



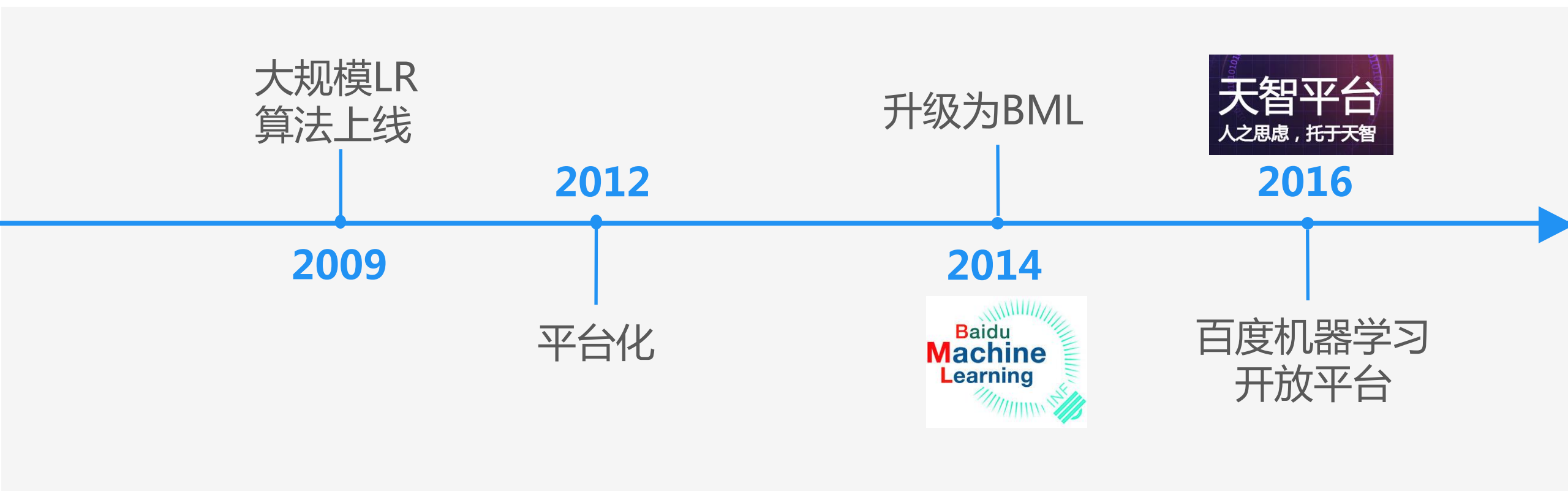
# 百度大规模推荐系统实践

百度基础架构部

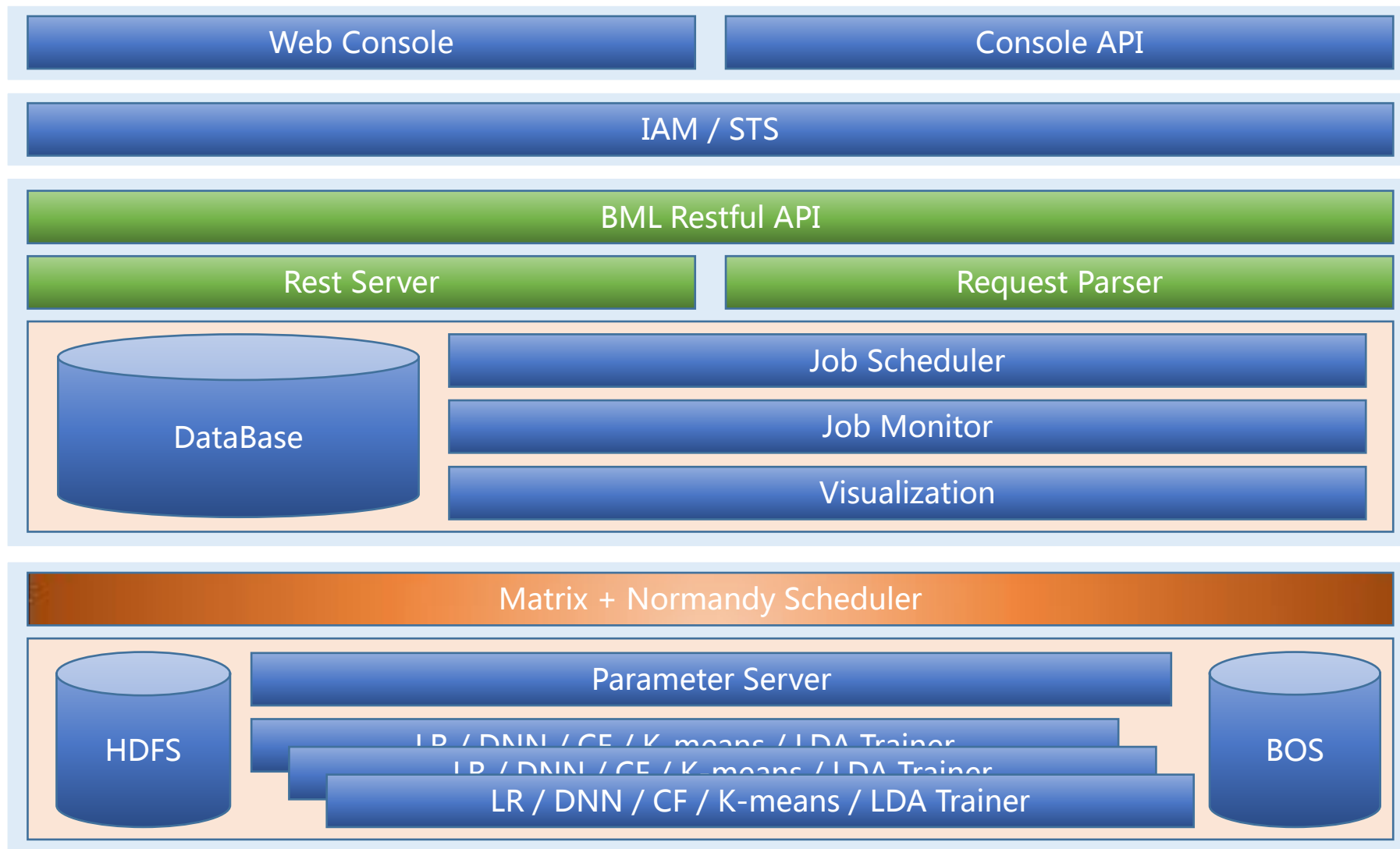
秦铎浩

- 百度机器学习开放平台历程
- 云端上的BML
- 大规模推荐系统实践

# 百度机器学习开放平台发展



# 平台整体架构



- 百度机器学习开放平台历程
- 云端上的BML
- 大规模推荐系统实践

# 云端上的BML



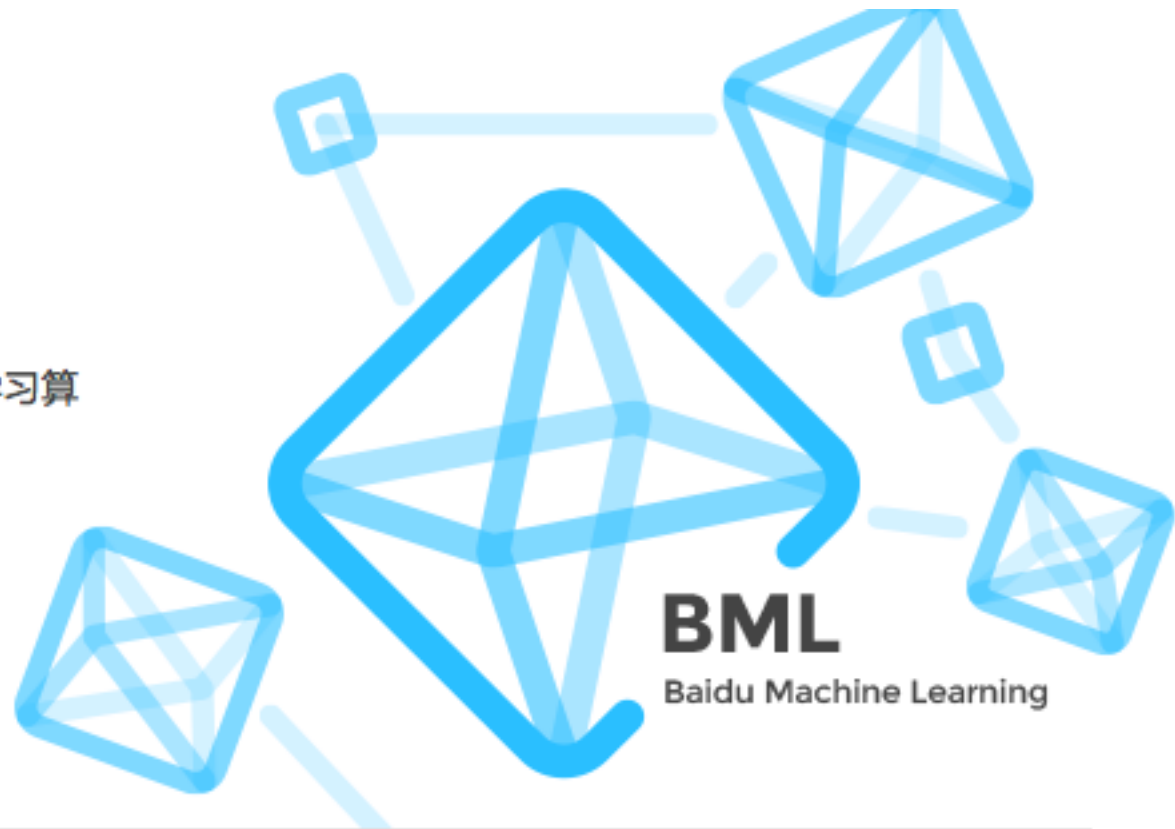
