USA.

T. Krishna was with the Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139 USA. He is now with the School of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA 30332 USA.

J. S. Emer is with the Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139 USA, and also with Nvidia Corporation, Westford, MA 01886 USA.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/JSSC.2016.2616357

0018-9200 © 2016 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See [http://www.ieee.org/publications\\_standards/publications/rights/index.html](http://www.ieee.org/publications_standards/publications/rights/index.html) for more information.

processing of CNNs has to not only provide high parallelism for high throughput but also optimize for the data movement of the entire system in order to achieve high energy efficiency. In addition, this optimization needs to adapt to the varying shapes of the high-dimensional convolutions in CNN.

To address these challenges, it is crucial to design a compute scheme, called a *dataflow*, that can support a highly parallel compute paradigm while optimizing the energy cost of data movement from both on-chip and off-chip. The cost of data movement is reduced by exploiting data reuse in a multilevel memory hierarchy, and the hardware needs to be reconfigurable to support different shapes. To further improve energy efficiency, data statistics can also be exploited. Specifically, CNN data contains many zeros. Techniques such as compression and data adaptive processing can be applied to save both memory bandwidth and processing power.

Previous work has proposed hardware designs for CNN acceleration [12]–[22]. However, most of them only have simulation results that are not verified by the measured results from fabricated chips; implementations using FPGA also do not reveal the actual throughput and energy efficiency of the architecture. A few efforts have demonstrated the measurement results of fabricated chips [23]–[25]. However, these works do not benchmark their implementations using widely used publicly available state-of-the-art CNNs, which is critical to the hardware evaluation. Specifically, Park et al. [23] propose a deep-learning processor for running both training and inference using a MIMD architecture, which was tested on a custom four-layer network using 5 × 5 filters. Cavigelli et al. [24] present a CNN accelerator for inference that is tested on the four-layer CNN using 7 × 7 filters. Sim et al. [25] demonstrate a CNN processor and only report the theoretical peak throughput along with the power measured on a CNN for the MNIST data set [26], which has storage and computation requirements that are orders of magnitude lower than the state-of-the-art CNNs. With the exception of [24], these works did not report the required DRAM bandwidth for the proposed compute schemes. It is not sufficient to look at the processor power alone, since DRAM access is one of the most important factors dictating the system energy efficiency.

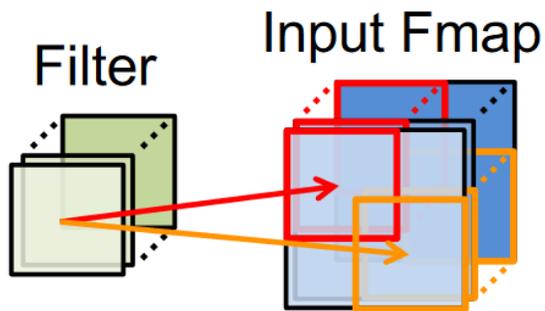
In this paper, we have implemented and fabricated a CNN accelerator, called Eyeriss, that can support high throughput CNN inference and optimizes for the energy efficiency of the entire system, including the accelerator chip and off-chip DRAM. It is also reconfigurable to handle different

- 卷积计算中的**数据复用模式**

- 通过**权重和激活值的复用**，可以显著降低访存需求，提升计算效率

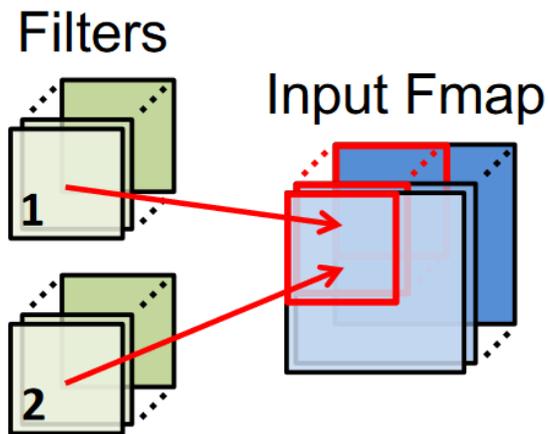
## 北京大学-智能硬件体系结构

卷积层



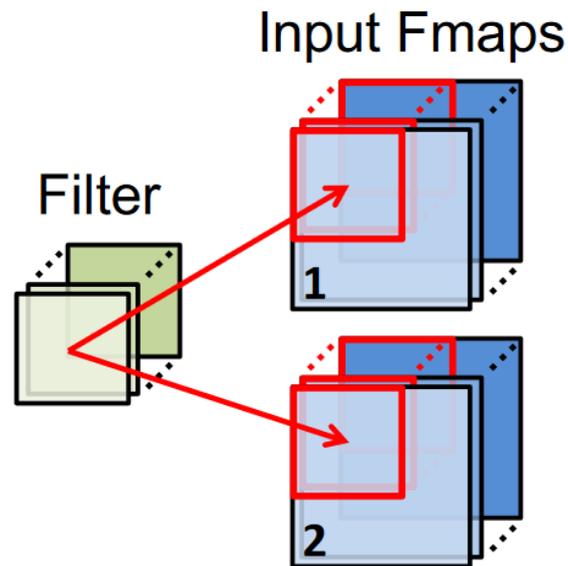
激活值和权重值的复用

卷积层和全连接层



激活值的复用

卷积层和全连接层

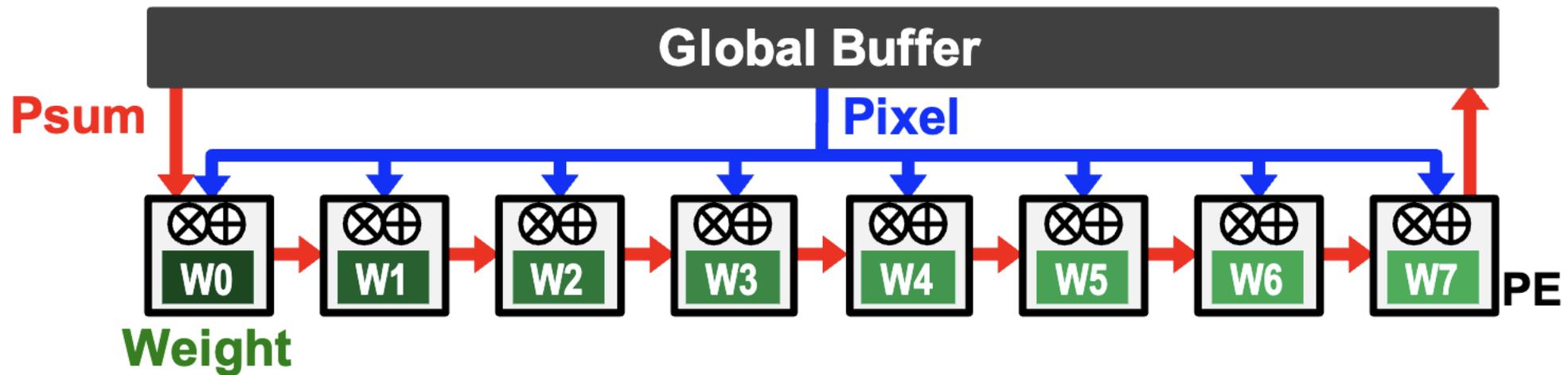


权重值的复用

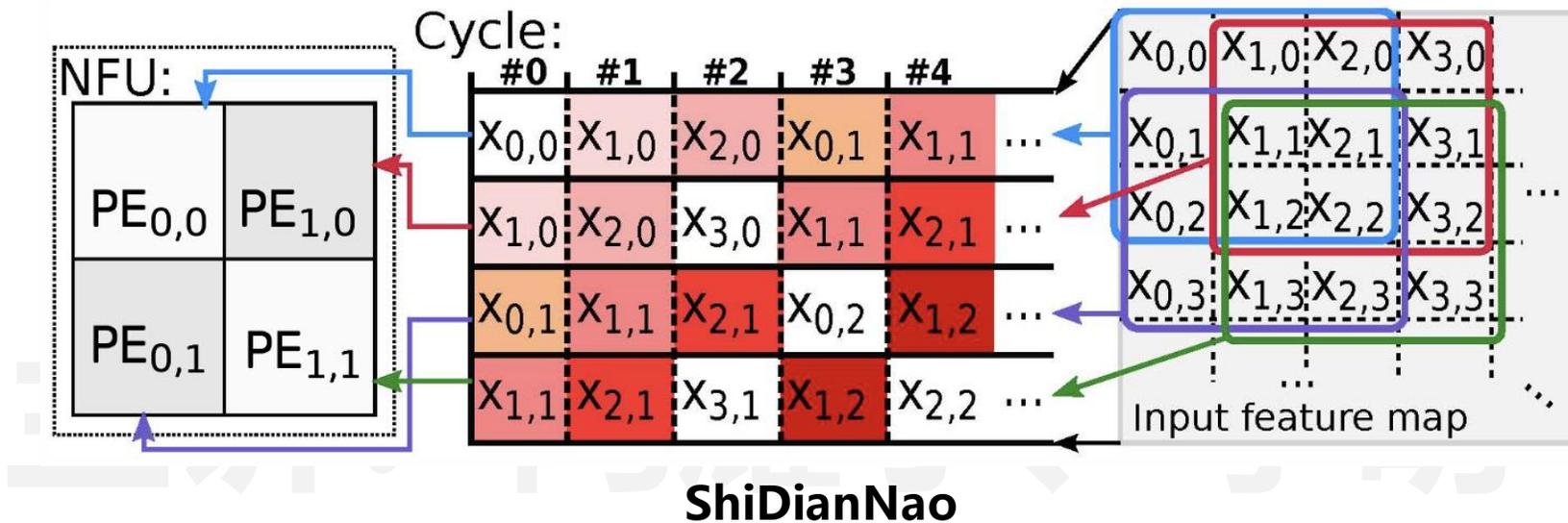
- Eyeriss把现有数据流加速器分为3类:

- **权重静态**, 代表工作包括TPU等, 运算过程中, 权重值被预先存储在PE中, 最大化权重复用

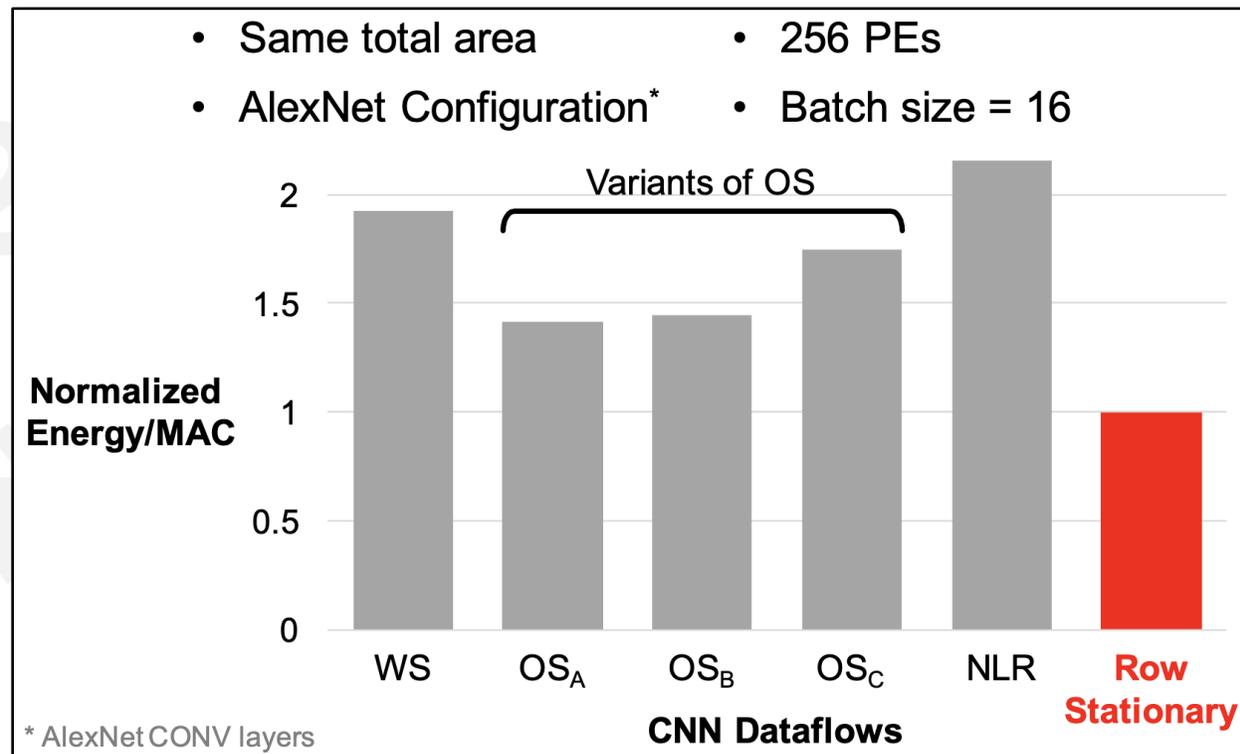
北京大学-智能硬件体系结构



- Eyeriss把现有数据流加速器分为3类：
  - **权重静态**，代表工作包括TPU等，权重值被预先存储在PE中，最大化权重复用
  - **输出静态**，代表工作包括ShiDianNao等，输出值停留在PE中被累加，最小化输出的读写开销



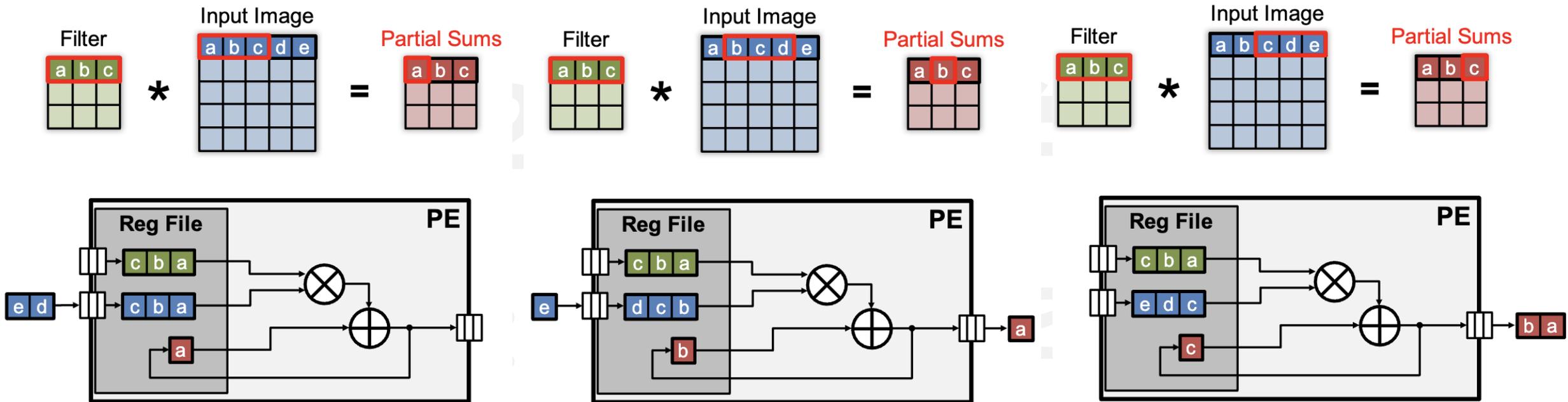
- Eyeriss把现有数据流加速器分为3类：
  - **权重静态**，代表工作包括TPU等，权重值被预先存储在PE中，最大化权重复用
  - **输出静态**，代表工作包括ShiDianNao等，输出值停留在PE中被累加，最小化输出的读写开销
  - **No local reuse**，代表工作包括DianNao、DaDianNao等



# AI加速器架构发展——Eyeriss

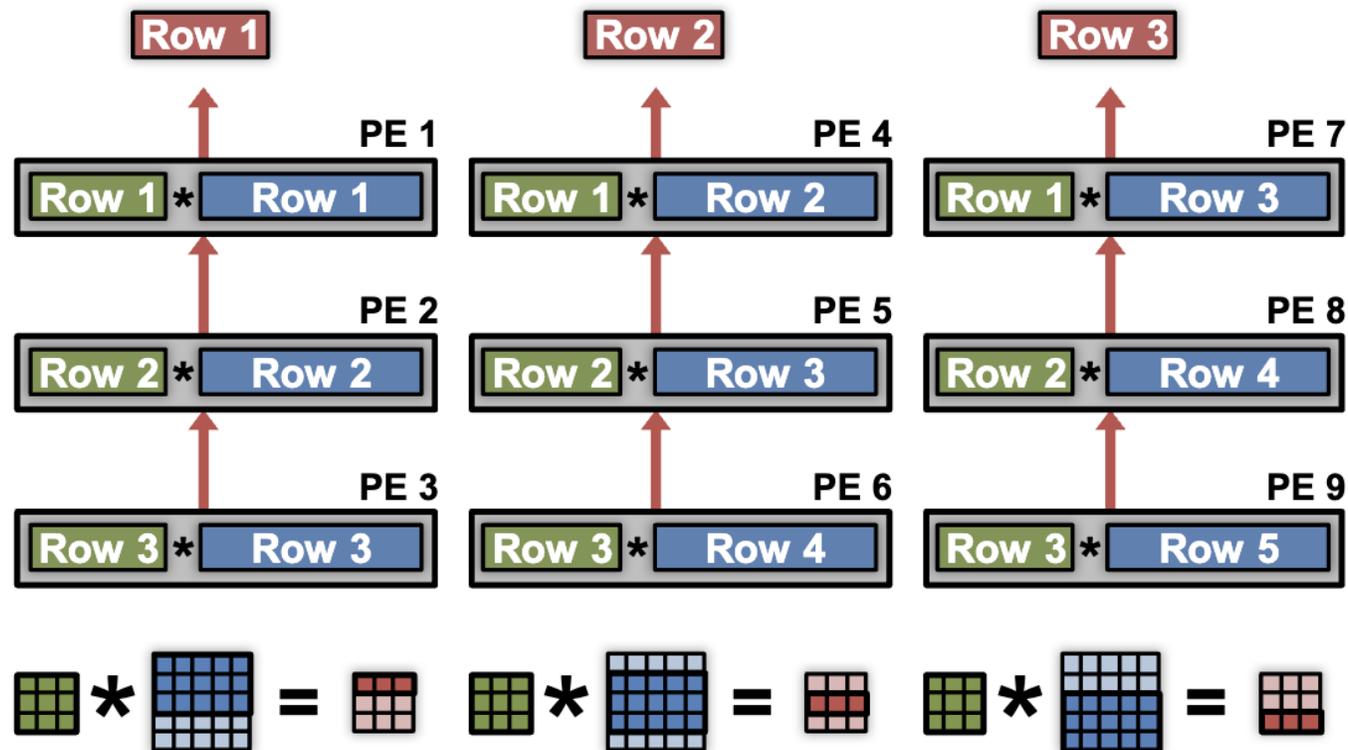
- 基于数据流分类，Eyeriss提出了row stationary的数据流
  - 计算单元内部，权重、输入和输出激活值均被复用

北京大学-智能硬件体系结构



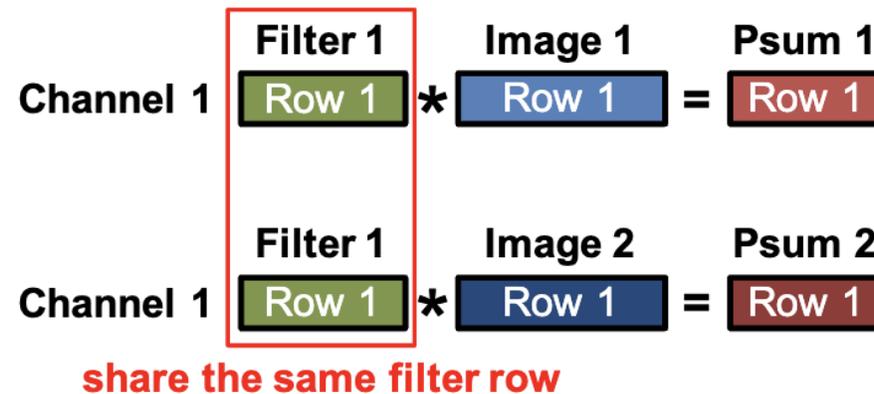
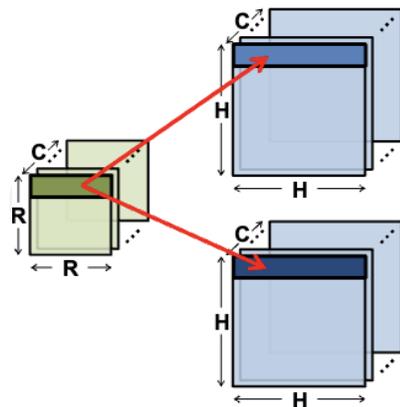
# AI加速器架构发展——Eyeriss

- 基于数据流分类，Eyeriss提出了row stationary的数据流
  - 计算单元内部，权重、输入和输出激活值均被复用
  - 不同计算单元间，进一步复用权重、输入和输出



# AI加速器架构发展——Eyeriss

- 基于数据流分类，Eyeriss提出了row stationary的数据流
  - 计算单元内部，权重、输入和输出激活值均被复用
  - 不同计算单元间，进一步复用权重、输入和输出
  - 进一步拓展到多输入图像

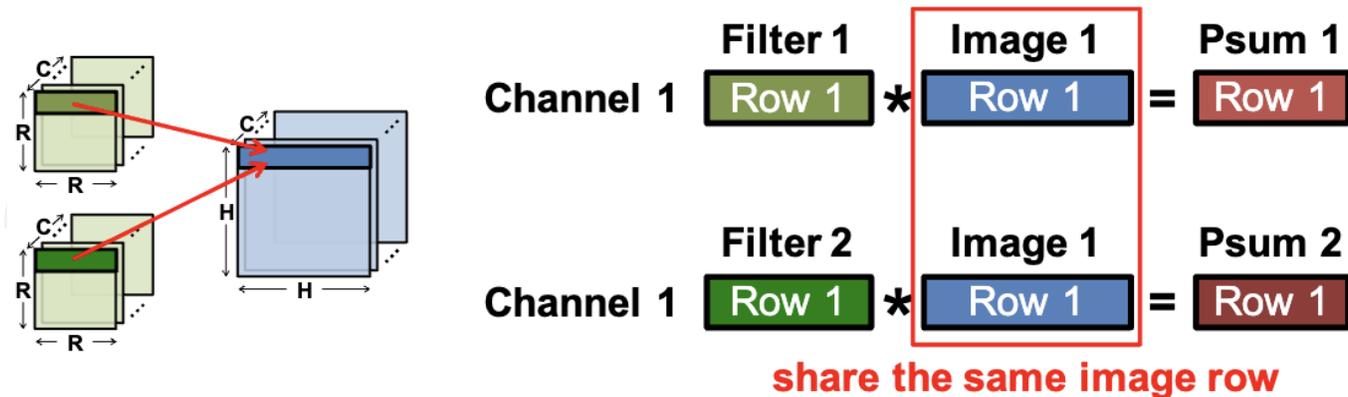


Processing in PE: concatenate image rows



# AI加速器架构发展——Eyeriss

- 基于数据流分类，Eyeriss提出了row stationary的数据流
  - 计算单元内部，权重、输入和输出激活值均被复用
  - 不同计算单元间，进一步复用权重、输入和输出
  - 进一步拓展到多输入图像、多卷积核

