



# 人工智能系统 System for AI

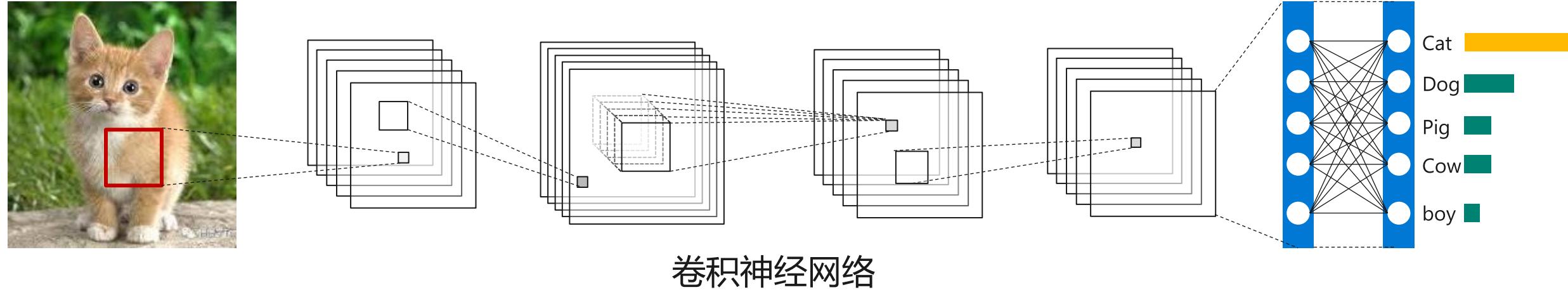
## 神经网络的稀疏化

Efficiency via compression and sparsity

# 主要内容

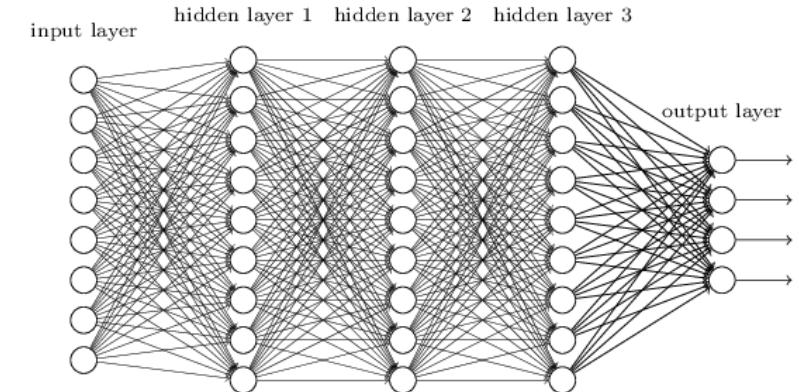
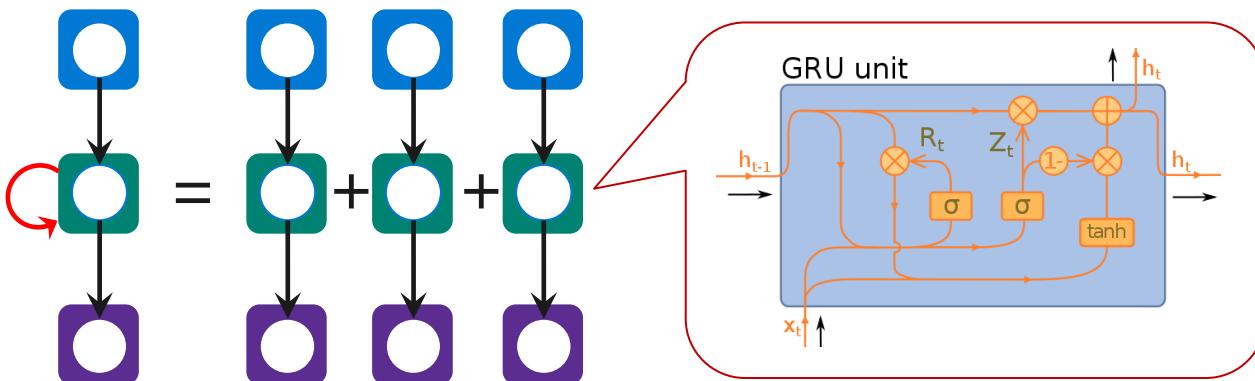
- 背景介绍
  - 模型的增长
  - 计算力的瓶颈
  - 模型压缩与稀疏的重要性
- 神经网络的稀疏性
  - 权重稀疏
  - 激活稀疏
  - 梯度稀疏
  - 稀疏与正则化
  - 量化

