



深度学习在CTR预估的应用

张俊林

新浪微博 AI Lab 资深算法专家

TABLE OF CONTENTES

当深度学习遇到CTR预估

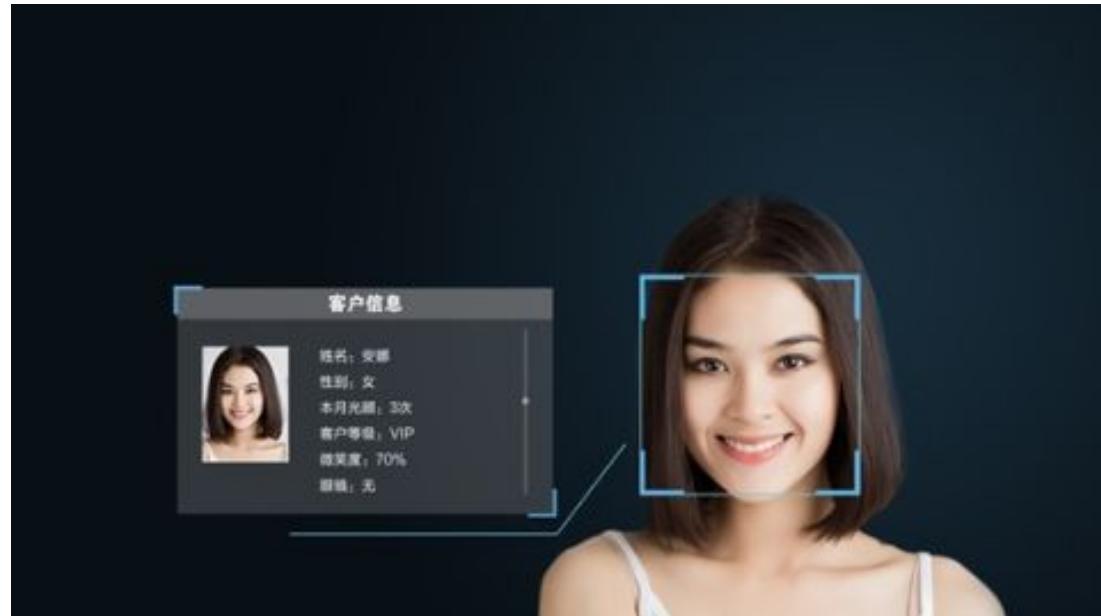
传统主流CTR预估方法

深度学习基础模型

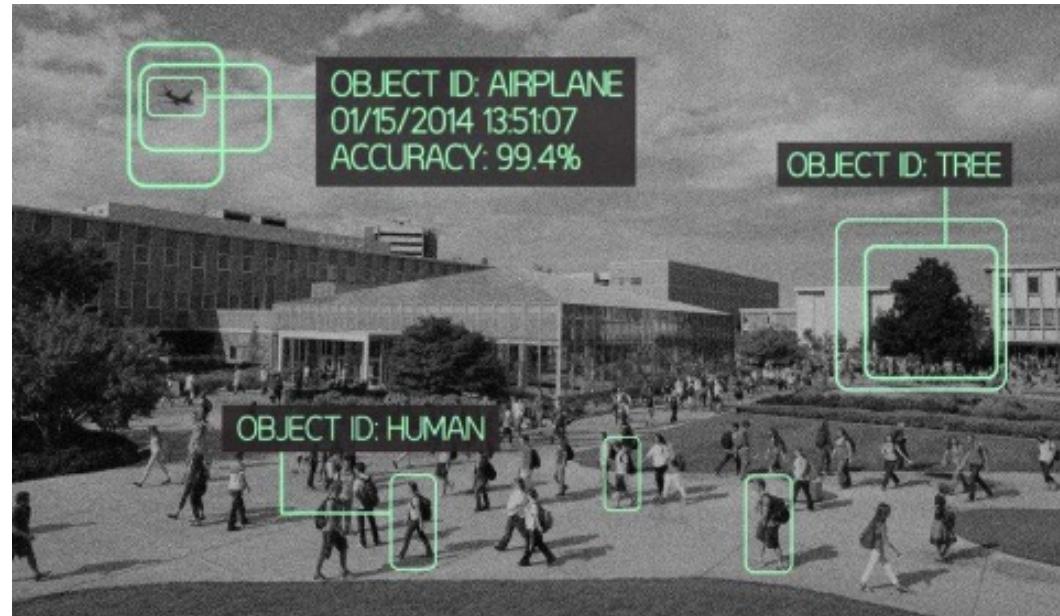
深度学习CTR预估模型

互联网公司深度学习CTR案例

深度学习:各个领域的成功



人脸识别



物体识别



语音识别



机器翻译



风格转换



图片生成

CTR任务的应用



计算广告

推荐系统

信息流排序

CTR任务例子

	Feature vector \mathbf{x}																Target y			
$\mathbf{x}^{(1)}$	1	0	0	...	1	0	0	0	...	0.3	0.3	0.3	0	...	13	0	0	0	0	...
$\mathbf{x}^{(2)}$	1	0	0	...	0	1	0	0	...	0.3	0.3	0.3	0	...	14	1	0	0	0	...
$\mathbf{x}^{(3)}$	1	0	0	...	0	0	1	0	...	0.3	0.3	0.3	0	...	16	0	1	0	0	...
$\mathbf{x}^{(4)}$	0	1	0	...	0	0	1	0	...	0	0	0.5	0.5	...	5	0	0	0	0	...
$\mathbf{x}^{(5)}$	0	1	0	...	0	0	0	1	...	0	0	0.5	0.5	...	8	0	0	1	0	...
$\mathbf{x}^{(6)}$	0	0	1	...	1	0	0	0	...	0.5	0	0.5	0	...	9	0	0	0	0	...
$\mathbf{x}^{(7)}$	0	0	1	...	0	0	1	0	...	0.5	0	0.5	0	...	12	1	0	0	0	...
	A	B	C	...	TI	NH	SW	ST	...	TI	NH	SW	ST	...	Time	TI	NH	SW	ST	...
	User				Movie					Other Movies rated					Last Movie rated					

CTR任务的特点

- 大量离散特征
- 大量高维度稀疏特征
- 特征工程：特征组合对于效果非常关键

当CTR预估遇到深度学习

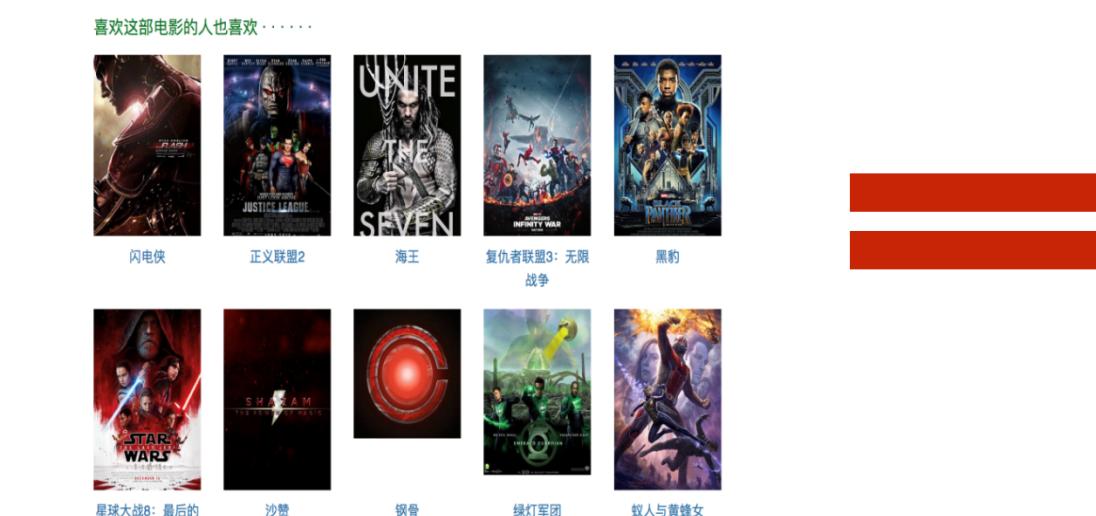
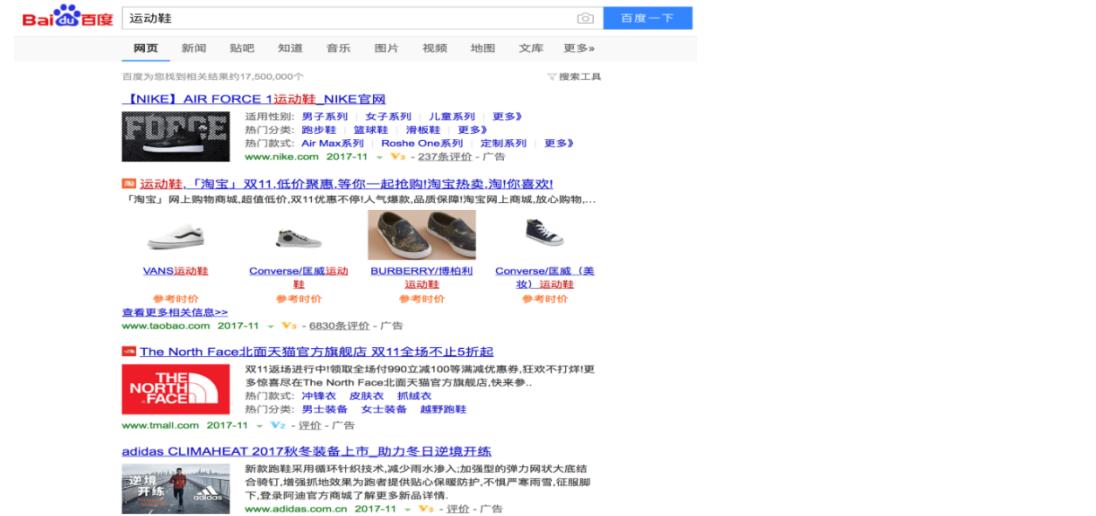
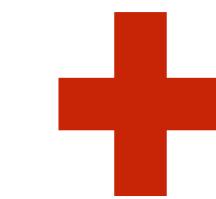
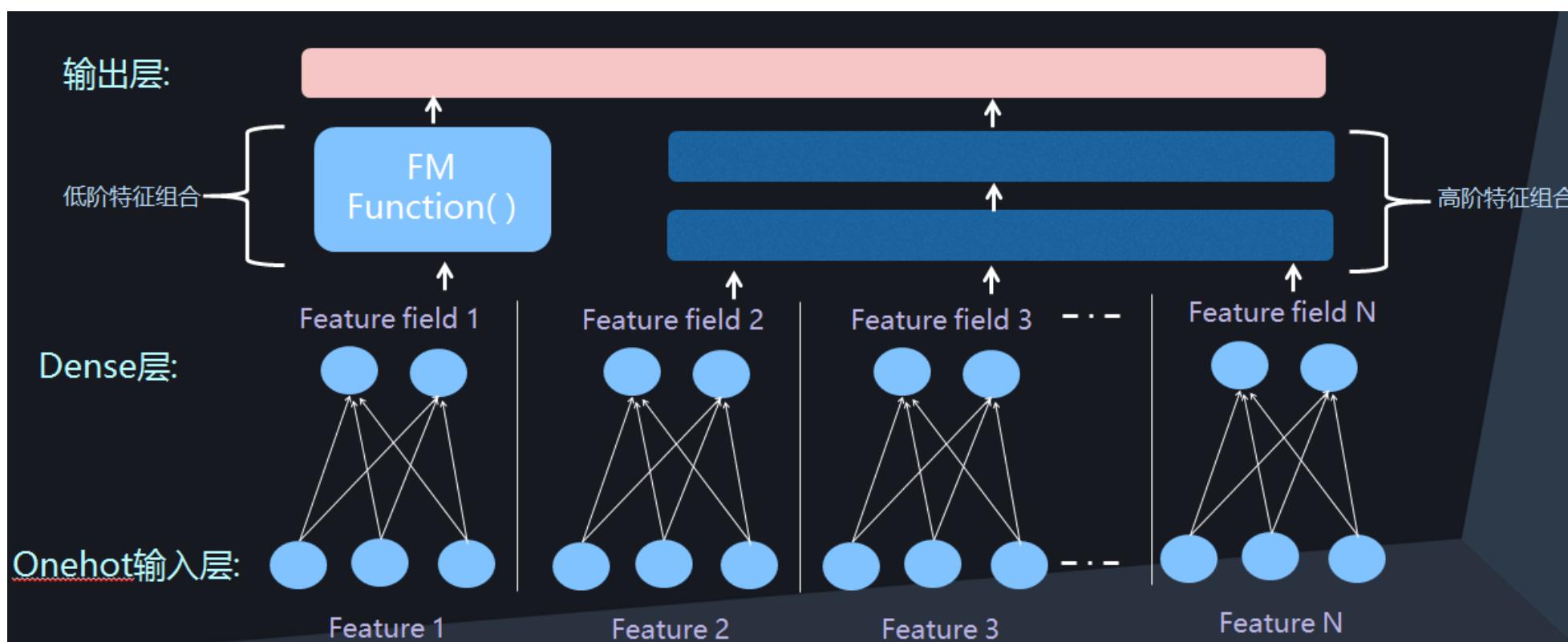


TABLE OF CONTENTES

当深度学习遇到CTR预估

传统主流CTR预估方法

深度学习基础模型

深度学习CTR预估模型

互联网公司深度学习CTR案例

线性模型：思路及问题

Linear: $\hat{y}(x) := w_0 + \sum_{i=1}^n w_i x_i$

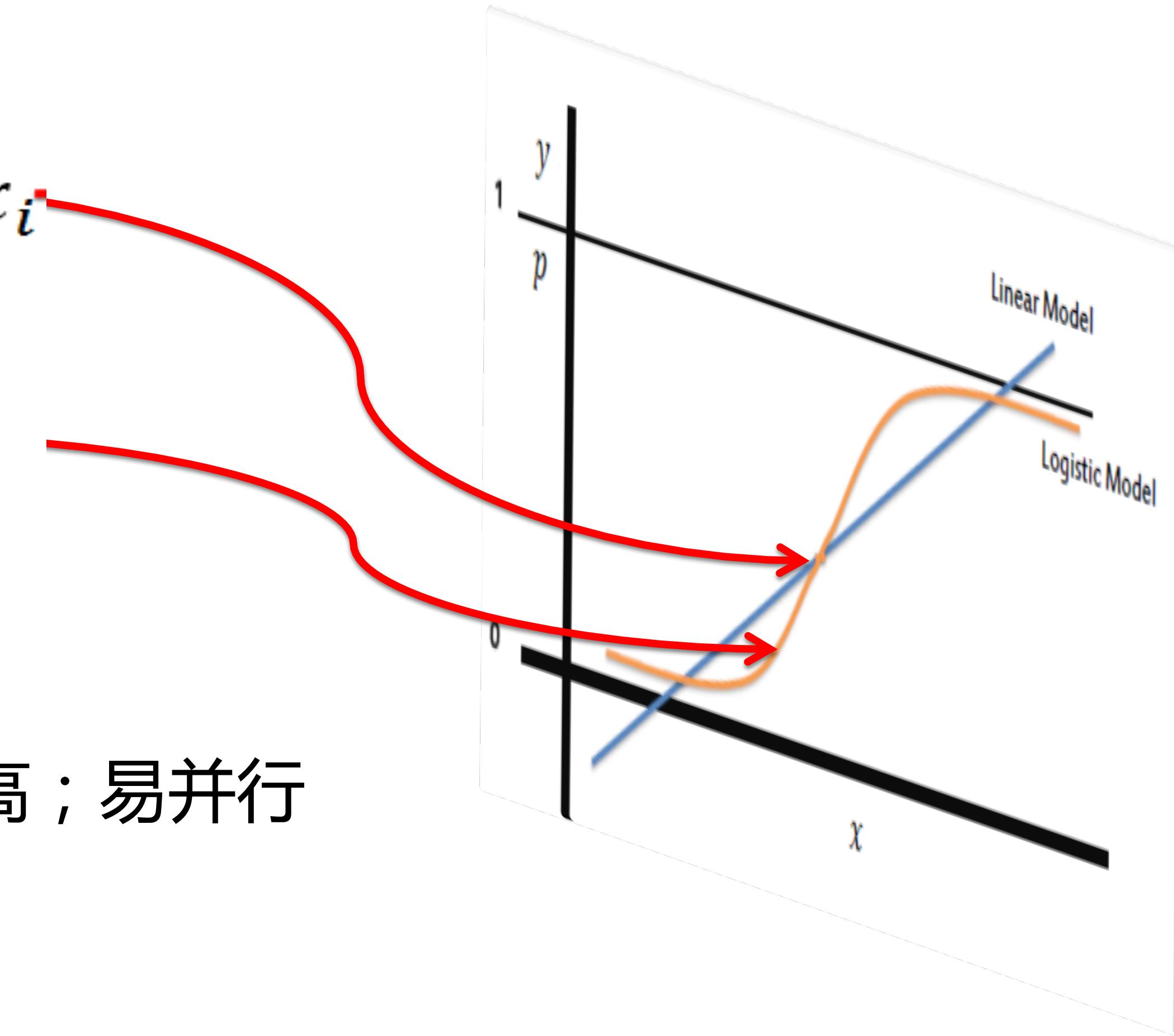
LR: $\hat{y}(x) = \frac{1}{1 + w_0 \exp(-w^T x)}$

优势：

简单；可解释；易扩展；效率高；易并行

缺点：

难以捕获特征组合



线性模型改进：加入特征组合

$$\text{改进版本 } \hat{y}(x) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n w_{i,j} x_i x_j$$

两两特征组合

优势：

直接将两两组合特征引入模型

缺点：

组合特征泛化能力弱

$w_{i,j} = 0 \quad if \text{ 在训练数据中 } x_i x_j = 0$

FM模型

$$\text{FM: } \hat{y}(x) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle v_i, v_j \rangle x_i x_j$$



$$\langle v_i, v_j \rangle := \sum_{f=1}^k v_{i,f} \cdot v_{j,f} =$$

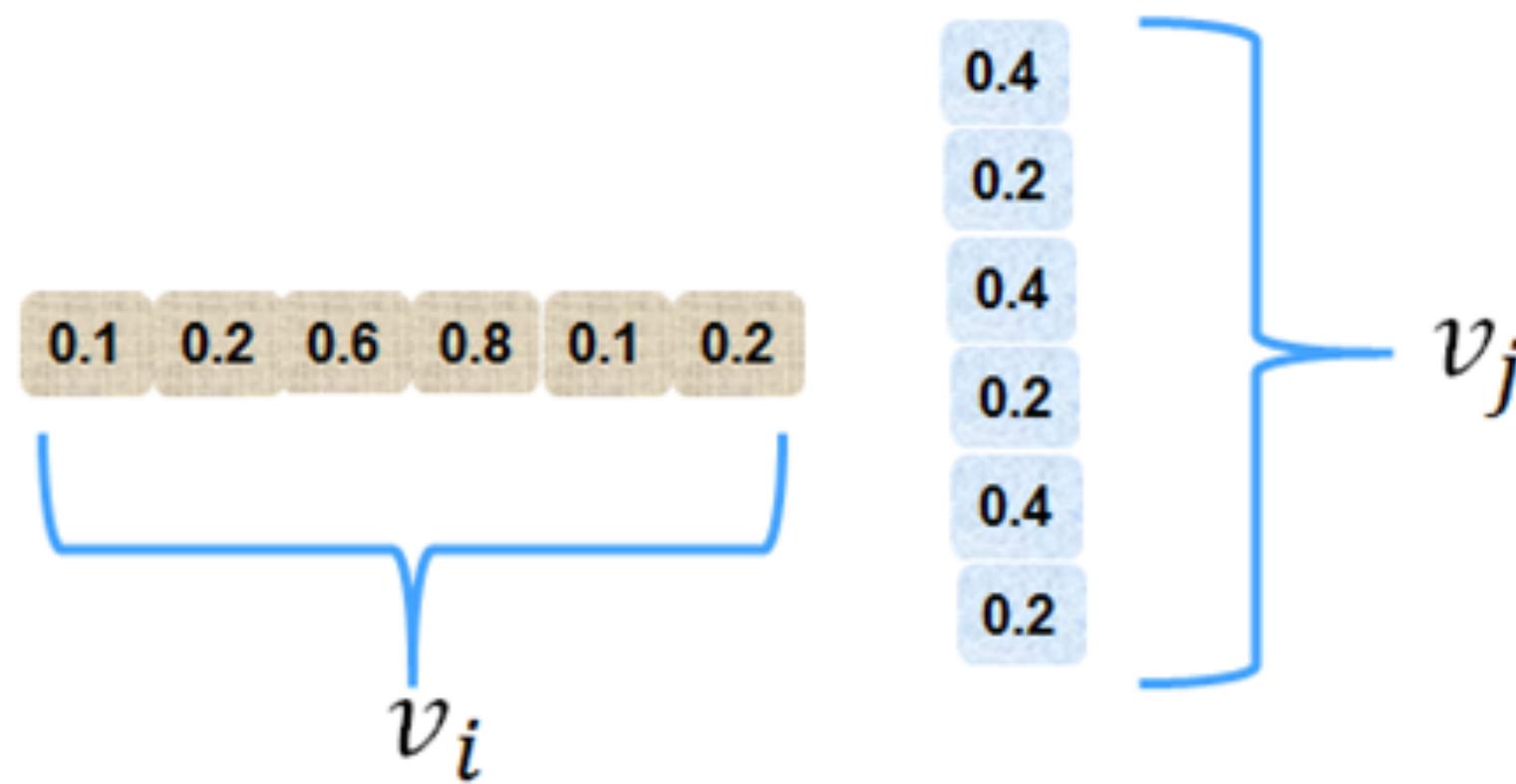
0.1	0.2	0.6	0.8	0.1	0.2
-----	-----	-----	-----	-----	-----

v_i

v_j

FM模型

$$\langle v_i, v_j \rangle := \sum_{f=1}^k v_{i,f} \cdot v_{j,f} =$$



v_1	0.3	0.2	0.6	0.8	0.1	0.2
v_2	0.1	0.8	0.6	0.8	0.4	0.6
v_3	0.4	0.2	0.7	0.2	0.1	0.2
v_4	0.1	0.2	0.6	0.8	0.5	0.2
	→					
v_{n-1}	0.3	0.2	0.6	0.8	0.1	0.2
v_n	0.5	0.8	0.9	0.8	0.4	0.6