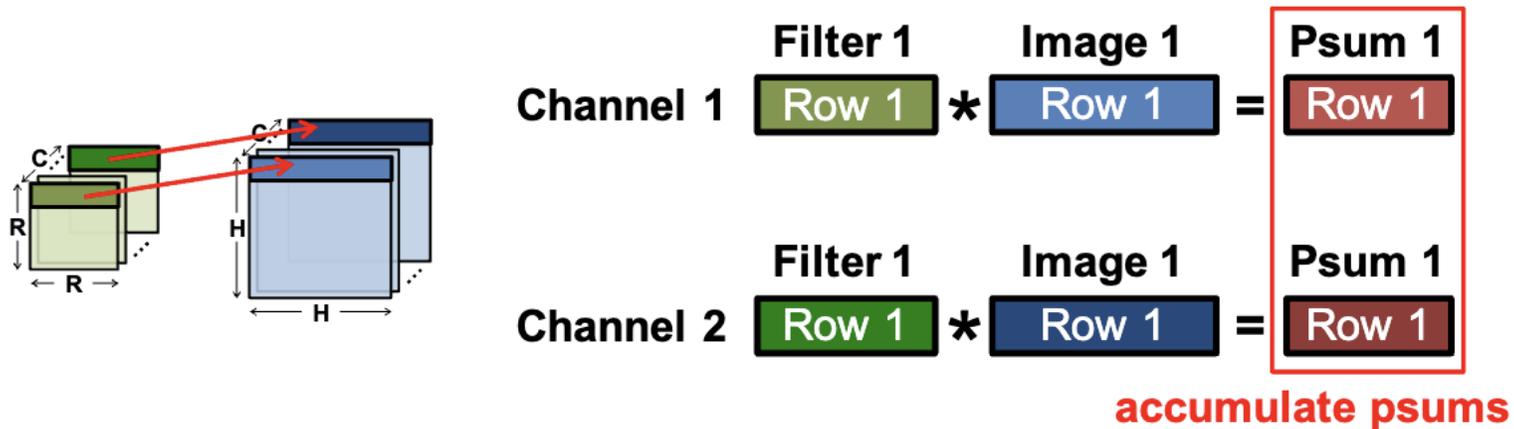


Processing in PE: interleave filter rows



# AI加速器架构发展——Eyeriss

- 基于数据流分类，Eyeriss提出了row stationary的数据流
  - 计算单元内部，权重、输入和输出激活值均被复用
  - 不同计算单元间，进一步复用权重、输入和输出
  - 进一步拓展到多输入图像、多卷积核，不同输入通道

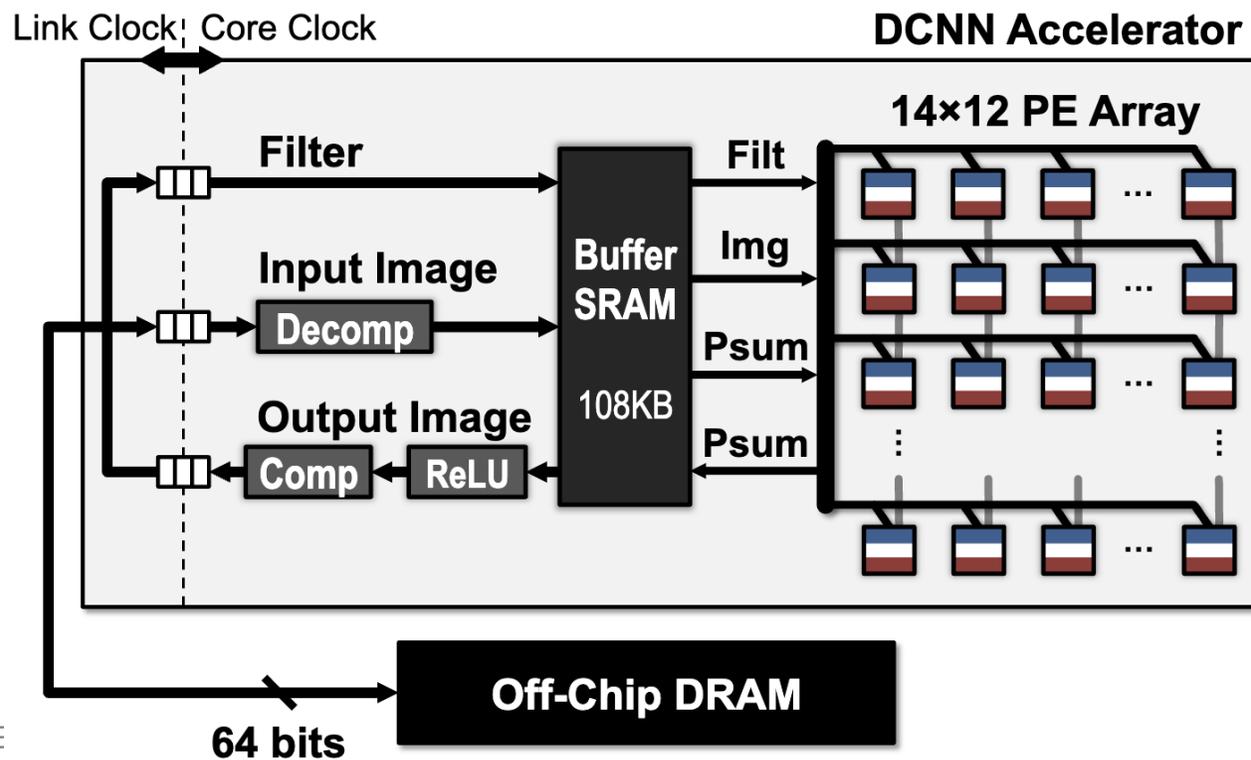


Processing in PE: interleave channels



## • Eyeriss的存储系统

- 计算单元内部需要交大的便签存储器，提升数据复用
- 采用了统一的便签存储器，容量很小，主要用于计算卷积层
- 插入了压缩和解压缩模块，能够利用输入和输出结果中的0元素

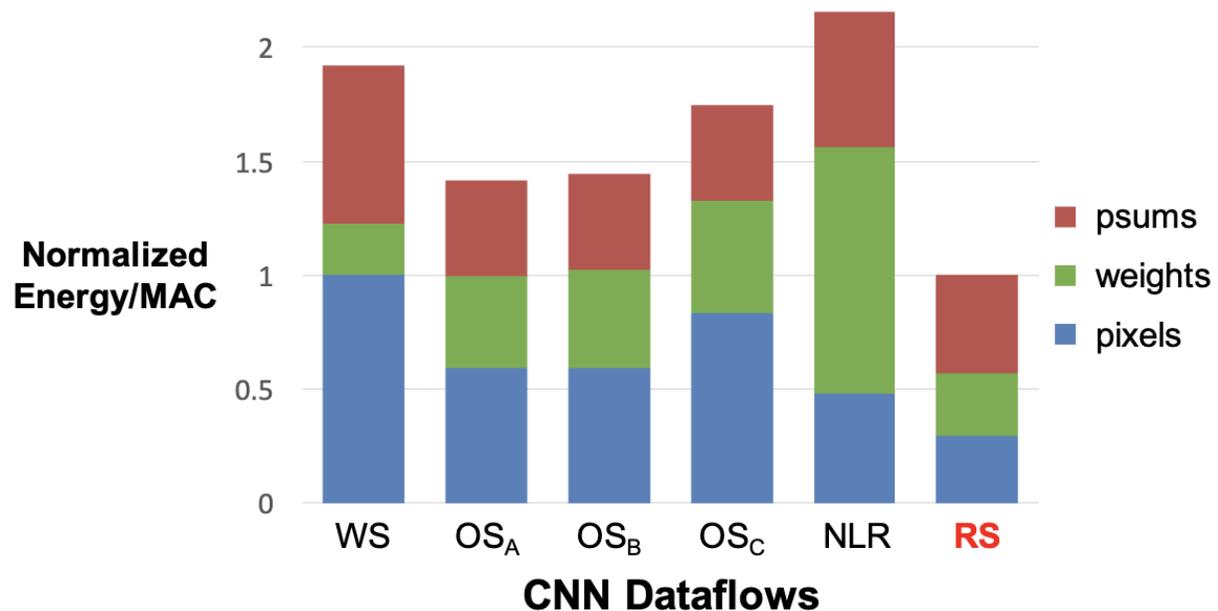
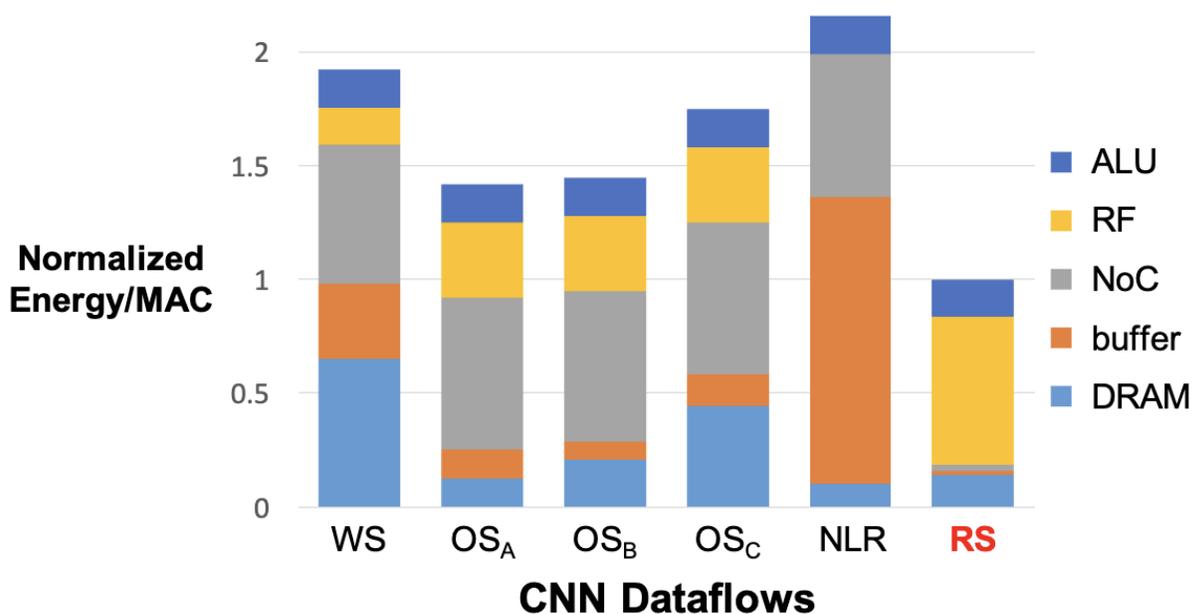


Technology	TSMC 65nm LP 1P9M
Core Area	3.5mm×3.5mm
Gate Count	1852 kGates (NAND2)
On-Chip Buffer	108 KB
# of PEs	168
Scratch Pad / PE	0.5 KB
Supply Voltage	0.82 – 1.17 V
Core Frequency	100 – 250 MHz
Peak Performance	33.6 – 84.0 GOPS (2 OP = 1 MAC)
Word Bit-width	16-bit Fixed-Point
Filter Size*	1 – 32 [width] 1 – 12 [height]
# of Filters*	1 – 1024
# of Channels*	1 – 1024
Stride Range	1–12 [horizontal] 1, 2, 4 [vertical]

# AI加速器架构发展——Eyeriss

- 通过将Eyeriss的功耗与其他数据流进行比较可以发现
  - RF的功耗显著提升（每位每个PE中的存储显著提升）
  - Buffer的功耗很小，因为进一步提升了输入、权重和输出数据的复用

北京大学-智能硬件体系结构



# AI加速器架构发展——Google TPU v1



- Google TPU论文发表于ISCA 2016，从2015年开始，TPU已经集成在了Google服务器中
- 只考虑神经网络推理加速，重点关注MLP、CNN、LSTM三类神经网络
- 相较于NV K80 GPU和Intel Haswell CPU, 15-30倍延迟降低、30-80倍能效降低

## In-Datcenter Performance Analysis of a Tensor Processing Unit™

Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Daniel Killebrew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Jonathan Ross, Matt Ross, Amir Salek, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Dan Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon

Google, Inc., Mountain View, CA USA

Email: {jouppi, cliffy, nishantpatil, davidpatterson}@google.com

To appear at the 44th International Symposium on Computer Architecture (ISCA), Toronto, Canada, June 26, 2017.

### Abstract

Many architects believe that major improvements in cost-energy-performance must now come from domain-specific hardware. This paper evaluates a custom ASIC—called a *Tensor Processing Unit (TPU)*—deployed in datacenters since 2015 that accelerates the inference phase of neural networks (NN). The heart of the TPU is a 65,536 8-bit MAC matrix multiply unit that offers a peak throughput of 92 TeraOps/second (TOPS) and a large (28 MIB) software-managed on-chip memory. The TPU’s deterministic execution model is a better match to the 99th-percentile response-time requirement of our NN applications than are the time-varying optimizations of CPUs and GPUs (caches, out-of-order execution, multithreading, multiprocessing, prefetching, ...) that help average throughput more than guaranteed latency. The lack of such features helps explain why, despite having myriad MACs and a big memory, the TPU is relatively small and low power. We compare the TPU to a server-class Intel Haswell CPU and an Nvidia K80 GPU, which are contemporaries deployed in the same datacenters. Our workload, written in the high-level TensorFlow framework, uses production NN applications (MLPs, CNNs, and LSTMs) that represent 95% of our datacenters’ NN inference demand. Despite low utilization for some applications, the TPU is on average about 15X - 30X faster than its contemporary GPU or CPU, with TOPS/Watt about 30X - 80X higher. Moreover, using the GPU’s GDDR5 memory in the TPU would triple achieved TOPS and raise TOPS/Watt to nearly 70X the GPU and 200X the CPU.

Index terms—DNN, MLP, CNN, RNN, LSTM, neural network, domain-specific architecture, accelerator

### 1. Introduction to Neural Networks

The synergy between the large data sets in the cloud and the numerous computers that power it has enabled a renaissance in machine learning. In particular, *deep neural networks* (DNNs) have led to breakthroughs such as reducing word error rates in speech recognition by 30% over traditional approaches, which was the biggest gain in 20 years [Dea16]; cutting the error rate in an image recognition competition since 2011 from 26% to 3.5% [Kri12] [Sze15] [He16]; and beating a human champion at Go [Sil16]. Unlike some hardware targets, DNNs are applicable to a wide range of problems, so we can reuse a DNN-specific ASIC for solutions in speech, vision, language, translation, search ranking, and many more.

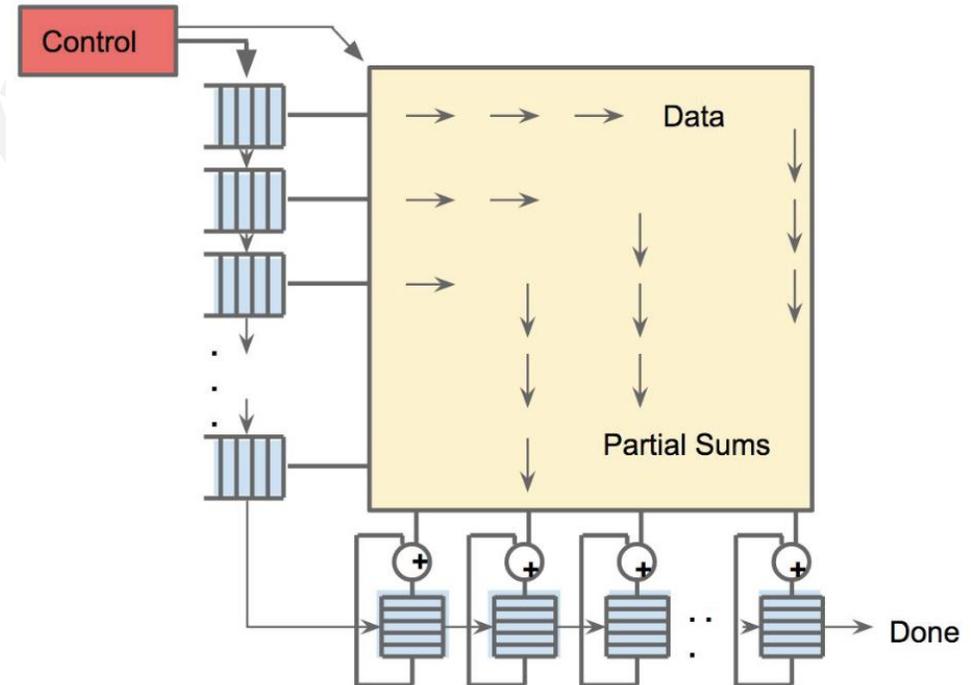
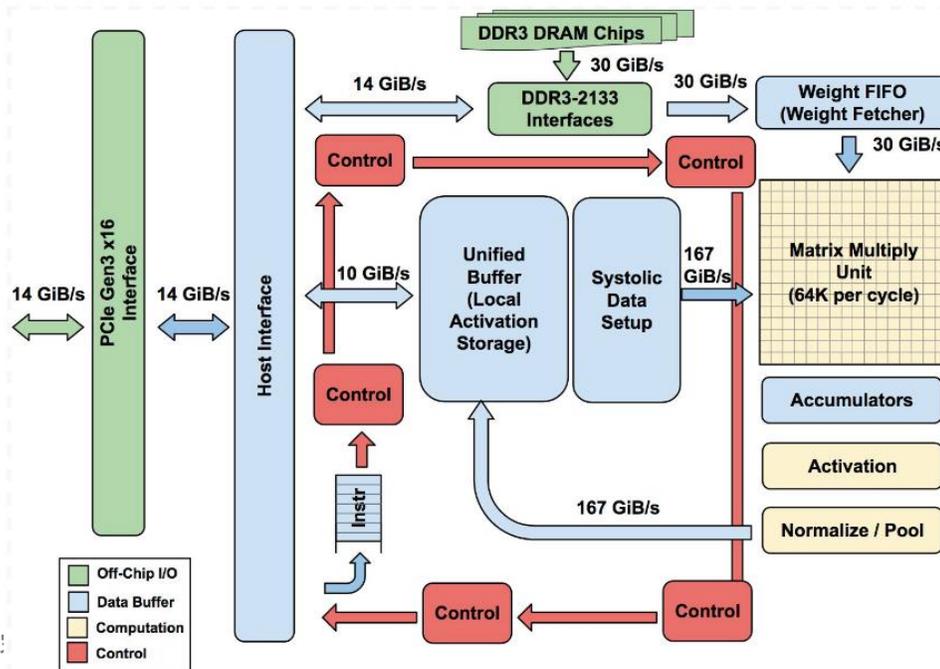
Neural networks (NN) target brain-like functionality and are based on a simple artificial neuron: a nonlinear function (such as  $\max(0, \text{value})$ ) of a weighted sum of the inputs. These artificial neurons are collected into layers, with the outputs of one layer becoming the inputs of the next one in the sequence. The “deep” part of DNN comes from going beyond a few layers, as the large data sets in the cloud allowed more accurate models to be built by using extra and larger layers to capture higher levels of patterns or concepts, and GPUs provided enough computing to develop them.