# 深度神经网络：多层的互联网络结构

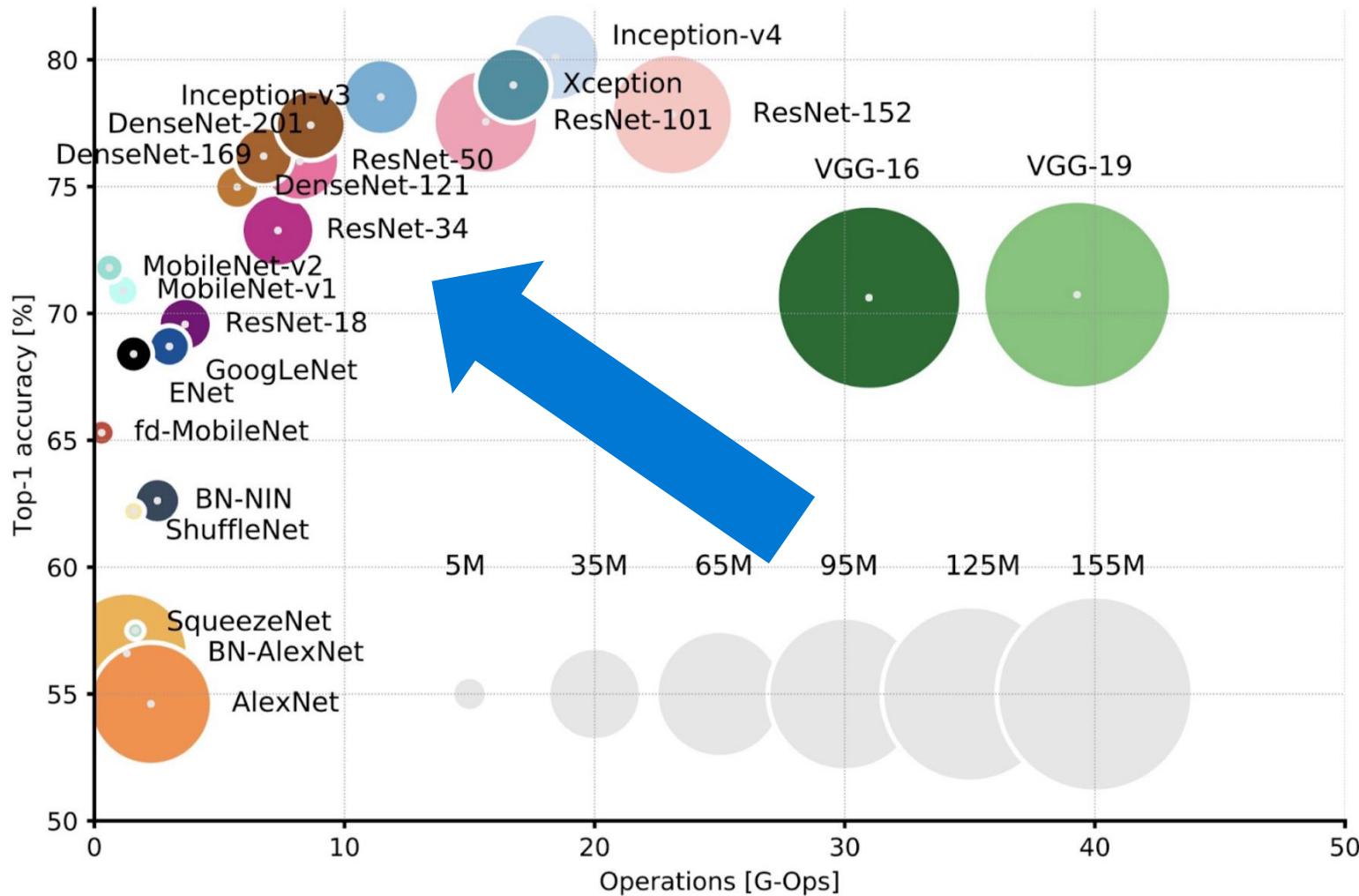


循环神经网络

全连接神经网络

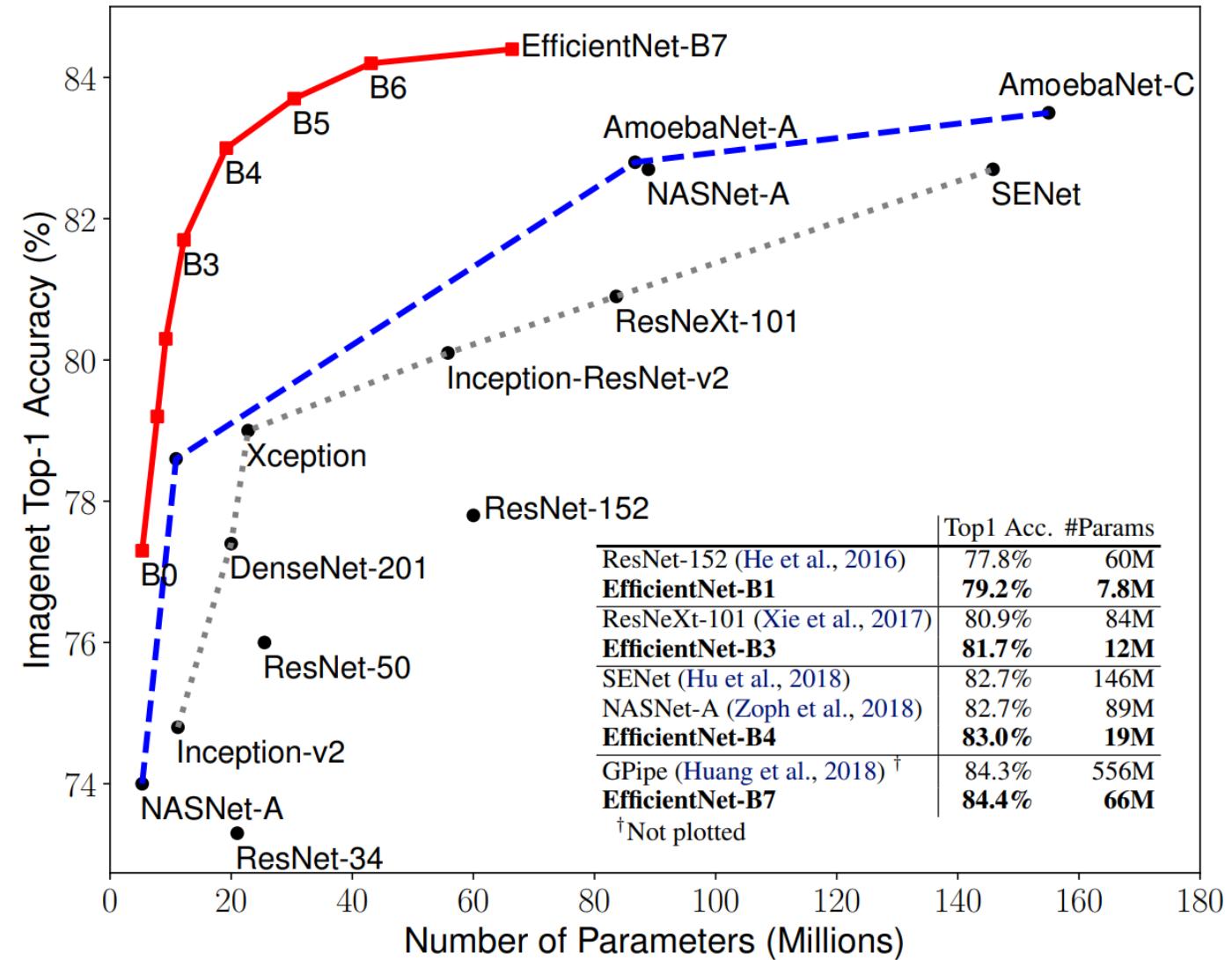


# 结构优化：寻找更高效的模型



# 增加模型维度是提高精度的主要手段

深度、宽度、分辨率



# 模型越来越大



BERT

*Language*

- 64 TPUv2 , • 8 P100
- 4 Days *or* • 365 Days
- 1000 GB , • 1000 GB



Wavenet

*Speech*

- 2 P100
- 6 Days
- 16 GB



Deformable CNN

*Vision*

- 8 P100
  - 10 Days
  - 64 GB
- 64 K80
  - 6 Days
  - 1500 GB

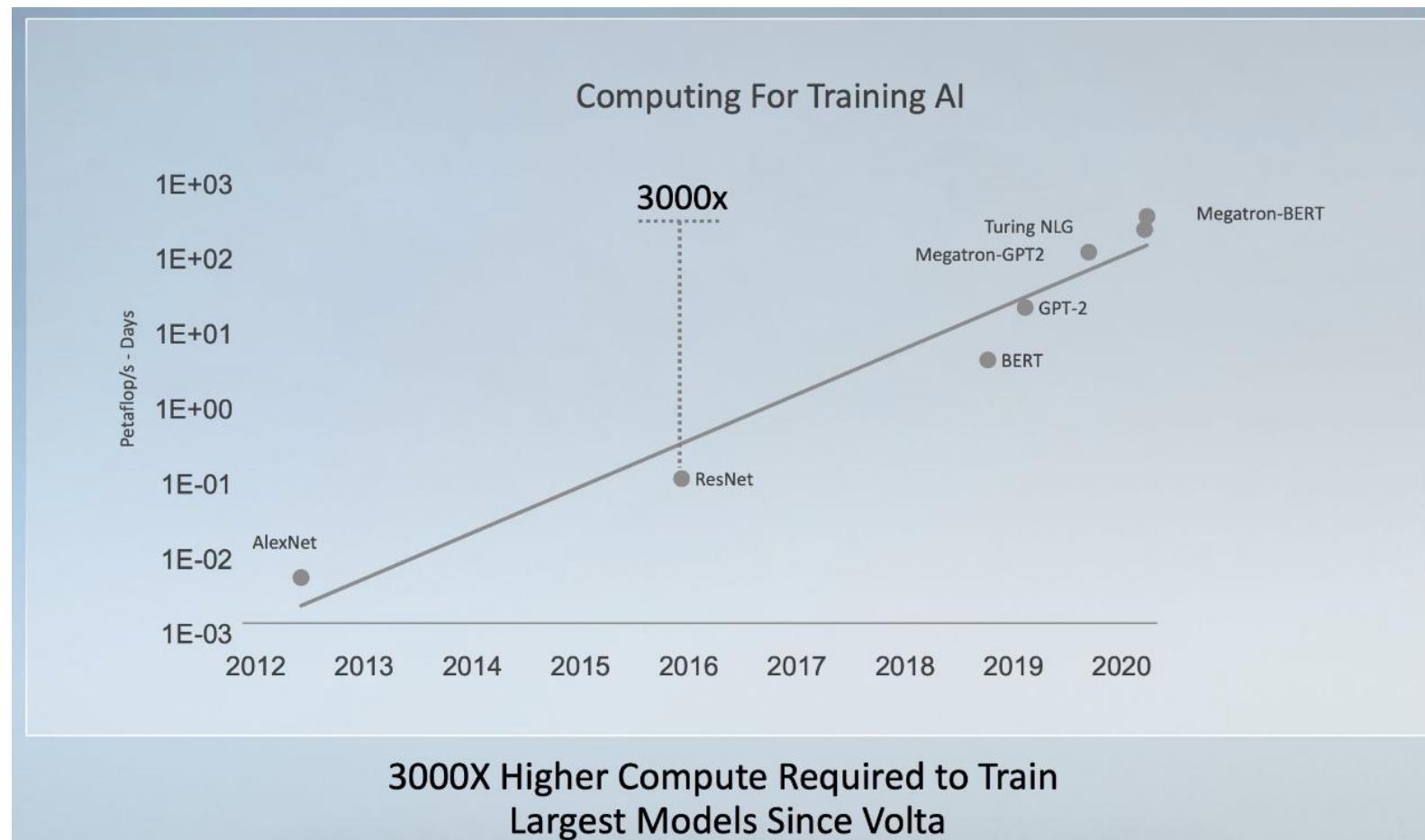


MoE

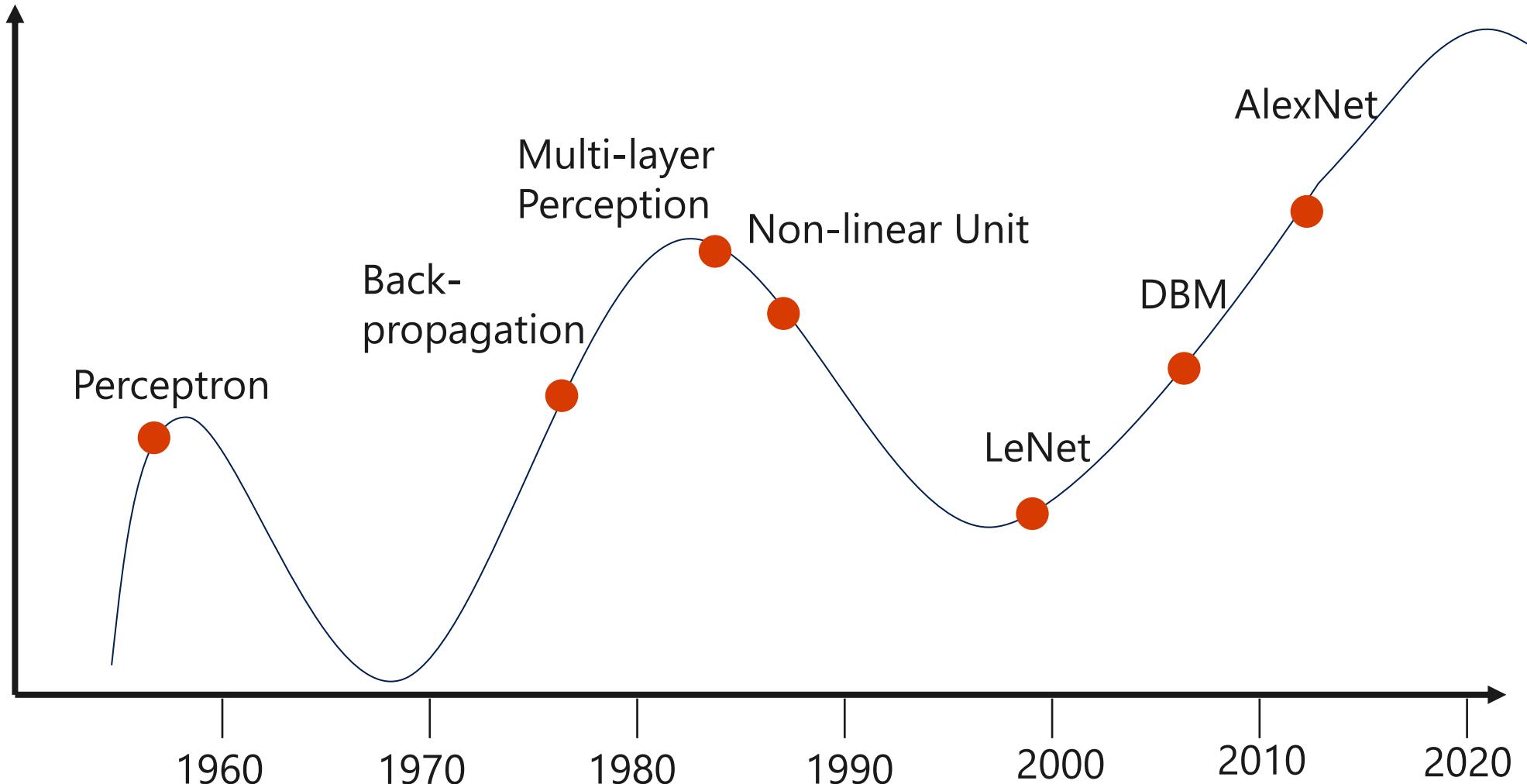
*Language*

# Volta GPU发布以来，大模型训练的算力要求增长了3000倍

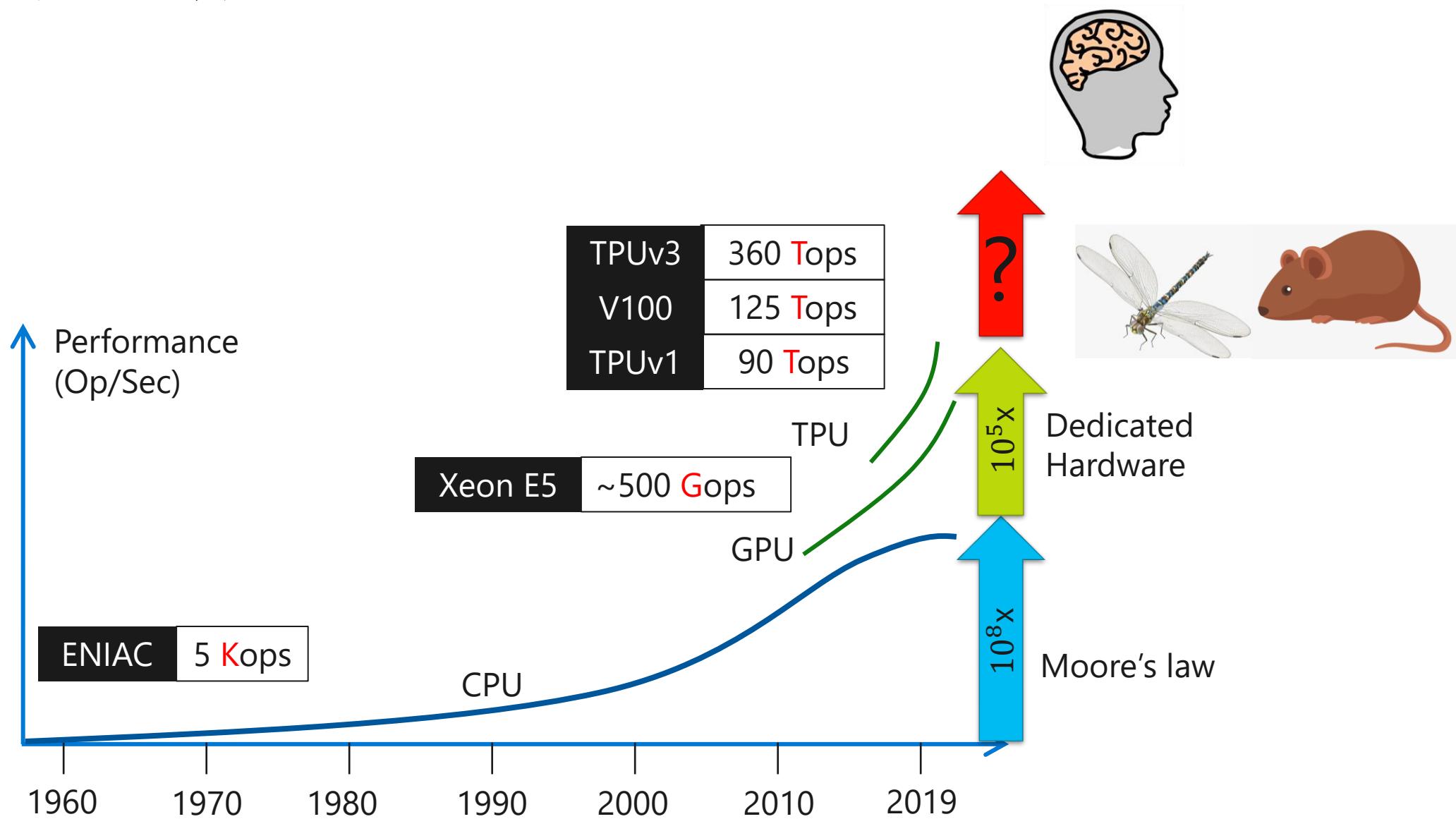
- GPT-3：一个亿人民币的模型



# 计算力是AI在二十一世纪初叶复兴的核心动力



# 计算力瓶颈

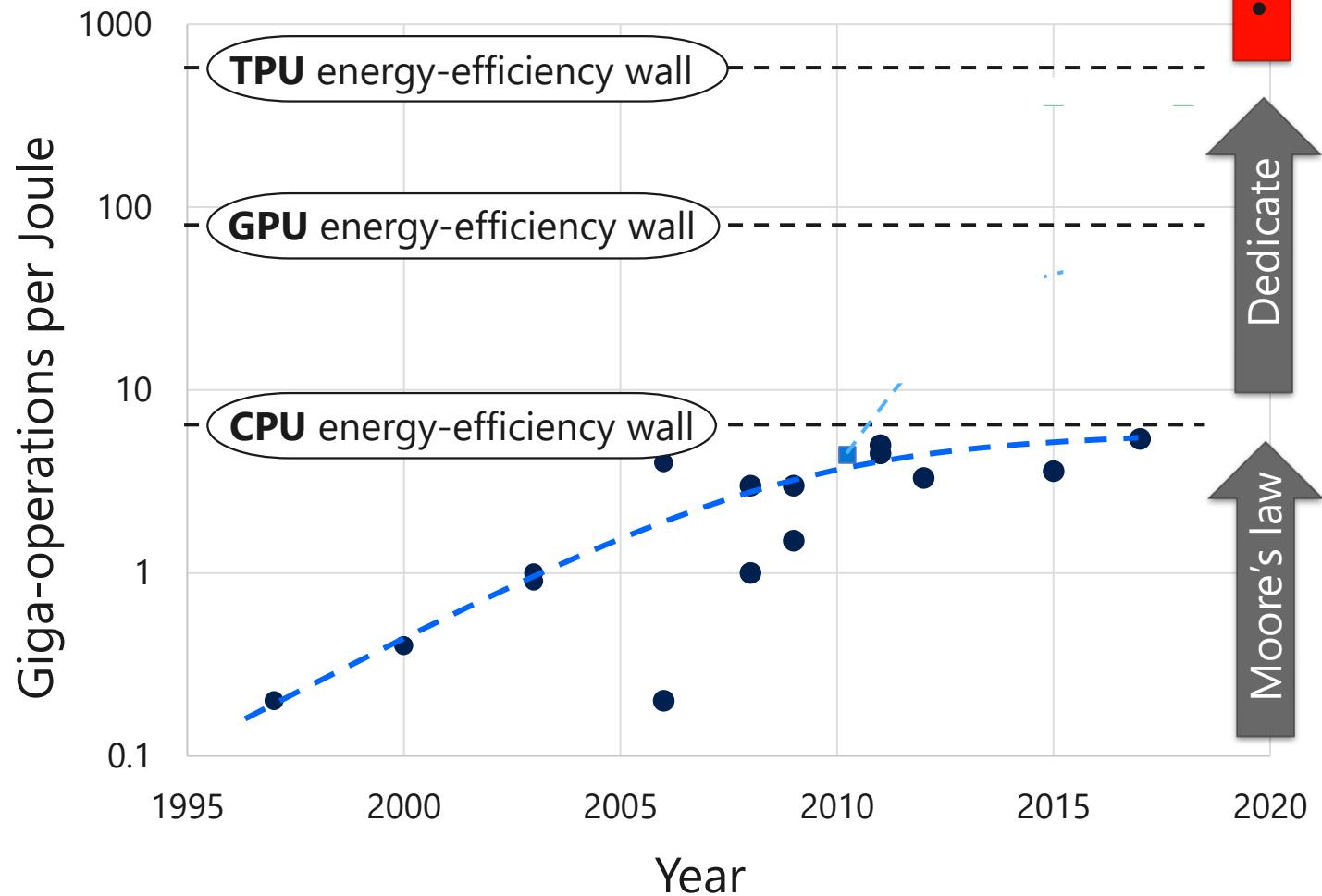


10

# 成本瓶颈



$\frac{10^{17} \sim 20}{20 \text{ Watt}}$

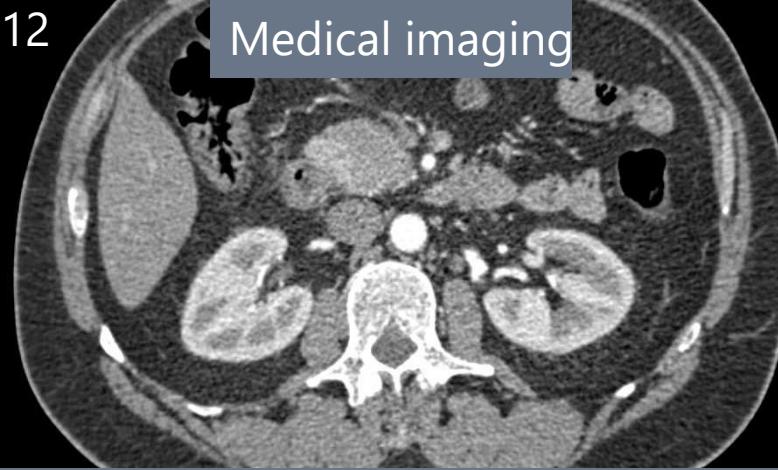


模型压缩与稀疏化  
(更高效的算法，软硬件协同)

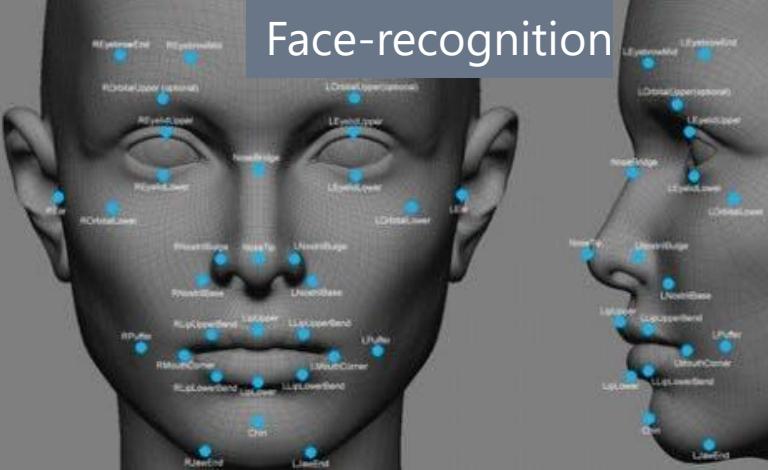
更大更强的处理器  
(锗基材料，量子计算)



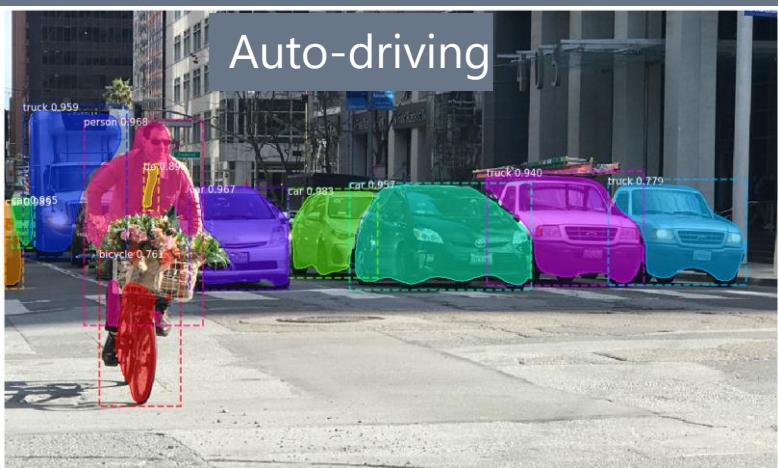
Medical imaging



Face-recognition



Auto-driving



Speech/Voice



LiDAR data

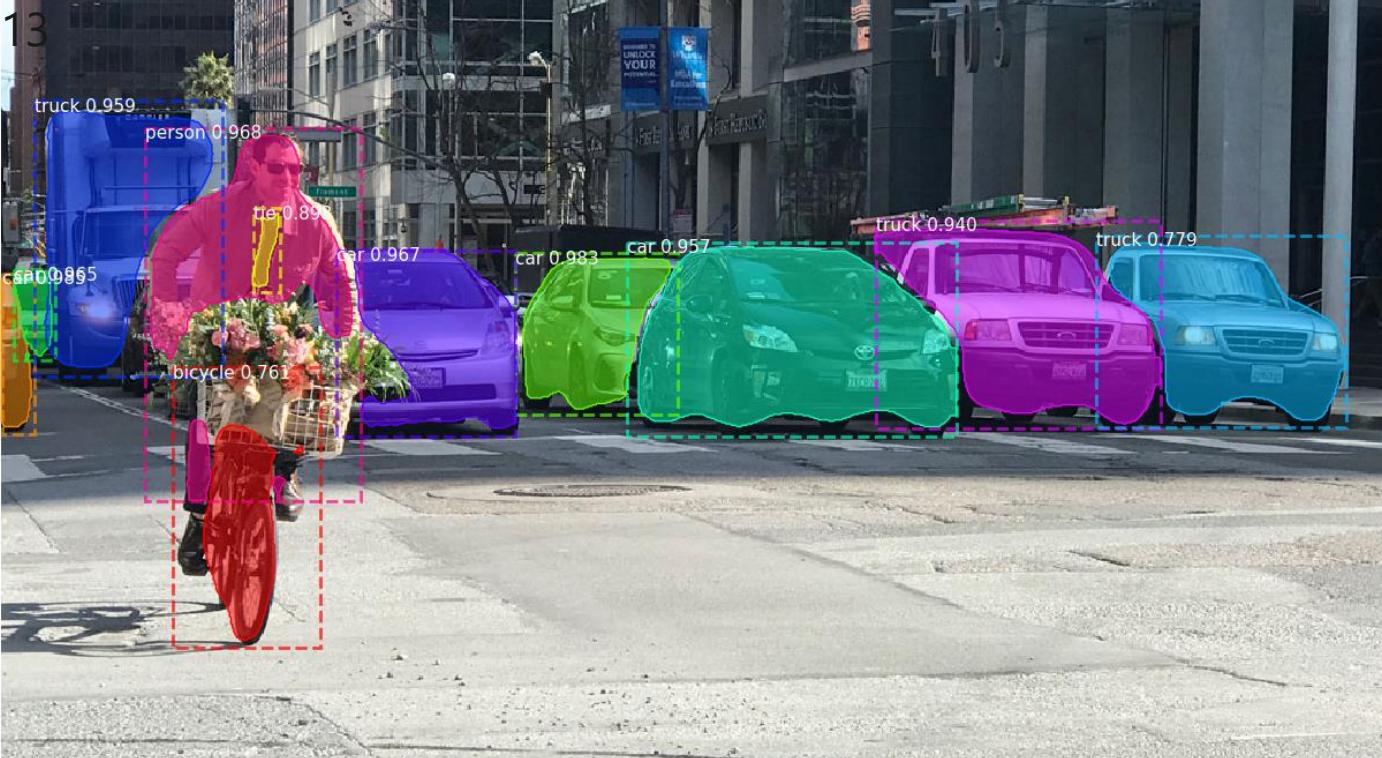


Social Networking



# 稀疏的来源： 高度结构化的数据

- 我们被无数的数据包围
- 每一种数据都有其内在结构性
- (自动地) 学习数据内生的结构性是人工智能算法的核心
- 数据的结构性带来信息表达中的稀疏性
- 高效的人工智能系统应该充分的利用稀疏性。



当人类看到眼前物体时，我们的神经系统：

- 不会处理所有像素；
- 关注主要物体；



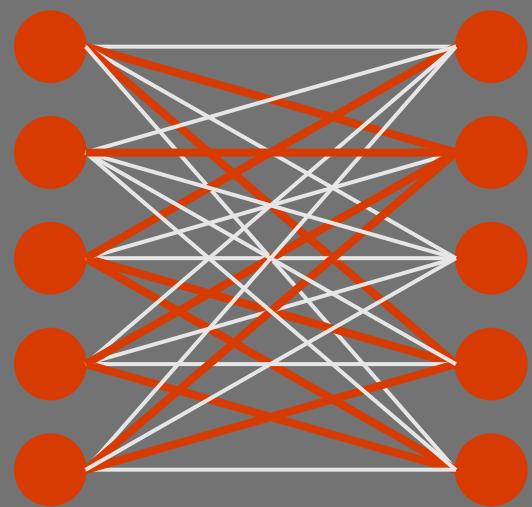
当人类识别一只猫时，我们：

- 不会仔细检查每一个毛发的纹理；
- 简单的几何边缘足以识别；

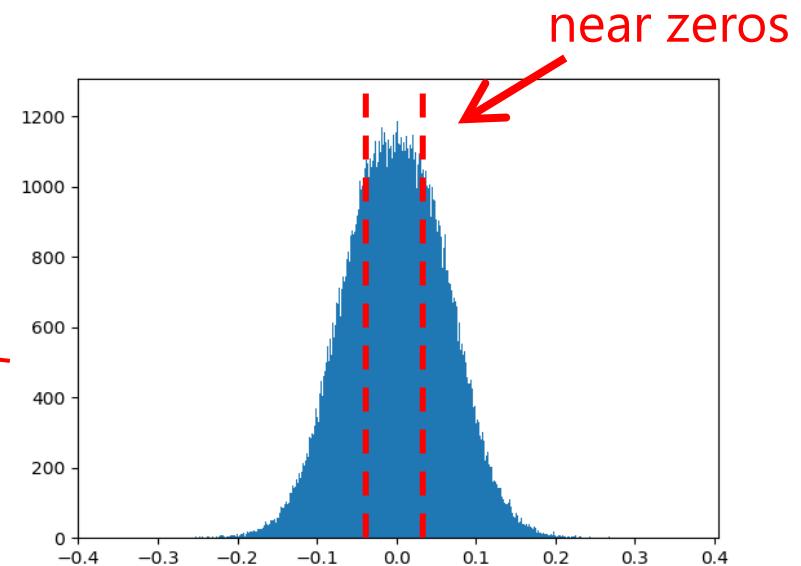
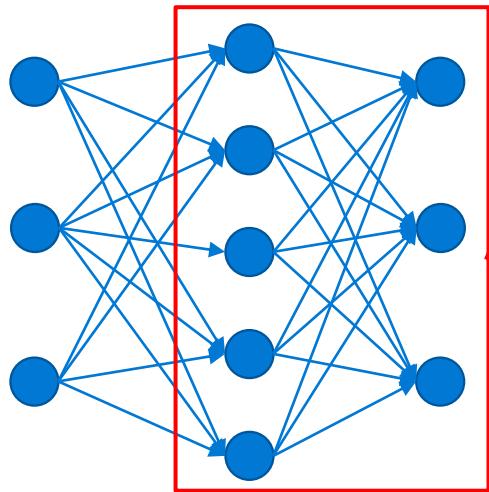
# 神经网络的稀疏性

Sparsity Types	Data Structure	Sources	Percentage of Sparsity	Speedup Potential																				
Static Sparsity	Weight	<p>Initial Weights of a layer</p> $f(w) = \begin{cases} w & \text{if }  w  > 0.5 \\ 0 & \text{if }  w  \leq 0.5 \end{cases}$ <p>Sparse Weights of a layer</p>	<p>Deep Compression AlexNet</p> <table border="1"> <thead> <tr> <th>Layer</th> <th>Sparsity (%)</th> </tr> </thead> <tbody> <tr><td>conv1</td><td>~10</td></tr> <tr><td>conv2</td><td>~60</td></tr> <tr><td>conv3</td><td>~65</td></tr> <tr><td>conv4</td><td>~65</td></tr> <tr><td>conv5</td><td>~60</td></tr> <tr><td>fc1</td><td>~95</td></tr> <tr><td>fc2</td><td>~95</td></tr> <tr><td>fc3</td><td>~90</td></tr> <tr><td>GeoMean</td><td>~65</td></tr> </tbody> </table>	Layer	Sparsity (%)	conv1	~10	conv2	~60	conv3	~65	conv4	~65	conv5	~60	fc1	~95	fc2	~95	fc3	~90	GeoMean	~65	1x~50x
Layer	Sparsity (%)																							
conv1	~10																							
conv2	~60																							
conv3	~65																							
conv4	~65																							
conv5	~60																							
fc1	~95																							
fc2	~95																							
fc3	~90																							
GeoMean	~65																							
Dynamic Sparsity	Activation	<p>Input features of ReLU layer</p> $f(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$ <p>Sparse input features of next layer</p>	<table border="1"> <thead> <tr> <th>Model</th> <th>Sparsity (%)</th> </tr> </thead> <tbody> <tr><td>Cifar-10</td><td>~35</td></tr> <tr><td>AlexNet</td><td>~65</td></tr> <tr><td>DeepComp</td><td>~50</td></tr> <tr><td>GoogleNet</td><td>~65</td></tr> <tr><td>VGG-16</td><td>~55</td></tr> <tr><td>ResNet-50</td><td>~55</td></tr> <tr><td>GeoMean</td><td>~50</td></tr> </tbody> </table>	Model	Sparsity (%)	Cifar-10	~35	AlexNet	~65	DeepComp	~50	GoogleNet	~65	VGG-16	~55	ResNet-50	~55	GeoMean	~50	1x~10x				
Model	Sparsity (%)																							
Cifar-10	~35																							
AlexNet	~65																							
DeepComp	~50																							
GoogleNet	~65																							
VGG-16	~55																							
ResNet-50	~55																							
GeoMean	~50																							
Gradient	<p>Error inputs to 2x2 Max Pooling layer</p> <p>Backward propagation</p> <p>Sparse error inputs to previous layer</p>	<p>CIFAR-10 DNN</p> <table border="1"> <thead> <tr> <th>Layer</th> <th>Sparsity (%)</th> </tr> </thead> <tbody> <tr><td>conv1</td><td>~85</td></tr> <tr><td>conv2</td><td>~85</td></tr> <tr><td>conv3</td><td>~75</td></tr> <tr><td>GeoMean</td><td>~80</td></tr> </tbody> </table>	Layer	Sparsity (%)	conv1	~85	conv2	~85	conv3	~75	GeoMean	~80	1x~10x											
Layer	Sparsity (%)																							
conv1	~85																							
conv2	~85																							
conv3	~75																							
GeoMean	~80																							

# 权重稀疏

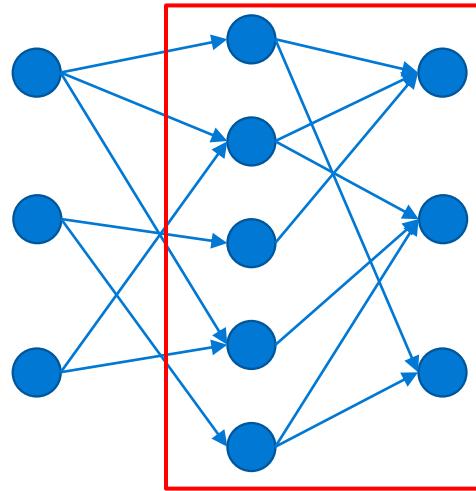


# 神经网络的权重 (weight) 稀疏

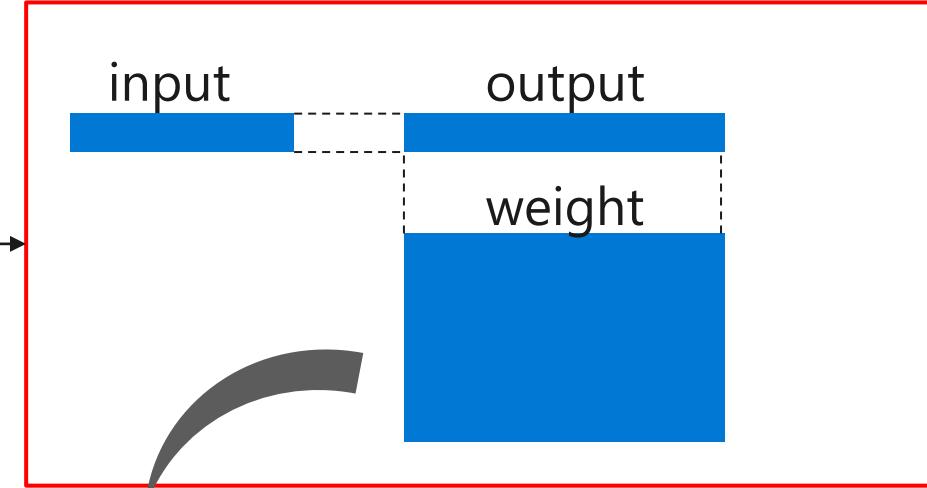


Data set: PTB, a language model  
1 Million training words,  
73k validation words,  
82k test words  
Model: 2-layer LSTM model,  
LSTM size = 1500

