

大规模异构计算

——深度学习背后的驱动力量

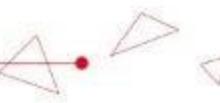
2019.12.05 颜深根 执行研究总监 商汤科技

目录

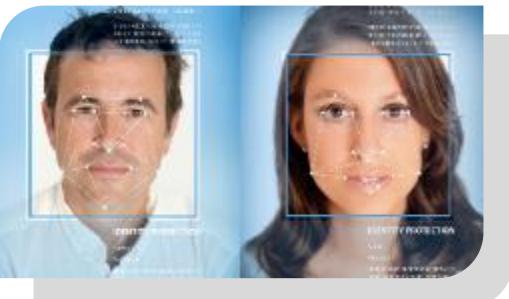
1 //
2 //
3 //
4 //
5 //
6 //

- AI复兴：深度学习
 - 计算优化：GPU
 - 计算以外：内存&访存
 - 携手共进：通信优化
 - 无所不能：云上AI
 - 轻松一刻：应用案例

各项人工智能技术取得突破：深度学习



人脸识别



语音识别



围棋博弈



图像分类



深度学习

深度学习让机器从大量资料中归纳总结经验
这一套技术体系让机器训练学会特定技能



自动驾驶

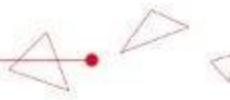


智能投顾、量化交易

语言理解

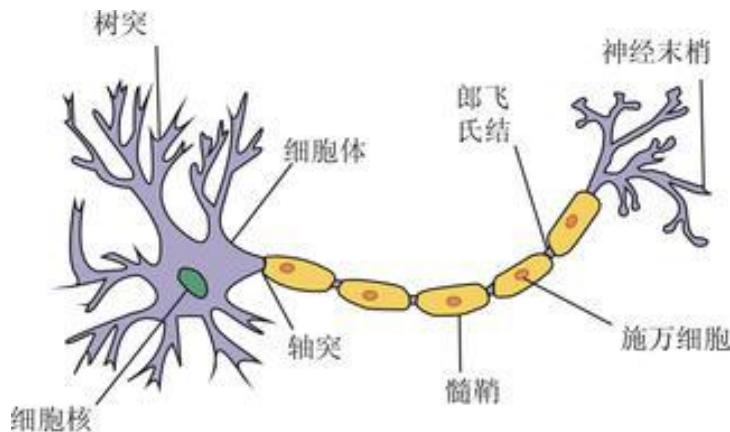


深度学习的基本思想

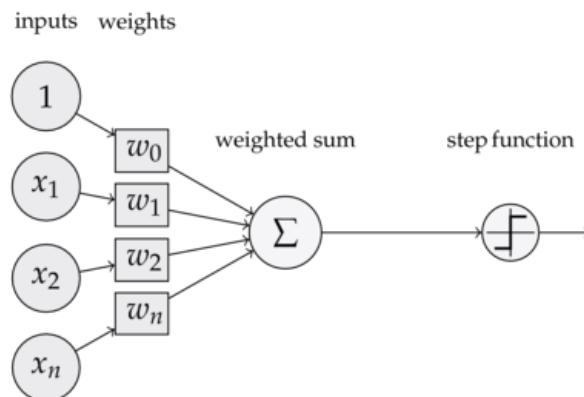


● 人工神经网络

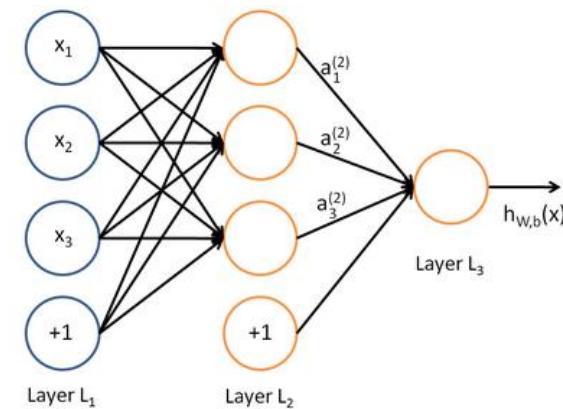
- 源于对人脑神经元运行机理的研究，始于20世纪40年代
- 将神经的突触连接结构用抽象化的模型表达
- 将多个单一神经元结构相互连接得到复杂的网络结构



生物神经元模型

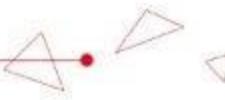


感知器(Perceptron)



多层感知器

BP(Back Propagation)神经网络



● Back Propagation原理

- Rumelhart, McClelland于1985年提出了BP网络的误差反向后传BP(Back Propagation)学习算法



- 利用输出后的误差来估计输出层的直接前导层的误差，再用这个误差估计更前一层的误差，如此一层一层的反传下去，就获得了所有其他各层的误差估计。⁵

深度学习：深层神经网络

● 深层神经网络

■ 普适逼近原理(Universal approximation theorem):

- 单隐藏层神经网络可以逼近任意函数

■ 神经网络的表达能力

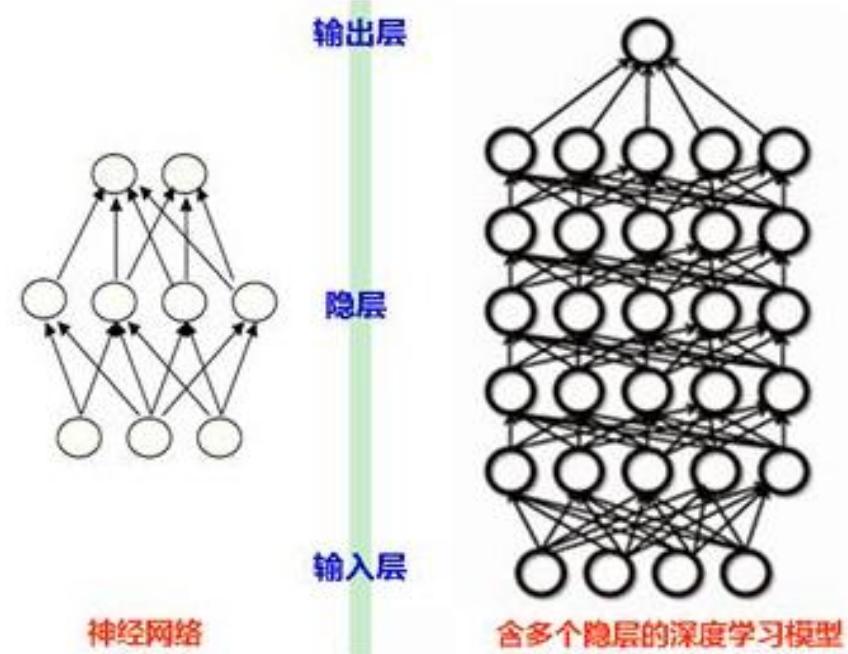
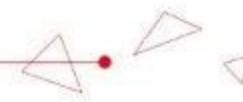
- 参数越多表达能力越强
- 层数越多表达能力越强
- 过拟合问题

■ 深层神经网络泛化能力更强，训练难度更大

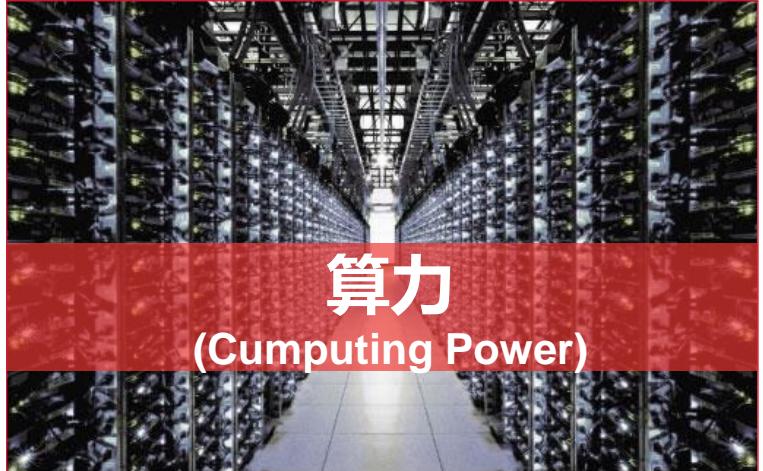
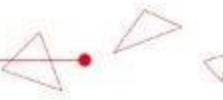
- 逐层训练，改进训练算法
- 增加训练样本，增加迭代次数

■ 深度学习

- 深层神经网络
- 通过海量计算能力，海量训练样本进行训练

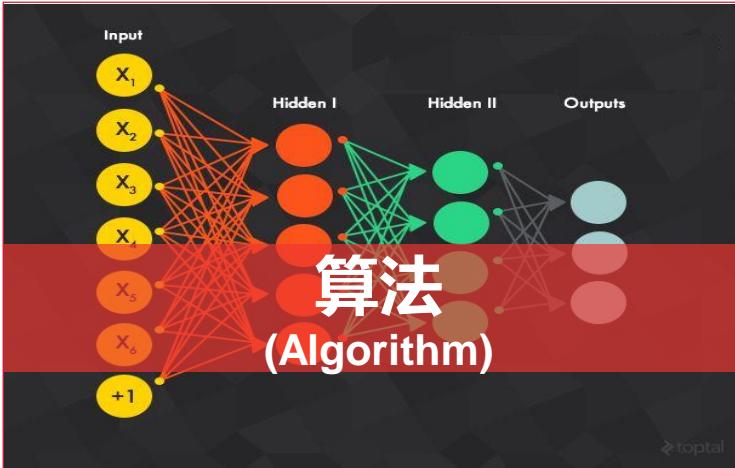


深度学习的特点



算力 (Computing Power)

大量高性能硬件组成的计算能力。
如：分布式GPU集群



算法 (Algorithm)

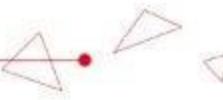
优秀的人工智能算法，比如现在最流行的深度学习相关的DNN、CNN、RNN、GAN...



数据 (Data)

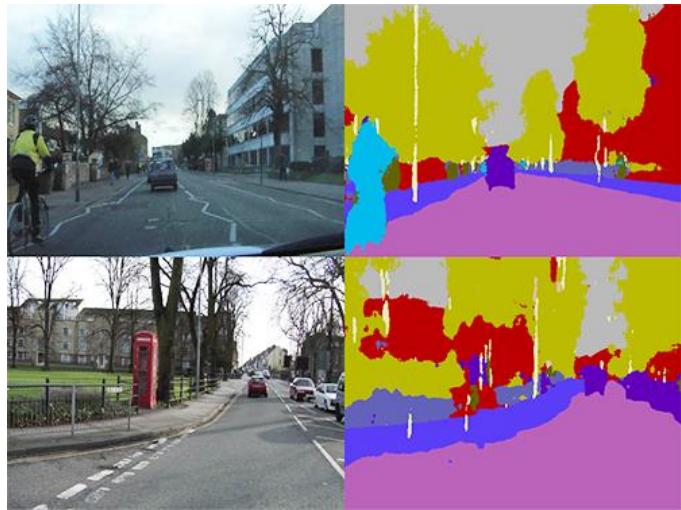
数据是驱动人工智能取得更好的识别率和精准度的核心因素

深度学习常用的几种模型

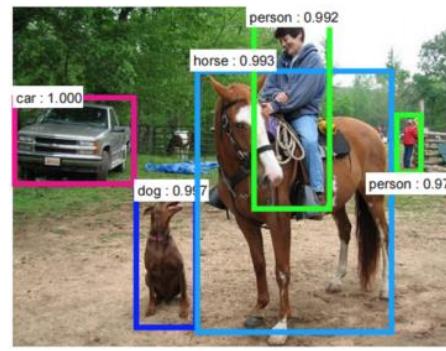


- 自动编码器：Auto-Encoder
- 卷积神经网络：Convolutional Neural Networks(CNN)
- 生成式对抗网络：Generative Adversarial Networks(GAN)
- 循环神经网络：Recurrent Neural Networks(RNN)
- 受限玻尔兹曼机：Restricted Boltzmann Machine(RBM)
- 深度信任网络：Deep Belief Networks(DBN)

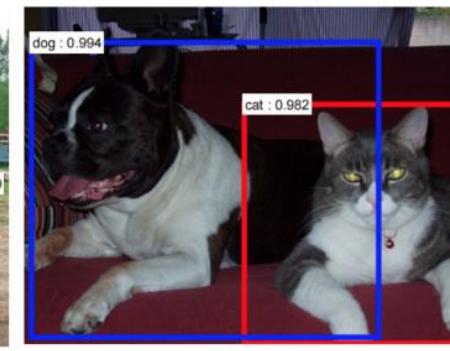
卷积神经网络



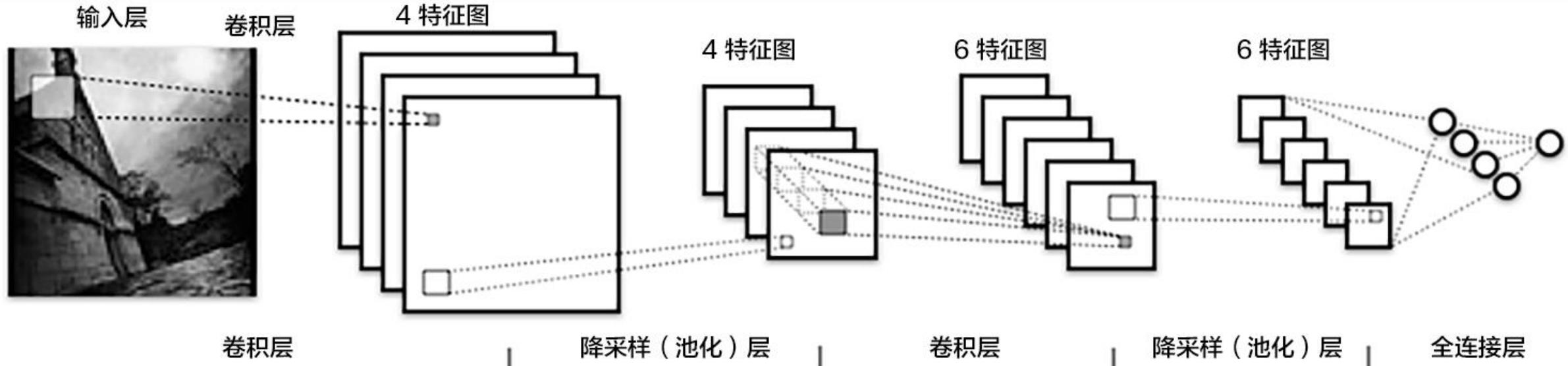
Segmentation



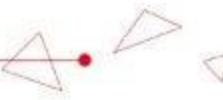
Detection



AutoDriving



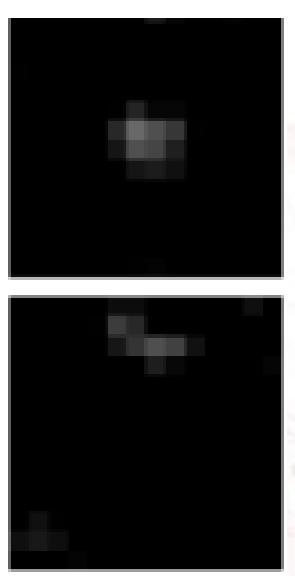
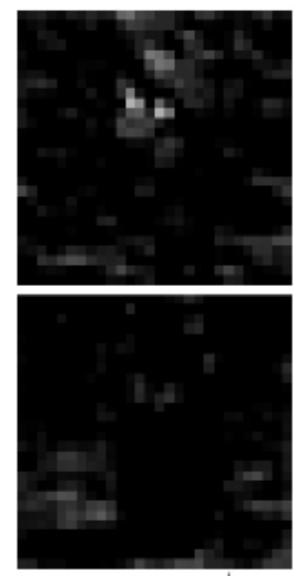
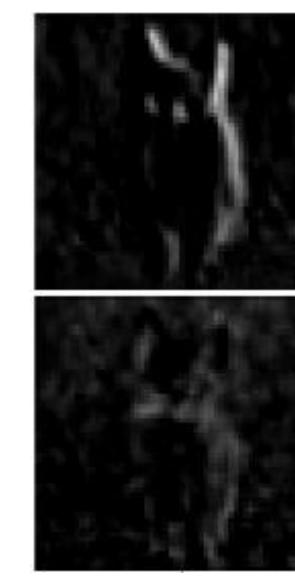
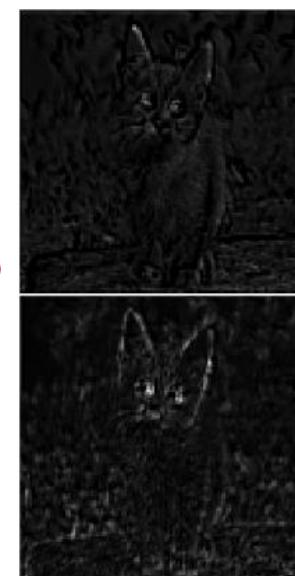
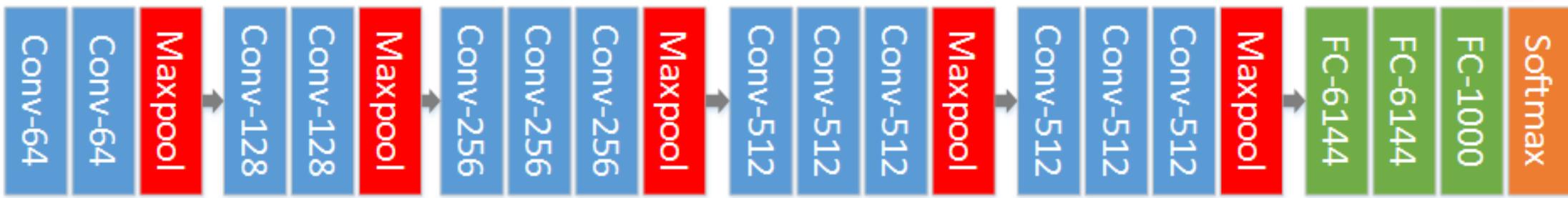
卷积神经网络的基本结构



- 数据预处理层
 - 数据读取，数据增强，数据采样等等
- 卷积层
 - 提取特征的滤波过程，在图像领域常表现为模板(mask)运算
- 激活函数
 - 增加网络的非线性性，拓展了网络的表达能力
 - 常见的激活函数有三种：Sigmoid函数，tanh函数，ReLU函数
- 池化层
 - 减少数据量，抑制过拟合，提高鲁棒性(max-pooling, average-pooling)
- 全连接层
 - 主要用于分类，现在很多网络也可以不需要全连接层
- Loss层
 - 与需要解决的问题相关，分类常用的Loss为Softmax

卷积神经网络：特征提取示例

输入



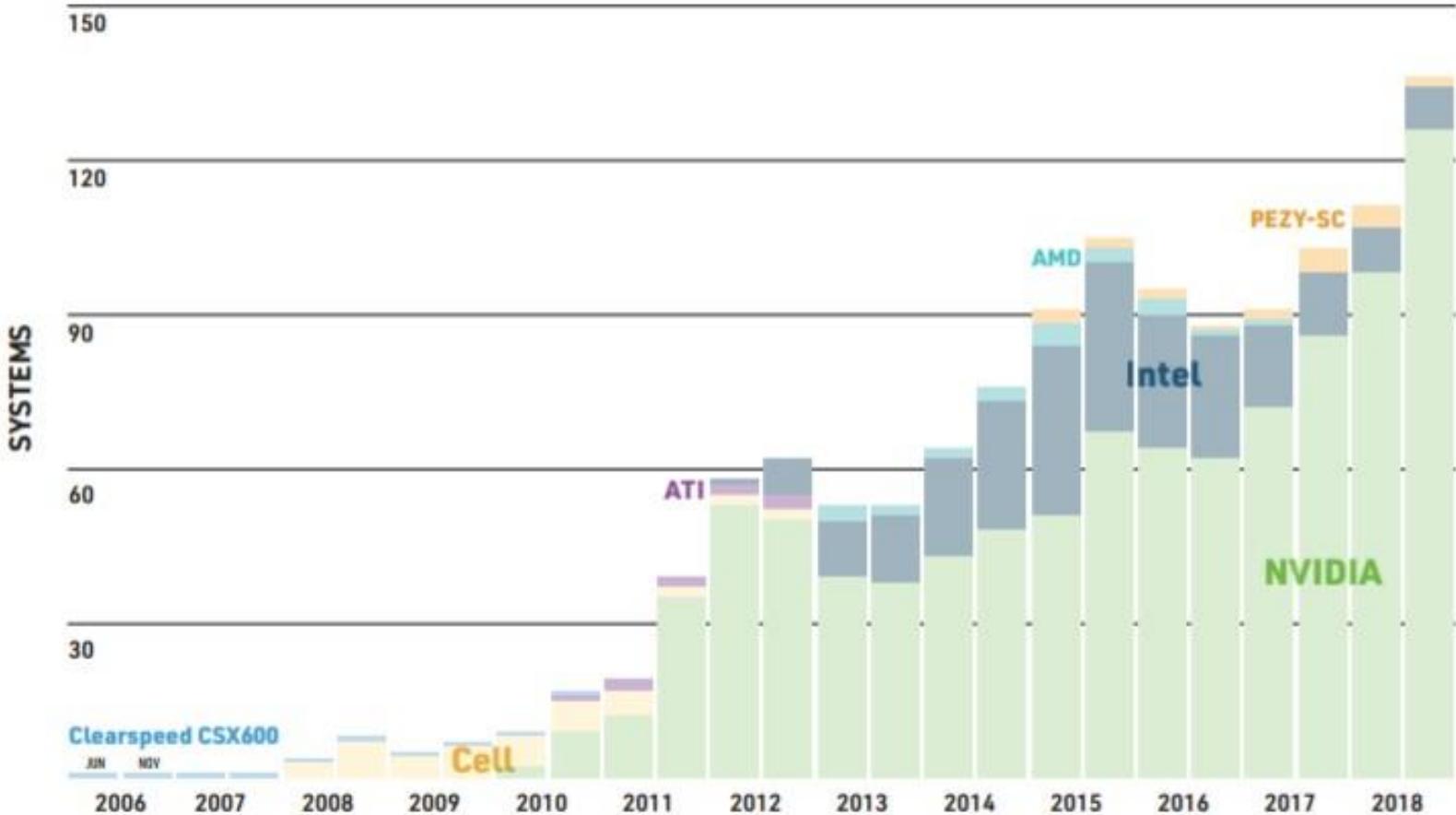
输出

目录

1 //
2 //
3 //
4 //
5 //
6 //

- AI复兴：深度学习
- 计算优化：GPU
- 计算以外：内存&访存
- 携手共进：通信优化
- 无所不能：云上AI
- 轻松一刻：应用案例

异构计算的兴起



近年HPC TOP500 厂家份额统计



CPU和GPU的特点

Nvidia Turing GPU架构

TU102 - FULL CONF

18.6 BILLION TRANSISTORS

SM

CUDA CORES

TENSOR CORES

RT CORES

GEOMETRY UNITS

TEXTURE UNITS

ROP UNITS

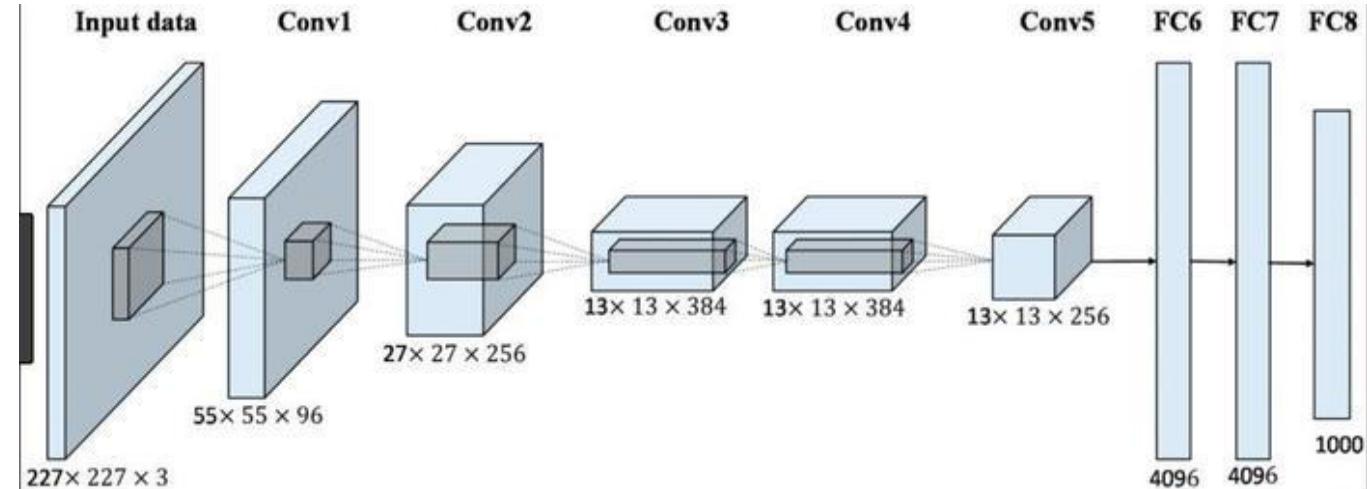
MEMORY

NVLINK CHANNELS



不同操作计算量统计

params	AlexNet	FLOPs
4M	FC 1000	4M
16M	FC 4096 / ReLU	16M
37M	FC 4096 / ReLU	37M
442K	Max Pool 3x3s2	
74M	Conv 3x3s1, 256 / ReLU	
1.3M	Conv 3x3s1, 384 / ReLU	112M
884K	Conv 3x3s1, 384 / ReLU	149M
	Max Pool 3x3s2	
	Local Response Norm	
307K	Conv 5x5s1, 256 / ReLU	223M
	Max Pool 3x3s2	
	Local Response Norm	
35K	Conv 11x11s4, 96 / ReLU	105M