## 百度机器学习 BML

Baidu Machine Learning

百度自主研发的新一代机器学习平台，基于百度内部应用多年的机器学习算法库，提供实用的行业大数据解决方案

立即申请公测

 帮助文档 >



<https://cloud.baidu.com/>

预处理

生成数据集

分类算法

推荐算法

聚类算法

主题模型

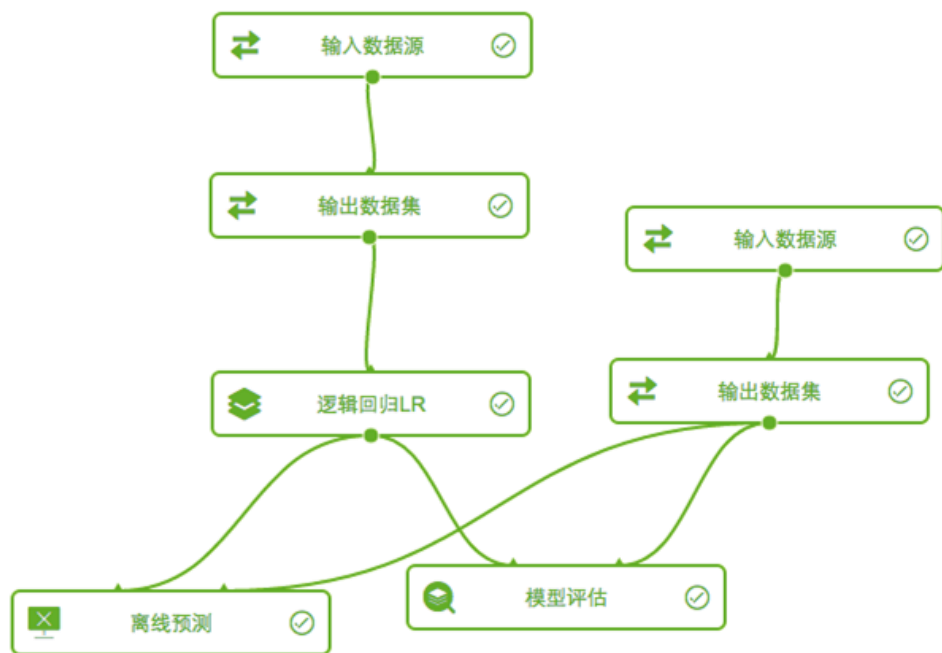
神经网络

在线学习

效果评估