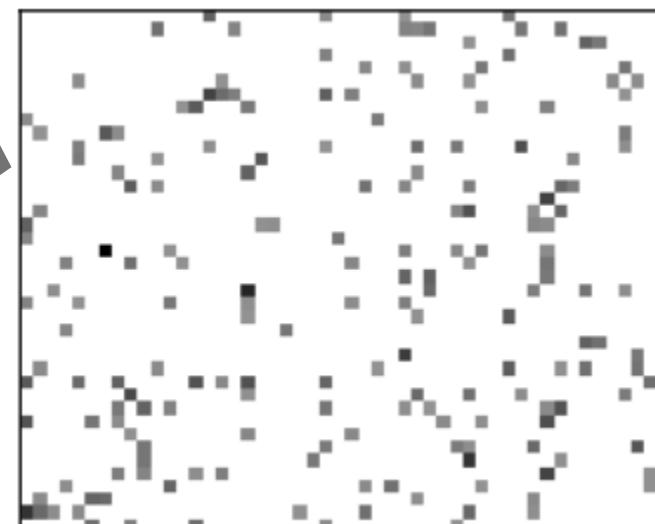
# 稀疏化的矩阵乘法表示



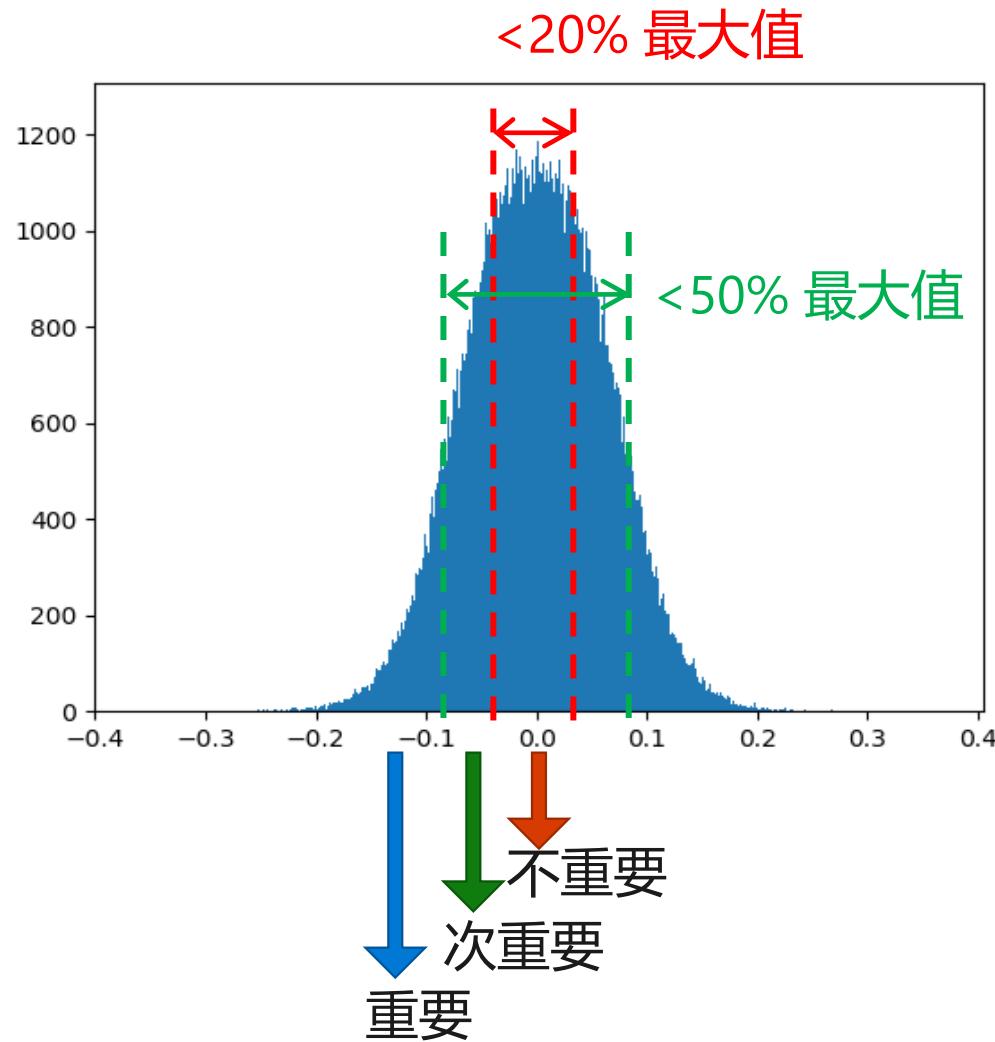
算子：  
矩阵乘法



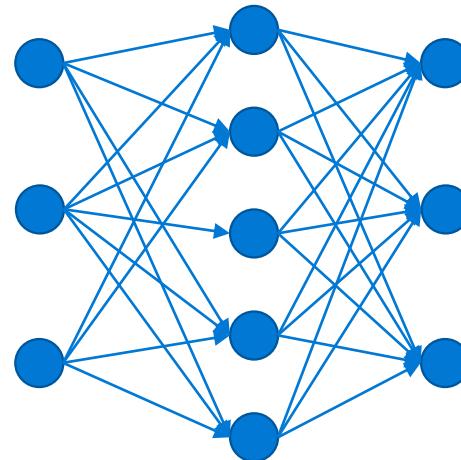
稀疏化



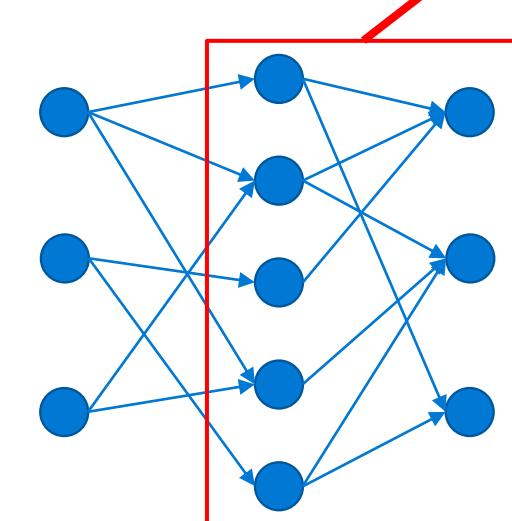
# 剪枝阈值



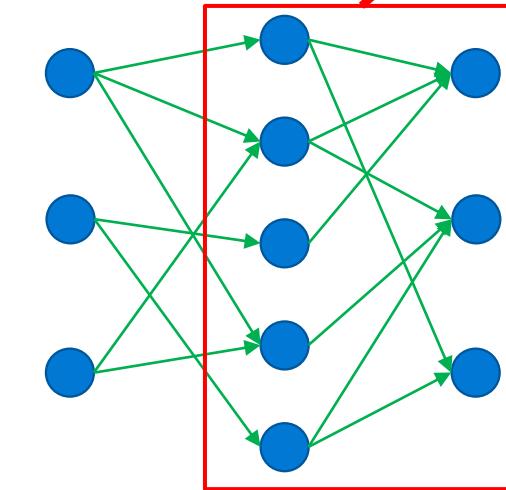
# 权重稀疏化的三个步骤



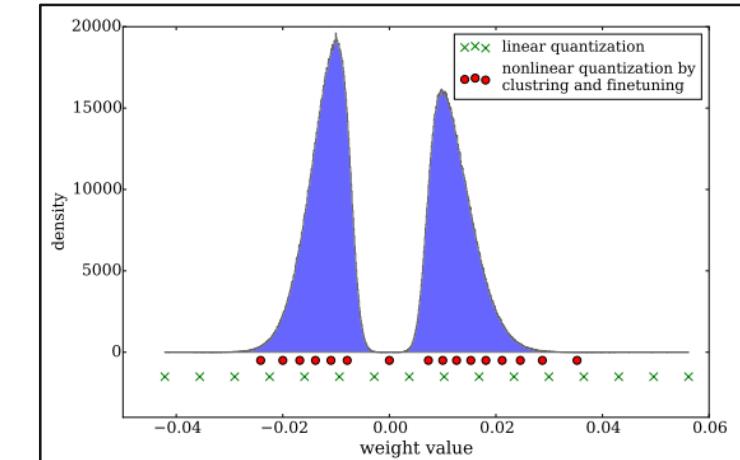
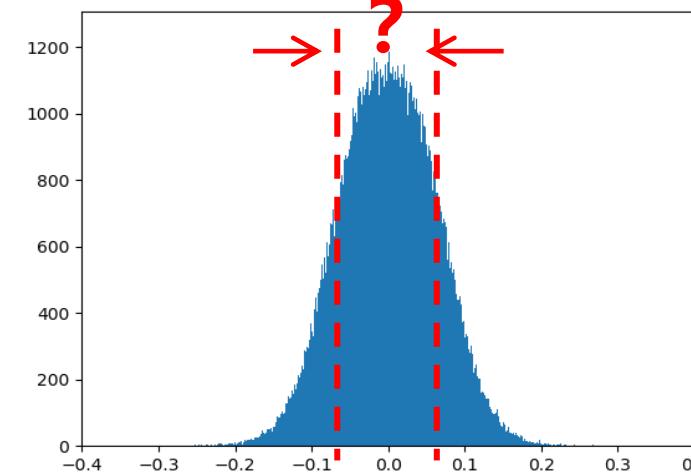
Train Connectivity



Prune Connections

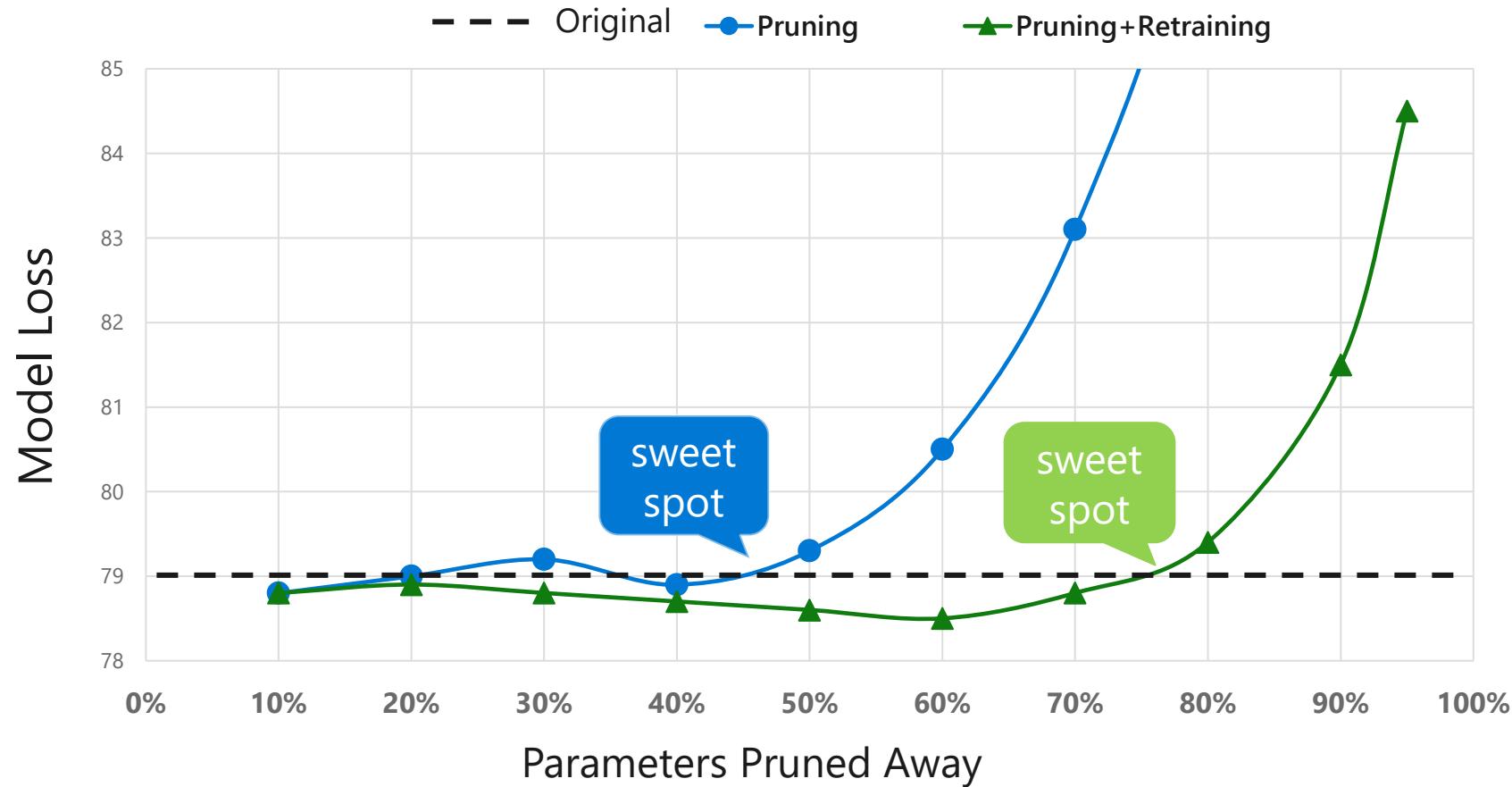


Fine-tuning



$\hat{w} \leftarrow \text{abs}(w)$   
if  $\hat{w}_i \leq \text{Threshold}$   
 $w_i \leftarrow 0$

# 通过精调（fine-tuning）提高模型稀疏度



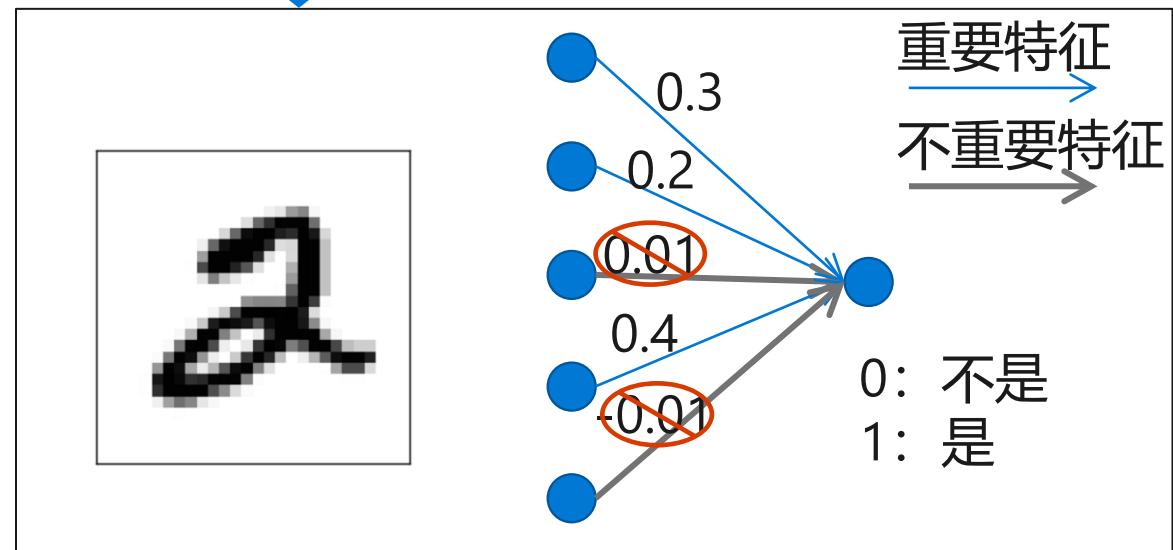
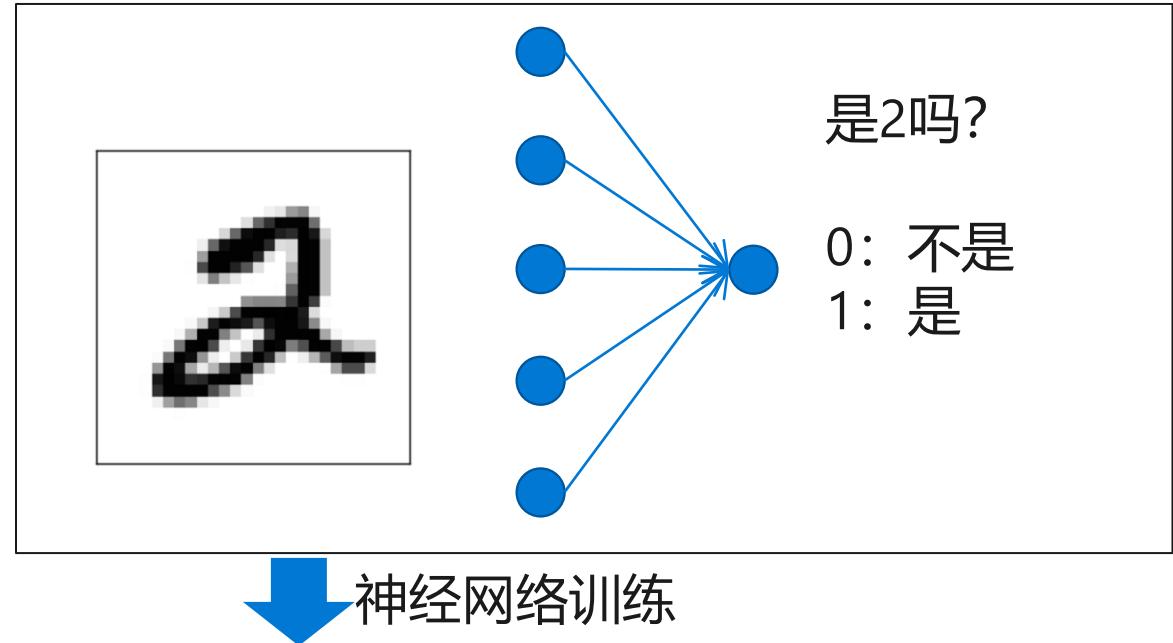
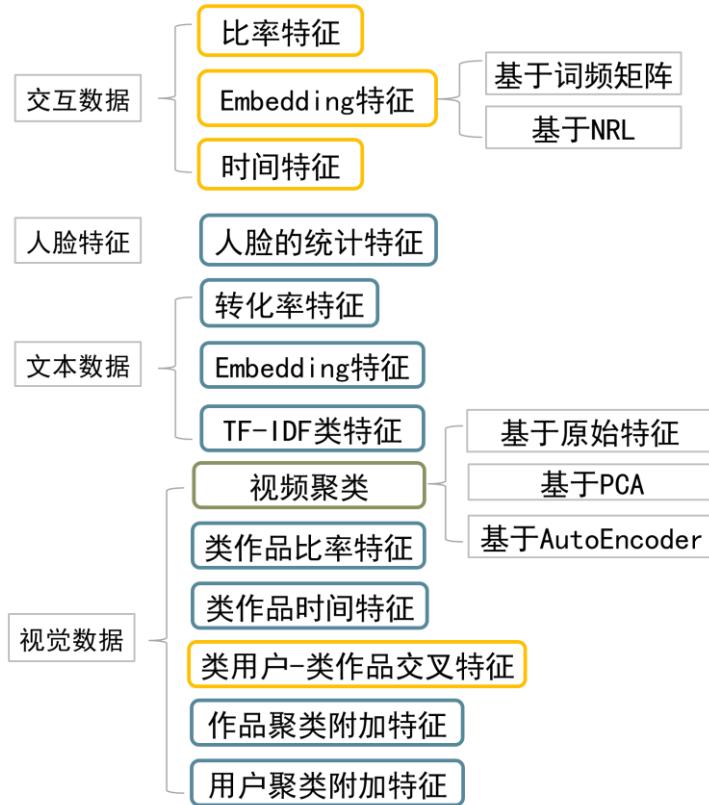
# 思考题