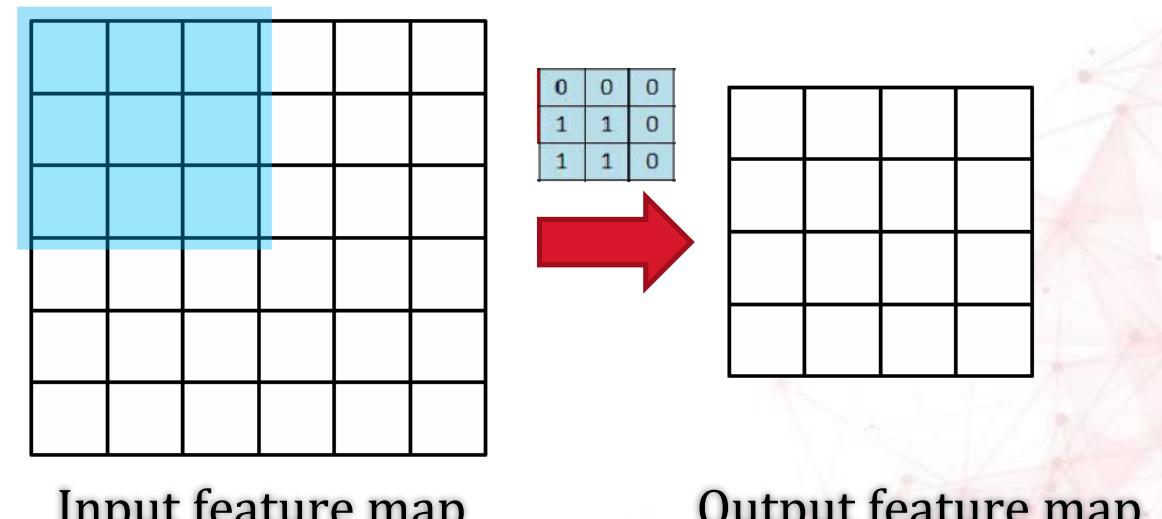
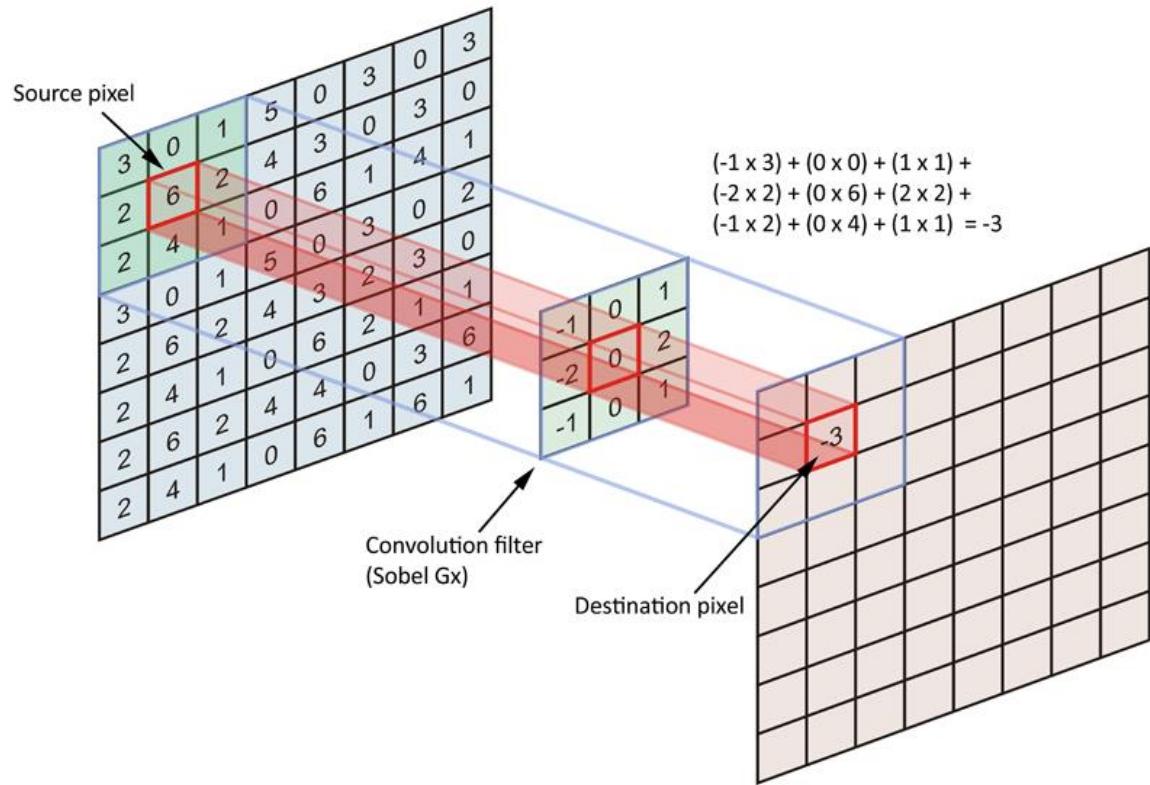
AlexNet	参数数量	参数占比	计算量	计算量占比
卷积层	~3M	5%	663M	92%
全链接层	~57M	95%	57M	8%

超过90%的计算量在卷积层

卷积计算方法

● 卷积计算公式

$$\text{Output}[i][j][n] = \sum_{m=0}^{M-1} \sum_{u=0}^{K-1} \sum_{v=0}^{K-1} \text{Input}[i * S + u][j * S + v][m] \times \text{Filter}[n][u][v][m]$$



Input feature map

Output feature map

整张图片输入输出（无Padding）

卷积计算方法



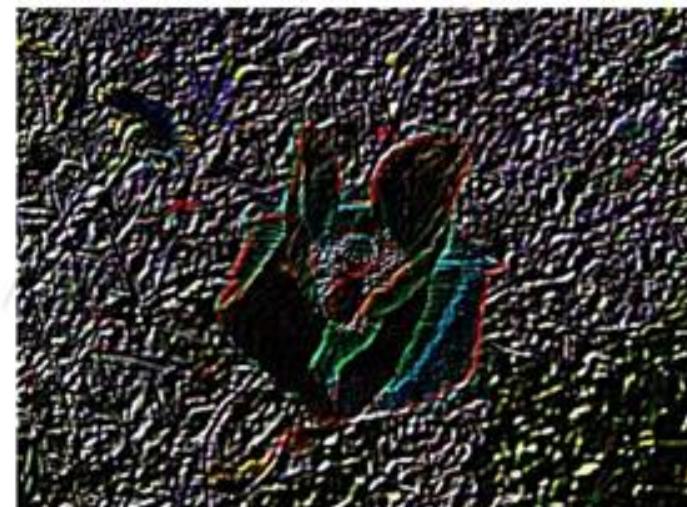
$$\ast \begin{bmatrix} -1 & -1 & -1 \\ -1 & 9 & -1 \\ -1 & -1 & -1 \end{bmatrix} =$$



示例1：图像锐化



$$\ast \begin{bmatrix} -1 & 0 & 0 & 0 & 0 \\ 0 & -2 & 0 & 0 & 0 \\ 0 & 0 & 6 & 0 & 0 \\ 0 & 0 & 0 & -2 & 0 \\ 0 & 0 & 0 & 0 & -1 \end{bmatrix} =$$



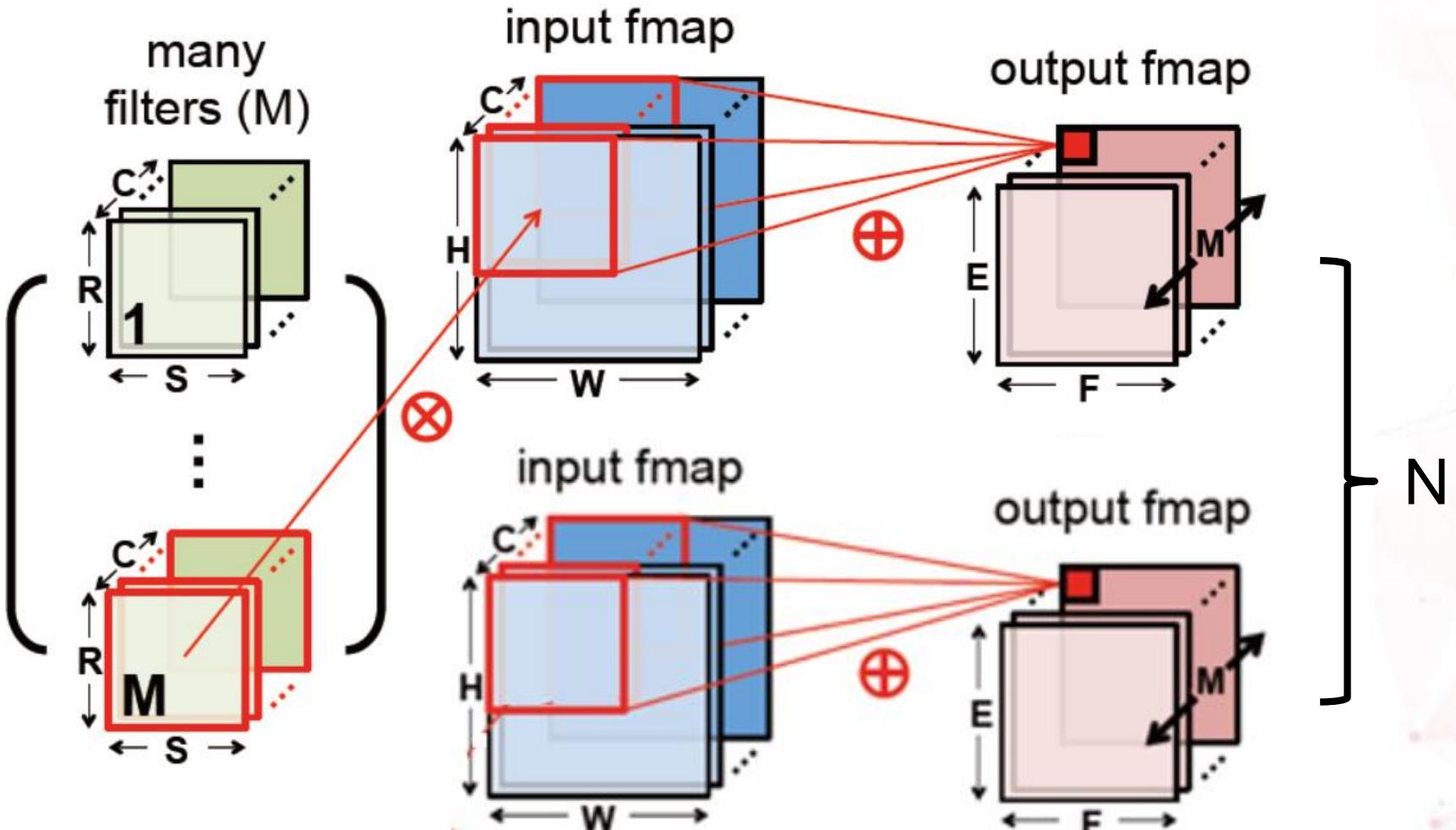
示例2：边缘检测

卷积神经网络中的卷积

- 输入 Feature Map
 - W: 宽度
 - H: 高度
 - C: 通道数
 - N: 样本数量

- Filter
 - S: 宽度
 - R: 高度
 - C: 通道数
 - M: 输出通道数

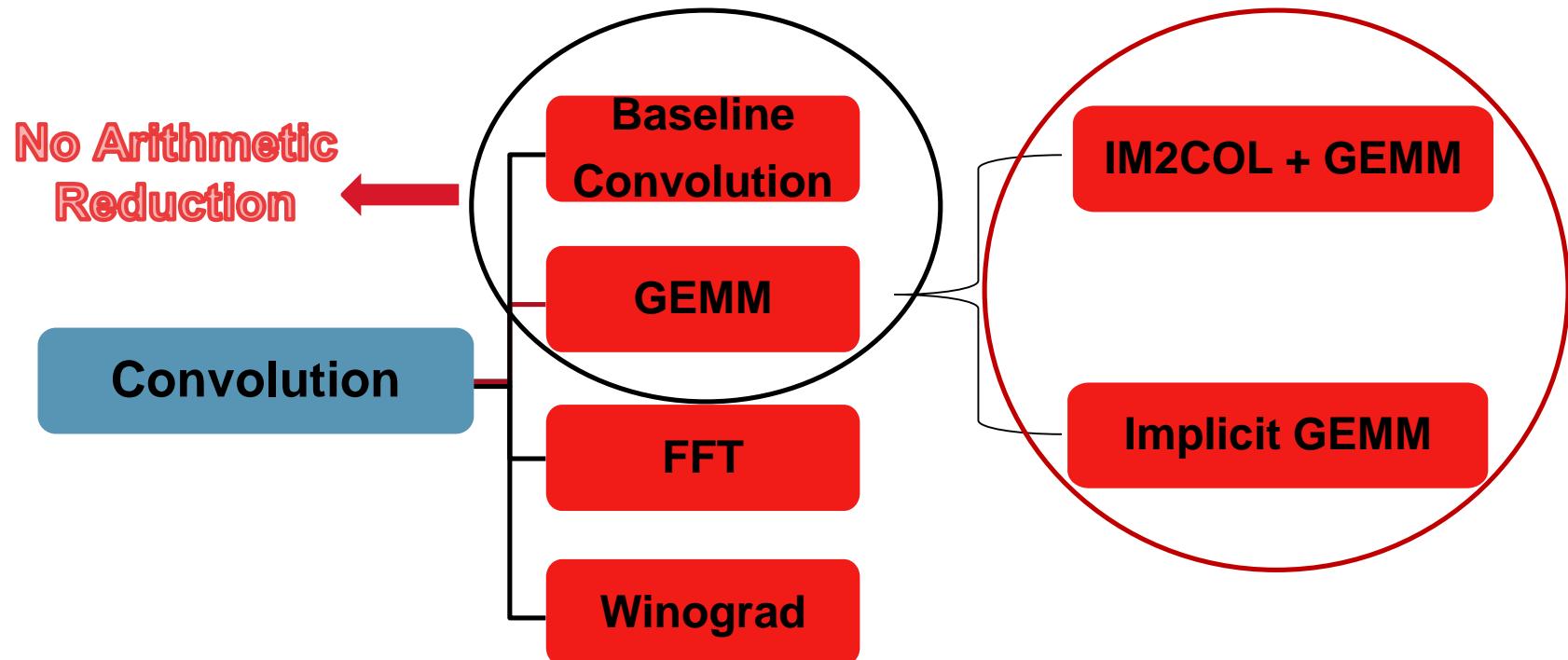
- 输出 Feature Map
 - F: 宽度
 - E: 高度
 - M: 通道数
 - N: 样本数量



$$\text{计算量} = 2 * F * E * S * R * C * M * N$$

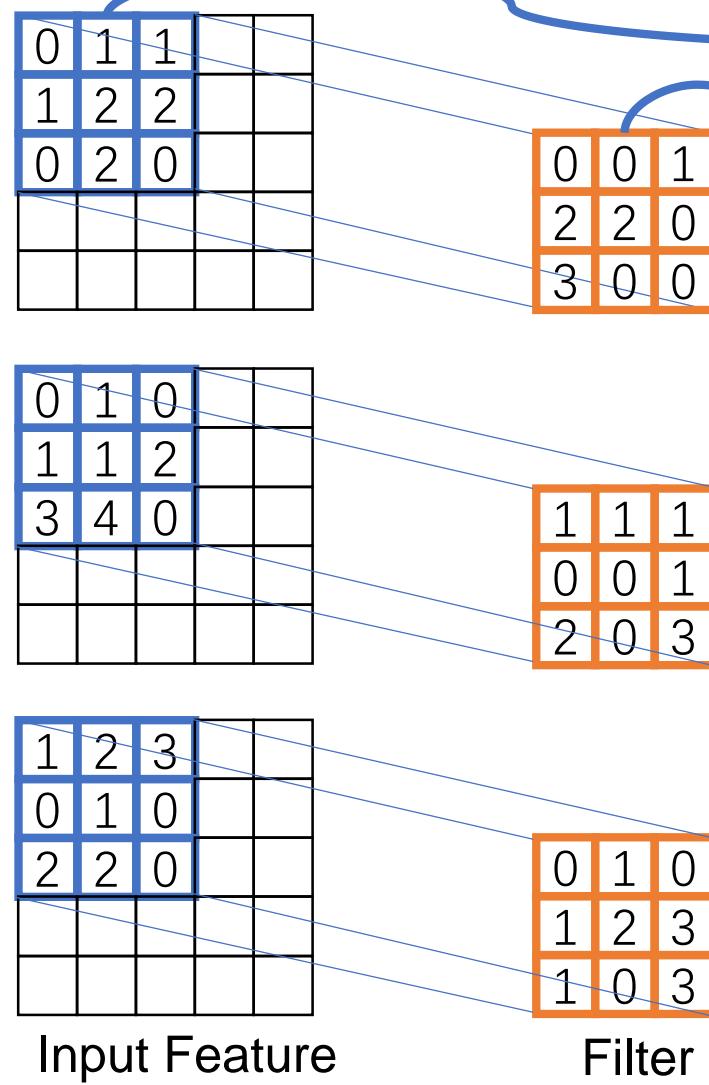
常用卷积计算方法

● 四种卷积计算方法

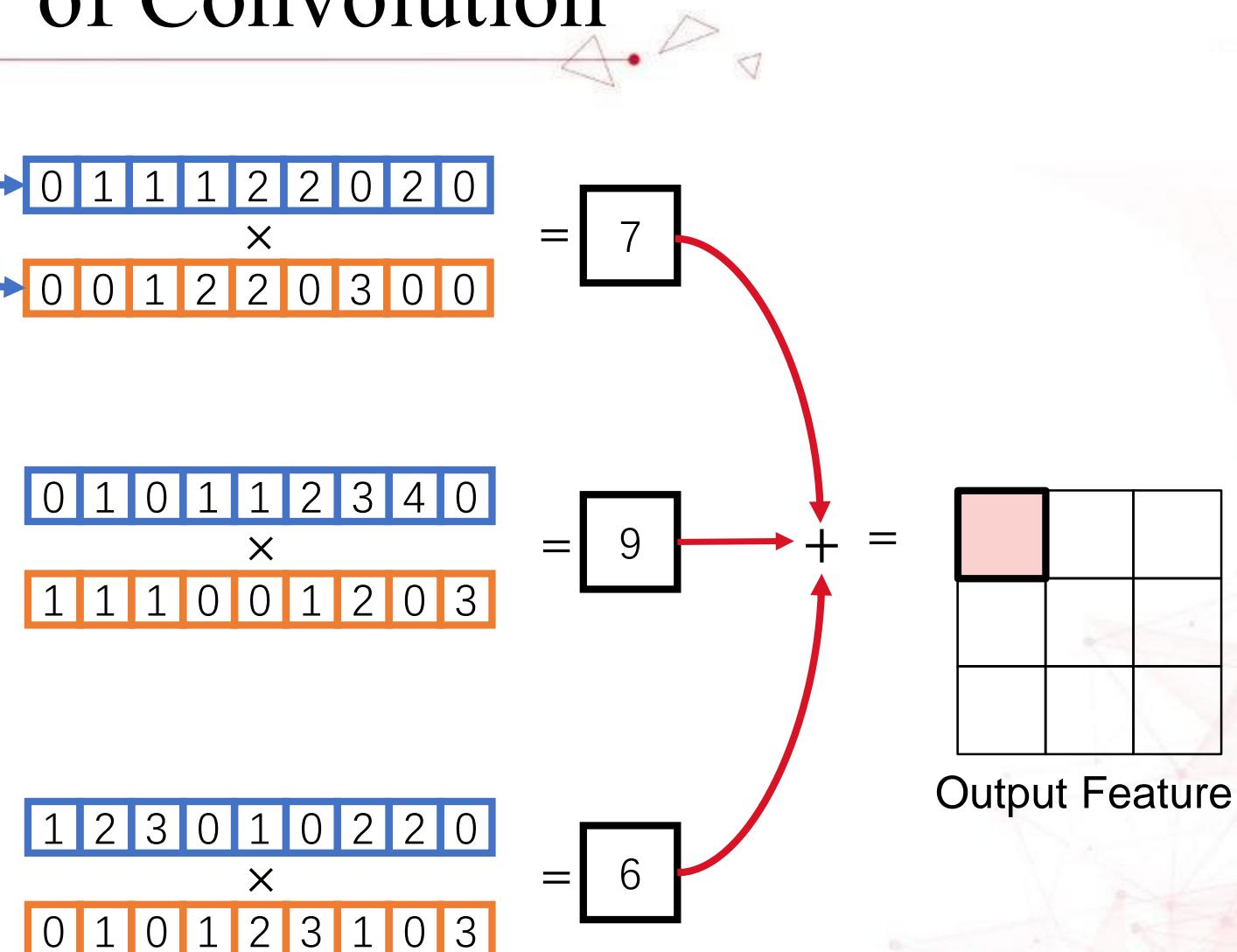


案例：GEMM优化

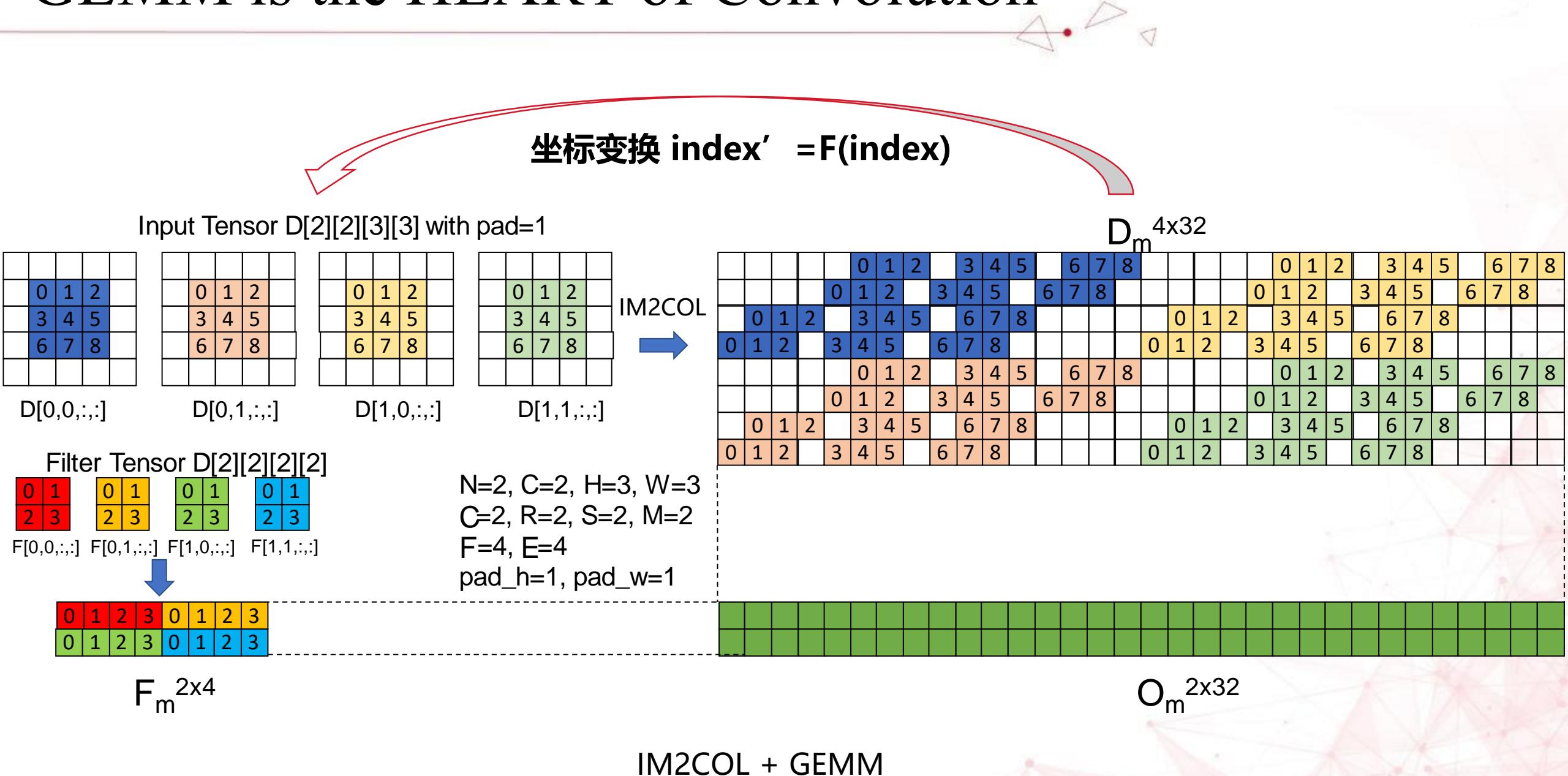
GEMM is the HEART of Convolution



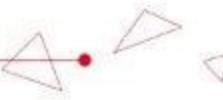
IM2COL + GEMM



GEMM is the HEART of Convolution



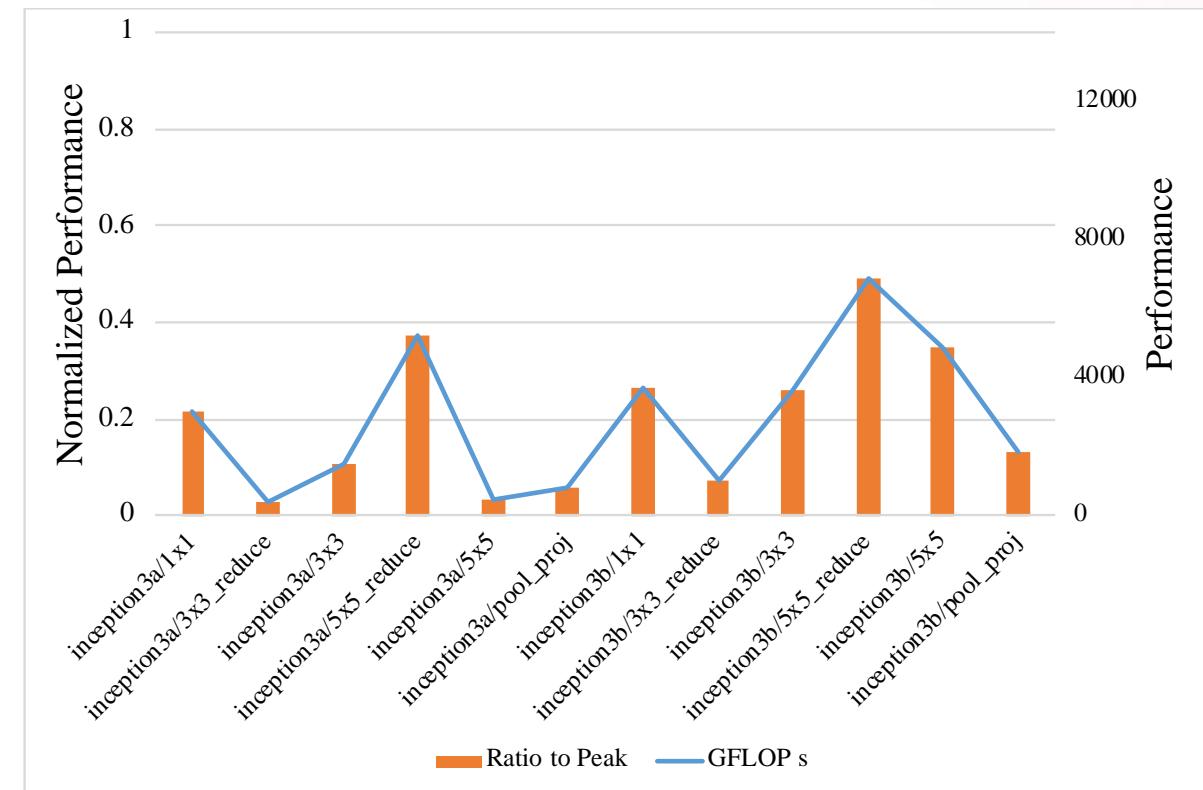
Motivation: Small Matrix Problem



Nvidia Tesla V100: FP32 Peak performance 14Tflops

	m	k	N (Batch=1)	GFLOPS	Ratio to Peak
inception3a_1	96	192	784	3011.18	0.215084
Inception3a_2	16	192	784	416.84	0.029774
Inception3a_3	64	192	784	1497.77	0.106984
Inception3a_4	128	864	784	5193.64	0.370974
Inception3a_5	32	400	256	470.9	0.033636
Inception3a_6	32	192	784	818.79	0.058485
Inception4a_1	128	256	784	3676.79	0.262628
Inception4a_2	32	256	784	1012.76	0.07234
Inception4a_3	128	256	784	3667.59	0.261971
Inception4a_4	192	1152	676	6910.33	0.493595
Inception4a_5	96	800	784	4888.3	0.349164
Inception4a_6	64	256	784	1845.13	0.131795