可视化

# 云端上的BML



## 评估结果

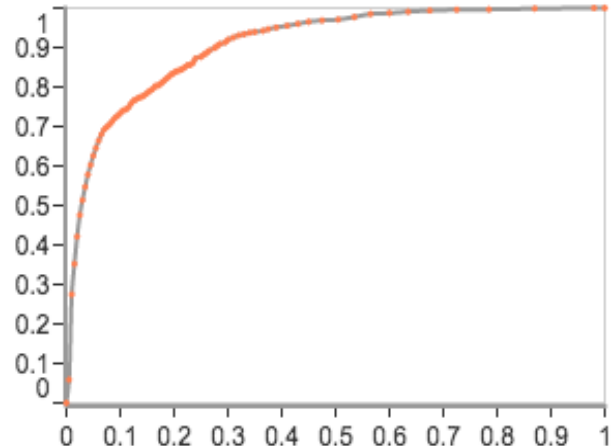
AUC= 0.91156

ROC曲线

PR曲线

## ROC曲线

True Positive Rate



False Positive Rate



- 百度机器学习开放平台历程
- 云端上的BML
- 大规模推荐系统实践

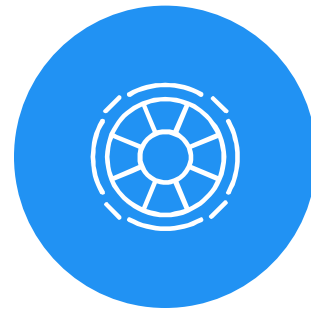