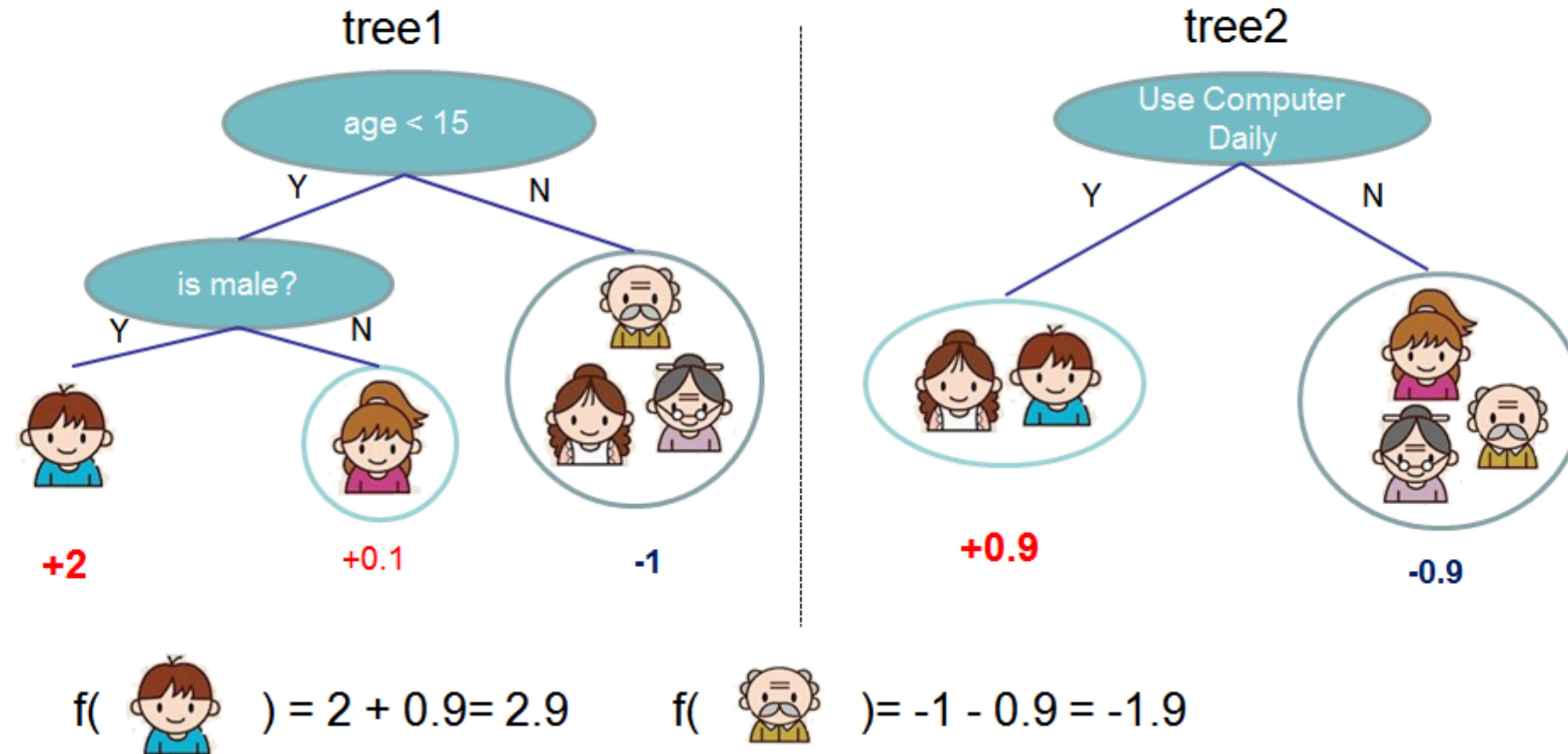
$w_{i,j} = \langle v_i, v_j \rangle \neq 0$
even if 训练数据中 $x_i x_j = 0$
only if 在训练数据中存在 k 使得 $x_i x_k \neq 0$

FM模型泛化能力强

GBDT模型

- Gradient Boosting Decision Tree: 迭代的决策树算法，多棵决策树构成，所有子树的决策值累加得出预测值；
- 很多名称，还被称为：MART (Multiple Additive Regression Tree) / GBRT (Gradient Boosting Regression Tree) / TreeLink等

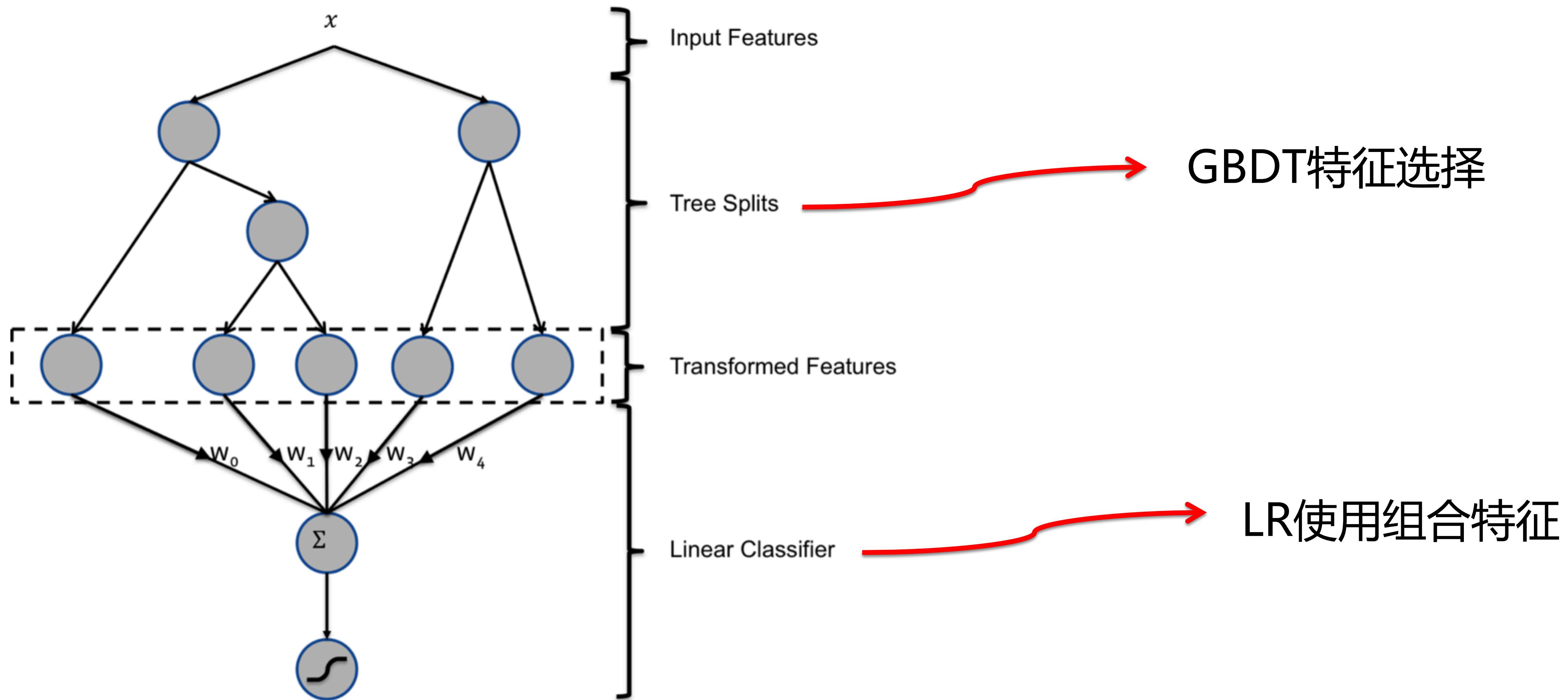
GBDT模型(例：预测是否喜欢打游戏)



LR+GBDT模型

- Facebook首先提出 (2014)
- 集成LR和GBDT各自的优势
 - GBDT发现有效的组合特征Feature Set
 - 将Feature Set引入LR模型中
- 目前广泛使用在各大互联网公司线上系统中

LR+GBDT模型



GBDT+FM模型

- 2014年提出（香港中文+百度）：GBFM模型
- 集成GBDT和FM各自的优势
 - GBDT发现最有效的组合特征Feature Set
 - 将Feature Set引入FM模型中
- GBFM：贪心策略选择组合特征

TABLE OF CONTENTES

当深度学习遇到CTR预估

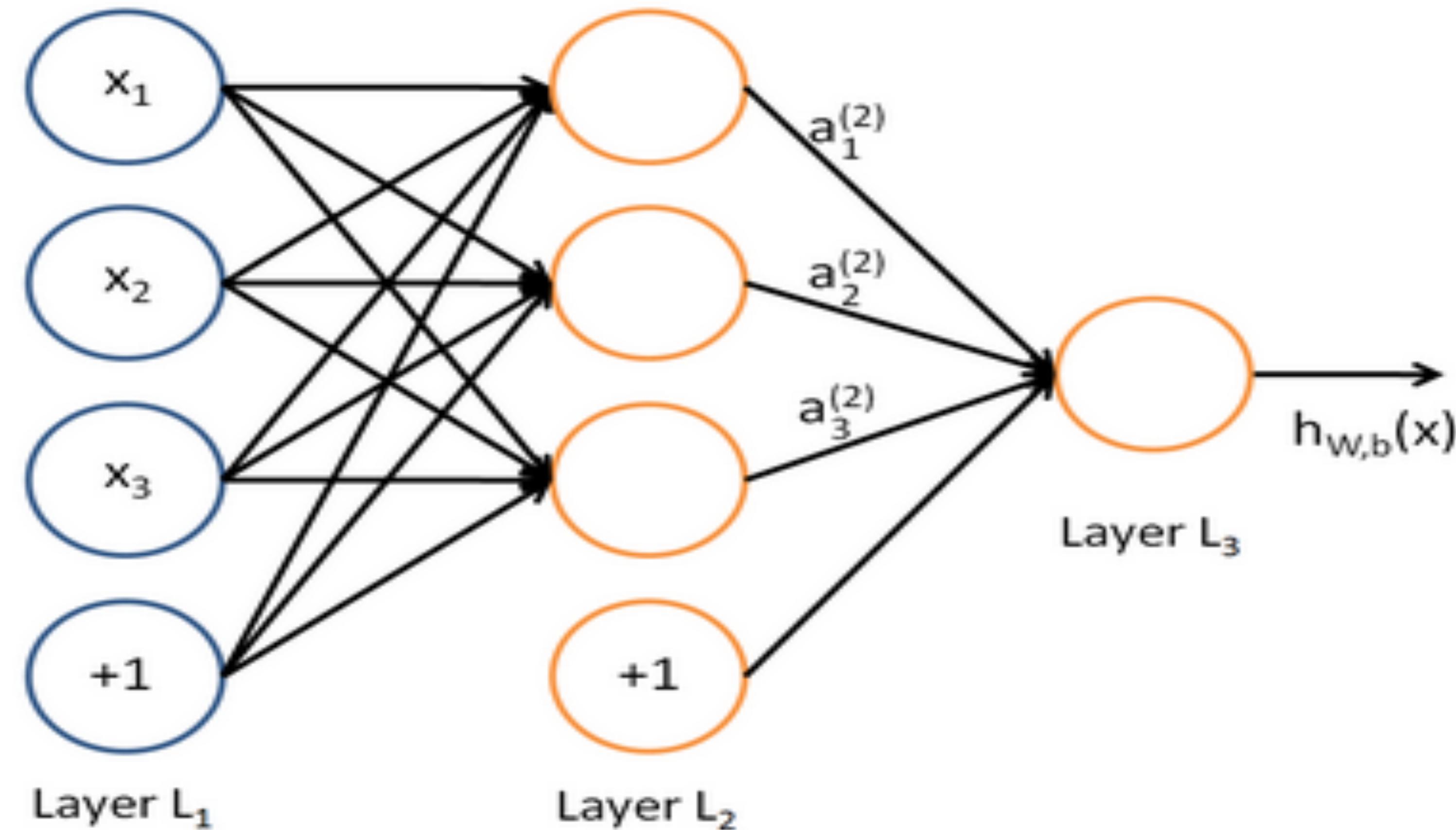
传统主流CTR预估方法

深度学习基础模型

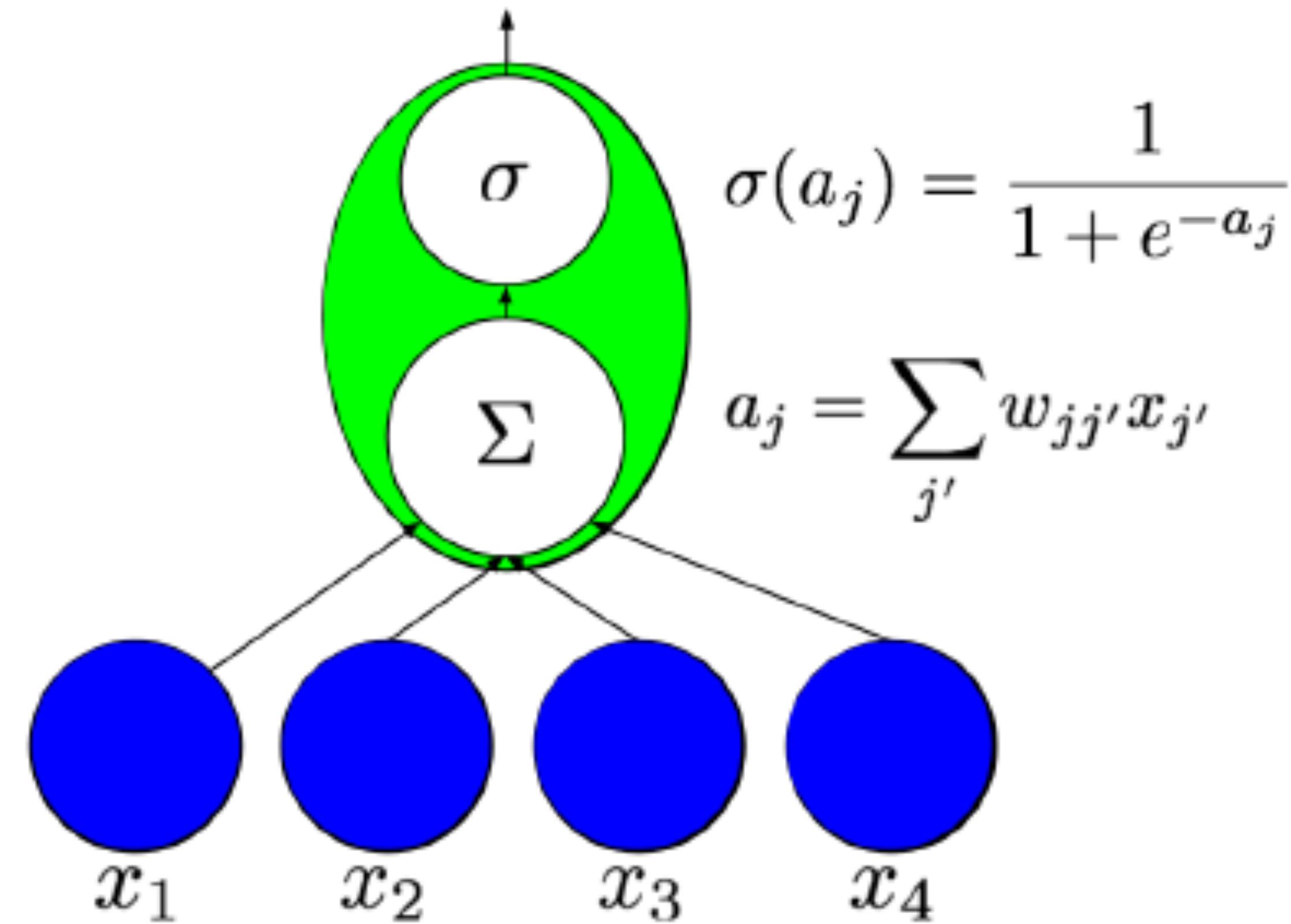
深度学习CTR预估模型

互联网公司深度学习CTR案例

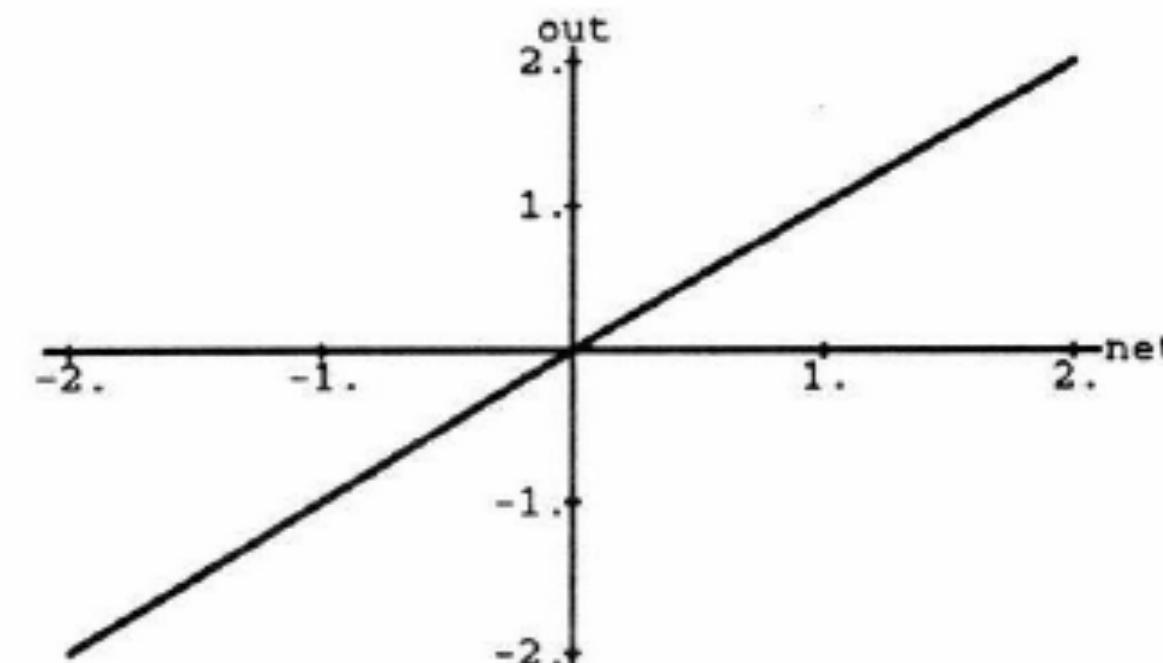
前向神经网络 (MLP)