- 权重稀疏的本质是什么？

# 从特征工程到深度学习

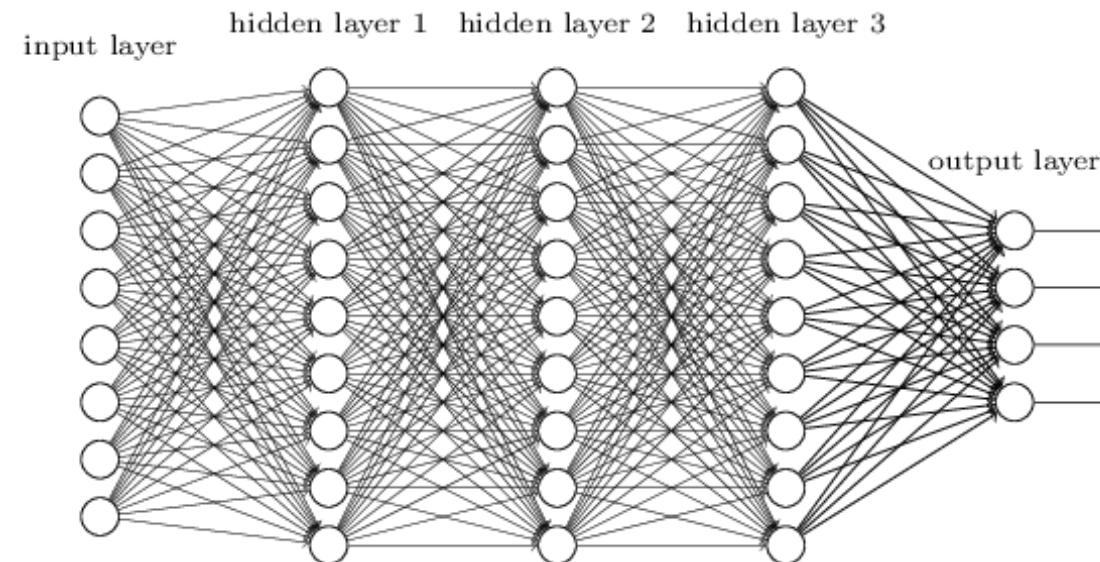
## 深度学习： 自动学习特征

### 传统方法： 人工选取特征

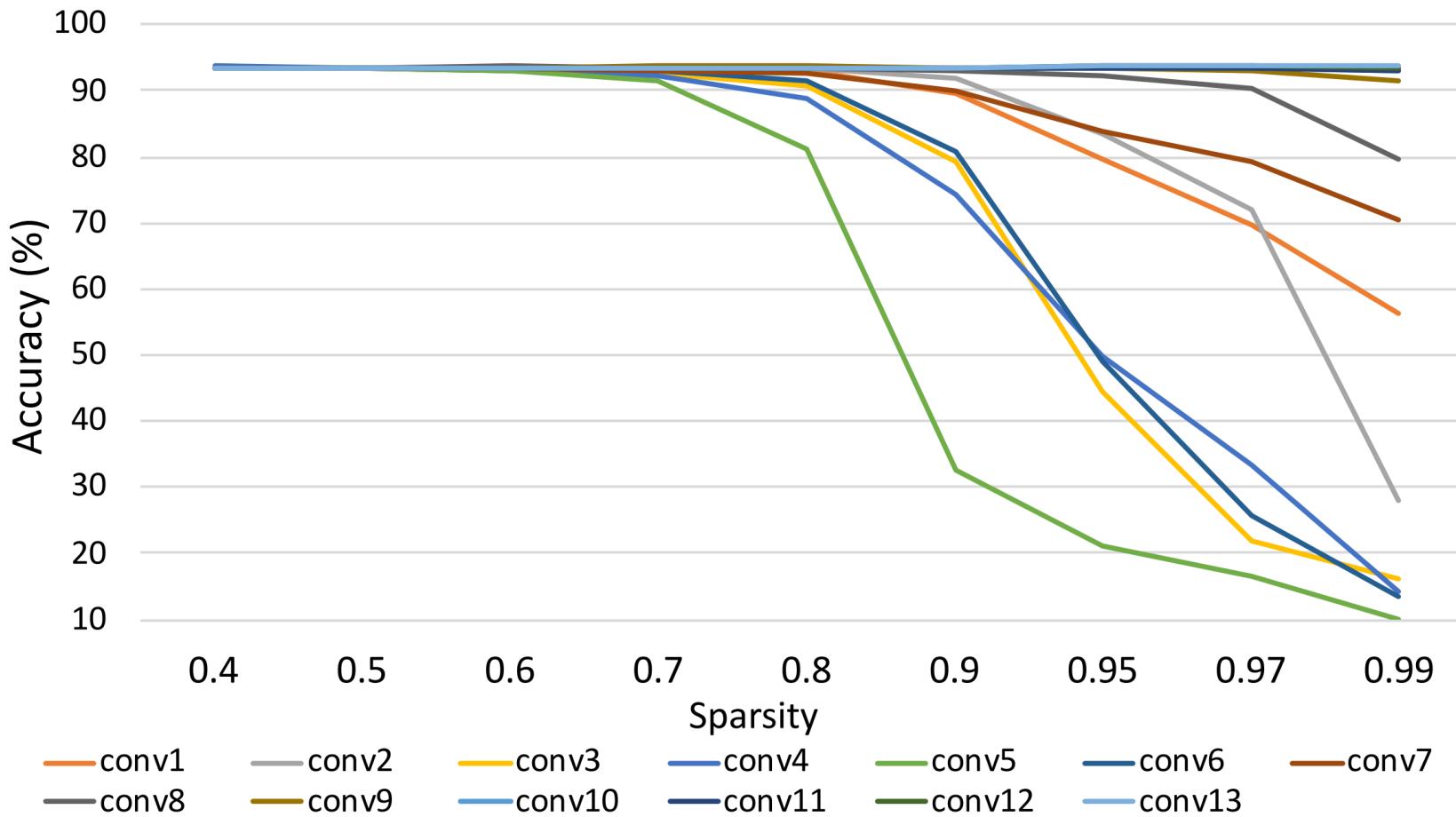


# 思考题

- 神经网络有许多层，每一层的稀疏度都一样吗？

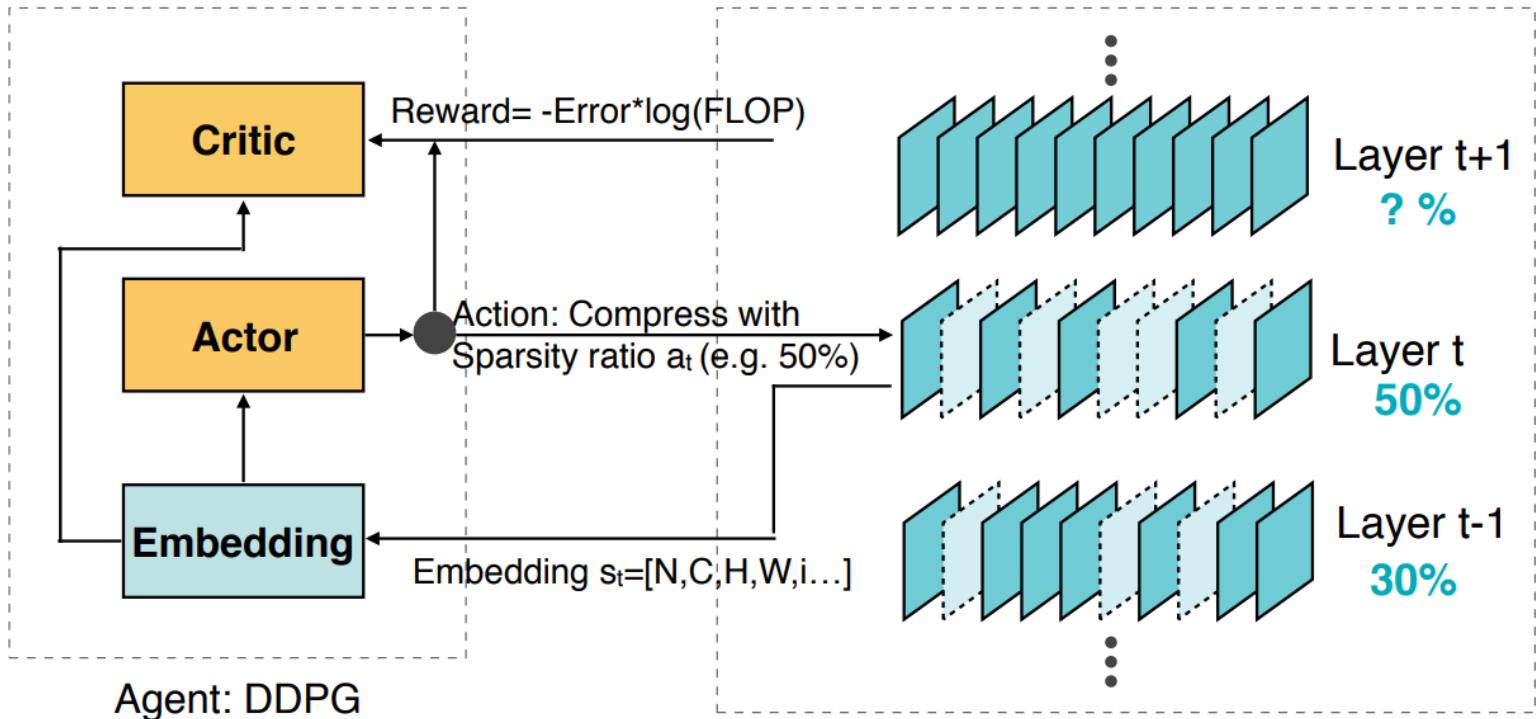
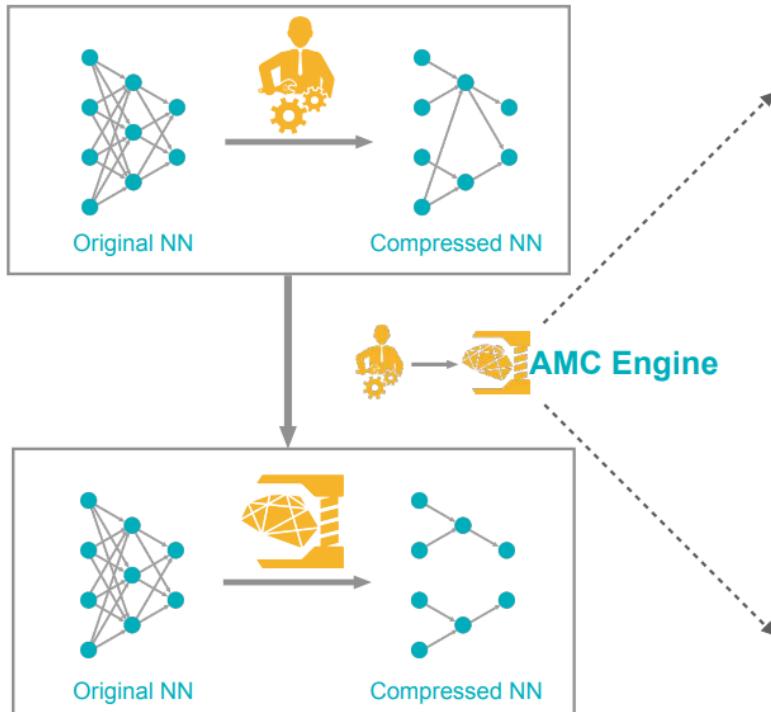


# 稀疏的敏感性：VGG-16



# 研究课题 1：各层最优稀疏度

Model Compression by Human:  
Labor Consuming, Sub-optimal



Model Compression by AI:  
Automated, Higher Compression Rate, Faster

Environment: Channel Pruning

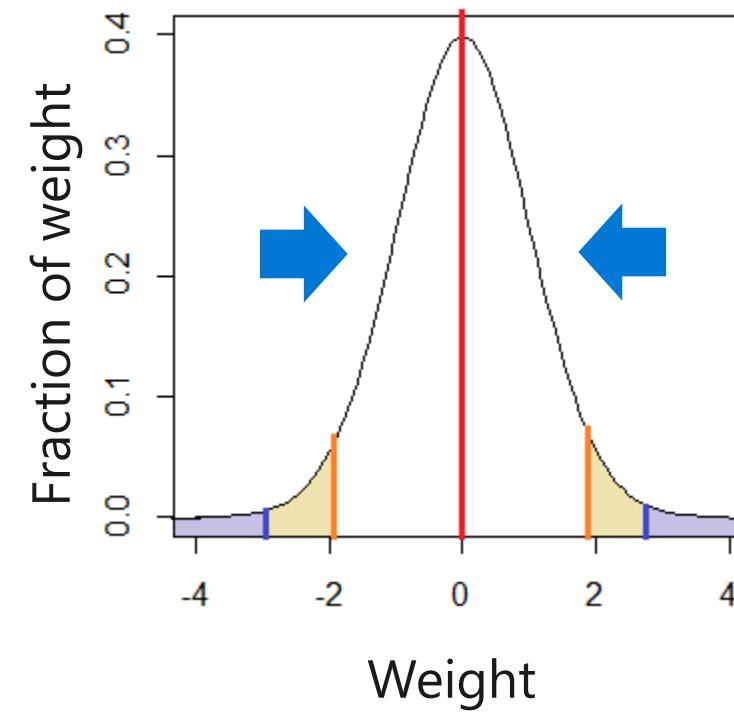
# 思考题

- 各层稀疏度和调参过程中的特征图数量的关系是什么？
- 如果某一层特别稀疏，是不是意味着该层特征图太多了？

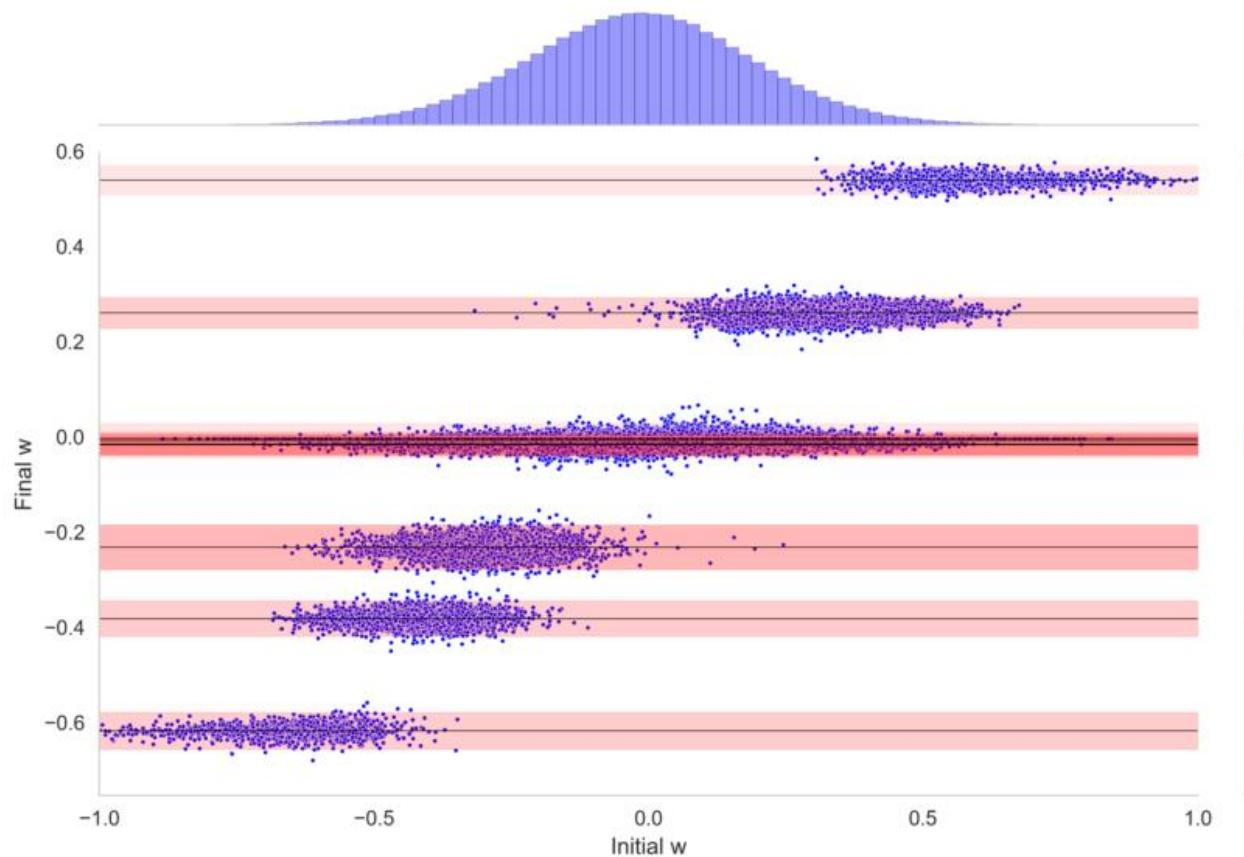
# 研究课题 2：稀疏与正则化

$$L1: \quad C = L(W, b) + \lambda \sum_w |w|$$

$$L2: \quad C = L(W, b) + \lambda \sum_w w^2$$



# 混合高斯模型



$$\mathcal{L}(q(\mathbf{w}), \mathbf{w}) = \underbrace{\mathbb{E}_{q(\mathbf{w})} [-\log p(\mathcal{D}|\mathbf{w})]}_{\mathcal{L}^E} + \underbrace{\text{KL}(q(\mathbf{w})||p(\mathbf{w}))}_{\mathcal{L}^C}$$

$$\mathbb{E}_{q(\mathbf{w})} [-\log p(\mathbf{w})] - H(q(\mathbf{w}))$$

$$= - \int_{\mathbb{R}^I} \mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma \mathbf{I}) \log \mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma \mathbf{I})$$

# 思考题

- 对稀疏来说，L1和L2正则化的区别是什么？
- 还有哪些正则化方法？

# 研究课题 3：硬件友好的稀疏化方法

1. Train Connectivity

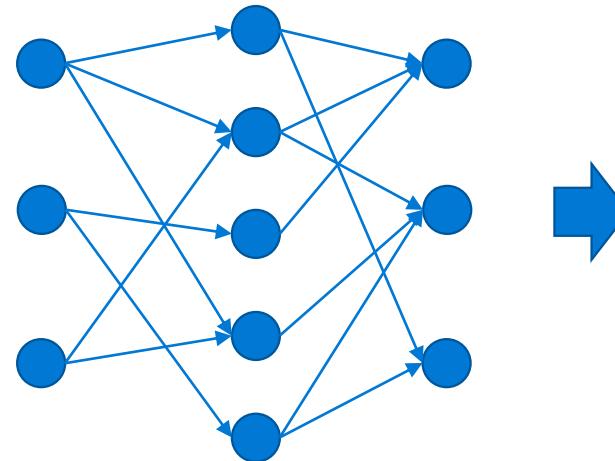
2. Prune Connections

$$\hat{w} \leftarrow abs(w)$$

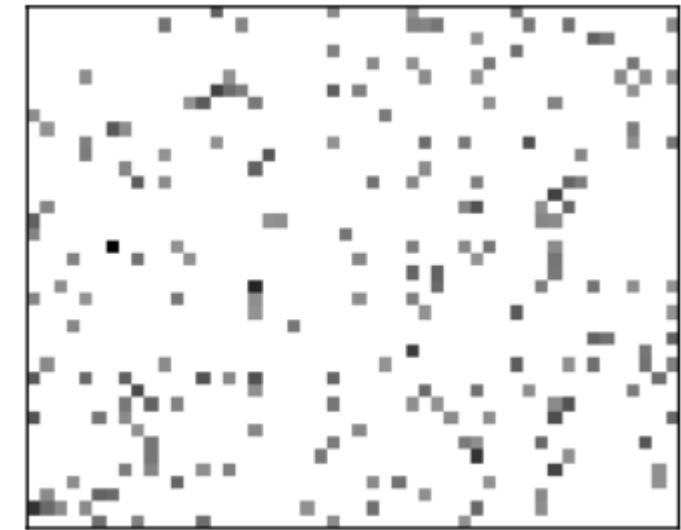
if  $\hat{w}_i \leq Threshold$

$$w_i \leftarrow 0$$

3. Train Weights

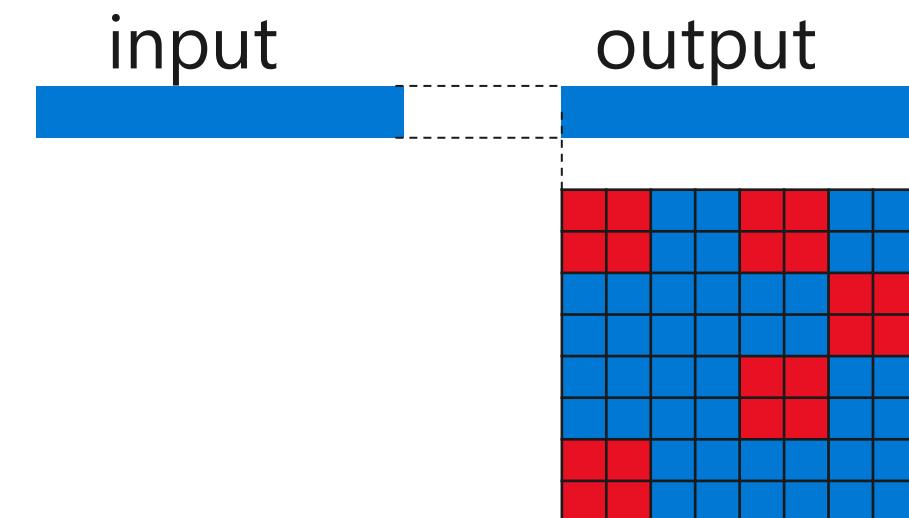
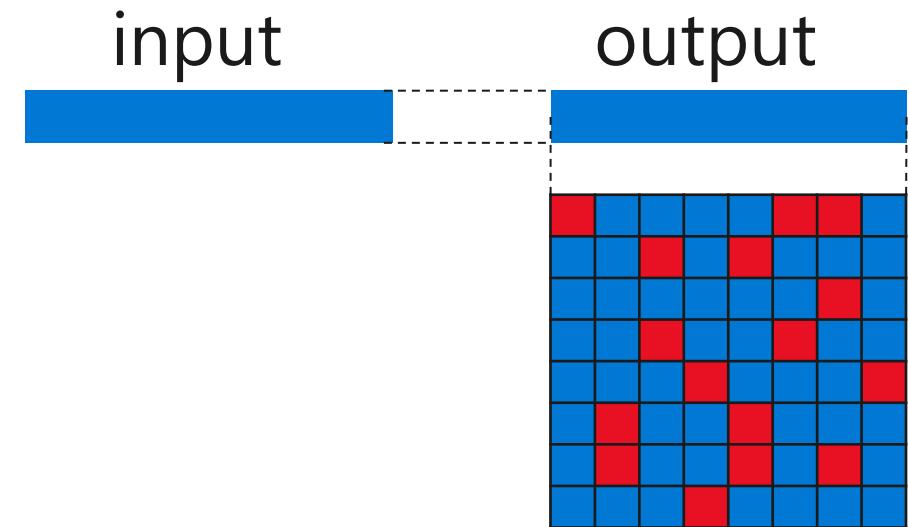
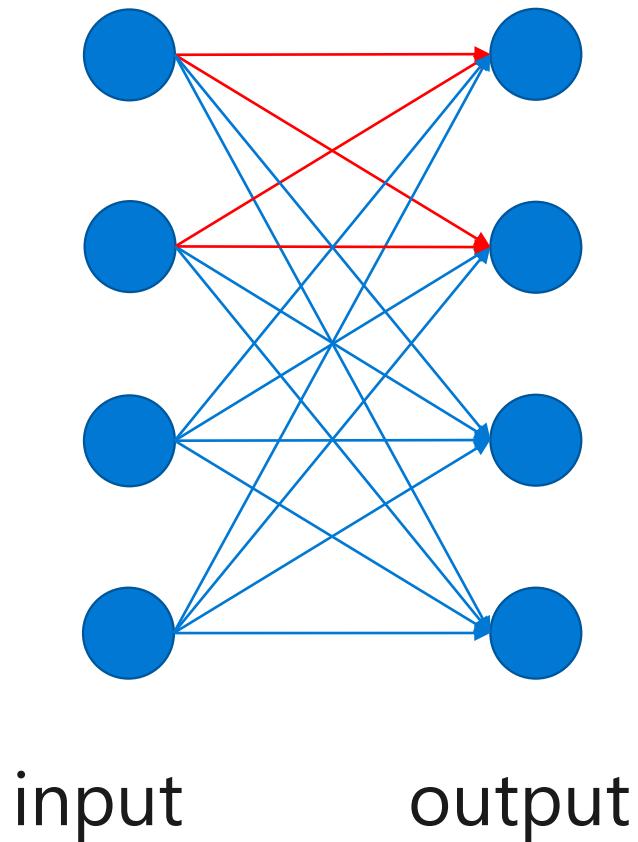