# 推荐系统必备



**干净的数据**



**高效的算法**



**完善的系统**

# 业务场景



菜品推荐

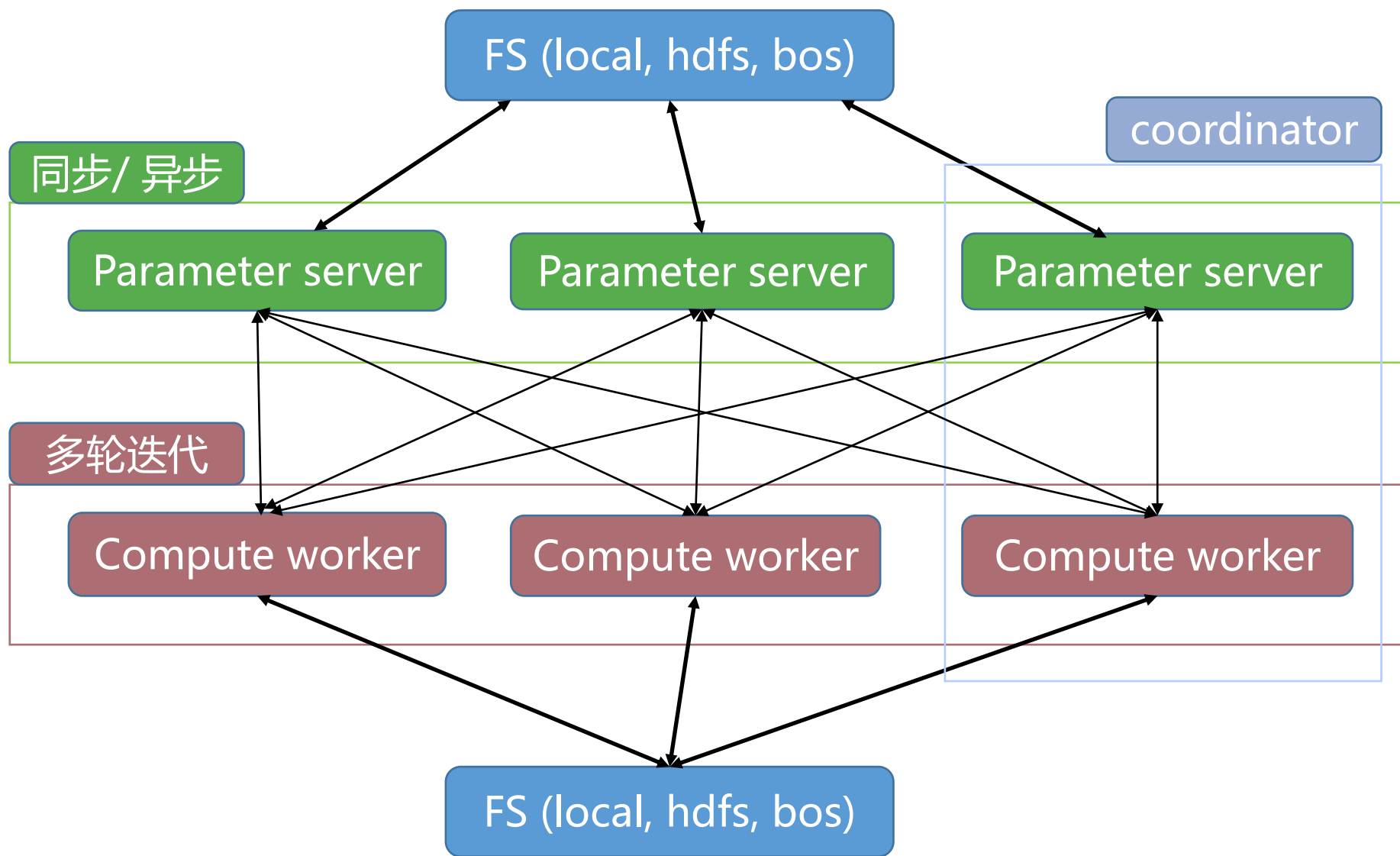


新闻推荐



搜索推荐

# ELF (Essential Learning Framework)

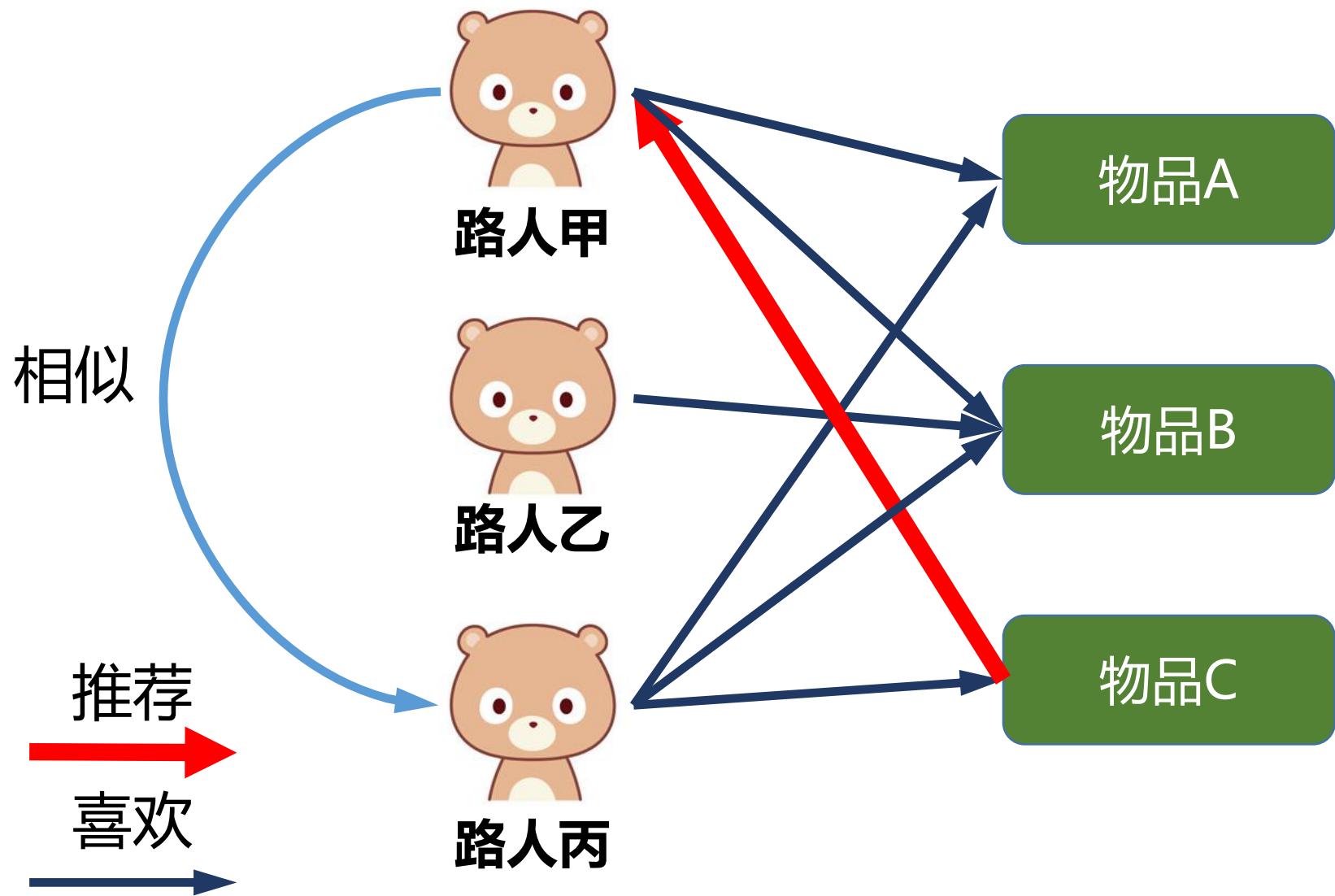


# ELF (Essential Learning Framework)

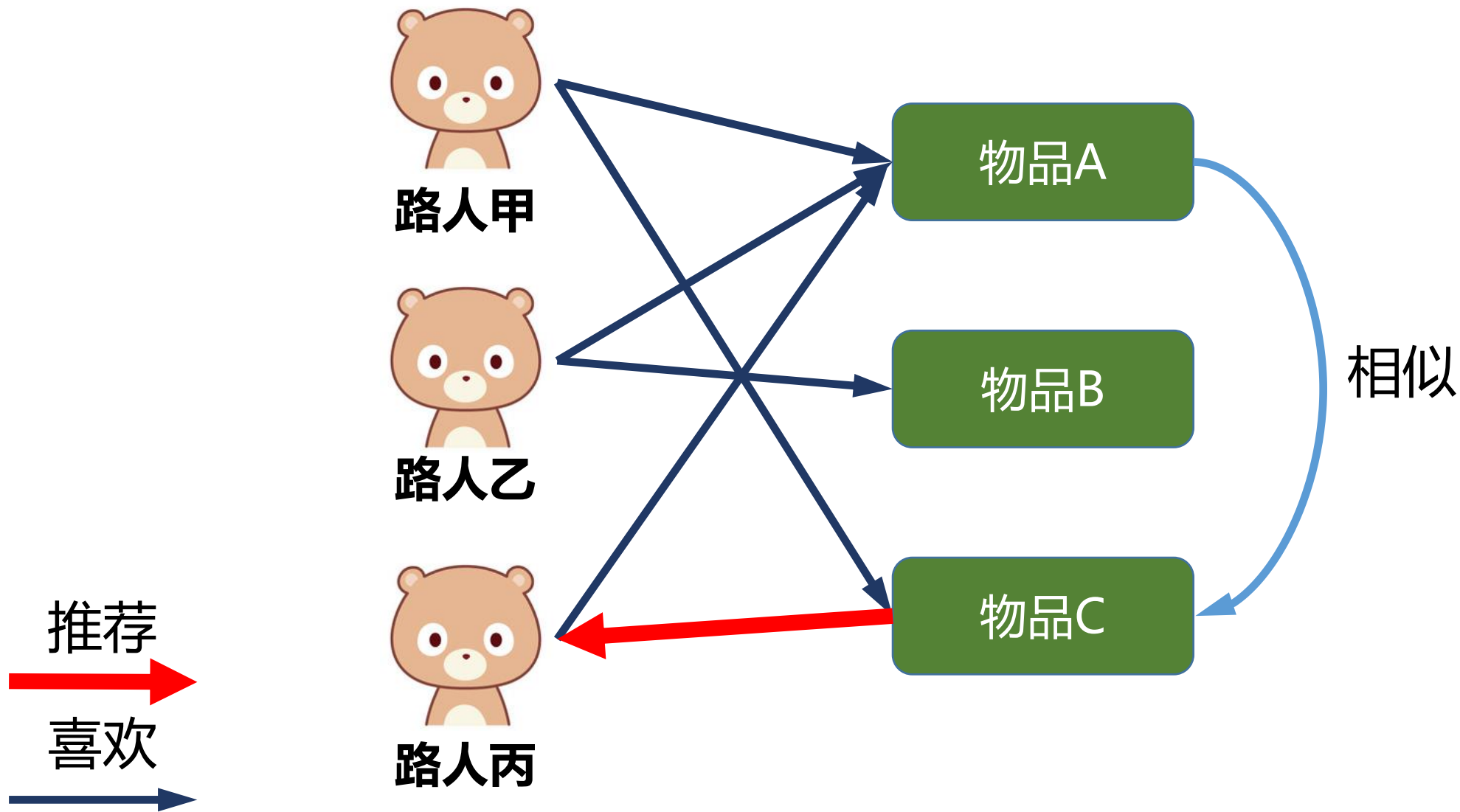


