

深度学习：零基础体验

PyTorch/Tensorflow的本质：大规模快速偏微分求导推理机器



TensorFlow

提纲

1. 求导方法

- 什么是导数,微积分
- 导数的传递(chain rule)
- 手算例子
- PyTorch/Tensorflow 计算图

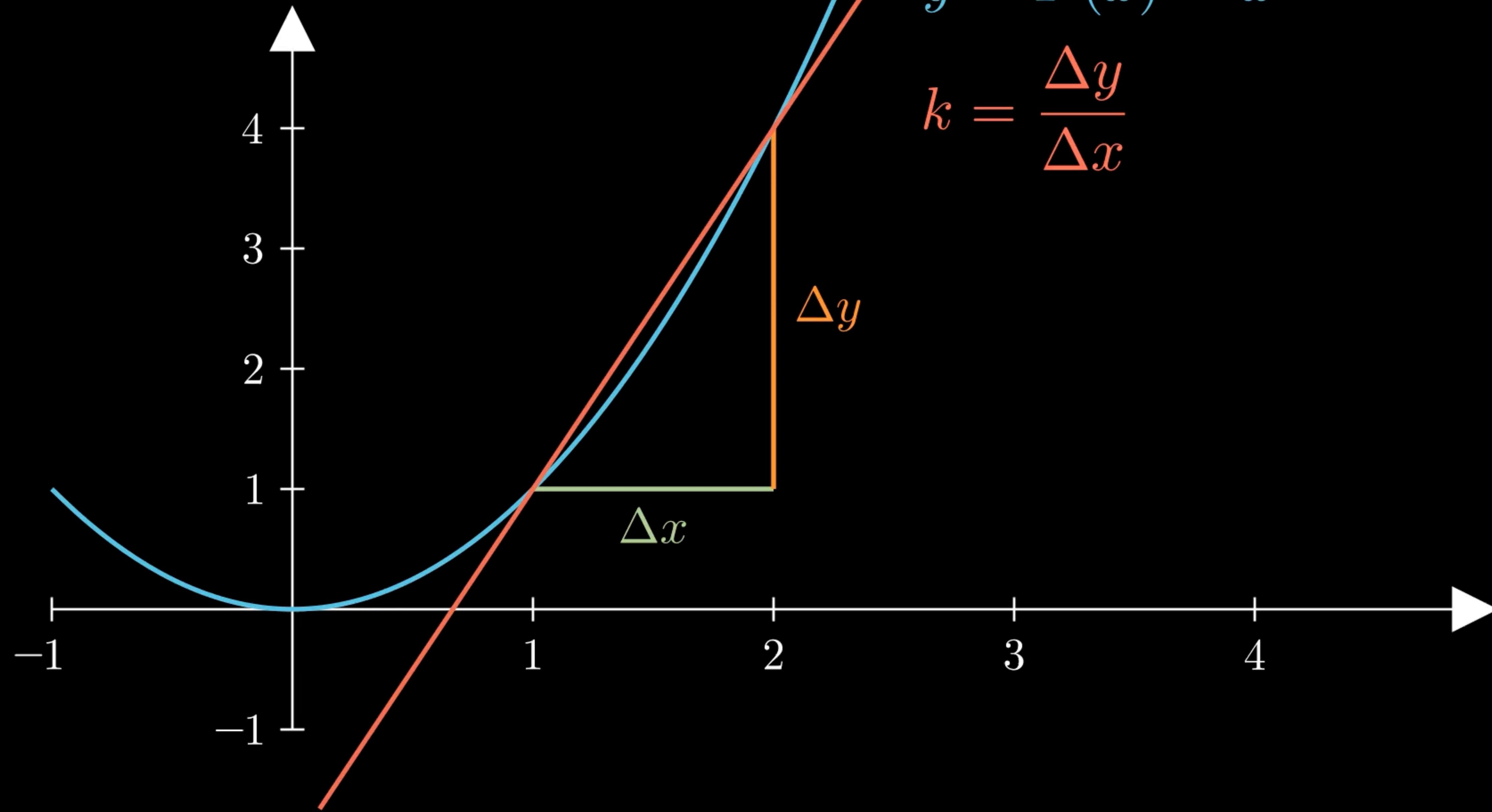
2. 深度学习的结构

- 全联接网络 (FCNN/ANN)
- 非线性变换
 - 单位变换/概率映射/损失函数
- 随机梯度下降 (SGD)
- 卷积神经网络CNN

3. 深度学习的目的：

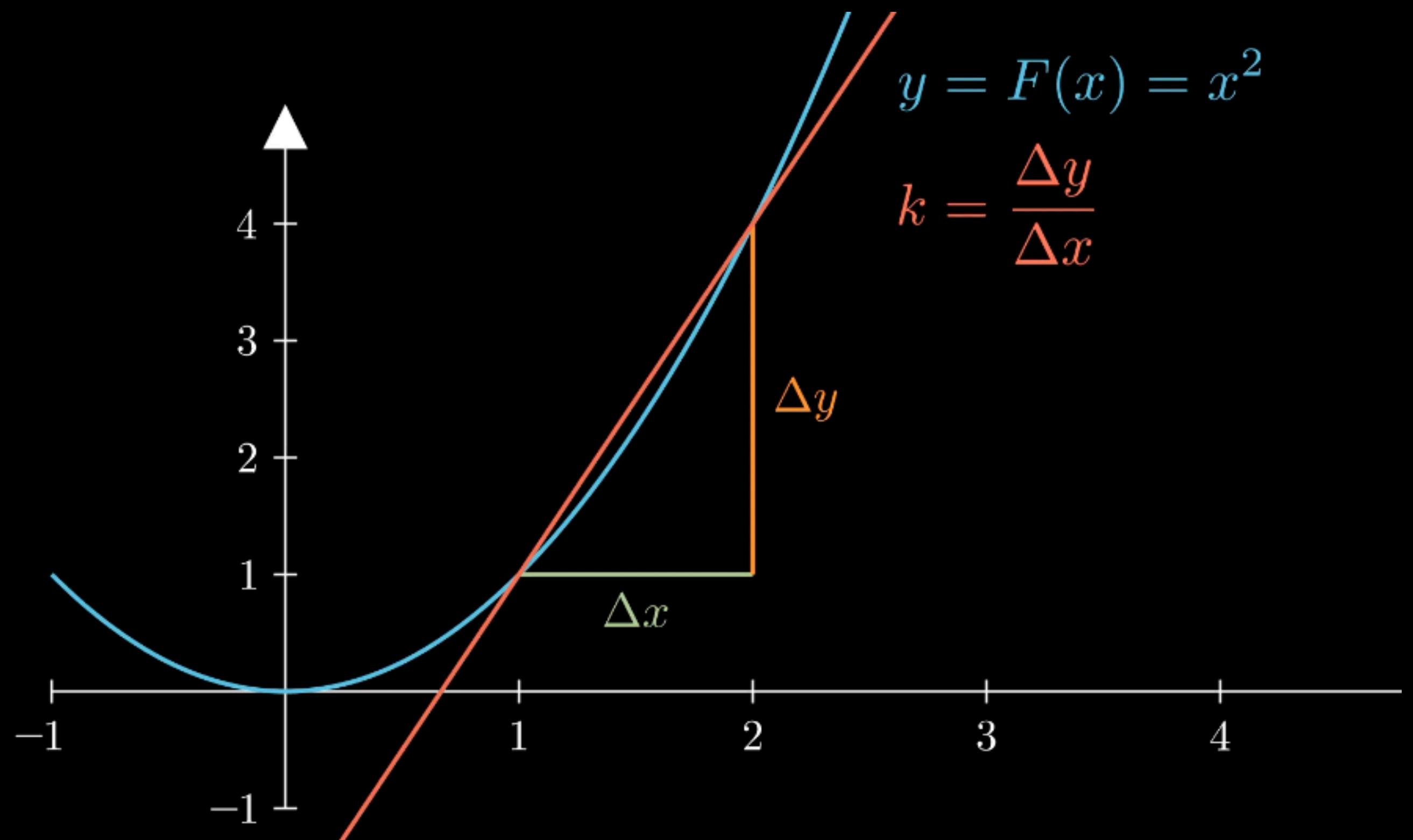
- 概率论视角
- 循环神经网络RNN
- 注意力机制(Attention mechanism)

导数的诞生

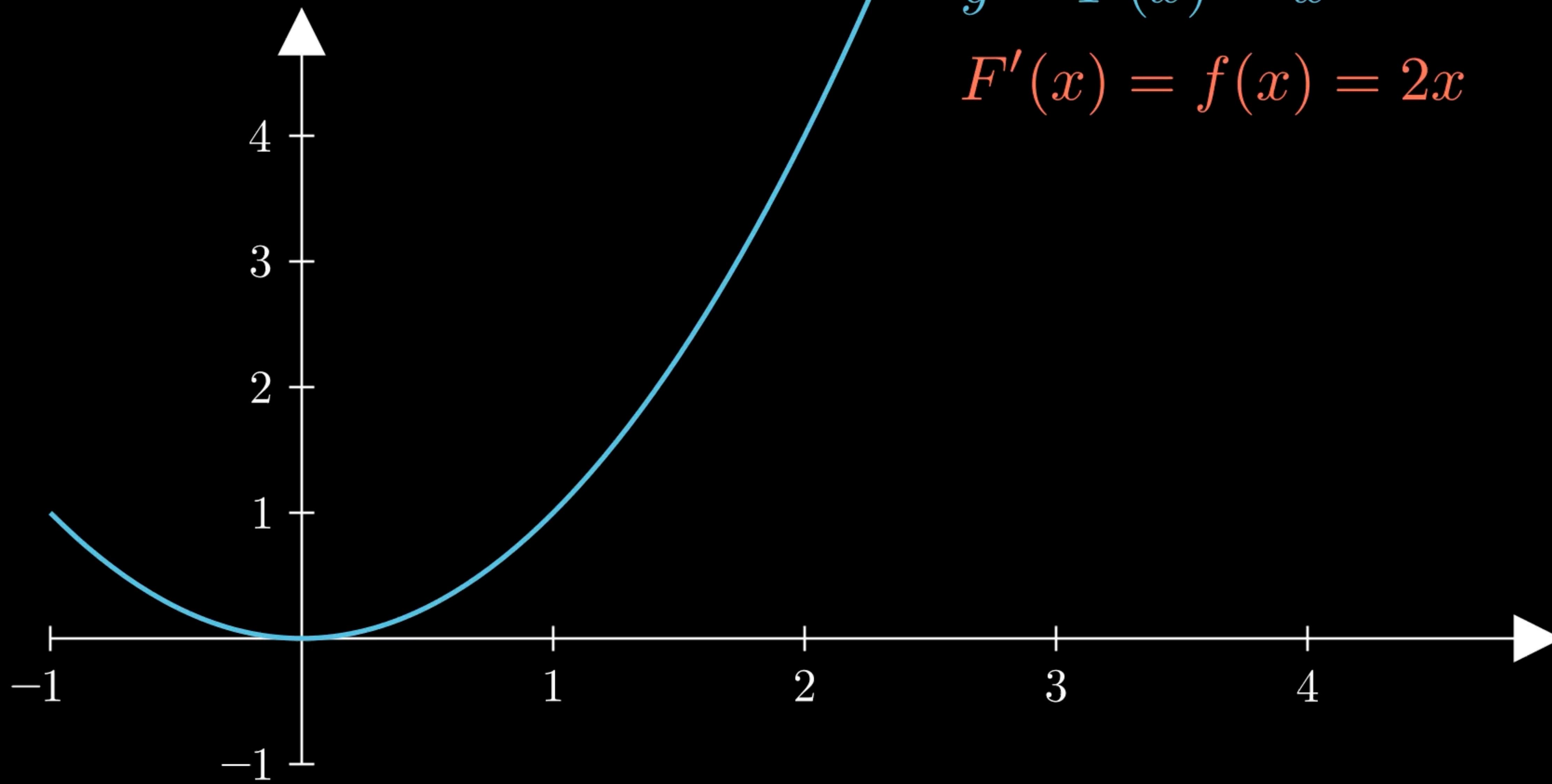


导数的诞生

$$\begin{aligned}f(x) &= F'(x) = y' \\&= \lim_{\Delta x \rightarrow 0} \frac{F(x + \Delta x) - F(x)}{\Delta x} \\&= \lim_{\Delta x \rightarrow 0} \frac{(x + \Delta x) \times (x + \Delta x) - F(x)}{\Delta x} \\&= \lim_{\Delta x \rightarrow 0} \frac{x^2 + 2x\Delta x + \Delta x^2 - x^2}{\Delta x} \\&= \lim_{\Delta x \rightarrow 0} (2x + \Delta x) \\&= 2x\end{aligned}$$



仅已知导函数 $f(x)$, 估计函数 $F(x)$



$$\Delta x = 1$$

x	$f(x) = 2x$	$\Delta y = 2x \cdot \Delta x$	$\hat{F}(x)$	$F(x) = x^2$
0	0	0	0	0
1	2	2	0	1
2	4	4	2	4
3	6	6	6	9

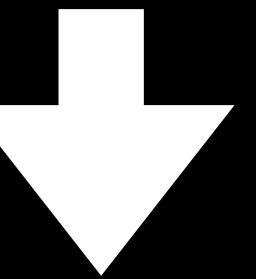
$$\Delta x = 0.5$$

x	$f(x) = 2x$	$\Delta y = 2x \cdot \Delta x$	$\hat{F}(x)$	$F(x) = x^2$
0	0	0	0	0
0.5	1	0.5	0	0.25
1	2	1	0.5	1
1.5	3	1.5	1.5	2.25
2	4	2	3	4
2.5	5	2.5	5	6.25
3	6	3	7.5	9