Sparsity = 50%~90%

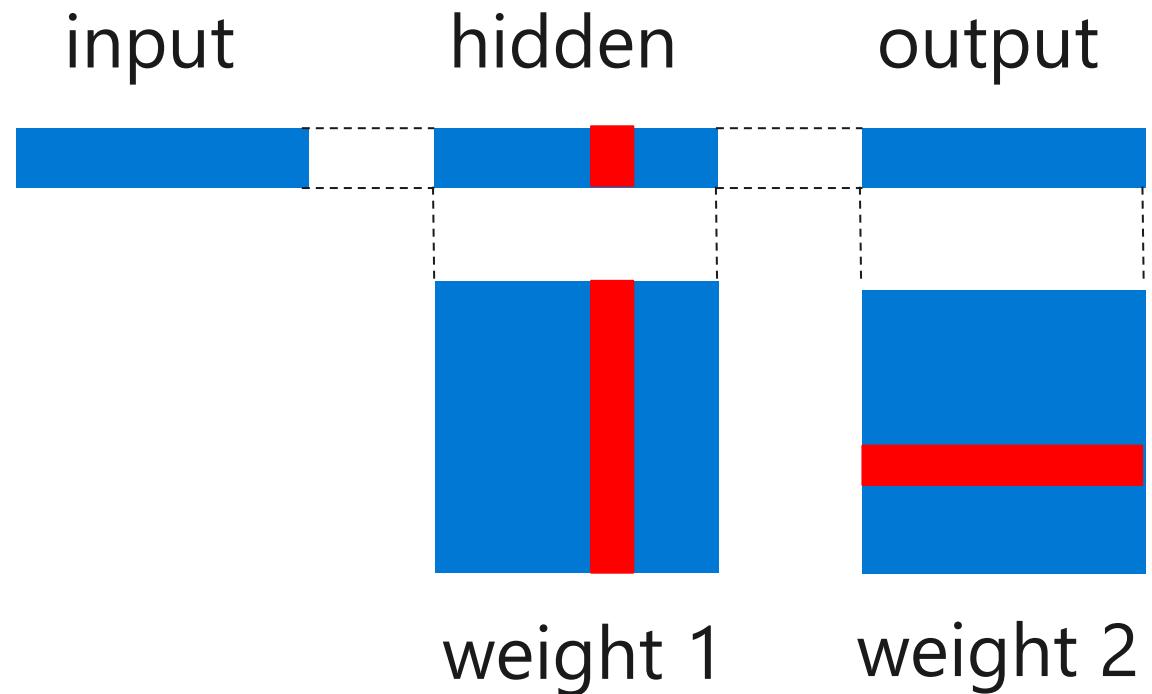
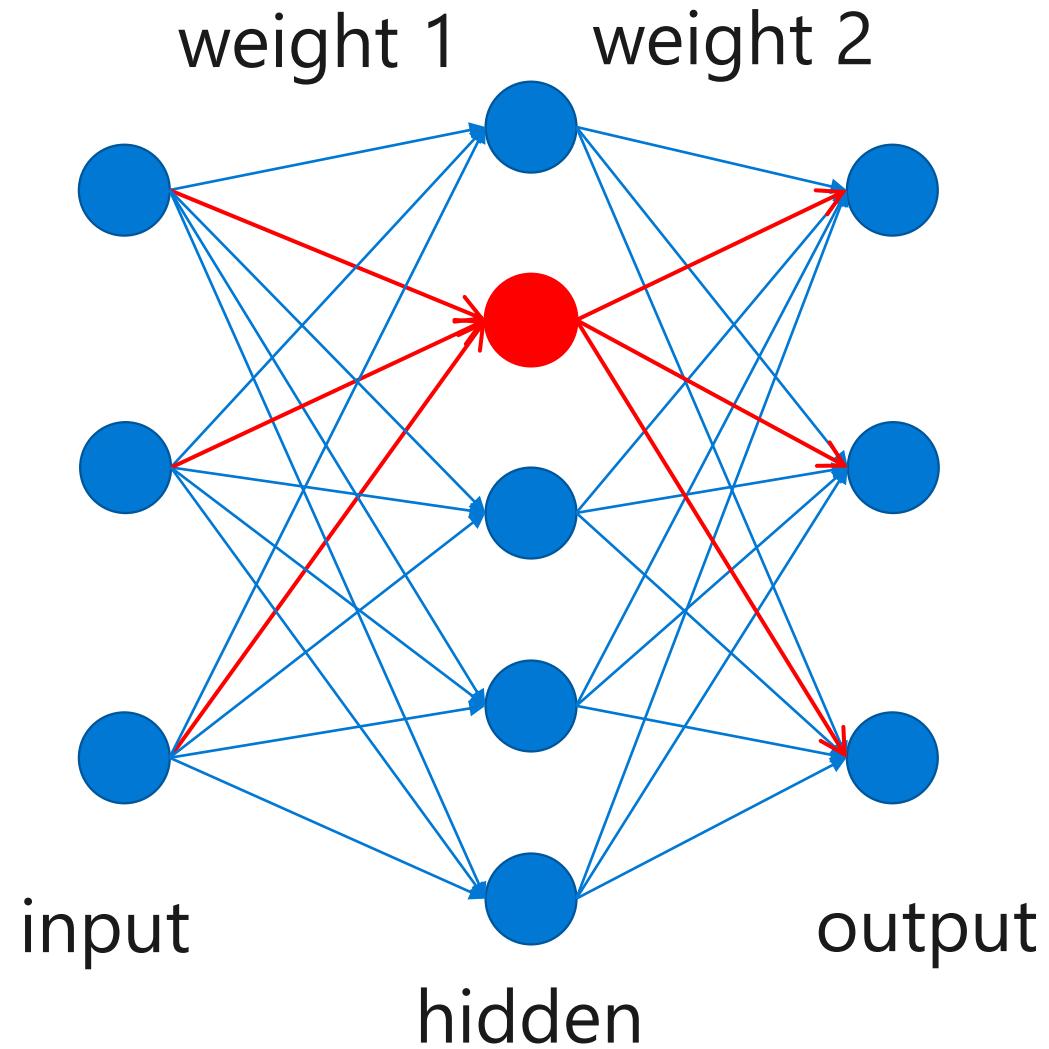


高度非规则的矩阵形状

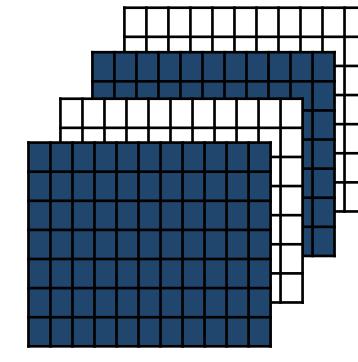
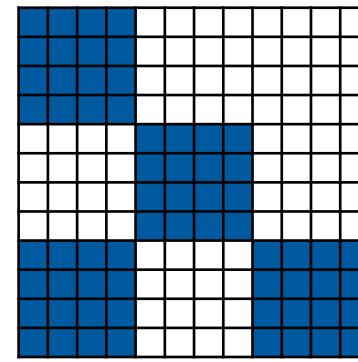
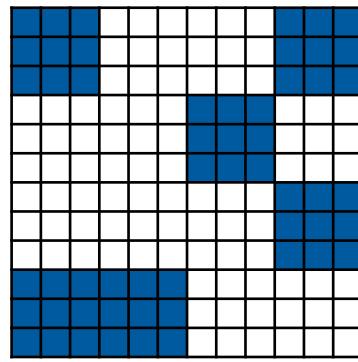
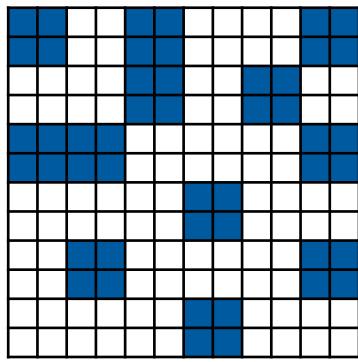
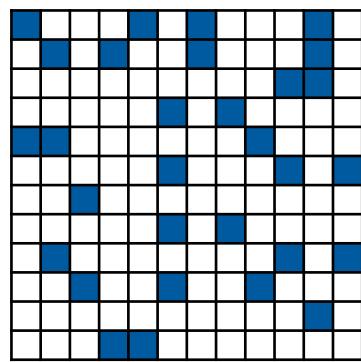
# 细粒度稀疏与粗粒度稀疏



# 通道剪枝 ( Channel Pruning )



# 加速与精度的权衡



细粒度

精度高

不规则

难以加速

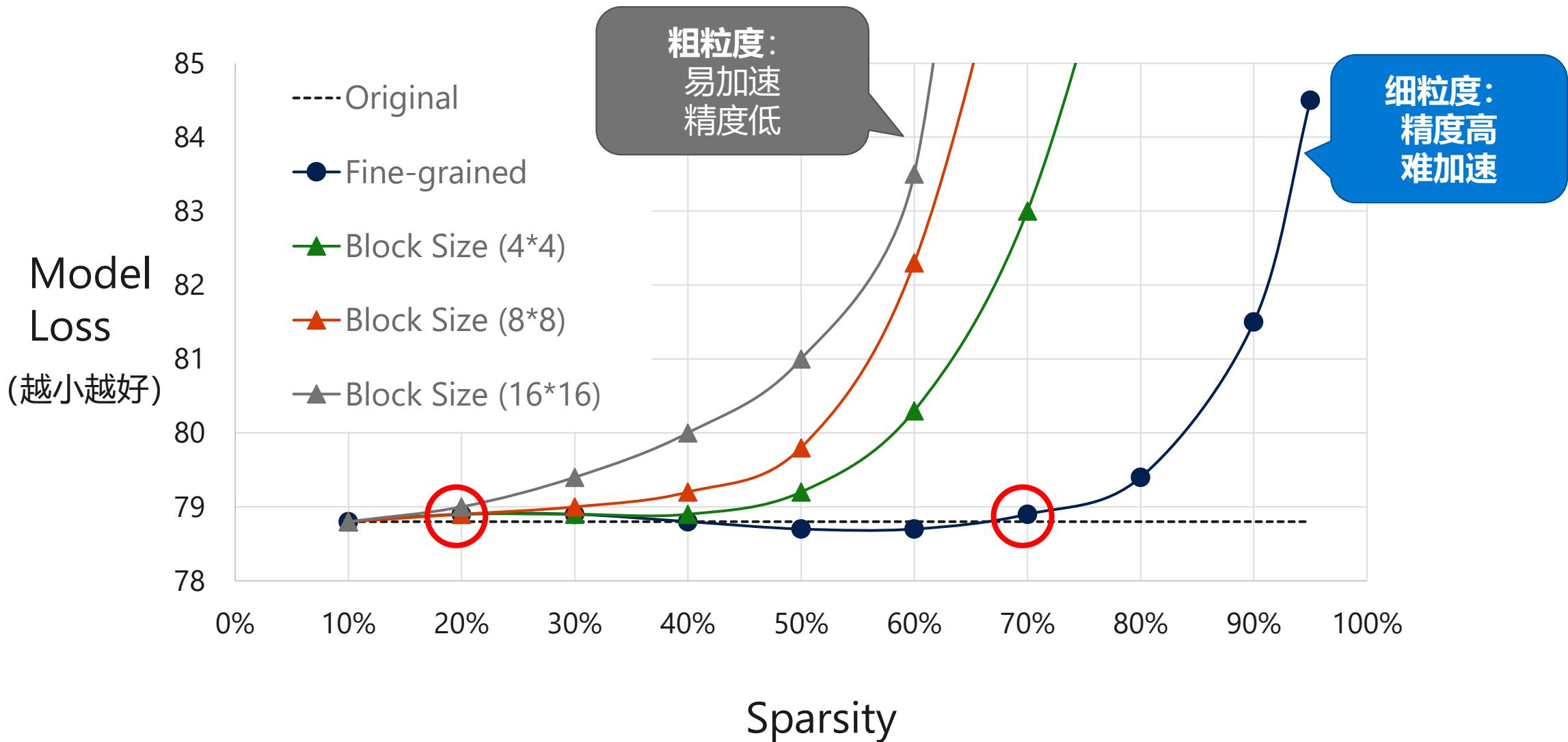
粗粒度

精度低

规则

容易加速

# 加速与精度的权衡



# FPGA'19 和 AAAI'19

0.2	0.1	0.2	-0.6	0.1	0.4	-0.1	0.6
0.4	-0.3	0.4	0.1	0.2	-0.4	0.1	0.5
0.7	-0.1	-0.3	0.1	0.5	-0.1	0.5	0.1
-0.1	0.6	-0.5	0.3	-0.4	-0.2	0.3	0.6

(a) Original Dense matrix

			-0.6		0.4		0.6
0.4		0.4			-0.4		0.5
0.7				0.5		0.5	
	0.6	-0.5	0.3	-0.4		0.3	0.6

(b) Unstructured sparse matrix  
by global pruning

		0.2	-0.6			-0.1	0.6
		0.4	0.1			0.1	0.5
0.7	-0.1				0.5	0.1	
-0.1	0.6				0.3	0.6	

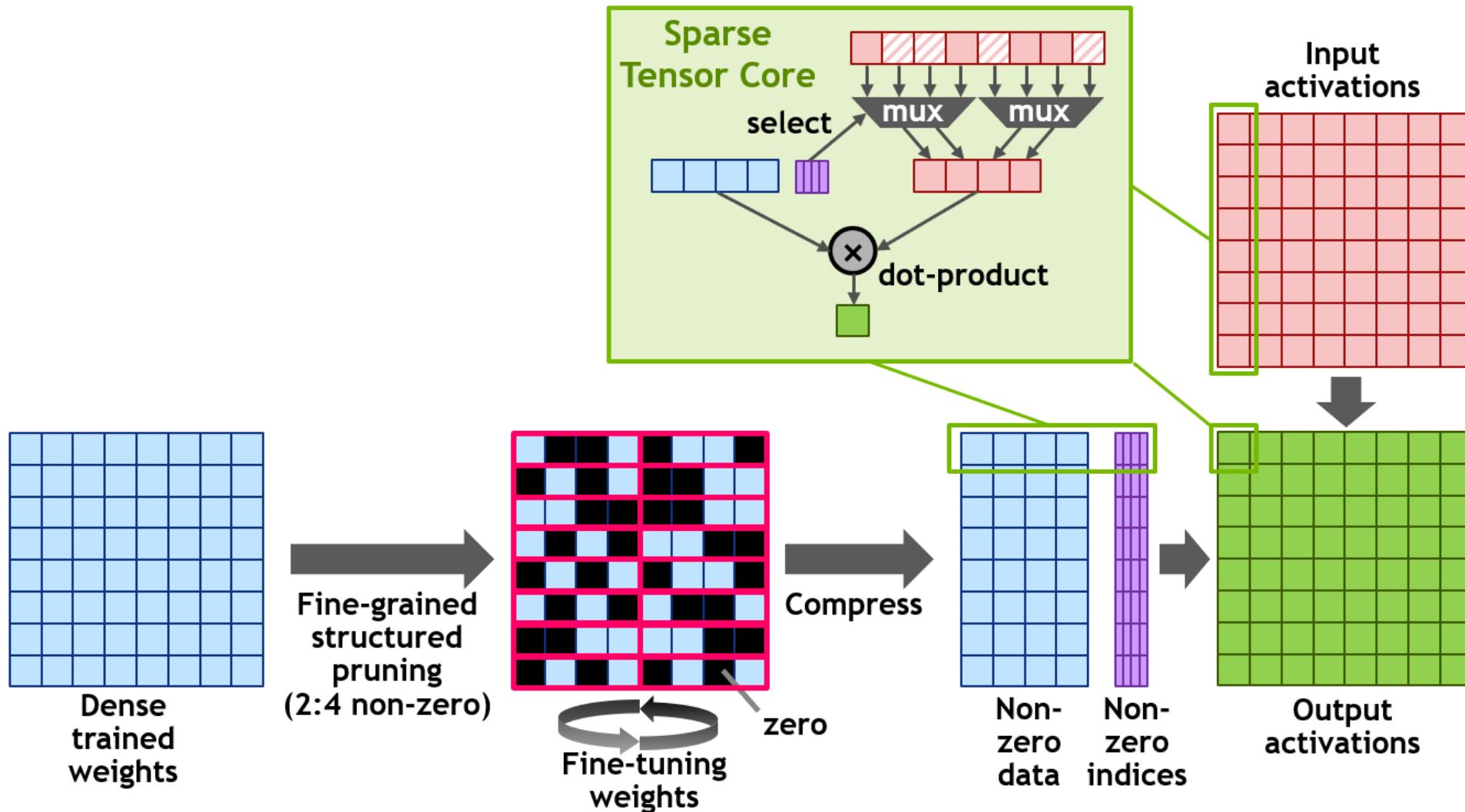
(c) Block sparse matrix by pruning 2x2  
blocks according to block average.

		0.2	-0.6			0.4		0.6
0.4		0.4				-0.4		0.5
0.7		-0.3			0.5		0.5	
	0.6	-0.5			-0.4			0.6

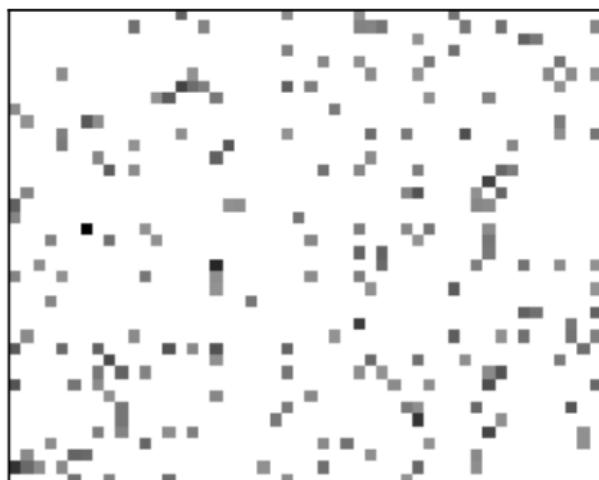
(d) Bank-balanced sparse matrix by  
local pruning inside each 1x4 bank

Efficient and effective sparse LSTM on FPGA with bank-balanced sparsity  
Balanced sparsity for efficient DNN inference on GPU

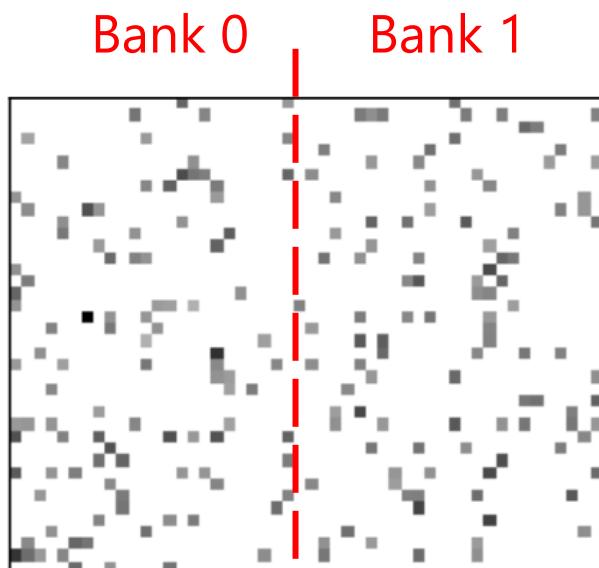
# 英伟达 A100 GPU: sparse tensor core



# 剪枝效果可视化



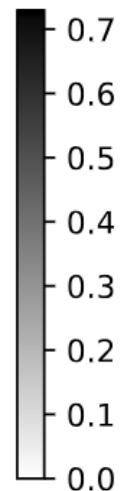
(a) 细粒度剪枝



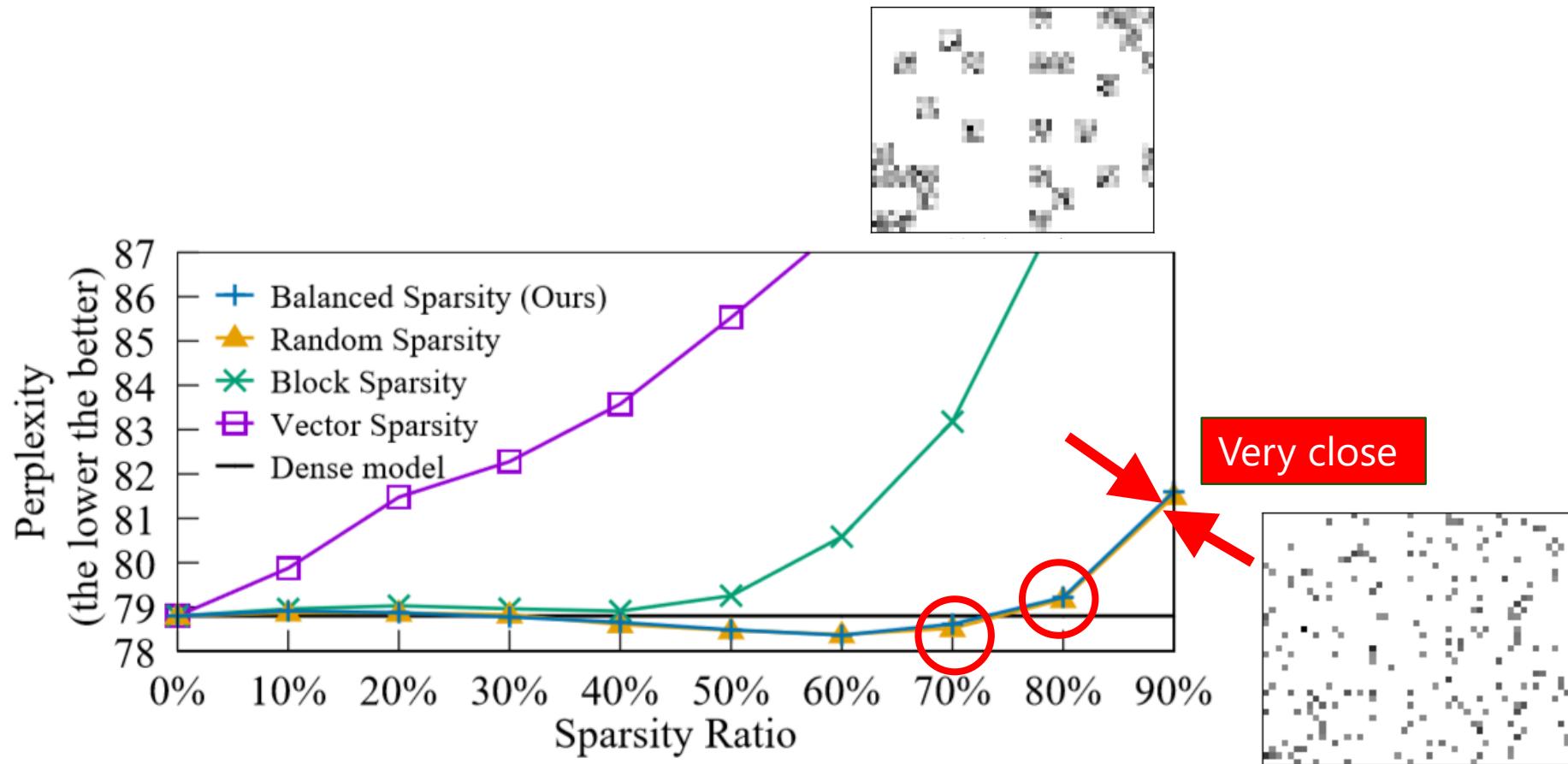
(b) BBS (我们的方法)