- 简单易用, 编写Async SGD LR仅需要200行代码
- 组件分布式多线程实现, 支持细粒度的线程控制
- 节点间通信依赖高效的baidu-rpc
- 深度优化hashtable, 专用于Parameter Server

# 协同过滤-基于user的CF



# 协同过滤-基于item的CF



- 收集用户偏好



隐式/显式反馈

Cosine相似度

$$sim(i, j) = \cos(\vec{i}, \vec{j}) = \frac{\vec{i} \cdot \vec{j}}{|\vec{i}| \times |\vec{j}|}$$

- 找到相似的用户或物品



欧几里得距离

$$d(i, j) = |\vec{i} - \vec{j}|$$
$$sim(i, j) = \frac{1}{1 + d(i, j)}$$

- 计算推荐

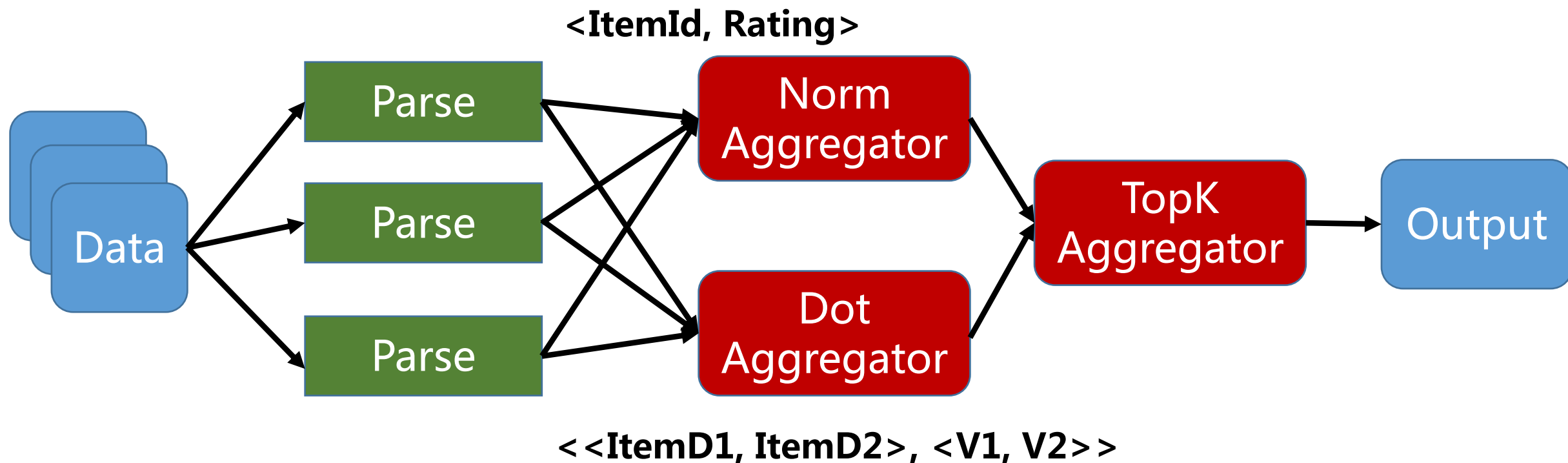


$$P_{u,i} = \frac{\sum_{all\ similar\ items, N} (S_{i,N} * R_{u,N})}{\sum_{all\ similar\ items, N} (|S_{i,N}|)}$$



	user1	user2	user3	user4	user5	user6
item1	5.0	2.0	2.0	5.0	4.0	0
item2	0	0	5.0	0	0	4.0
item3	3	2.5	0	0	3.0	0
item4	2.5	0	0	3.0	2.0	2.0
item5	0	2.0	4.0	4.5	0	0
item6	0	0	0	0	3.5	3.5
item7	0	0	5.0	0	0	4.0