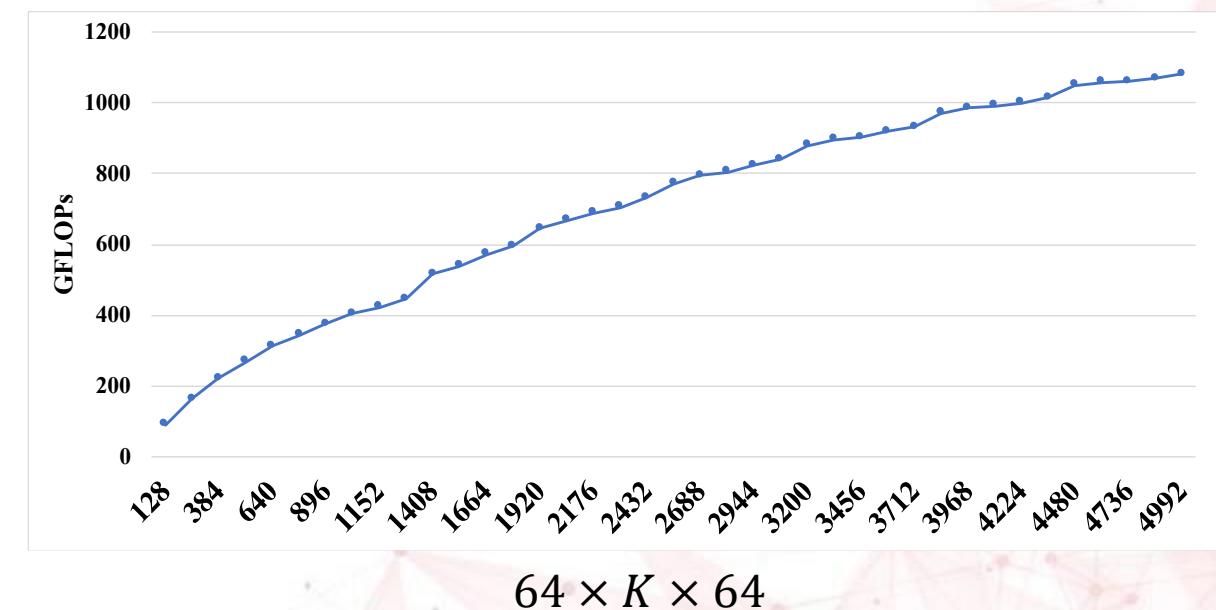
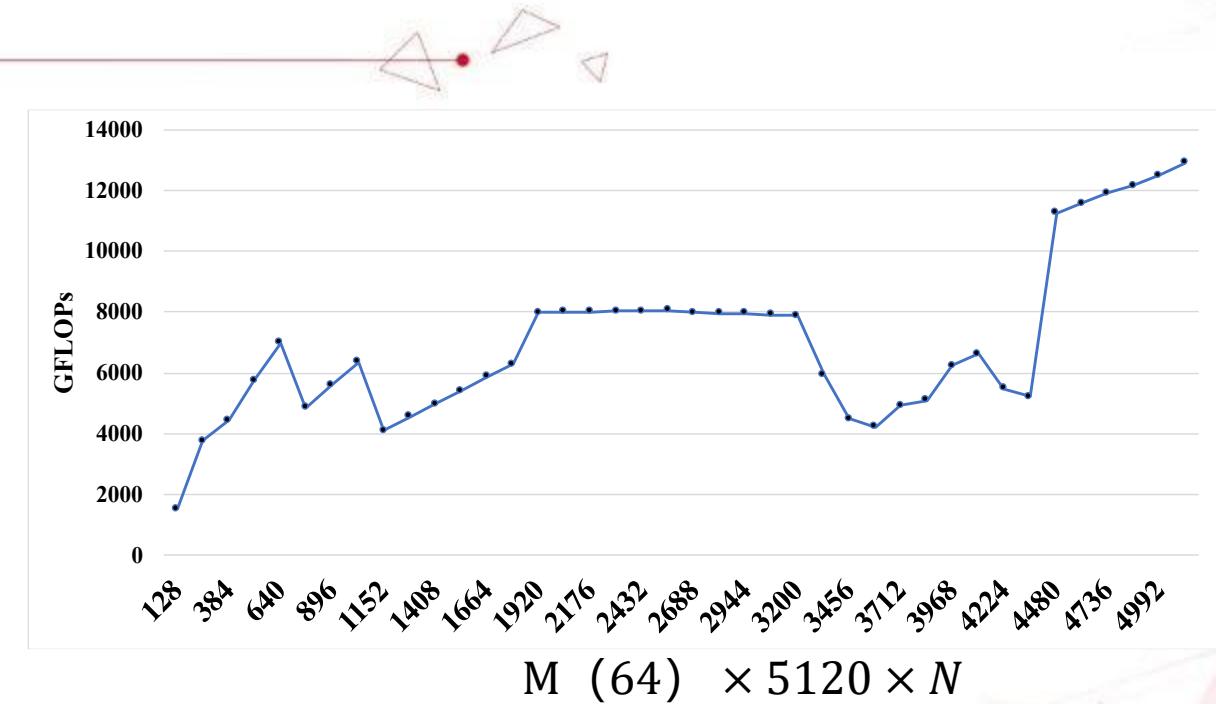
(a) Matrix Size of some convolution



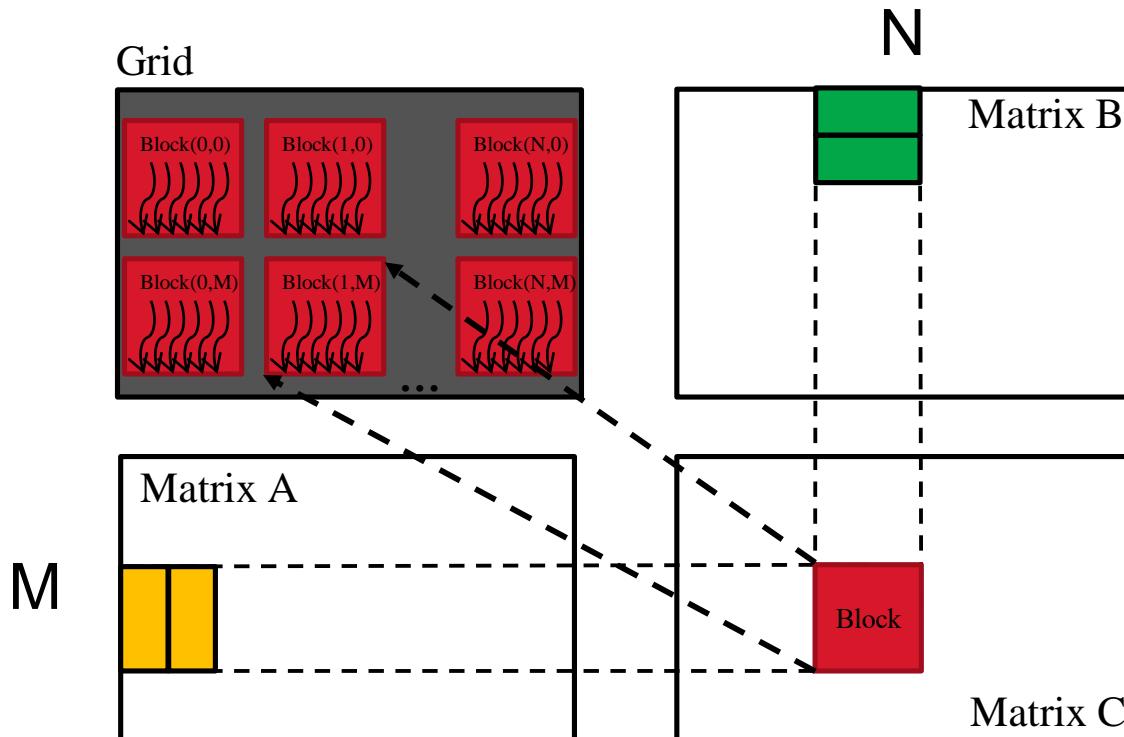
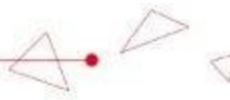
(b) Performance and Normalized Performance

Issues on Small Matrix

- M and N determine tile size and the number of tiles
 - impact the thread-level parallelism (**TLP**)
 - Tiling Strategy
- K determines the workload of each tile
 - Impact on instruction-level parallelism (**ILP**) within a single thread
 - Software Pipelining



Tiling Strategy



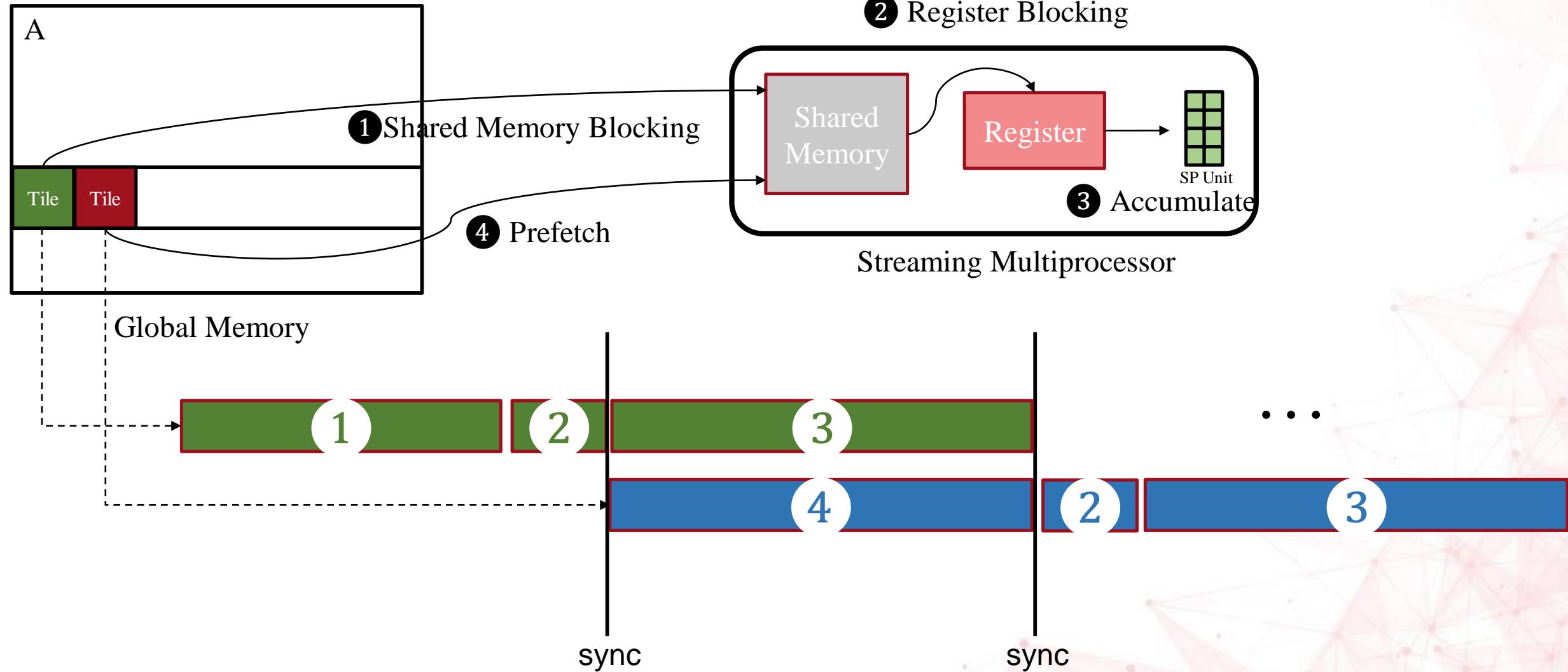
K
matrix multiplication

Tiling Strategy	BY	BX	Thread
Small	16	16	32
Medium	32	32	64
Large	64	64	64
Tall	128	64	128
Wide	64	128	128
Huge	128	128	256

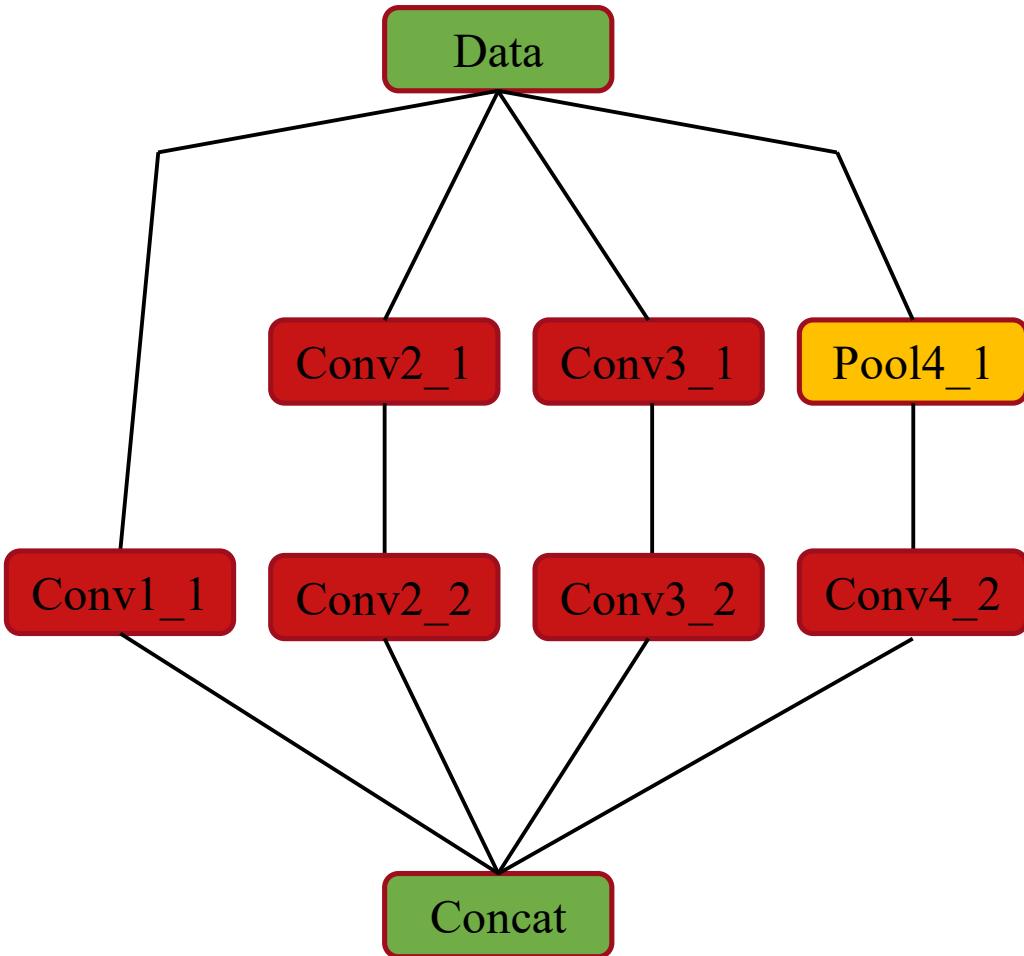
$$\text{计算访存比} = \frac{M \times N \times K}{(M+N) \times K + M \times N}$$

Prefers large tile size to hide the memory latency.

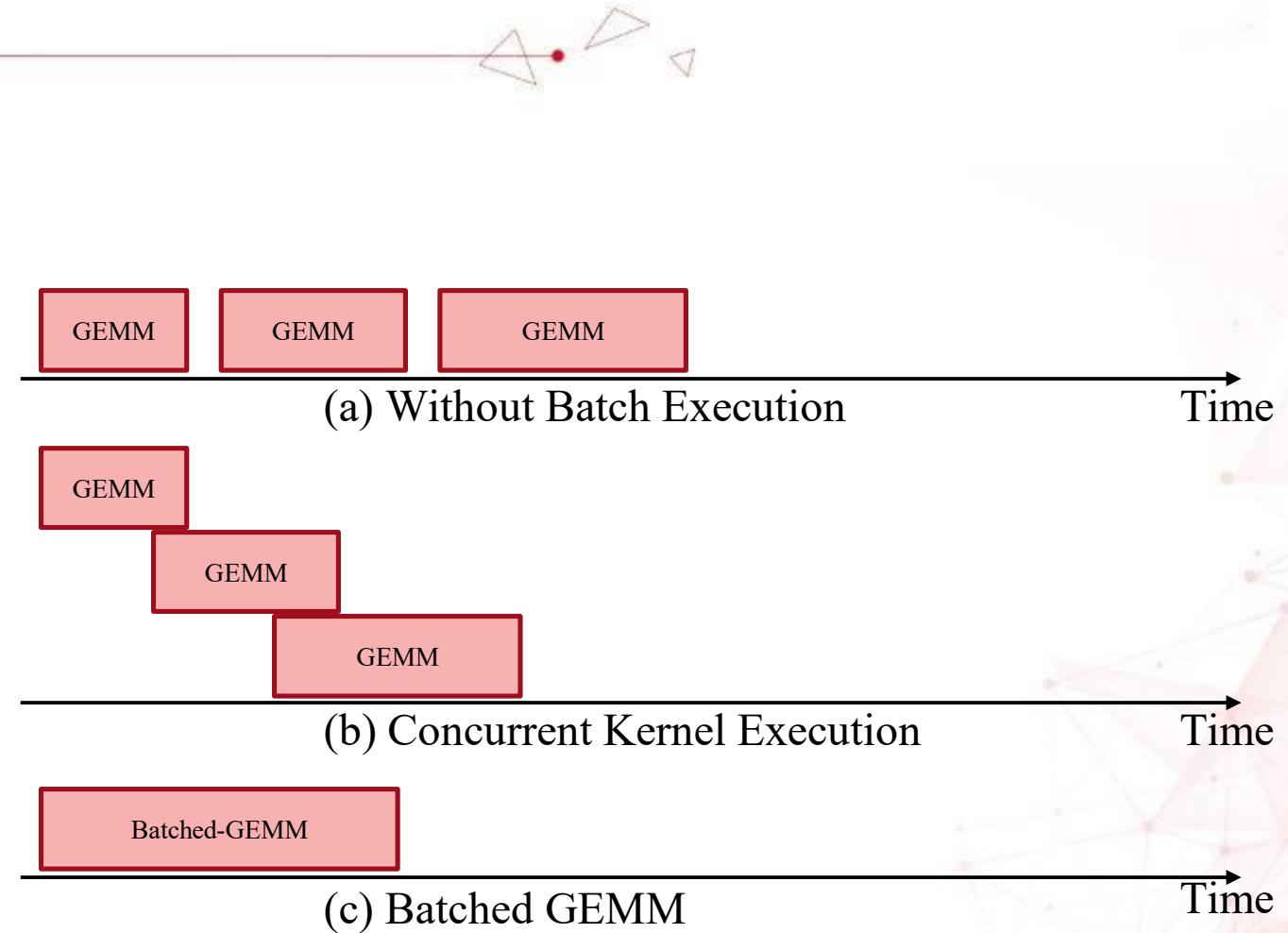
Software Pipelining



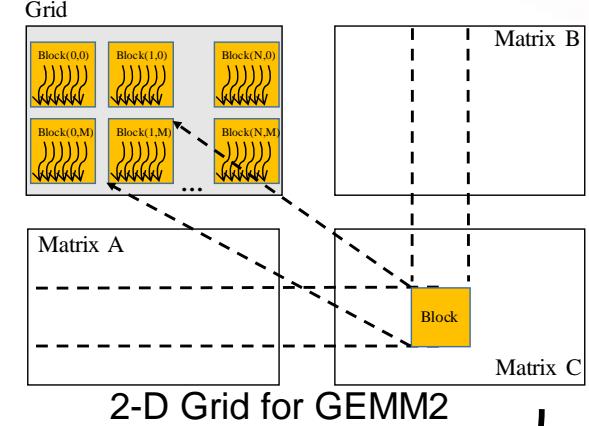
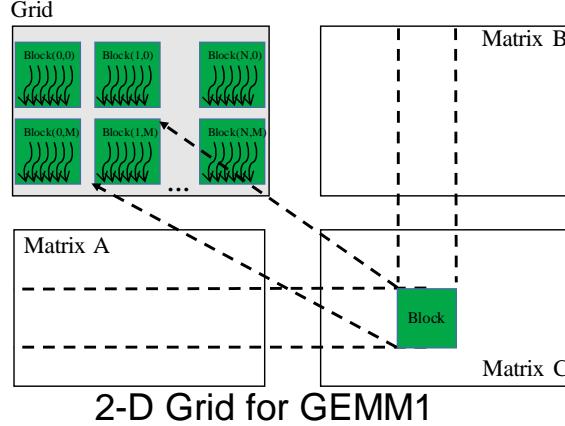
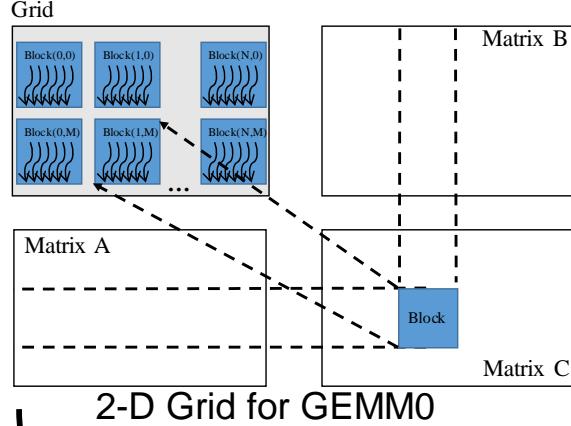
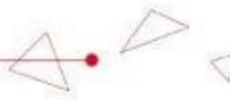
Batched GEMM



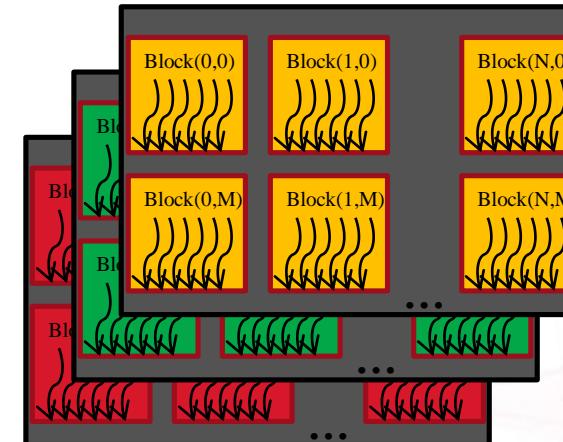
A fan-structure in Google-Net



Batched GEMM

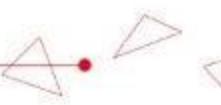


Merge three 2D CUDA kernels into a 3D CUDA kernel

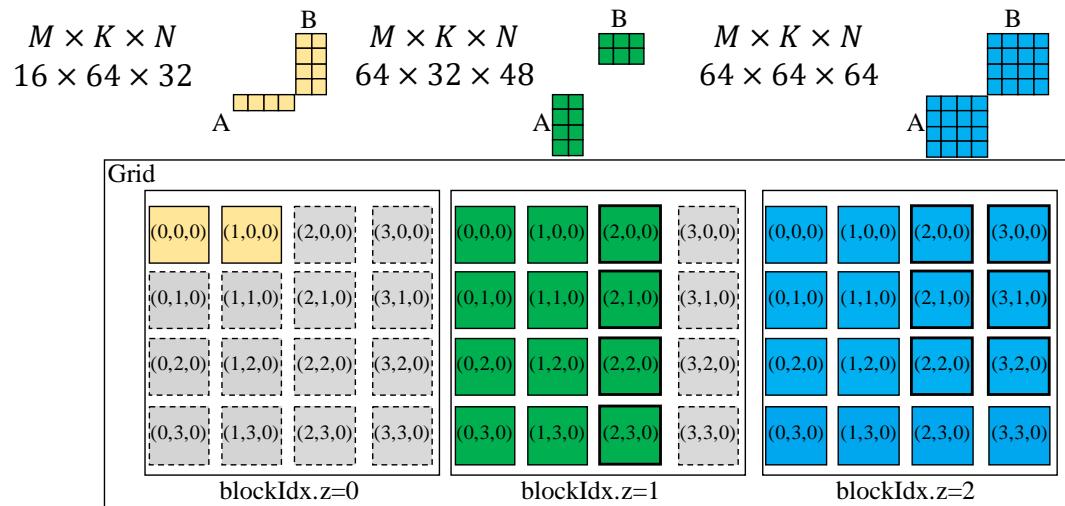


3-D Grid for Batched GEMM

Challenges of Batched GEMM

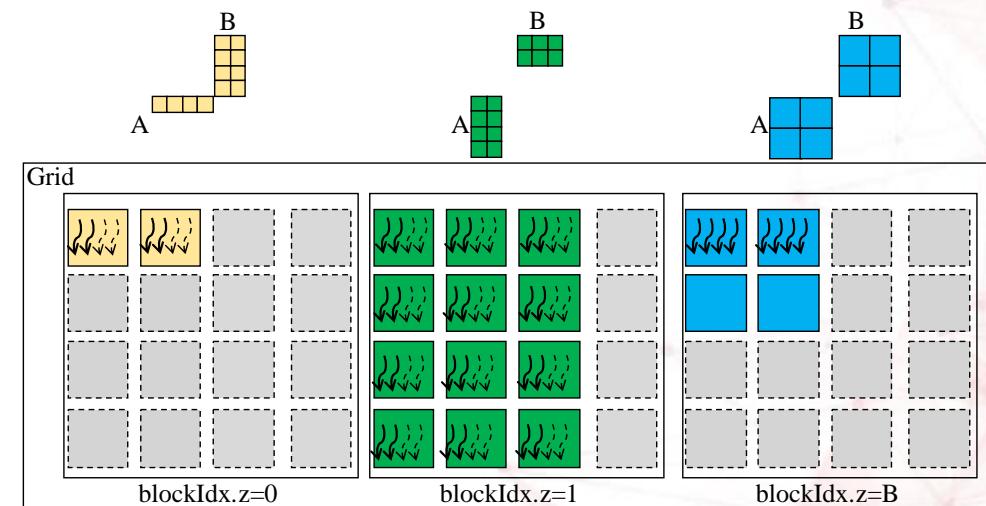


- How to choose proper tile size for each GEMM?
- Thread Under-utilized?