The two phases of NN are called *training* (or learning) and *inference* (or prediction), and they refer to development versus production. The developer chooses the number of layers and the type of NN, and training determines the weights. Virtually all training today is in floating point, which is one reason GPUs have been so popular. A step called *quantization* transforms floating-point numbers into narrow integers—often just 8 bits—which are usually good enough for inference. Eight-bit integer multiplies can be 6X less energy and 6X less area than IEEE 754 16-bit floating-point multiplies, and the

Name	Layers					Nonlinear function	Weights	TPU Ops / Weight Byte	TPU Batch Size	% of Deployed TPUs in July 2016
	FC	Conv	Vector	Pool	Total					
MLP0	5				5	ReLU	20M	200	200	61%
MLP1	4				4	ReLU	5M	168	168	
LSTM0	24		34		58	sigmoid, tanh	52M	64	64	29%
LSTM1	37		19		56	sigmoid, tanh	34M	96	96	
CNN0		16			16	ReLU	8M	2888	8	5%
CNN1	4	72		13	89	ReLU	100M	1750	32	

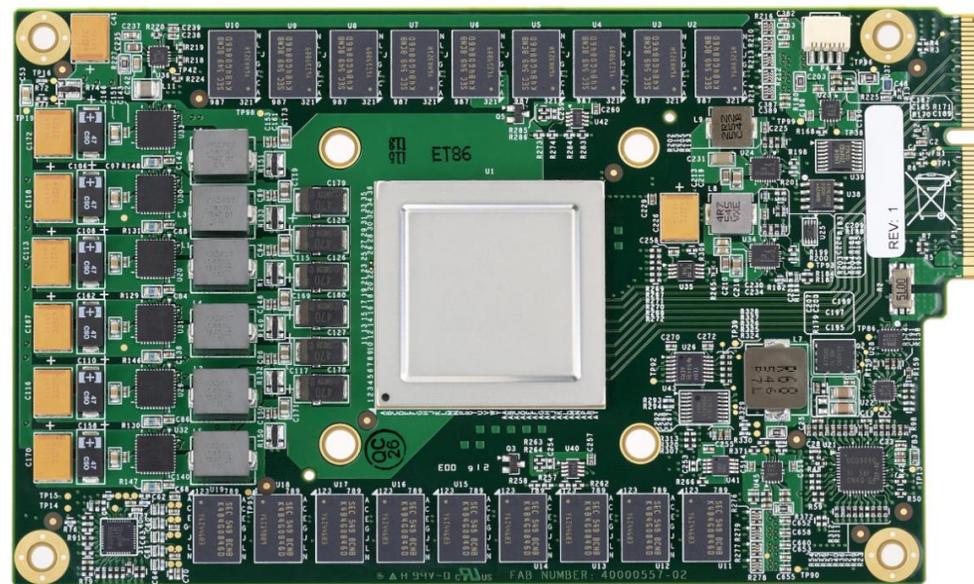
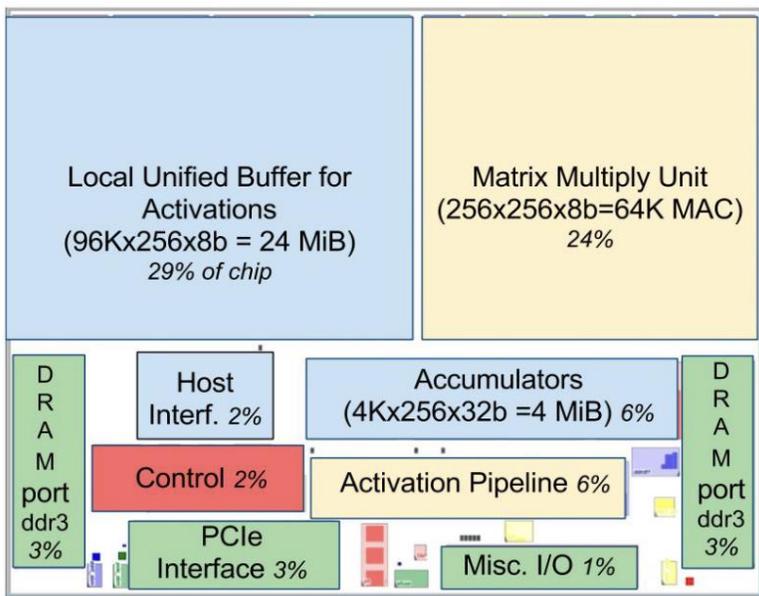
## • TPU设计参数

- 矩阵乘法单元:  $256 \times 256$ , 支持稠密矩阵乘法 (不支持稀疏计算), double buffering, 每周期产生256个partial sum, 支持8比特权重, 8/16比特输入
- 累加单元: 4 MB缓存 ( $4096 \times 256 \times 32b$ ), double buffering
- 片上存储: MMU存储64KB权重, FIFO深度为4, 激活存储为24MB



# AI加速器架构发展——Google TPU

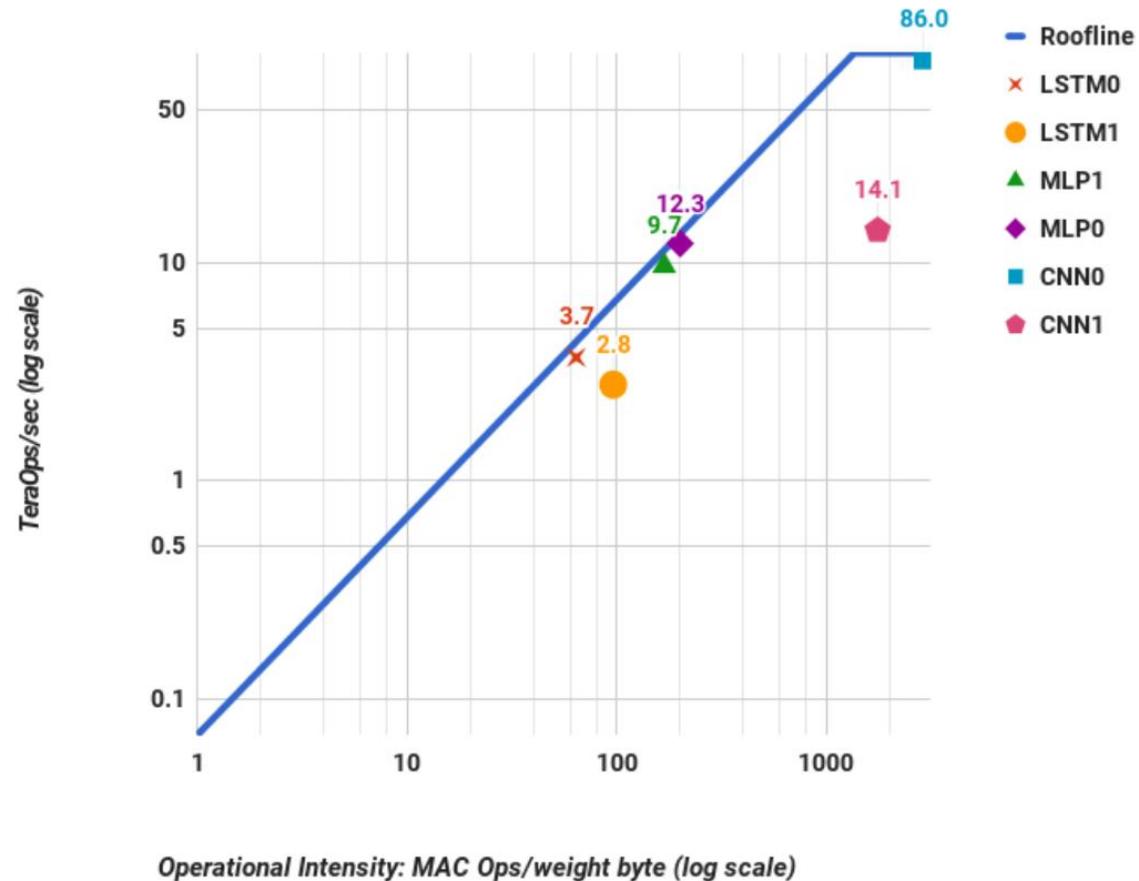
- 在TPU的版图中，数据通路面积占比为67%，I/O面积占比10%，控制占比2%
- 采用CISC指令，其中5个重要指令包括
  - Read\_Host\_Memory, Read\_Weights, MatrixMultiply/Convolve, Activate, Write\_Host\_Memory



- 基于Roofline模型的系统瓶颈评估
- TPU Roofline:
  - 每Byte的权重读取需要1350乘加计算平摊, 才能使得系统不会成为访存瓶颈
  - 访存瓶颈的操作: LSTM、MLP
  - 计算瓶颈的操作: 卷积

主讲: 陶耀

TPU Log-Log



- Google TPU v1的性能分析
  - 针对MLP、LSTM，模型权重的加载成为主要瓶颈
  - 实际算力与峰值算力之间存在显著差距

Name	Layers					Nonlinear function	Weights	TPU Ops / Weight Byte
	FC	Conv	Vector	Pool	Total			
MLP0	5				5	ReLU	20M	200
MLP1	4				4	ReLU	5M	168
LSTM0	24		34		58	sigmoid, tanh	52M	64
LSTM1	37		19		56	sigmoid, tanh	34M	96
CNN0		16			16	ReLU	8M	2888
CNN1	4	72		13	89	ReLU	100M	1750

Application	MLP0	MLP1	LSTM0	LSTM1	CNN0	CNN1	Mean	Row
Array active cycles	12.7%	10.6%	8.2%	10.5%	78.2%	46.2%	28%	1
Useful MACs in 64K matrix (% peak)	12.5%	9.4%	8.2%	6.3%	78.2%	22.5%	23%	2
Unused MACs	0.3%	1.2%	0.0%	4.2%	0.0%	23.7%	5%	3
Weight stall cycles	53.9%	44.2%	58.1%	62.1%	0.0%	28.1%	43%	4
Weight shift cycles	15.9%	13.4%	15.8%	17.1%	0.0%	7.0%	12%	5
Non-matrix cycles	17.5%	31.9%	17.9%	10.3%	21.8%	18.7%	20%	6
RAW stalls	3.3%	8.4%	14.6%	10.6%	3.5%	22.8%	11%	7
Input data stalls	6.1%	8.8%	5.1%	2.4%	3.4%	0.6%	4%	8
TeraOps/sec (92 Peak)	12.3	9.7	3.7	2.8	86.0	14.1	21.4	9

**Table 3. Factors limiting TPU performance of the NN workload based on hardware performance counters. Rows 1, 4, 5, and 6 total 100% and are based on measurements of activity of the matrix unit. Rows 2 and 3 further break down the fraction of 64K weights in the matrix unit that hold useful weights on active cycles. Our counters cannot exactly explain the time when the matrix unit is idle in row 6; rows 7 and 8 show counters for two possible reasons, including RAW pipeline hazards and PCIe input stalls. Row 9 (TOPS) is based on measurements of production code while the other rows are based on performance-counter measurements, so they are not perfectly consistent. Host server overhead is excluded here. CNN1 results are explained in the text.**

- 不同于TPU v1只支持推理，TPU v2需要**同时支持训练**
- 相较于推理，支持训练的难度更大：
  - 更多计算（和类型）：反向传播、转置、求导等操作
  - 更多存储：中间计算需要保存，反向传播中会用到
  - 更高精度：INT8无法满足训练需求
  - 更高的可编程性：不同的优化器等
  - 更难并行计算：

Build it quickly

Achieve high performance...

...at scale

...for new workloads out-of-the-box

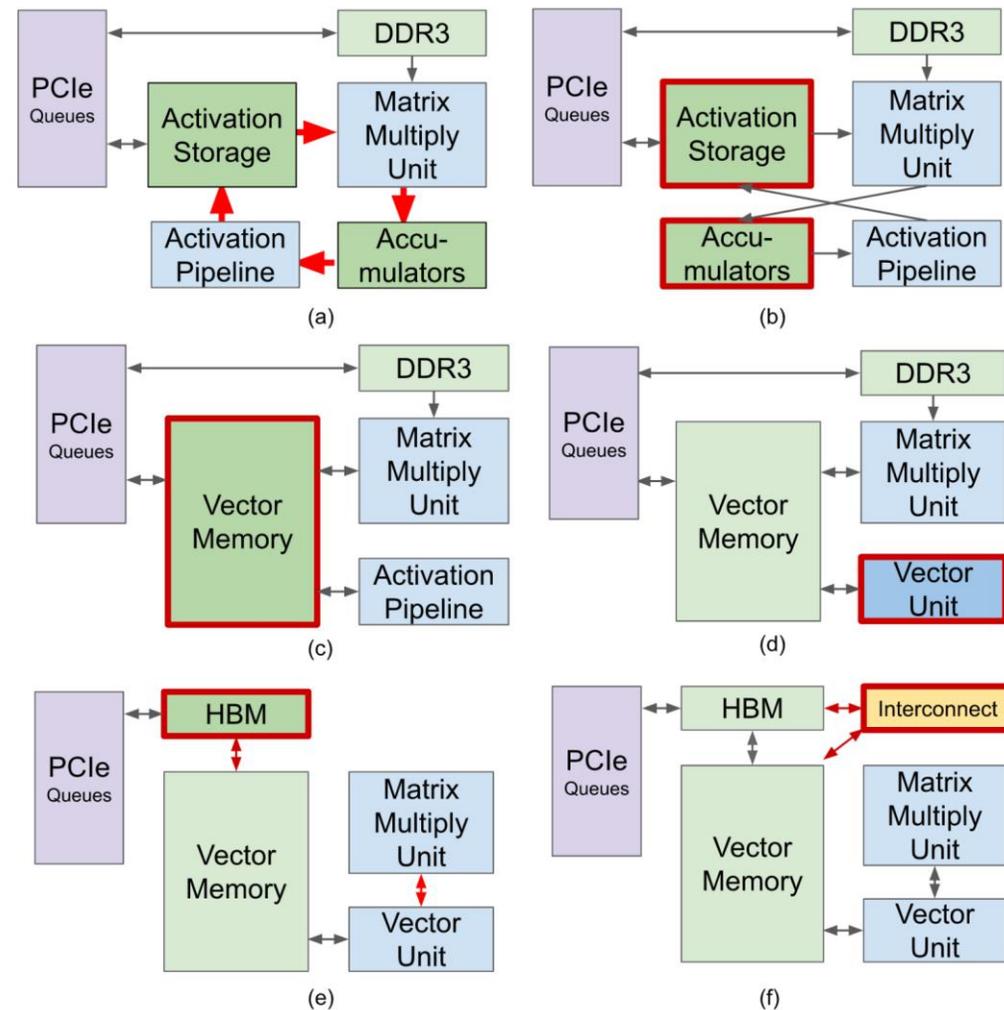
...all while being cost effective

主讲：陶耀宇、李明

# AI加速器架构发展——Google TPU v2

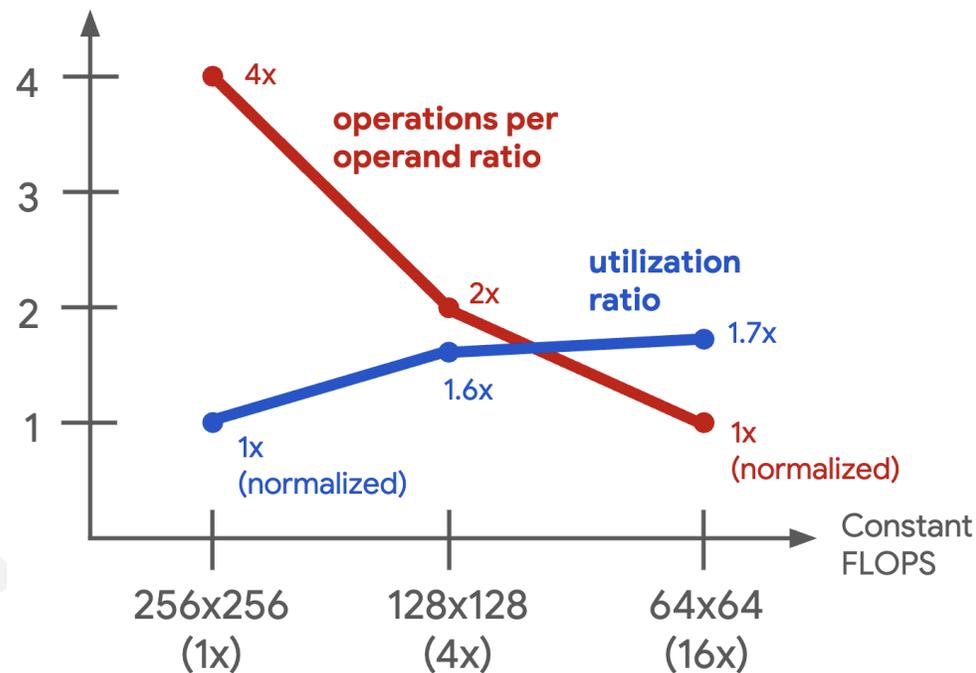
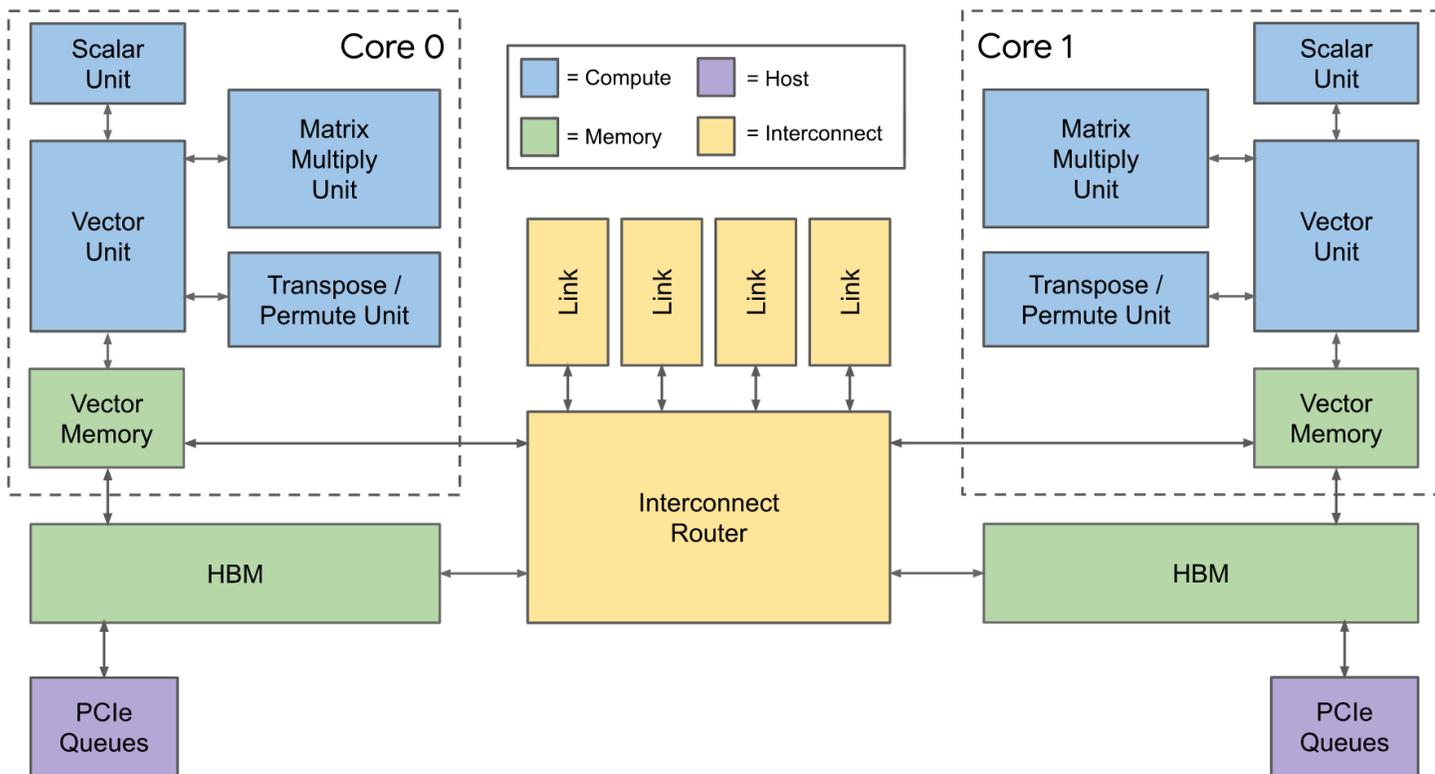
- 相较于TPU v1, TPU v2的**核心改动**包括
  - 将累加便签存储器、激活值便签存储器合并
  - 将累加器改为更为灵活、可编程的向量处理单元
  - 将MMU改为向量单元的协处理器, 简化访存需求
  - 将DDR3改为HBM, 提供更高访存带宽
  - 增加多卡互联, 提升多卡扩展效率

主讲：陶耀宇



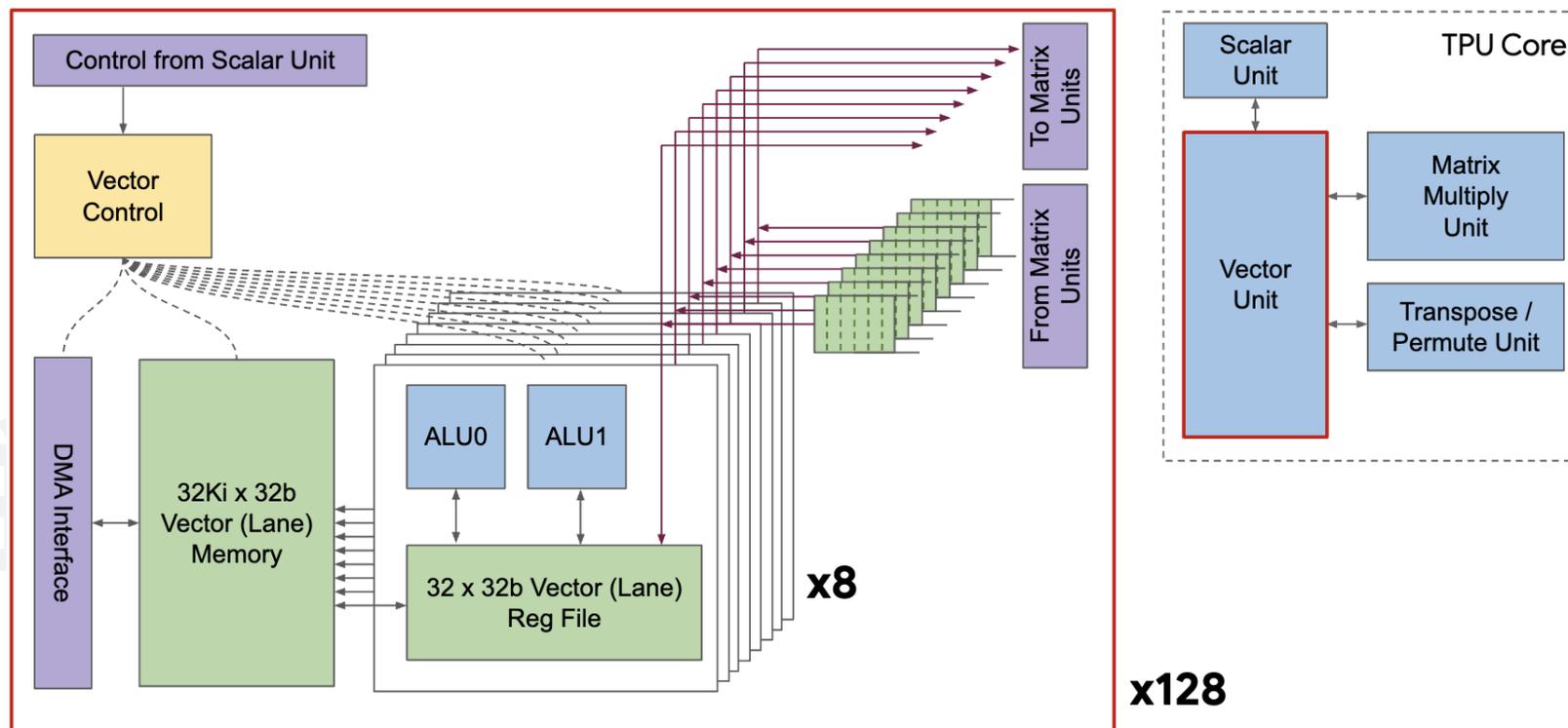
# AI加速器架构发展——Google TPU v2

- 采用双核架构，ISA支持矩阵、向量和标量计算，注重通用性
- 矩阵计算单元采用128 x 128的脉动阵列，支持bf16 (s1e8m7) 乘法和FP32累加
- 为什么要选择128 x 128的脉动阵列，而不是TPU v1中的256 x 256？