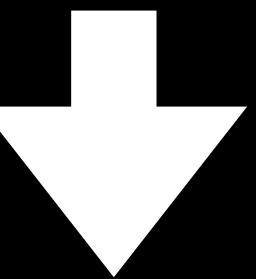
$$\Delta x = 0.01$$

	A	B	C	D
1	x	dy	$F(x)_{\text{hat}}$	$F(x)$
297	2.95	0.059	8.673	8.7025
298	2.96	0.0592	8.732	8.7616
299	2.97	0.0594	8.7912	8.8209
300	2.98	0.0596	8.8506	8.8804
301	2.99	0.0598	8.9102	8.9401
302	3	0.06	8.97	9

$$F(N) \approx \sum_{i=0}^m f(x_i) \cdot \Delta x, \quad \Delta x = N/m$$

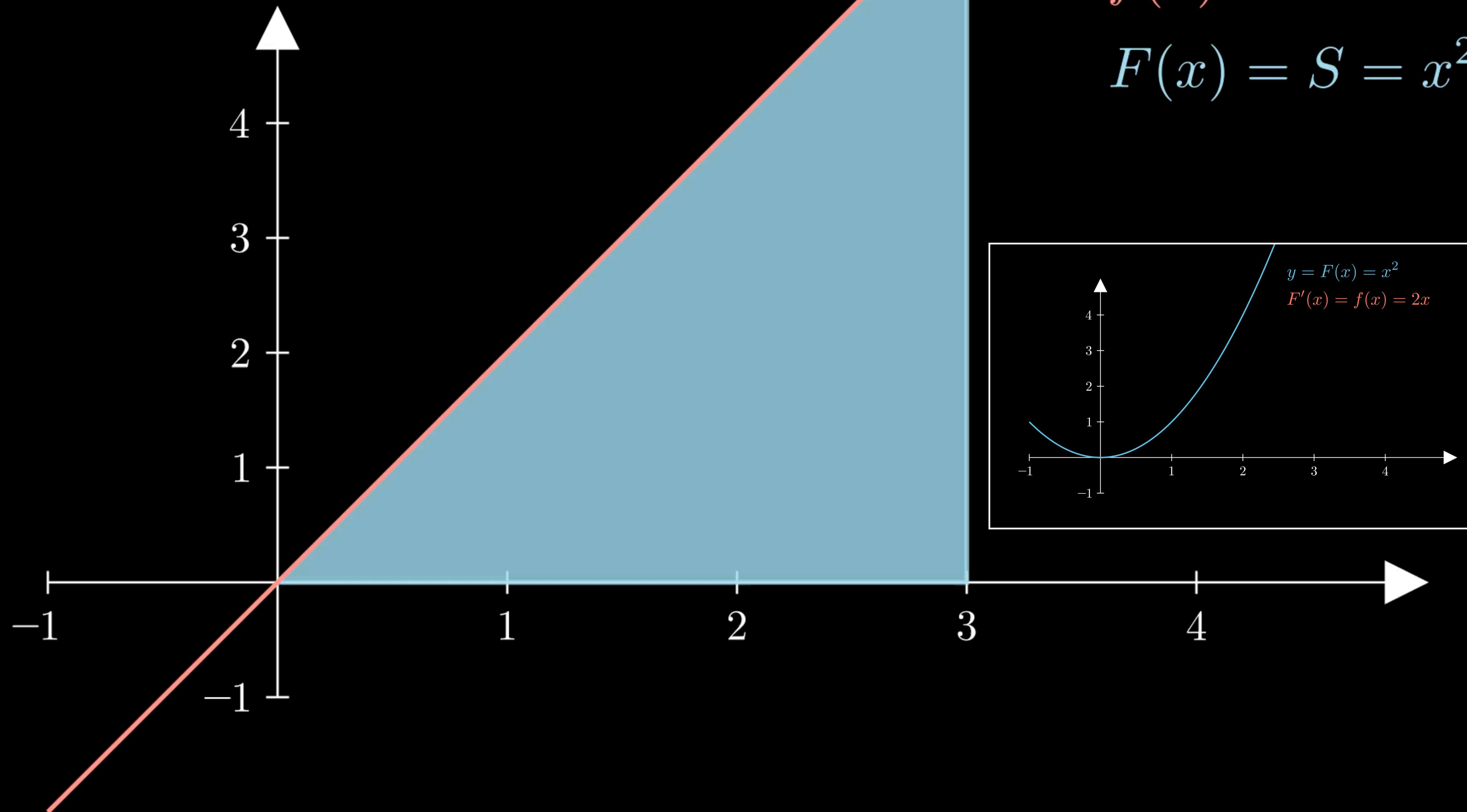


$$F(N) = \sum_{i=0}^{m \rightarrow \infty} f(x_i) \cdot \Delta x, c = \Delta x = N/m \rightarrow 0$$



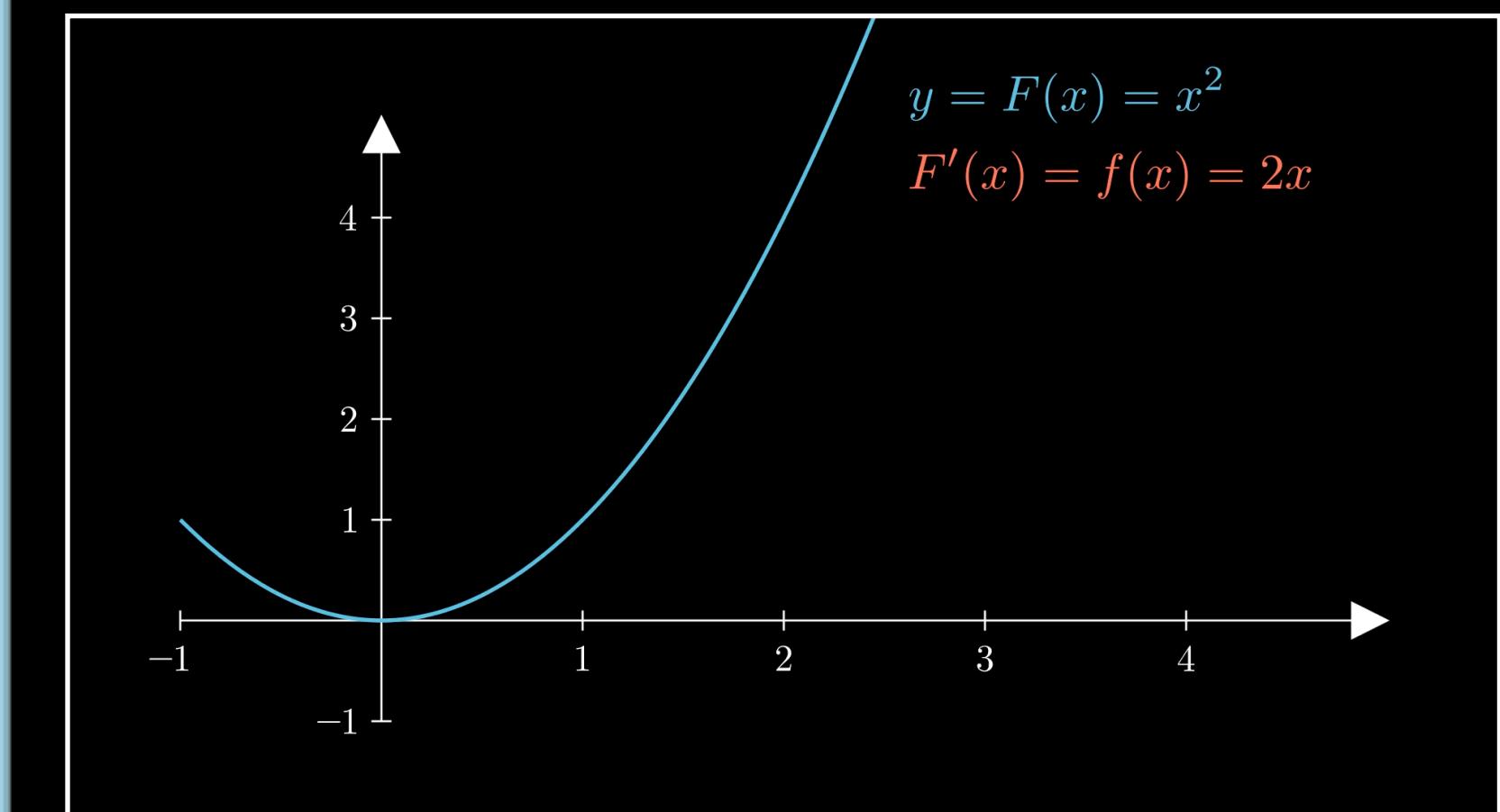
$$\int_{x=0}^N f(x) dx = F(N) - F(0)$$

更常见理解：曲线下面积(Area under curve)

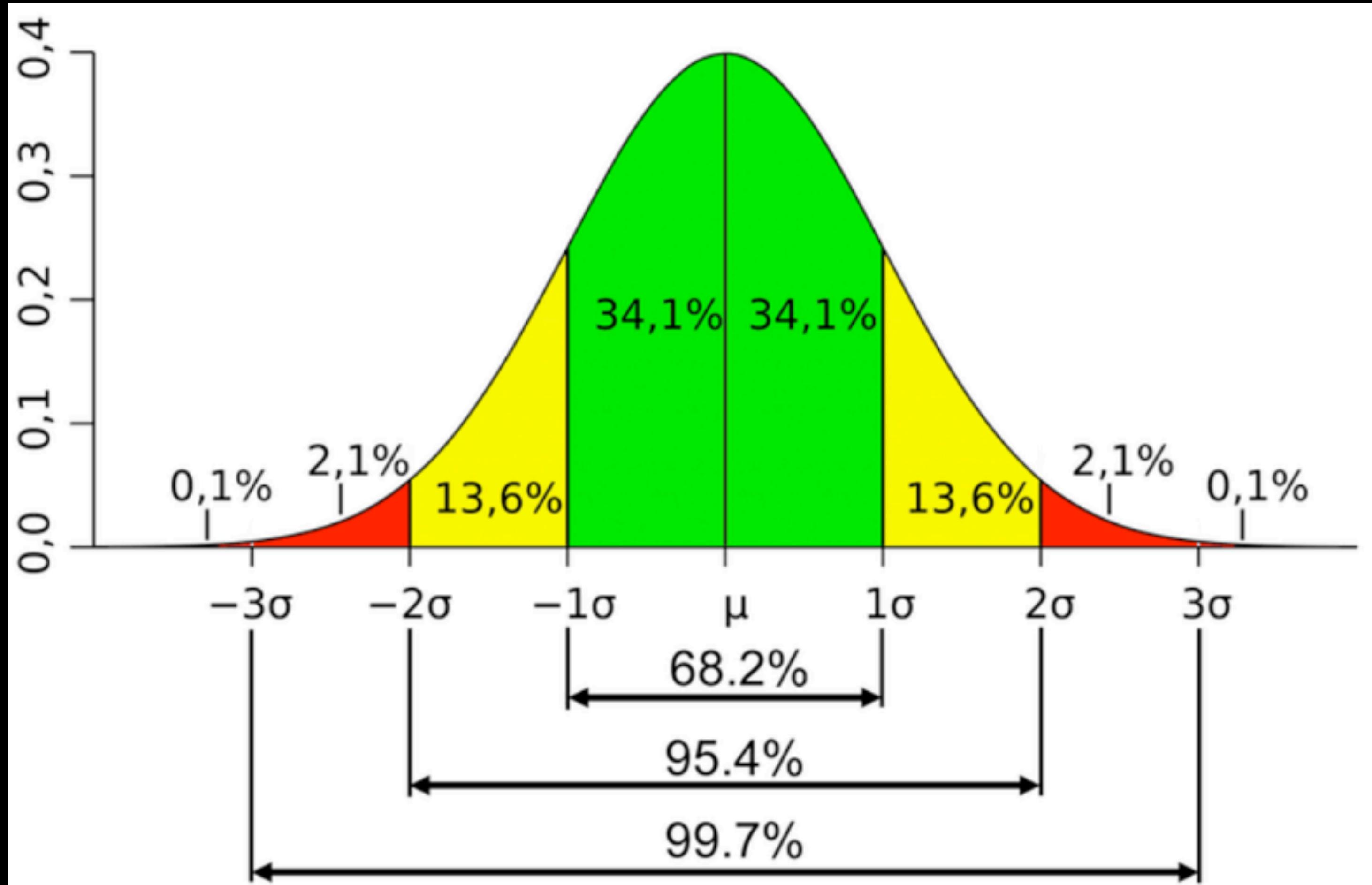


$$f(x) = 2 \times x$$

$$F(x) = S = x^2$$



更常见理解：曲线下面积(Area under curve)



$$\int_{x=x_l}^{x_u} f(x) dx = F(x_u) - F(x_l)$$

标准正态分布：概率密度函数(pdf)
积分得到概率

一些常见的 导数公式

Basic Properties and Formulas

If $f(x)$ and $g(x)$ are differentiable functions (the derivative exists), c and n are any real numbers,

$$1. (c f)' = c f'(x)$$

$$5. \frac{d}{dx}(c) = 0$$

$$2. (f \pm g)' = f'(x) \pm g'(x)$$

$$6. \frac{d}{dx}(x^n) = n x^{n-1} - \textbf{Power Rule}$$

$$3. (f g)' = f' g + f g' - \textbf{Product Rule}$$

$$7. \frac{d}{dx}(f(g(x))) = f'(g(x))g'(x)$$

$$4. \left(\frac{f}{g}\right)' = \frac{f' g - f g'}{g^2} - \textbf{Quotient Rule}$$

This is the **Chain Rule**

Common Derivatives

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(\csc x) = -\csc x \cot x$$

$$\frac{d}{dx}(a^x) = a^x \ln(a)$$

$$\frac{d}{dx}(\sin x) = \cos x$$

$$\frac{d}{dx}(\cot x) = -\csc^2 x$$

$$\frac{d}{dx}(\mathbf{e}^x) = \mathbf{e}^x$$

$$\frac{d}{dx}(\cos x) = -\sin x$$

$$\frac{d}{dx}(\sin^{-1} x) = \frac{1}{\sqrt{1-x^2}}$$

$$\frac{d}{dx}(\ln(x)) = \frac{1}{x}, \quad x > 0$$

$$\frac{d}{dx}(\tan x) = \sec^2 x$$

$$\frac{d}{dx}(\cos^{-1} x) = -\frac{1}{\sqrt{1-x^2}}$$

$$\frac{d}{dx}(\ln|x|) = \frac{1}{x}, \quad x \neq 0$$

$$\frac{d}{dx}(\sec x) = \sec x \tan x$$

$$\frac{d}{dx}(\tan^{-1} x) = \frac{1}{1+x^2}$$

$$\frac{d}{dx}(\log_a(x)) = \frac{1}{x \ln a}, \quad x > 0$$

Derivative Rules

Product Rule

乘法求导规则

$$y = f(x) \cdot g(x)$$

$$y' = f'(x) \cdot g(x) + f(x) \cdot g'(x)$$

Ex

$$y = x^4 \cdot \sin(x)$$

$$f(x) = x^4 \quad g(x) = \sin(x)$$

$$f'(x) = 4x^3 \quad g'(x) = \cos(x)$$

$$y' = 4x^3 \cdot \sin(x) + x^4 \cdot \cos(x)$$

Ex

$$y = 2x^2 \cdot \ln(x)$$

$$f(x) = 2x^2 \quad g(x) = \ln(x)$$

$$f'(x) = 4x \quad g'(x) = \frac{1}{x}$$

$$y' = 4x \cdot \ln(x) + 2x^2 \cdot \frac{1}{x}$$

$$y' = 4x \cdot \ln(x) + 2x$$

Quotient Rule

除法求导规则

$$y = \frac{f(x)}{g(x)} = \frac{hi}{lo}$$

$$y' = \frac{lo \cdot dhi - hi \cdot dlo}{(lo)^2}$$

Ex

$$y = \frac{e^x}{4x^3}$$

$$hi = e^x \quad lo = 4x^3$$

$$dhi = e^x \quad dlo = 12x^2$$

$$y' = \frac{4x^3 \cdot e^x - e^x \cdot 12x^2}{(4x^3)^2}$$

Ex

$$y = \frac{5x}{\sin(x)}$$

$$hi = 5x \quad lo = \sin(x)$$

$$dhi = 5 \quad dlo = \cos(x)$$

$$y' = \frac{\sin(x) \cdot 5 - 5x \cdot \cos(x)}{(\sin(x))^2}$$

Chain Rule

链式求导规则

$$y = f(g(x))$$

$$y' = f'(g(x)) \cdot g'(x)$$

Ex

$$y = \cos(3x^2)$$

$$f(g(x)) = \cos(3x^2) \quad g(x) = 3x^2$$

$$f'(g(x)) = -\sin(3x^2) \quad g'(x) = 6x$$

$$y' = -\sin(3x^2) \cdot 6x$$

Ex

$$y = e^{2x^3+4}$$

$$f(g(x)) = e^{2x^3+4} \quad g(x) = 2x^3 + 4$$

$$f'(g(x)) = e^{2x^3+4} \quad g'(x) = 6x^2$$

$$y' = e^{2x^3+4} \cdot 6x^2$$

求导：链式法则 (chain rule)