(c) 粗粒度剪枝



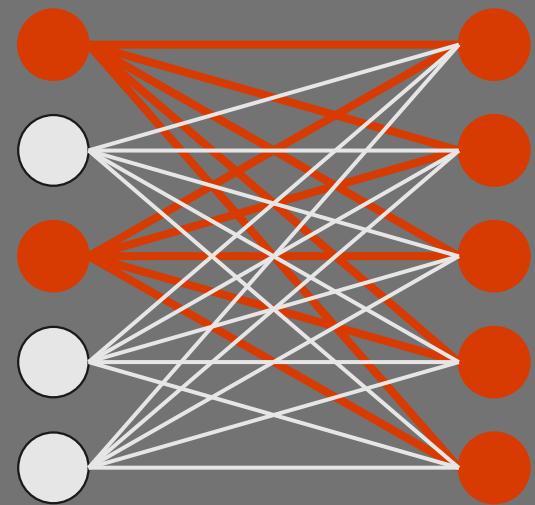
# 模型准确度对比



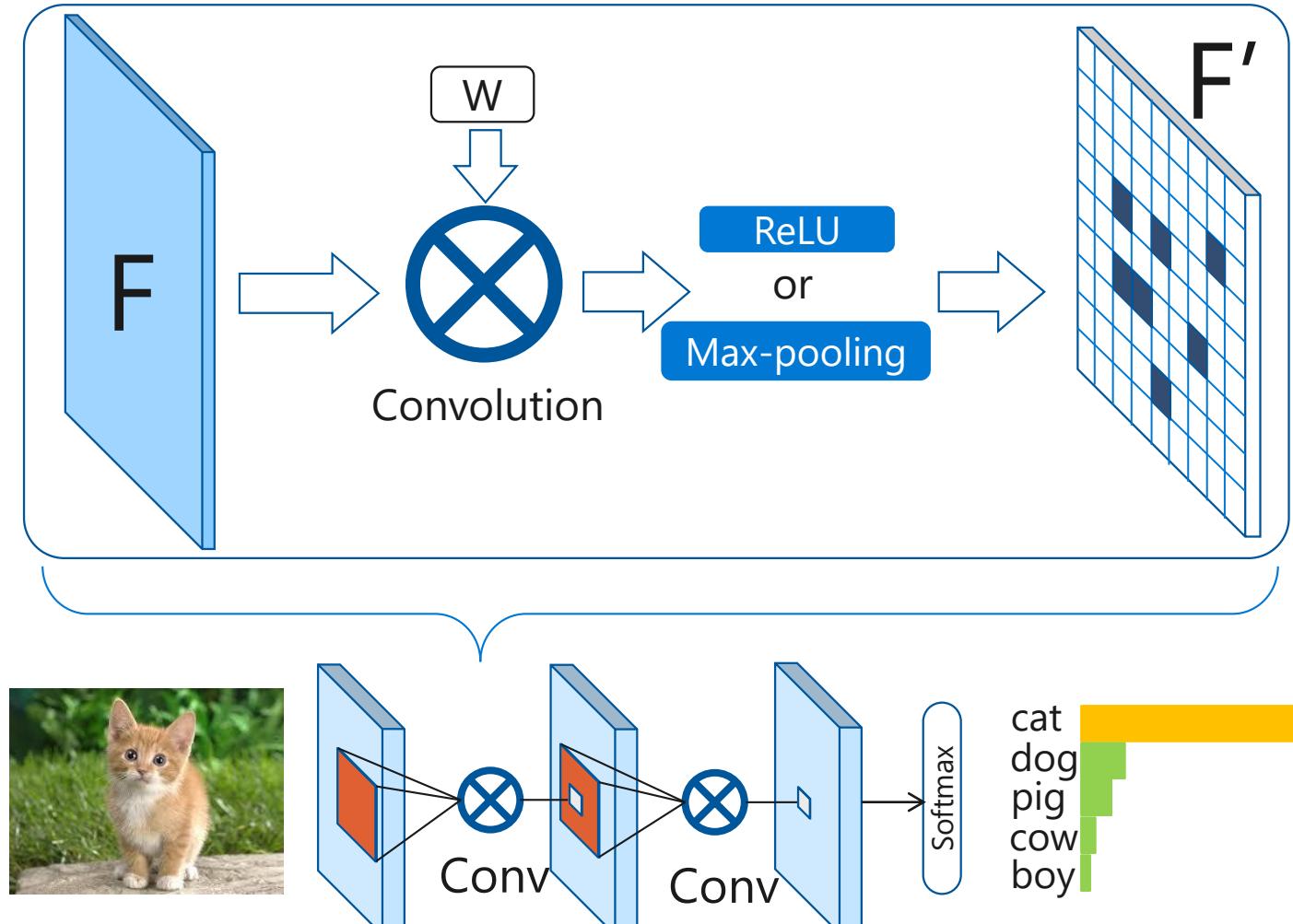
# 思考题

- 卷积神经网络上，可以有多少种不同粒度的剪枝？
- 剪枝和架构搜索（architecture search）的关系是什么？

# 激活稀疏

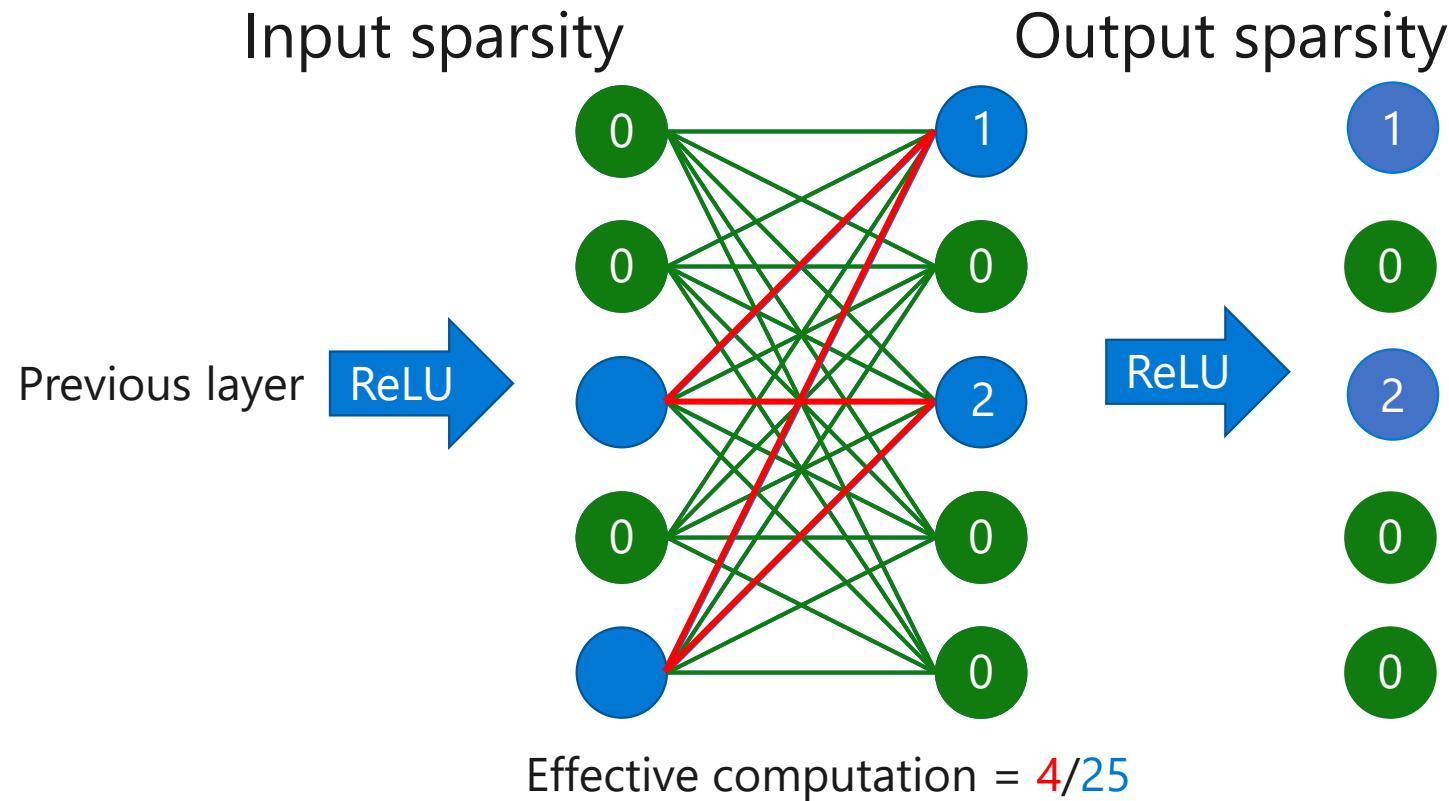


# 神经网络的激活 (activation) 稀疏



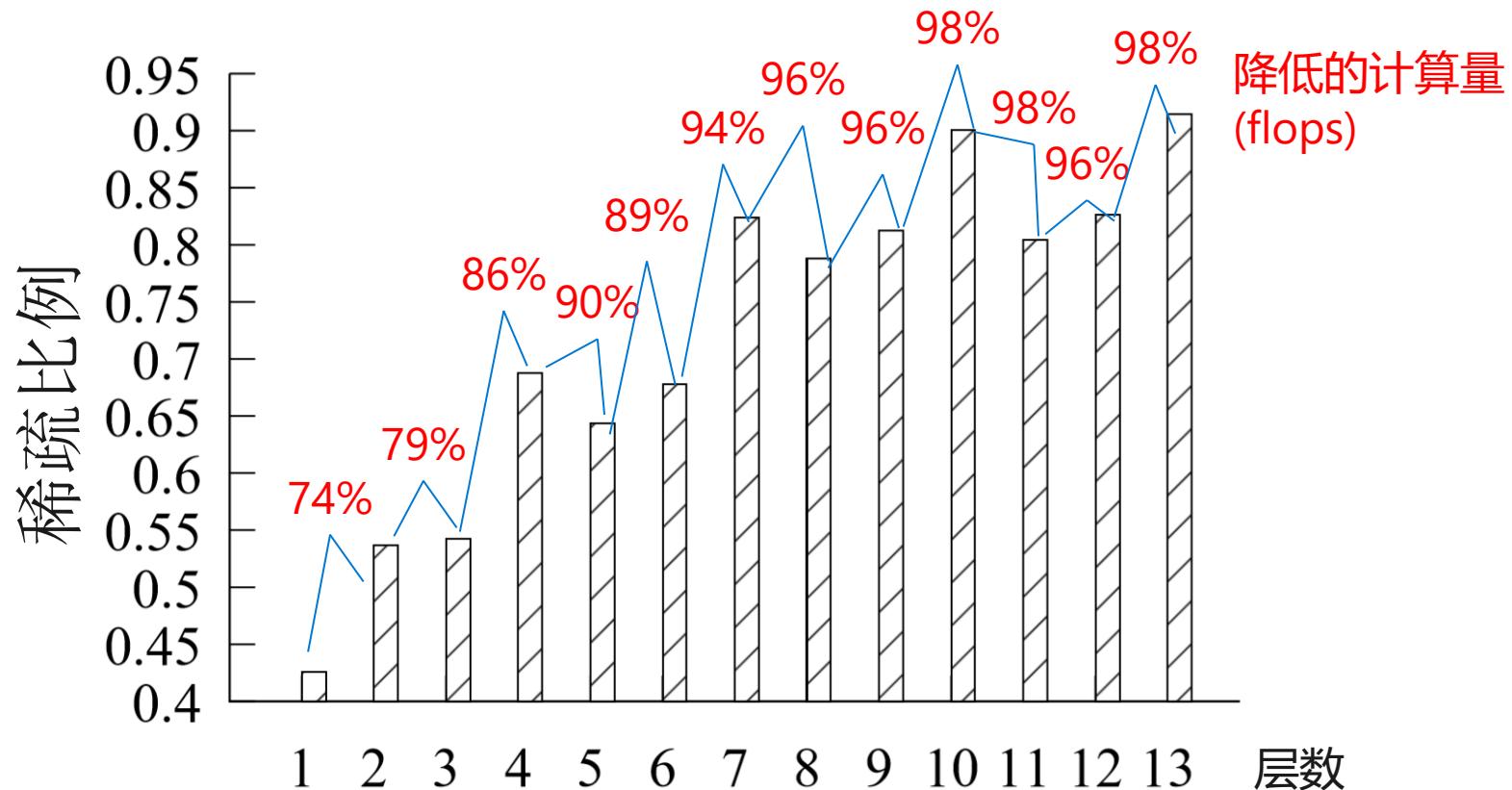
- ReLU
  - $y = \max(0, x)$
- Max-pooling
  - $y = \max(x_i \mid i=\{1,2,\dots,n\})$

# 激活稀疏的加速机会

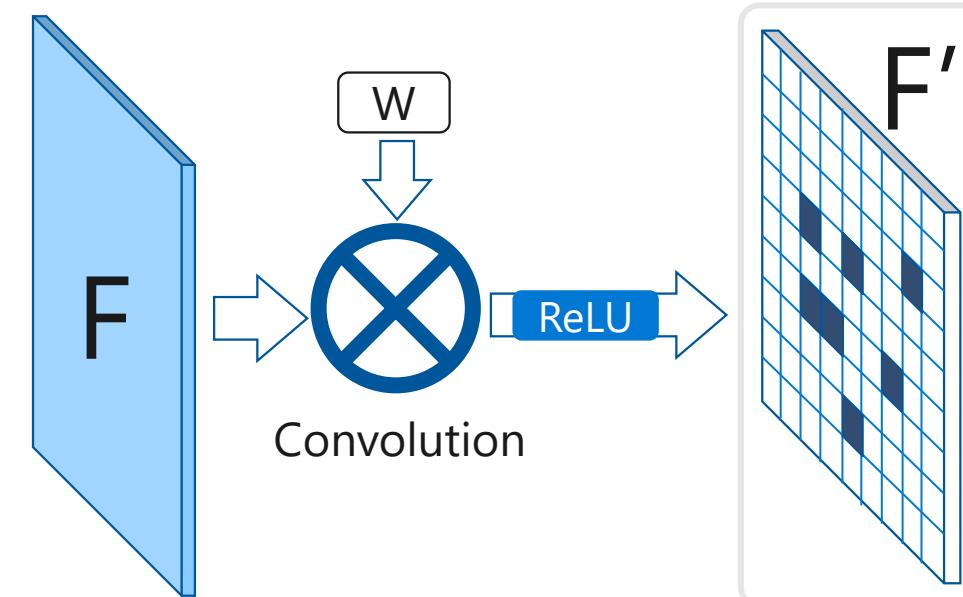


# Resnet-16的激活稀疏

- 稀疏度 : 40% ~ 90%
- 稀疏性同时存在于输入和输出特征图



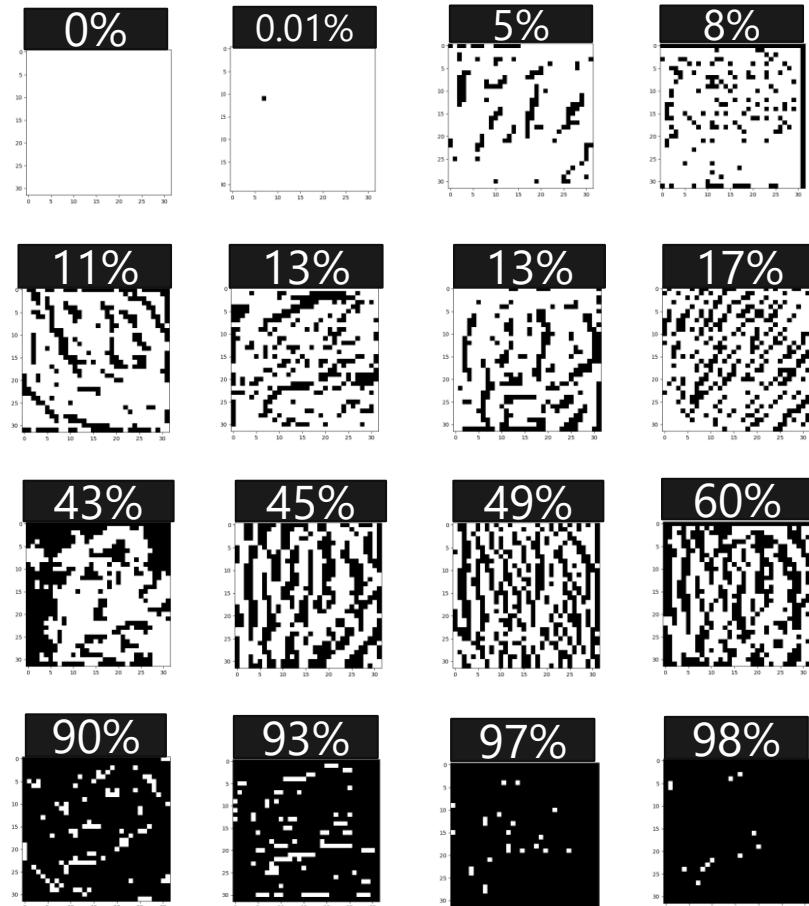
# 真实的例子： ResNet-16 layer 2



输入特征图

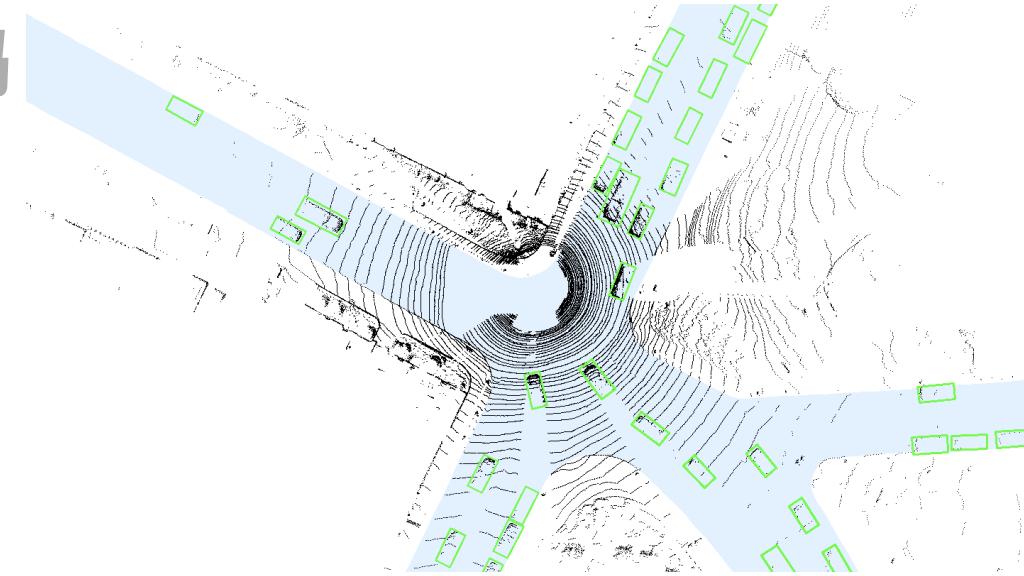
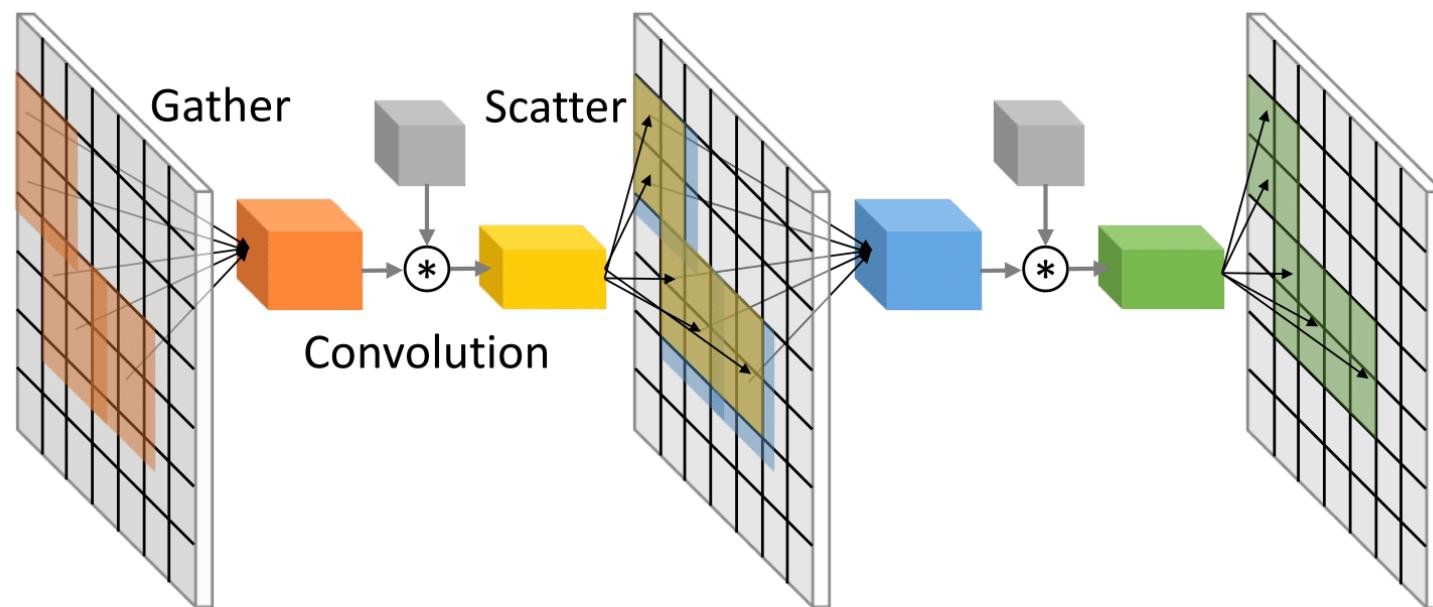
输出特征图