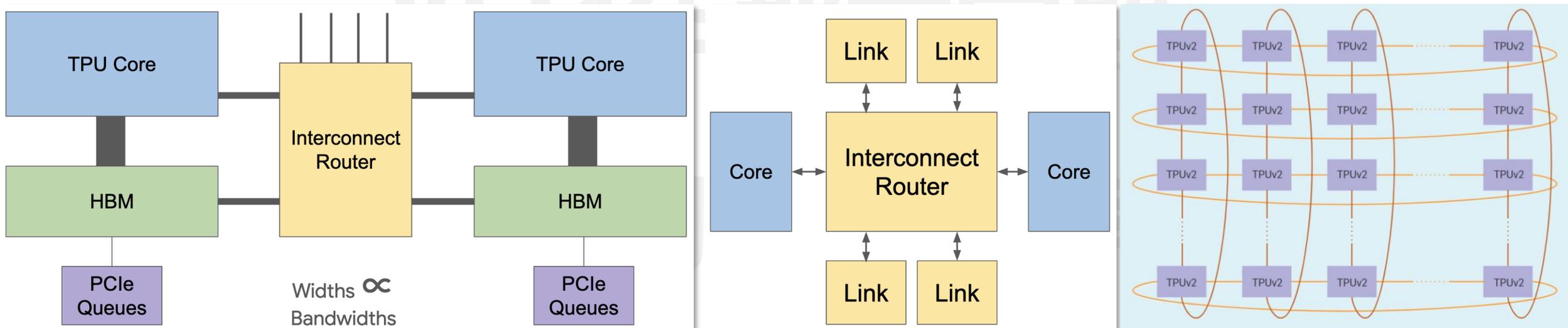
- 向量计算单元的功能：
  - 8个向量计算核，每个核每个周期处理128个输入
  - 为矩阵计算单元提供输入

北京大学-智能硬件体系结构



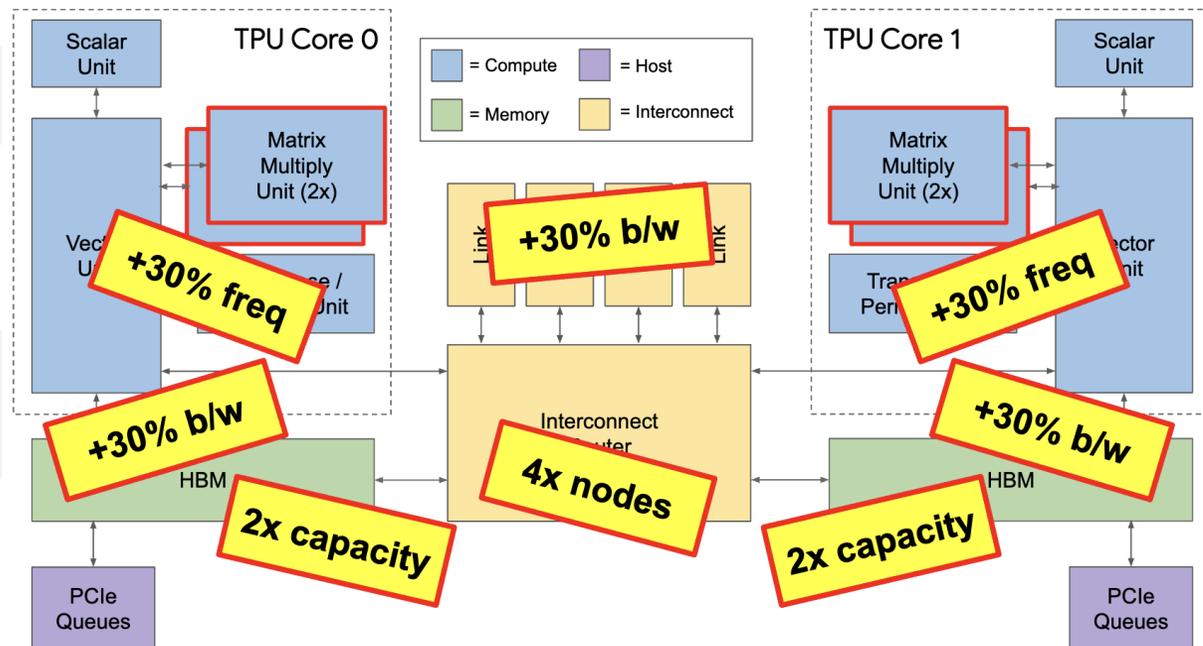
- TPU v2的存储系统

- 片上存储约为32MB
- 16GB HBM, 带宽约为600 GB/s (TPU v1采用DDR3, 带宽约为30 GB/s)
- 片上路由有4个链路, 每个链路带宽为500 Gbps, 采用2维圆环 (Torus) 的拓扑结构
- 最多支持256个TPU的拓展



# AI加速器架构发展——Google TPU v3

- 在TPU v2的基础上，**TPU v3**进行了进一步的改进了提升，包括
  - 矩阵乘法单元扩大2倍
  - 通过工程优化，将时钟频率从700 MHz提升到了940 MHz
  - 提升HBM带宽30%，扩大HBM容量2倍，提升互连带宽30%，达到每个链路650 Gb/s
  - 最大支持1024个TPU的扩展



# 目录

CONTENTS



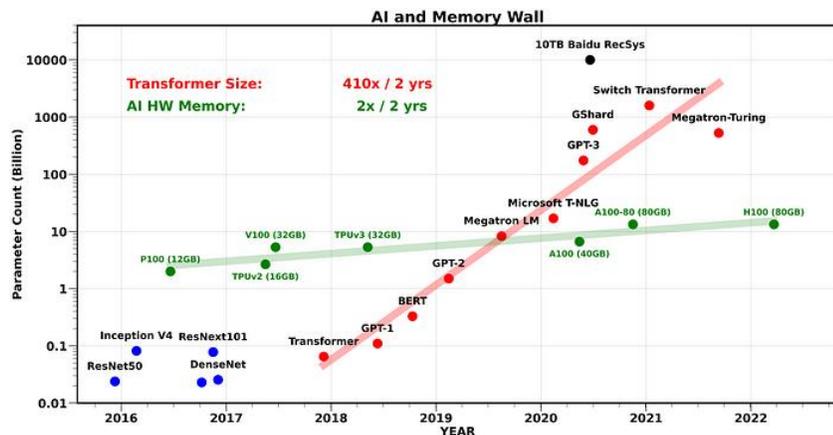
01. 典型AI芯片架构
02. AI芯片软硬件协同设计
03. AI大模型芯片设计
04. 存算一体AI芯片架构

- 我们介绍了人工智能处理器芯片的基本组成和代表性架构
  - 基本组成：计算单元、访存、互连通信
  - 代表性架构：DianNao系列、Eyeriss系列、Google TPU系列
- 接下来我们将介绍如何通过神经网络/处理器架构的协同设计和优化，进一步提升处理器芯片效率
  - 神经网络模型**量化**与混合精度人工智能处理器
  - 神经网络模型**稀疏化**与人工智能处理器架构

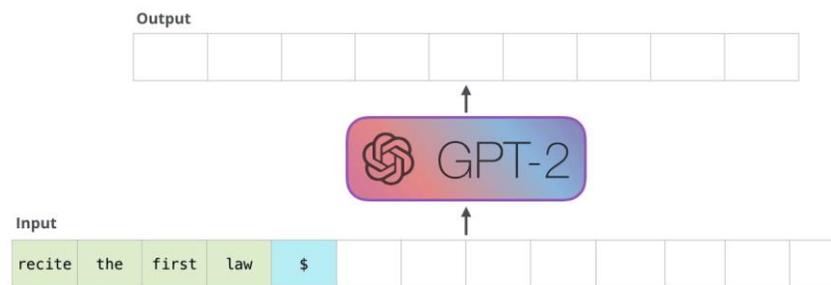
主讲：陶耀宇、李萌

- 高效人工智能模型计算面临高存储、高带宽、动态计算图、部署平台复杂多样的挑战

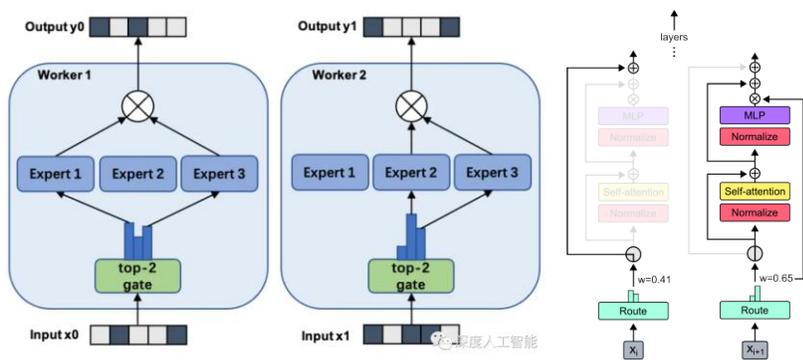
模型参数指数级增长，造成单模型存储需求 > 1 TB



大模型自回归解码，导致平均访存带宽需求 > 1 TB/s



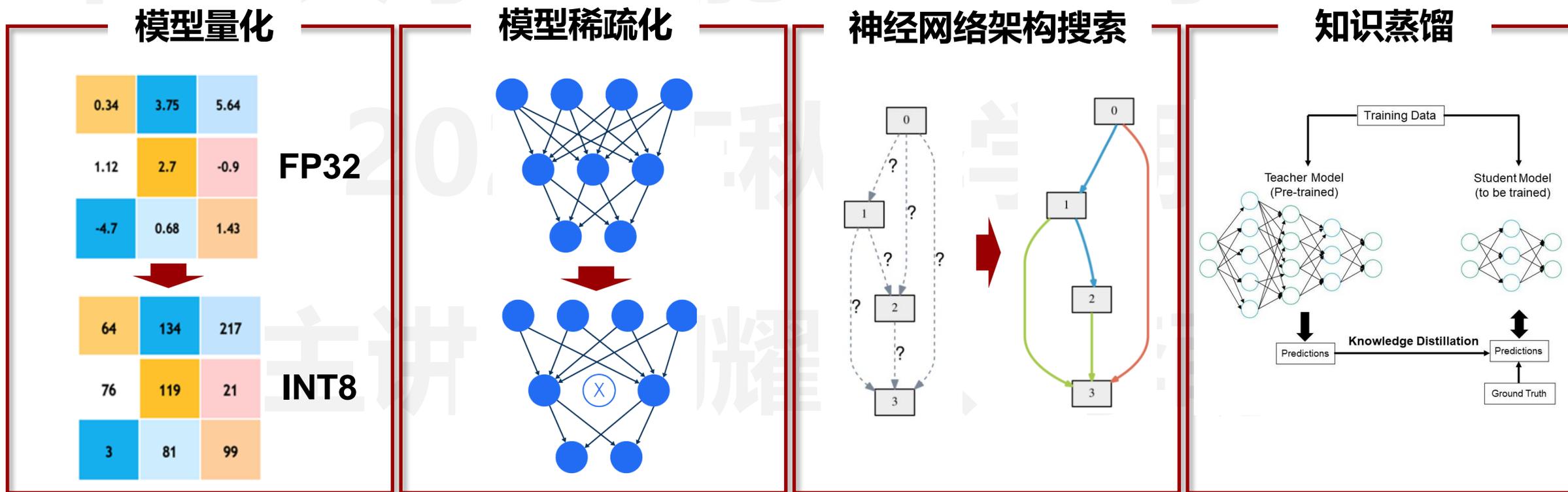
动态计算图导致传统静态数据流优化失效



部署平台的计算、存储能力存在显著差距

	Cloud AI	Mobile AI	Tiny AI
Memory	24 GB	4 GB	500 KB
Storage	~ TB/PB	~100s GB	~1s MB
# params	> 70 B	3 - 13 B	~1s M

- 为了进一步提升人工智能模型的推理效率，采用算法和AI处理器协同设计和优化方法
  - 神经网络模型**量化**、**稀疏化**、**神经网络架构搜索**、**知识蒸馏**等
- 本次课我们将以神经网络模型量化、稀疏化为例，介绍算法和加速器的协同设计



# 基于模型量化的软硬件协同设计

- 量化是将数据从**连续值集**映射为**离散值集**的过程
- 模型量化能够利用神经网络对于量化噪声的容错能力，显著提升模型推理计算效率
  - 无论是输入图像，还是神经网络本身都对于量化噪声具有较高的容忍能力

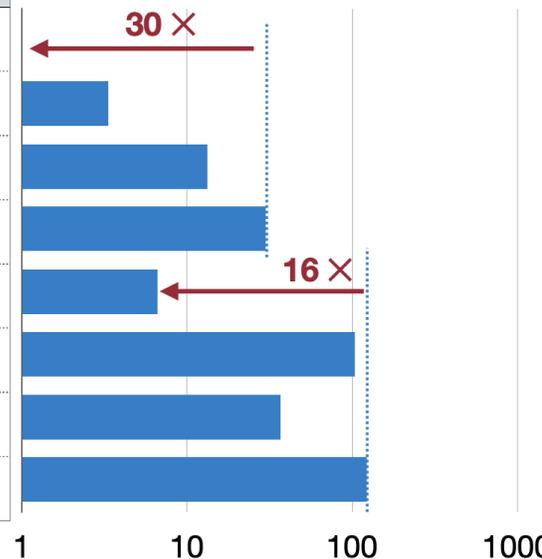
2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
-0.91	1.92	0	-1.03
1.87	0	1.53	1.49



1	-2	0	-1
-1	-1	-2	1
-2	1	-1	-2
1	-1	0	0

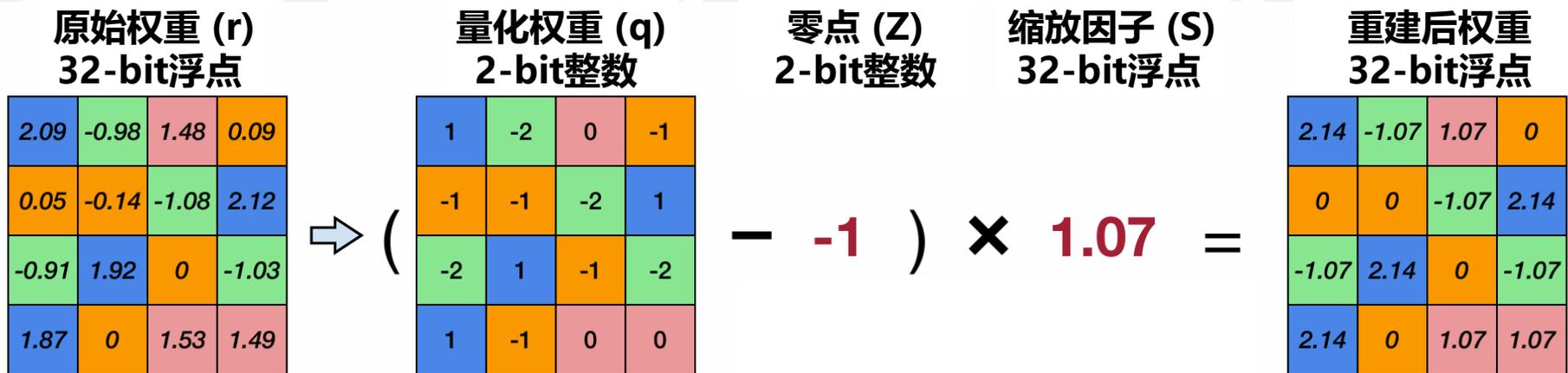
$- (-1) \times 1.07$

Operation	Energy [pJ]
8 bit int ADD	0.03
32 bit int ADD	0.1
16 bit float ADD	0.4
32 bit float ADD	0.9
8 bit int MULT	0.2
32 bit int MULT	3.1
16 bit float MULT	1.1
32 bit float MULT	3.7



Rough Energy Cost For Various Operations in 45nm 0.9V

- 量化过程中的主要参数：缩放因子、零点、量化后数值
  - 静态量化：量化参数在推理前确定，推理时保持不变
  - 动态量化：量化参数在推理时计算得到



2-bit整数表示方法

Binary	Decimal
01	1
00	0
11	-1
10	-2

量化误差

-0.05	0.09	0.41	0.09
0.05	-0.14	-0.01	-0.02
0.16	-0.22	0	0.04
-0.27	0	0.46	0.42

- 权重的量化与激活值的量化对于计算的影响不同

北京大学 - 智能硬件体系结构

$$Y = WX$$

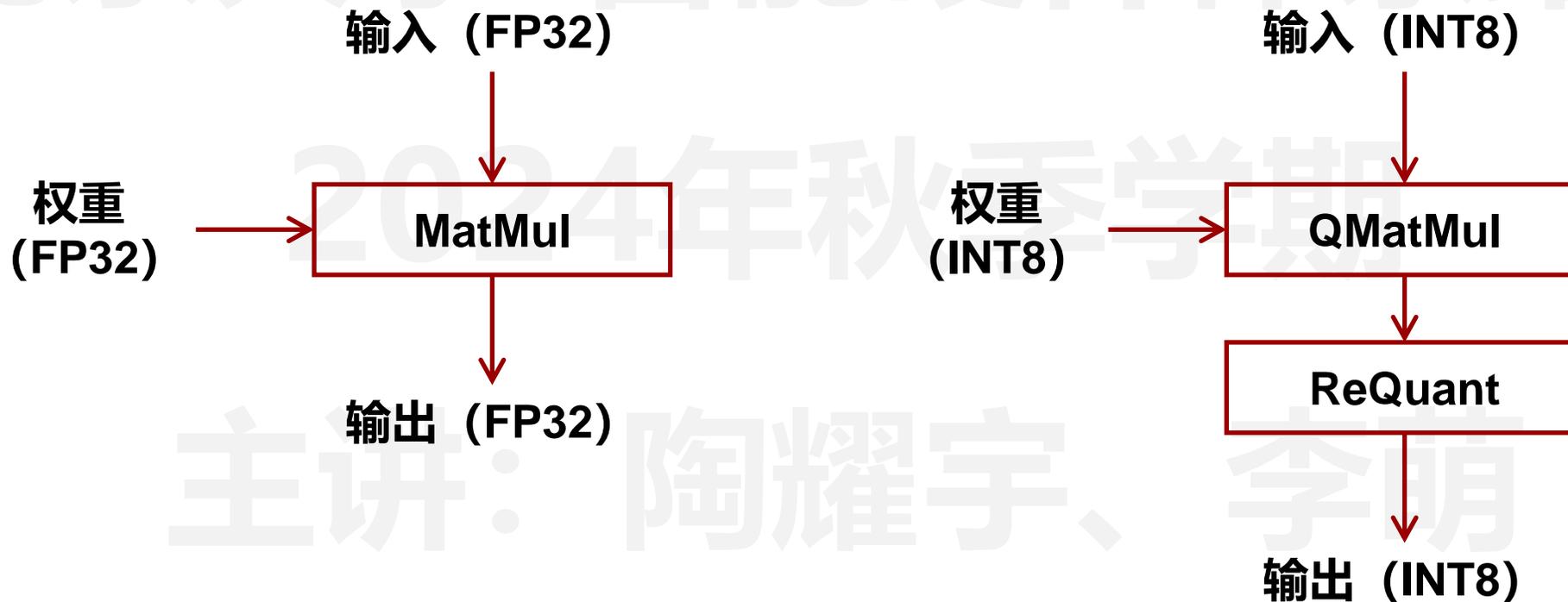
$$S_Y (q_Y - Z_Y) = S_W (q_W - Z_W) \cdot S_X (q_X - Z_X)$$

$$q_Y = \frac{S_W S_X}{S_Y} (q_W - Z_W)(q_X - Z_X) + Z_Y$$

$$q_Y = \frac{S_W S_X}{S_Y} (q_W q_X - Z_W q_X - Z_X q_W + Z_W Z_X) + Z_Y$$

**QMatMul**  $Z_W$  通常为0 提前计算

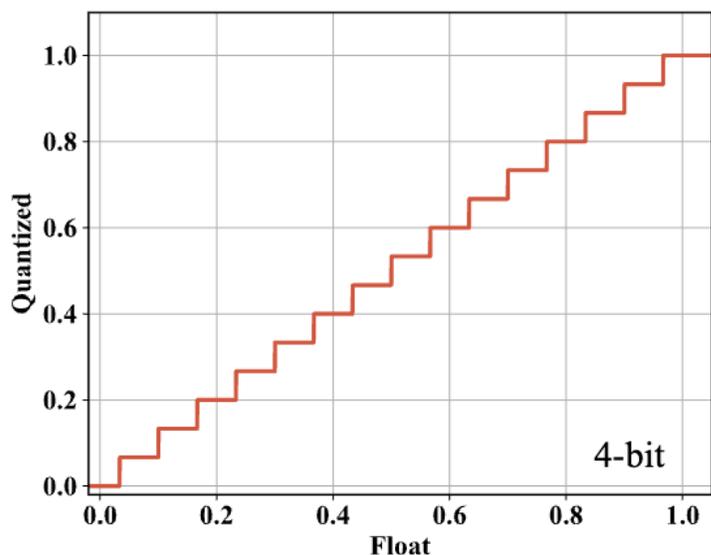
- 权重的量化与激活值的量化对于计算的影响不同
- 量化神经网络推理：引入新的**重量化**计算，将高比特精度计算结果重量化为低比特精度输出
- 思考：QMatMul和MatMul的计算有什么区别？



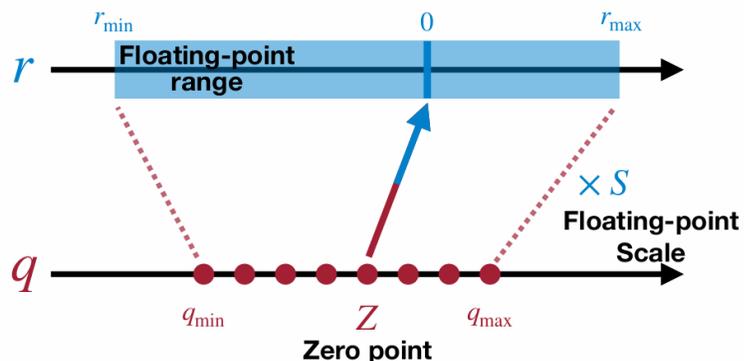
# 均匀量化与非均匀量化

- **均匀量化**: 输入到输出的映射为**线性函数**, 因此均匀间隔的输入产生均匀间隔的输出
  - 针对均匀量化, QMatMul和MatMul的计算, 除比特精度外, 基本相同