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

$$z = \log(x_1 \cdot x_2) \cdot \sin x_2$$

$$\frac{\partial z}{\partial x_1} = \frac{x_2}{x_1 x_2} \cdot \sin x_2 + \log(x_1 x_2) \cdot 0 = \frac{\sin x_2}{x_1}$$

$$\frac{\partial z}{\partial x_2} = \frac{x_1}{x_1 x_2} \cdot \sin x_2 + \log(x_1 x_2) \cdot \cos x_2 = \frac{\sin x_2}{x_2} + (\cos x_2) \cdot \log(x_1 x_2)$$

手算结果（基本上是高中数学内容）

Take derivatives (求导)

```
In[2]:= f[z] := Log[x1 * x2] * Sin[x2]
```

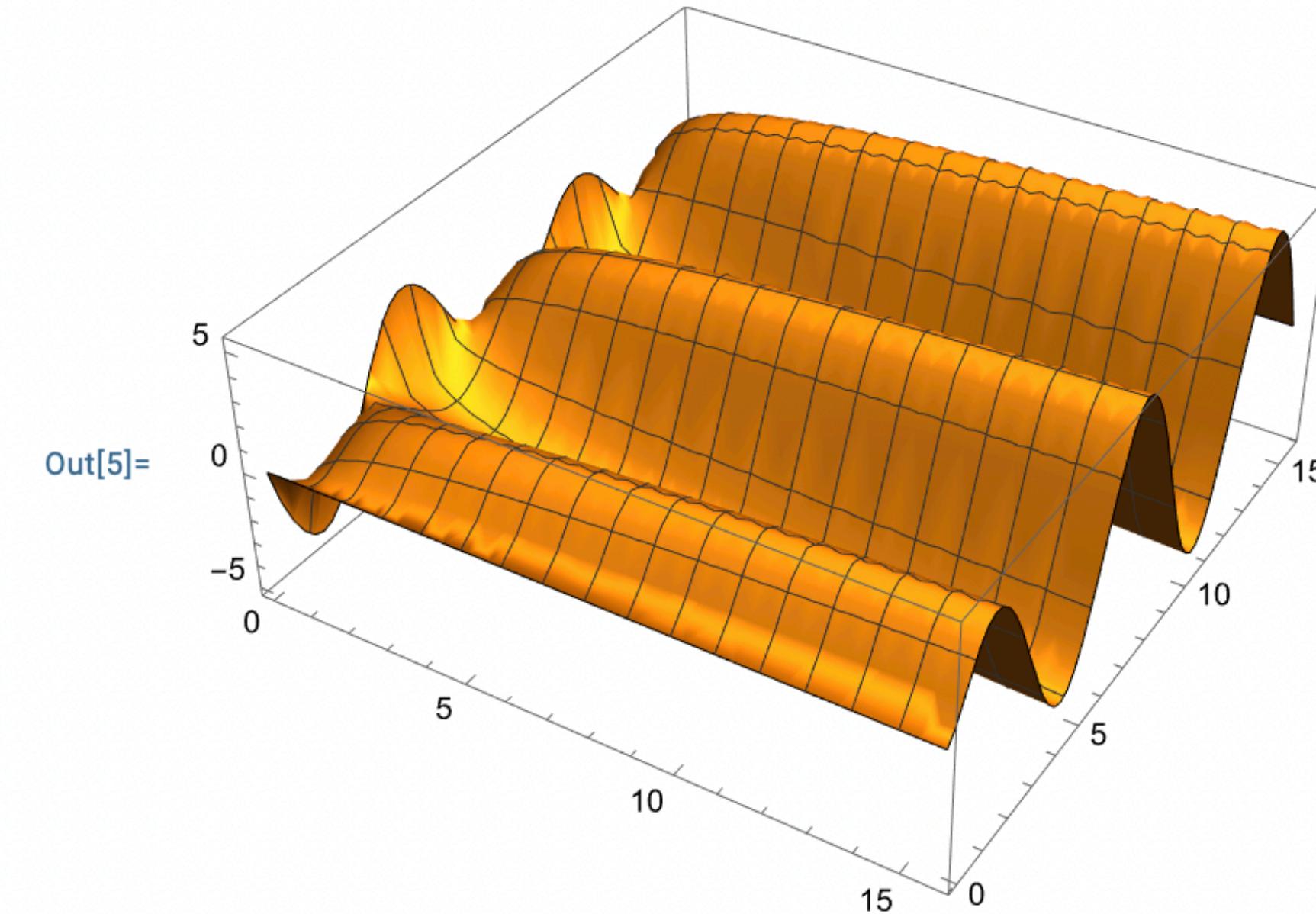
```
In[3]:= D[f[z], x1]
```

$$\text{Out}[3]= \frac{\text{Sin}[x2]}{x1}$$

```
In[4]:= D[f[z], x2]
```

$$\text{Out}[4]= \text{Cos}[x2] \text{Log}[x1 x2] + \frac{\text{Sin}[x2]}{x2}$$

```
In[5]:= Plot3D[Log[x1 * x2] * Sin[x2], {x1, 0.01, 5*Pi}, {x2, 0.01, 5*Pi}]
```



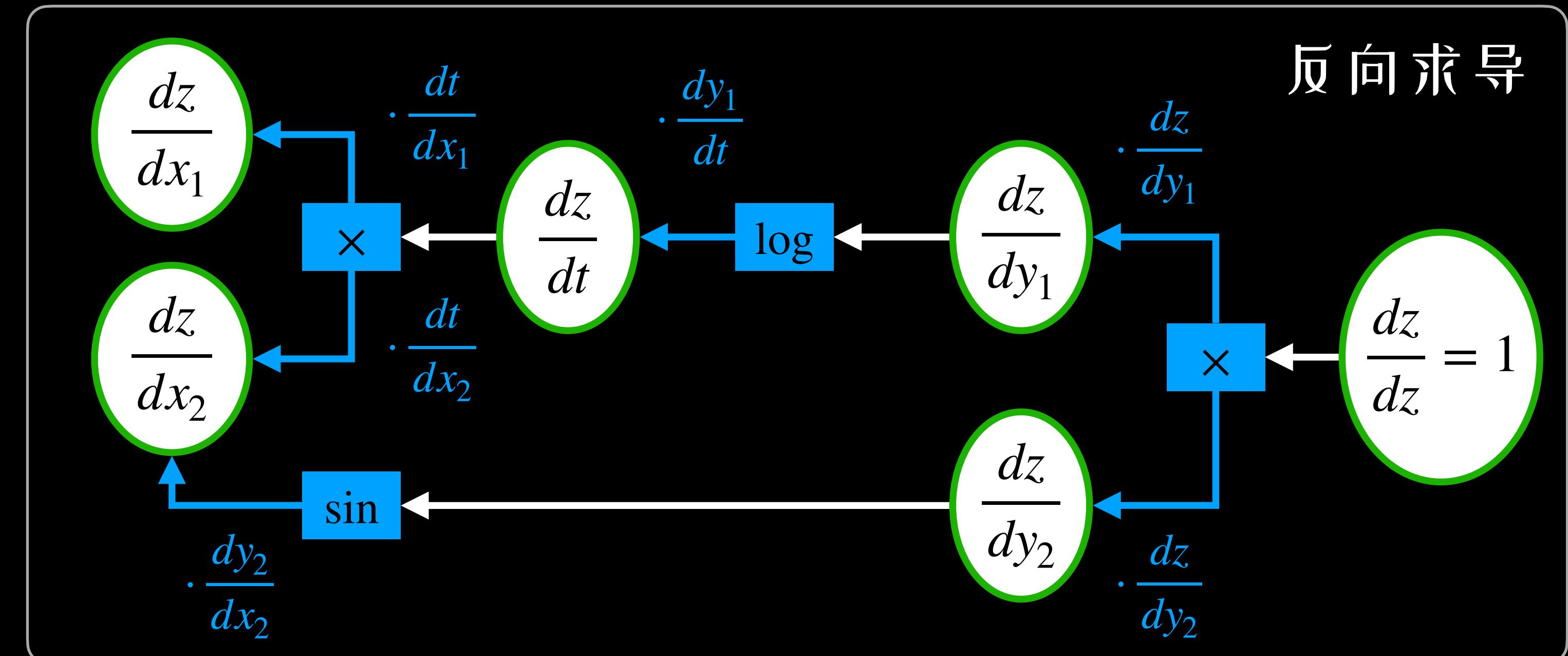
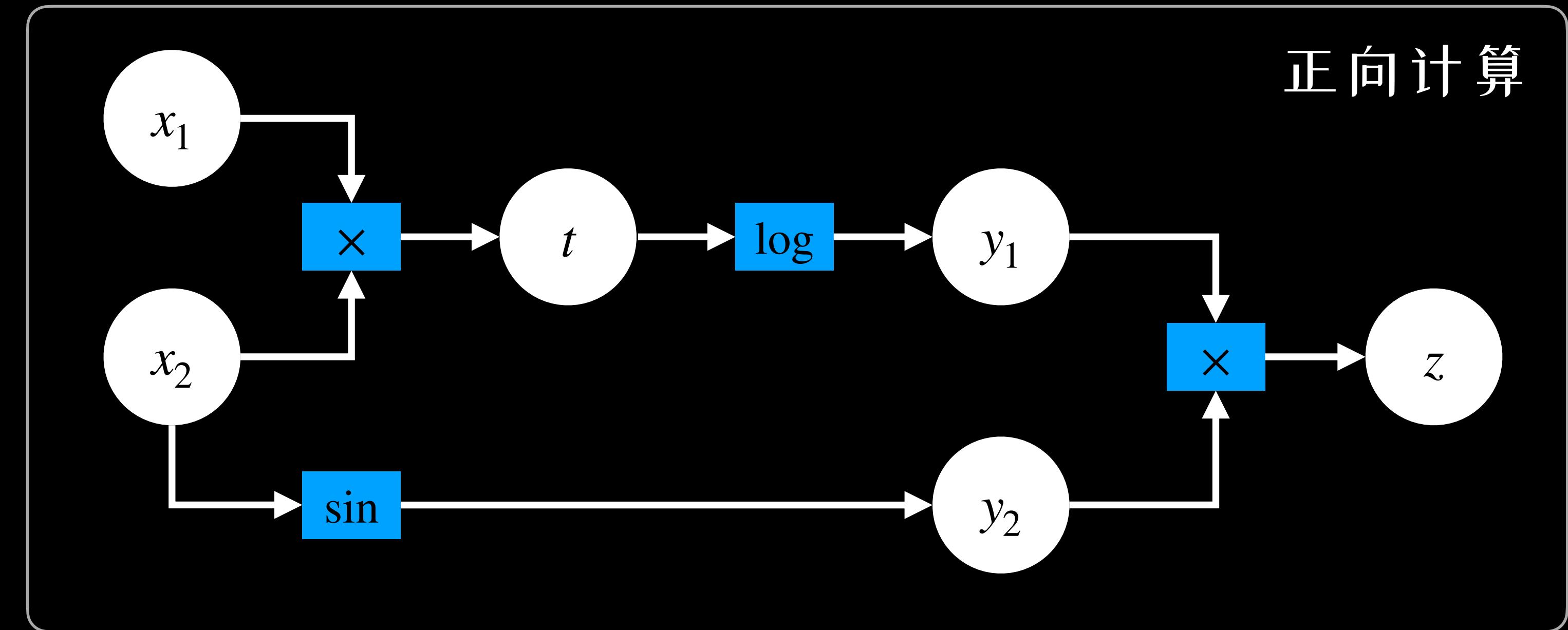
Mathematica求导结果