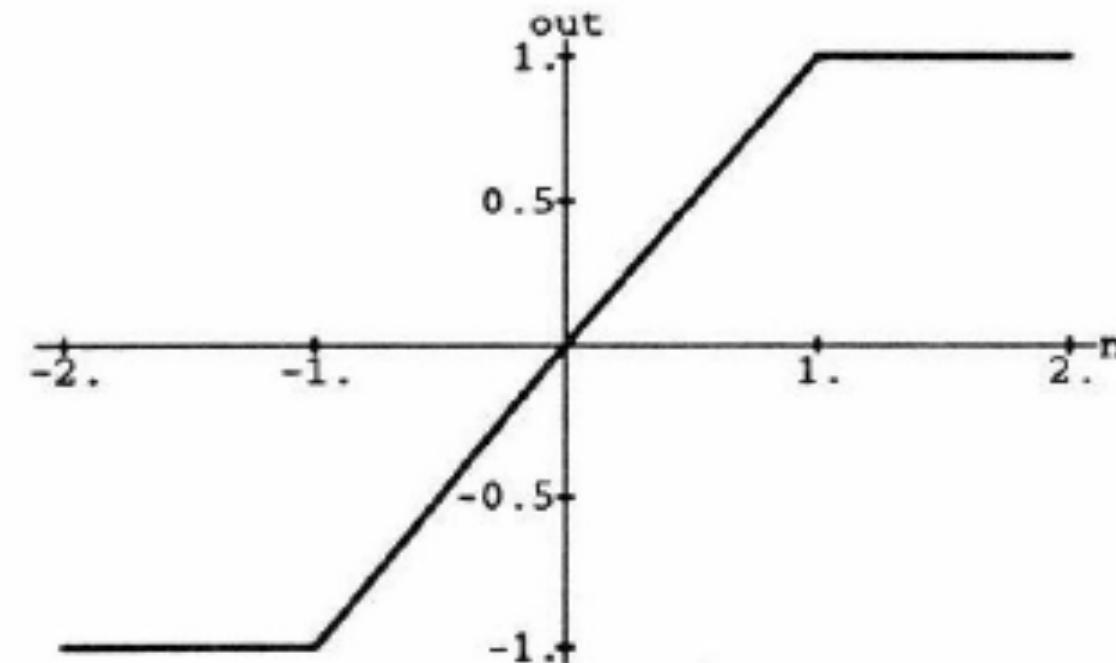
前向神经网络 (MLP) : 隐层节点



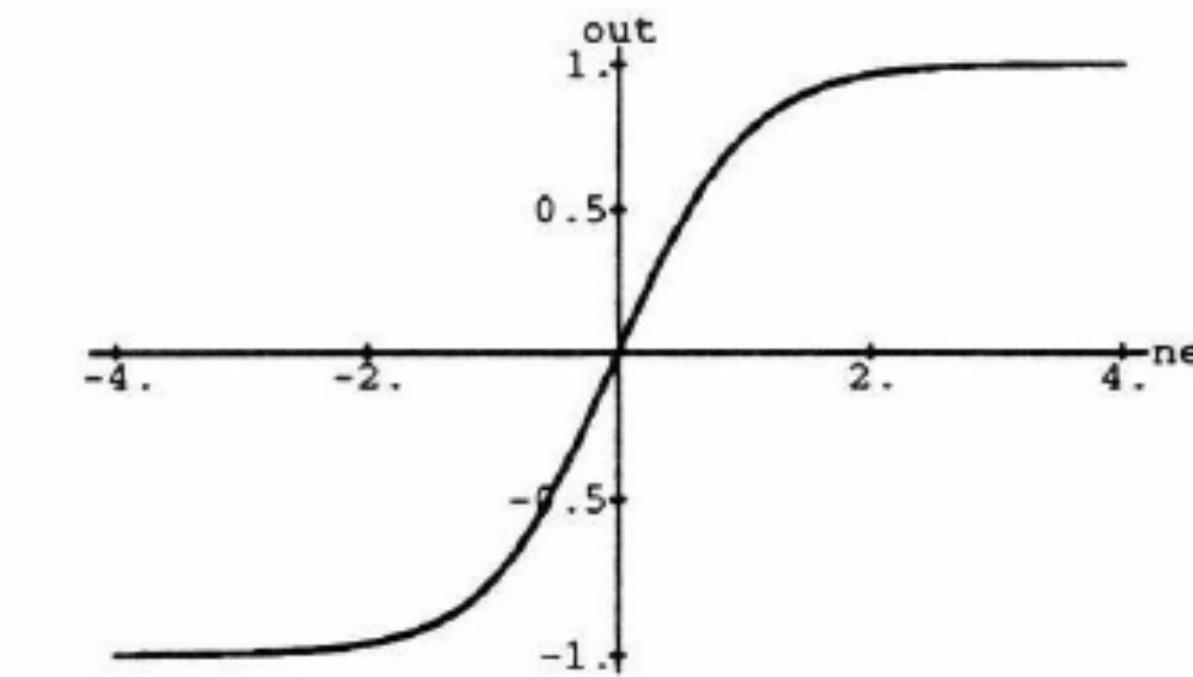
前向神经网络 (MLP) : 激活函数



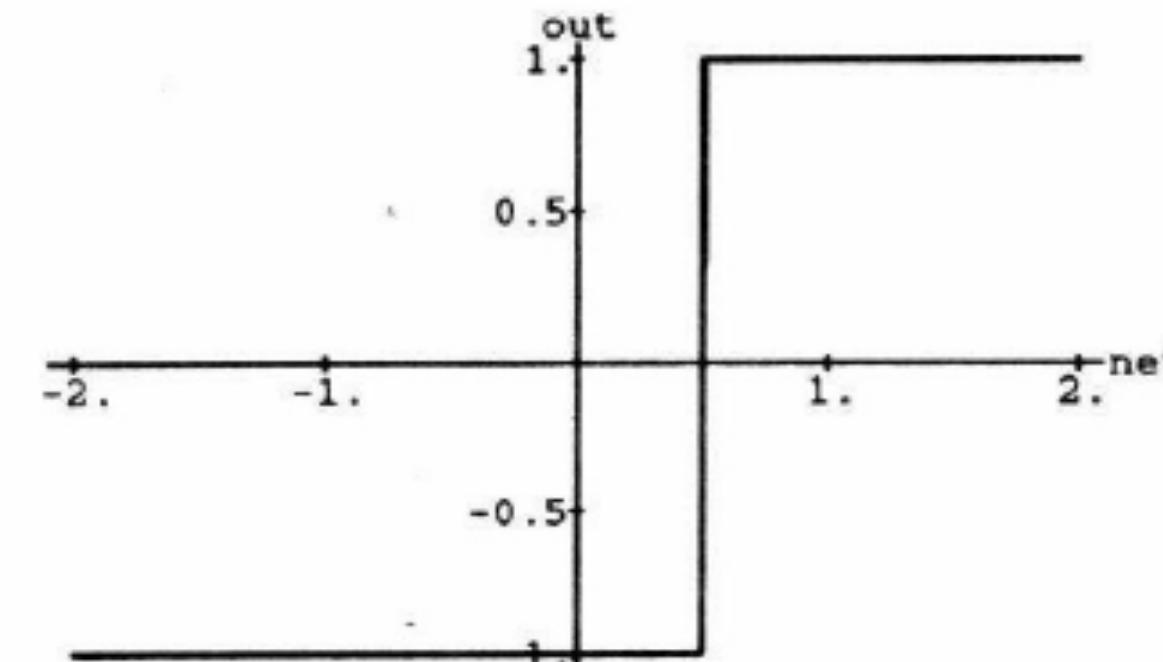
linear



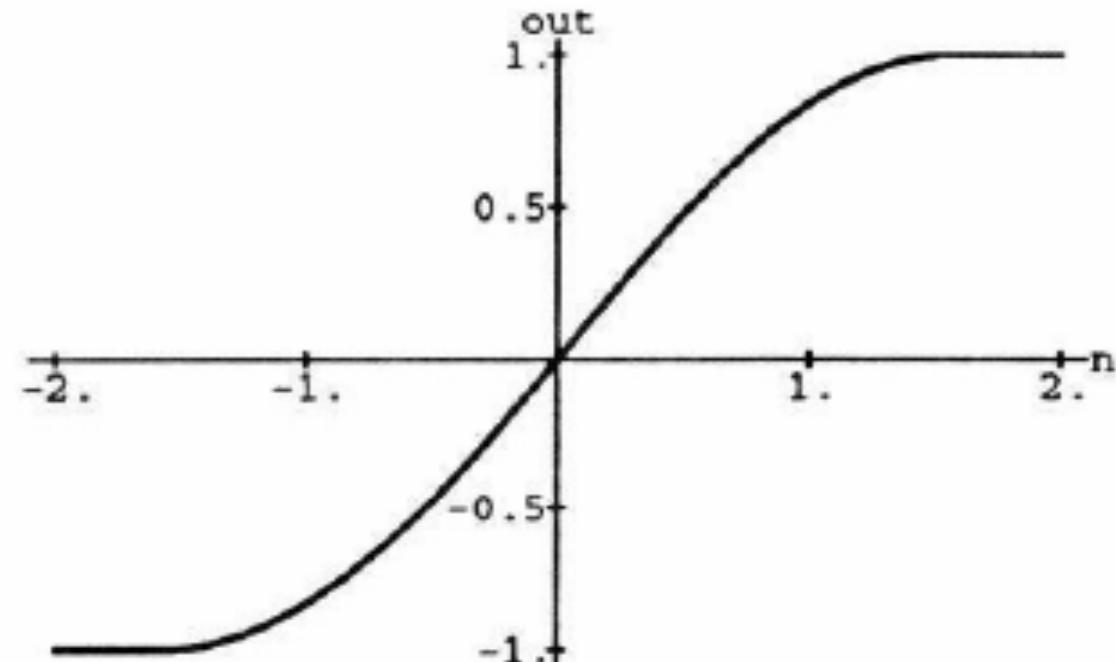
piecewise linear



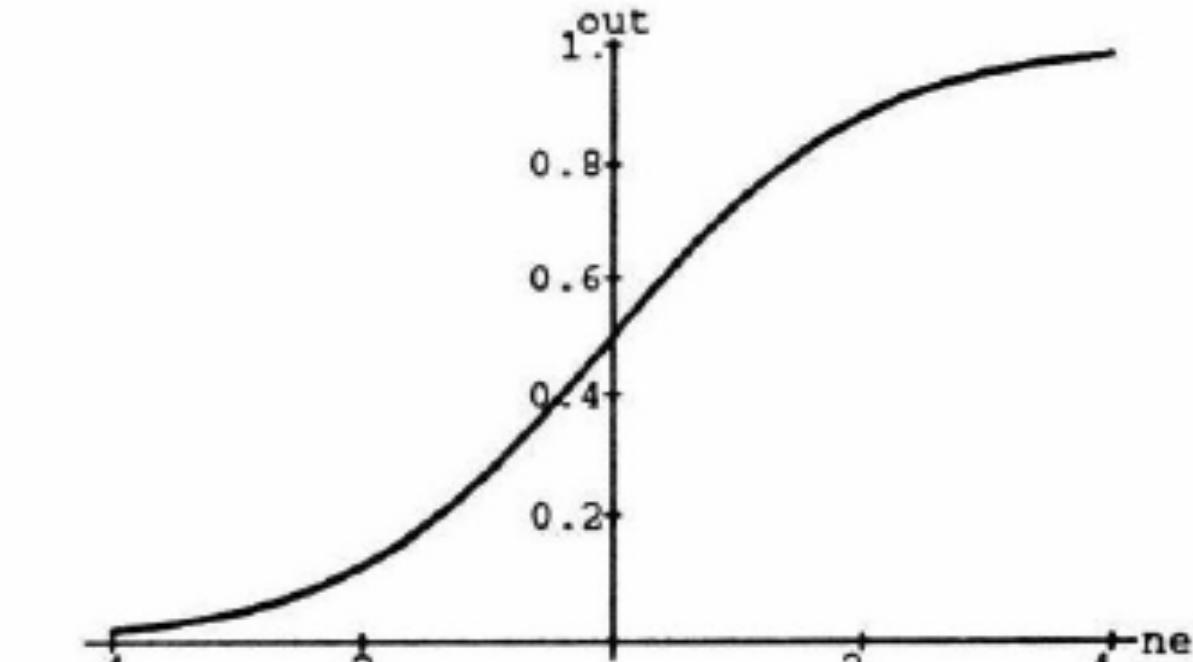
$\tanh(x)$



threshold

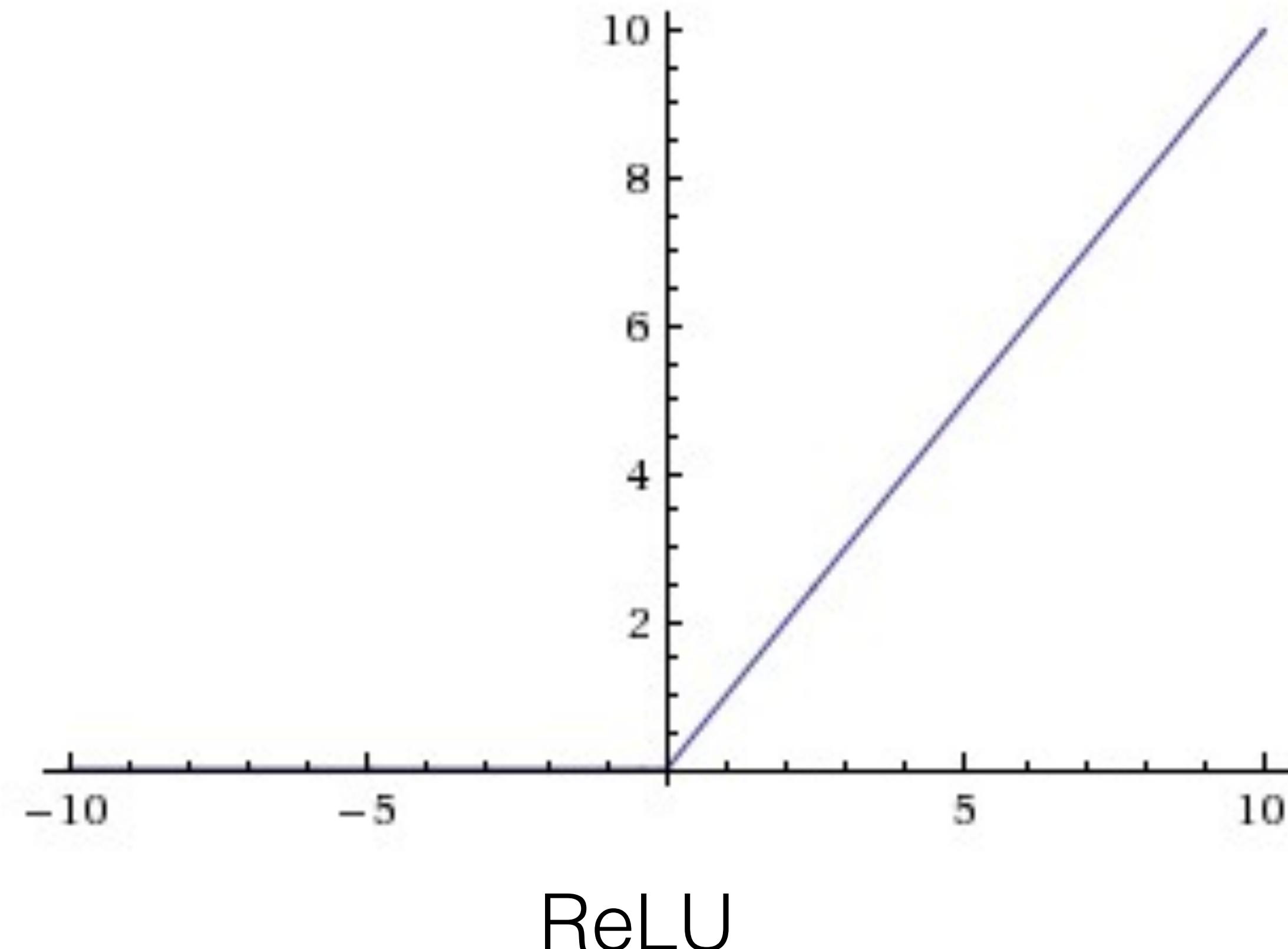


$\sin(x)$ till saturation

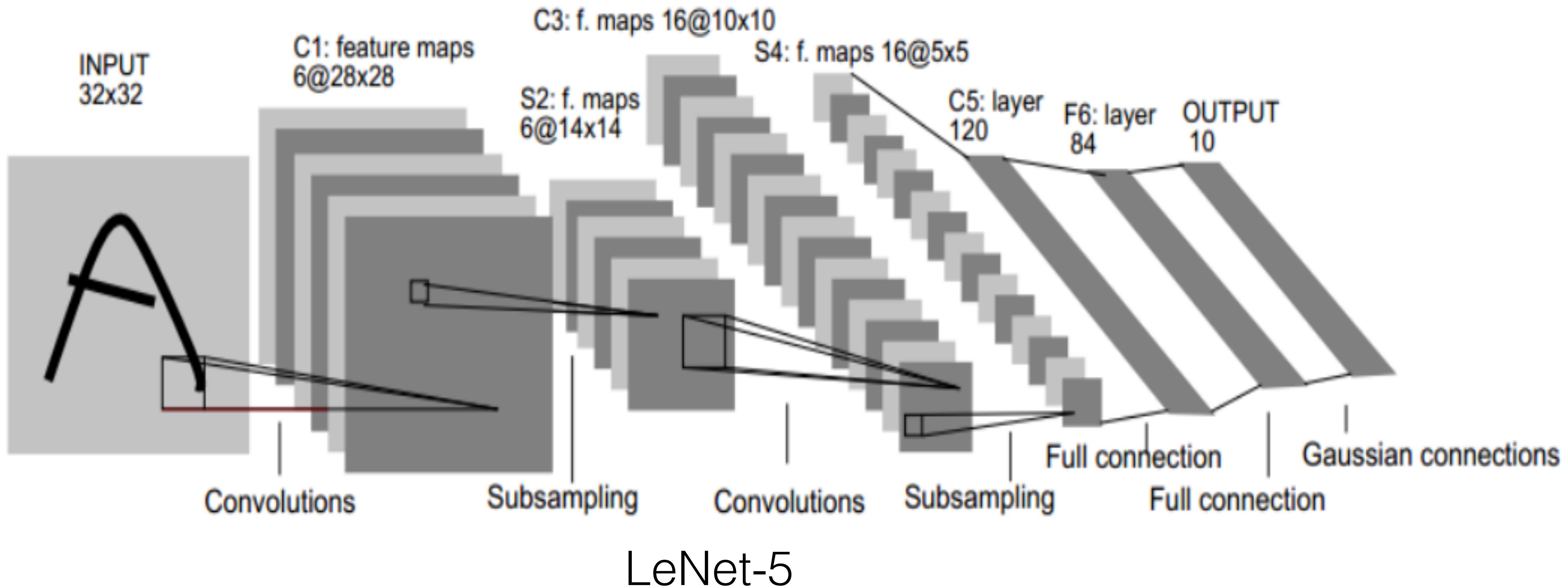


logistic $1/(1+\exp(-x))$

前向神经网络 (MLP) : 激活函数



CNN : LeNet-5整体结构



CNN: 卷积层

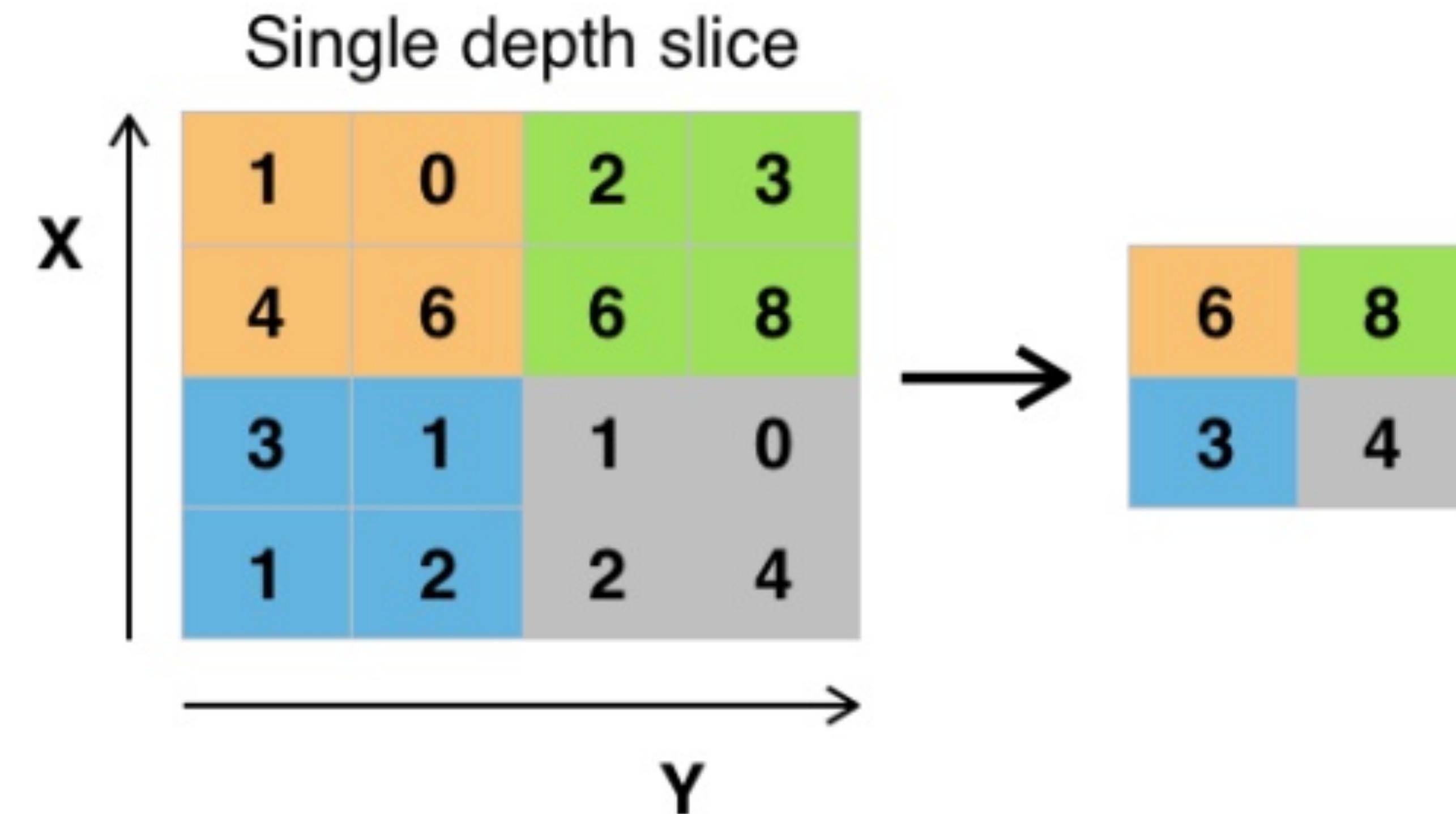
1 <small>$\times 1$</small>	1 <small>$\times 0$</small>	1 <small>$\times 1$</small>	0	0
0 <small>$\times 0$</small>	1 <small>$\times 1$</small>	1 <small>$\times 0$</small>	1	0
0 <small>$\times 1$</small>	0 <small>$\times 0$</small>	1 <small>$\times 1$</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