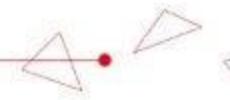
(a) Each GEMM uses the same Tile size 16×16 .

Tiling Strategy	BY	BX	Thread
Small	16	16	32
Medium	32	32	64
Large	64	64	64
Tall	128	64	128
Wide	64	128	128
Huge	128	128	256



(b) The third tile size 32×32 .

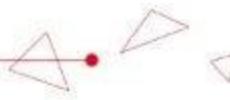
Tiling Engine



- Design Tiling Strategy dedicated for Batched GEMM scenario

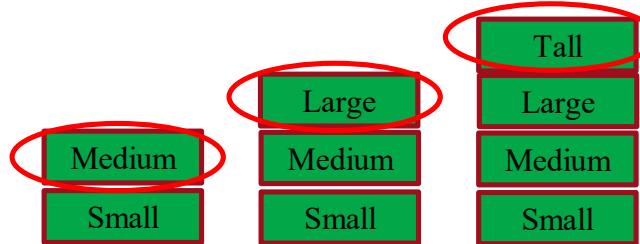
Tiling Strategy	BY	BX	256-Thread (Sub-Tile per Thread)	128-Thread (Sub-Tile per Thread)
Small	16	16	1x1	2x1
Medium	32	32	2x2	4x2
Large	64	64	4x4	8x4
Tall	128	64	8x4	8x8
Wide	64	128	4x8	8x8
Huge	128	128	8x8	16x8

Tiling Engine



$M \times N \times K$ GEMM 1 GEMM 2 GEMM 3
 $32 \times 64 \times 16$ $64 \times 64 \times 32$ $128 \times 64 \times 32$

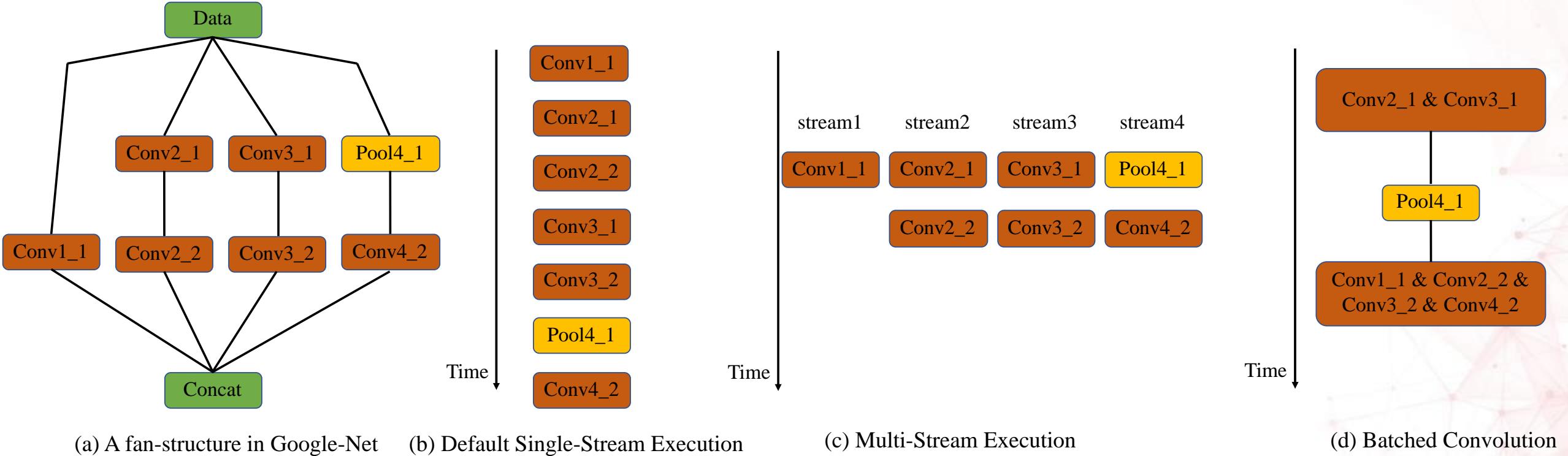
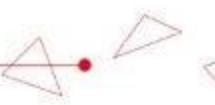
1. Obtain applicable tiling strategies for each GEMM
2. First try the small tiling strategy to guarantee TLP
3. Try the larger tiling strategy
4. The bottom of each bucket is the selected tiling strategies



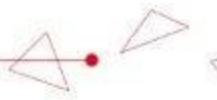
Tiling Strategy	BY	BX
Small	16	16
Medium	32	32
Large	64	64
Tall	128	64
Wide	64	128
Huge	128	128

TLP
TLP > Threshold? Larger: END

Evaluation – Google-Net



Evaluation – Real-world Case Study



- Platform
 - NVIDIA V100 GPU with CUDA 9.0
 - Batch size =1

Implementation	Time (ms)	Speedup
cuDNN without stream	3.18	1X
cuDNN with stream	2.41	1.31X
Batched GEMM	2.01	1.58X

目录

1 //
2 //
3 //
4 //
5 //
6 //

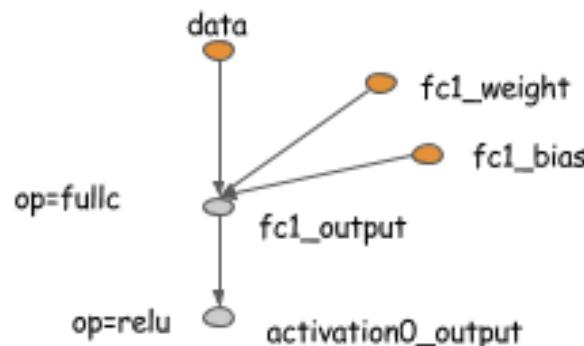
- AI复兴：深度学习
- 计算优化：GPU
- 计算以外：内存&访存
- 携手共进：通信优化
- 无所不能：云上AI
- 轻松一刻：应用案例

计算以外：内存优化

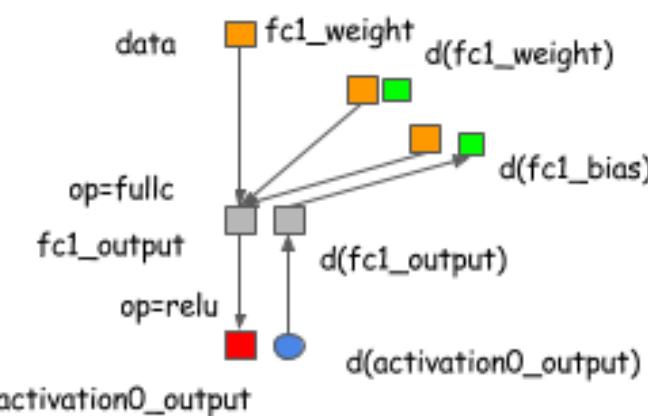
■ 深度神经网络的计算过程可以用一个图来进行描述

- 分为前向计算和反向传播两个过程

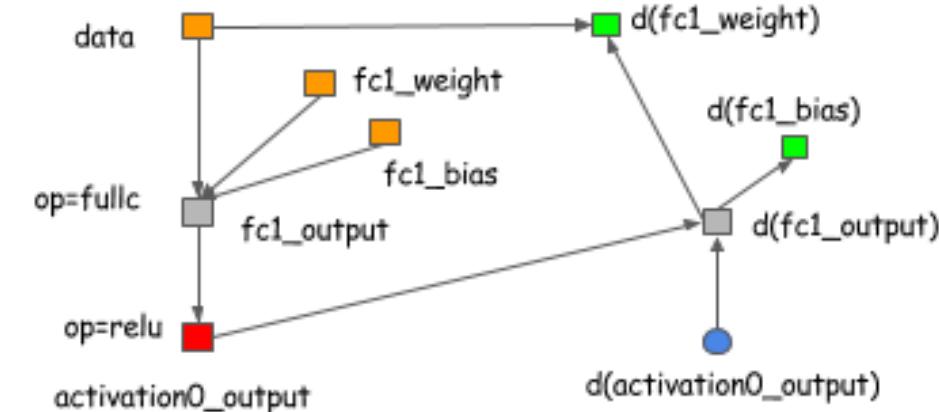
Forward Graph (Configuration)



Backprop on Same Graph



Explicit Backward Graph



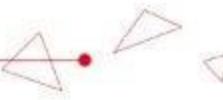
- argument symbol nodes
- internal symbol nodes

- argument arrays
- output arrays

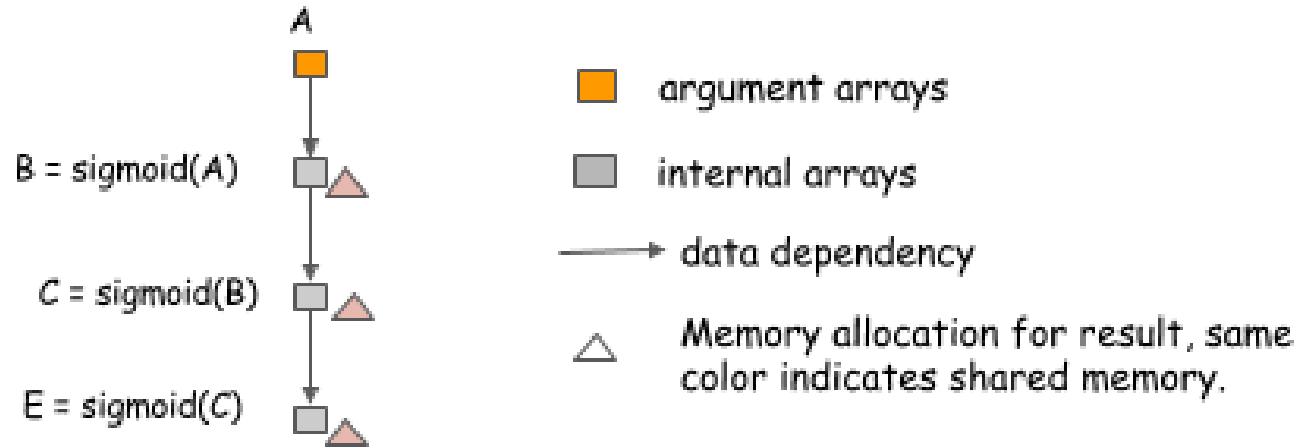
- argument gradient holders
- internal arrays

- backward input nodes
- d(?) gradient node of ?

计算以外：内存优化

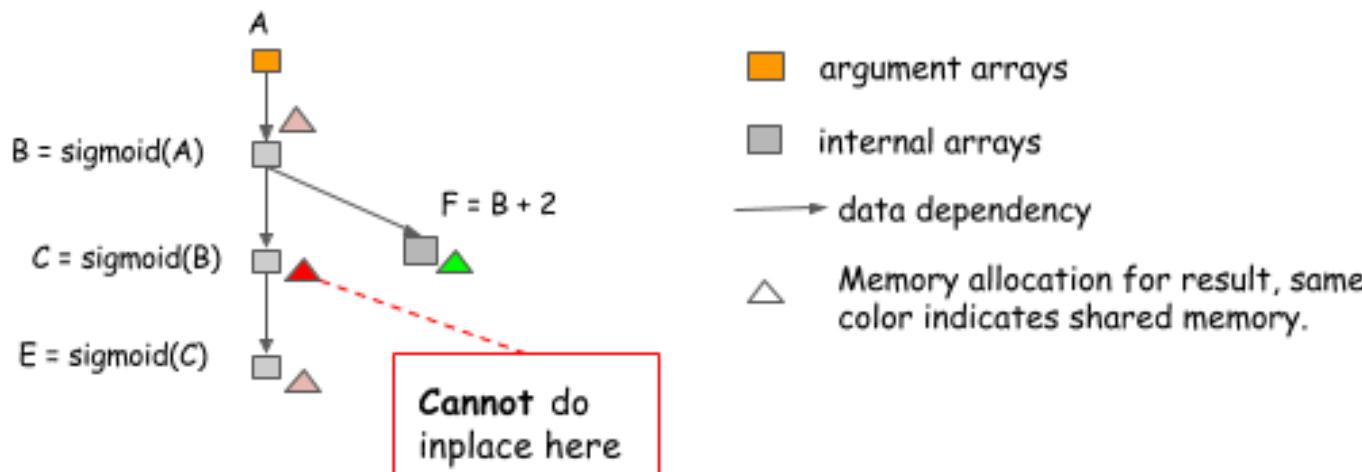


■ 优化方法1：In-place Operation



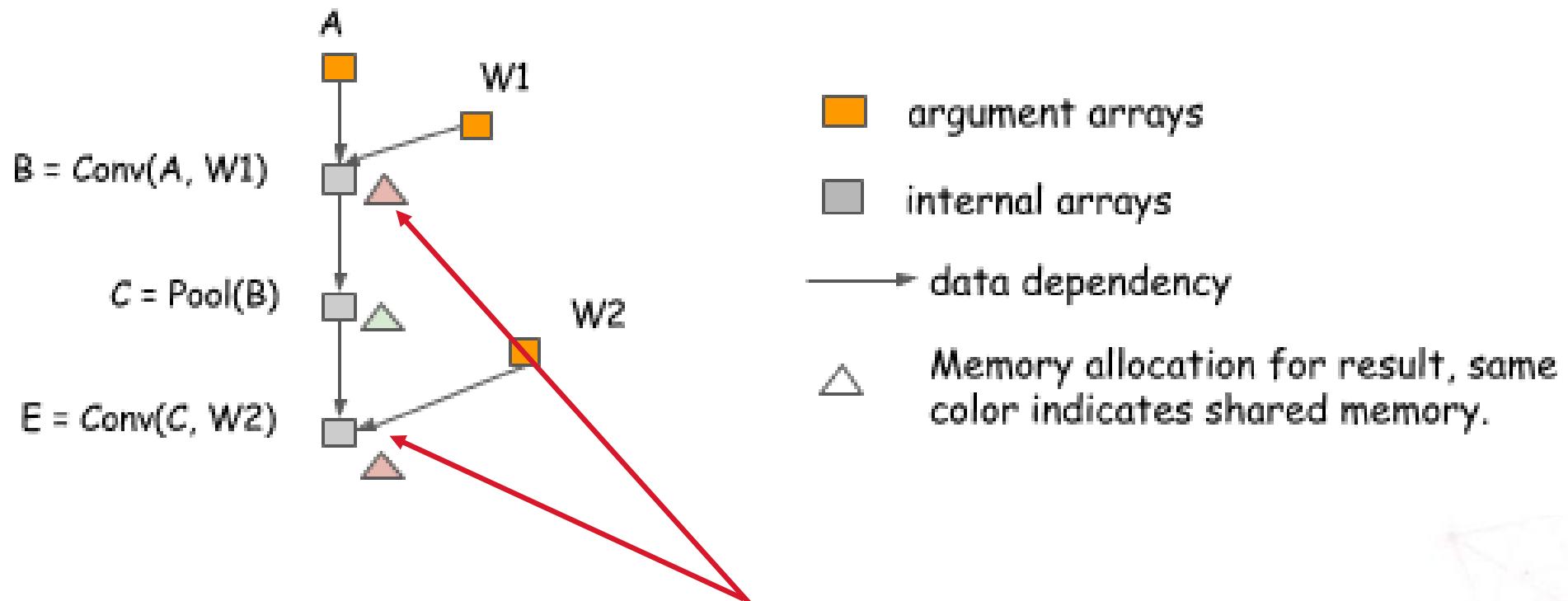
● 激活函数层往往可以使用In-Place操作

● 带有分支的激活函数不能使
用In-Place操作



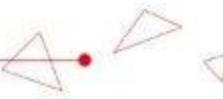
计算以外：内存优化

■ 优化方法2：Memory Sharing

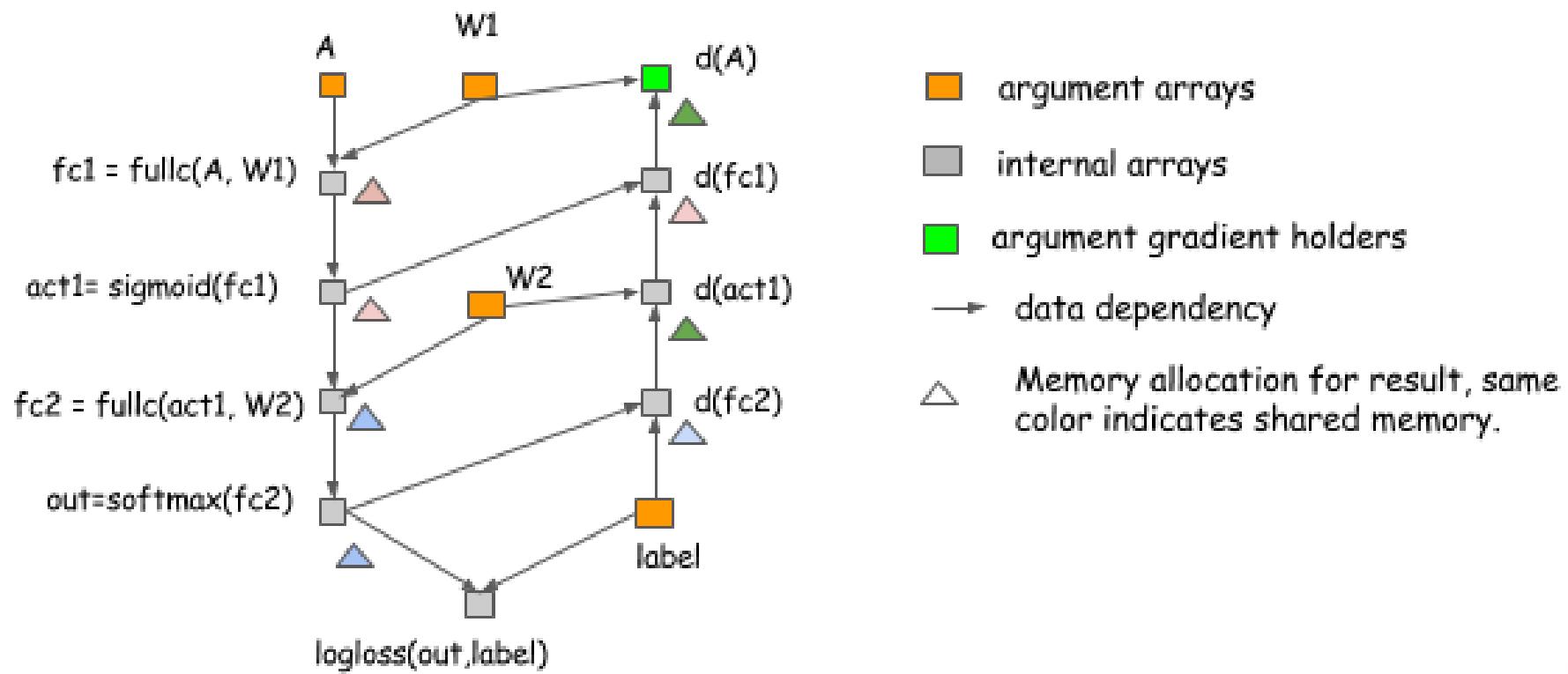


在计算E的时候，B的内存可以释放

计算以外：内存优化



■ 实际案例：In-Place Operation & Memory Sharing



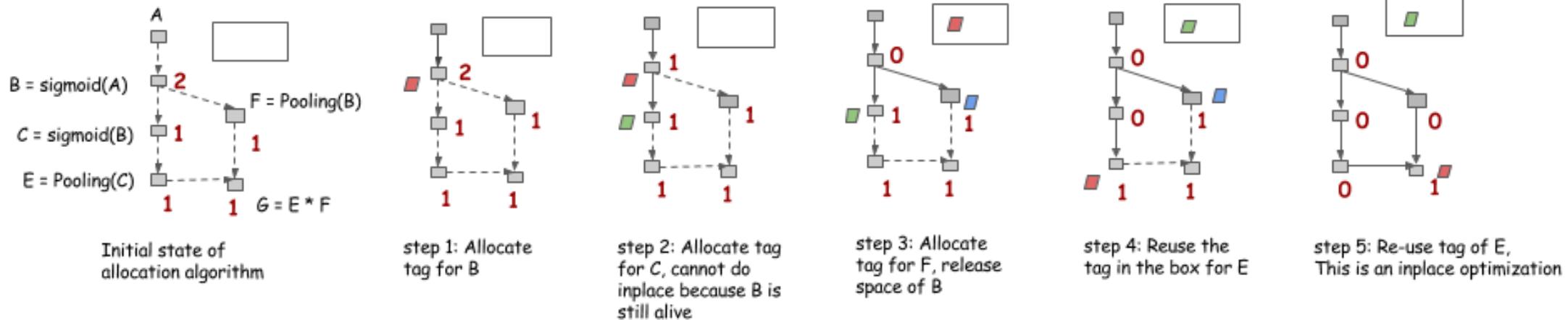
相同颜色的三角形复用内存

计算以外：内存优化

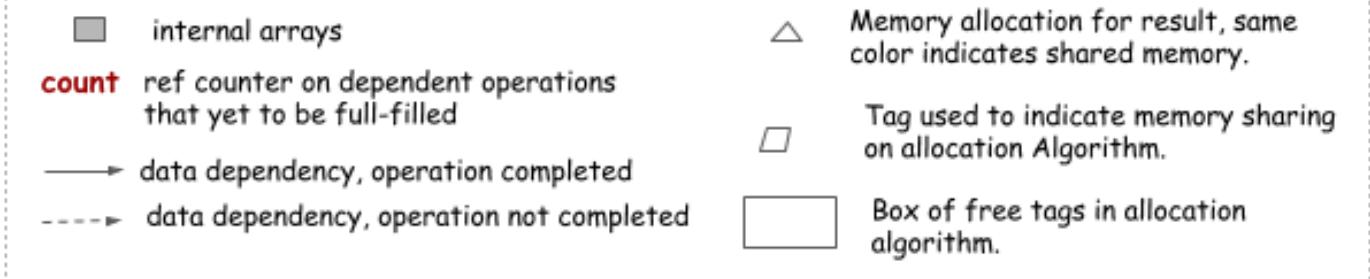
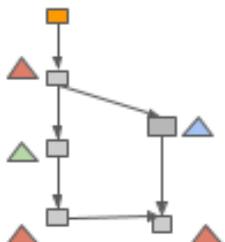
■ 如何实现？