北京大学-智能硬件体系结构



(a) Uniform Quantization



$$r_{max} = S (q_{max} - Z)$$

$$r_{min} = S (q_{min} - Z)$$

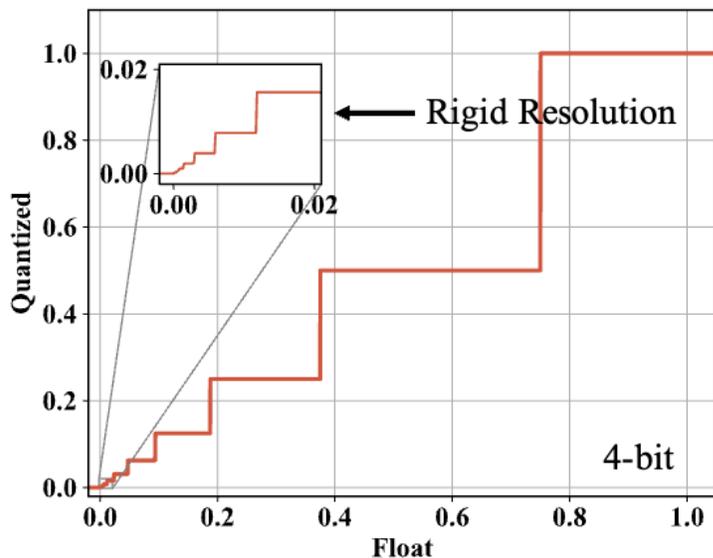


$$r_{max} - r_{min} = S (q_{max} - q_{min})$$

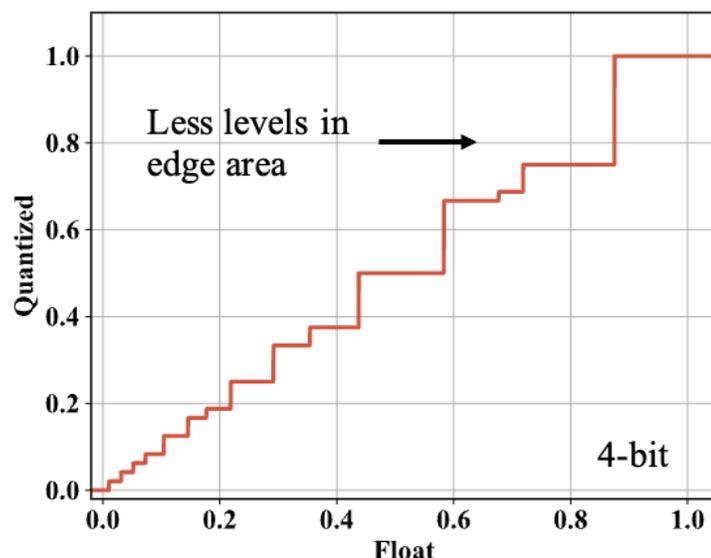
$$S = \frac{r_{max} - r_{min}}{q_{max} - q_{min}}$$

- 均匀量化：输入到输出的映射为线性函数，因此均匀间隔的输入产生均匀间隔的输出
- 非均匀量化：输入到输出的映射为**非线性函数**，常见算法包括Power of Two (PoT)、APoT等
- 思考：均匀和非均匀量化还有哪些类别？

北京大学-智能硬件体系结构



(b) Power-of-Two Quantization

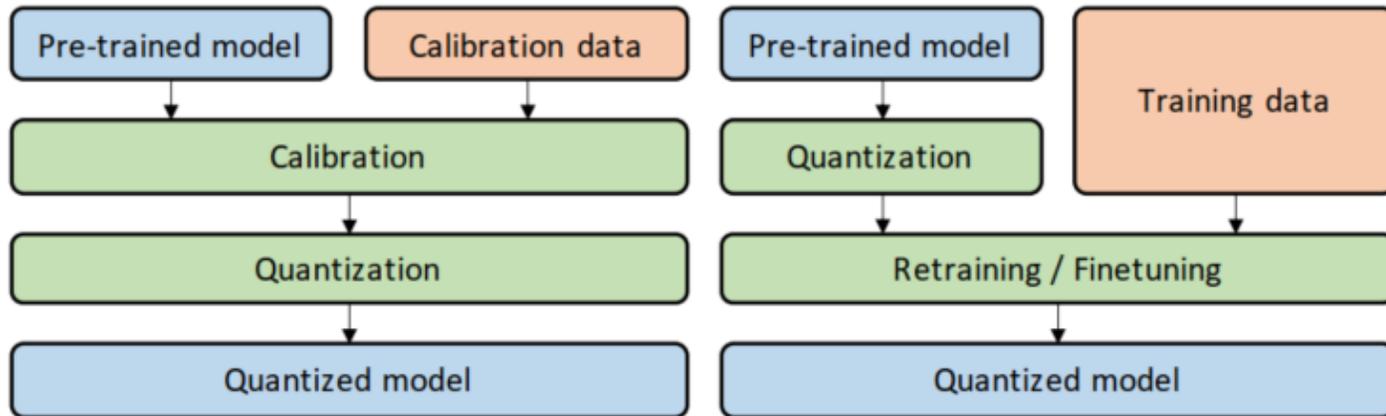


(c) Additive Power-of-Two Quantization

**PoT量化：乘法可以用移位进行计算，显著降低计算复杂度，但是，对于较高比特数，PoT量化提升有限**

**APoT量化：乘法通过移位和加法实现，计算效率受APoT的项数决定**

- **训练后量化 (post training quantization)** : 模型训练完成后对权重、激活值进行量化, 其中, 激活值的量化, 需要借助校准数据
- **量化感知训练 (quantization-aware training)** : 将训练过的模型量化后又再进行重训练, 或者直接训练量化模型



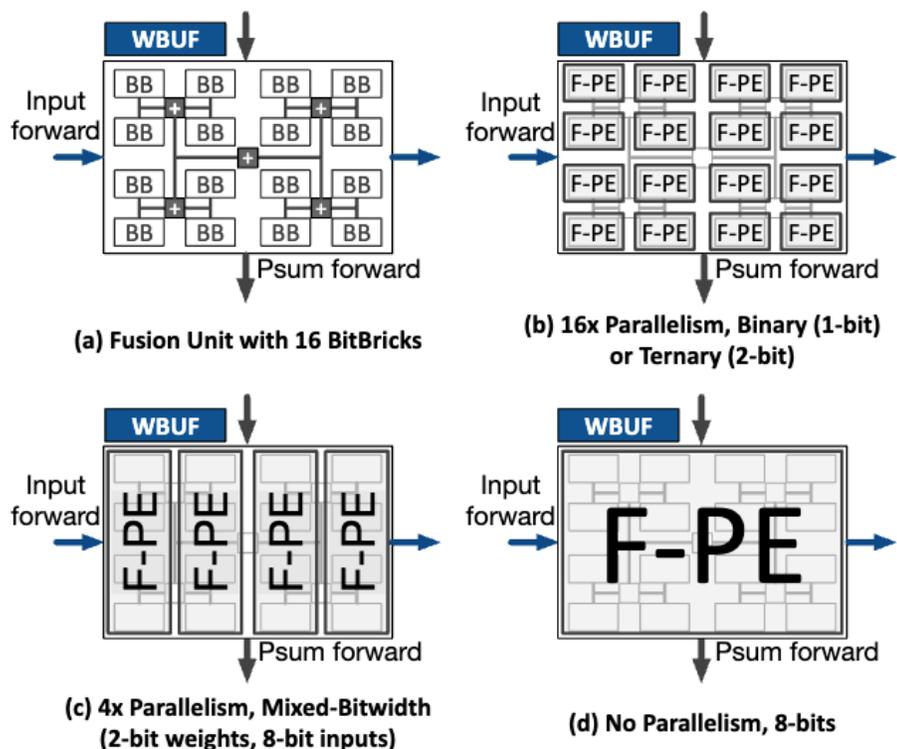
**PTQ运行时间短, 数据需求小, 但是模型推理准确率下降较大**

**QAT准确率高, 但是训练时间长, 数据要求高, 在特定任务上可能过拟合**

- 量化神经网络想要真正实现计算能效、吞吐等的提升，离不开**AI处理器**的支持
  - 例如支持低比特计算的计算单元（乘法器、累加器等）、支持低比特数据的存储方式等等
- 因此，神经网络量化往往需要配合专用硬件（或核函数）设计，形成软硬件协同设计和优化
- 本节课，我们重点介绍3个例子
  - BitFusion：支持混合精度计算的AI处理器
  - OliVe：面向低比特大模型的AI处理器
  - ANT：基于定制化数据类型的AI处理器

主讲：陶耀宇、李萌

- Bit Fusion文章发表于ISCA 2018，核心在于高效支持混合精度神经网络推理
- 核心观察：不同神经网络模型以及同一神经网络模型中的不同层，对于量化噪声的容错能力不同
  - 采用混合精度量化可以在保持模型精度的同时，最大化压缩模型
- 文章提出Bit Fusion架构，重点通过可重构的计算单元，实现不同计算精度的灵活支持



- Fusion单元由16个BitBricks组成，每个BitBrick每周期可以支持2比特计算
- 多个BitBricks组合，能够实现W4A4、W2A8、W8A2以及W8A8的计算

Hardik Sharma et al., *Bit Fusion: Bit-Level Dynamically Composable Architecture for Accelerating Deep Neural Network*, ISCA 2018



- ANT发表于MICRO 2022, 提出了新的数据类型, 适应神经网络推理需求
- 核心观察: 1. 神经网络中不同的tensor分布各不相同; 2. 同一个tensor内部对于接近0的值或特别大的值不需要使用高精度表示; 使用现有的INT或Float格式量化难以适应这两种变化
- 提出了一种新的数据类型flint, 适合表示高斯分布的数据

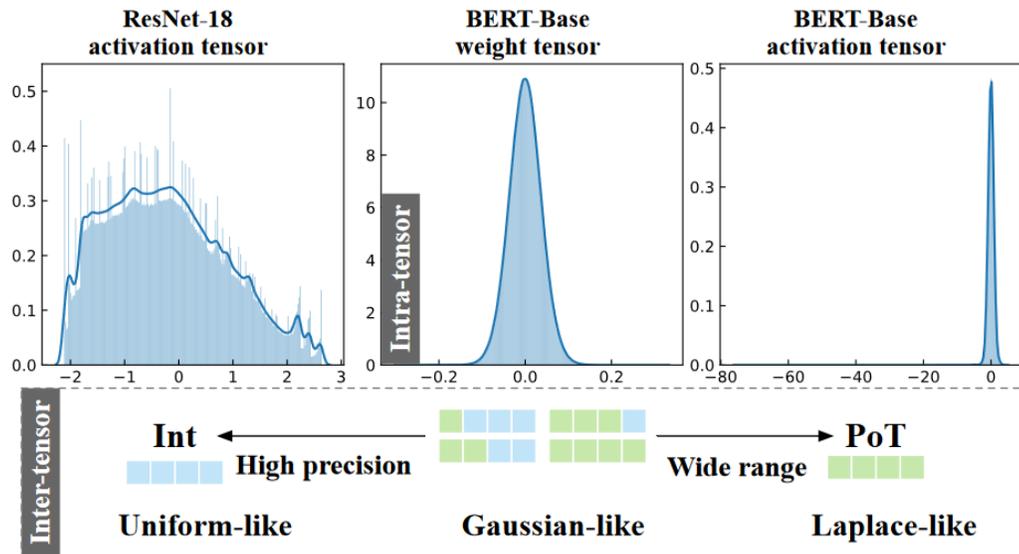
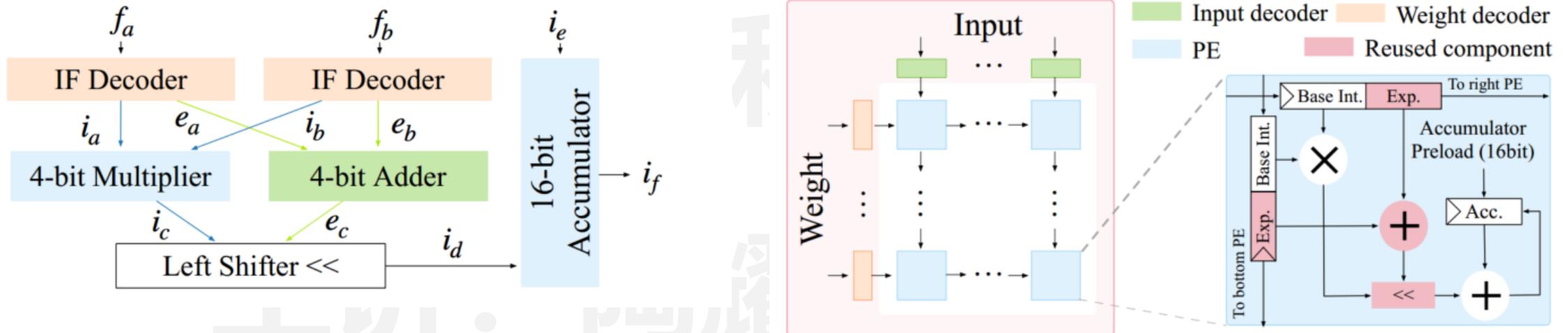


Figure 1: Intra-tensor and inter-tensor adaptivity.

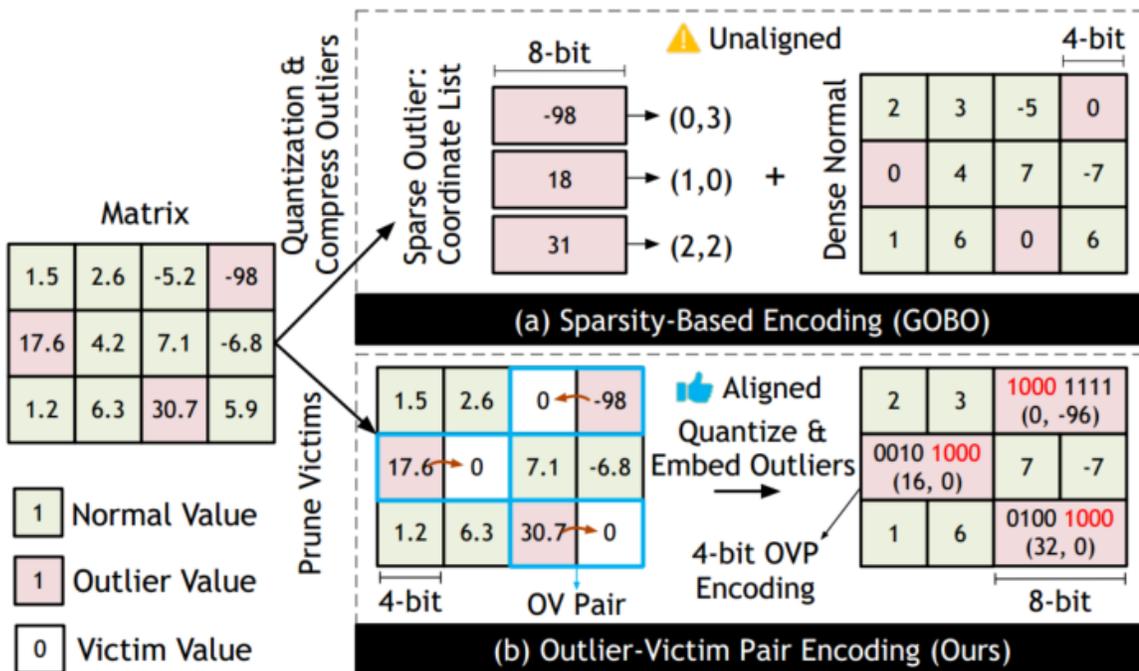
Bits	Exponent Value	Fraction Value	Value in Decimal
0000	-	0	0
0001	$1 - 1 = 0$	1	$2^0 \times 1 = 1$
001x	$2 - 1 = 1$	1, 1.5	2, 3
01xx	$3 - 1 = 2$	1, 1.25, 1.5, 1.75	4, 5, 6, 7
11xx	$4 - 1 = 3$	1, 1.25, 1.5, 1.75	8, 10, 12, 14
101x	$5 - 1 = 4$	1, 1.5	16, 24
1001	$6 - 1 = 5$	1	32
1000	$7 - 1 = 6$	1	$2^6 \times 1 = 64$

Table II: The value table of 4-bit unsigned flint with the exponent bias of  $-1$ . The blue numbers are the first-one-encoded exponent and “x” is mantissa with value of 0 or 1.

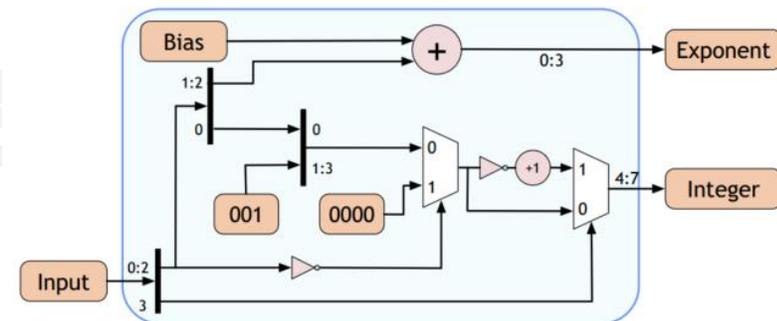
- 为了解决tensor之间数据分布不同的问题，通过校准数据集选择量化后均方误差最小的数据类型
  - 直觉上int类型适合均匀分布，Float/PoT类型适合拉普拉斯分布，Flint类型适合高斯分布
- 设计了可以支持以上**四种格式的处理单元（PE）**，主要包括解码器，乘法器，加法器和累加器
- 将设计的PE集成到了**脉动阵列**形成新的架构



- OliVe发表于ISCA 2023, 针对大模型中的**离群值** (outliers) 进行设计
- 核心观察: 大语言模型中存在一些outliers, 这些outliers很重要, 但旁边的正常值 (称为victims) 不重要, 因此可以牺牲victims进而用更多比特表示outliers
- 对于outliers, 提出了新的数据格式**Abfloat**, 通过自适应的bias使编码的值表示比正常值更大的数

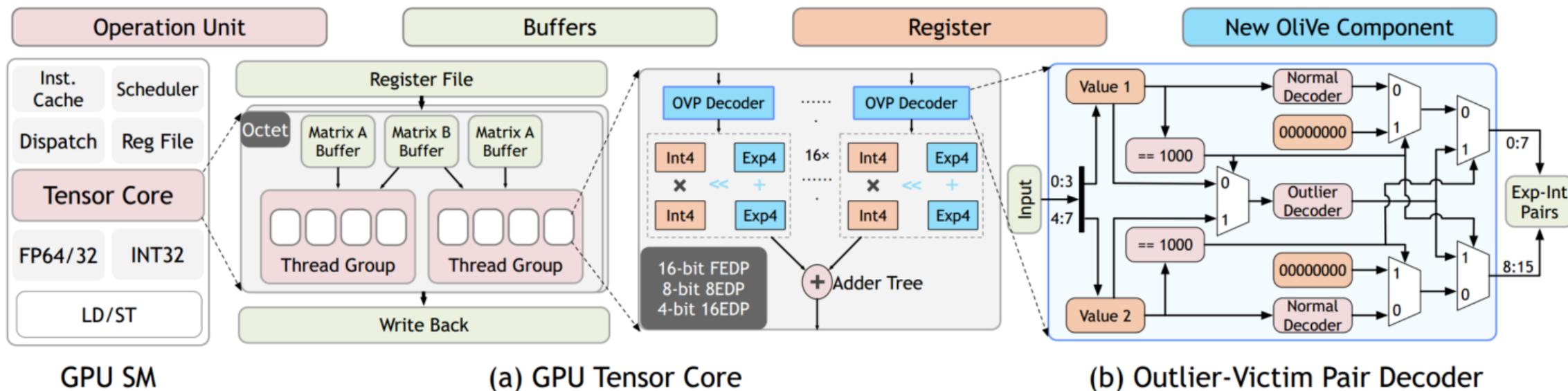


Binary	Exponent	Integer	Real Value
000	0	0	0
001	0	3	$3 \times 2^0 = 3$
01x	1	2, 3	$2 \times 2^1 = 4, 3 \times 2^1 = 6$
10x	2	2, 3	$2 \times 2^2 = 8, 3 \times 2^2 = 12$
11x	3	2, 3	$2 \times 2^3 = 16, 3 \times 2^3 = 24$



- 解码完成后，文章设计了乘累加单元电路（MAC Unit），并将他们集成到了GPU tensor core以及脉动阵列中

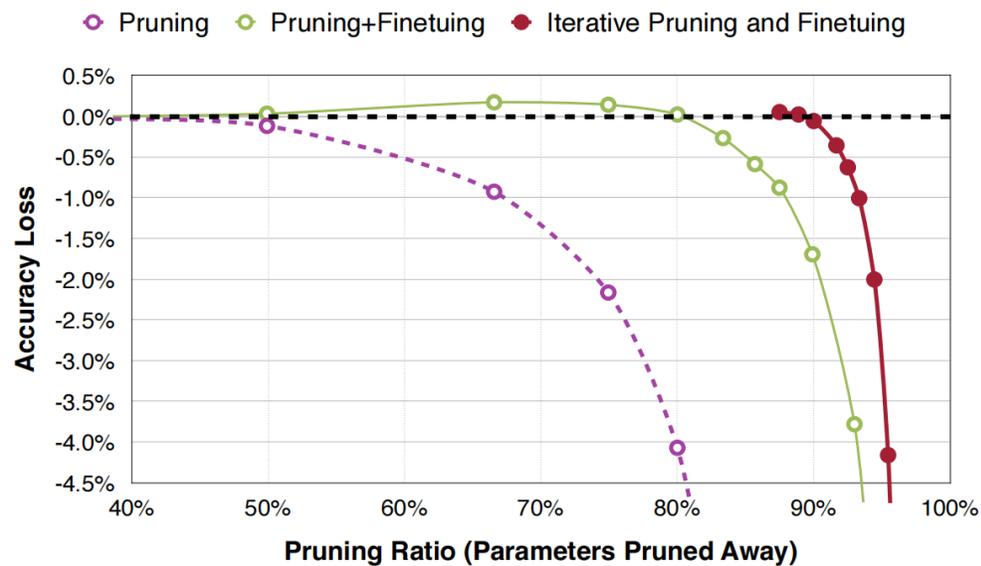
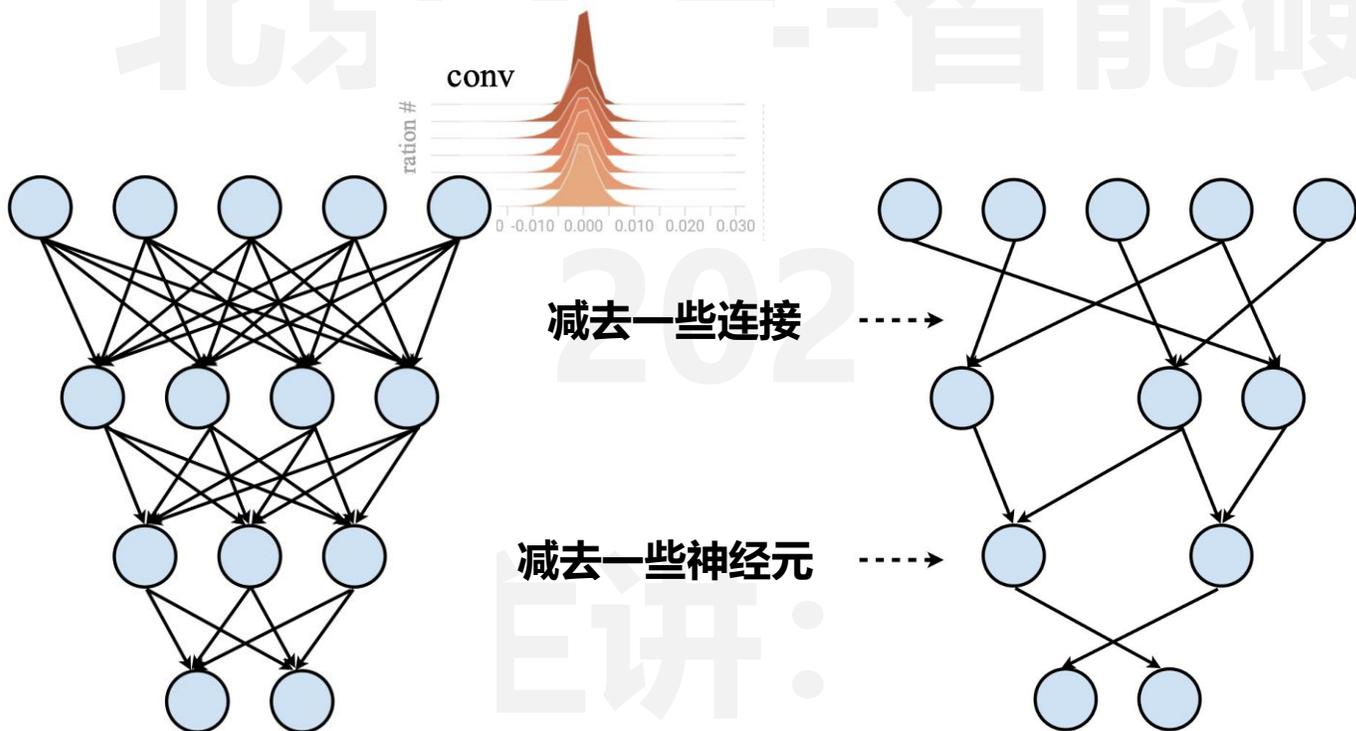
北京大学 智能硬件系统结构