PyTorch/Tensorflow: 通过计算图推导/求导

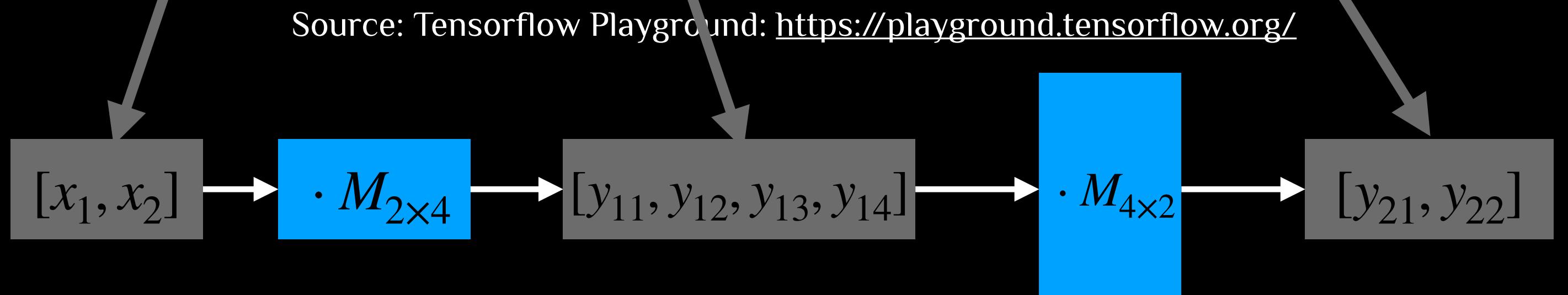
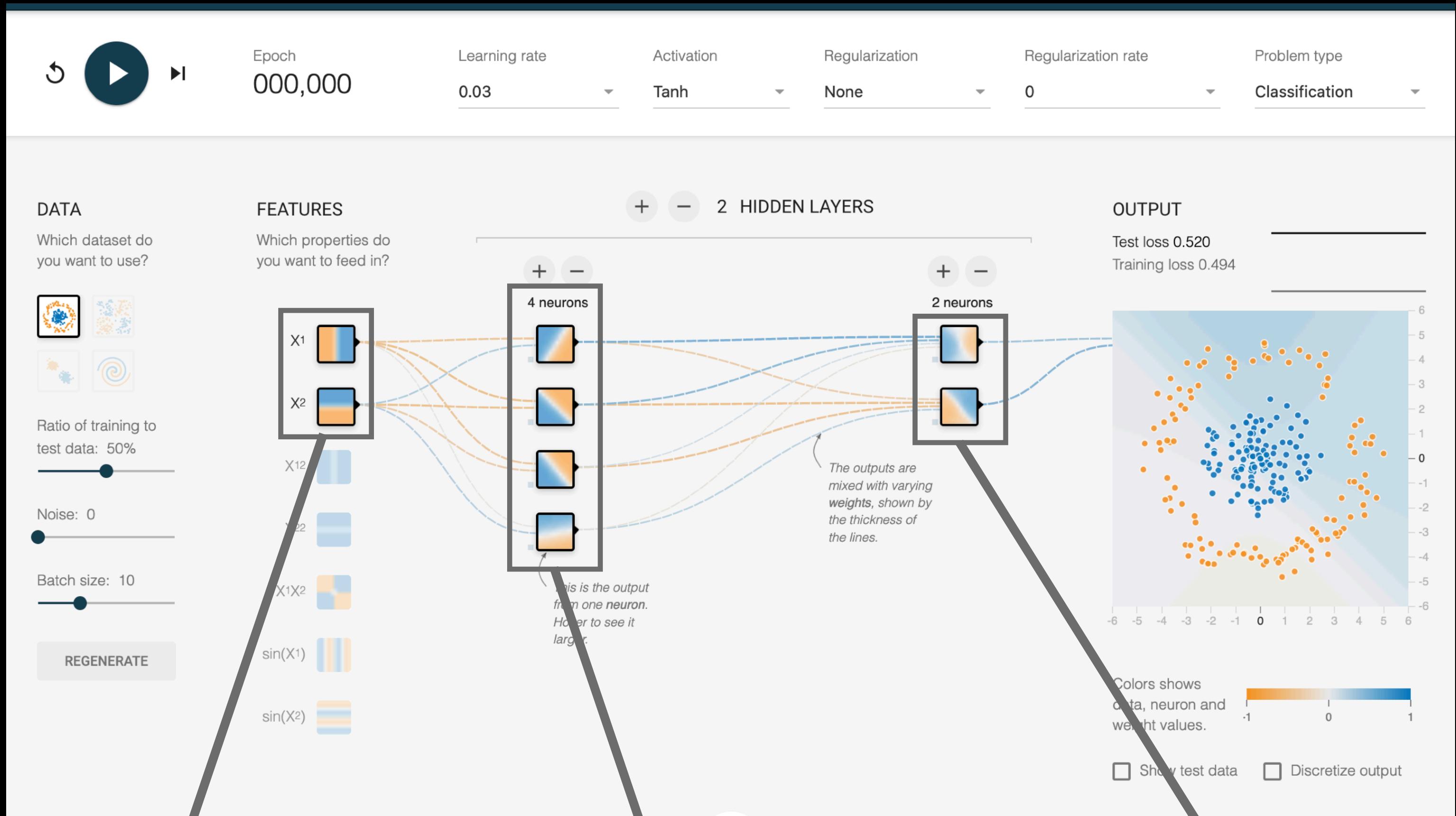
$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

$$z = \log(x_1 \cdot x_2) \cdot \sin x_2$$

1. 每种运算的导数是确定的
2. 通过计算图，实现连续求导
3. 一遍运算，求出所有偏导
4. 需要带入具体数值（不做符号推导，只是传递导数数值）



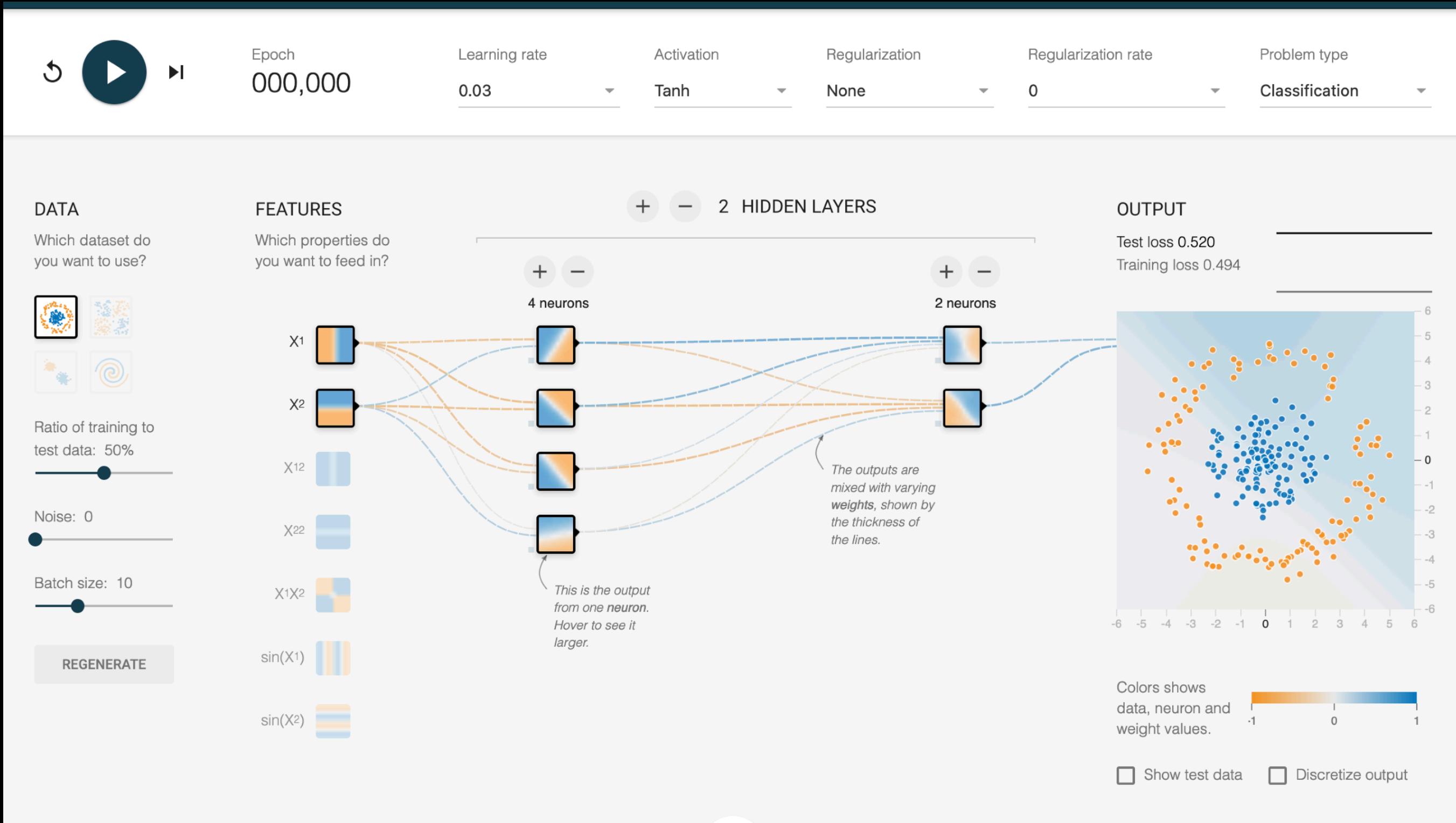
深度学习：线性变换(全连接网络/FCNN)



分类器

1. 输入 : (x_1, x_2)
2. 输出 : $y = -1 \text{ or } 1$

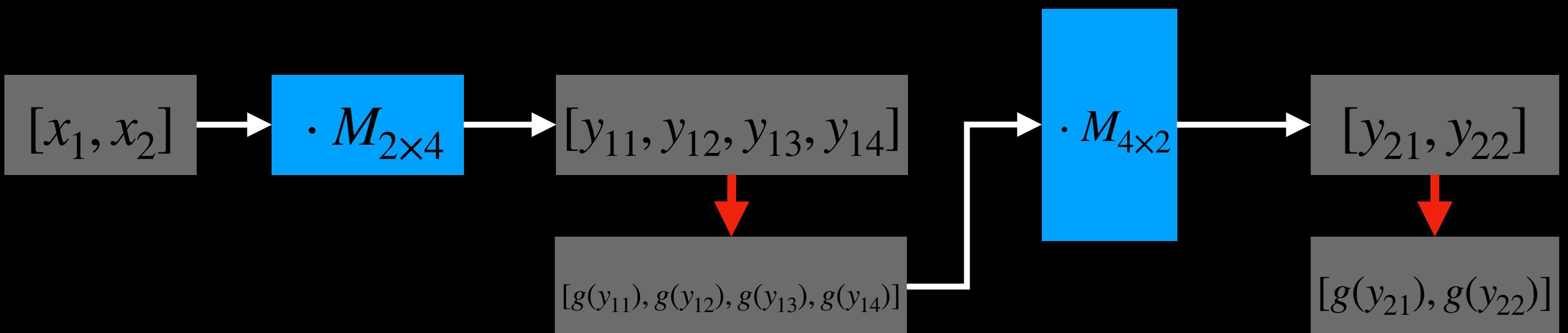
深度学习：线性变换(全连接网络/FCNN)



$$M_{2 \times 4} \cdot M_{4 \times 2} = M_{2 \times 2}$$

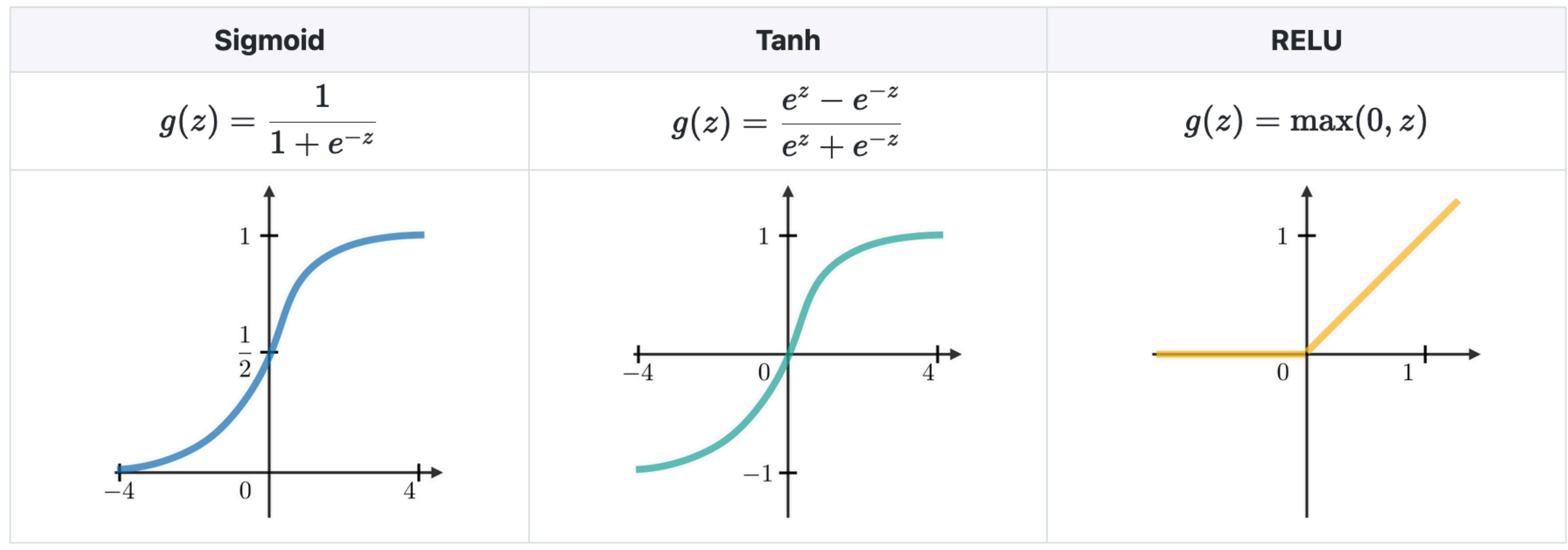
多个矩阵相乘 等价于一个矩阵
所以，需要在那个矩阵运算后添
加非线性变换

Source: Tensorflow Playground: <https://playground.tensorflow.org/>



深度学习：非线性变换（单个神经元的变换）

Commonly used activation functions – The most common activation functions used in RNN modules are described below:

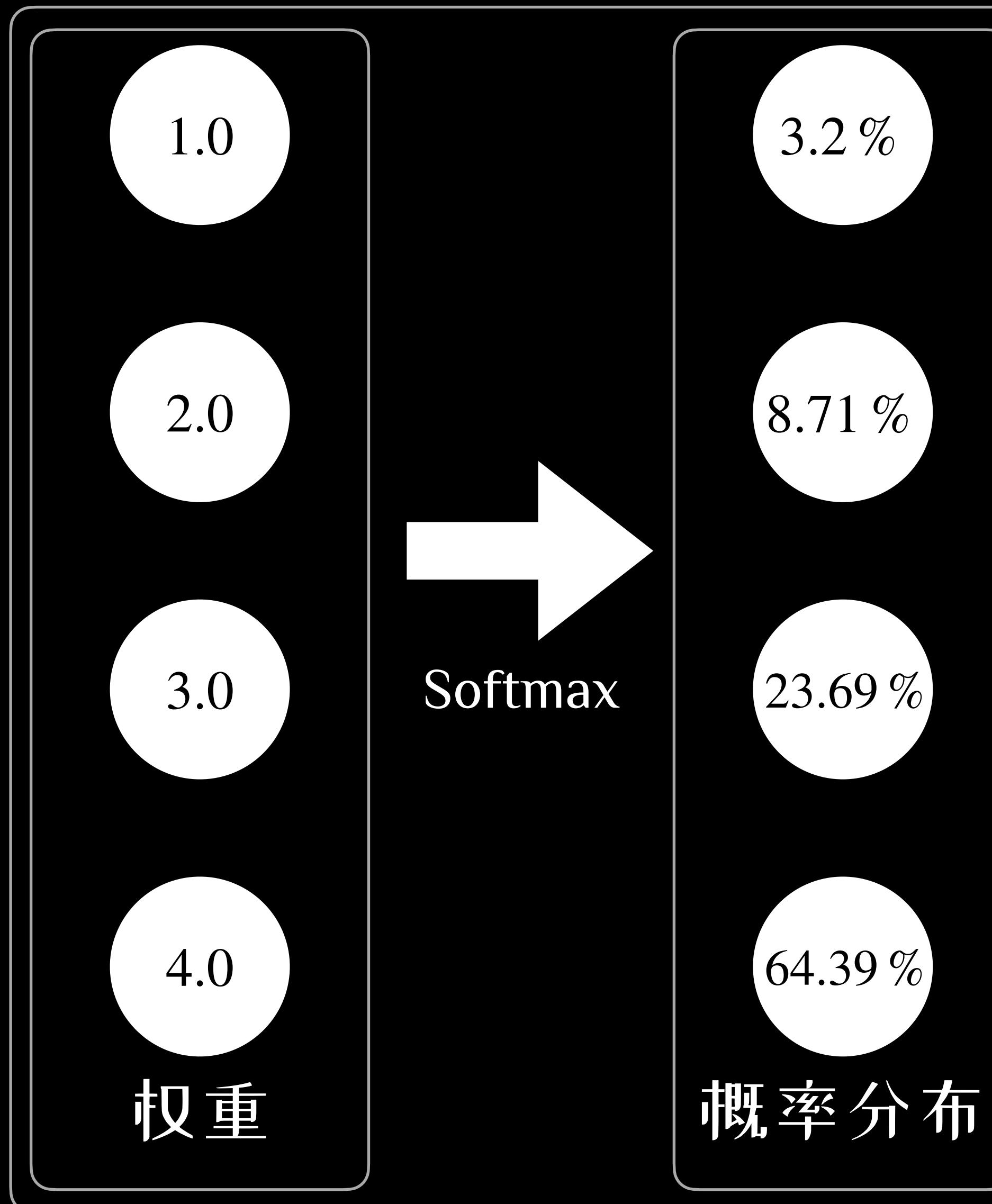


$$\frac{dg(z)}{dz} = g(z)(1 - g(z))$$

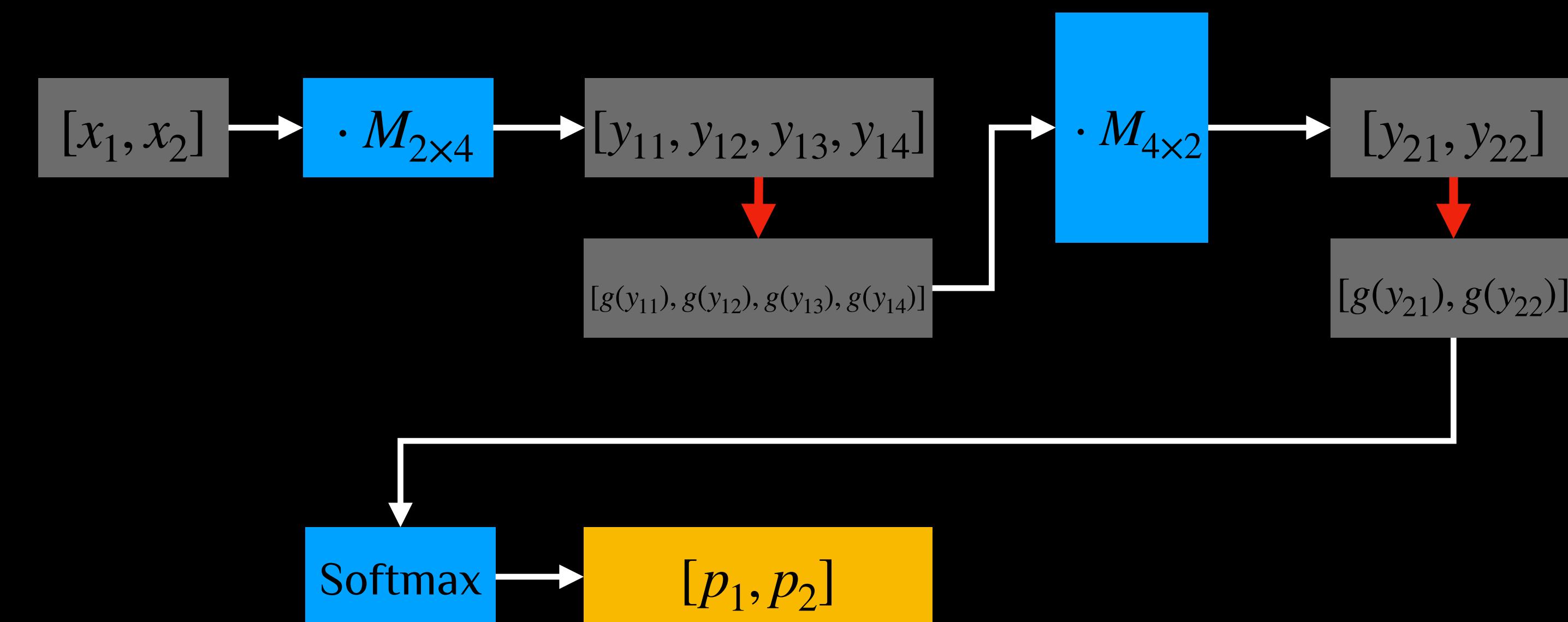
$$\frac{dg(z)}{dz} = 1 - (g(z))^2$$

$$\frac{dg(z)}{dz} = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases}$$

深度学习：非线性变换（Softmax单层神经元变换为概率）



$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad \text{for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K.$$



深度学习: 非线性变换 (损失函数)

Cross Entropy Loss

$$l = - \sum_{i=1}^M t_i \log p_i, M \text{ 个标注类型}$$

$$L = \sum_{j=1}^N l, N \text{ 条训练数据}$$

样本1

[0,1]

真实概率

[0.1,0.9]

模型预测概率

$$l = - (0 \cdot \log[0.1] + 1 \cdot \log[0.9]) = 0.1053$$

样本2

[0,1]

真实概率

[0.9,0.1]

模型预测概率

$$l = - (0 \cdot \log[0.9] + 1 \cdot \log[0.1]) = 2.3026$$

梯度学习: 非线性变换 (损失函数)

$$l = - \sum_{i=1}^M t_i \log p_i, M \text{ 个标注类型}$$

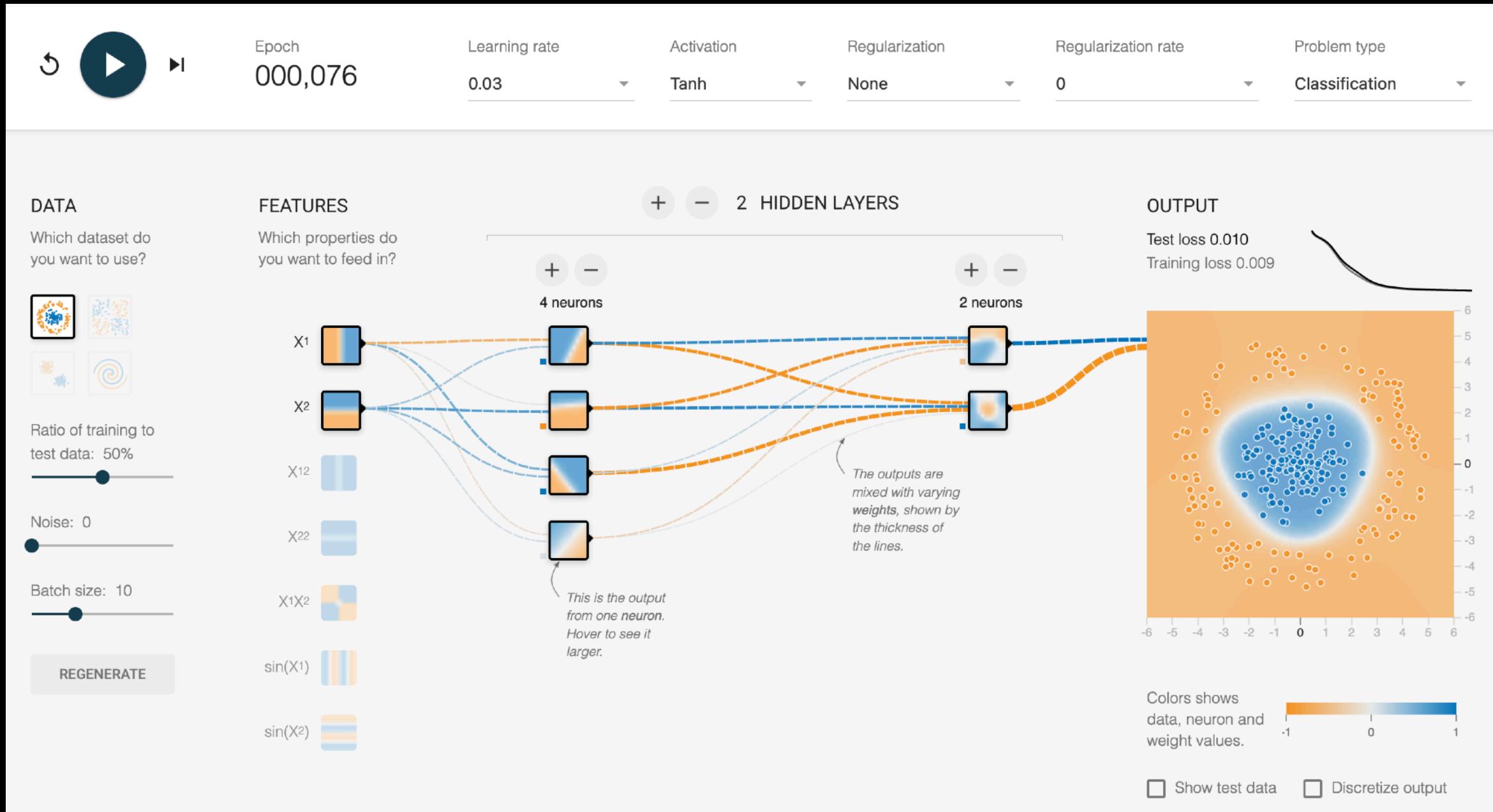
$$L = \sum_{j=1}^N l, N \text{ 条训练数据}$$

$$\frac{\partial L}{\partial w_k}, \forall k \in \{1, \dots, K\}$$

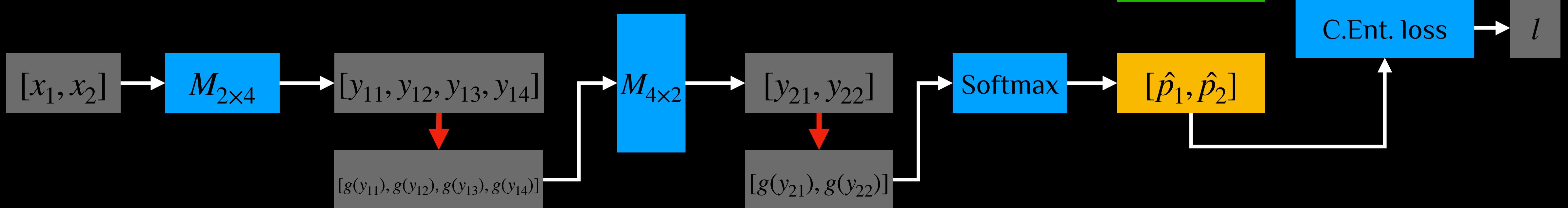
梯度下降更新参数(Gradient Descent)

$$w_k := w_k - \gamma \frac{\partial L}{\partial w_k} \rightarrow \begin{cases} n = N, \text{G.D.} \\ n = \text{B.S.} = 4, \text{S.G.D.} \end{cases}$$

深度学习：线性变换(全连接网络/FCNN)



训练几秒钟后...

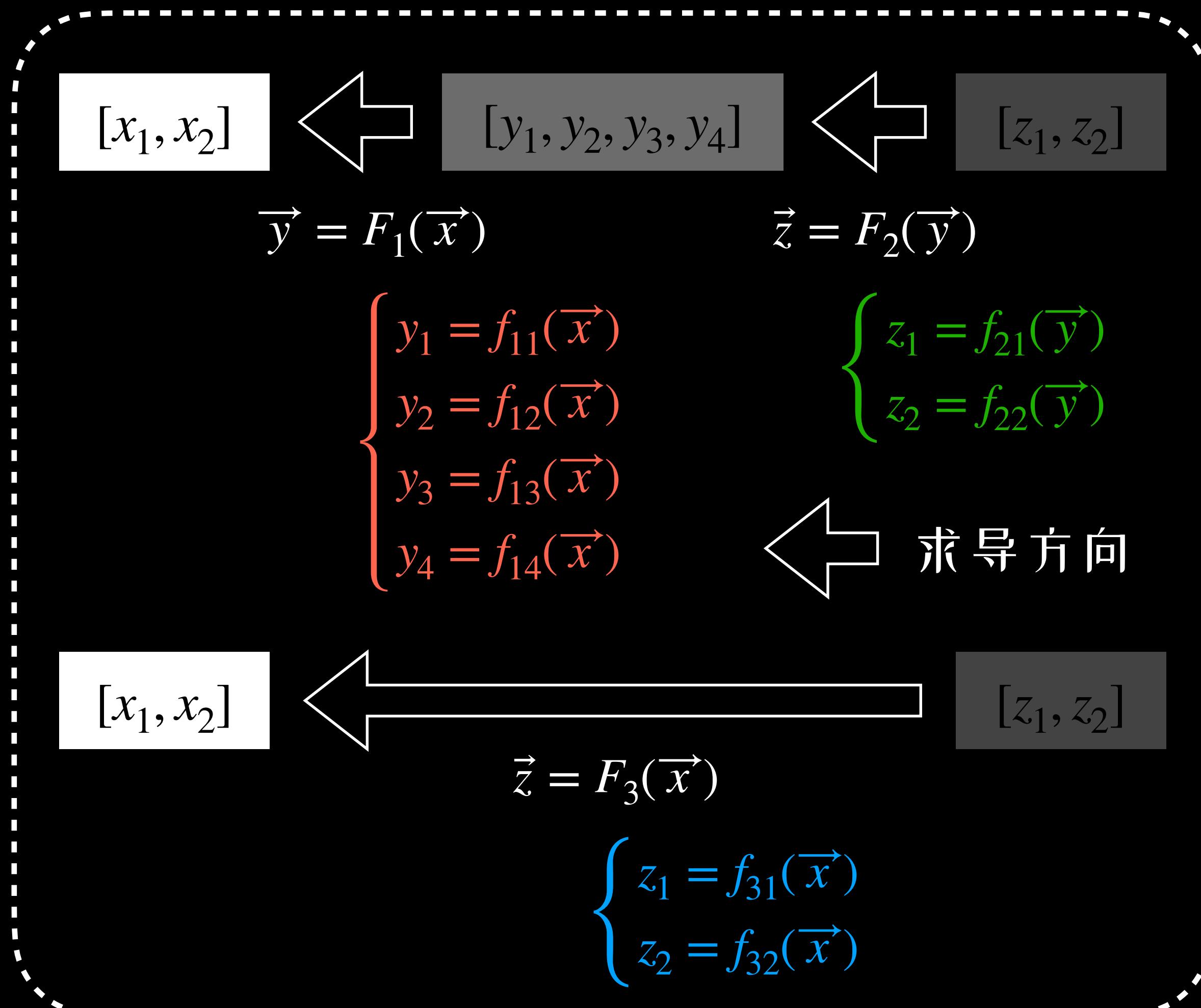


对应的大致计算图

需要训练的模型参数
(矩阵权重)：
 $2 \times 4 + 4 \times 2 = 16$

向量化/矩阵化/张量化(vectorization)