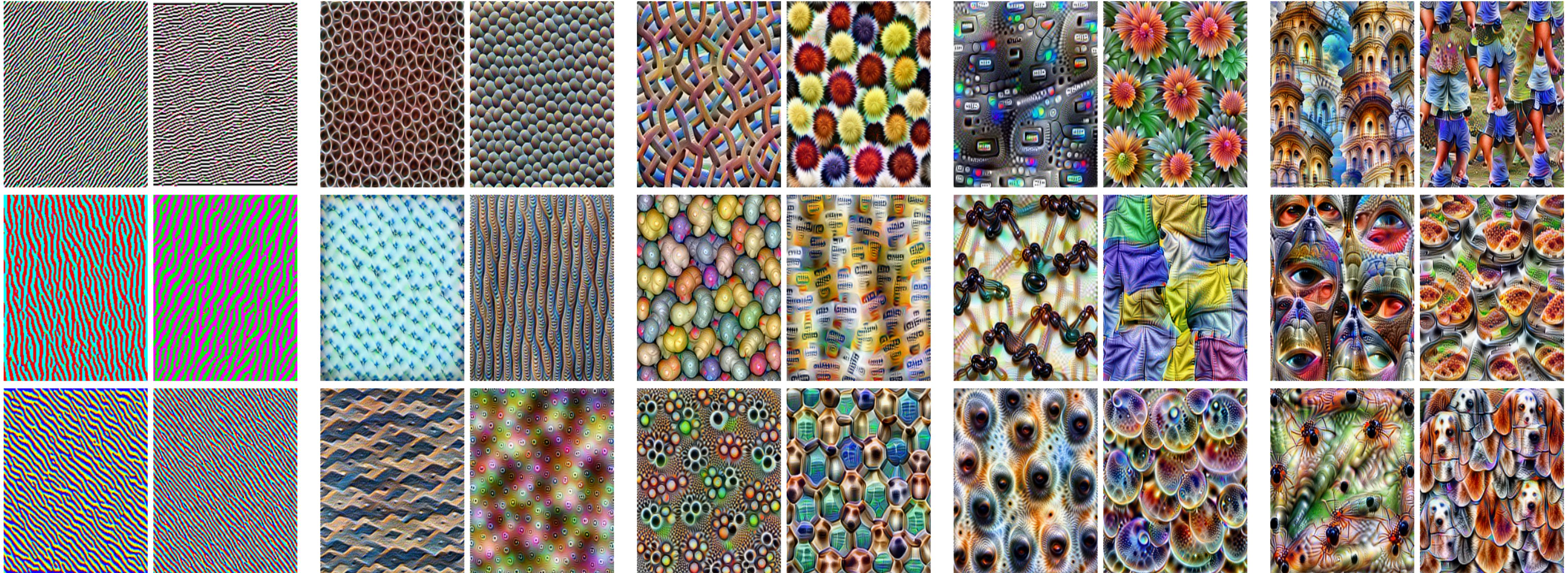
4		

Convolved Feature

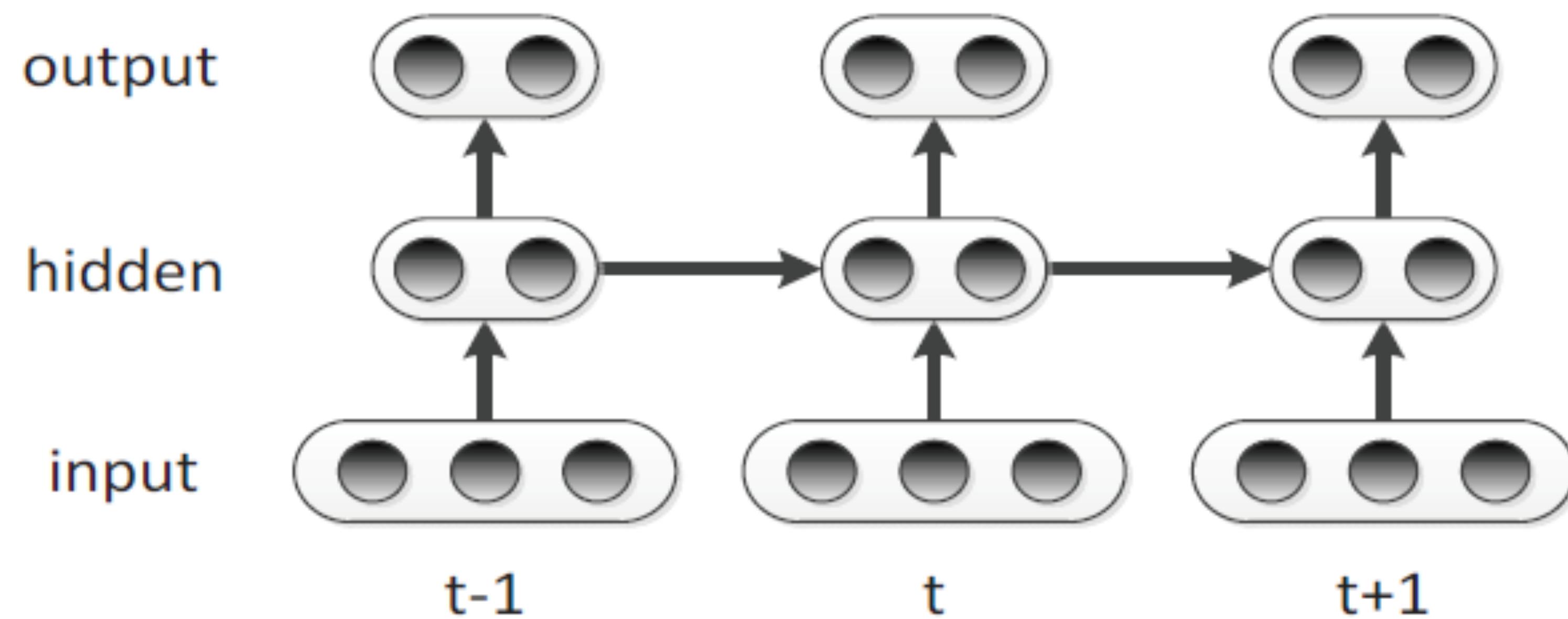
CNN:MaxPooling层



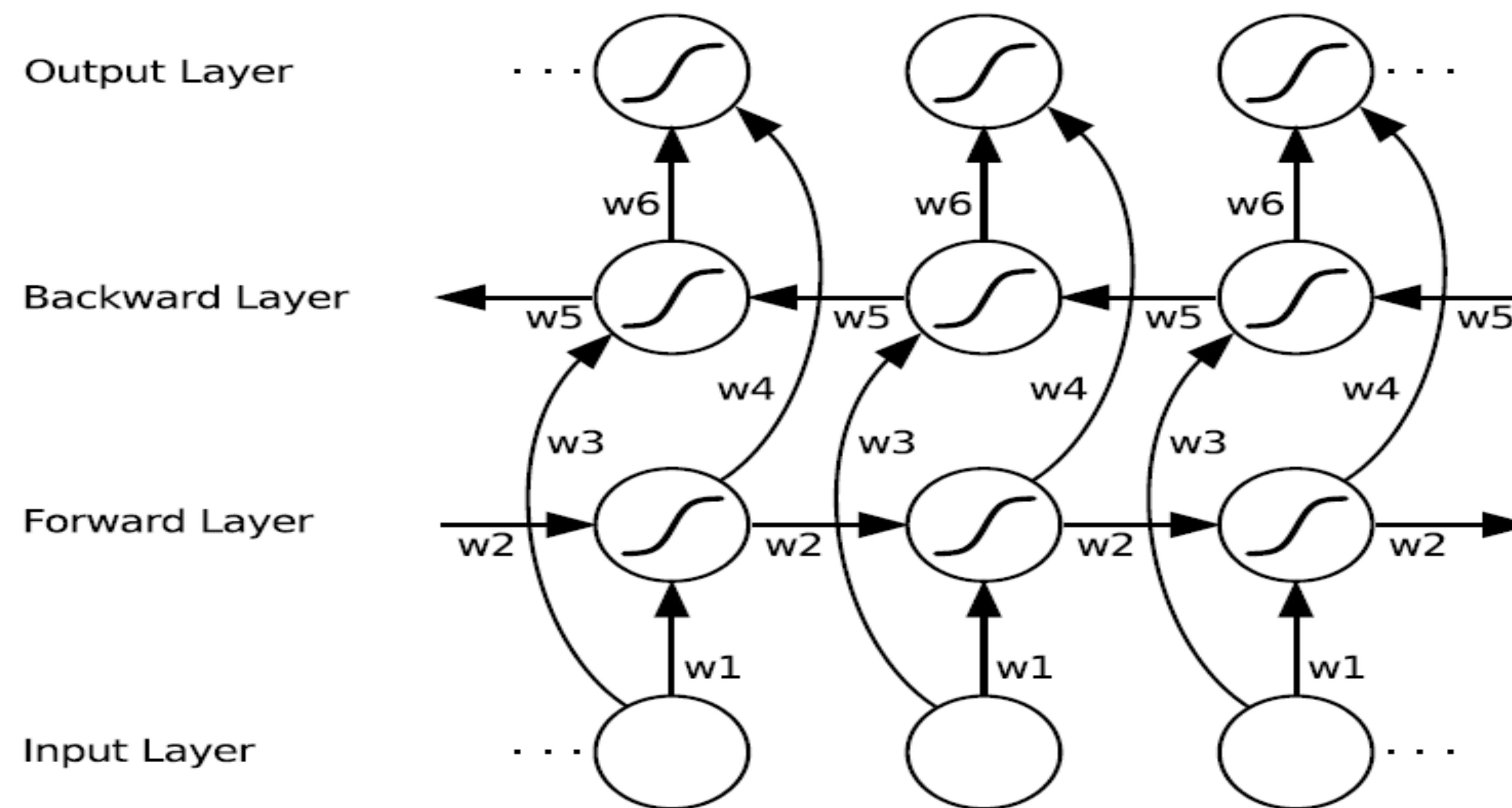
CNN: 学到了什么特征？



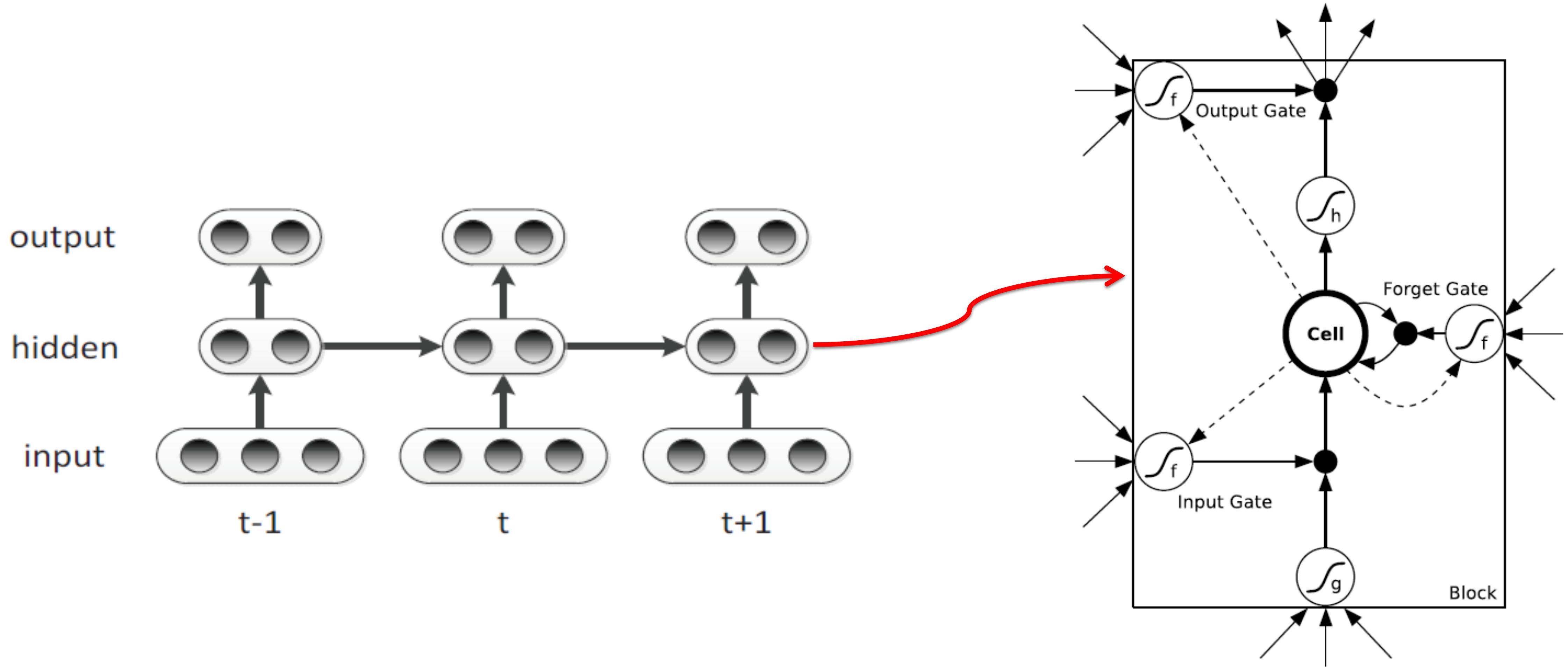
RNN



双向RNN



LSTM



双向深度LSTM

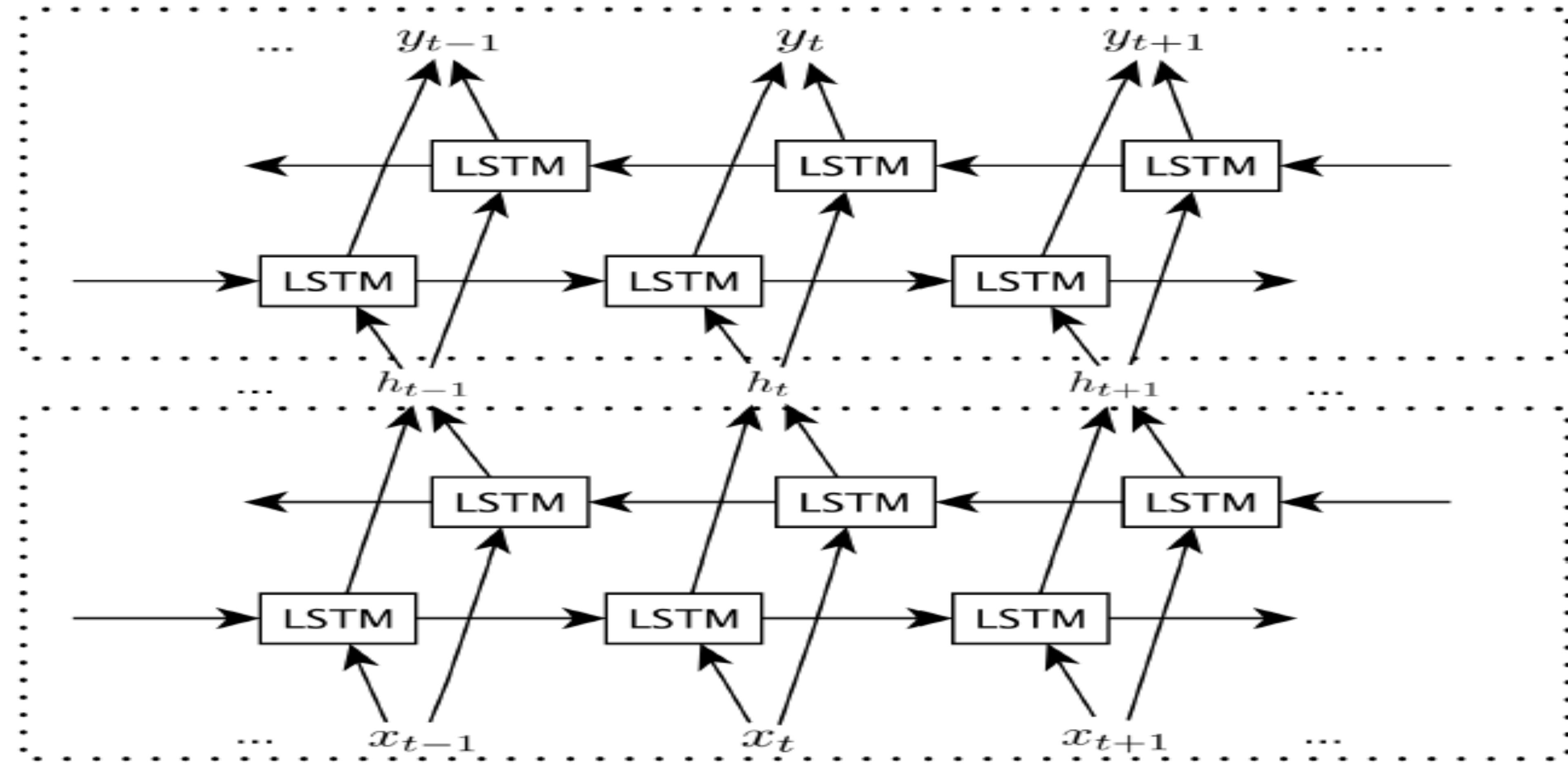


TABLE OF CONTENTES

当深度学习遇到CTR预估

传统主流CTR预估方法

深度学习基础模型

深度学习CTR预估模型

互联网公司深度学习CTR案例

深度学习CTR模型要解决的几个关键问题

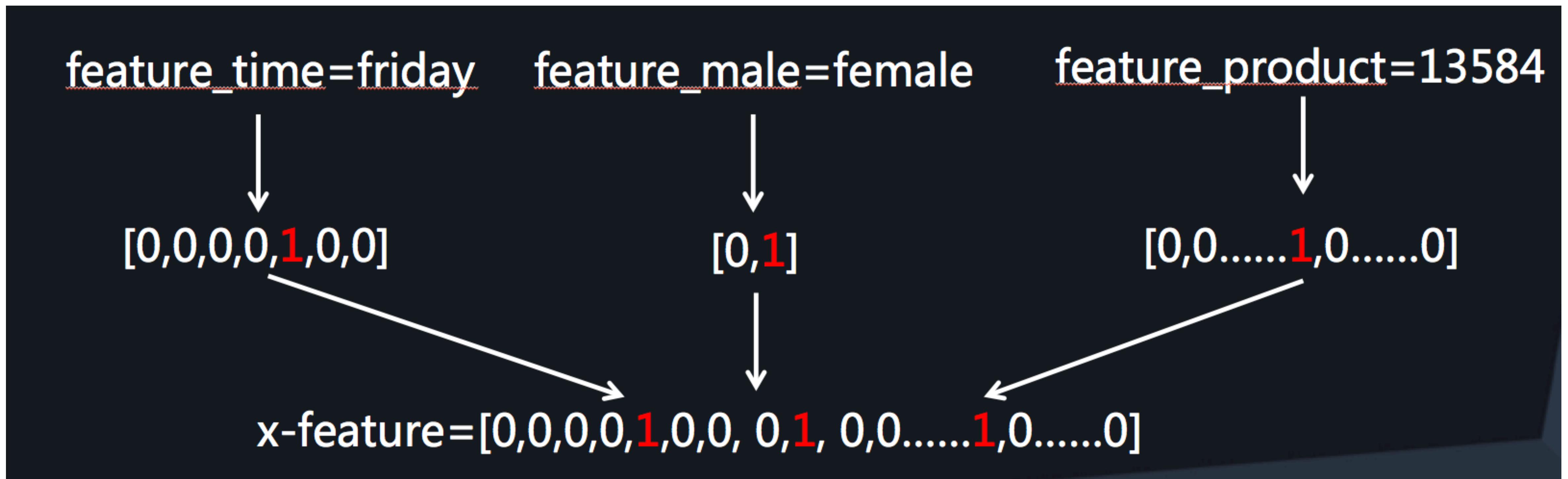
- CTR任务特点：大量离散特征的表示问题
- CTR任务特点：如何快速处理大量高维度稀疏特征？（OneHot 2 Dense）
- 特征工程：如何从手工到自动？（深度学习的优势）
- 特征工程：如何捕获和表达两两组合特征？（FM机制神经网络化）
- 特征工程：如何捕获和表达多组组合特征？（利用Deep网络）