工研·阿雁子、子明

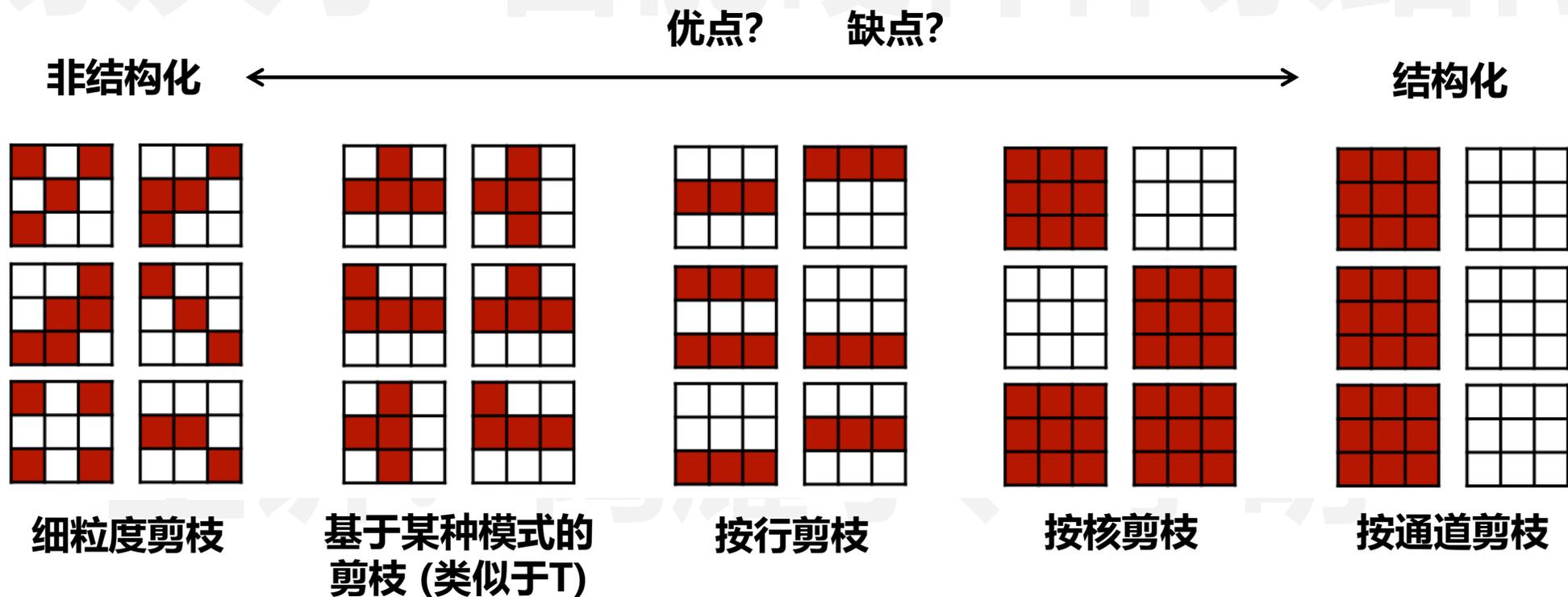
- 除了对于量化噪声的容错性，神经网络往往呈现显著的**稀疏性**，即对特定神经元或者模型权重进行剪枝，不会影响其推理准确率

北京大学 - 智能硬件体系结构



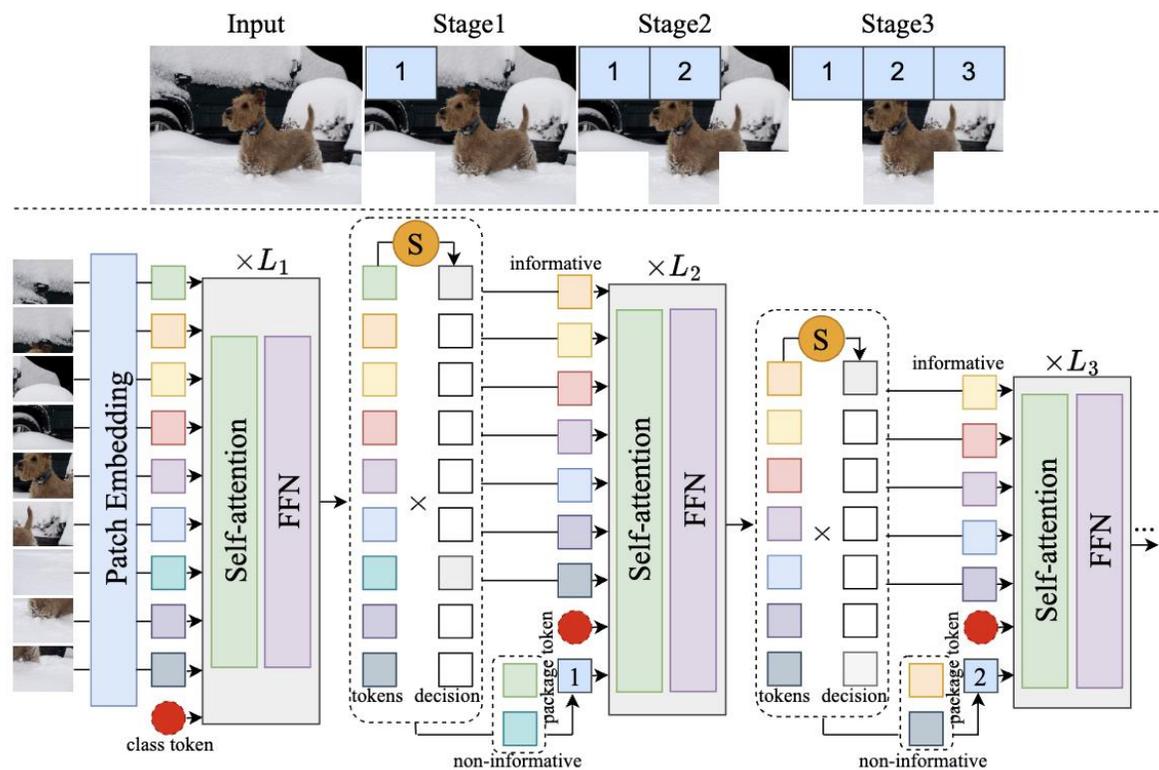
- 依据剪枝粒度不同，神经网络剪枝可以分为**结构化**和**非结构化剪枝**
- 细粒度非结构化剪枝更加灵活，往往准确率更高（或稀疏度更高），但是硬件加速更加困难
- 粗粒度结构化剪枝则更加硬件友好，但是灵活性受限

北京大学-智能硬件体系结构



- 依据剪枝对象不同，神经网络剪枝可以分为**权重剪枝**与**激活值剪枝**
- 权重剪枝适用于CNN、RNN、Transformer等，通常为**静态剪枝**
- 激活值剪枝则更常见于Transformer模型或基于ReLU的CNN模型中，通常为**动态剪枝**

北京大学-智能硬件体系结构



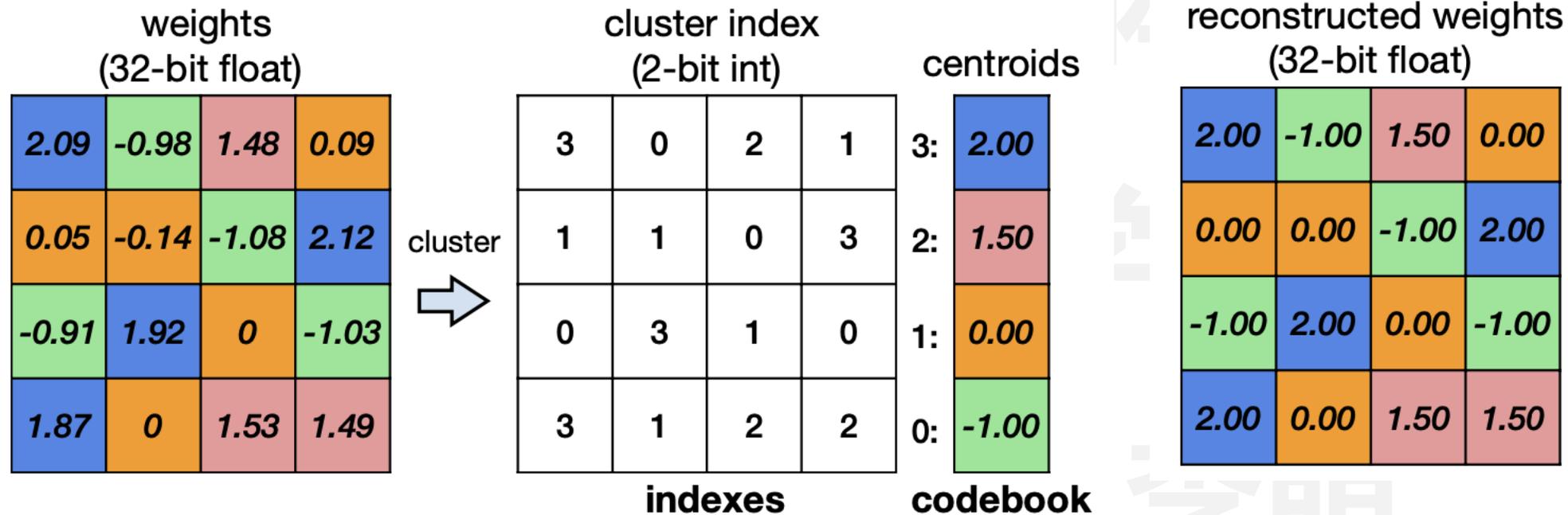
春季学期

Peiyan Dong et al., *HeatViT: Hardware-Efficient Adaptive Token Pruning for Vision Transformers*, HPCA 2023

李、李萌

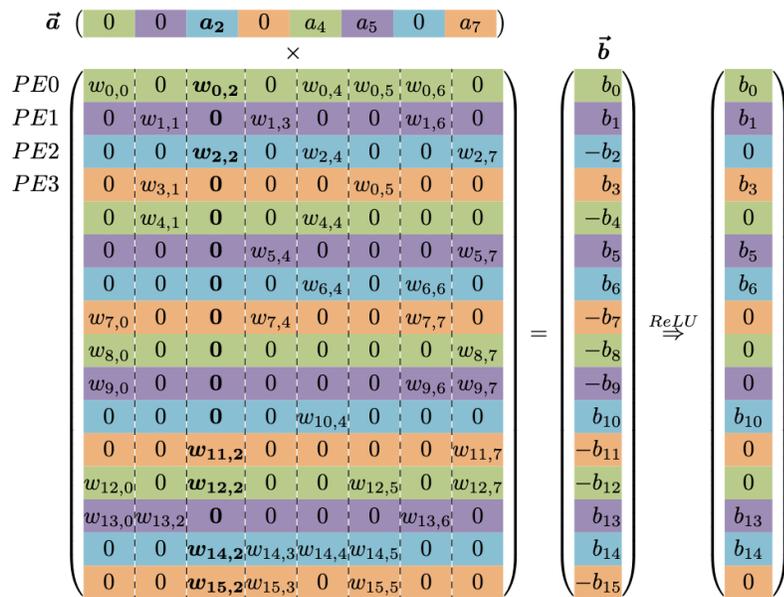
- 稀疏神经网络想要真正实现计算能效、吞吐等的提升，同样需要底层AI处理器的支持
  - 特别是对于**稀疏激活值**，导致计算呈现显著动态性
  - 例如动态剪枝单元、动态计算单元等
- 因此，神经网络量化往往需要配合专用硬件（或核函数）设计，形成软硬件协同设计和优化
- 本节课，我们重点介绍3个例子
  - EIE、SNAP：同时考虑权重和激活稀疏的AI处理器架构
  - Nvidia sparse tensor core：商用稀疏AI处理器架构
- 一个很好的Tutorial：Sparse Tensor Accelerator Modeling Tutorial @ ISCA 2021

- 文章发表于ISCA 2016，针对基于量化和剪枝的压缩模型，提出了高效推理引擎EIE



Song Han et al., *EIE: efficient inference engine on compressed deep neural network*, ISCA 2016

- 文章发表于 **ISCA 2016**，针对基于量化和剪枝的压缩模型，提出了高效推理引擎EIE
- 核心观察：缺少有效处理**不规则存储**的稀疏数据的硬件单元
  - 静态的权重稀疏+动态的激活稀疏
- 文章提出高效推理引擎EIE，实现了高效的稀疏矩阵向量乘法



- 权重稀疏通过 Index + value 存储为 Compressed sparse column (CSC) 格式
- 激活稀疏通过只向PE广播非零激活值实现
- 权重共享通过将权重固定为16个值并用4位索引查找实现

Figure 2. Matrix  $W$  and vectors  $a$  and  $b$  are interleaved over 4 PEs. Elements of the same color are stored in the same PE.

- **Act Queue平衡负载**: 非零激活值对应的一系列权重中的非零数量不同导致不同PE间负载不平衡
- **Leading Non-zero Detection Node**: 检测非零激活值并广播到所有PE

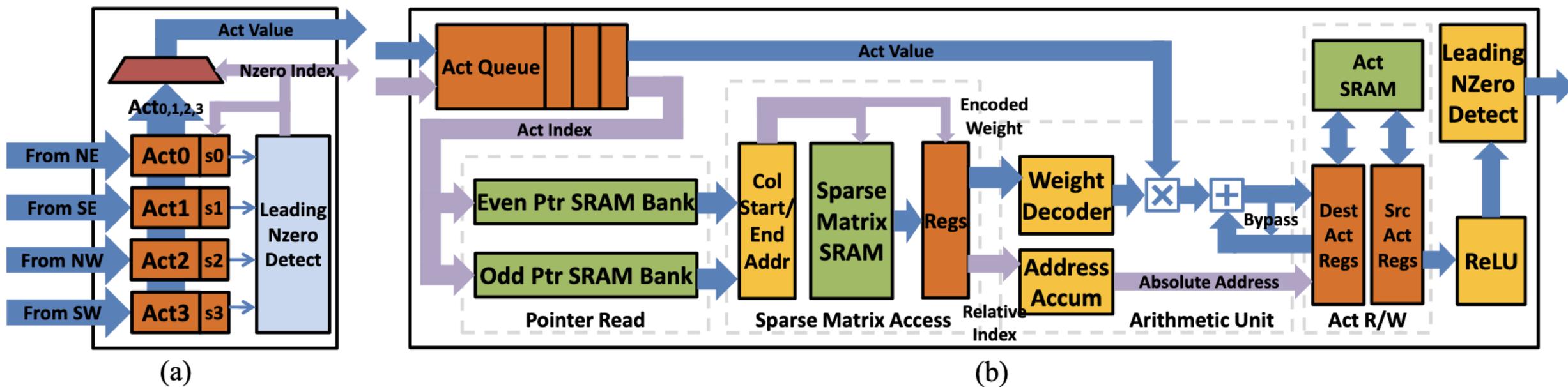
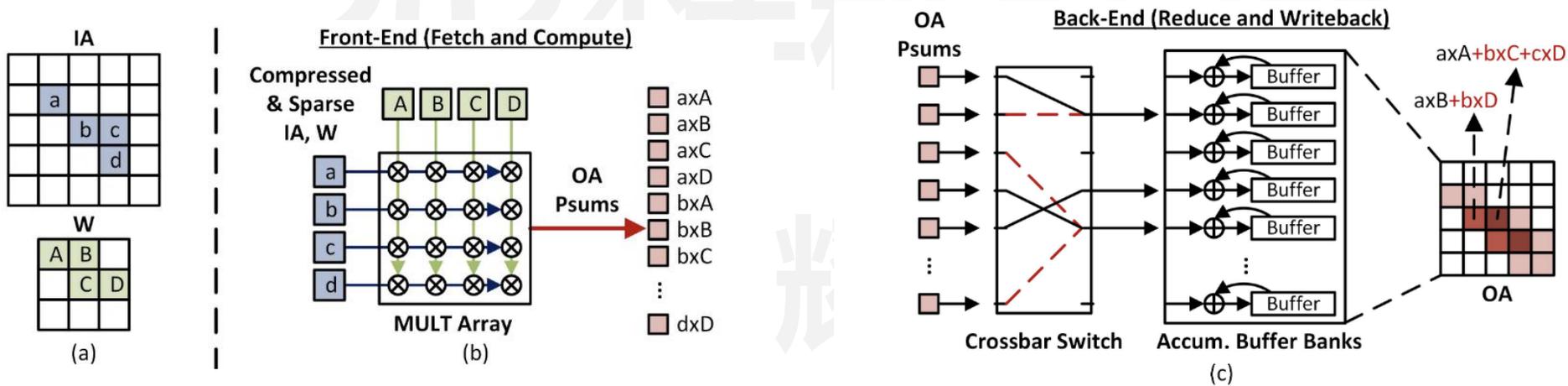


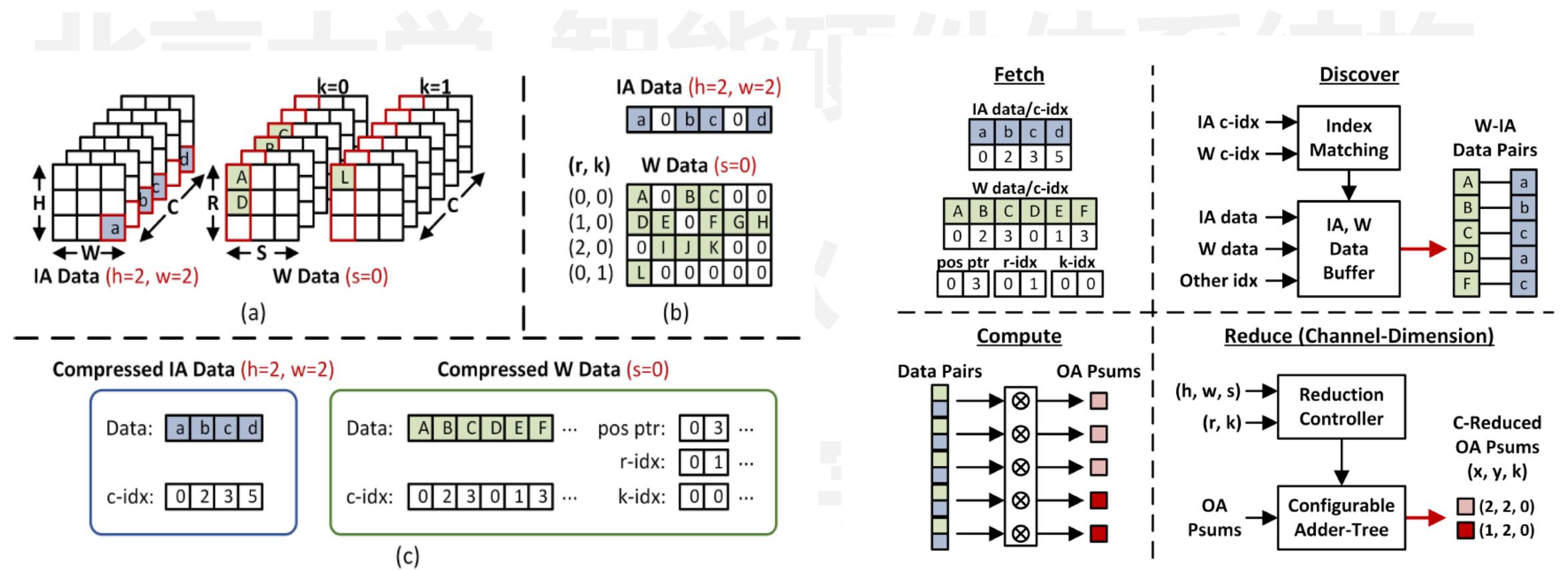
Figure 4. (a) The architecture of Leading Non-zero Detection Node. (b) The architecture of Processing Element.