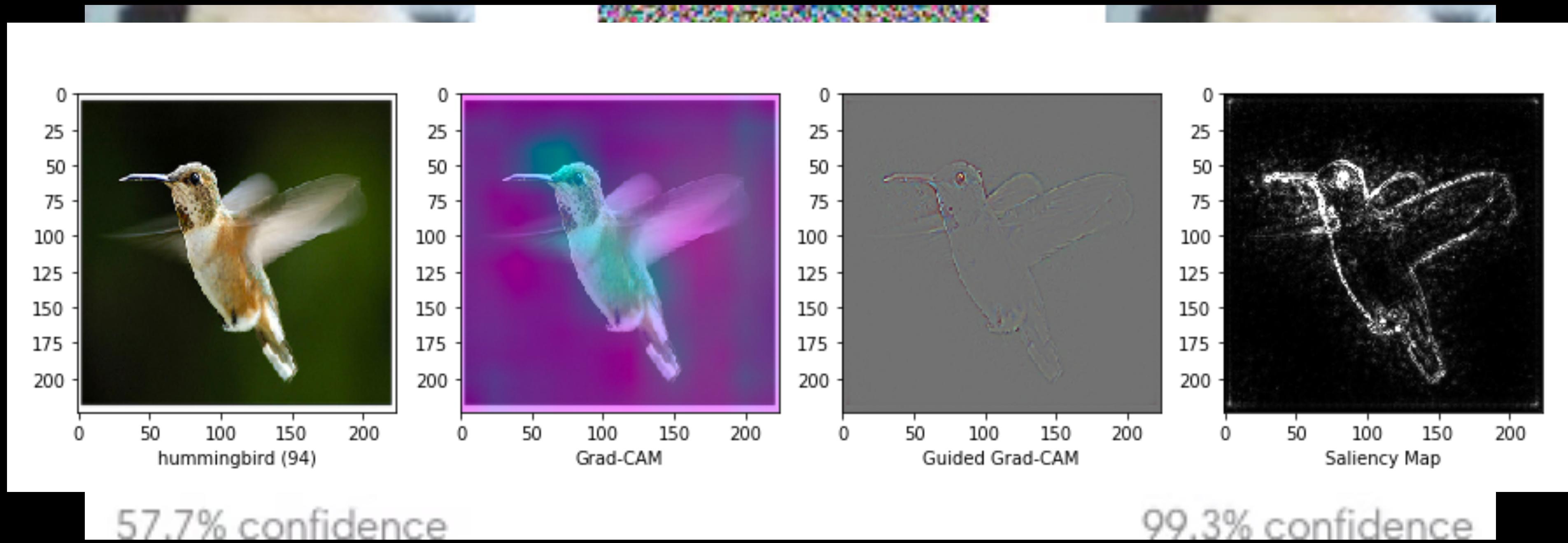
- Jacobian Matrix (每层之间导数传递遵循矩阵“乘法”运算法则)



$$\begin{array}{c}
 \left[\begin{array}{cc} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} \\ \frac{\partial y_3}{\partial x_1} & \frac{\partial y_3}{\partial x_2} \\ \frac{\partial y_4}{\partial x_1} & \frac{\partial y_4}{\partial x_2} \end{array} \right] \cdot \left[\begin{array}{cc} \frac{\partial z_1}{\partial y_1} & \frac{\partial z_1}{\partial y_2} \\ \frac{\partial z_2}{\partial y_1} & \frac{\partial z_2}{\partial y_2} \end{array} \right] = \left[\begin{array}{cc} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_1}{\partial x_2} \\ \frac{\partial z_2}{\partial x_1} & \frac{\partial z_2}{\partial x_2} \end{array} \right]
 \end{array}$$

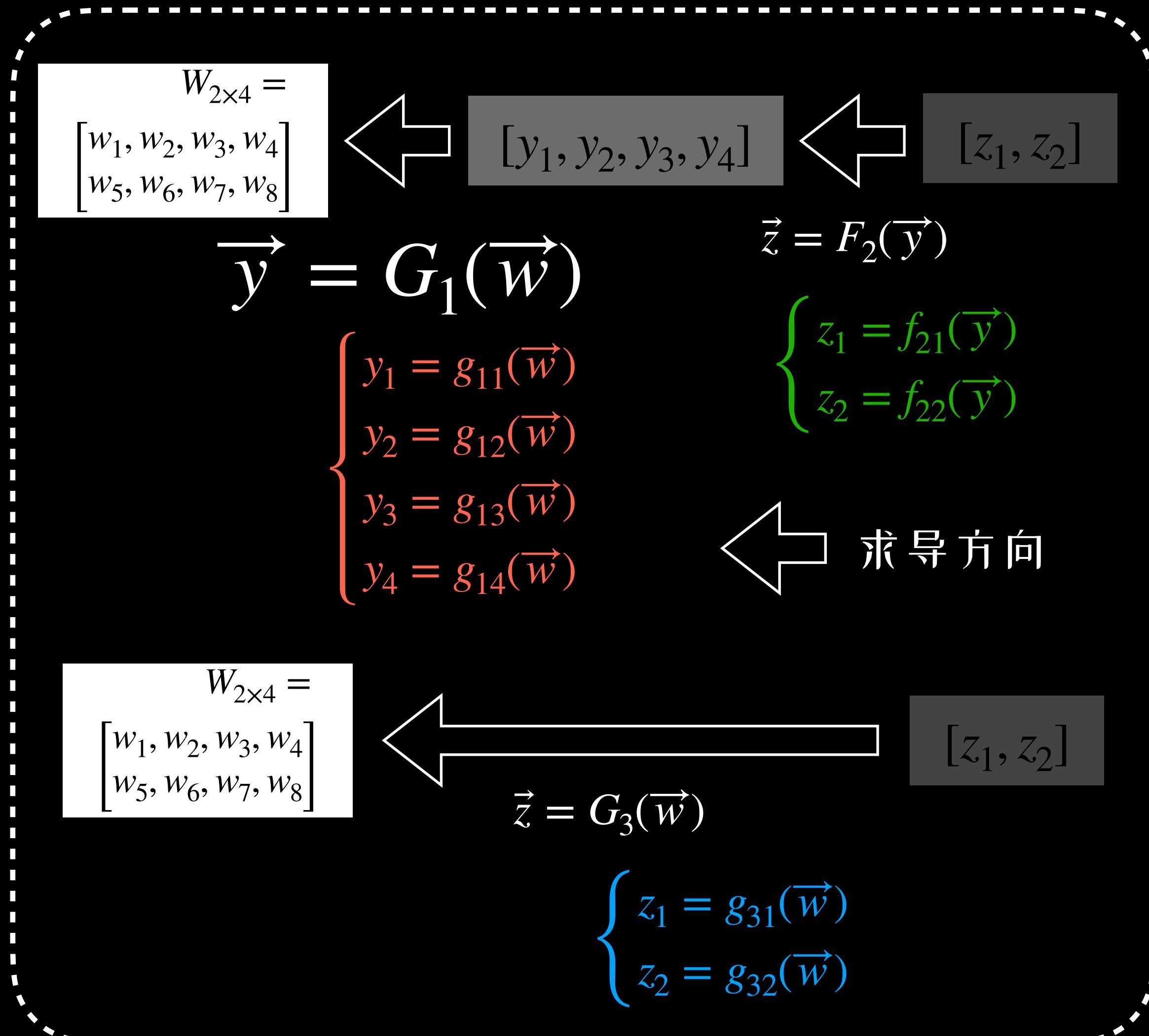
向量化/矩阵化/张量化(vectorization)

- 这里是对“数据” / “神经元”进行求导
- 应用1 攻击模型 Adversarial Attack
- 应用2 解释模型 saliency map (突出“导数更大”的图片像素)



向量化/矩阵化/张量化(vectorization)

- Jacobian Matrix (每层之间导数传递遵循矩阵“乘法”运算法则)



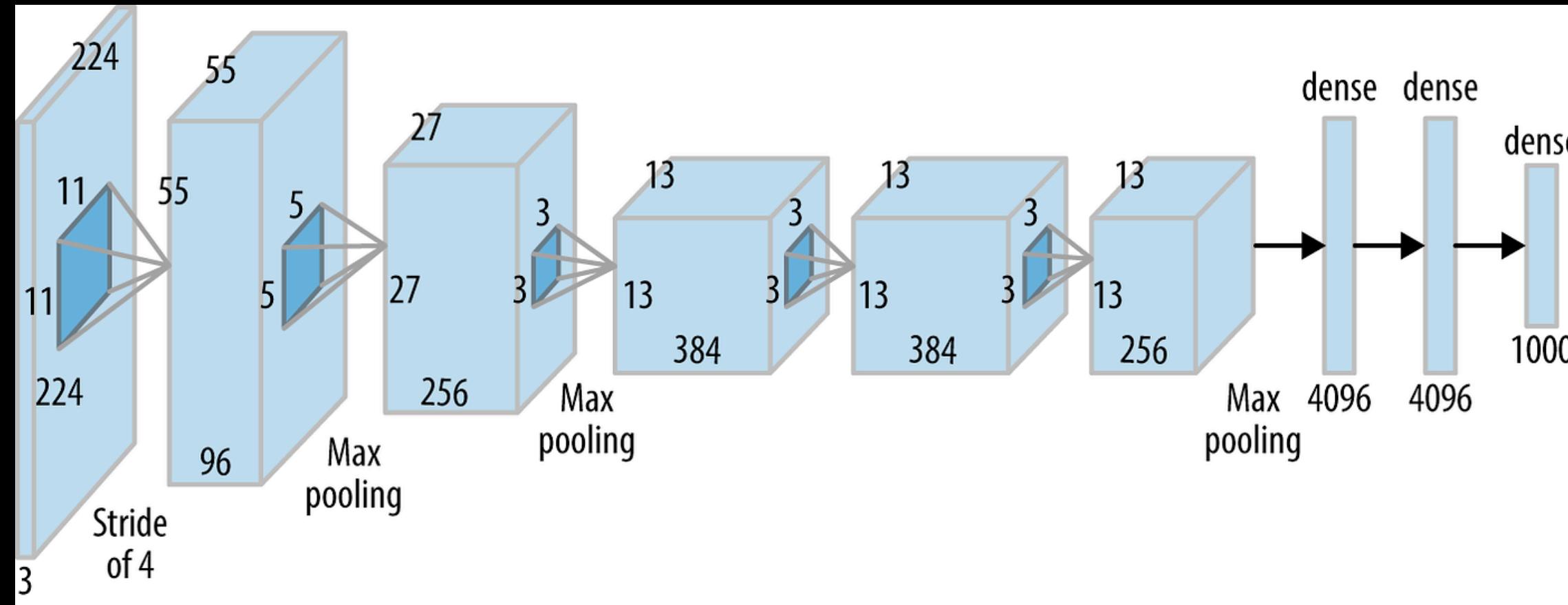
$$\left[\begin{array}{cccc} \frac{\partial z_1}{\partial y_1} & \frac{\partial z_1}{\partial y_2} & \frac{\partial z_1}{\partial y_3} & \frac{\partial z_1}{\partial y_4} \\ \frac{\partial z_2}{\partial y_1} & \frac{\partial z_2}{\partial y_2} & \frac{\partial z_2}{\partial y_3} & \frac{\partial z_2}{\partial y_4} \end{array} \right] \cdot \left[\begin{array}{cccc} \frac{\partial y_1}{\partial w_1} & \frac{\partial y_1}{\partial w_2} & \cdots & \frac{\partial y_1}{\partial w_8} \\ \frac{\partial y_2}{\partial w_1} & \frac{\partial y_2}{\partial w_2} & \cdots & \frac{\partial y_2}{\partial w_8} \\ \frac{\partial y_3}{\partial w_1} & \frac{\partial y_3}{\partial w_2} & \cdots & \frac{\partial y_3}{\partial w_8} \\ \frac{\partial y_4}{\partial w_1} & \frac{\partial y_4}{\partial w_2} & \cdots & \frac{\partial y_4}{\partial w_8} \end{array} \right] = \left[\begin{array}{cccc} \frac{\partial z_1}{\partial w_1} & \frac{\partial z_1}{\partial w_2} & \cdots & \frac{\partial z_1}{\partial w_8} \\ \frac{\partial z_2}{\partial w_1} & \frac{\partial z_2}{\partial w_2} & \cdots & \frac{\partial z_2}{\partial w_8} \end{array} \right]$$

- 训练中，对模型参数求导也是非常类似的
- 我的观点：“自动求导机”
- 通过计算图，可以求二阶导/多阶导
(也很多基于二阶导数的方法)

卷积神经网络 (CNN)