CTR任务中的特征类型

- 连续特征
 - 收入，身高，体重.....
 - 适合DNN处理
- 离散特征
 - 职业，性别，毕业学校.....
 - 不适合DNN处理

离散特征如何让DNN可以处理？

- 直观思路：离散特征使用Onehot表达



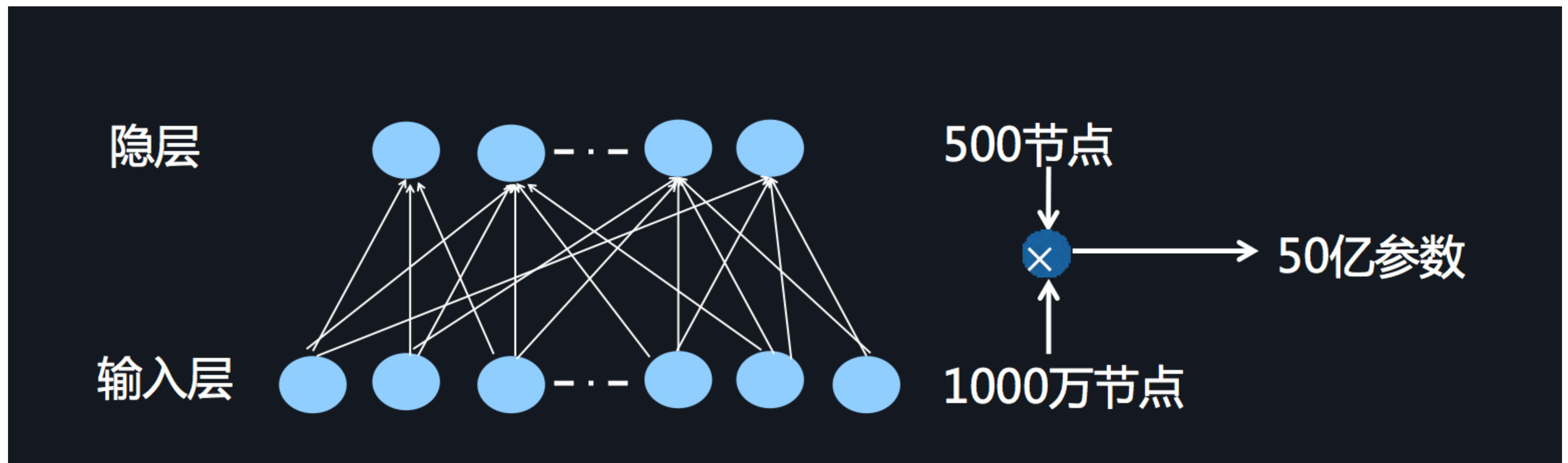
离散特征如何让DNN可以处理？

- 直观思路：离散特征使用Onehot表达



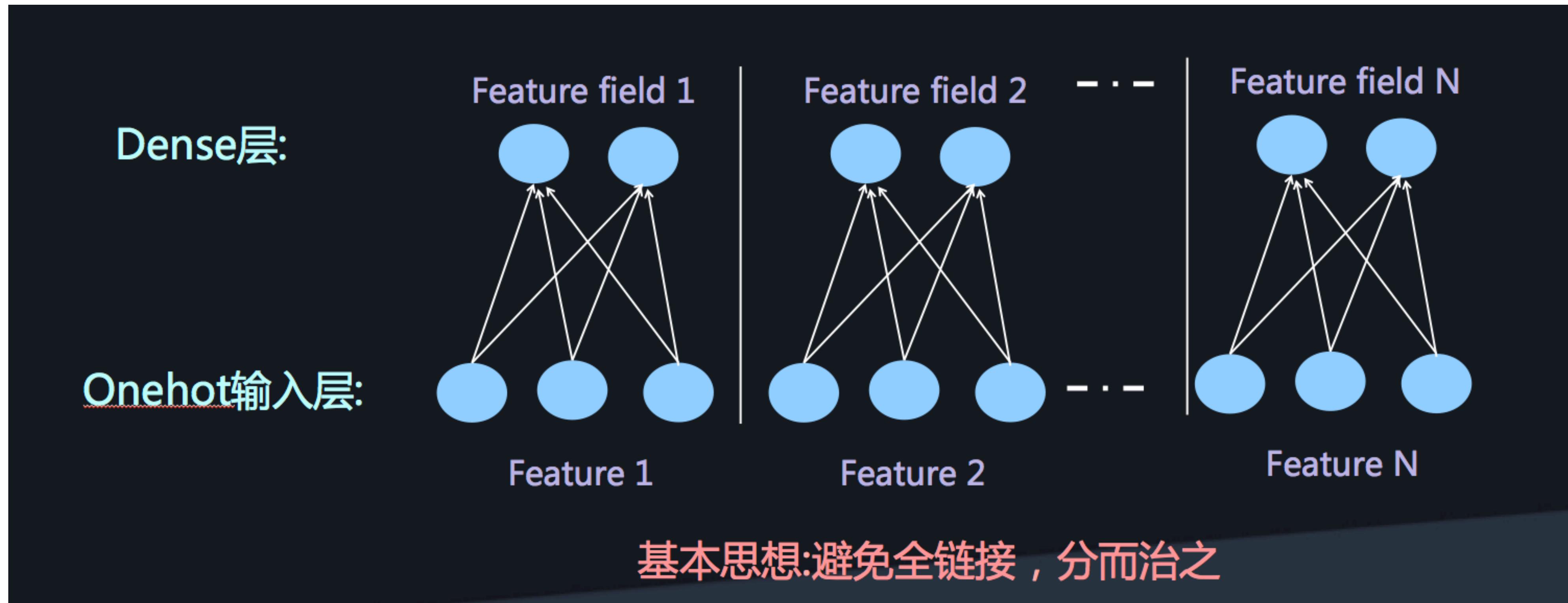
离散特征如何让DNN可以处理？

- Onehot作为DNN输入的问题：CTR预估任务里不可行

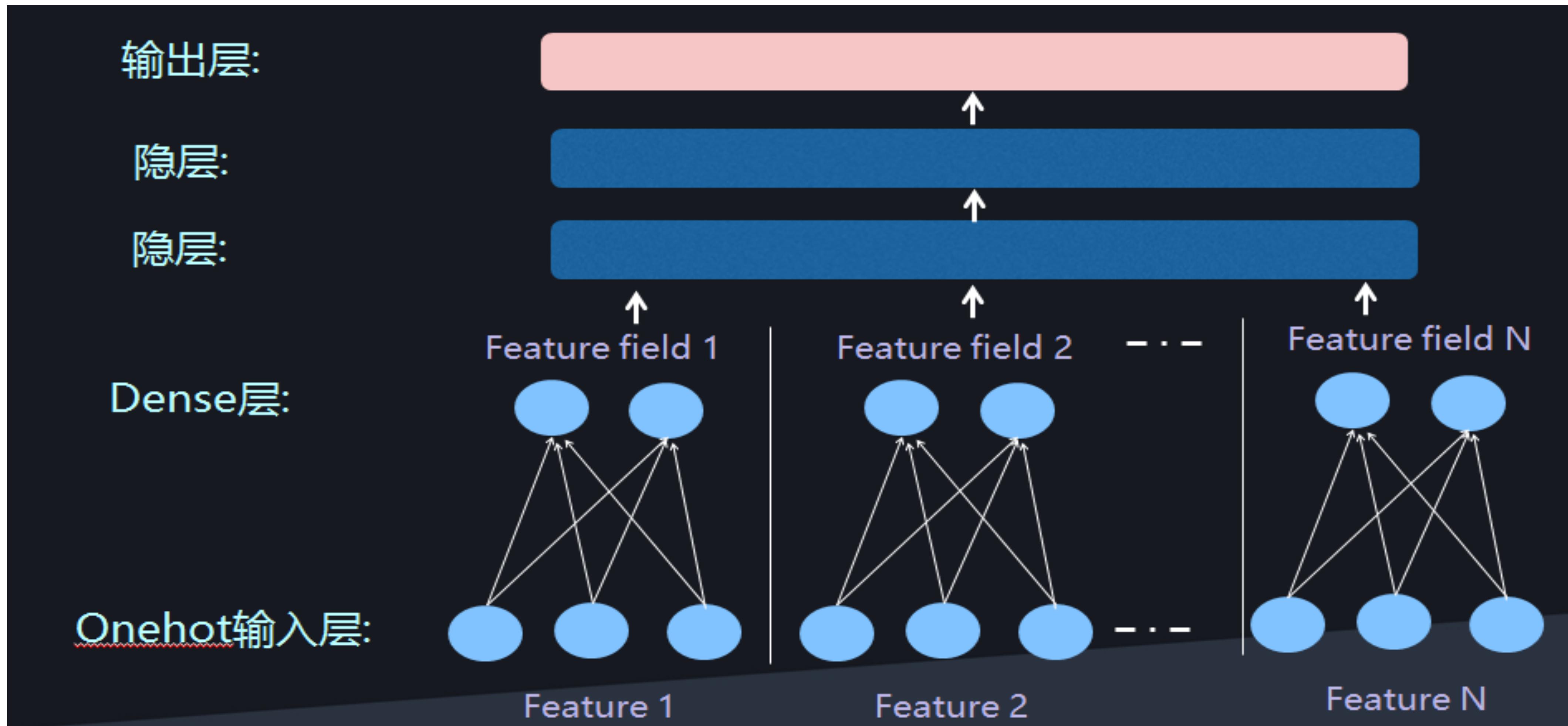


离散特征如何让DNN可以处理？

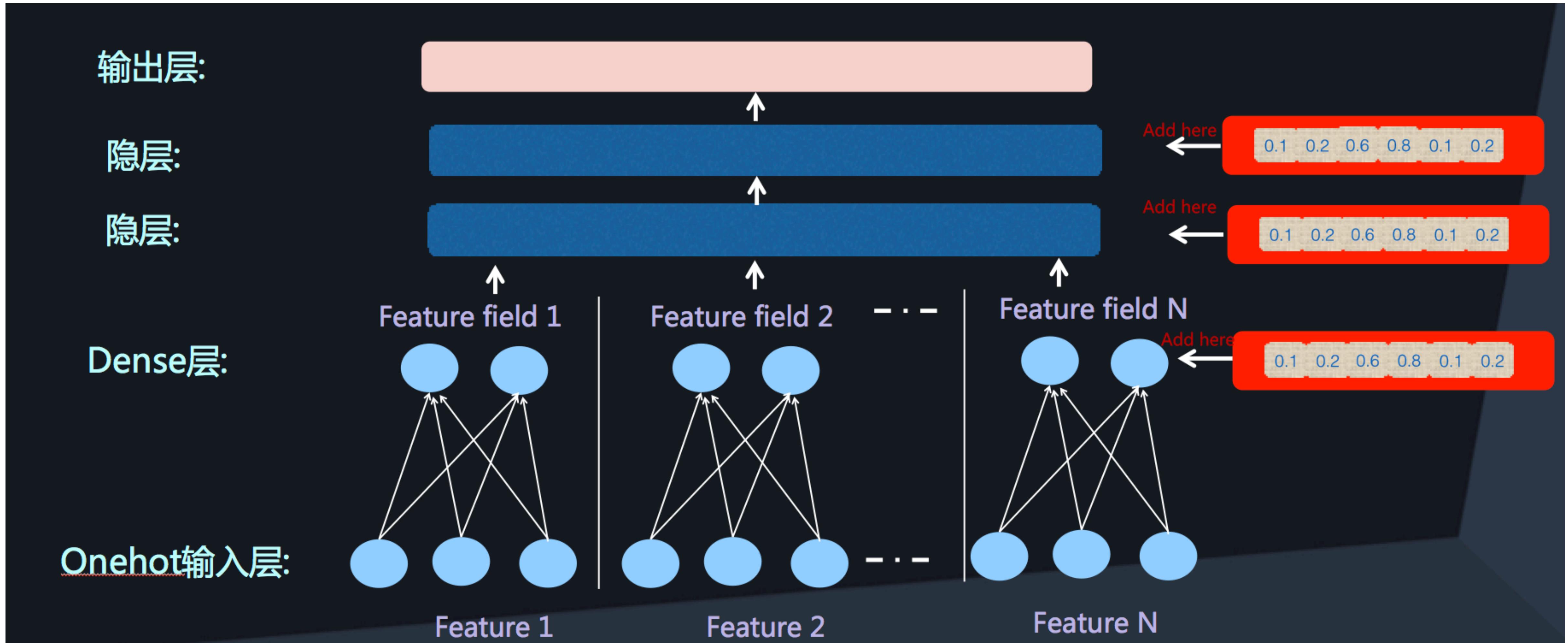
- 解决思路：从OneHot到Dense Vector



形成DNN结构

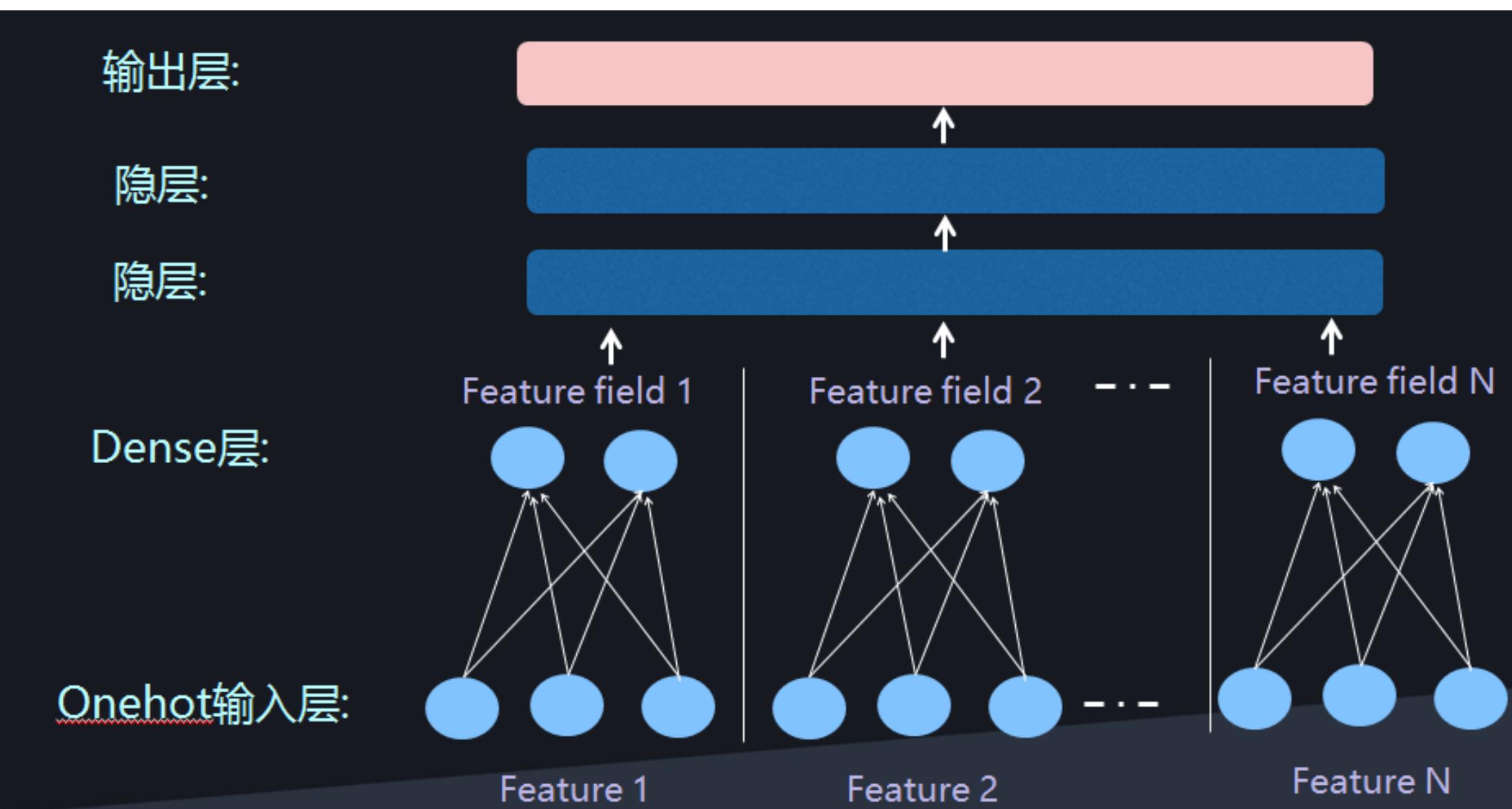


再加入连续特征



这是通用的深层模型结构

这就是FNN模型：Factorisation-Machine Supported Neural Networks



Wide&Deep模型的Deep部分：相同的结构

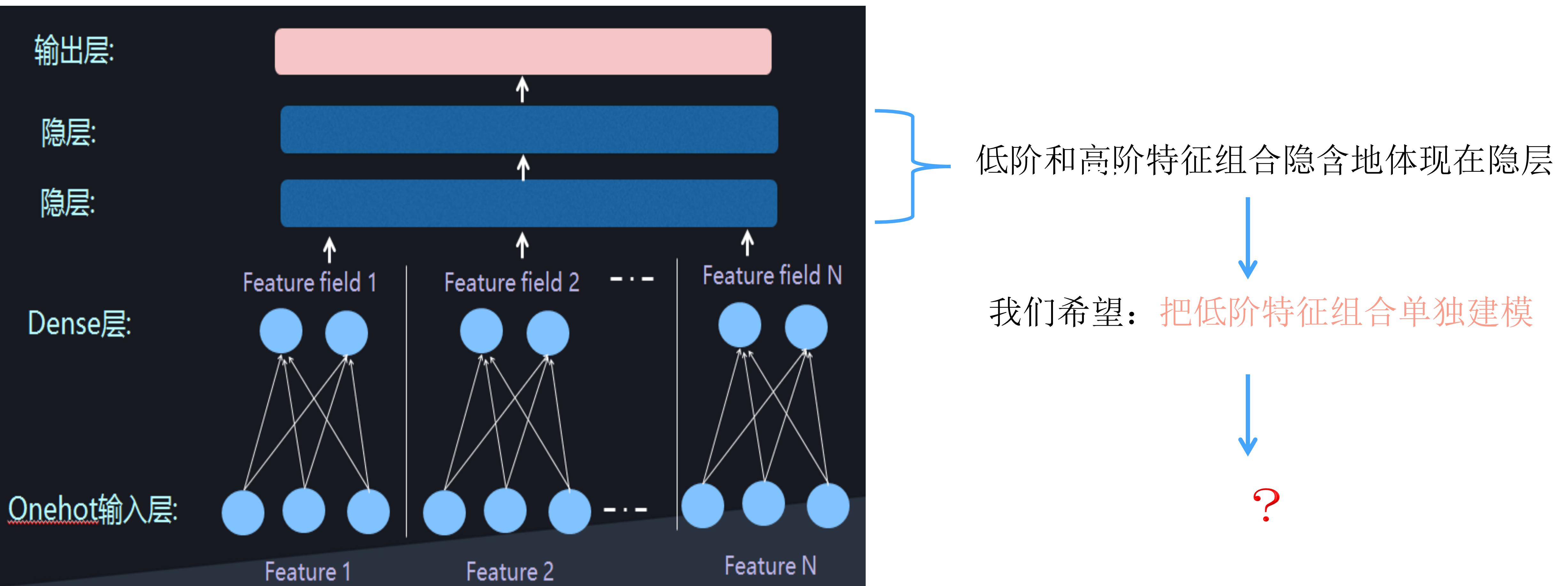


很多其它改进模型的Deep部分：相同的结构



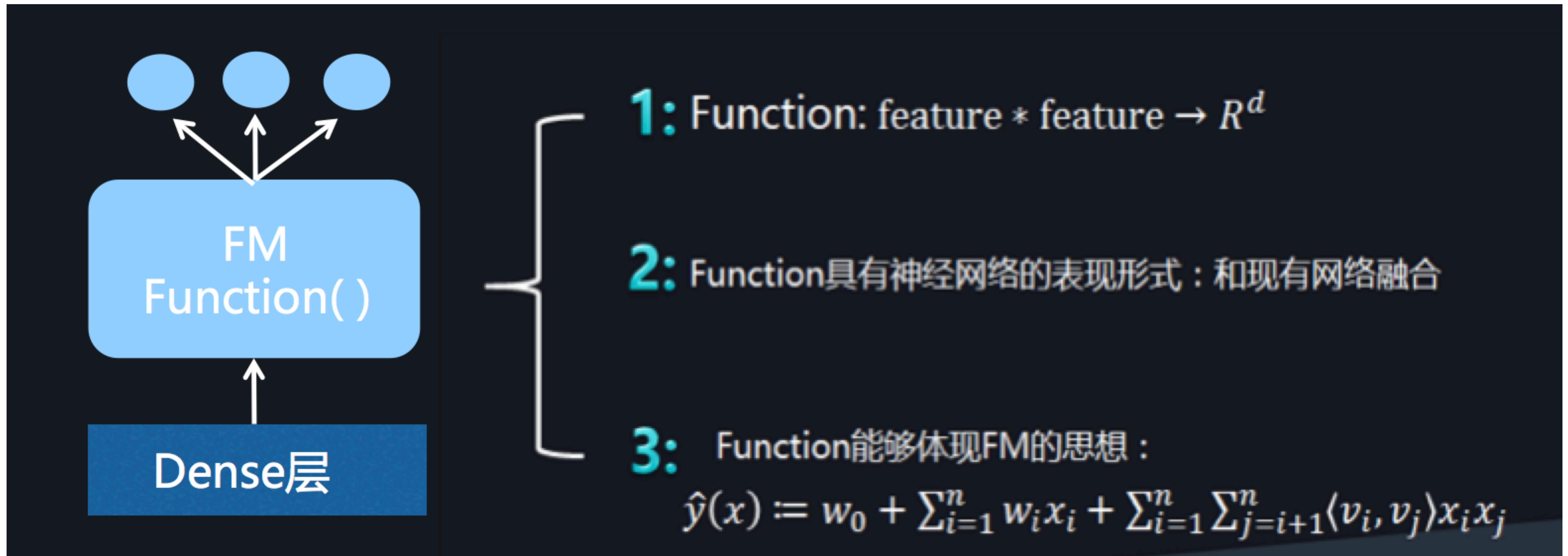
几乎所有DL+CTR模型的输入部分:这种Onehot2Dense映射

DNN输入问题解决了，但是.....