- 内存申请算法，类似编译器中寄存器分配

Allocation Steps



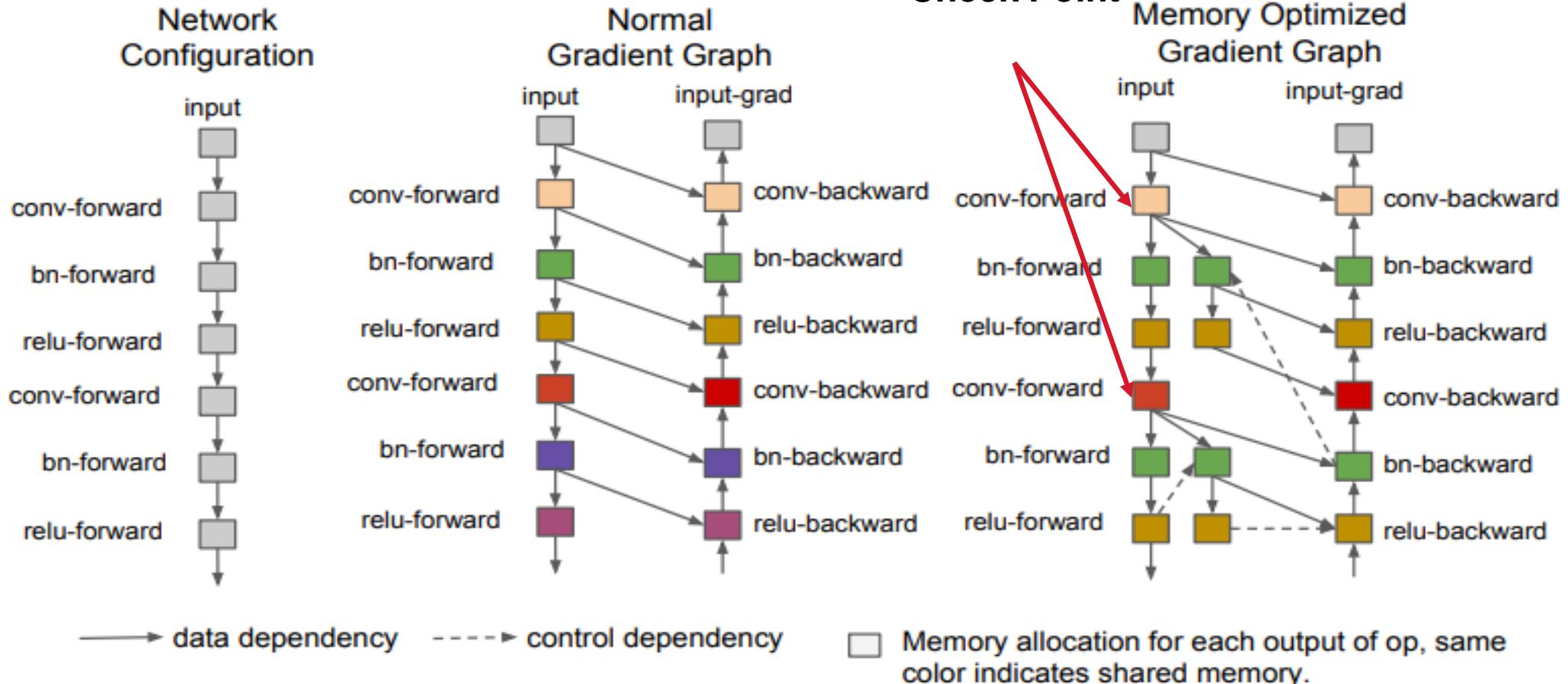
Final Memory Plan



采用简单计数器，标记每一块内存之后需要被用到的次数

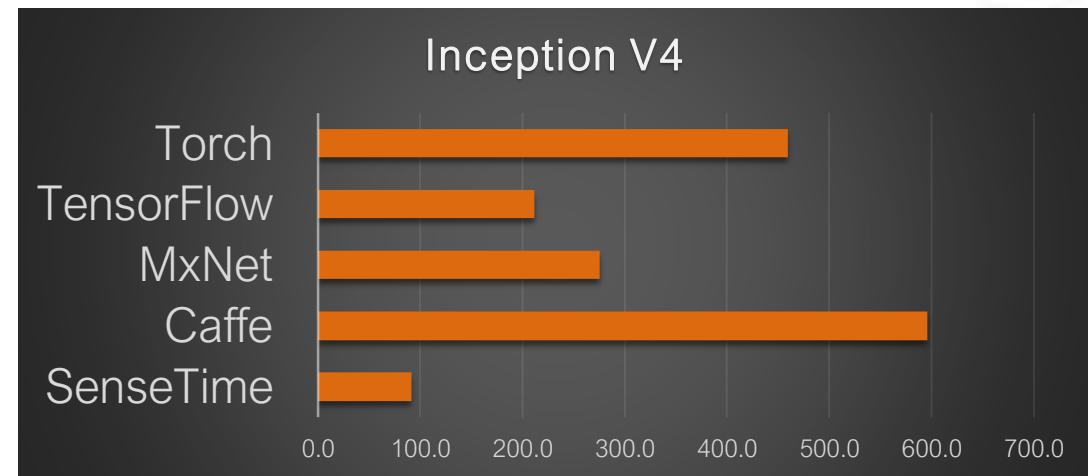
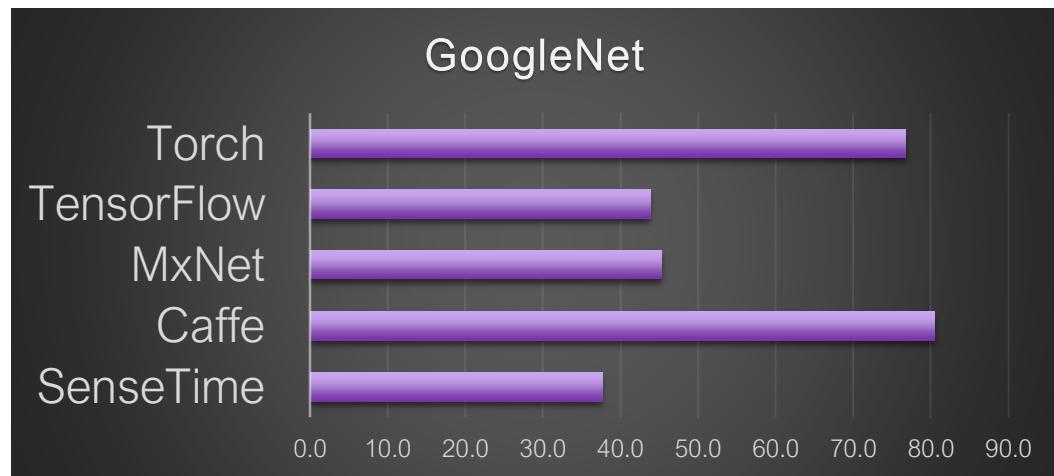
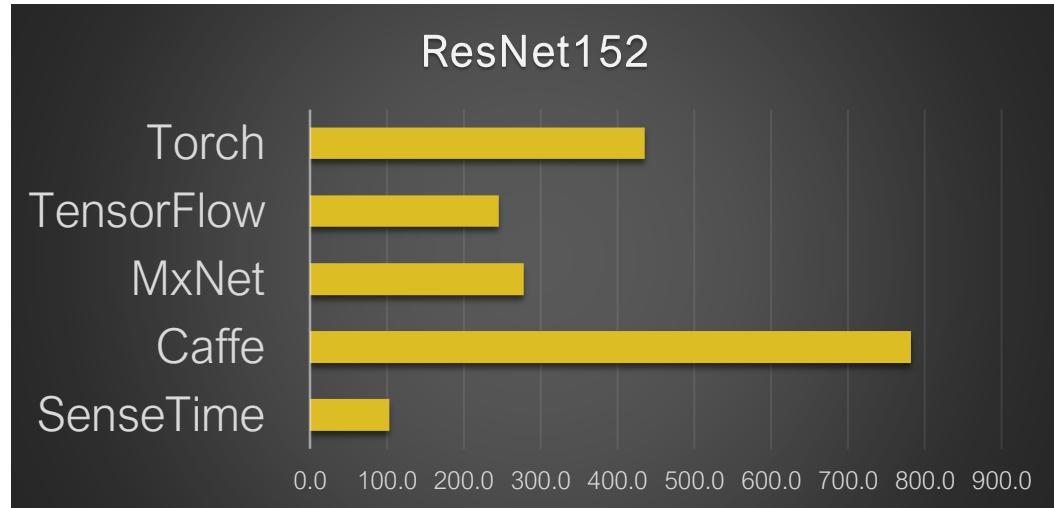
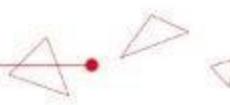
计算以外：内存优化

■ 优化方法3：Computation Vs. Memory



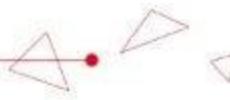
通过重复计算一次Forward，可以获得次线性的内存消耗

计算以外：内存优化

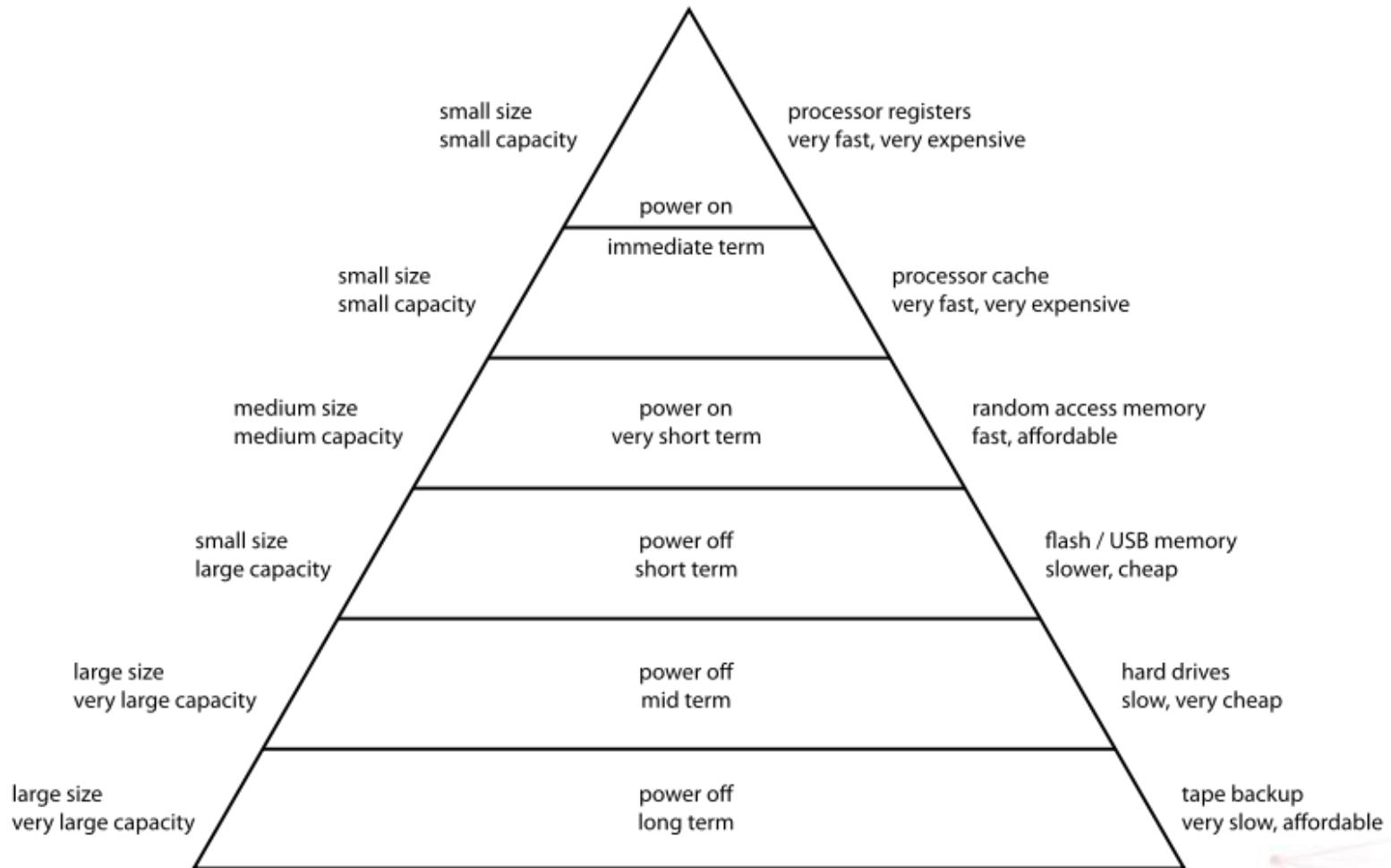


Memory efficiency, lower is better (MxNet不确认是否是最优版本)

计算以外：访存优化

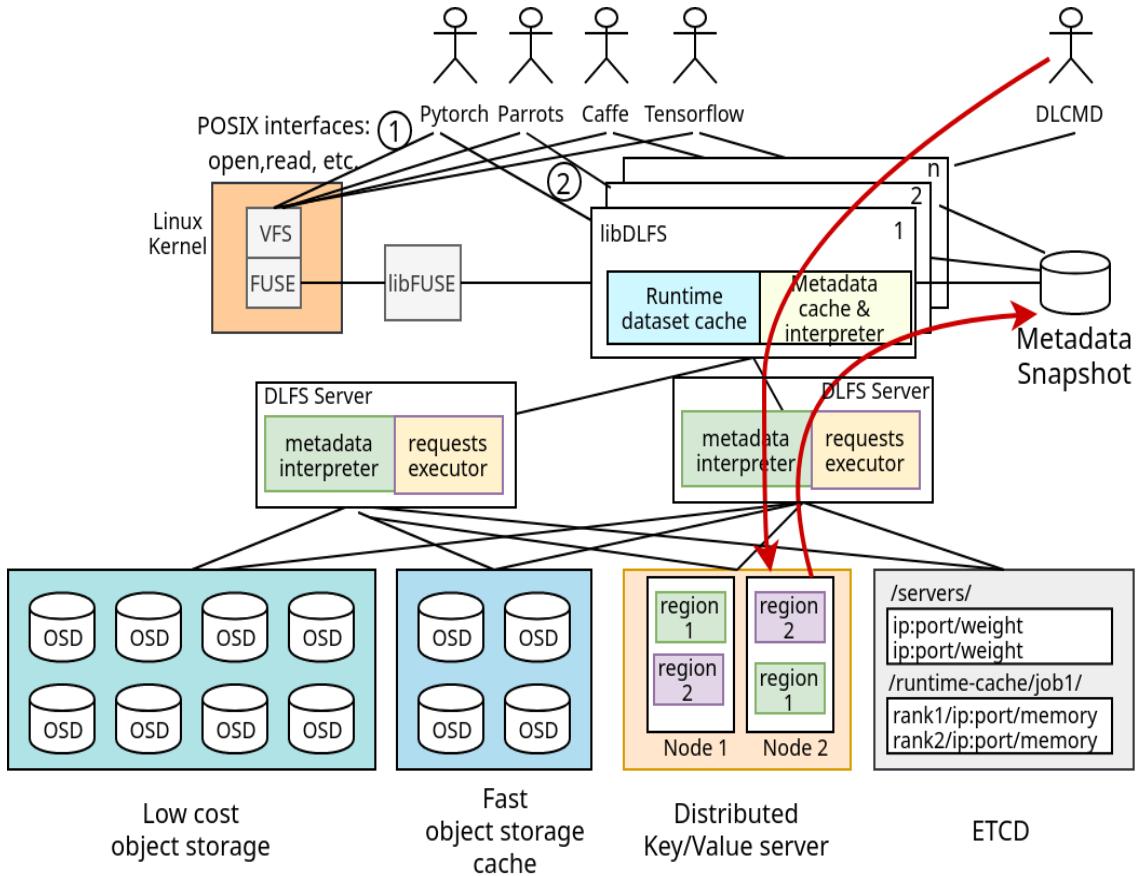


Computer Memory Hierarchy



计算以外：访存优化

System Overview

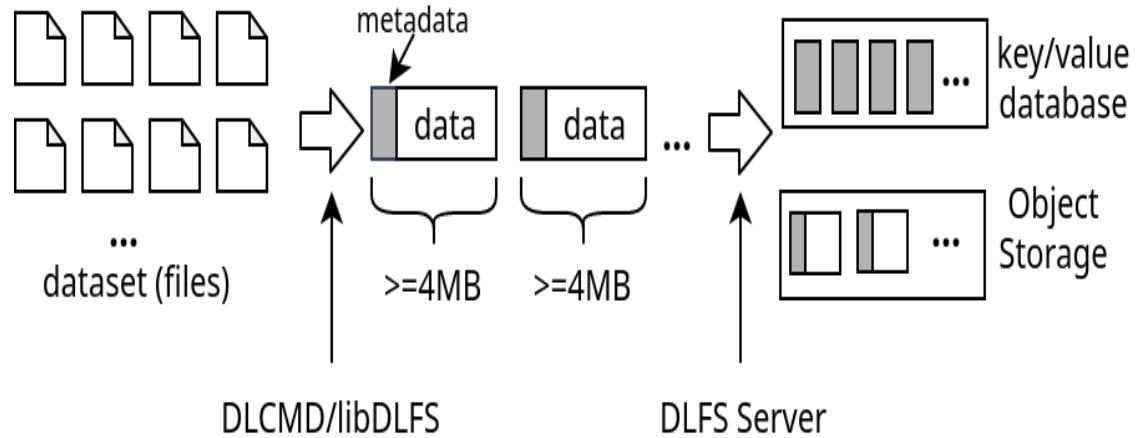


Overview of DLFS System

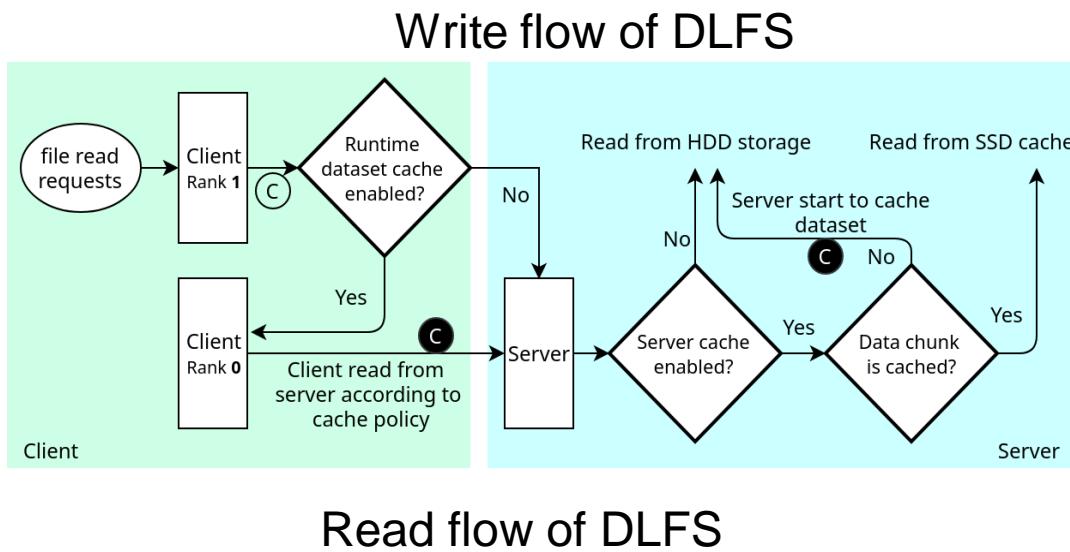
- Use multiple data storage backends to reduce hardware cost
- Server is stateless and has no dependencies between each other, free to scale up
- Metadata is stored in distributed key/value server. Support two types of key/value server (NVME-based Tikv and memory-based Redis).
- Users have two ways to access the data in our system: proprietary APIs and POSIX (via FUSE - Filesystem in Userspace)
- Client has runtime dataset caching system to cache job dataset among client instances
- Metadata can be materialized (snapshot) locally to bypass the key/value server

高性能存储

Write/Read dataflow from DLFS client to DLFS server

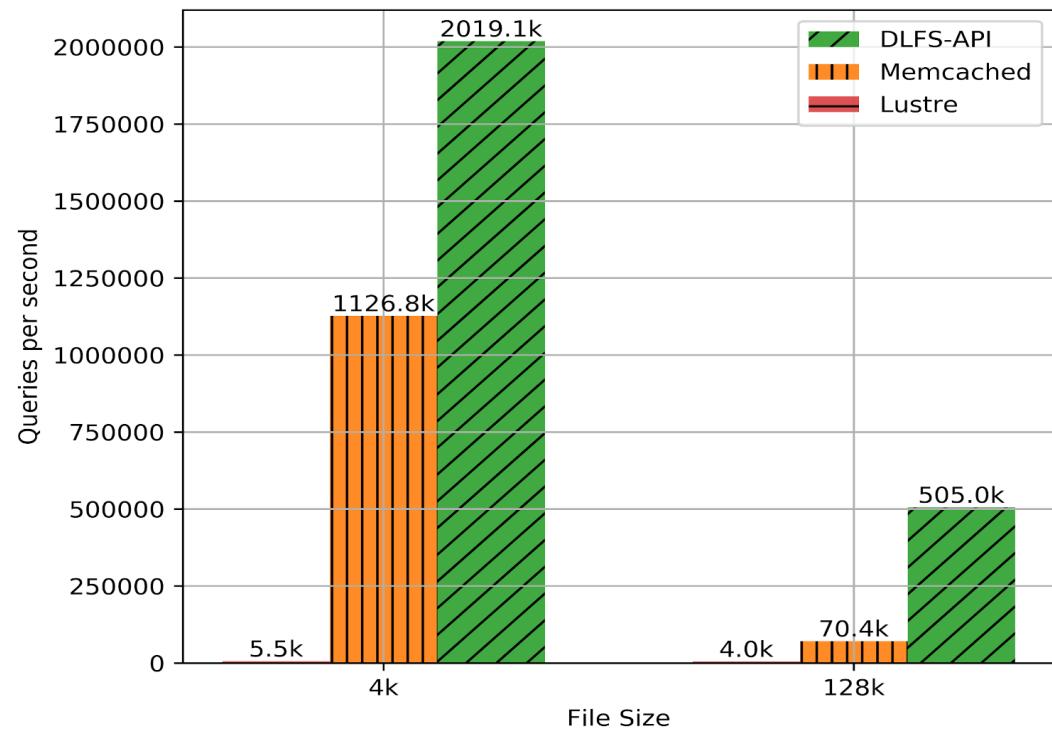


- File writes in batches
- File aggregates to large data chunks with metadata include in the head
- Runtime dataset cache works on metadata snapshots
- There are multiple client instances on a machine. To reduce the number of network connections, we only select one client instance to perform the runtime dataset cache
- The client cache used for block-based shuffle read method and the FUSE interface

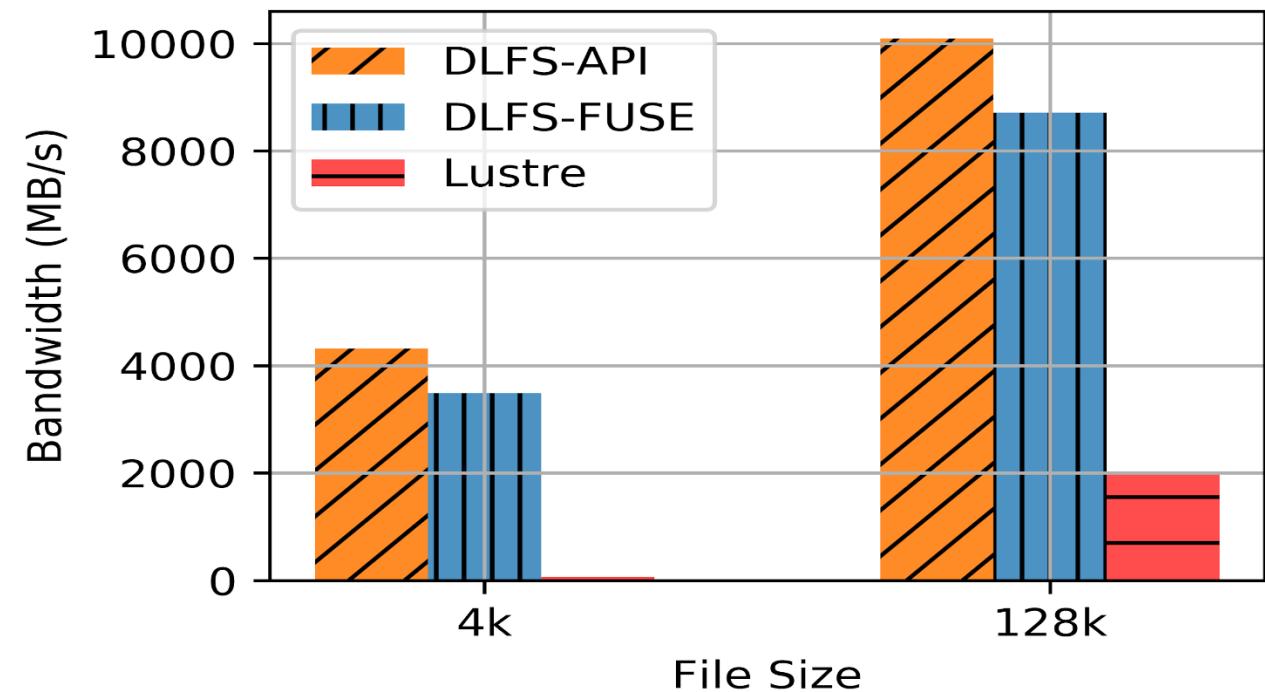


高性能存储|Evaluation

Points: DLFS writes in batches, merges small files into large data chunks
DLFS API Vs. Lustre vs. Memcached



QPS with 64 threads (on 4 nodes)