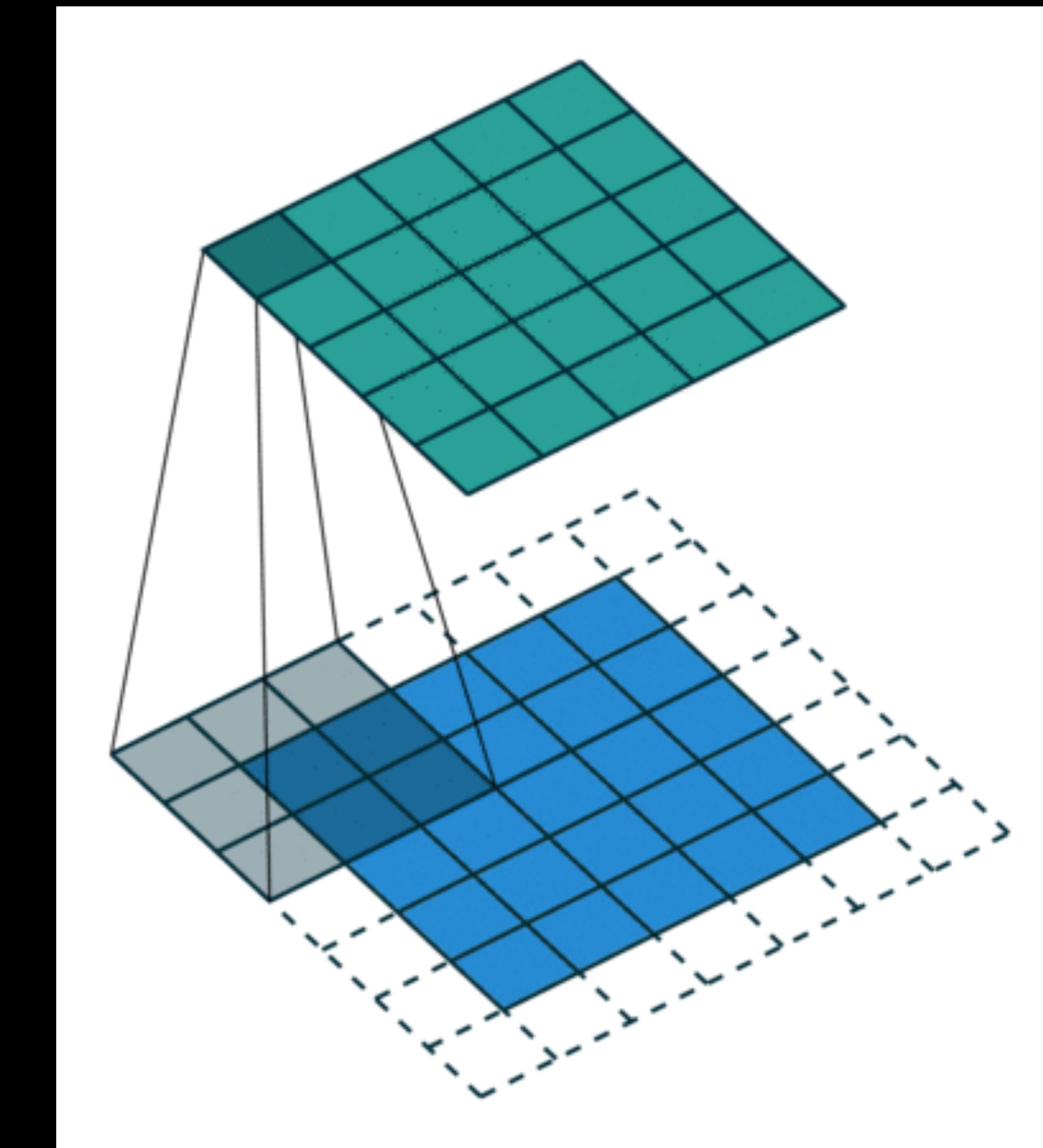
Convolutional Neural Network

- FCNN参数量很大
- 神经元互相影响，一定阶段后难以训练
- CNN使用卷积运算
 - 降低参数量，高效提取图片信息



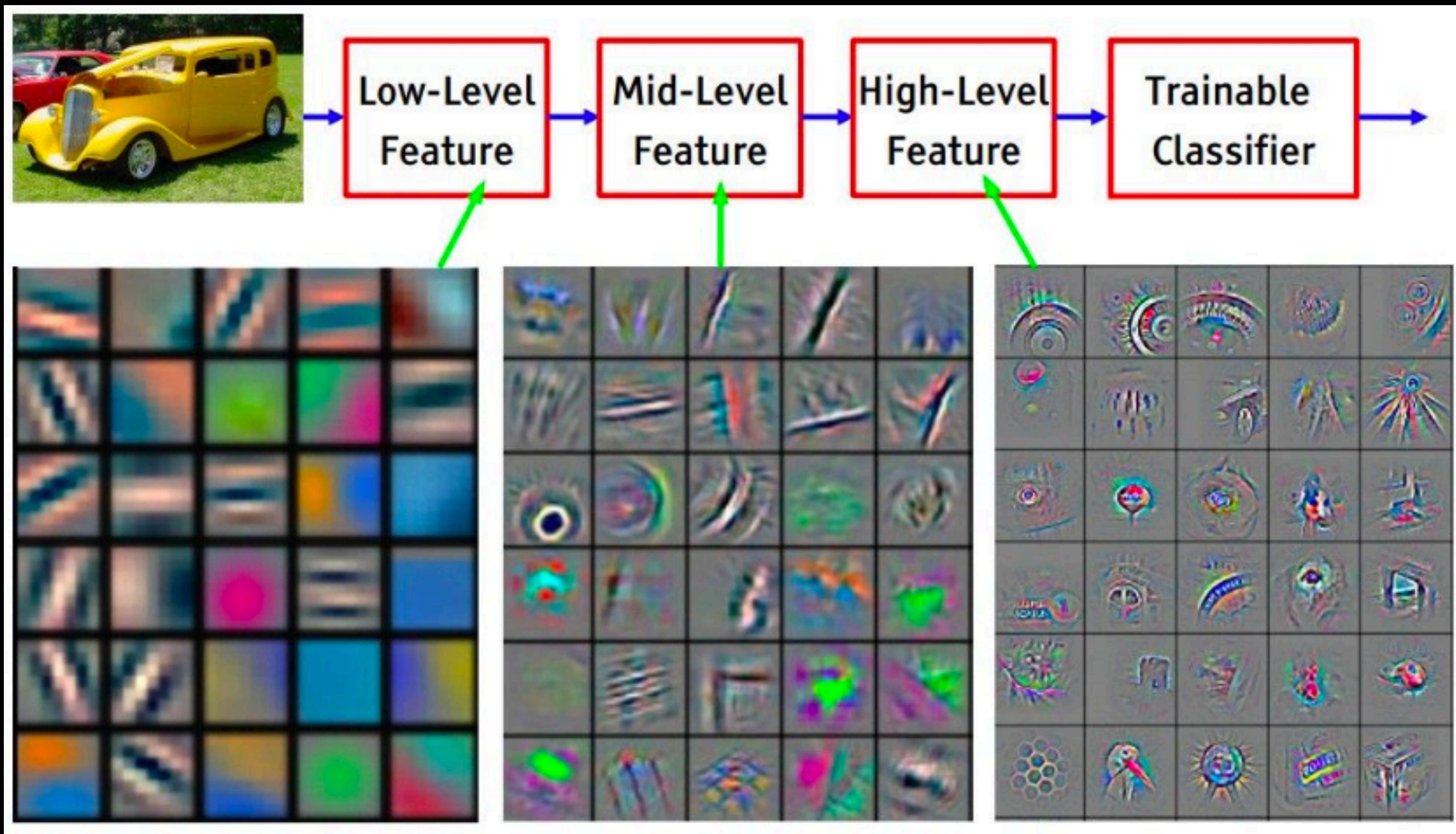
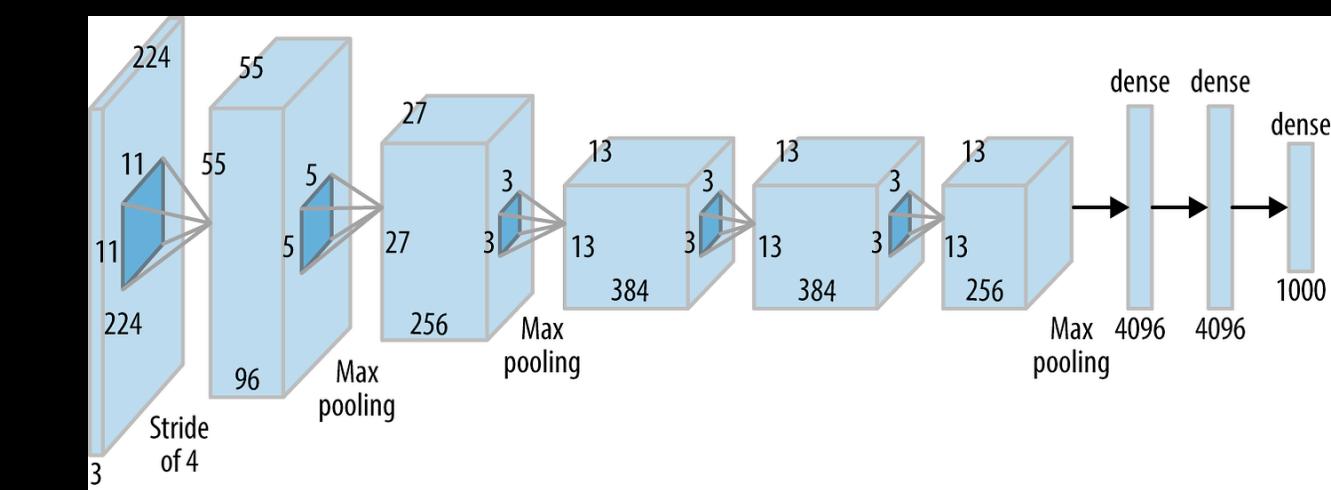
AlexNet 结构

Source : https://www.researchgate.net/figure/Simplified-illustration-of-the-AlexNet-architecture_fig2_329790469



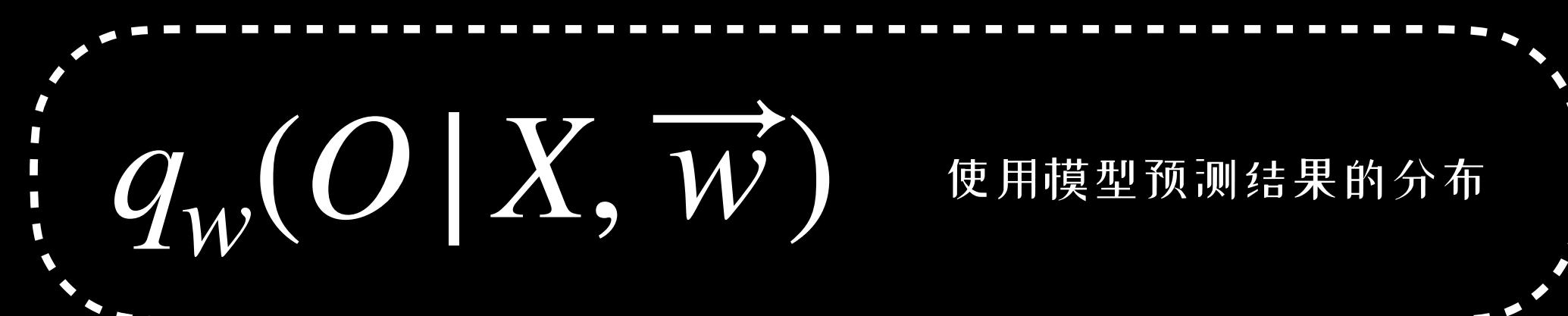
卷积神经网络 (CNN)

Convolutional Neural Network



*统计学视角：训练超高维空间中的概率分布

- 每个输入数据属于一个 m 维空间
 - $\vec{x}_i \in X \subset R^m$
- 每个数据对应的结果也在这个 k 维空间
 - $\vec{y}_i \in Y \subset R^k$
- 模型输出结果也属于一个 k 维空间
 - $\vec{o}_i \in O \subset R^k$
- 模型参数属于一个 n 维空间 $\vec{w} \in R^n$
- 问题：使用概率论语言（概率分布）描述机器学习（监督学习）模型？



$$\begin{aligned} & \arg \min_w D_{KL}(q_w || p_t) \\ & \approx \arg \min_w \sum_{i=1}^N L_w(f_w(\vec{x}_i), \vec{y}_i) \end{aligned}$$