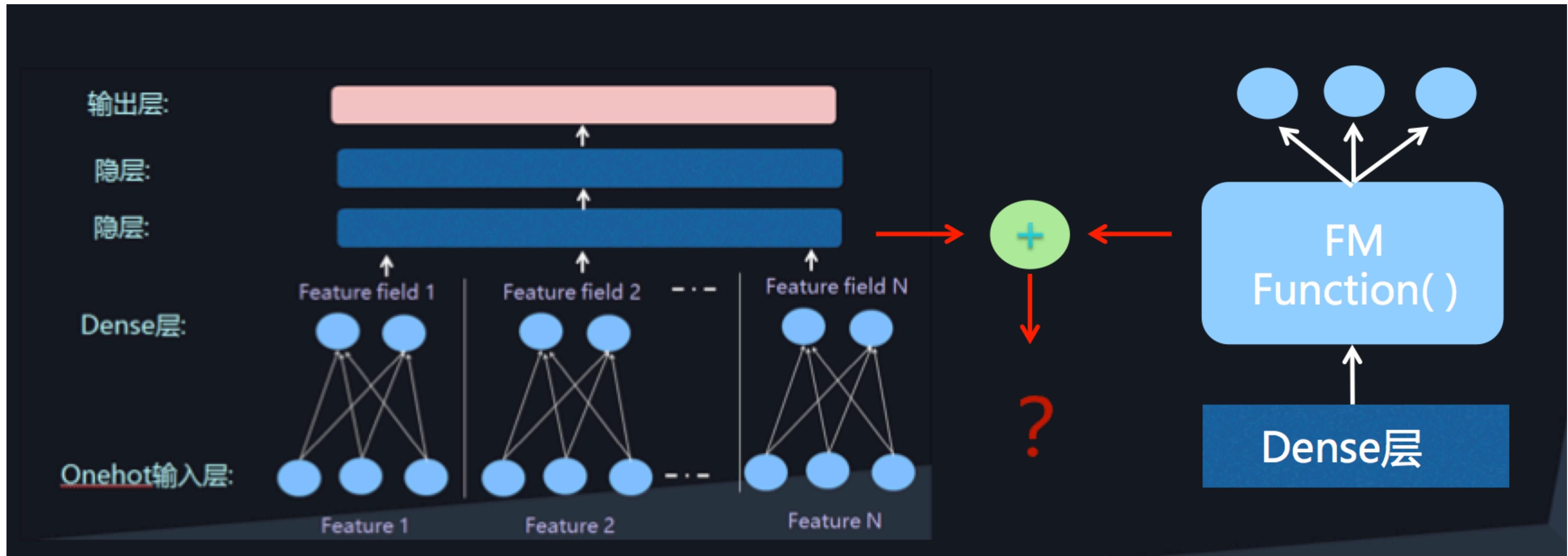


把低阶特征组合单独建模

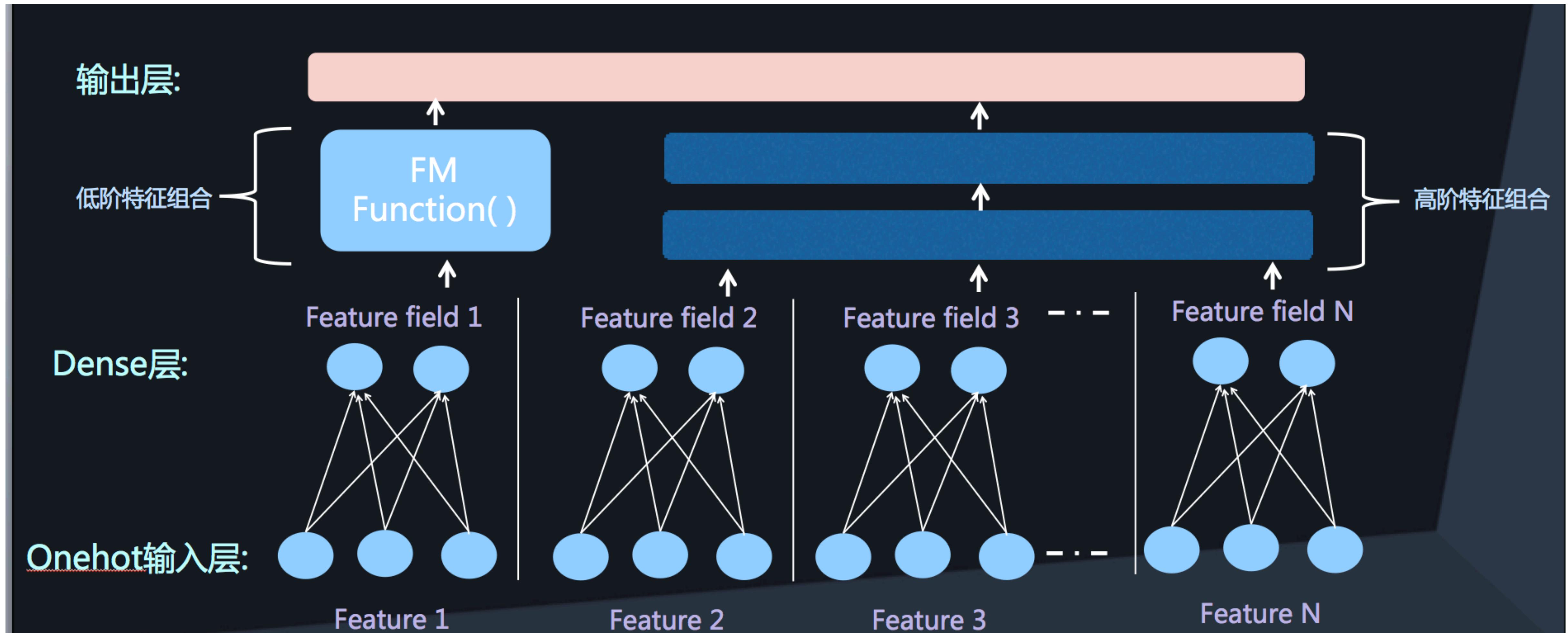
- 首先需要：定义一个神经网络版的低阶特征组合模型



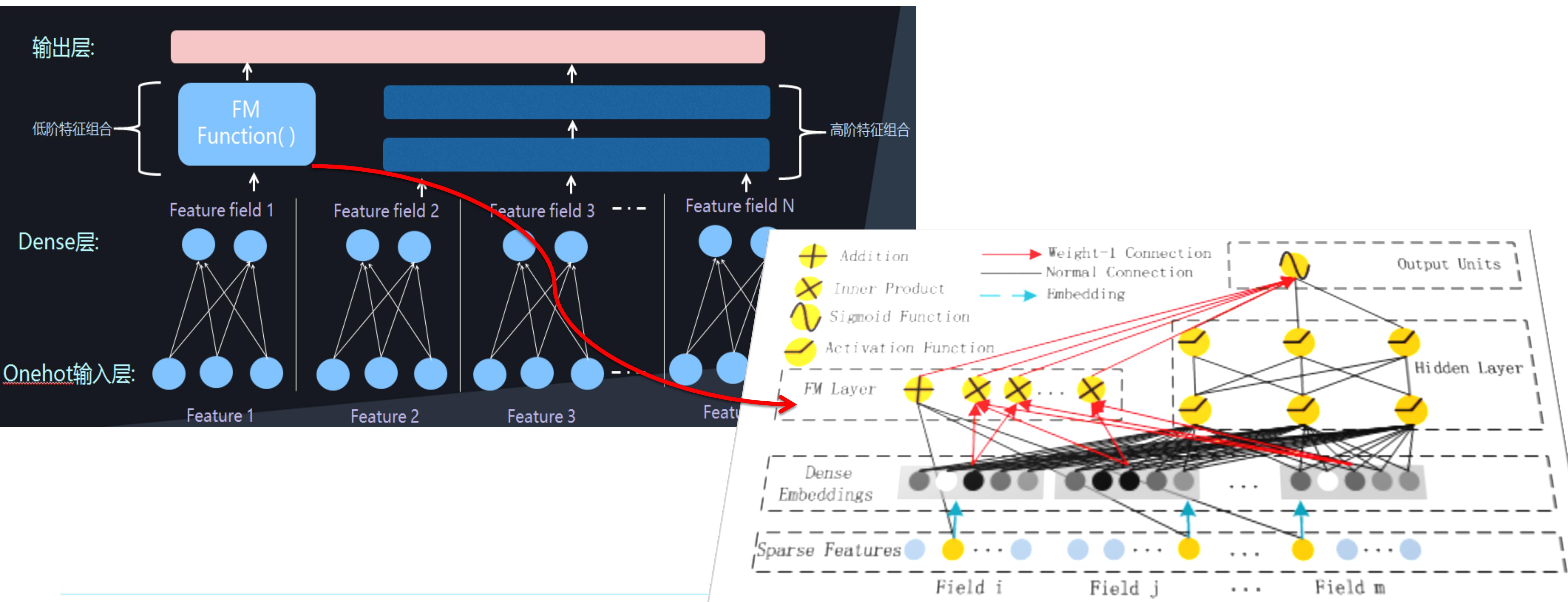
把低阶特征组合模型插入网络结构中



典型网络融合结构之一：并行结构

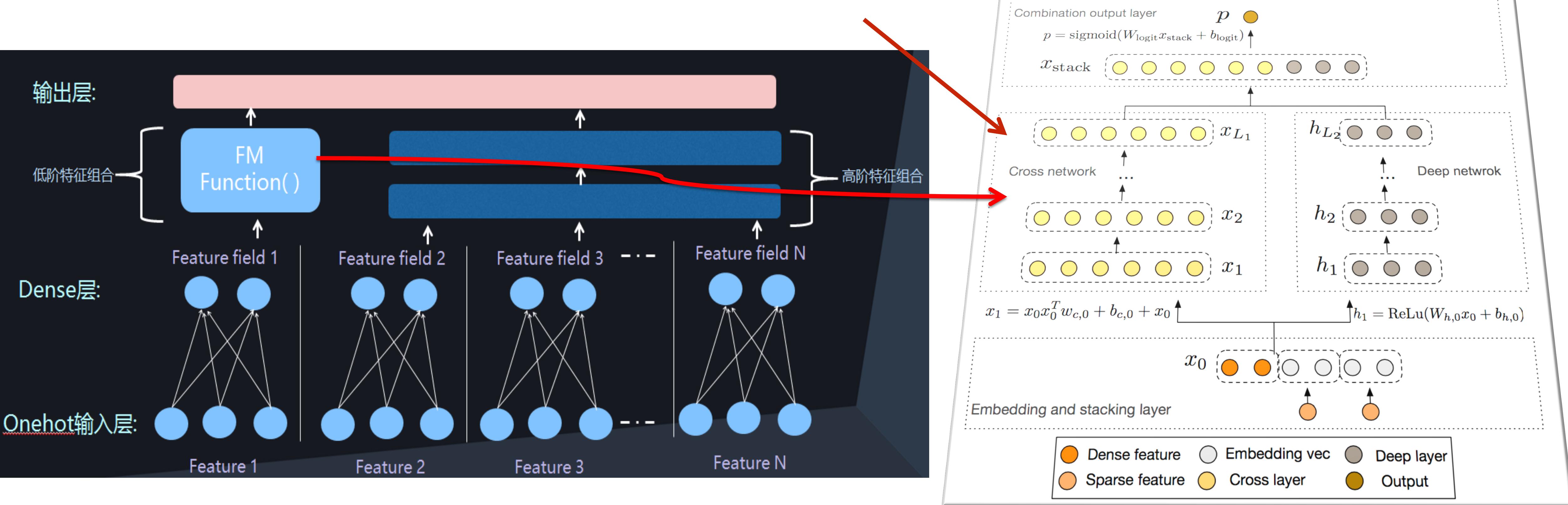


并行结构实例：DeepFM模型

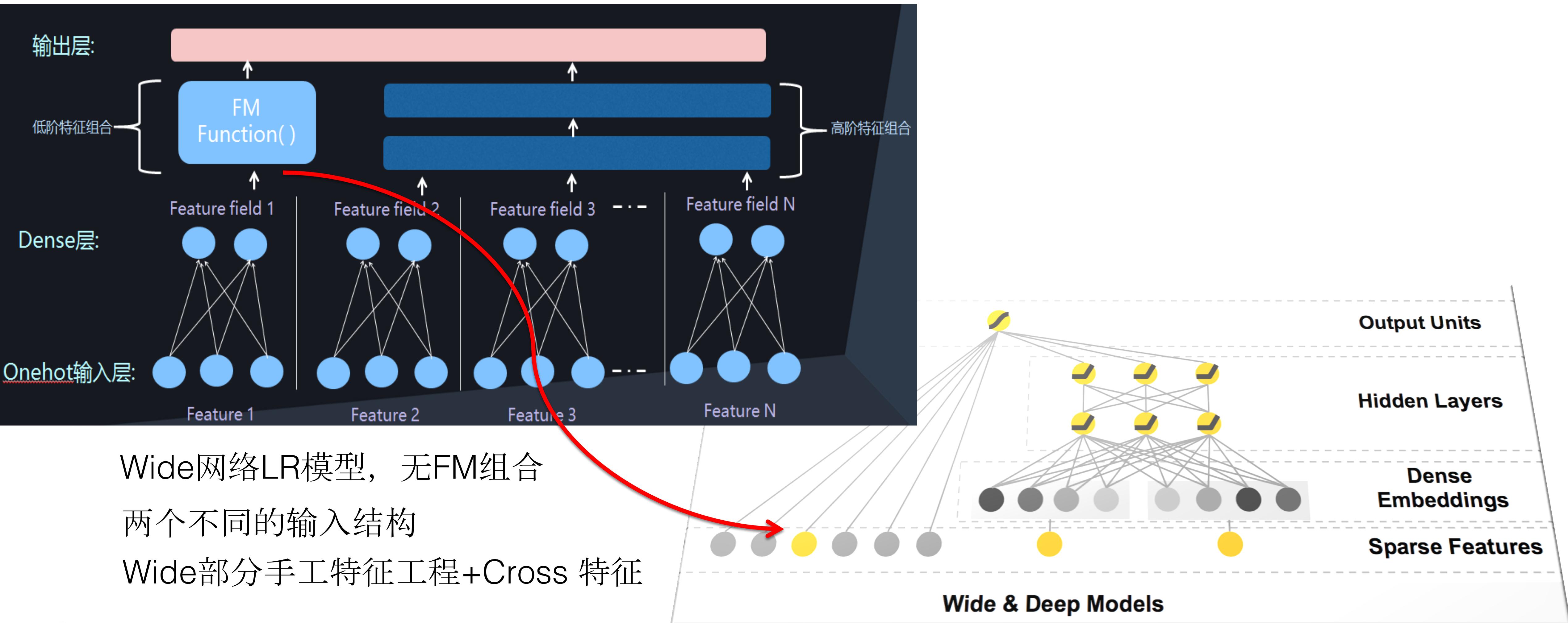


并行结构实例：Deep&Cross模型

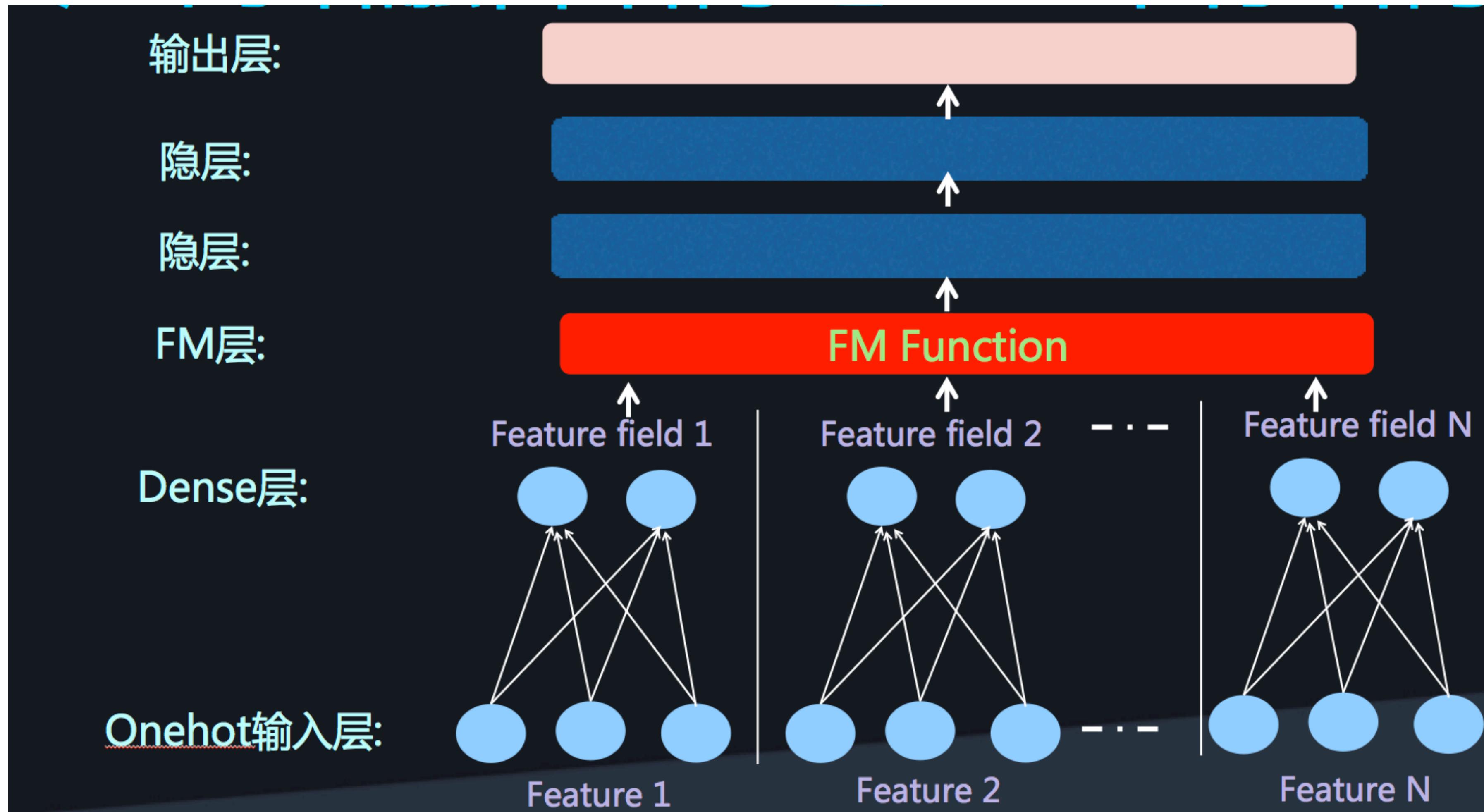
FM Function=Cross Network=ResNet



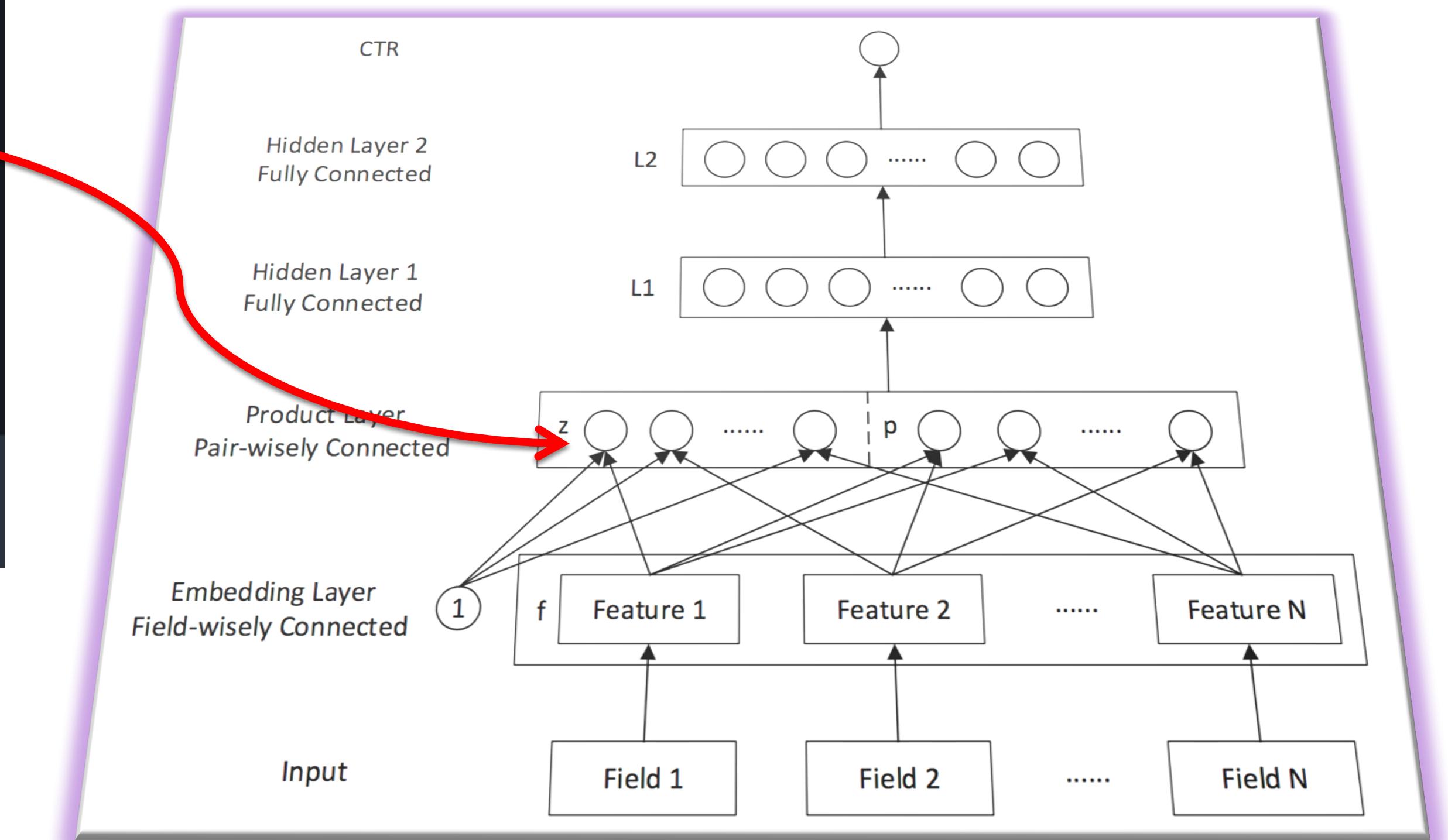
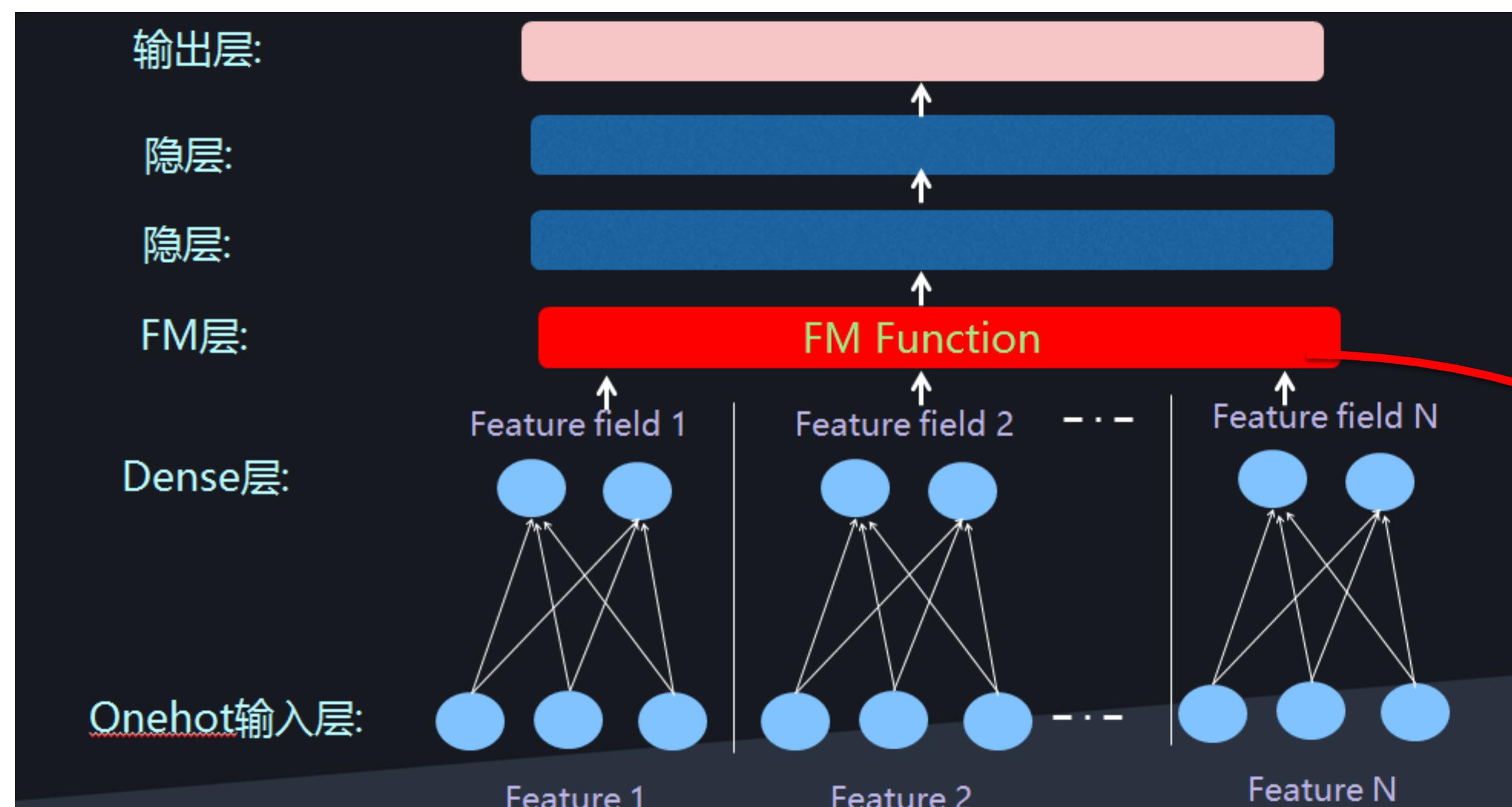
并行结构实例：Wide&Deep模型