特征图上的激活稀疏度不是均匀分布的



# 研究课题 4：结构化特征图稀疏

SBnet: 借助外部信息给特征图施加方形掩码

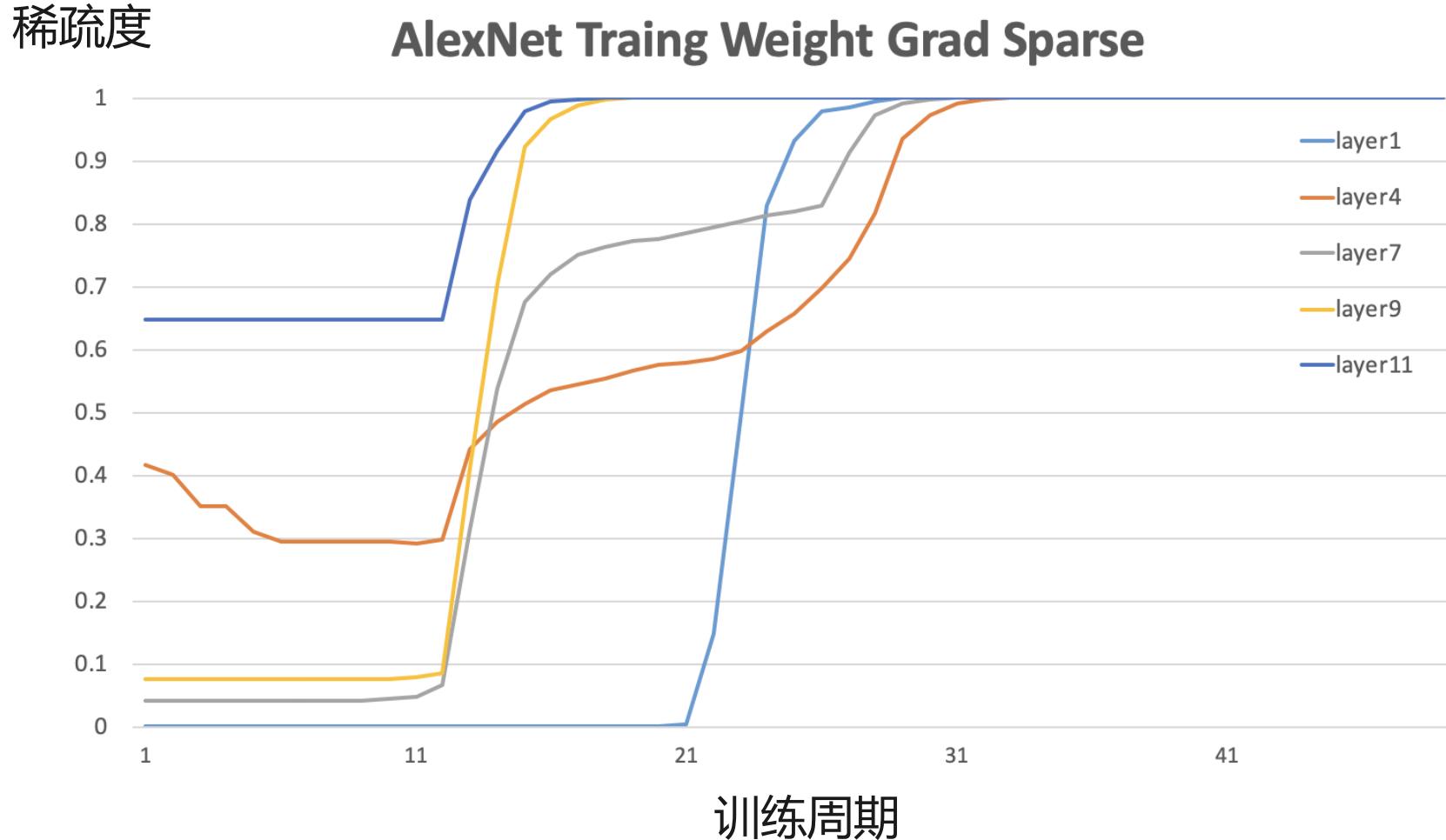


## 思考题：

- 从数据依赖角度看，权重稀疏与激活稀疏的关系是什么？
- 能否把权重稀疏和激活稀疏结合考虑取得更大的加速比？

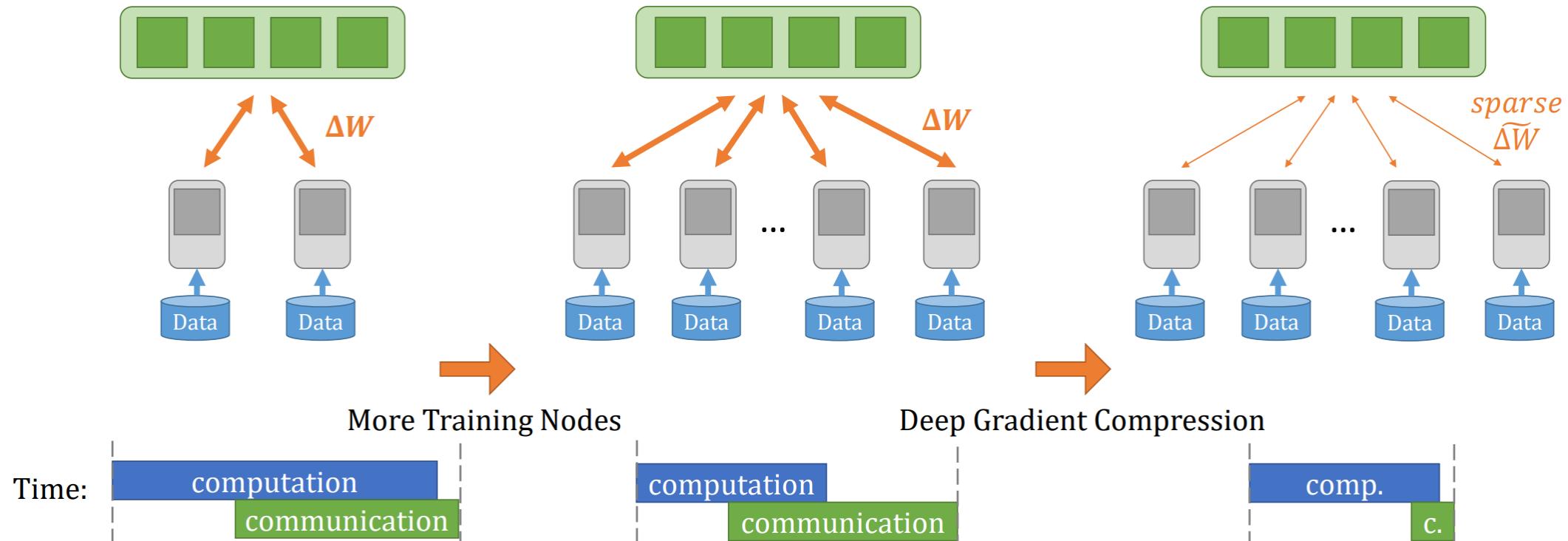
梯度稀疏

# 梯度稀疏：AlexNet



# 研究课题 5：加速训练

通过梯度稀疏降低分布式训练的通信代价



[] Deep Gradient Compression: Reducing the Communication Bandwidth for Distributed Training

# 量化/低比特运算

# 整形与浮点

整形

符号位	数值位
1-bit	31-bit

最小值:  $-2^{31}$

最大值:  $2^{31} - 1$

浮点  
(单精度)

符号位	指数位	尾数位
1-bit	8-bit	23-bit

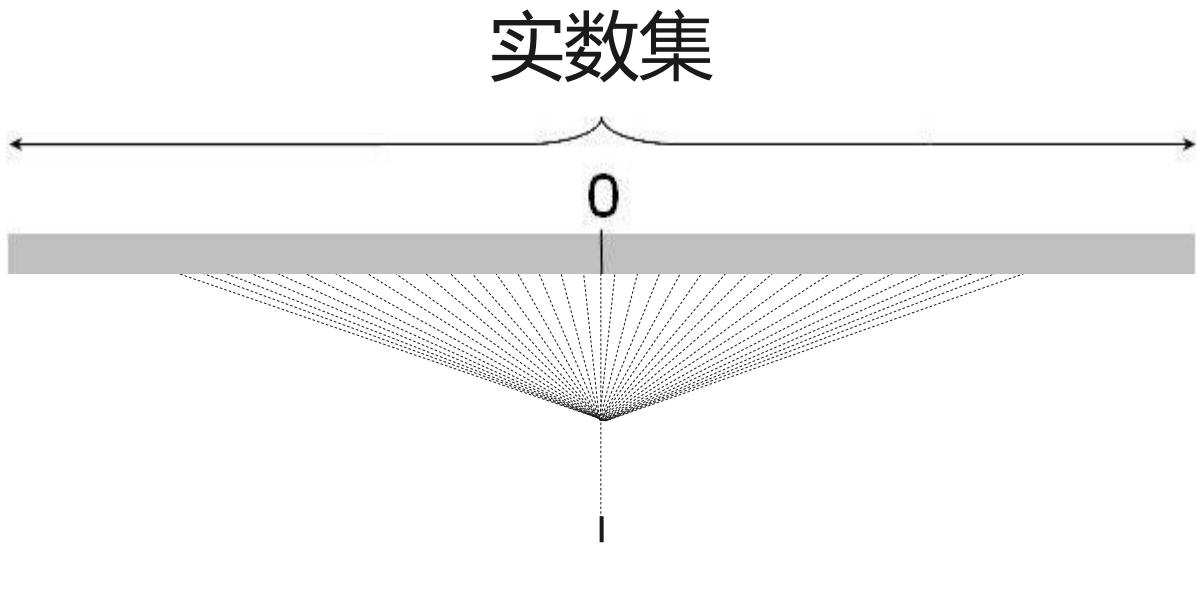
最小值:

$$-2^{-126} \times 1.0$$

最大值:

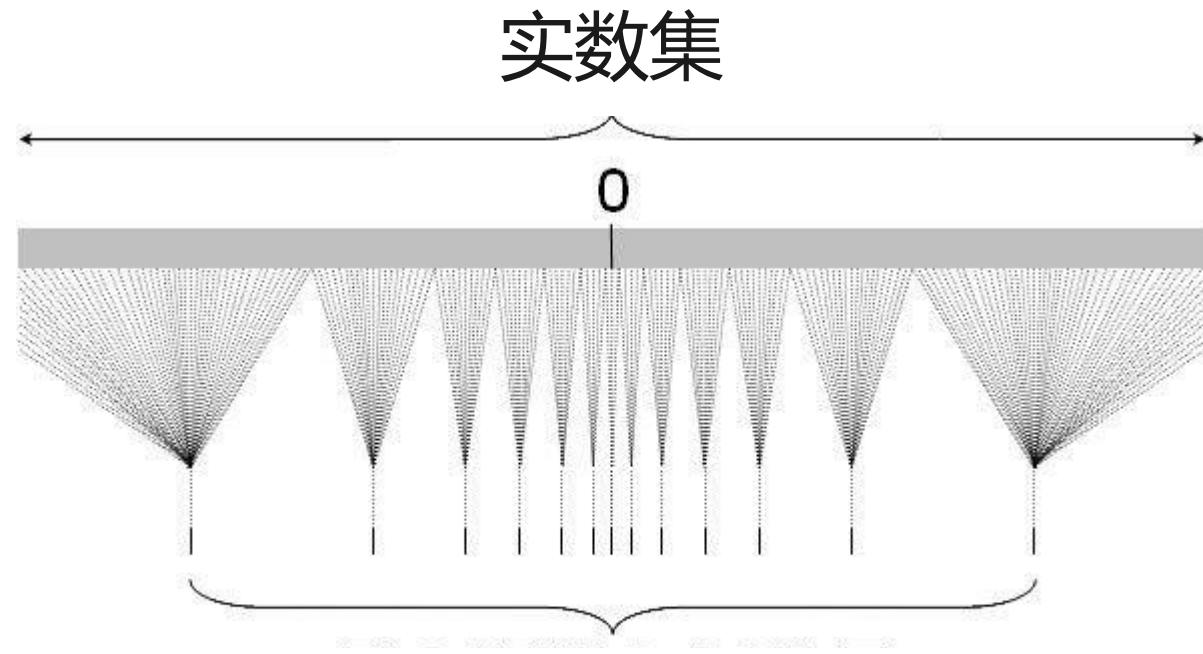
$$+2^{127} \times 1.11111111111111111111111$$