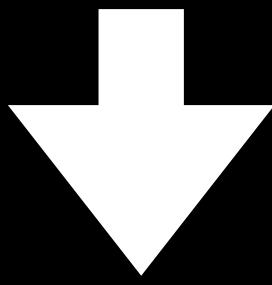
KL Divergent : 衡量两个分布不不同

- 核心思想：只要维度够高，通过概率分布组合可以得到任意复杂的分类器
- 观点：记忆的过程是在高维空间中的记录/索取“高维向量”的能力

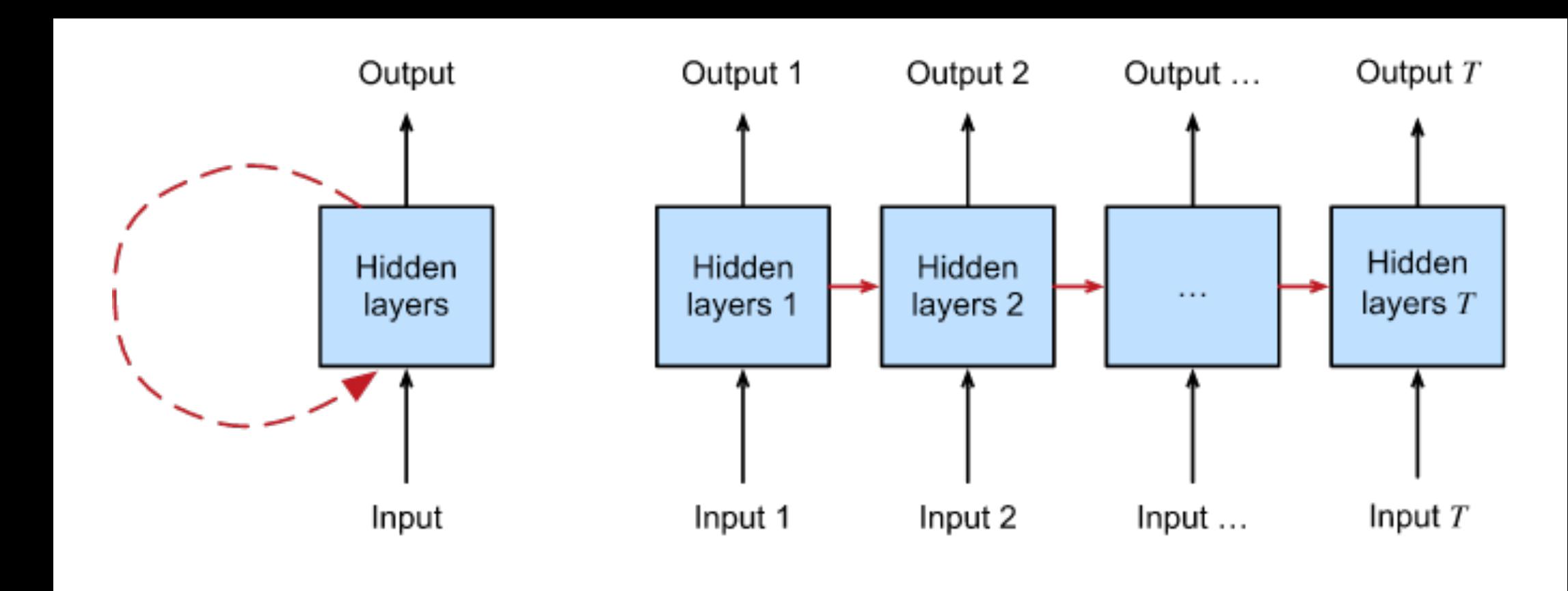
Recurrent Neural Network

循环神经网络

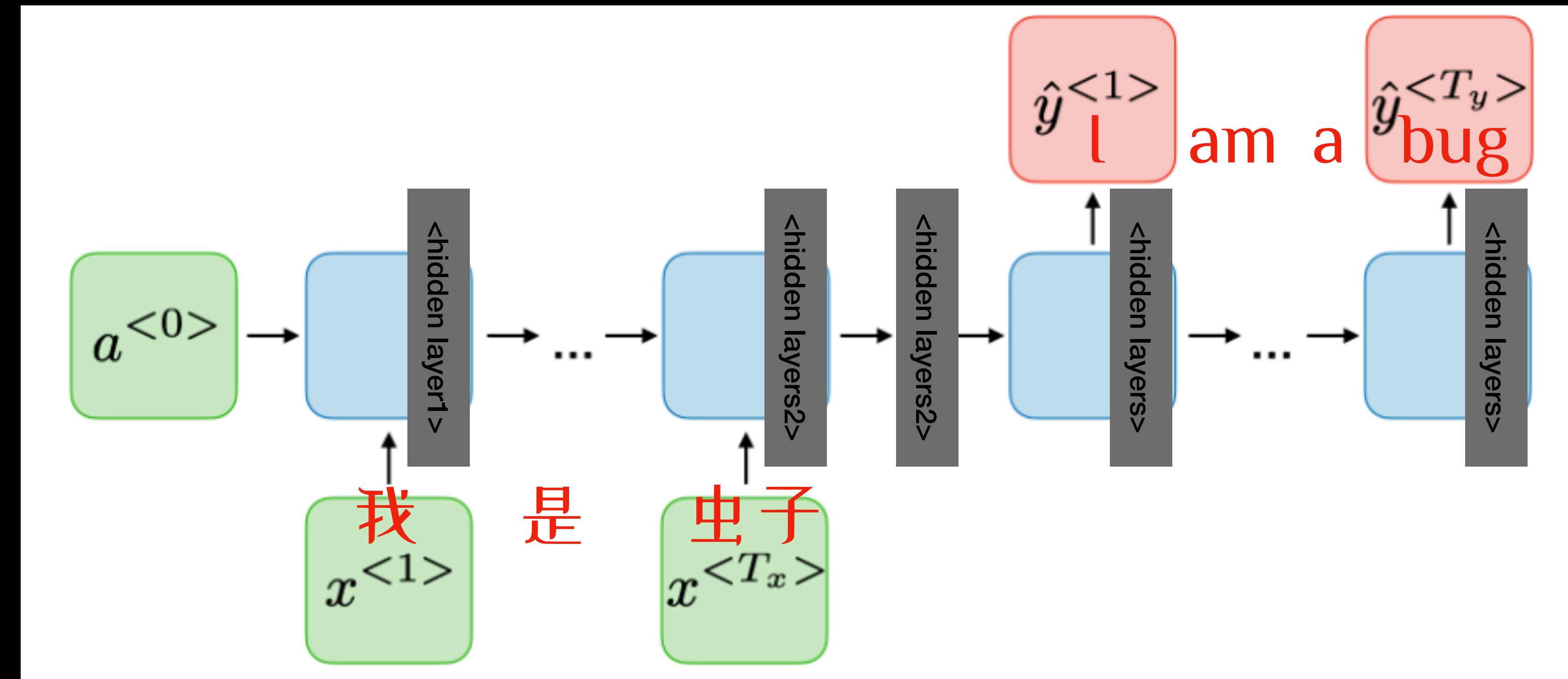
我是虫子.



I am a bug.

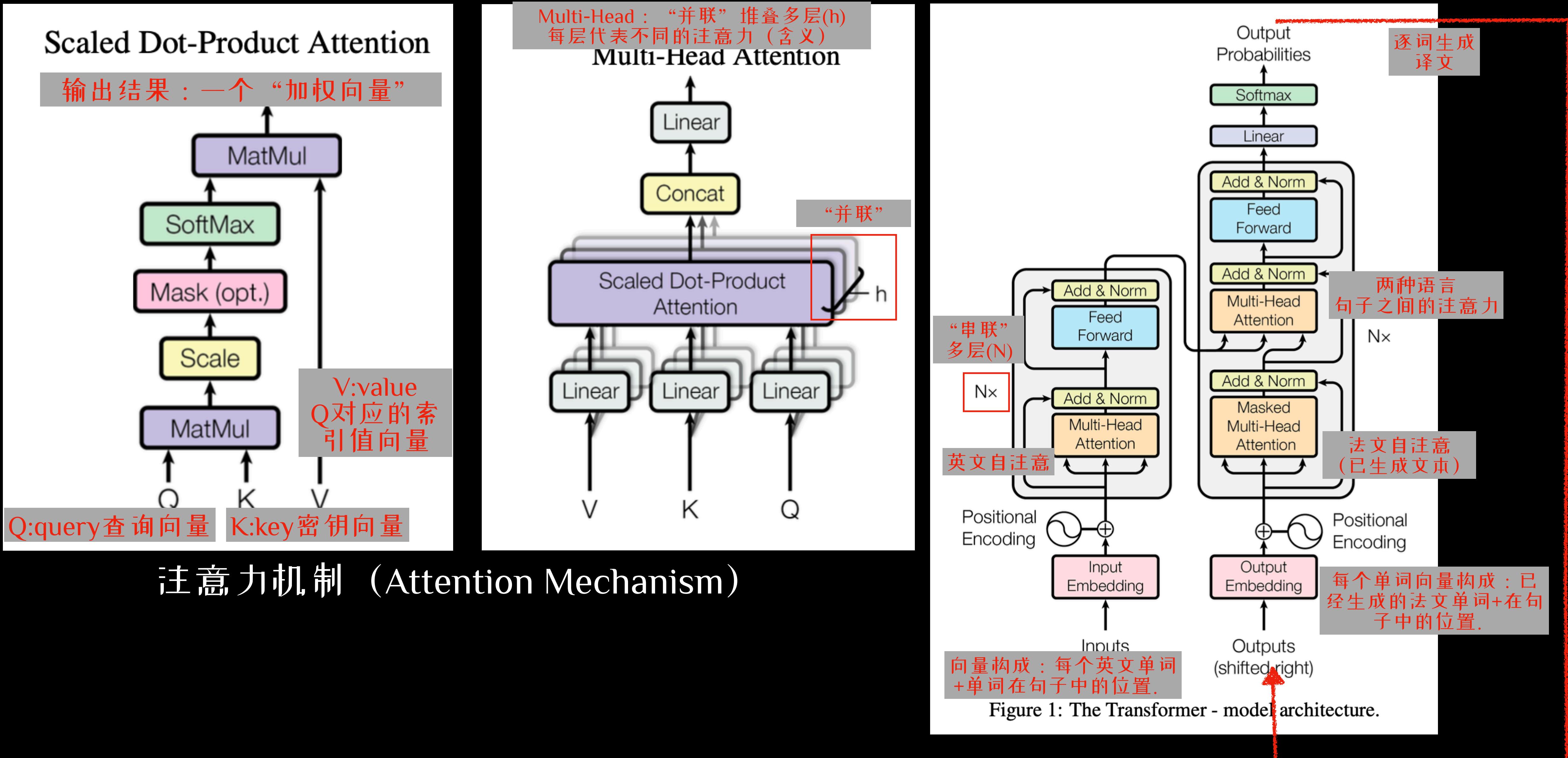


循环神经网络

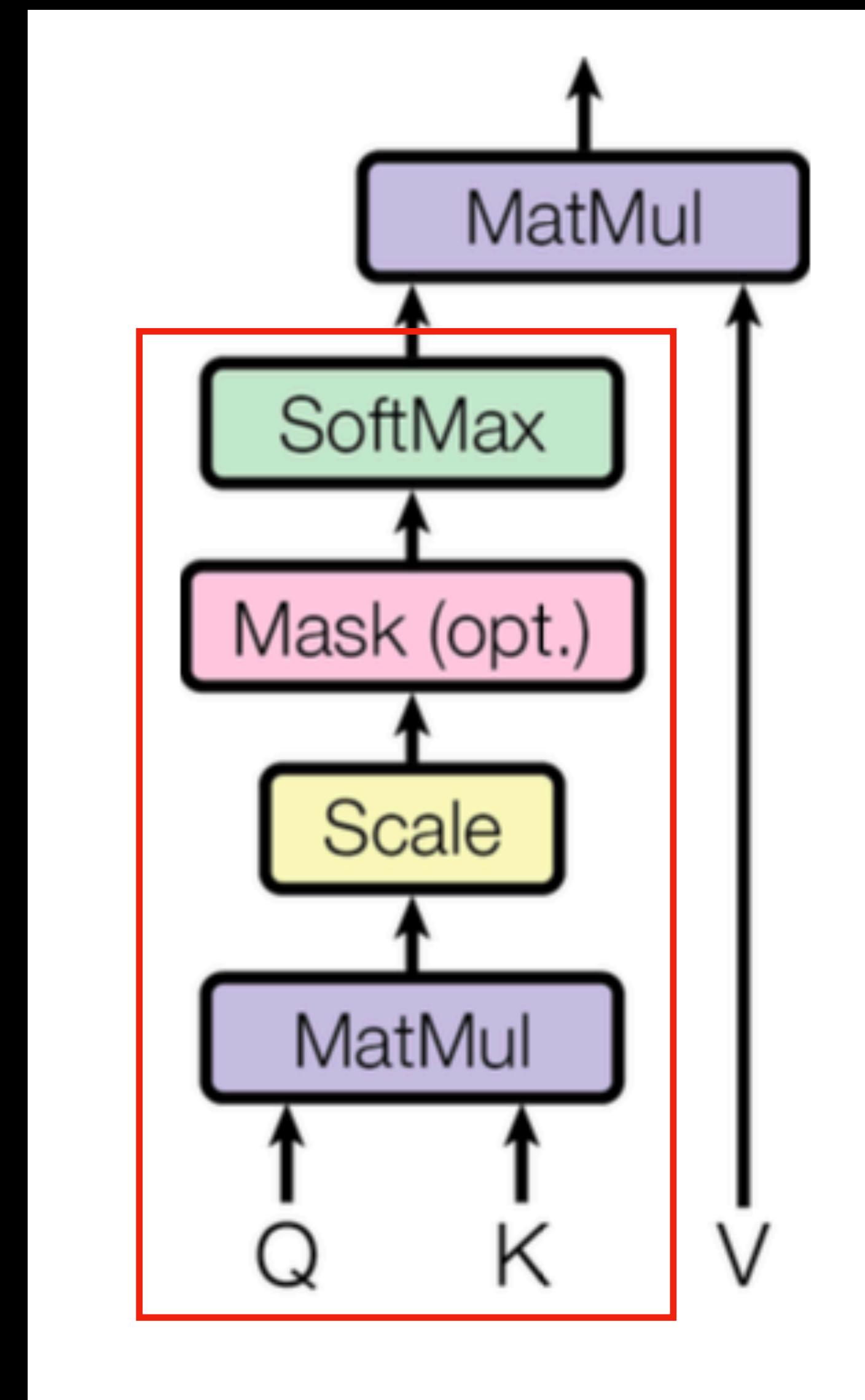


Many-to-Many结构 (用作语言翻译/对话)

注意力机制：《Attention is All You Need》



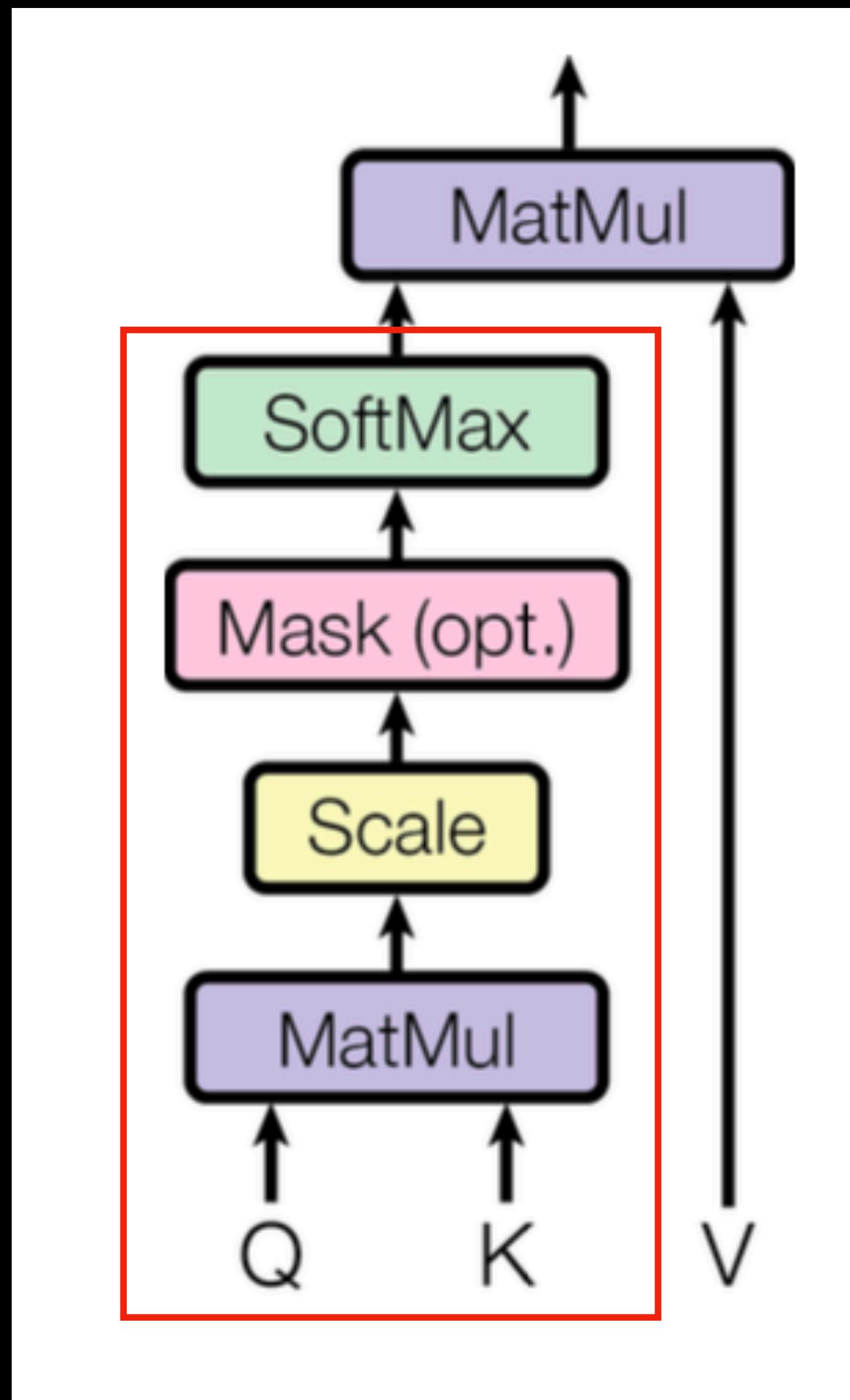
自注意力（可能的主谓关系）



可能的词义注意力(概率)

$\overrightarrow{Q_3}$ 田	0.01	0.05	0.94
$\overrightarrow{Q_2}$ 昆	0.00	0.95	0.05
$\overrightarrow{Q_1}$ 狗	0.99	0.00	0.01

$\overrightarrow{K_1}$ l $\overrightarrow{K_2}$ am $\overrightarrow{K_3}$ bug



Transformer论文中模型展现的一些注意力模型

Figure 5: Many of the attention heads exhibit behaviour that seems related to the structure of the sentence. We give two such examples above, from two different heads from the encoder self-attention at layer 5 of 6. The heads clearly learned to perform different tasks.