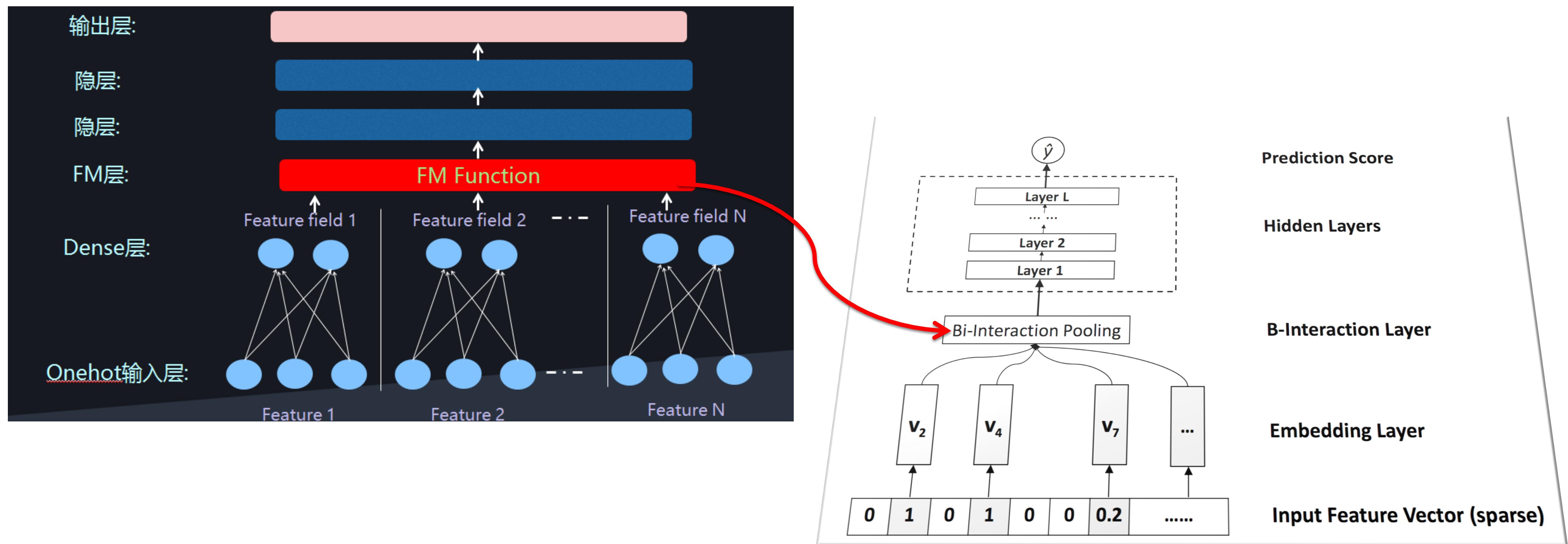
典型网络融合结构之二：串行结构



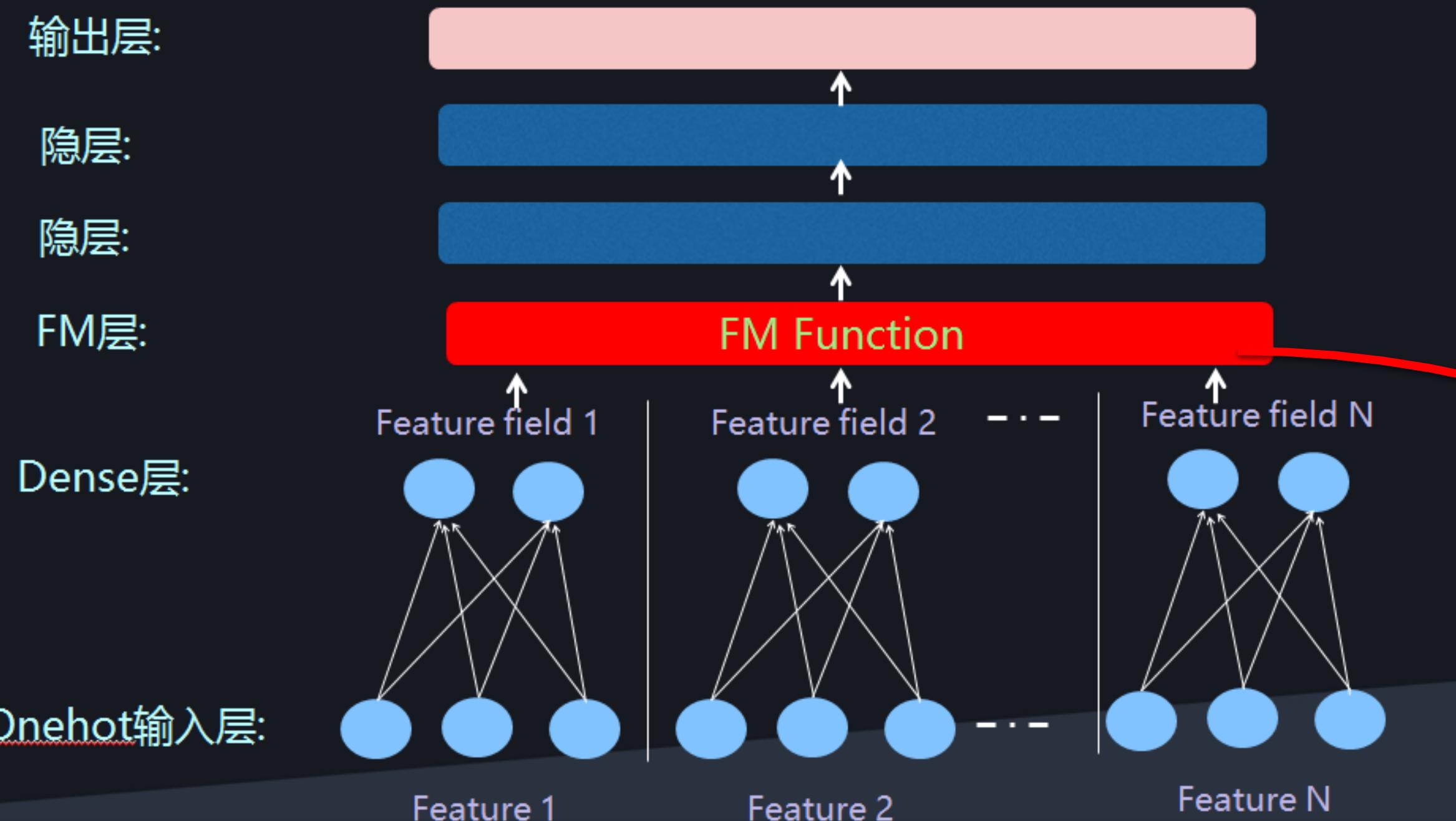
串行结构实例：PNN模型



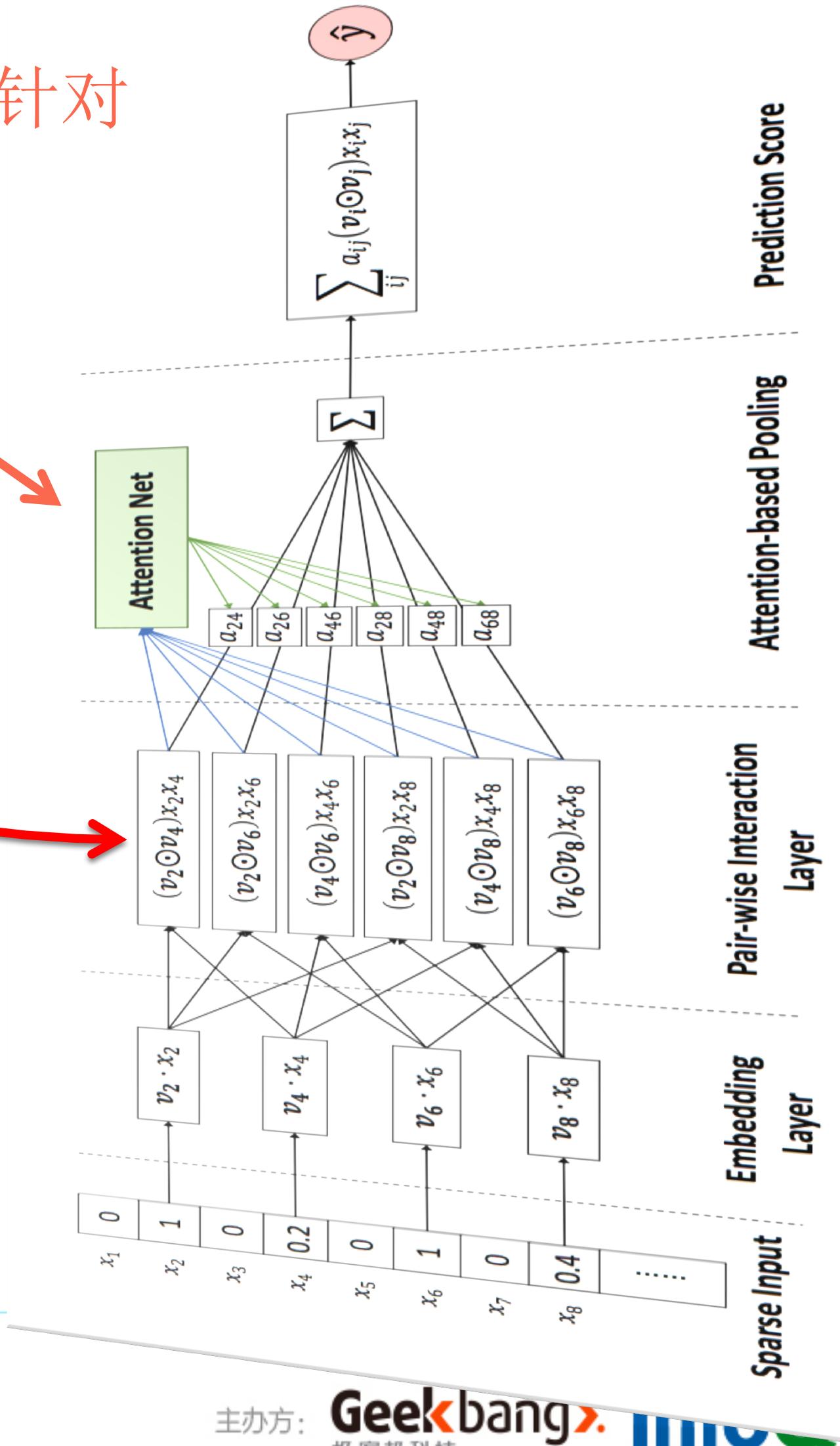
串行结构实例：NFM模型



串行结构实例：AFM模型



相比NFM模型只是多了针对组合特征的Attention



模型训练与优化：Dense层的预训练

1: Wide & Deep模型Dense层需要预训练



用FM初始化Onehot到Dense层的映射性能明显提升

2: 类似的FNN等无明显FM结构的模型Dense层需要预训练

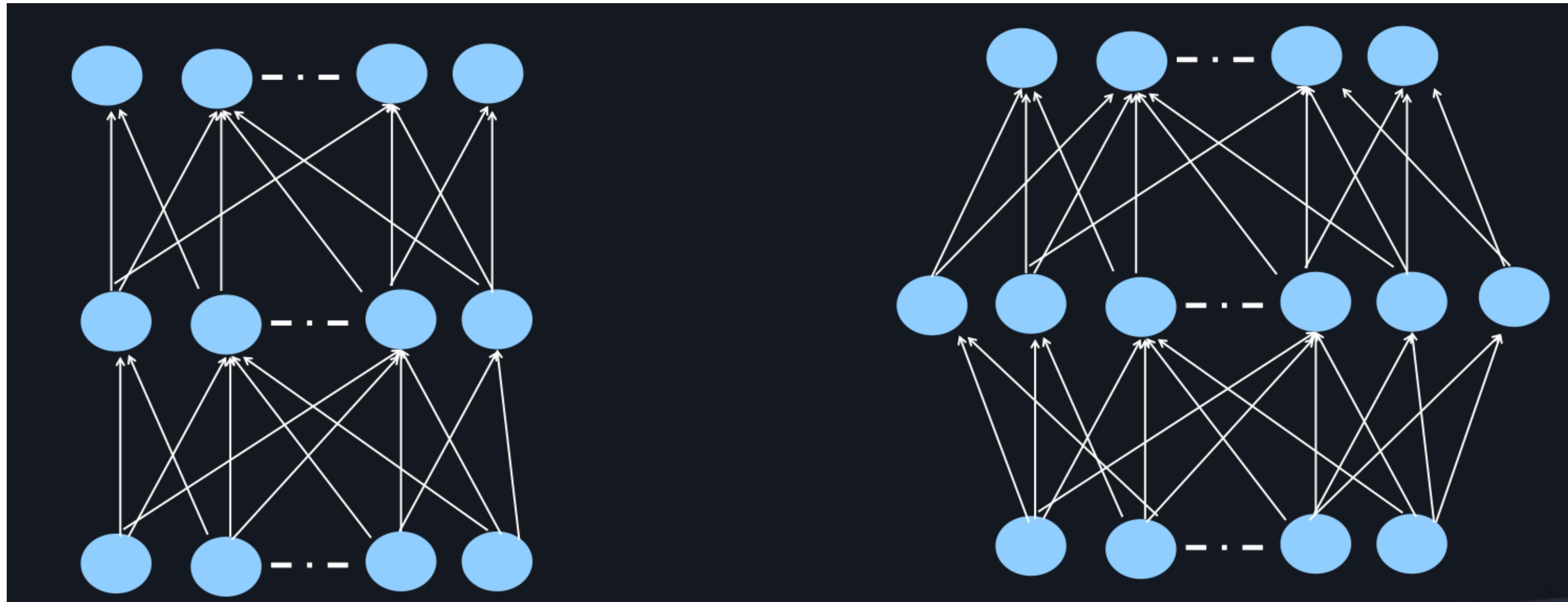


3: 串行结构的模型Dense层需要预训练



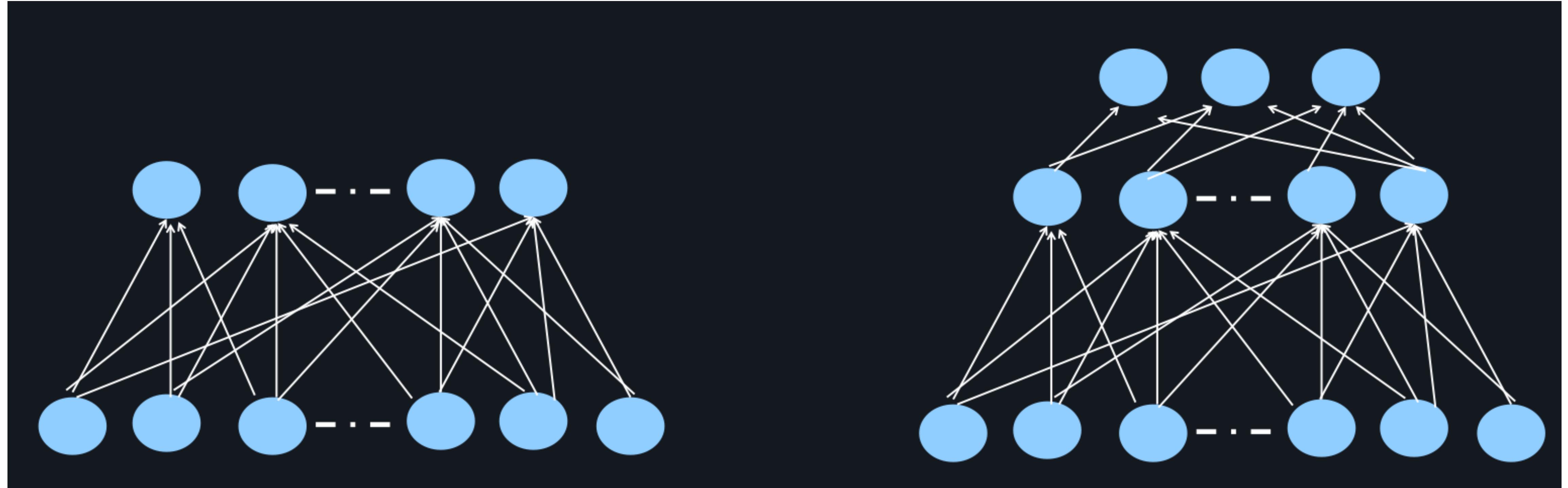
?

模型训练与优化：Deep网络隐层网络结构



平行结构或者菱形结构效果较好

模型训练与优化：Deep网络隐层的层深



两层或三层

TABLE OF CONTENTES

当深度学习遇到CTR预估

传统主流CTR预估方法

深度学习基础模型

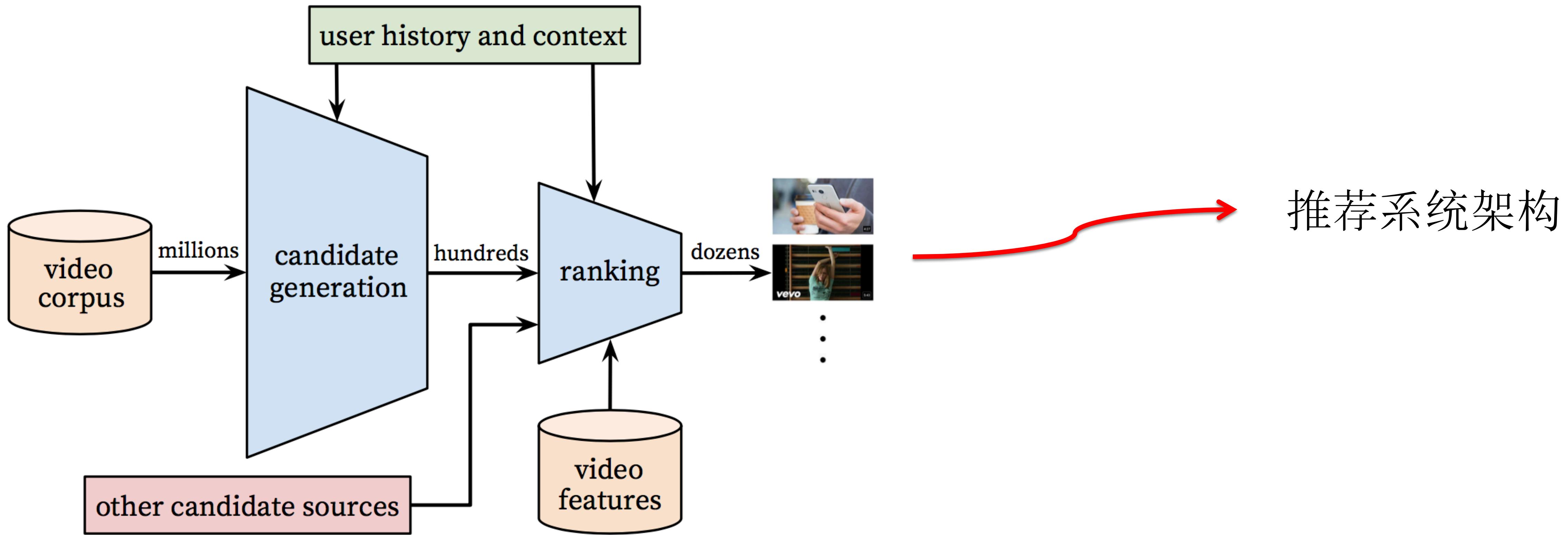
深度学习CTR预估模型

互联网公司深度学习CTR案例

Google : Youtube视频推荐

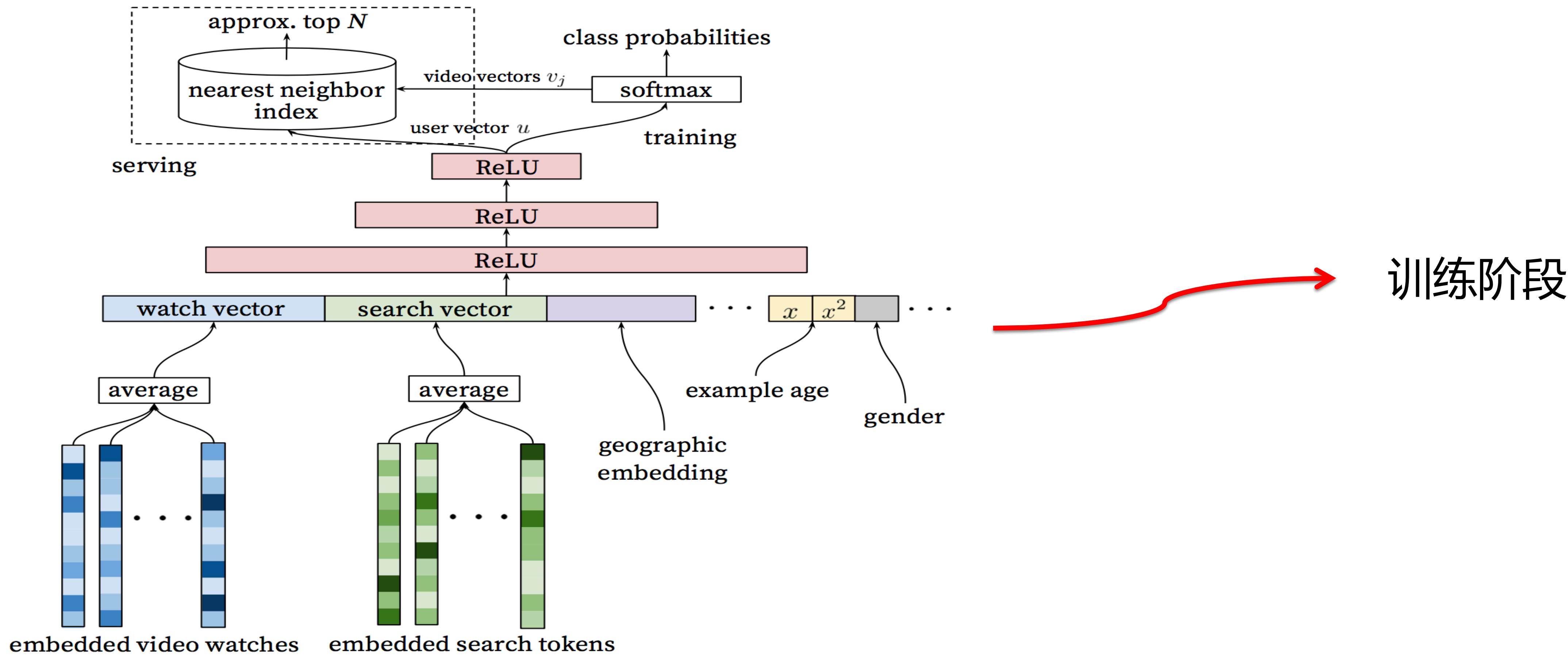
- Google 2016年工作
 - 论文：Deep Neural Networks for YouTube Recommendations
 - 已上线，效果优于传统模型
- 主要关注
 - 视频推荐
 - 深度神经网络的应用

Google : Youtube视频推荐

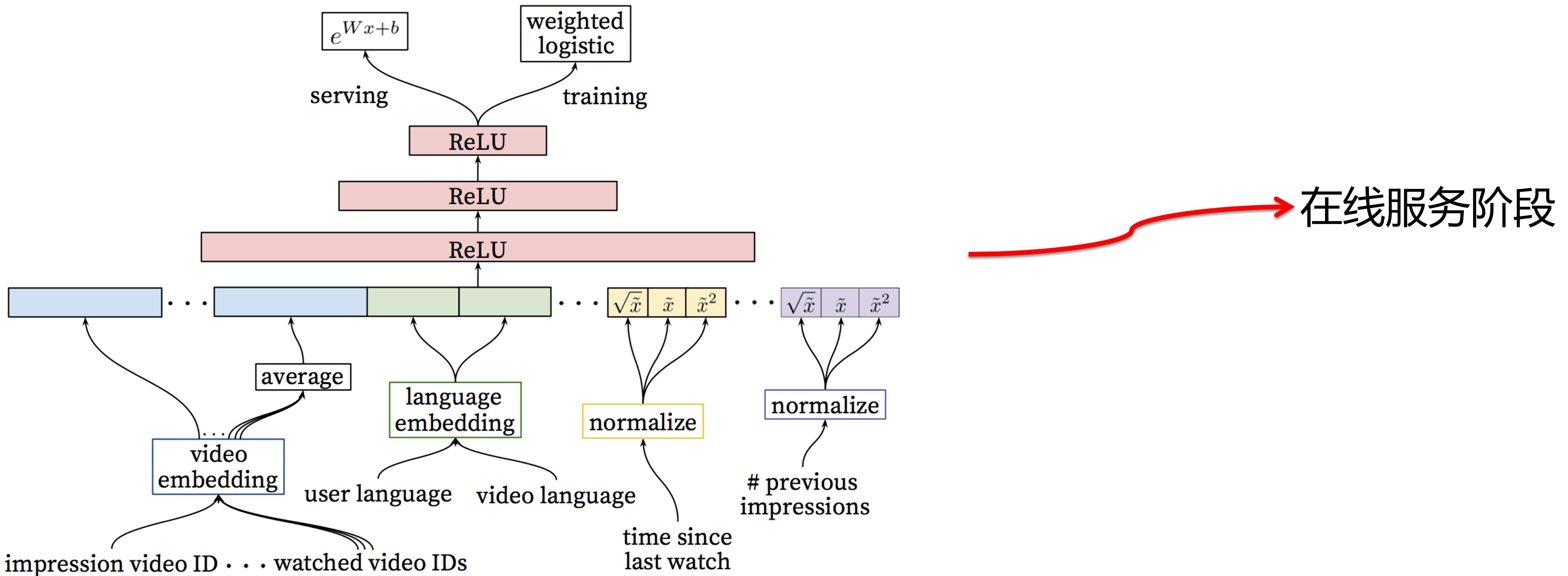


推荐系统架构

Google : Youtube视频推荐



Google : Youtube视频推荐



Google : Youtube视频推荐

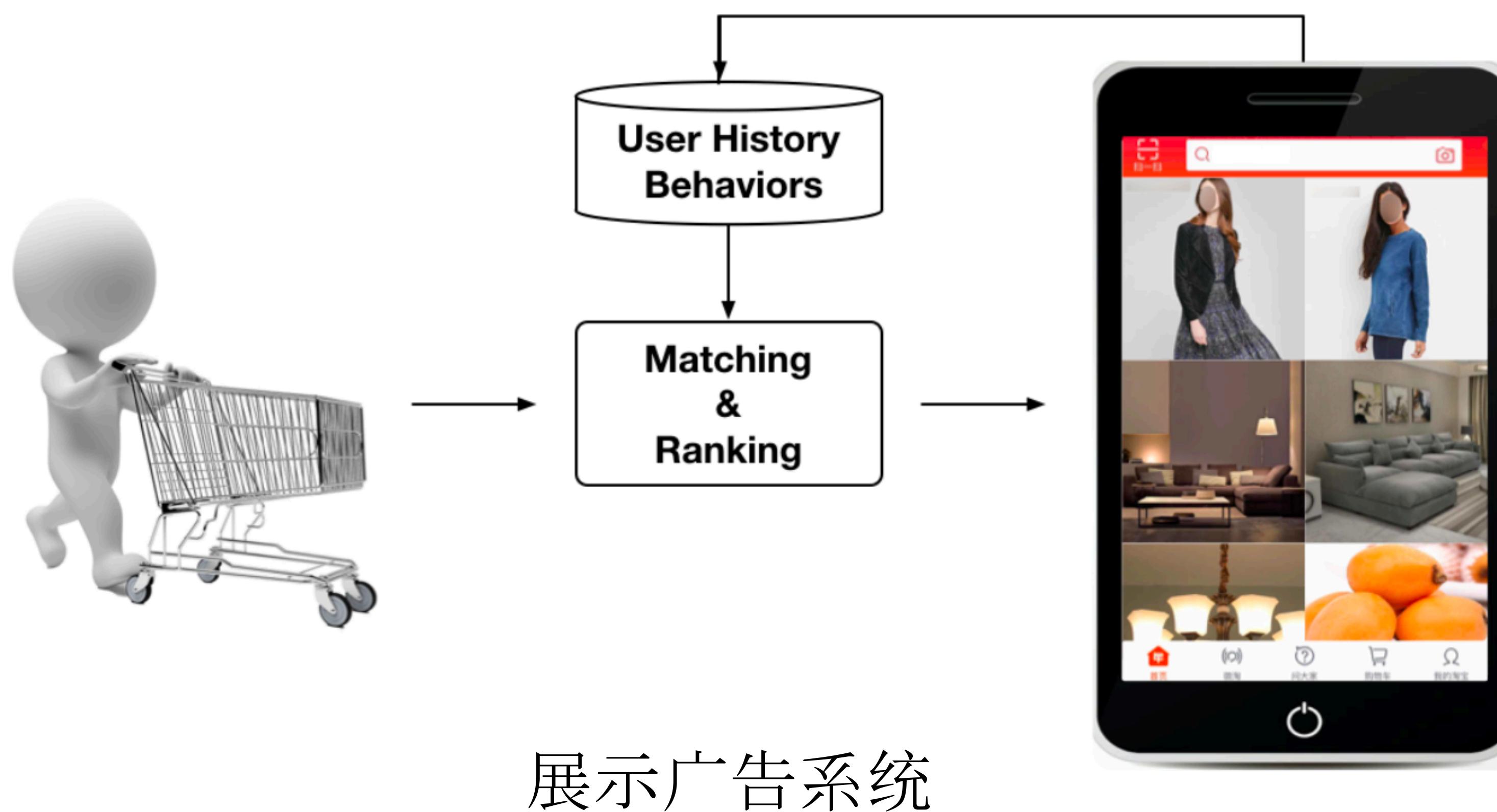
Hidden layers	weighted, per-user loss
None	41.6%
256 ReLU	36.9%
512 ReLU	36.7%
1024 ReLU	35.8%
512 ReLU → 256 ReLU	35.2%
1024 ReLU → 512 ReLU	34.7%
1024 ReLU → 512 ReLU → 256 ReLU	34.6%