- 文章发表于JSSCC 2021，通过channel-first数据流对稀疏网络模型进行了硬件加速
- 核心观察：数据稀疏使得网络推理更加高效，但现有稀疏计算的数据流仍面临以下挑战

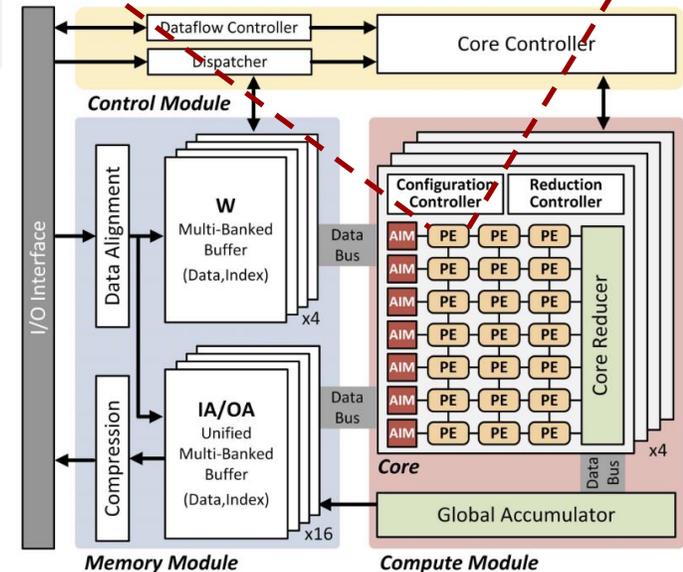
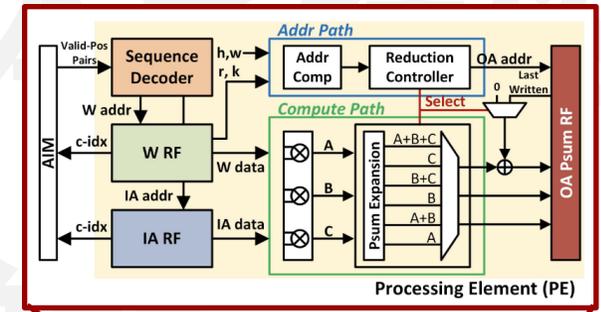
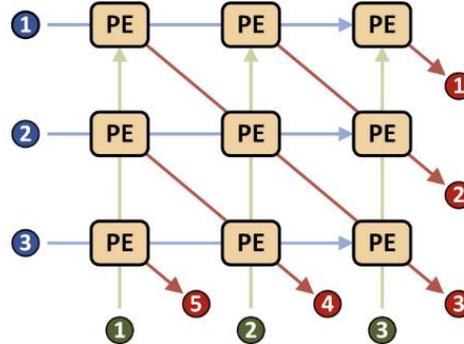
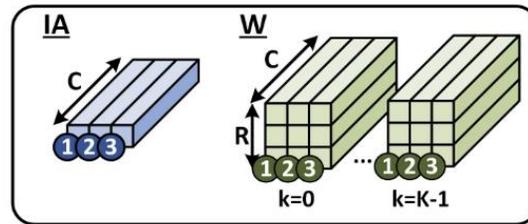
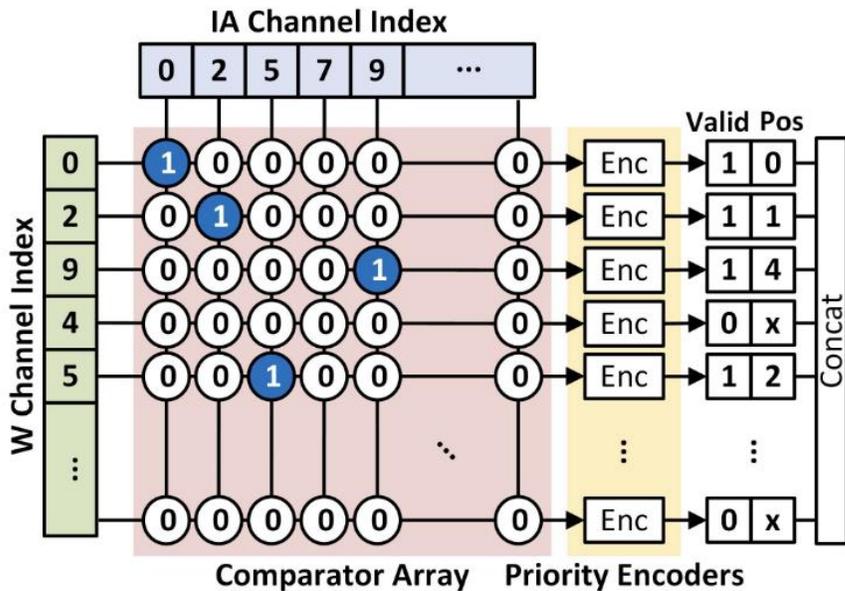
- Front-end：读取的W-IA pairs数量不足导致计算单元利用率低
- Back-end：计算结果写回的地址冲突
- 灵活性：对于不同kernel大小和层类型的的支持性不足



- **Channel-first数据流**: 同一个channel的可以任意配对, 产生的乘法结果都是有用的
- 通过可配置加法树实现先**归约加和**, 然后更新Psum, 从而减少Psum更新时的地址冲突

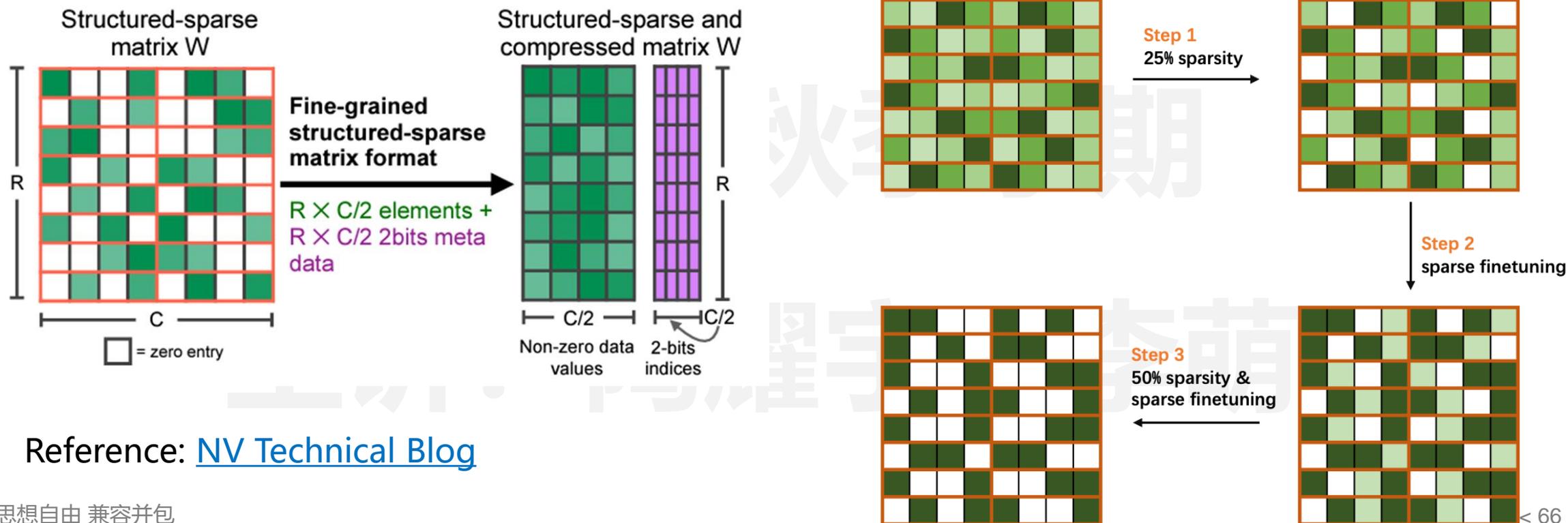


- Front-end: Index Matching单元提取足够数量的 W-IA pairs, 提高乘法阵列利用率
- Back-end: 两级 partial sum 压缩, PE-level 和 core-level, 减少内存访问冲突
- 灵活性: core-level的压缩可以支持不同层的计算



# 神经网络量化与AI处理器——NV Sparse TensorCore

- 2020年，英伟达推出Ampere架构，支持稀疏张量计算
- Ampere架构支持固定的2:4稀疏模式，即每4个元素中，2个为0，并且能够实现2倍的加速比
- 该稀疏模式可以拓展为更加普适的N:M稀疏，即每M个元素中，N个为0



Reference: [NV Technical Blog](#)

# 目录

CONTENTS



01. 典型AI芯片架构
02. AI芯片软硬件协同设计
03. AI大模型芯片设计
04. 存算一体AI芯片架构

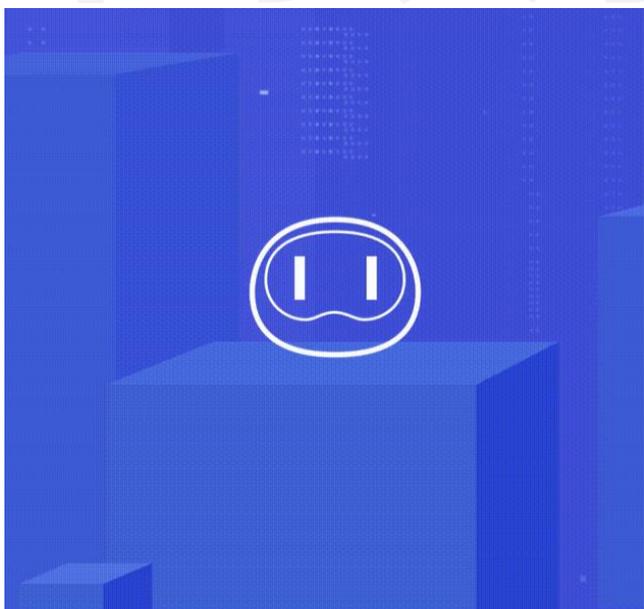
主流大模型

主干网络

Transformer

聊天机器人

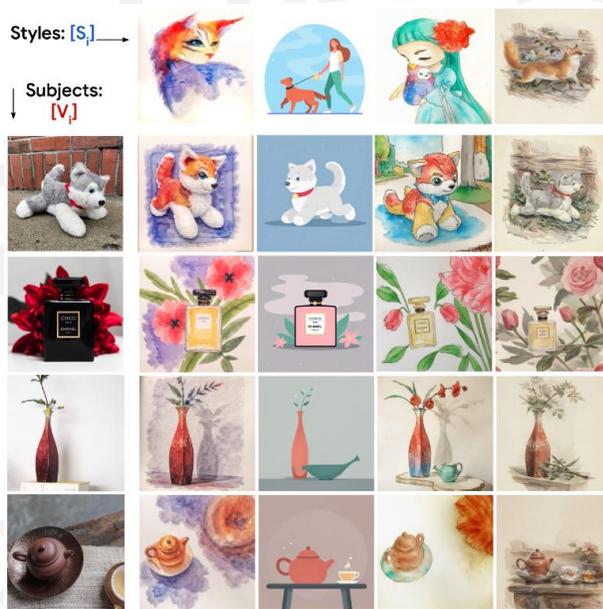
实时聊天 + 性格设定



<https://ai.baidu.com/unit/v2/static/socialbot>

图像合成

更改图片风格



Shah, Viraj, et al. arXiv:2311.13600 (2023).

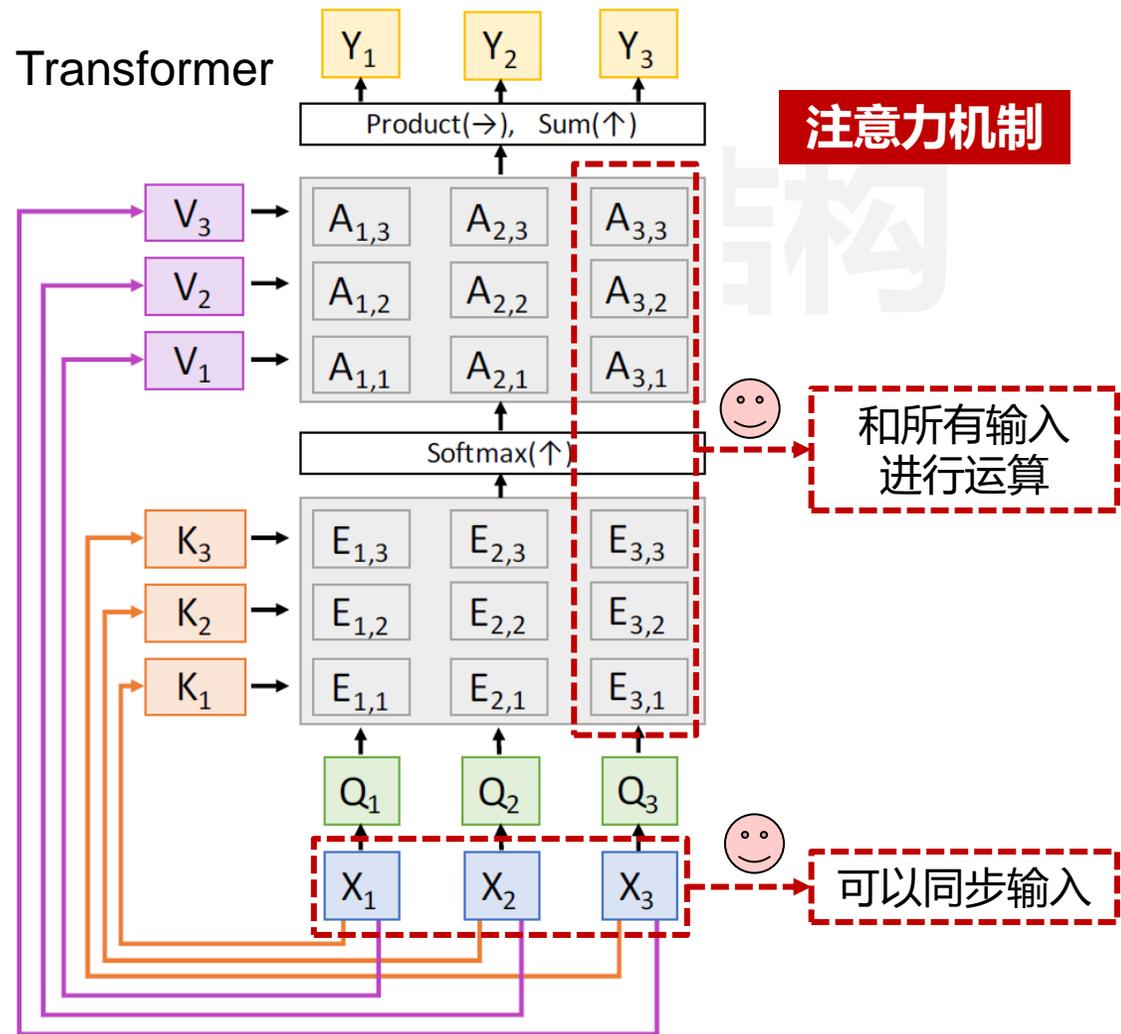
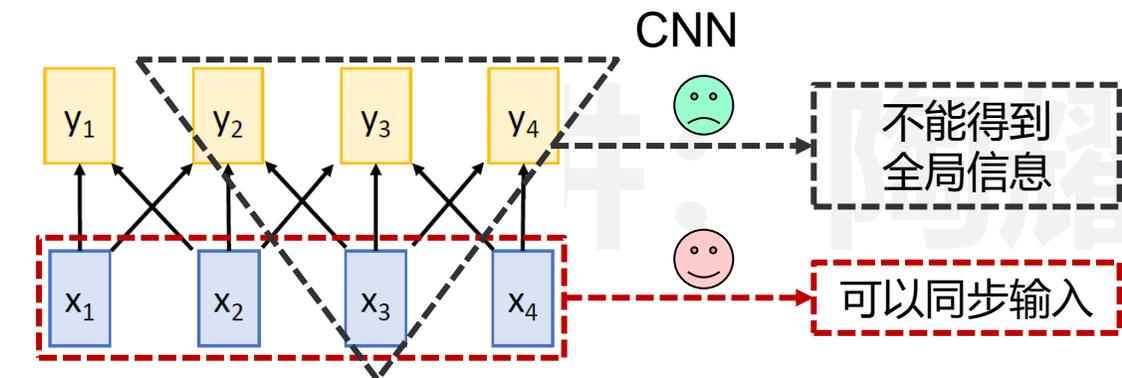
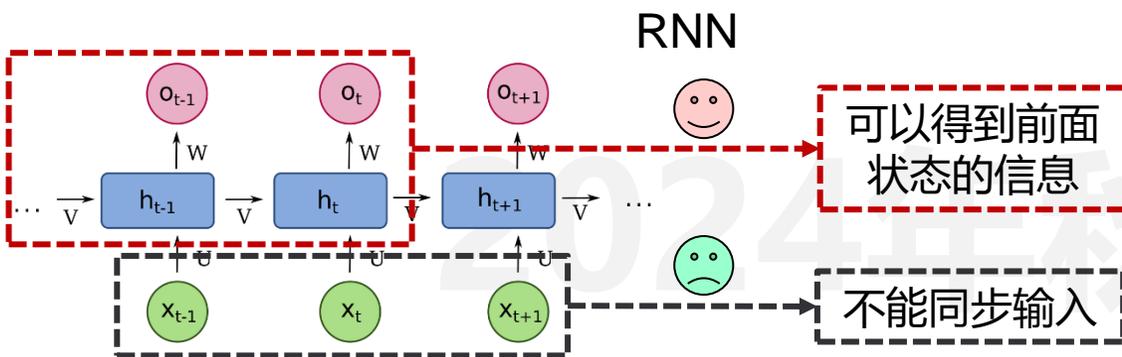
文本视频关联

根据文字更改视频内容

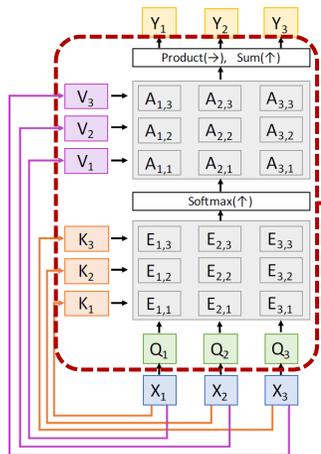


<https://pika.art>

指标	RNN	CNN	Transformer
数据特征	有序序列	多维网格	一组向量
长序列处理	✓	X	✓
可并行性	X	✓	✓



# AI大模型加速芯片架构



时间复杂度

$$O(n^2 \times d)$$

$n$ : 序列长度

$d$ : 维度

序列长度(32K)增大

任务种类增多 → 字典维度增大

注意力模块  
延时增加

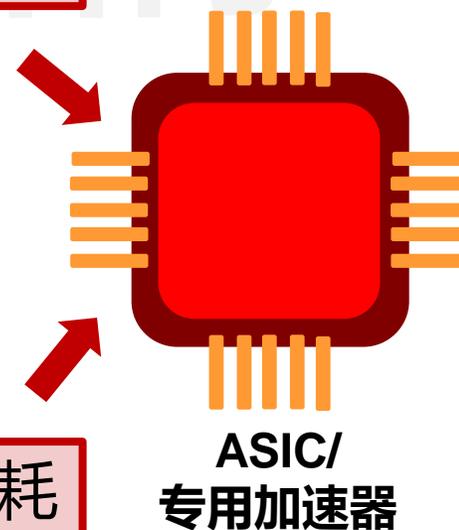
高延时

训练延时随模  
型增大增加

输入维度增大  
导致能耗增加

高能耗

能耗随模型  
增大增加

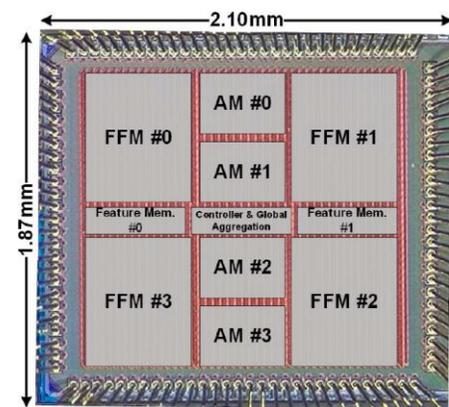
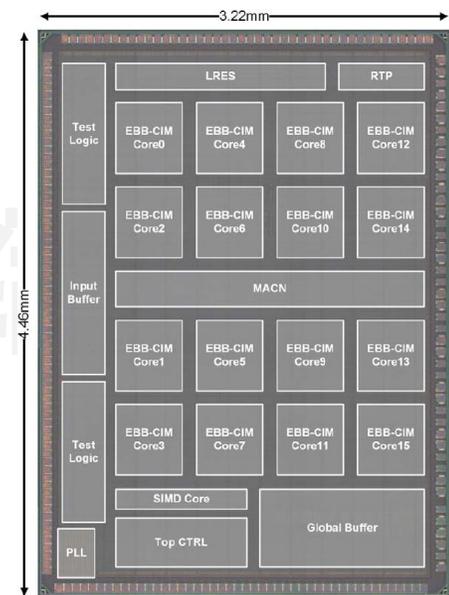


模型名称	训练时间(GPU-hours)	总能耗(MWh)	碳排放量(tCO <sub>2</sub> eq)
OPT-175B	809472	365	137
BLOOM-175B	1082880	475	183
LLaMA	7B	82432	36
	13B	135168	59
	33B	530432	233
	65B	1022362	449
LLaMA-2	7B	184320	83
	13B	368640	162
	34B	1038336	396
	70B	1720320	748

From LLaMA/LLaMA 2, 数据收集于A100-80GB

# AI大模型加速芯片架构

工作	单位	存储介质	存算模式	加速模型	流片/仿真
ReTransformer	杜克大学	ReRAM	存内计算	Transformer	仿真
ISSCC'2023	复旦大学	SRAM	近存计算	Transformer	流片
ReBERT	KAIST	ReRAM	存内计算	BERT	仿真
TranCIM	清华大学	SRAM	近存计算	BERT	流片
X-Former	普渡大学	SRAM/RRAM	近存/存内计算	BERT	仿真
P3ViT	澳门大学	SRAM	存内计算	ViT	流片
H3DAtten	佐治亚理工	SRAM/RRAM	近存/存内计算	Swin Transformer	仿真
MuTCIM	清华大学	SRAM	近存计算	MML Transformer	流片
TransPIM	加利福尼亚大学	DRAM	近存计算	BERT+GPT	仿真



Liu, Shiwei, et al. /ISSCC 2023

Tu, Fengbin, et al. /ISSCC 2023

高延时&高能耗

ASIC/专用加速器

挑战1

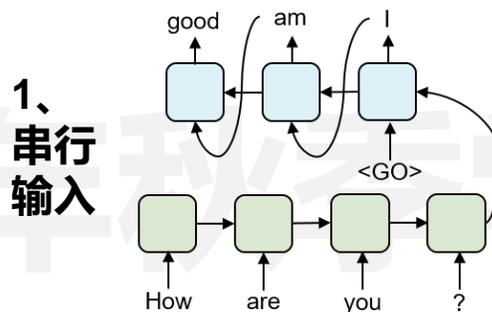
模型参数量大



模型级加速  
减少计算量

挑战2

数据调度和传统深度神经网络不同



模块级加速  
高效数据调度

挑战3

缺乏对专有算子的支持

主要线性算子	非线性算子
矩阵向量乘法(MVM)	GELU
矩阵矩阵乘法(MatMul)	LayerNorm
残差相加	Softmax
输入和权重不固定	动态归一化