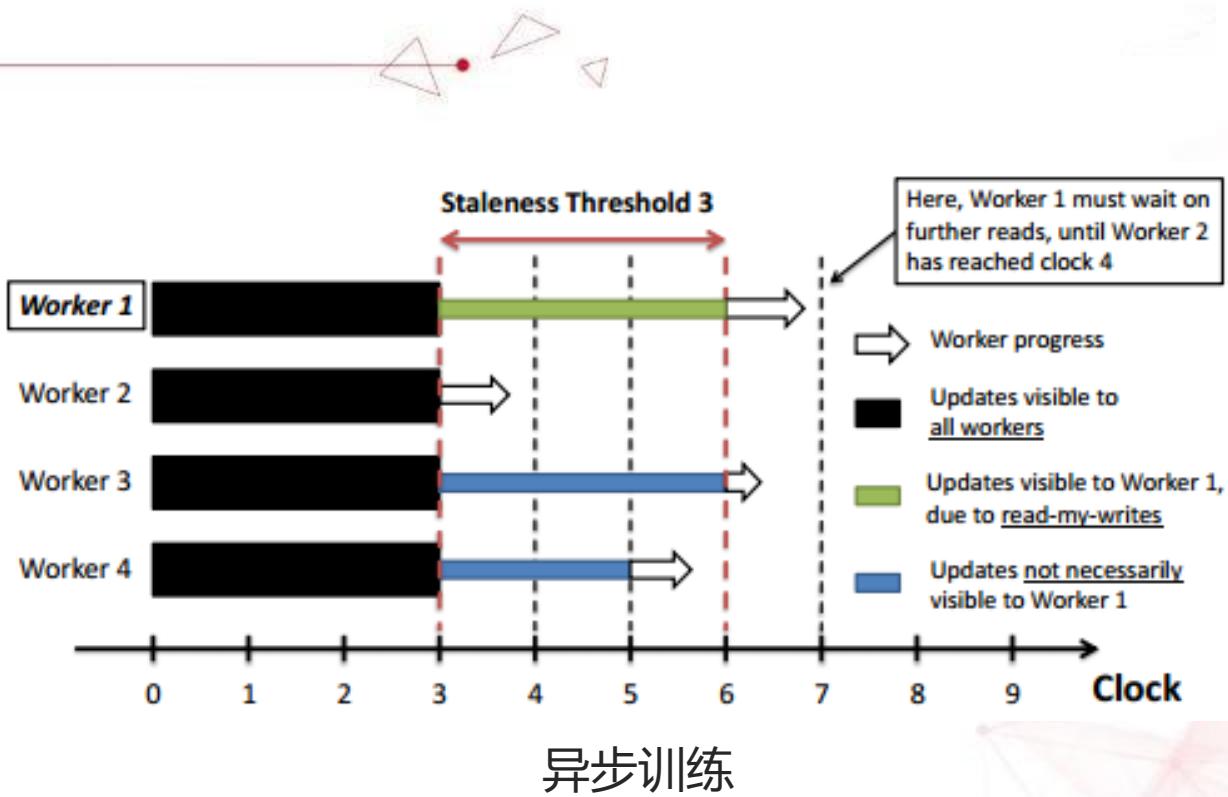
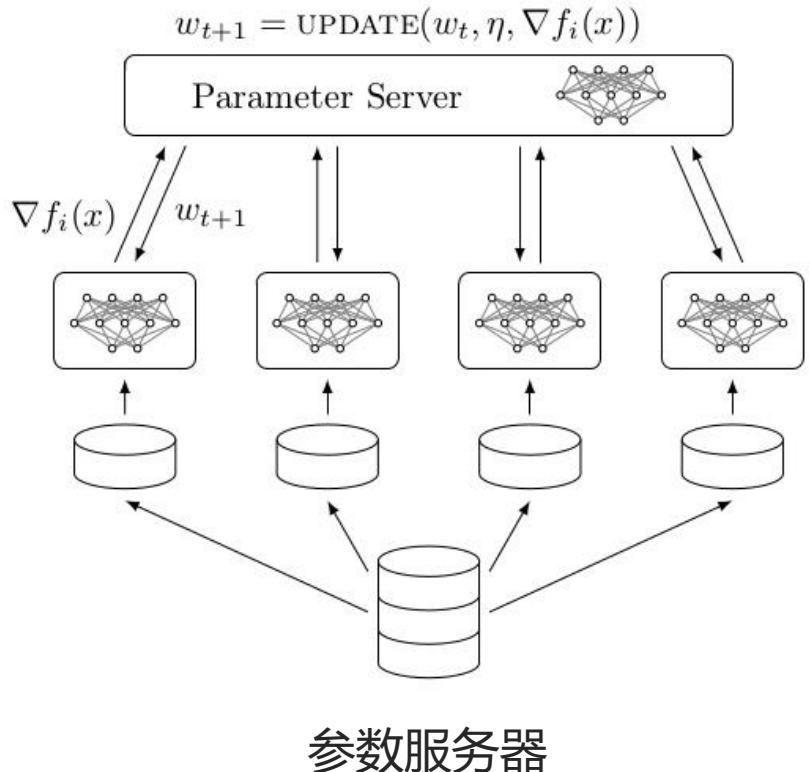
Bandwidth on 4 nodes (higher is better)

目录

1 //
2 //
3 //
4 //
5 //
6 //

- AI复兴：深度学习
- 计算优化：GPU
- 计算以外：内存&访存
- 携手共进：通信优化
- 无所不能：云上AI
- 轻松一刻：应用案例

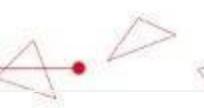
模型同步：参数服务器



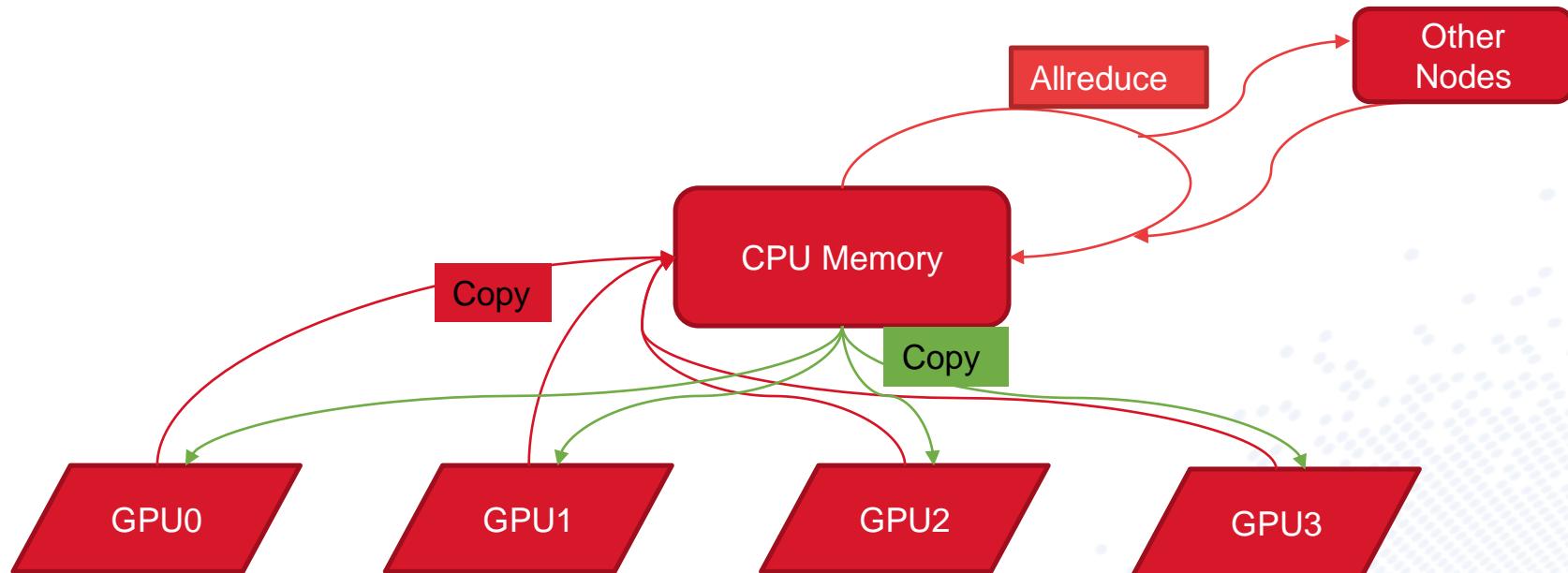
参数服务器的好处：

1. 服务器分为Server和Work，二者之间采用Put和Get操作通信，无需显示同步
2. 模型存放在Server上，不受单机内存限制，方便做模型并行，支持大模型

模型同步: Allreduce



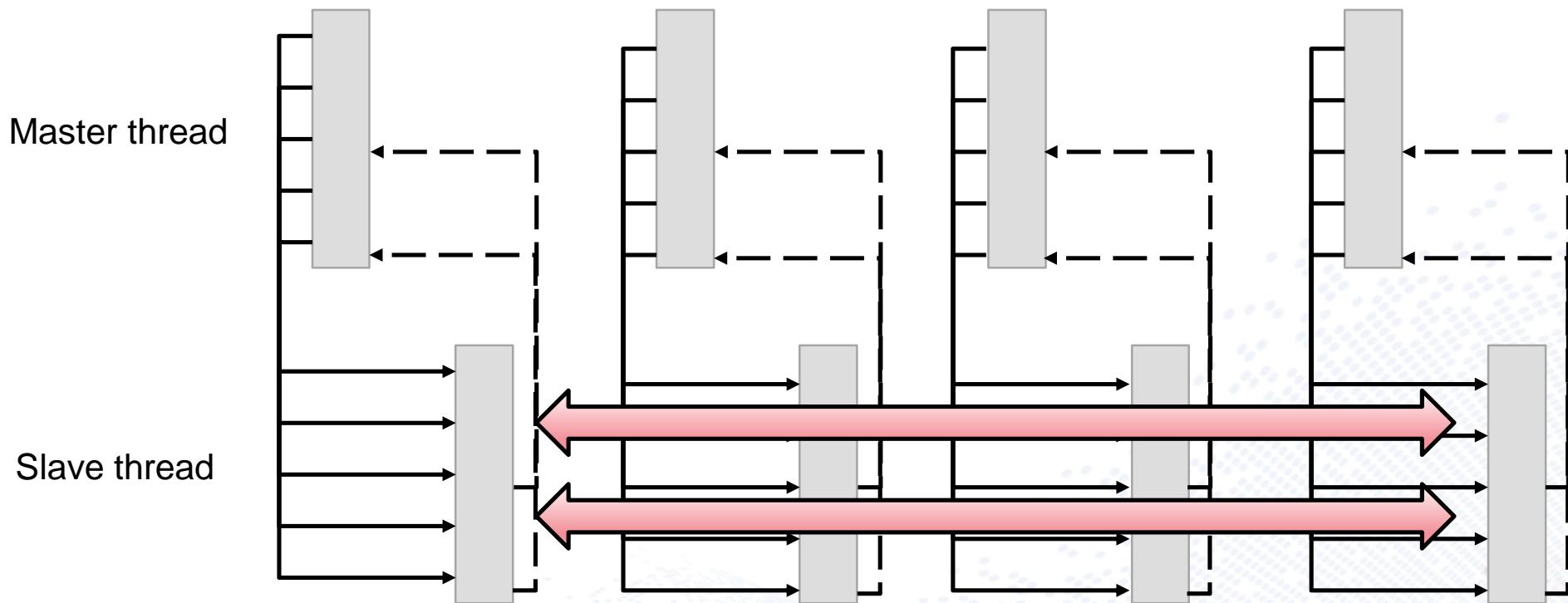
- Support Multi-GPUs and Multi-Nodes
- Three procedures: Copy, Allreduce, Copy



Copy: PCI-E

MPI Allreduce: CPU Allreduce

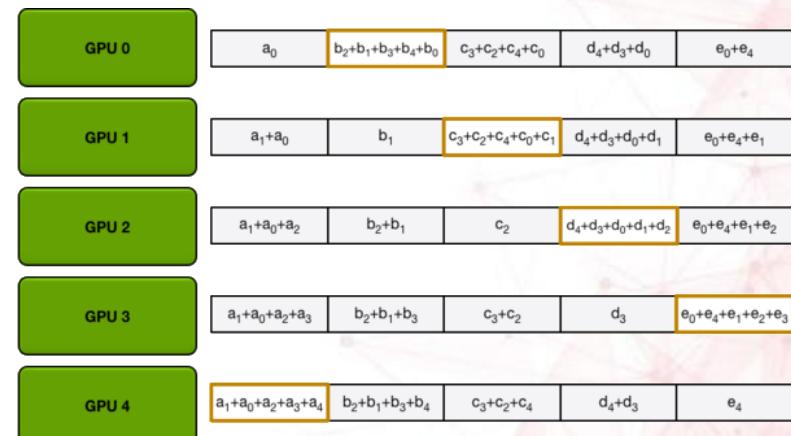
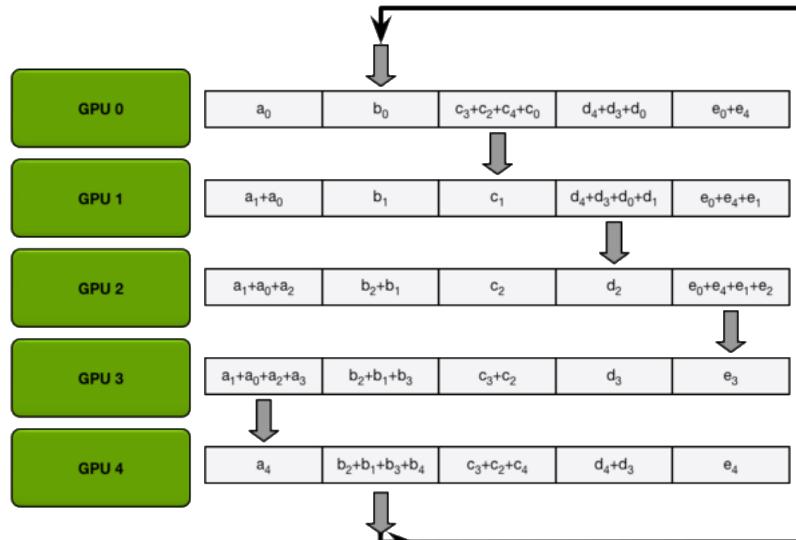
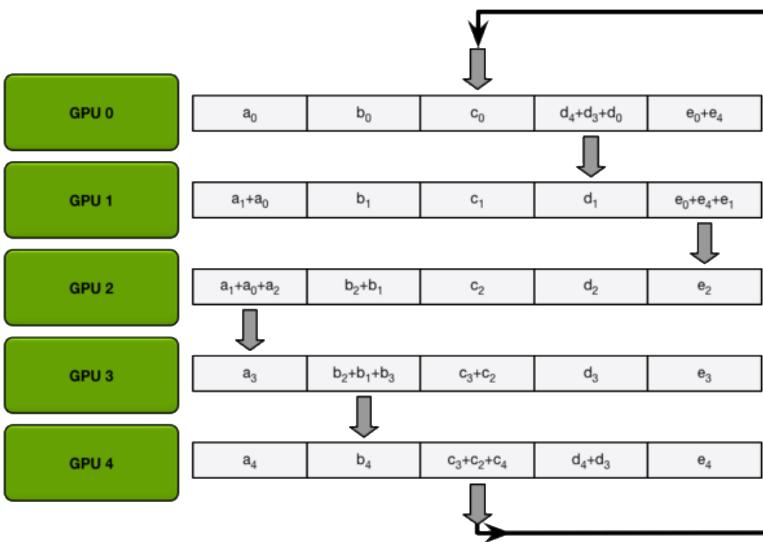
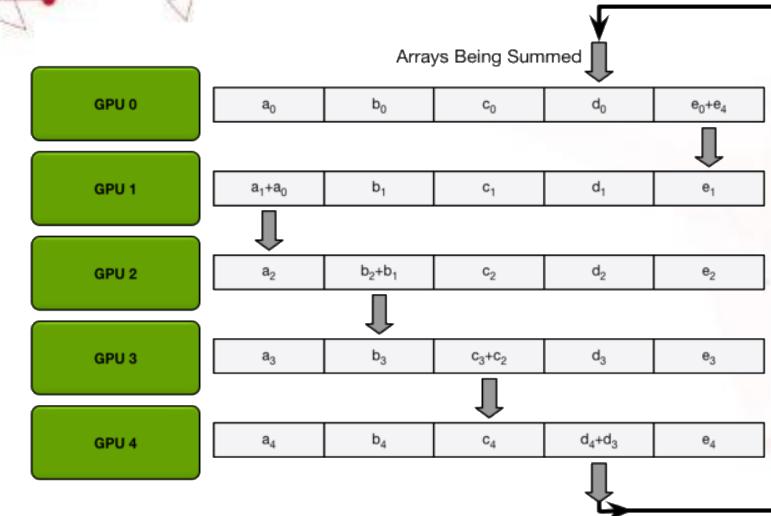
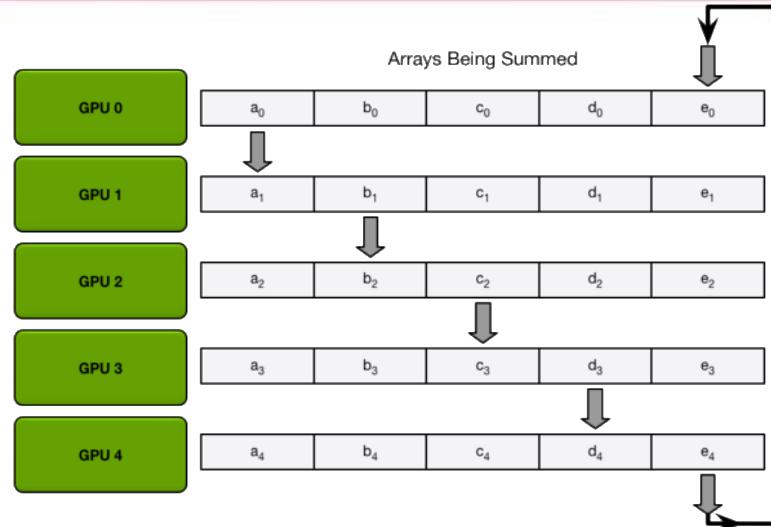
■ Task queue



Ring-based Allreduce

Ring Allreduce流程

1. scatter-reduce
2. allgather

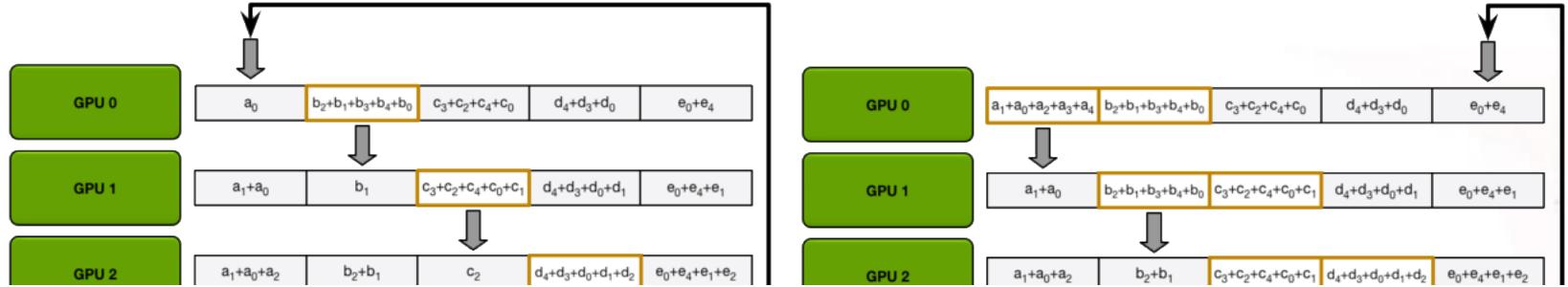


Ring-based Allreduce

Ring Allreduce流程

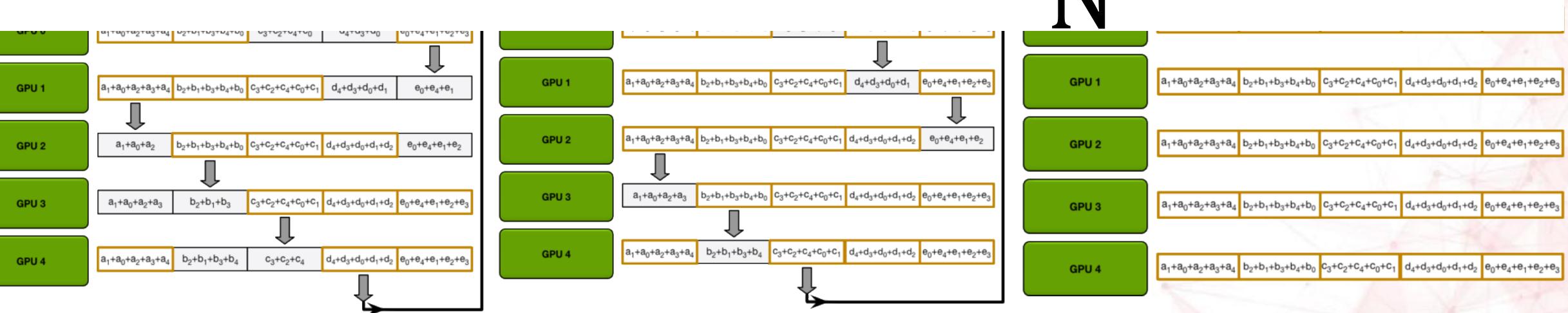
1. scatter-reduce

② 

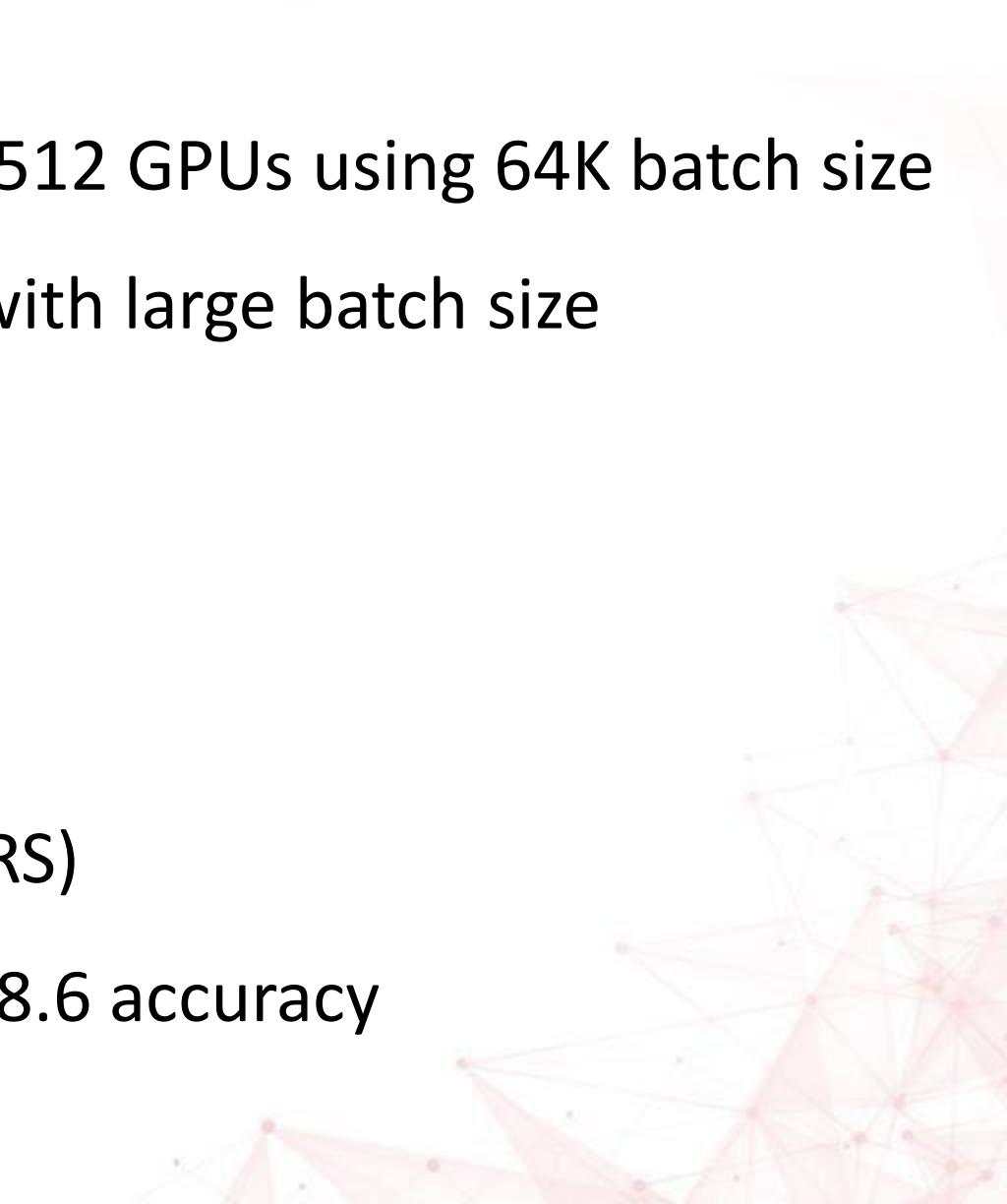


M

$= 2 \frac{M}{N} * (N - 1)$



Scale out distributed DL on a cluster with >512 GPUs



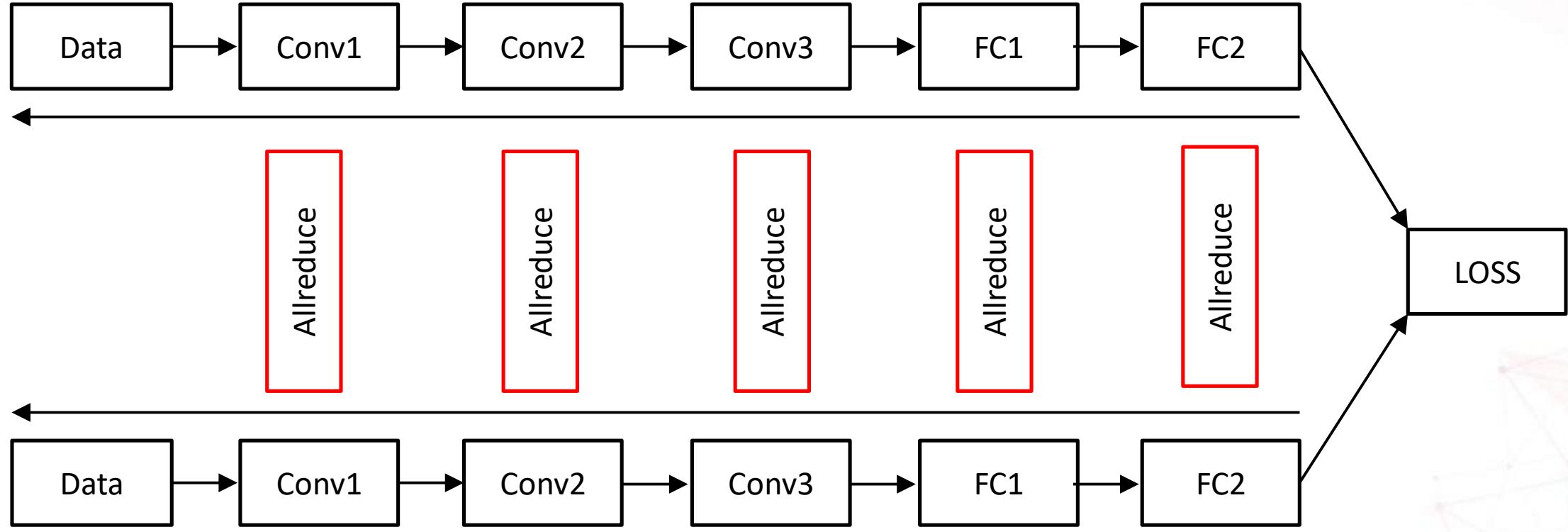
Target: Train AlexNet on ImageNet with 512 GPUs using 64K batch size