参数对性能影响

阿里巴巴 : Deep Interest Network

- 阿里妈妈2017年工作
 - 论文: Deep Interest Network for Click-Through Rate Prediction
 - 已上线, 效果明显优于传统模型
- 主要关注
 - 兴趣的多样性
 - 局部激活 (Local Activation)

阿里巴巴 : Deep Interest Network



阿里巴巴 : Deep Interest Network

Feature Category	Feature Name	Dimemision	Type	#Nonzero Ids/Sample
User Profile Features	gender	2	one-hot	1
	age_level	~ 10	one-hot	1

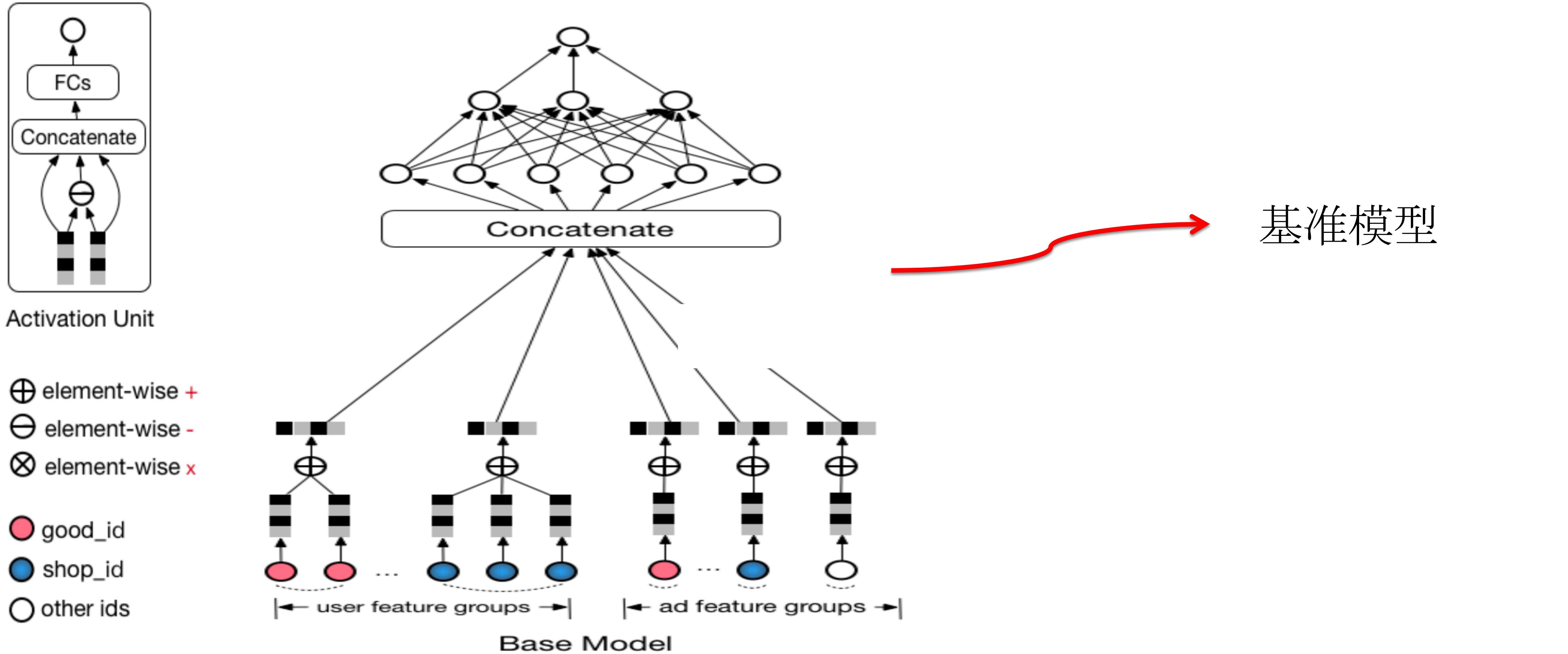
User Behavior Features	visited good_ids	~ 10^9	multi-hot	~ 10^3
	visited shop_ids	~ 10^7	multi-hot	~ 10^3
	visited cate_ids	~ 10^4	multi-hot	~ 10^2

Ad Features	good_id	~ 10^7	one-hot	1
	shop_id	~ 10^5	one-hot	1
	cate_id	~ 10^4	one-hot	1

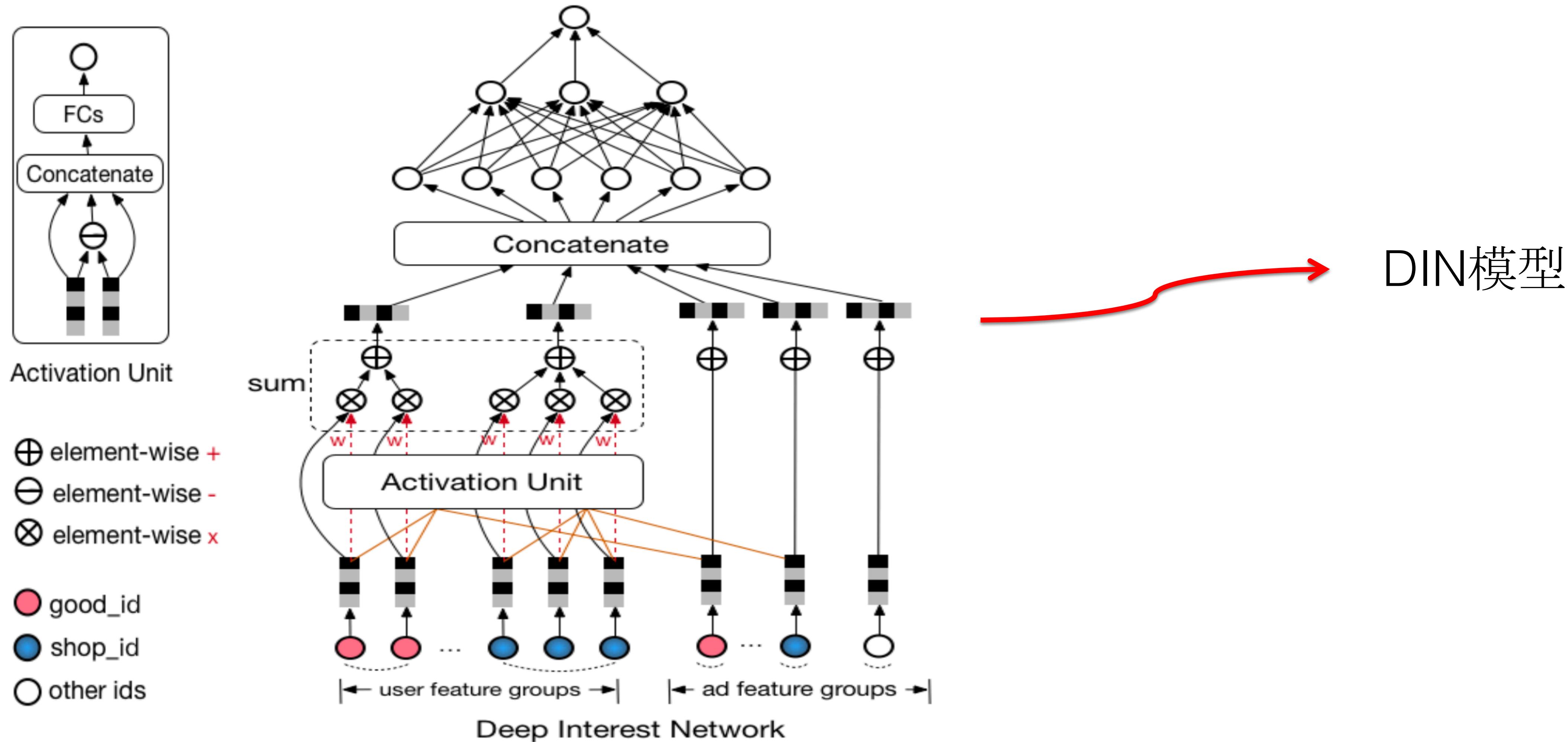
Scene Features	pid	~ 10	one-hot	1
	time	~ 10	one-hot	1

使用特征

阿里巴巴 : Deep Interest Network



阿里巴巴 : Deep Interest Network



阿里巴巴 : Deep Interest Network

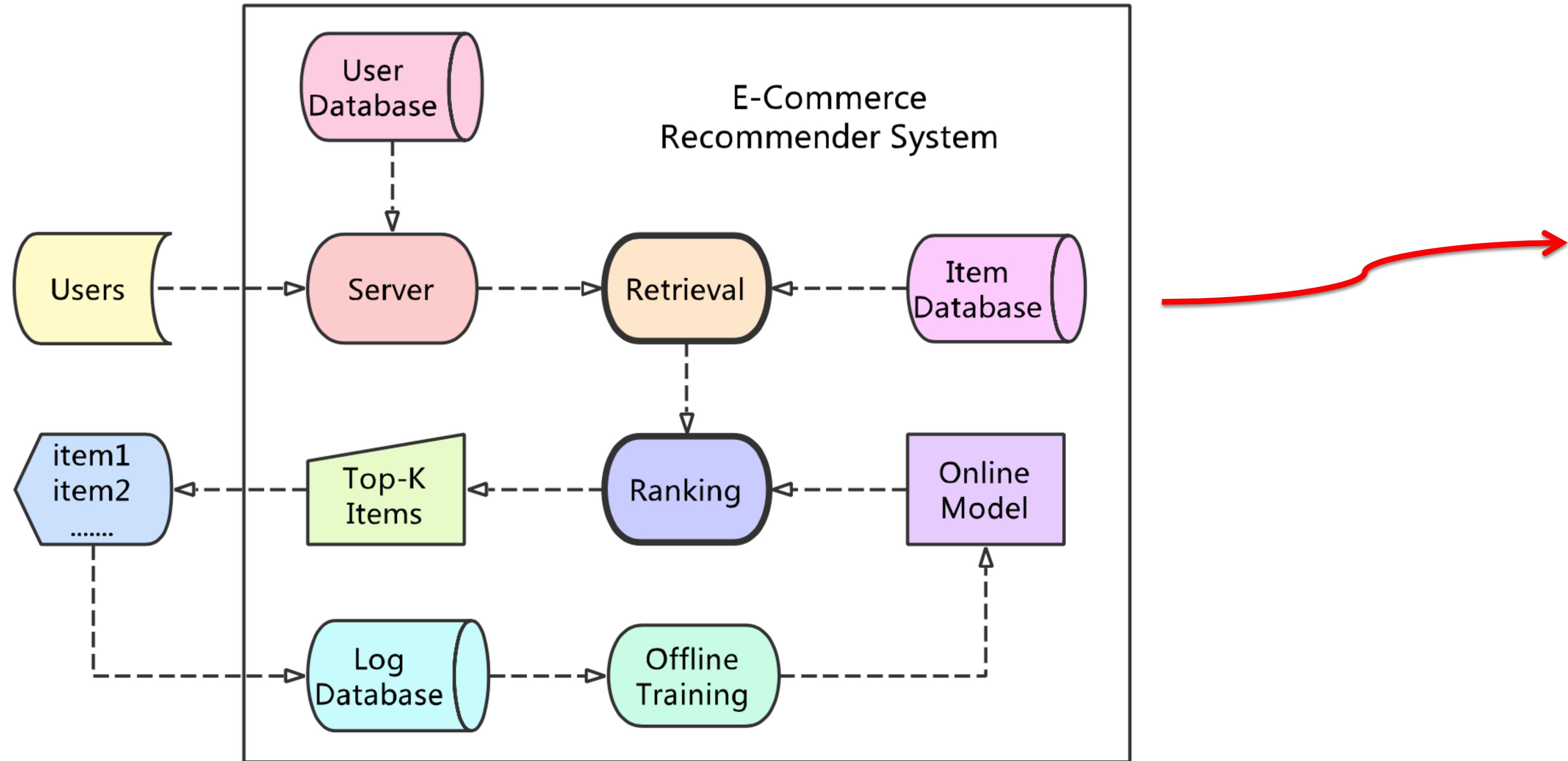


Local Activation

京东商城：Telepath

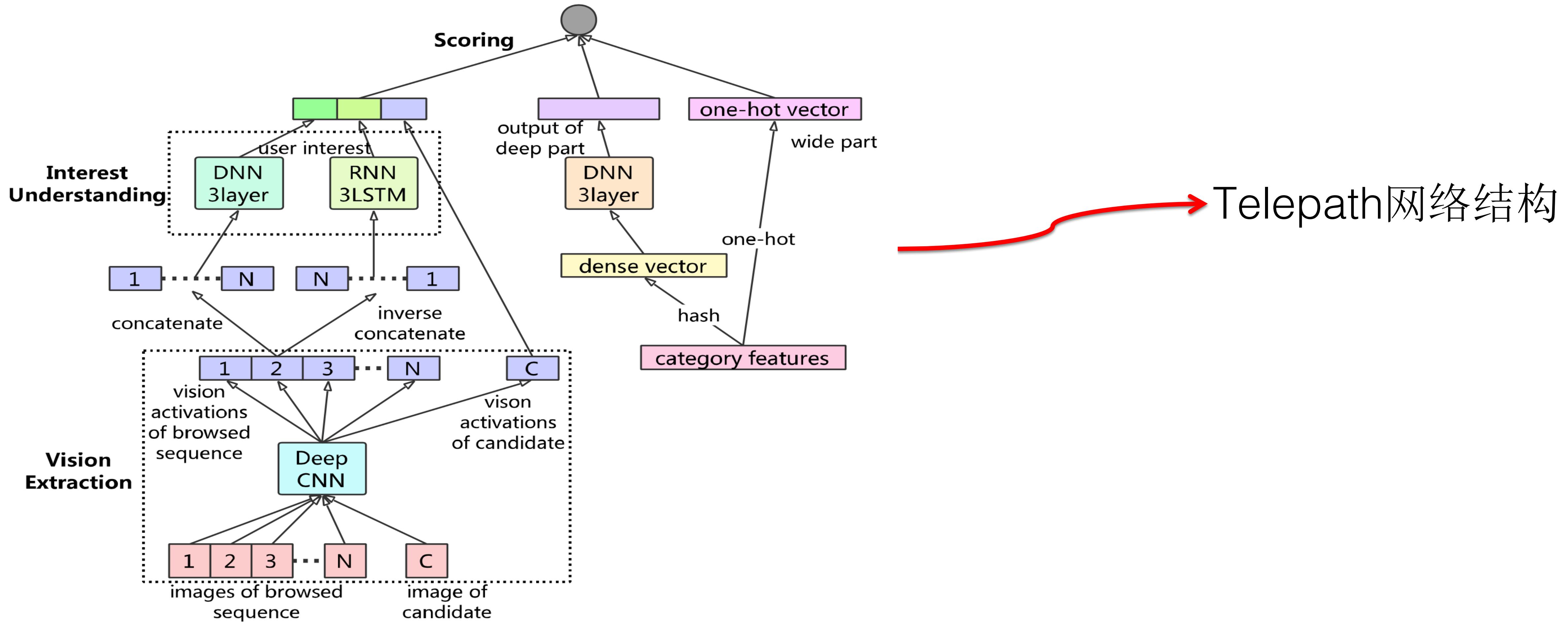
- 京东2017年工作
 - 论文：Telepath: Understanding Users from a Human Vision Perspective in Large-Scale Recommender Systems
 - 已上线，效果明显优于传统模型
- 主要关注
 - 用于推荐和广告
 - 融合图片和行为等各种信息

京东商城：Telepath

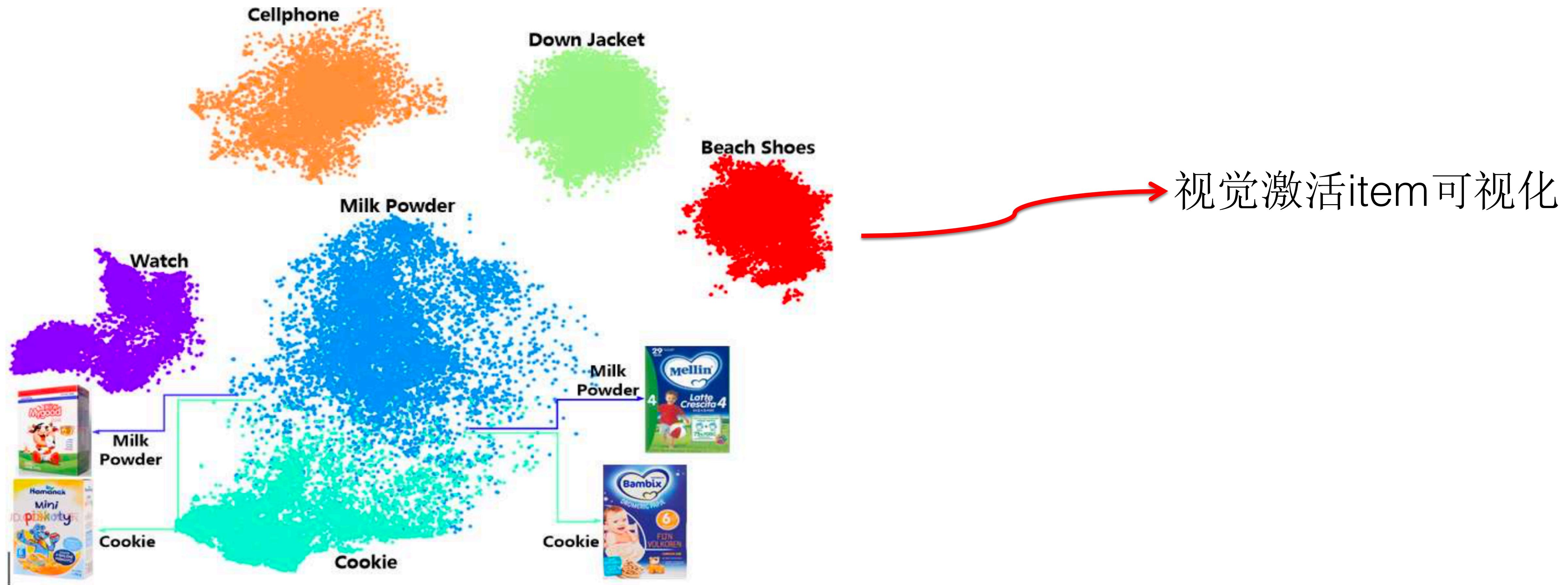


推荐系统架构

京东商城：Telepath

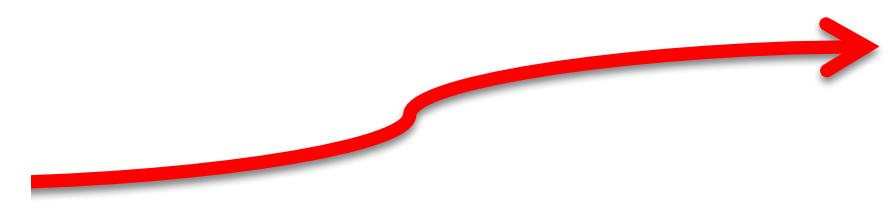


京东商城：Telepath



京东商城：Telepath

Date	Day1	Day2	Day3	Day4
CTR	+5.15%	+8.07%	+10.5%	+6.15%
GMV	+ 126.48%	+9.1%	+18.4%	-19.24%
ROI	+129.53%	+14.35%	+14.2%	-17.44%
Date	Day5	Day6	Day7	Average
CTR	+4.63%	+2.11%	+9.48%	+6.58%
GMV	+8.53%	+143.09%	+ 145.74%	+61.72%
ROI	+9.17%	+161.36%	+147.79%	+65.57%



上线对比效果

Thanks!