算子级加速  
加速专有算子

## 基本结构:

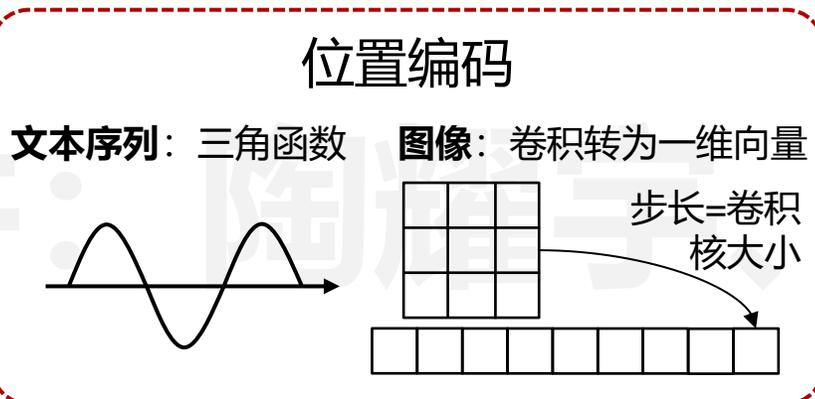
- 多个编码器层-多个解码器层

## 输入序列:

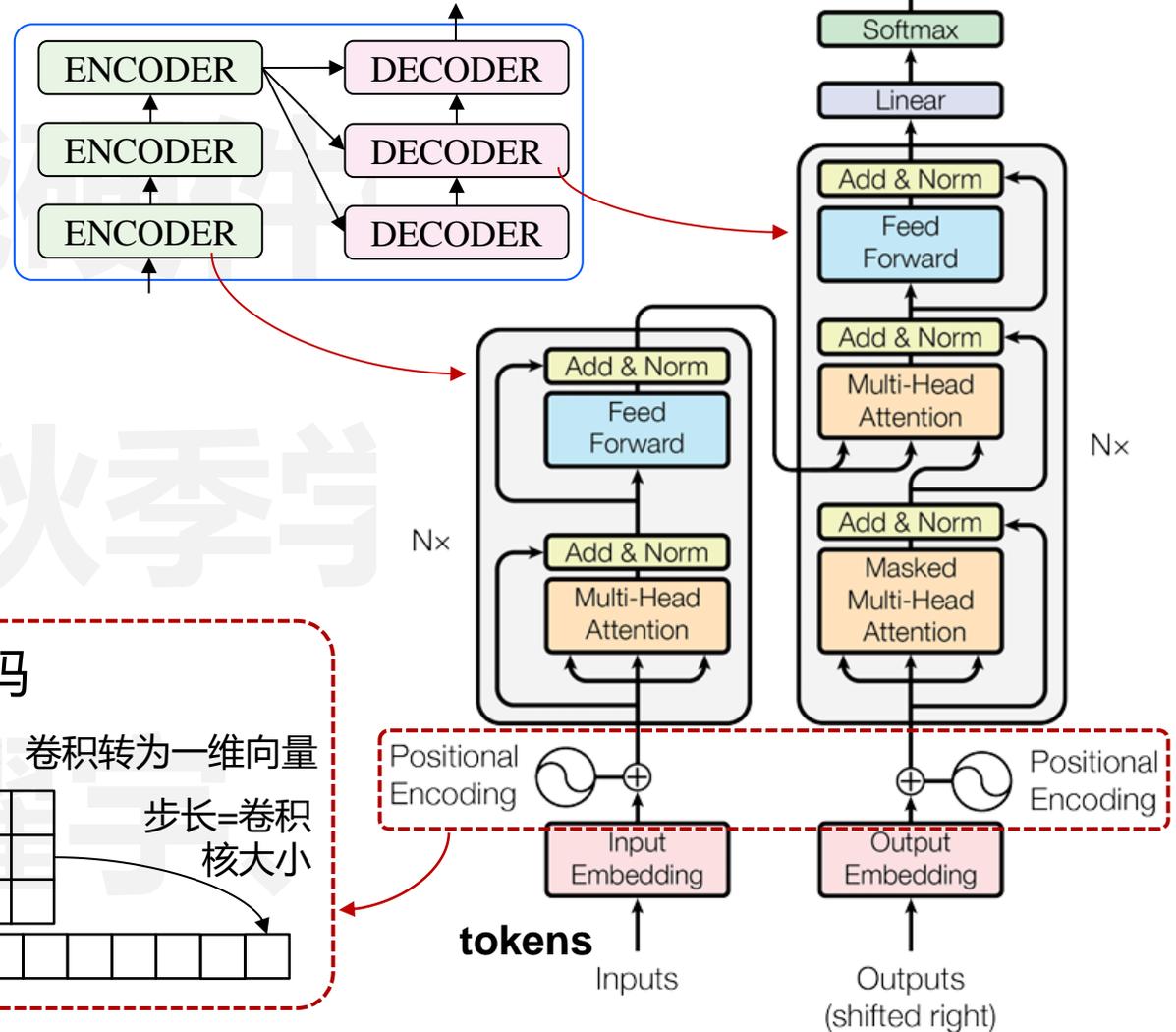
- 由一组token组成
- 每个token为一个向量表示一个词根
- 序列开头[CLS]——句子向量表示
- 输入模型前经过embedding层

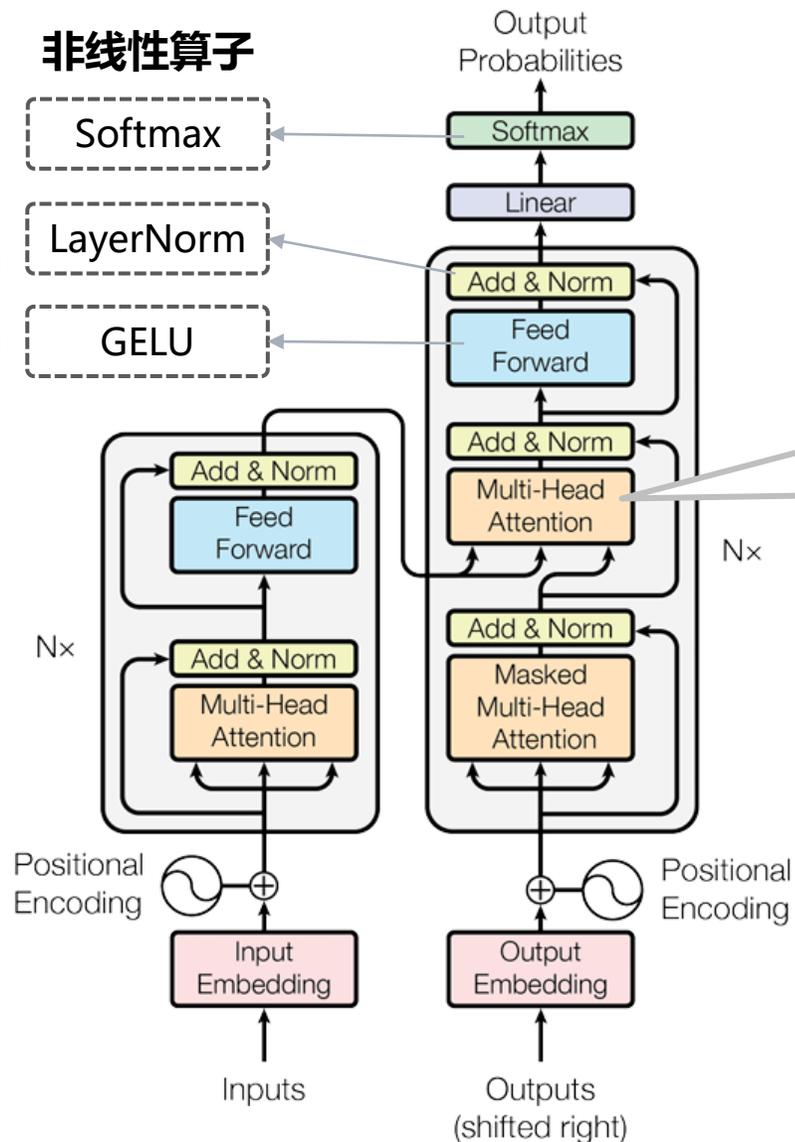
## embedding:

- 输入编码: 字典信息
- 位置编码: 位置信息



基本结构: 多编码器-多解码器

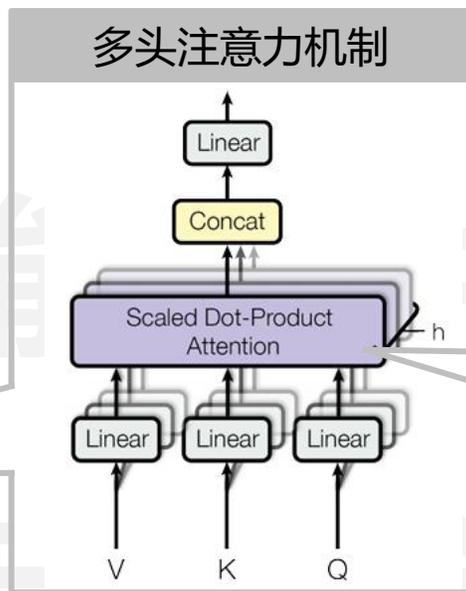




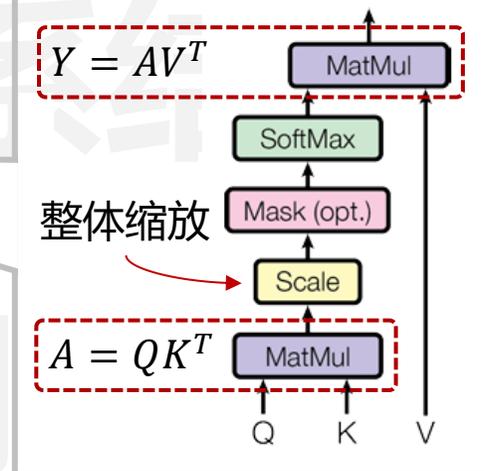
## 非线性算子

- Softmax
- LayerNorm
- GELU

## 多头注意力机制



## 缩放点乘注意力机制

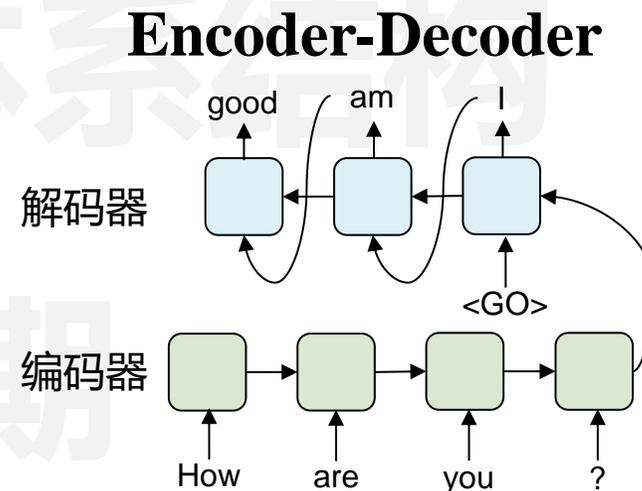
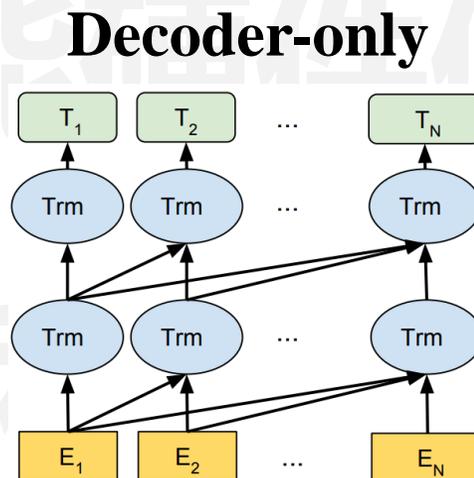
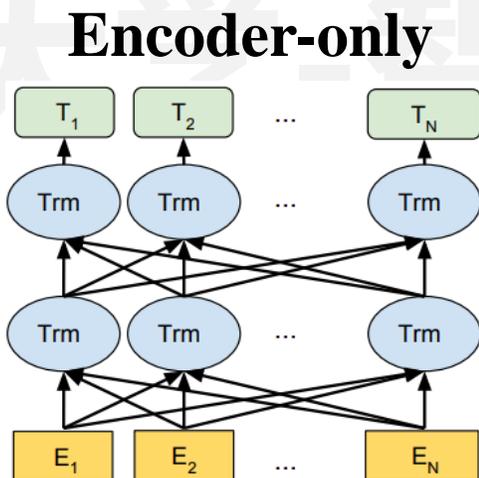


层级	模块级	主要线性算子	非线性算子
embedding	注意力模块	矩阵向量乘法(MVM)	GELU
编码器	前馈神经网络	矩阵矩阵乘法(MatMul)	LayerNorm
解码器	残差模块	残差相加	Softmax

## • 根据网络结构可分为：

- Encoder-only; Decoder-only; Encoder-Decoder

低延时要求  
输入依赖于上一次的输出  
流水设计  
分阶段并行设计



常见模型	BERT, RoBERTa	GPT, LLAMA	T5, BART
适合任务	分类任务	问答任务	翻译任务
输入输出特征	所有token并行输入	输入依赖于上一次的输出	输出强依赖于输入
适配加速器特征	<b>高吞吐</b>	<b>低延时</b>	<b>低延时</b>
适配并行设计	<b>流水并行设计</b>	<b>分阶段并行设计</b>	<b>分阶段并行设计</b>

- 模型级加速

- 模型压缩

- 稀疏性处理

- 模块级加速

- 注意力模块加速

- 前馈神经网络加速

- 算子级加速

- 线性算子加速

- 非线性算子加速

北京大学-智能硬件体系结构

2024年秋季学期

主讲：陶耀宇、李萌

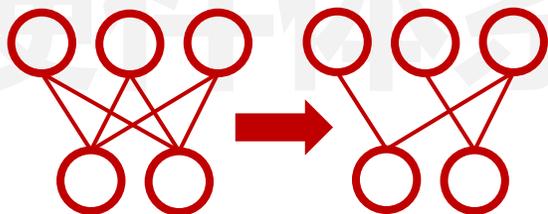
## • 模型压缩：硬件部署前算法优化

① **模型量化**



- 对**激活/权重/输出**进行量化
- 降低模型**权重存储**需求
- 降低硬件**计算精度**要求

② **静态剪枝**



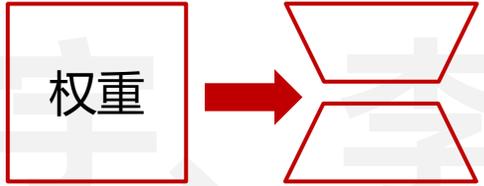
- 删去**不重要的**注意力维度/头
- 降低硬件**计算量**

③ **权重复用**



- **重复利用**网络中的某一部分
- 降低模型**权重存储**需求
- 降低**硬件设计**需求

④ **低秩分解**



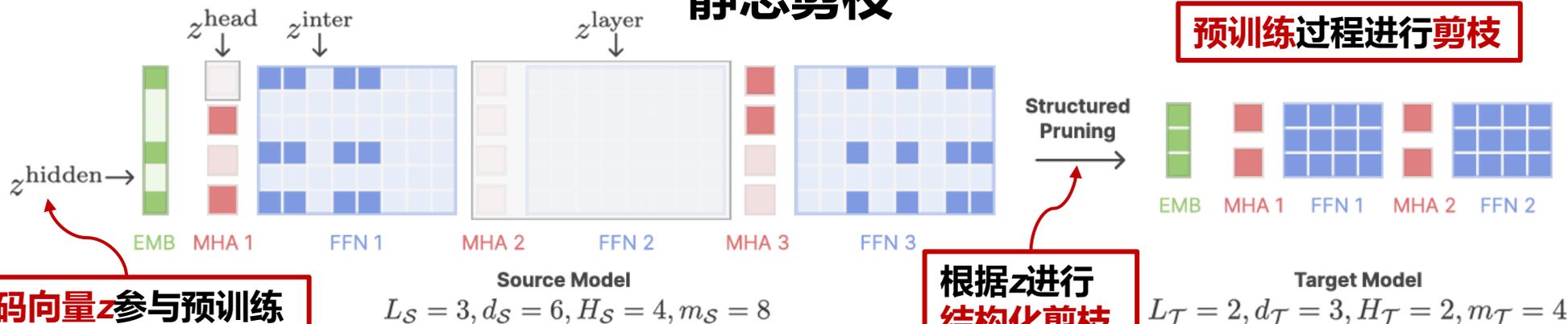
- 将权重矩阵拆分为几个低秩矩阵的组合
- 降低模型**训练存储**需求

## 模型量化

Q-BERT	TernaryBERT	I-BERT	BinaryBERT	PTQ4ViT	FQ-ViT	BiT
INT8	三值	INT8	二值	INT8	INT4	二值

目前部分运算比如全连接层可实现**三值/二值**精度下降**5%**以内，**INT8无损**

## 静态剪枝



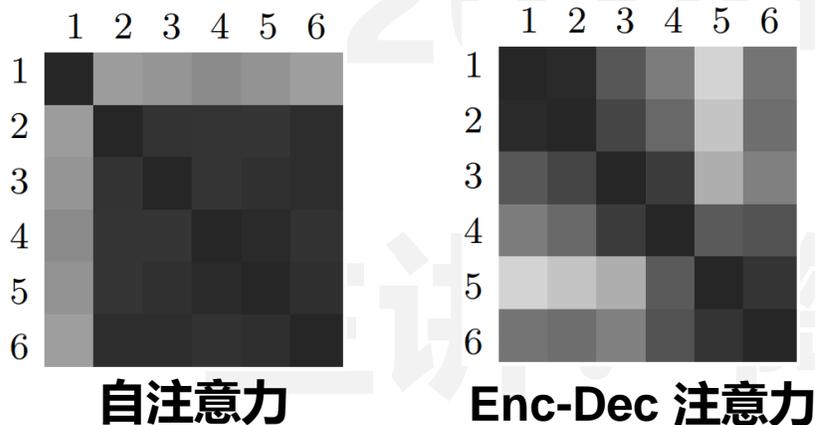
Xia, Mengzhou, et al. arXiv:2310.06694 (2023)

## • 权重复用

### 注意力模块权重相似

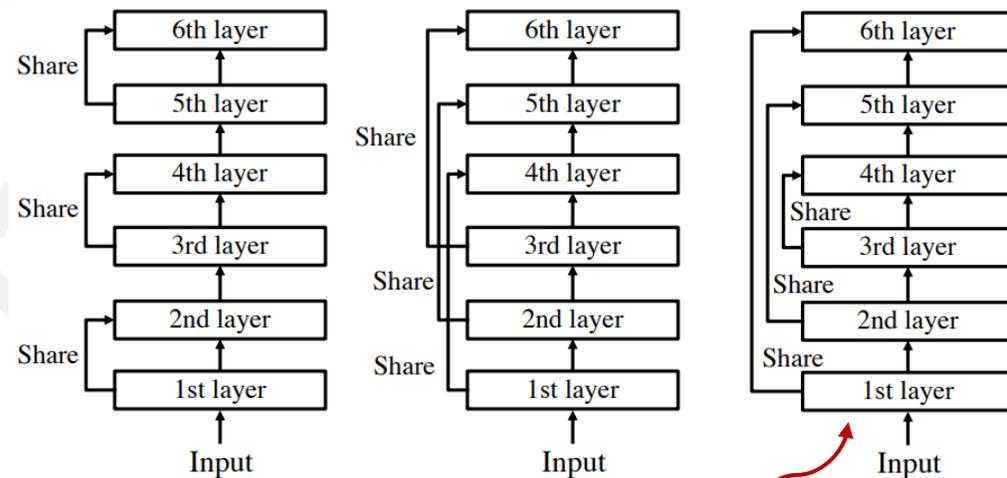
1、不同**维度**之间权重相似  
多个维度学习同样的东西

2、不同**层**之间权重相似  
不同层的权重的**JS**散度  
颜色越深代表相似度越高



### 层间/维度间权重复用

根据**权重相似性**进行复用



**ALBERT**  
——全层复用

低层有更多参数冗余

## • 低秩分解

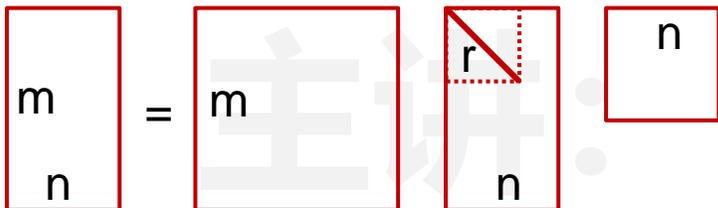
- 适用于**微调训练**过程
- 训练中的权重更新为**低秩**矩阵，可以用两个更小矩阵来表示
- 拆分更新权重降低**训练显存**

### 矩阵低秩分解

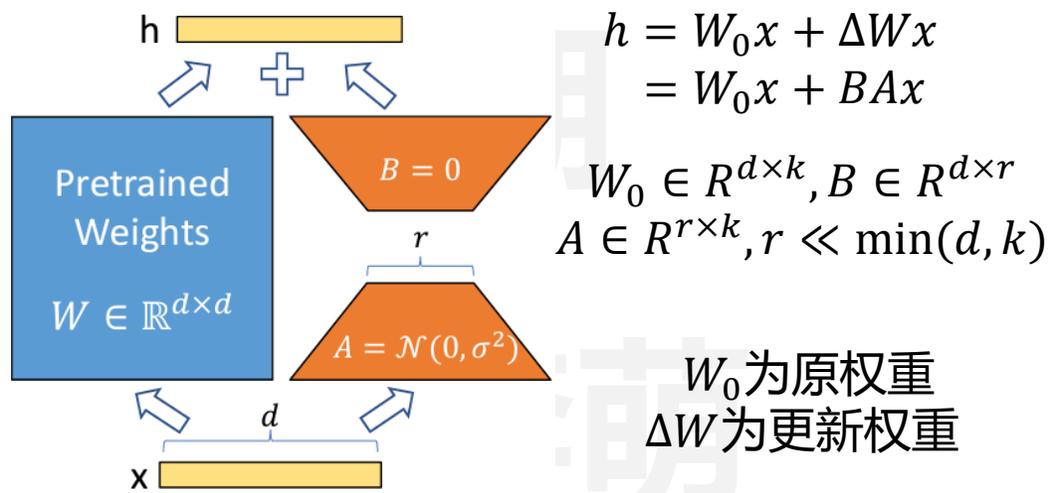
□ 奇异值分解SVD

$$A = U\Sigma V^T$$

仅有 $r$ 宽度的对角线上为特征值



### 微调训练应用



- **模型级加速**

- 模型压缩
- 稀疏性处理

- **模块级加速**

- 注意力模块加速
- 前馈神经网络加速

- **算子级加速**

- 线性算子加速
- 非线性算子加速

北京大学-智能硬件体系结构

2024年秋季学期

主讲：陶耀宇、李萌