Problem: Guarantee the model quality with large batch size

Solutions:

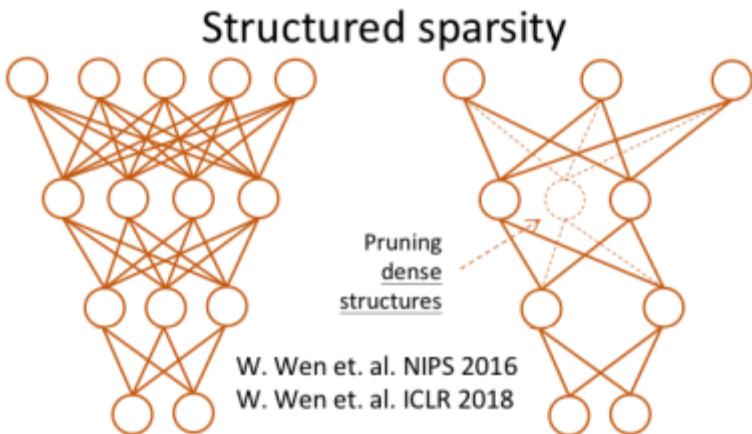
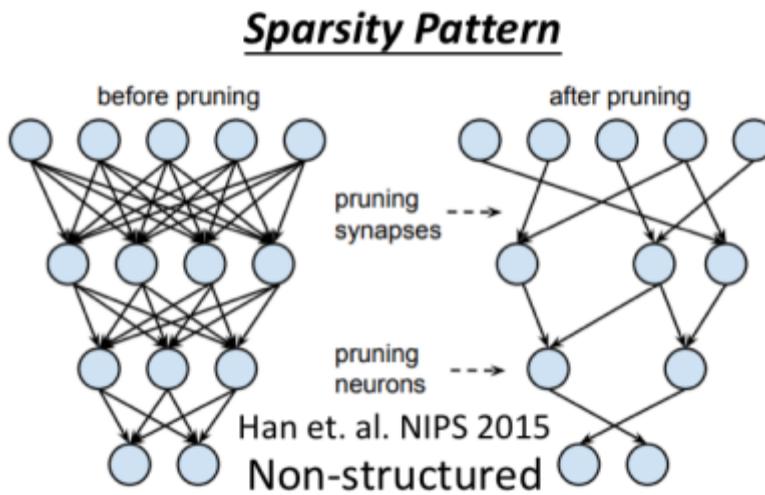
- Linearly scaling up learning rate (LR)
- Warmup LR
- Layer-wise Adaptive Rate Scaling (LARS)
- AlexNet: 8K (2017) -> 64K (2018), > 58.6 accuracy

Improve network performance with block-based communication design

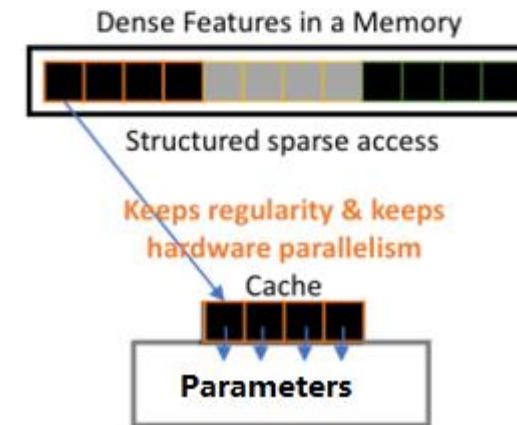
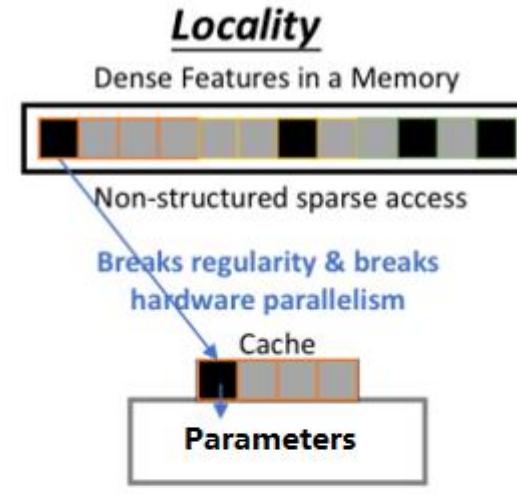


- Layer-based allreduce communication operation

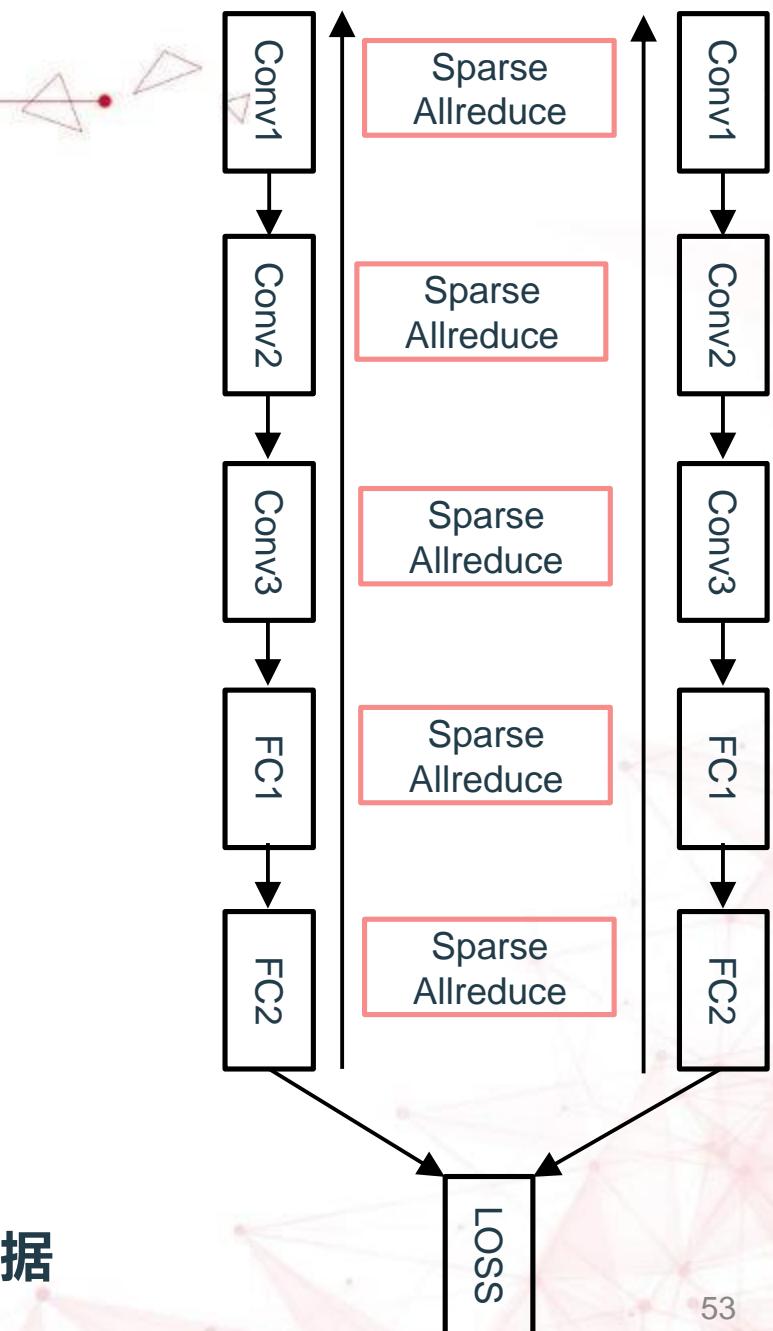
Sparse block-based communication design



结构化Sparse示意图



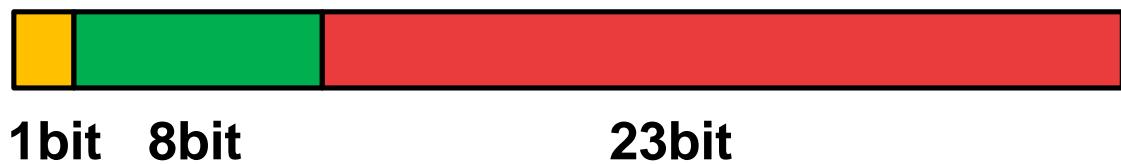
结构化Sparse后通信的数据



数据类型：单精度和半精度



32 bit = float, single precision



可表示的最大值

3.402823×10^{38}

非规格化最小值

2^{-149}

16 bit = half, half precision



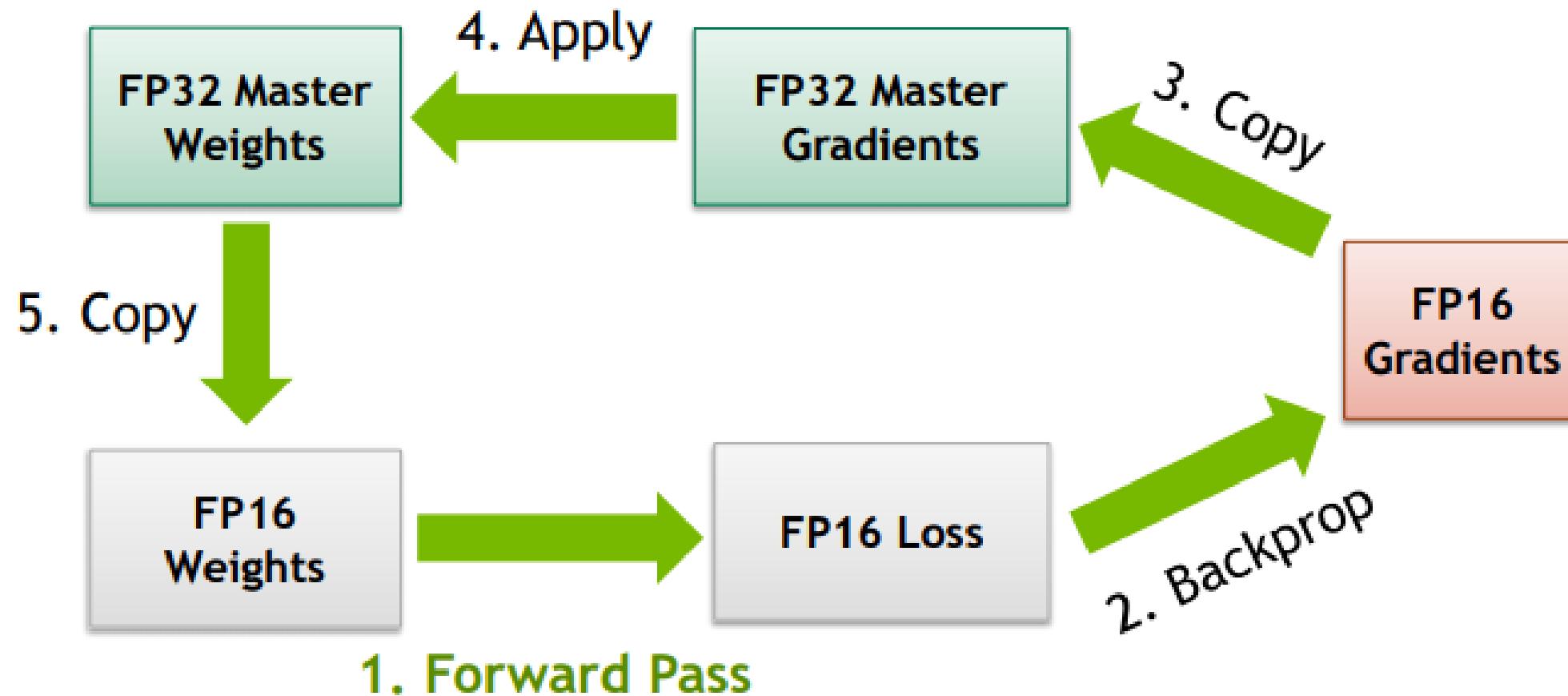
可表示的最大值

65504

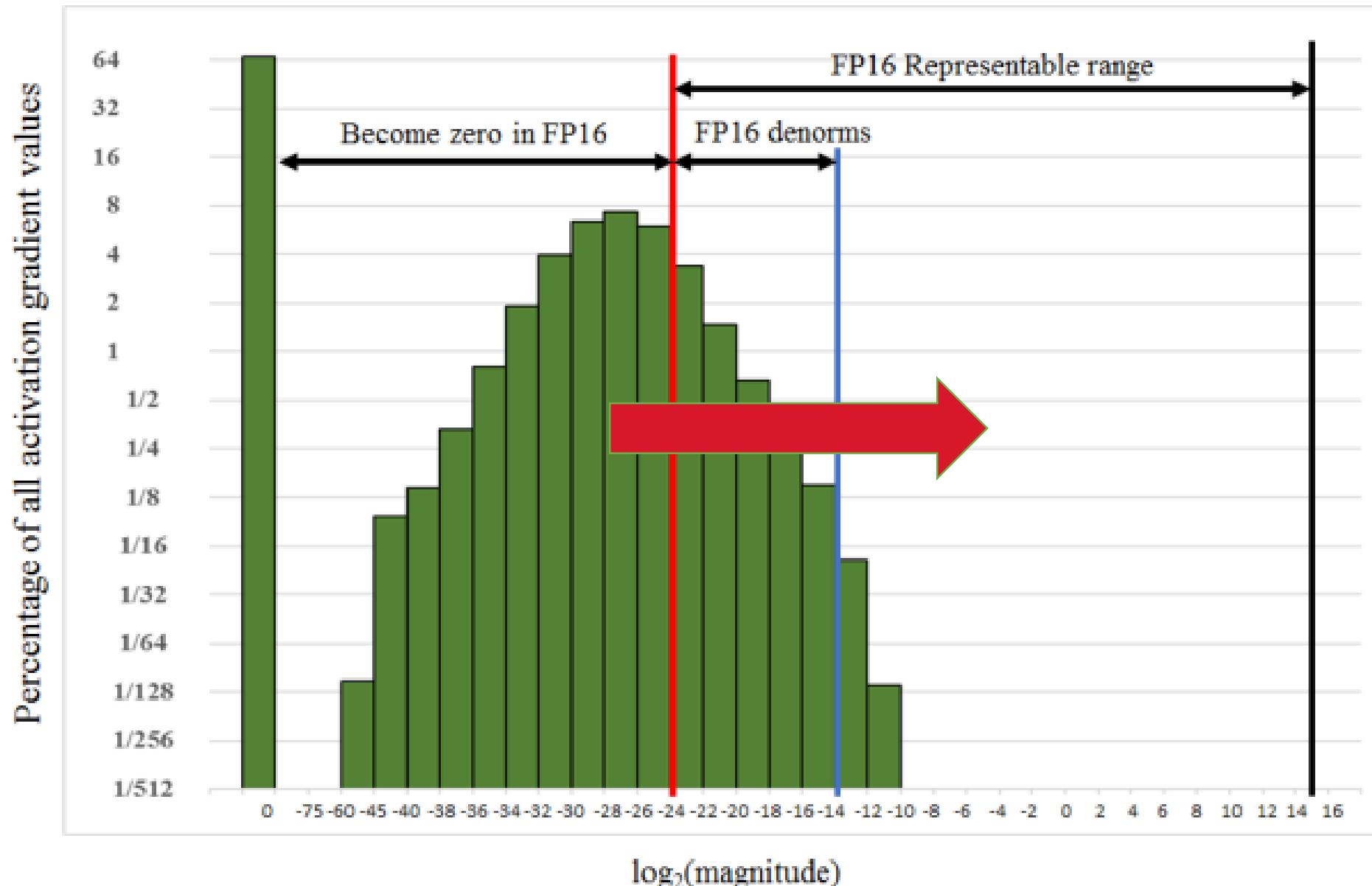
非规格化最小值

2^{-24}

混合精度解决方案: FP32 MASTER WEIGHTS

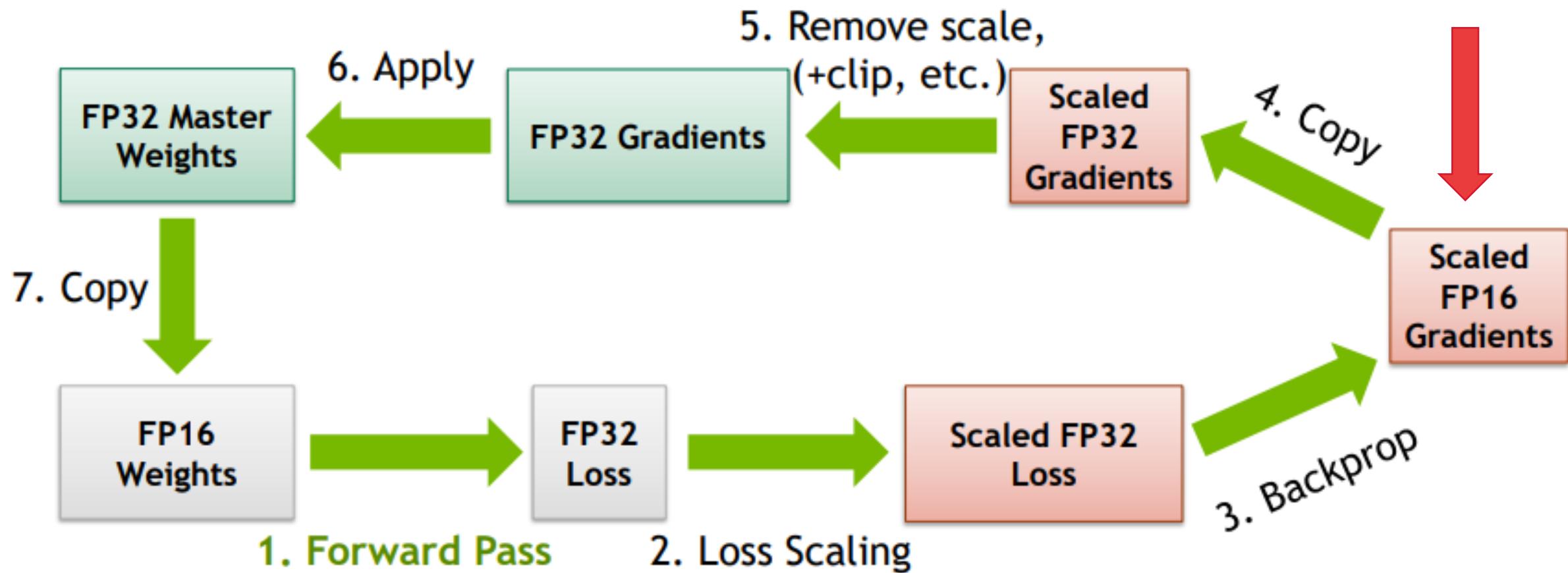


混合精度的问题：导数下溢

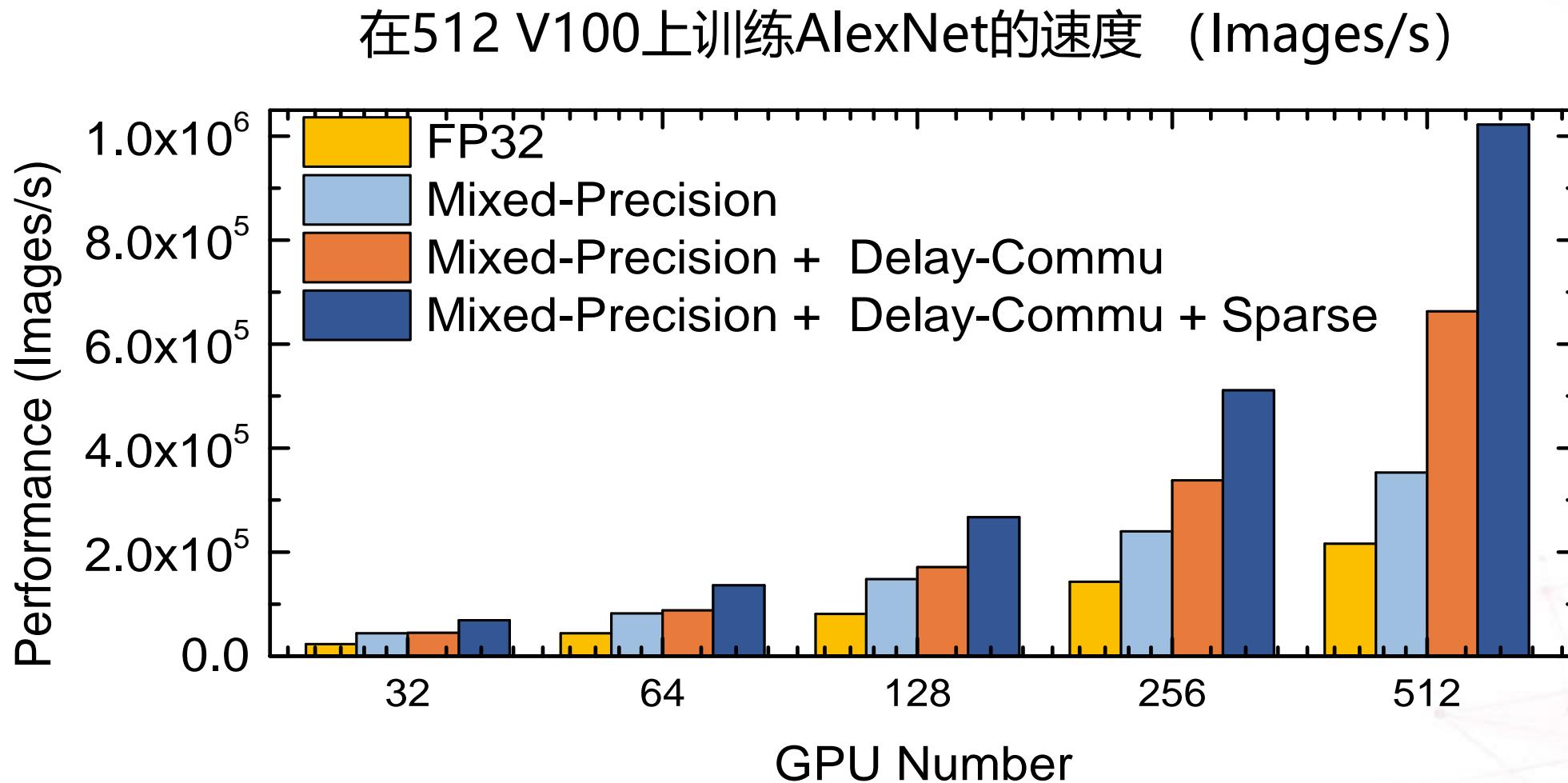


混合精度解决方案：LOSS SCALING

需要通信的部分



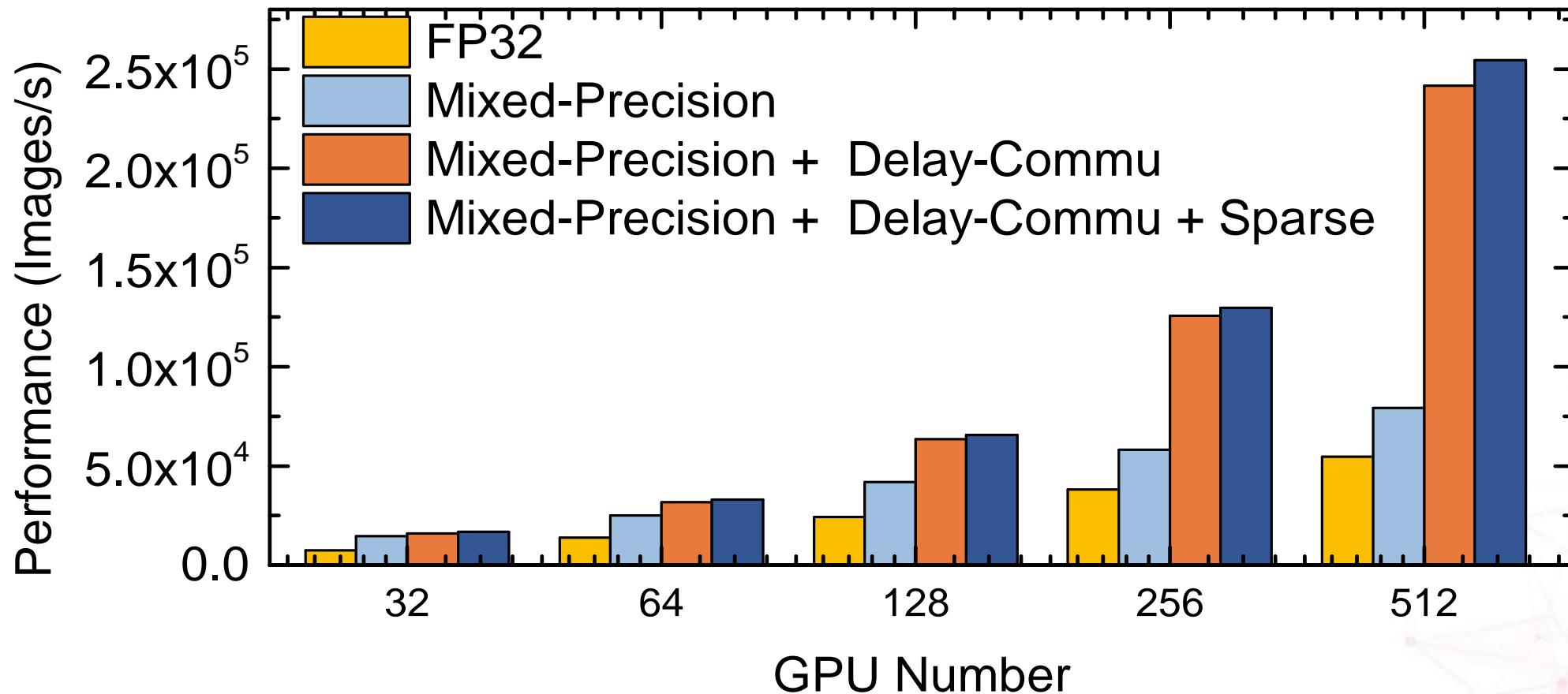
ImageNet training with AlexNet and Resnet-50