# 整形：线性采样



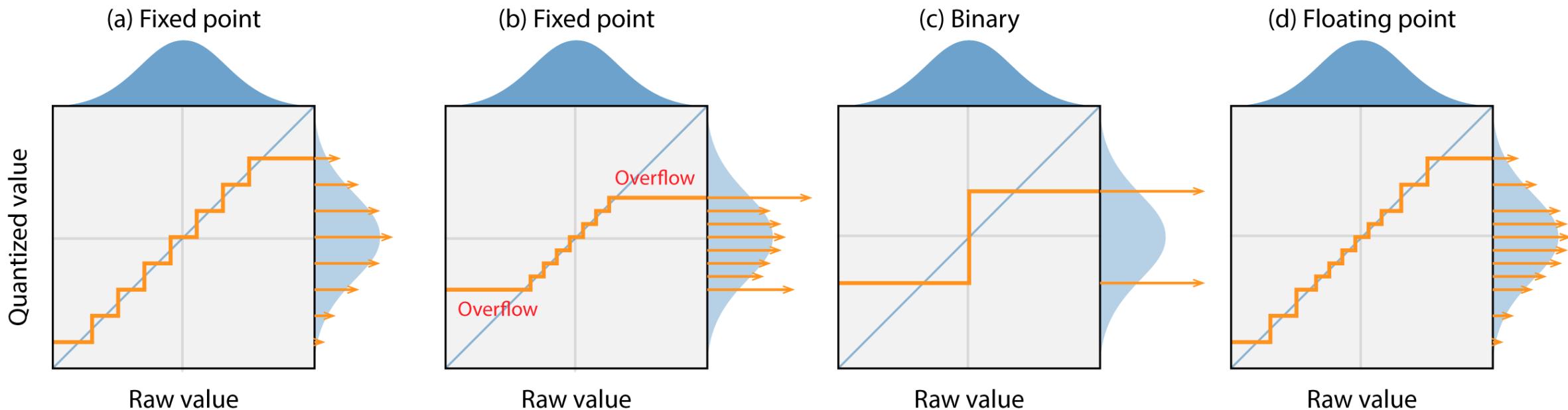
整数表达

# 浮点：指数分段采样

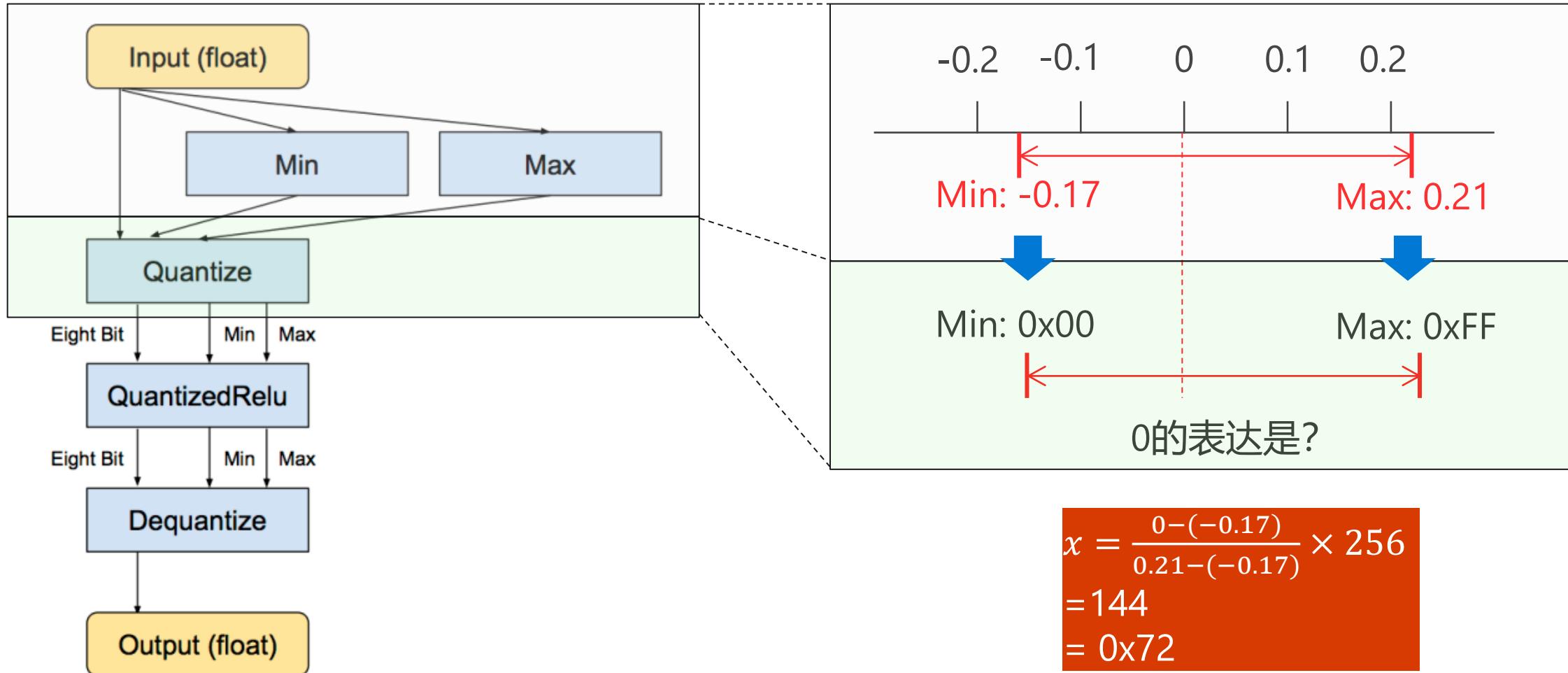


浮点数表达

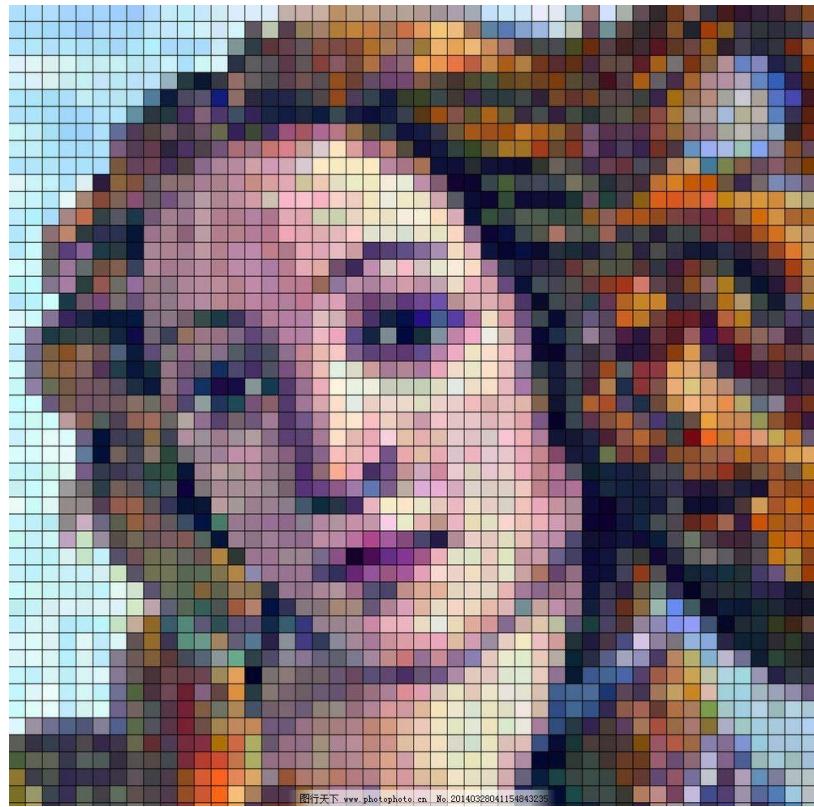
# 不同数据格式：不同采样精度



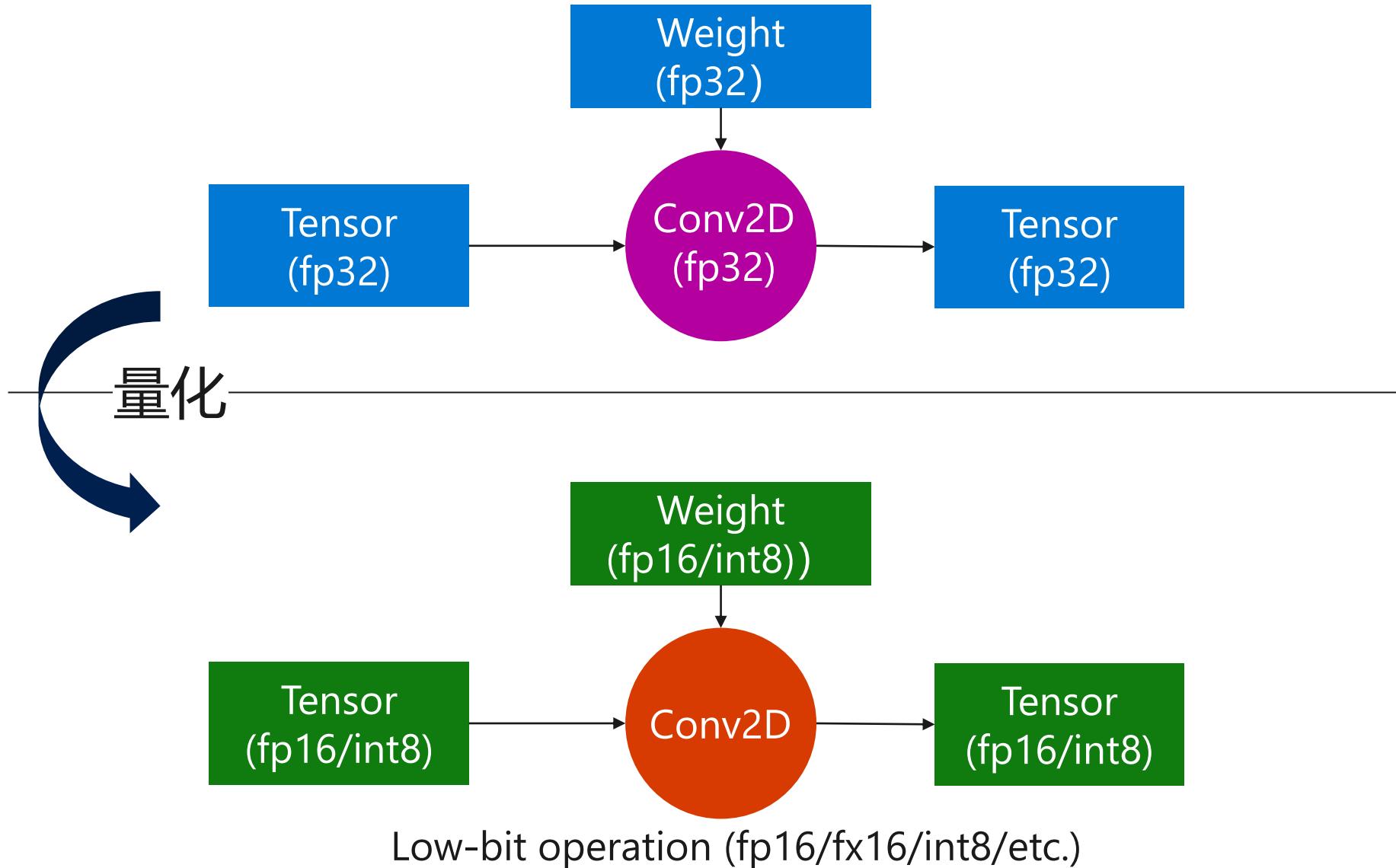
# Tensorflow中的定点化方法



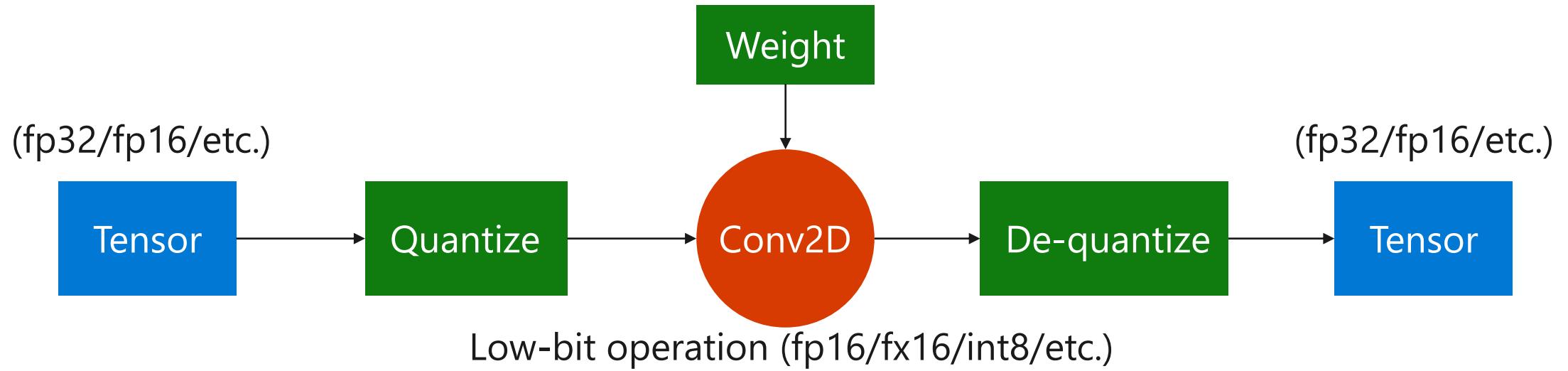
# 量化 (quantization) : 低精度运算



# 量化 (quantization) : 低精度运算



# 激活量化与权重量化



# 研究课题 6：量化位宽与精度格式 三值网络

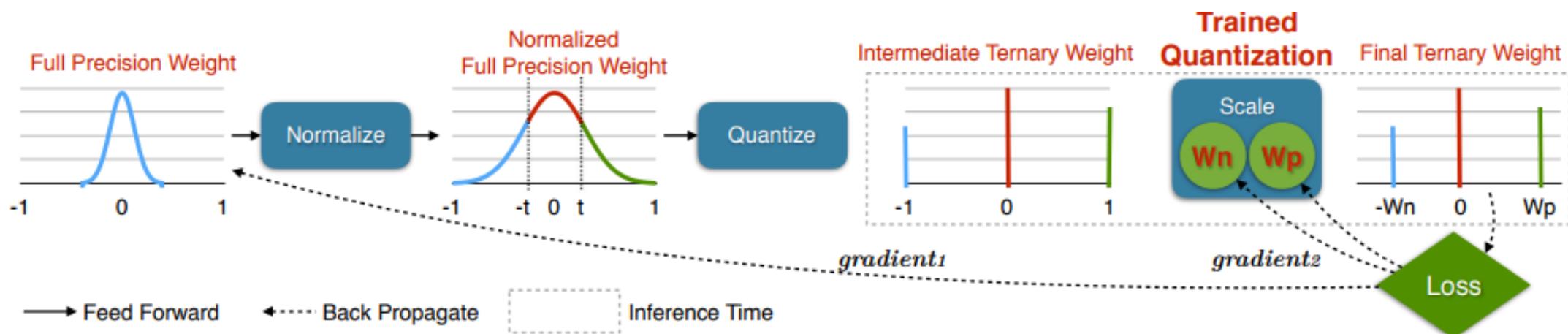
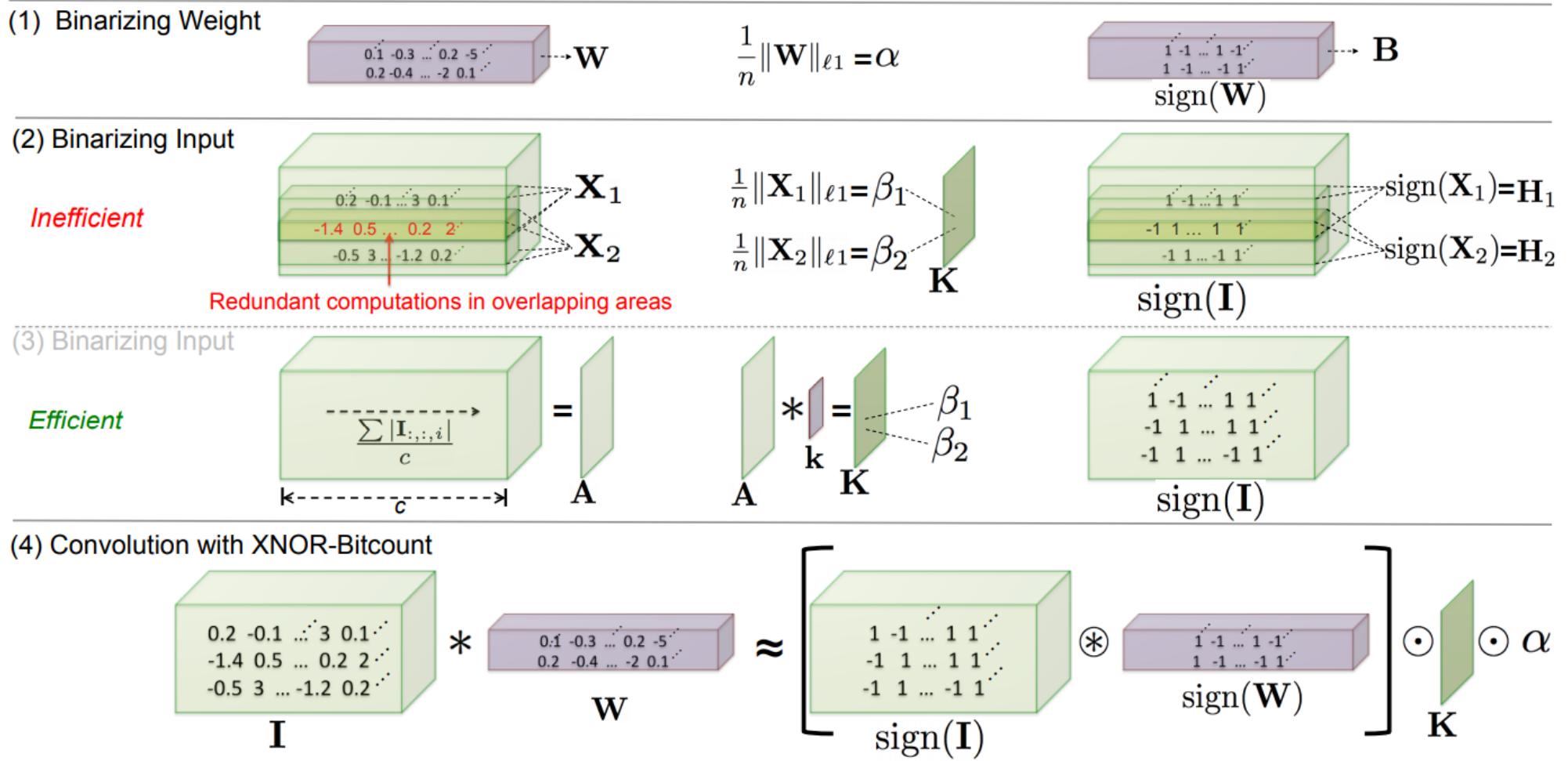


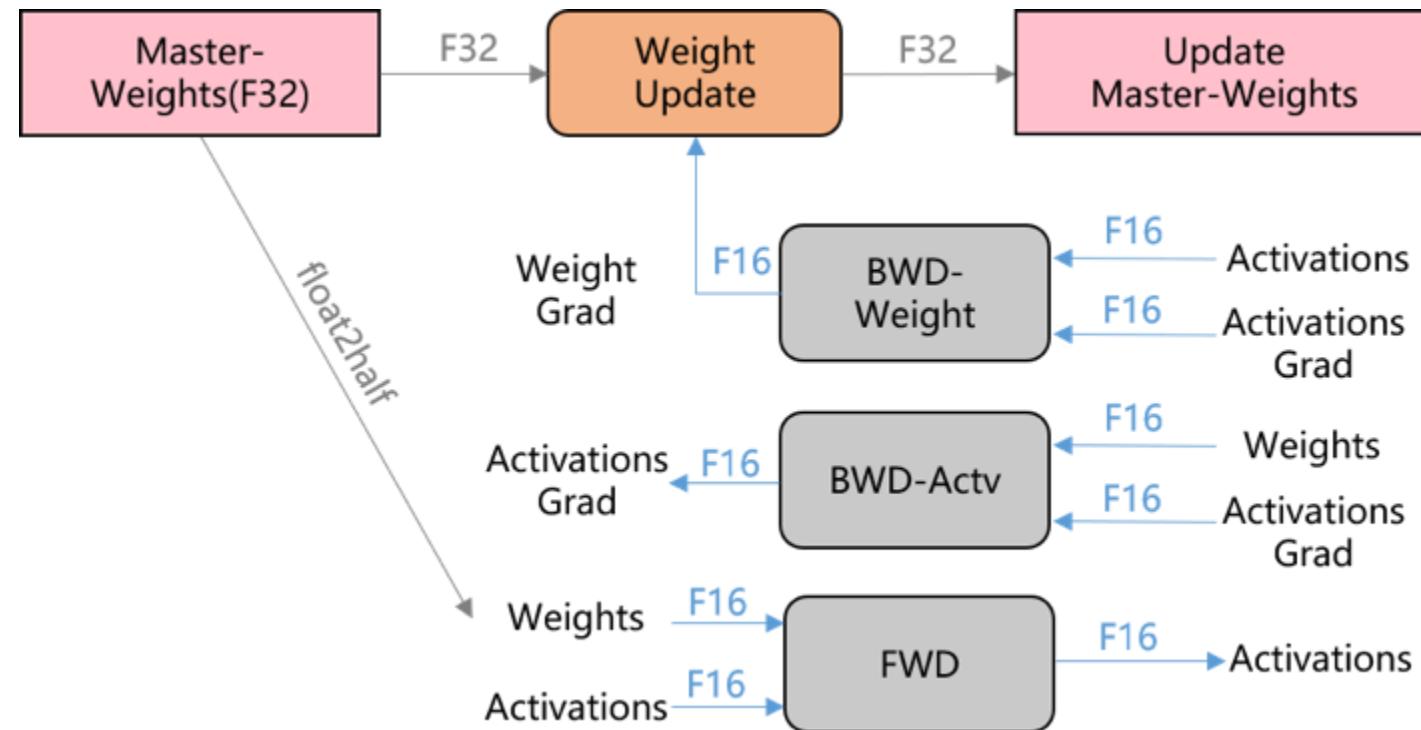
Figure 1: Overview of the trained ternary quantization procedure.

quantized ternary weights:  $w_l^t = \begin{cases} W_l^p : \tilde{w}_l > \Delta_l \\ 0 : |\tilde{w}_l| \leq \Delta_l \\ -W_l^n : \tilde{w}_l < -\Delta_l \end{cases}$

# 二值网络



# 研究课题 7：混合精度训练



# 思考题：

- Tensorflow中的 8bit 定点化数据如何做乘法运算的呢？

## 实数运算

$$\text{向量 A: } [-0.1, 0.1, 0.2] \times \text{向量 B: } [-2, 1, 3] = \text{向量 C: } [0.2, 0.1, 0.6]$$

## Tensorflow 8-bit 运算

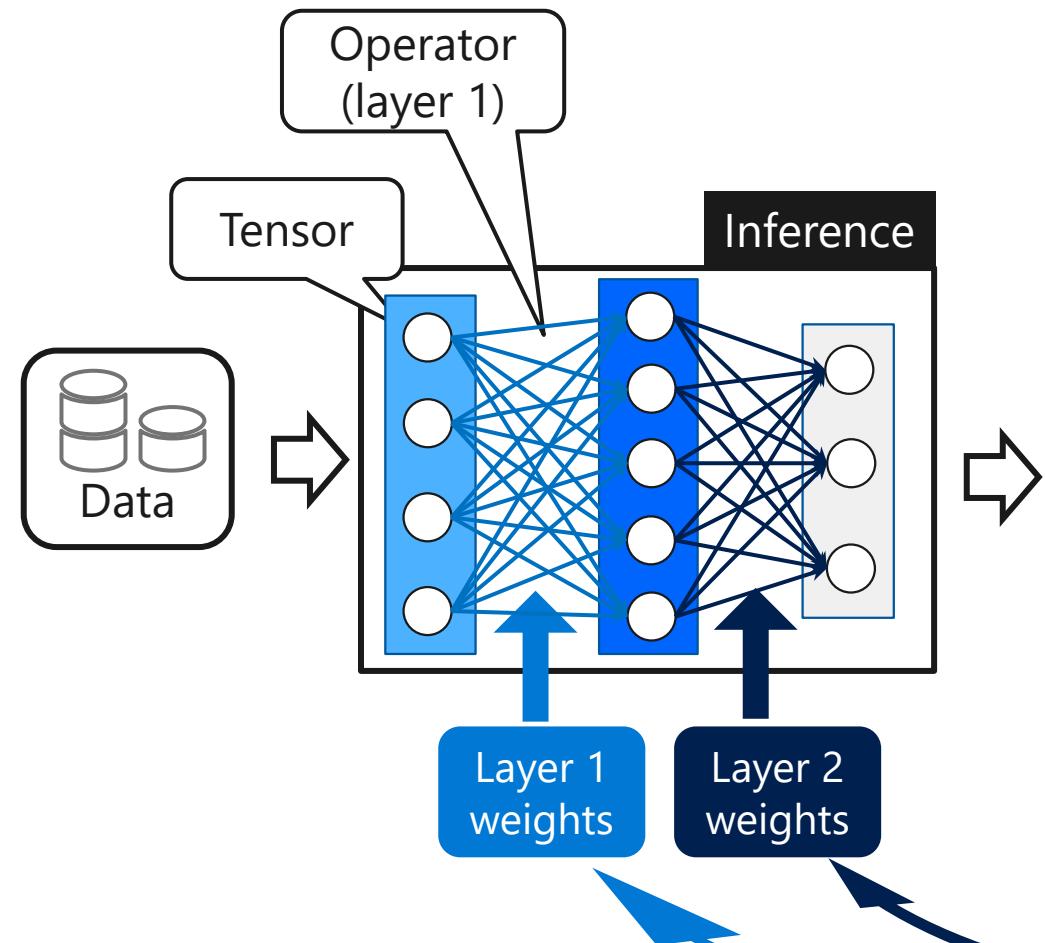
$$\text{向量 A: } [?, ?, ?] \times \text{向量 B: } [?, ?, ?] = \text{向量 C: } [?, ?, ?]$$

# 课程总结

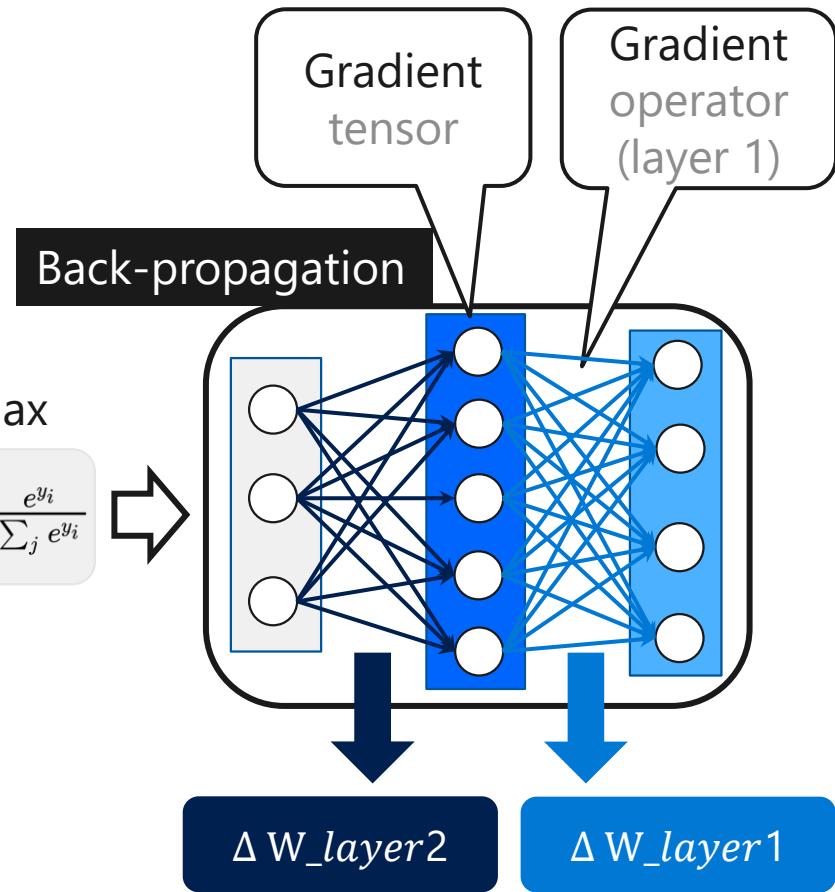
- 模型压缩和稀疏化是提升DNN模型运行效率的重要手段
- 稀疏化：权重稀疏，激活稀疏，梯度稀疏
- 量化：低精度运算
- 还有许多问题等待进一步的研究
  - 模型压缩和稀疏化的最优解尚待研究
  - 定制化硬件设计以充分发挥稀疏性和量化的优势

# Backup

# 回顾：DNN训练步骤



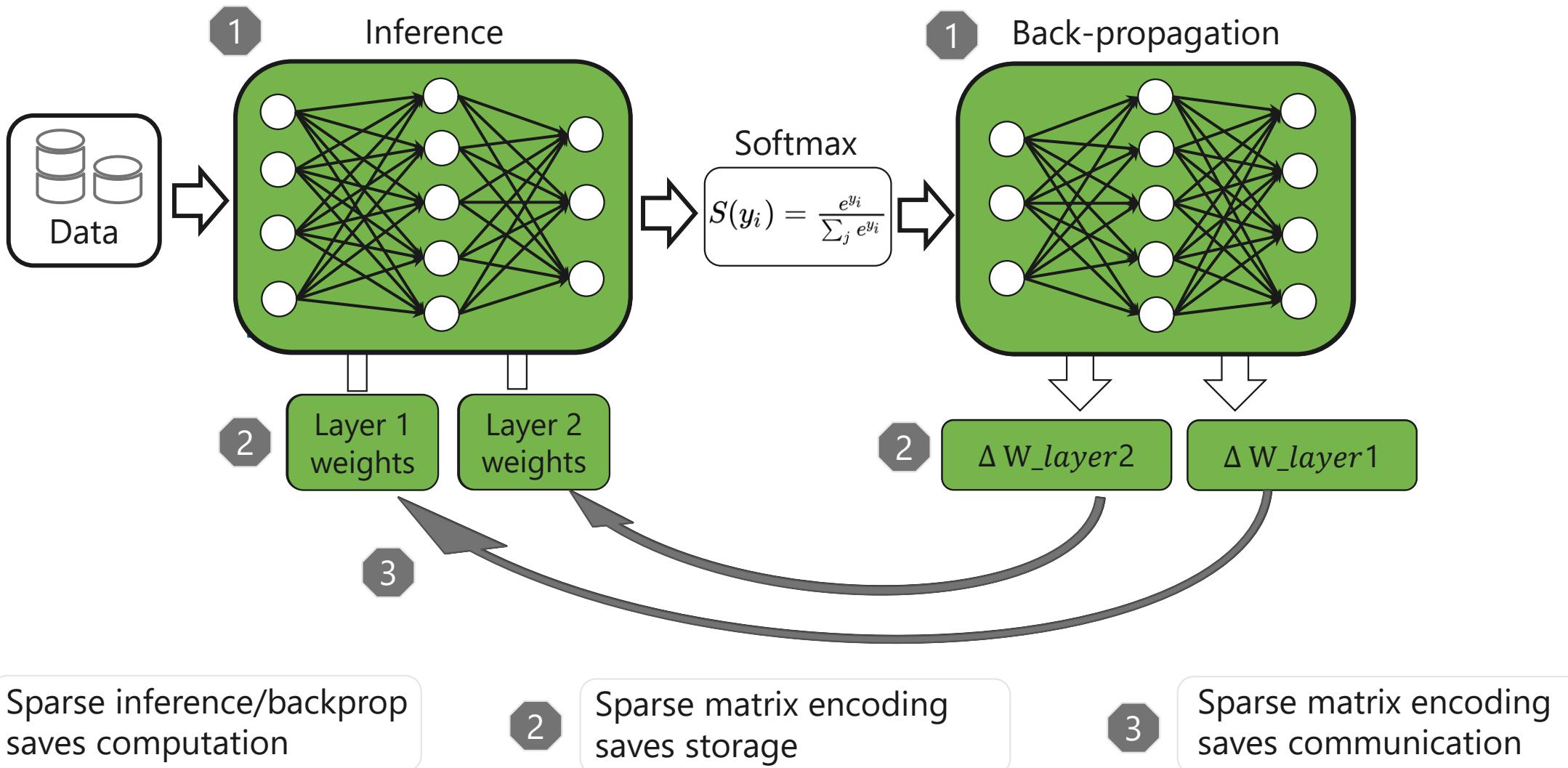
Auto-differentiate by framework (TF/Pytorch/etc.)



1 Iteration: a basic learning step

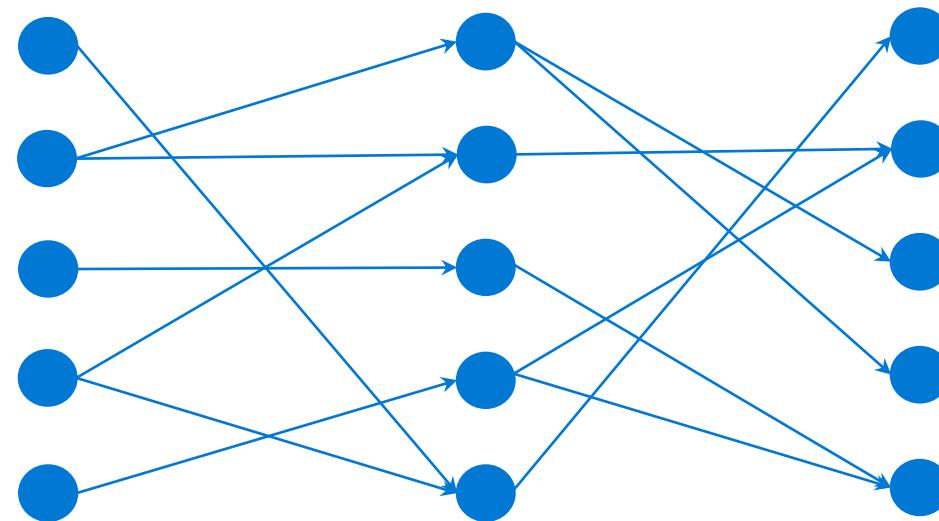
Training: 100~100000 Iterations

# 稀疏存在于神经网络训练的每一个阶段



# 前向计算：静态稀疏

neuron pruning is conducted offline



# 训练：动态稀疏