150秒内完成ImageNet/AlexNet 92 Epoch 训练，精度达到58.3%

ImageNet training with AlexNet and Resnet-50

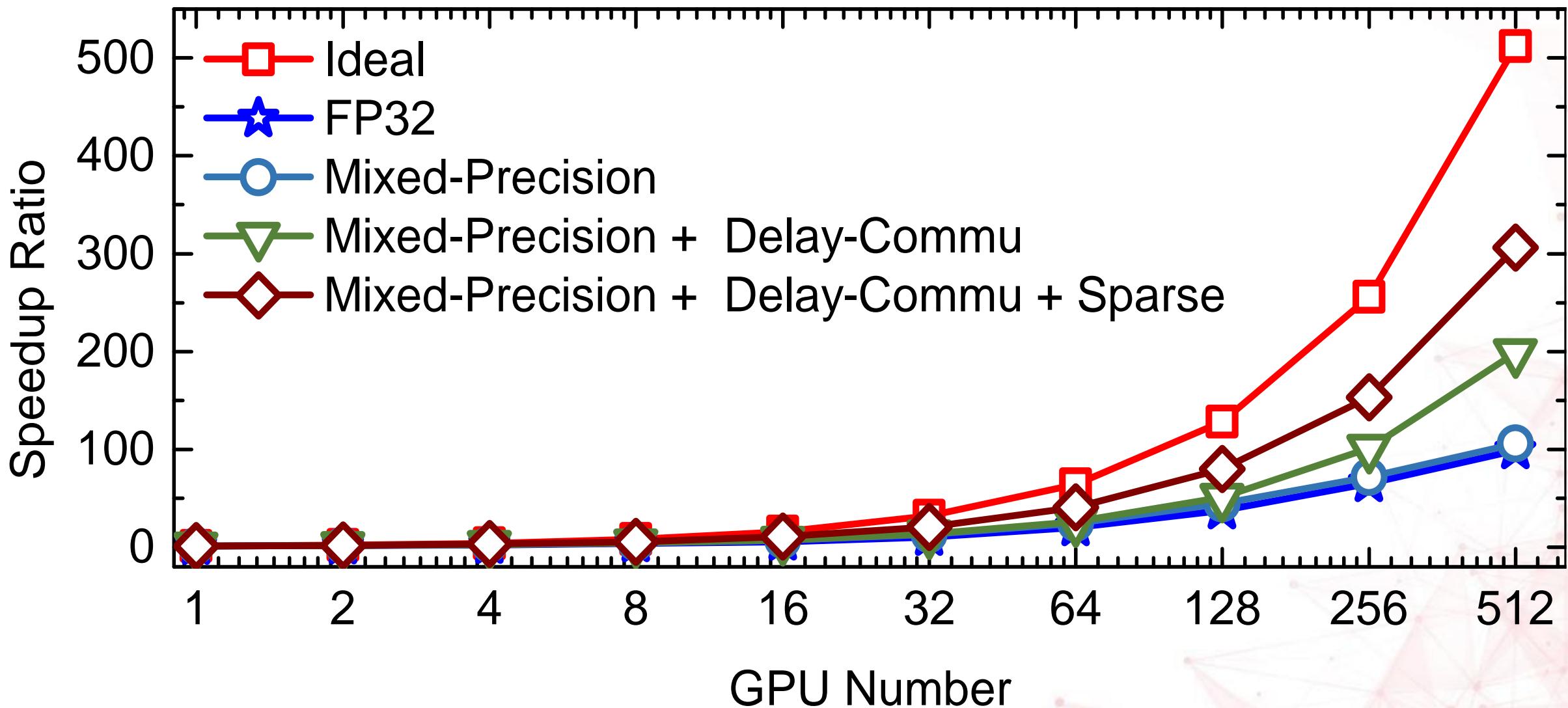
在512 V100上训练ResNet-50的速度 (Images/s)



480秒内完成ImageNet/ResNet 90 Epoch 训练，精度达到75%

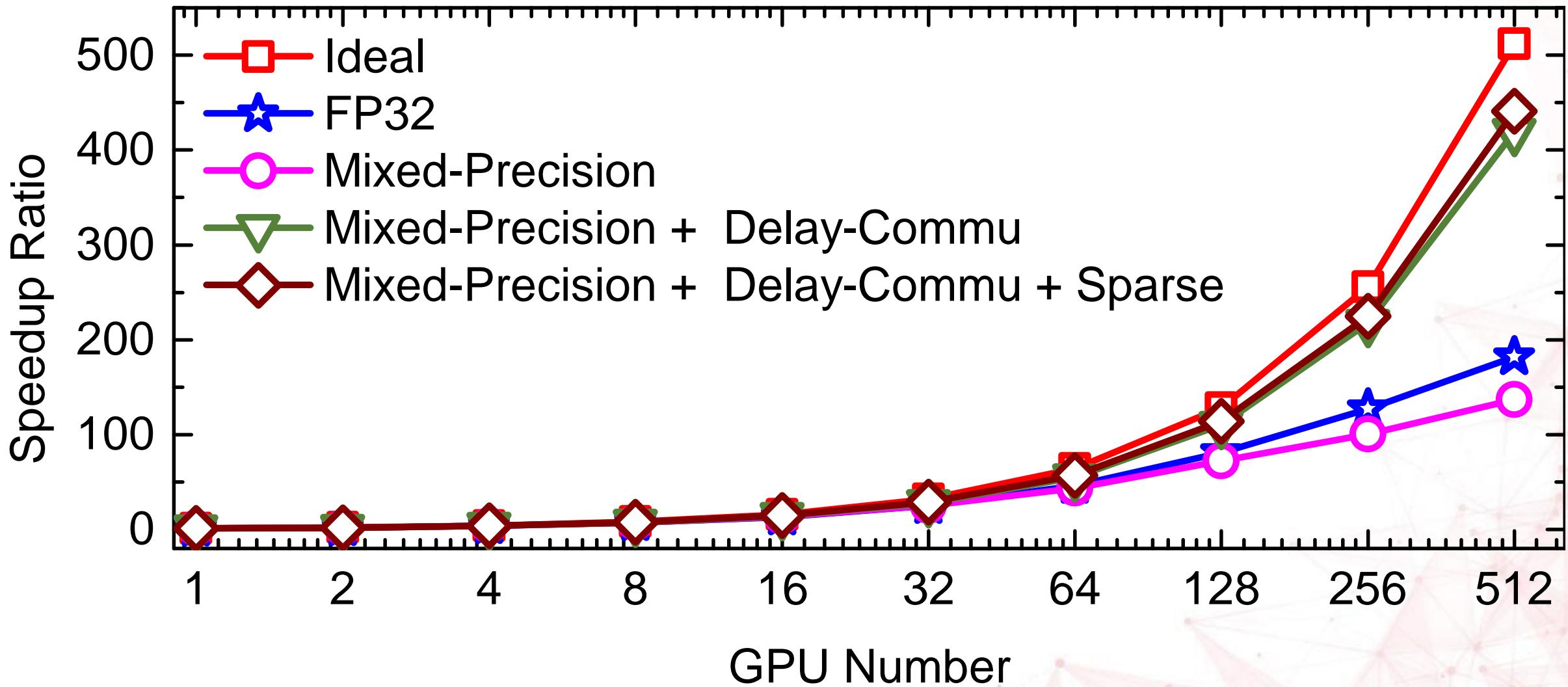
ImageNet training with AlexNet and Resnet-50

512 V100上训练AlexNet加速比



ImageNet training with AlexNet and Resnet-50

512 V100上训练ResNet-50加速比

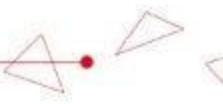


目录

1 //
2 //
3 //
4 //
5 //
6 //

- AI复兴：深度学习
- 计算优化：GPU
- 计算以外：内存&访存
- 携手共进：通信优化
- 无所不能：云上AI
- 轻松一刻：应用案例

Parrots综合算力平台



Parrots定位

完全商汤自研，提供**综合算力平台**，支撑AI研发和行业赋能

应用场景



安防



互娱



手机



自动驾驶



教育



医疗图像

...

算法/工具链

检测

分类

跟踪

关键点

分割

定点化

....

算力



综合算力平台

GPU/CPU



QUALCOMM

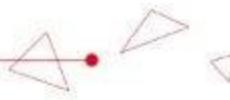
Storage



中科曙光
SUGON
HUAWEI

Interconnect





主要竞争对手趋势分析



Facebook



Google



Baidu

动态图支持

- PyTorch 用作前端语言

- 引入Eager模式

- 自定义Fluid语法，使用python编写 Fluid Program

计算性能

- Caffe2 用作底层执行

- 针对常用case优化Eager性能

- 使用c++内核执行Fluid Program

模型部署

- Scripts: 将python代码编译成中间表示

- 使用AutoGraph将python动态逻辑编译成静态图代码

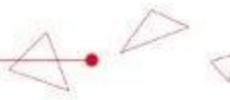
- Fluid Program做为中间表示，支持直接部署

特色/生态

- Glow: 跨平台底层引擎

- Lite: 端上部署套件
- TFX: 机器学习pipeline

- Cloud: 支持浏览器开发应用
- EDL: 虚拟化部署、调度训练



设计目标

运行效率

→ PaddlePaddle



- 先定义计算图、再执行
- 脚本编写复杂

Parrots

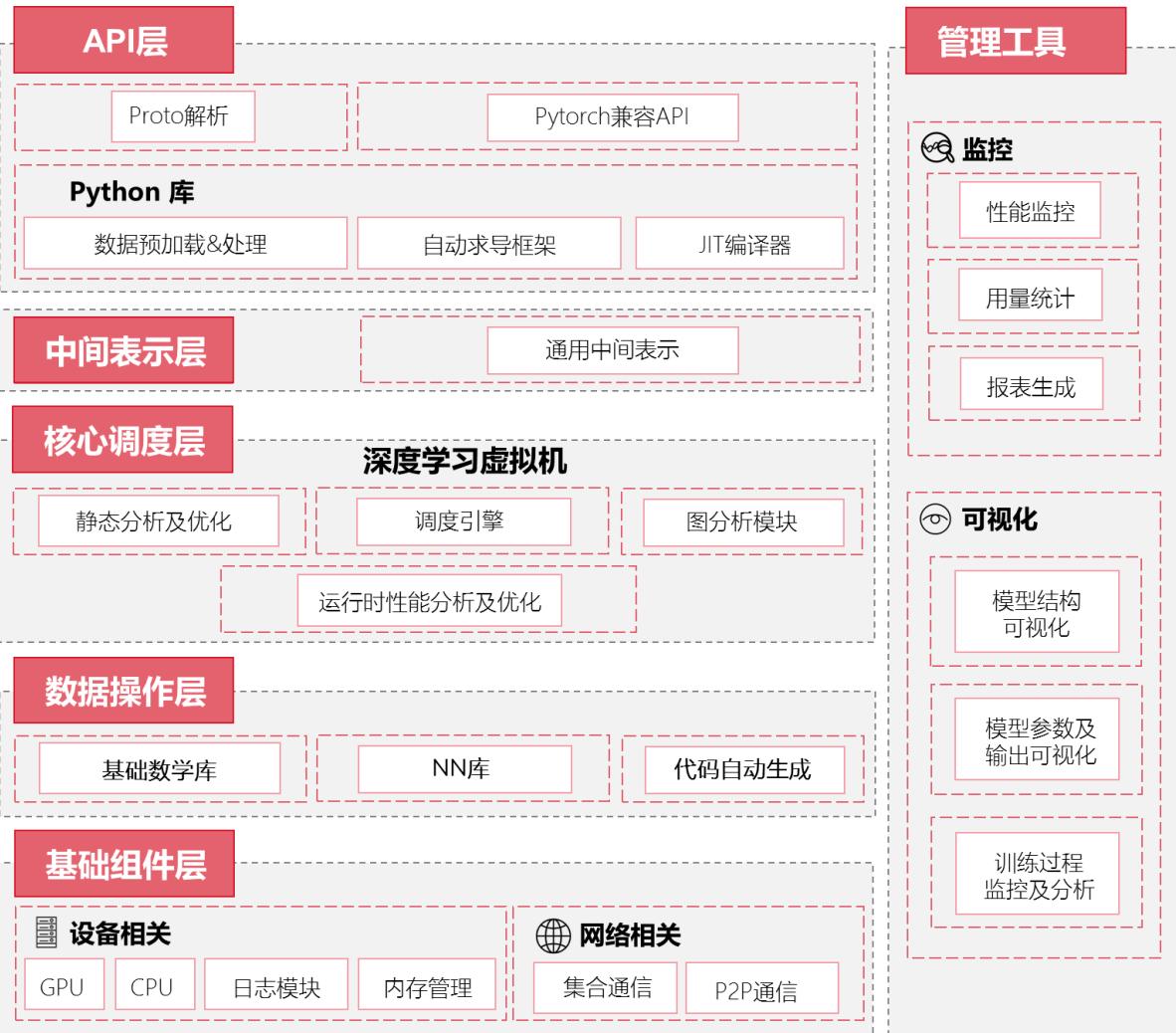
- 兼具现有框架优点
- 面向业务场景优化
- 凝聚商汤及合作伙伴的力量、构建行业品牌
- 无缝衔接研究到生产

○ PyTorch

- 分布式训练性能受限
- python核心——性能优化存在天花板

生产力

新一代Parrots



系统架构

● 主要特性



兼容Pytorch

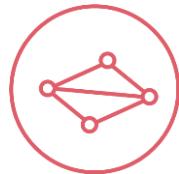
一键切换直接source



线性拓展

128卡加速比超过90%

● 关键技术



动态图并行

支持灵活定义模型结构



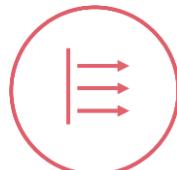
JIT编译

自动优化训练脚本



非阻塞执行引擎

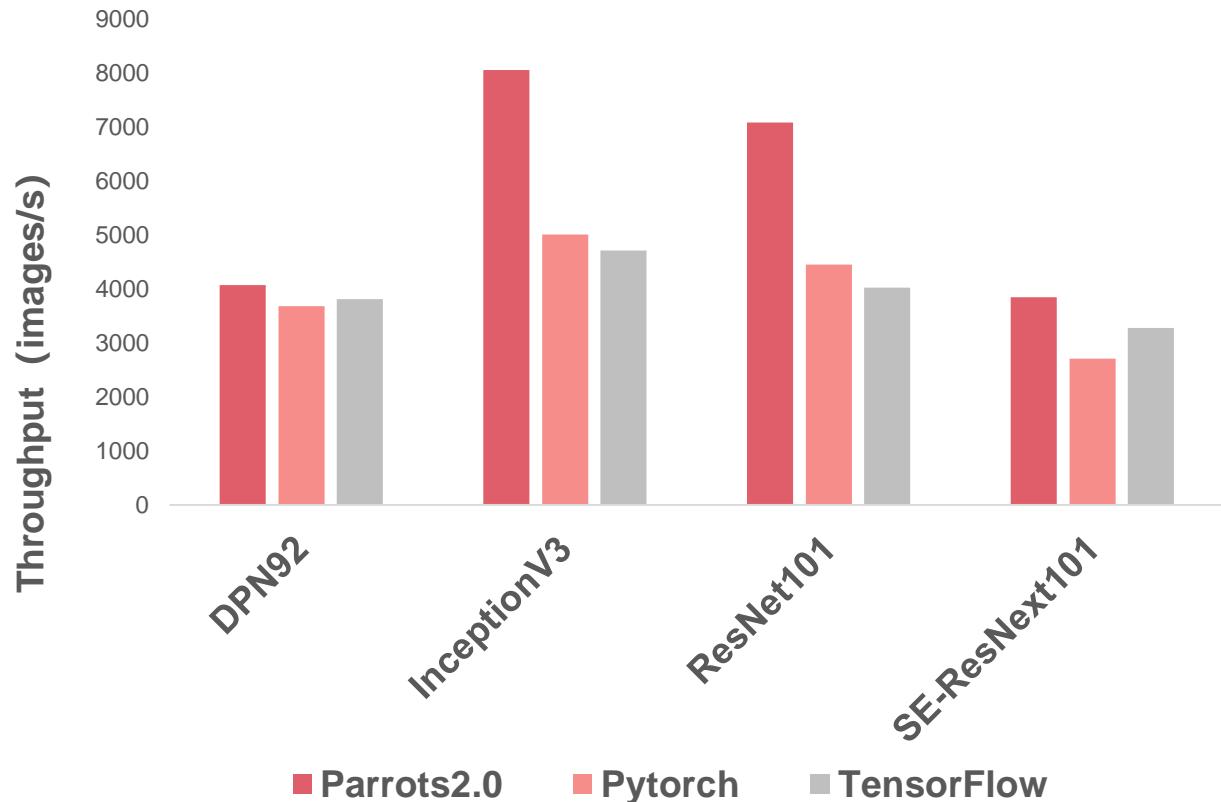
运算速度不受限于python



多stream并行

最大限度利用显卡算力

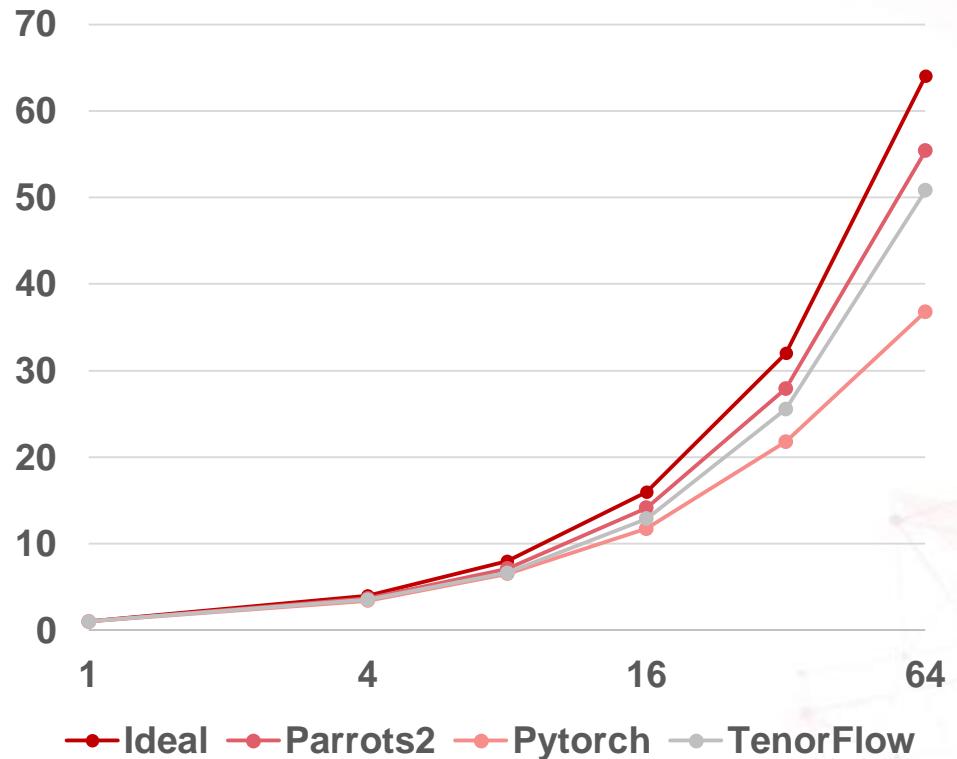
新一代Parrots



Training throughput on 64GPUs (images/s)

CPU: Intel(R) Xeon(R) CPU E5-2620 v4 @ 2.40GHz

GPU: NVIDIA Titan XP



Speedup on InceptionResNetV2 (Batchsize = 32)

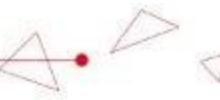
CPU: Intel(R) Xeon(R) CPU E5-2680 v4 @ 2.40GHz

GPU: NVIDIA GeForce GTX 1080 Ti

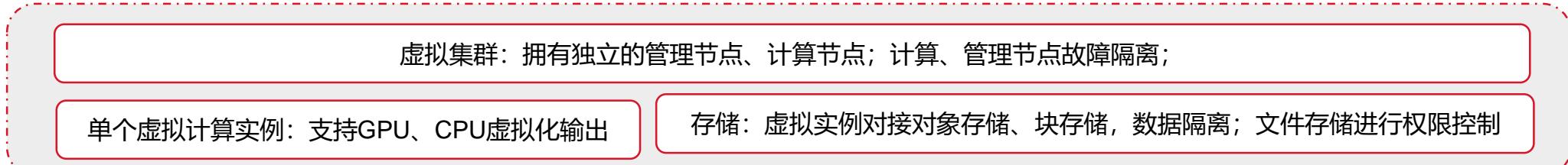
AI云原型系统概览



AI云原型系统概览



服务输出



虚拟抽象



存储资源



硬件基础



目录

1 //
2 //
3 //
4 //
5 //
6 //

- AI复兴：深度学习
- 计算优化：GPU
- 计算以外：内存&访存
- 携手共进：通信优化
- 无所不能：云上AI
- 轻松一刻：应用案例

轻松一刻：应用案例

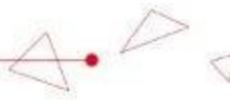


Image and Video Editing

1



De-hazing

- Restore pictures impacted by fog or haze to original appearance through industry leading image processing techniques

2



Super-Resolution

- Improve the resolution of photos without impacting quality, allowing clearer and stunning visual effects

3



Restoration (Fast De-noising)

- Achieve fast de-noising of images on mobile devices using leading algorithm for image de-noising

4



Dark Light Enhancement

- Improve overall brightness and definition of images or video being shot under dark light environment

5



Single Picture HDR

- Use one single picture to deliver impressive pro-level HDR effect shooting with mobile devices

6



De-focus Repair

- Leading image processing algorithm to deliver rapid repair and restoration of defocused pictures

7



Video Stabilization

- Get rid of the jitter of the video so that the user can control the magnitude of stabilization

8



Re-focus Ability

- Dual camera-based 3D visual technology to deliver the function of focusing after shooting & aperture after shooting

9



Filtering

- Deliver various special image effects. Currently capable of delivering over 70 filter effects and 6 artistic filtering effects

轻松一刻：应用案例

The screenshot displays the SenseFace surveillance system interface. At the top, there are tabs for Monitor, History, Task, Target, Statistic, and Tools. The date and time are shown as Apr 3 2016 11:56:16.

Task List: Shows a hierarchical list of tasks under Control Task 1, including local-video-1 and various surveillance tasks from task 1 to task 18.

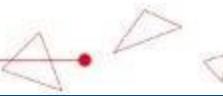
Monitor: Displays a live video feed from a camera labeled 08-13-015 at 16:12:52. The video shows several people on an escalator. A sign on the escalator reads "禁止跨越" (No crossing).

Captured: Shows a grid of 12 captured faces with their respective timestamps: 11:56:17, 11:56:16, 11:56:15, 11:56:08, 11:56:05, 11:55:57, 11:55:47, 11:55:47, 11:55:47, 11:55:47, 11:55:46, and 11:55:46.

Target: Shows three target monitoring panels for different dates and video feeds. Each panel includes a captured face, a target face with a confidence score, and a zoomed-in view of the target.

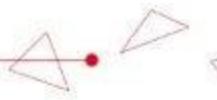
- 2016-04-03 11:55:39: local-video-1. Captured face. Target: chenzhao... (92.6%).
- 2016-04-03 11:54:56: local-video-1. Captured face. Target: suzhekun (91.7%).
- 2016-04-03 11:53:51: local-video-1. Captured face. Target: chenzhao... (93.9%).

轻松一刻：应用案例



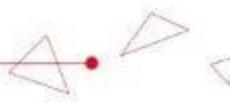
人群分析系统 Crowd Analysis System

轻松一刻：应用案例



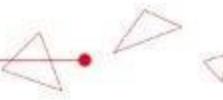
Enabling cars to sense the surrounding environment and recognizing objects

轻松一刻：应用案例



夜晚、大雾、恶劣天气等困难复杂场景对算法的高要求

轻松一刻：应用案例



谢谢！