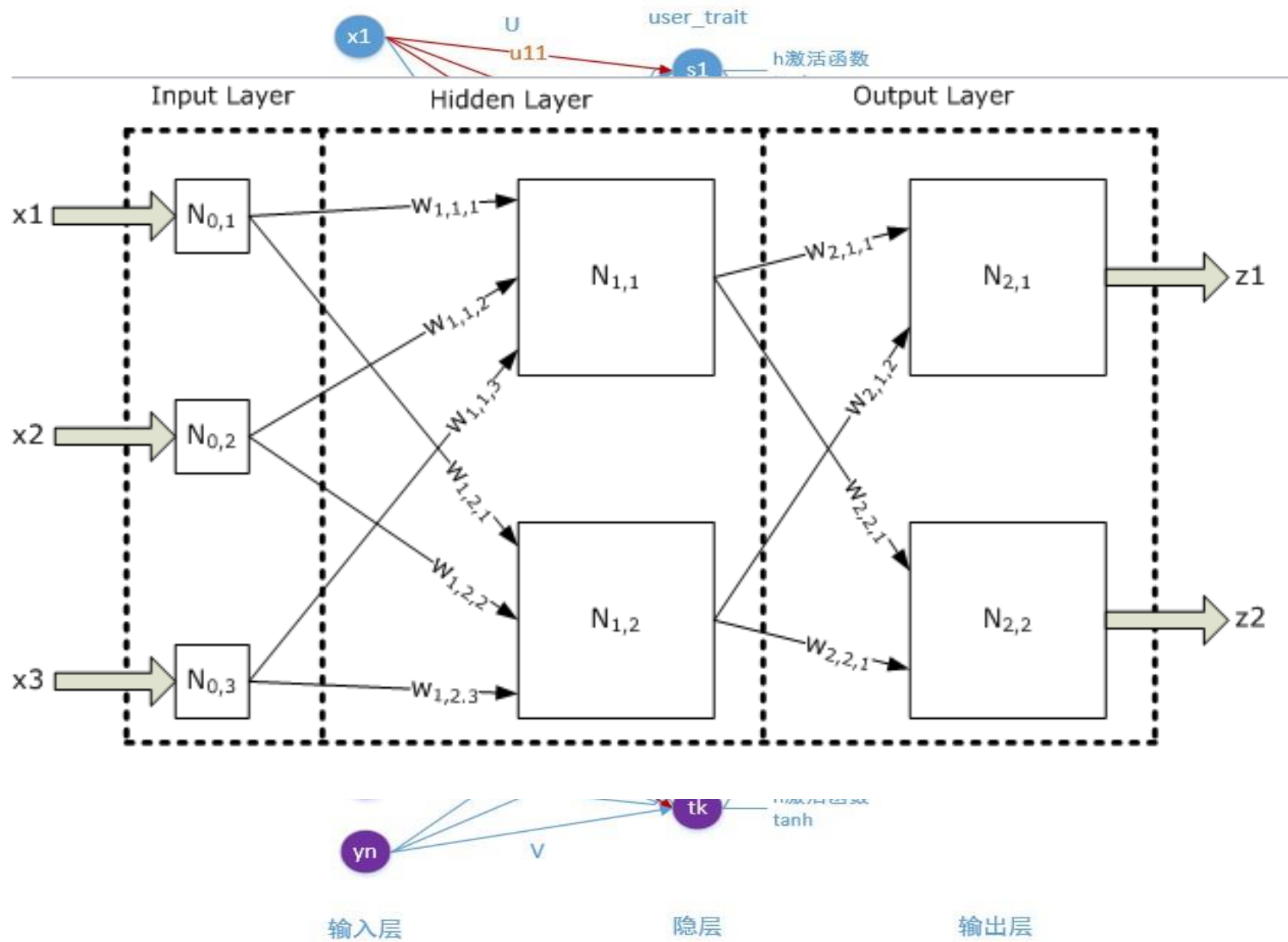
- 当item数量很大时，计算量会大到难以计算
- 没有利用到矩阵稀疏的性质



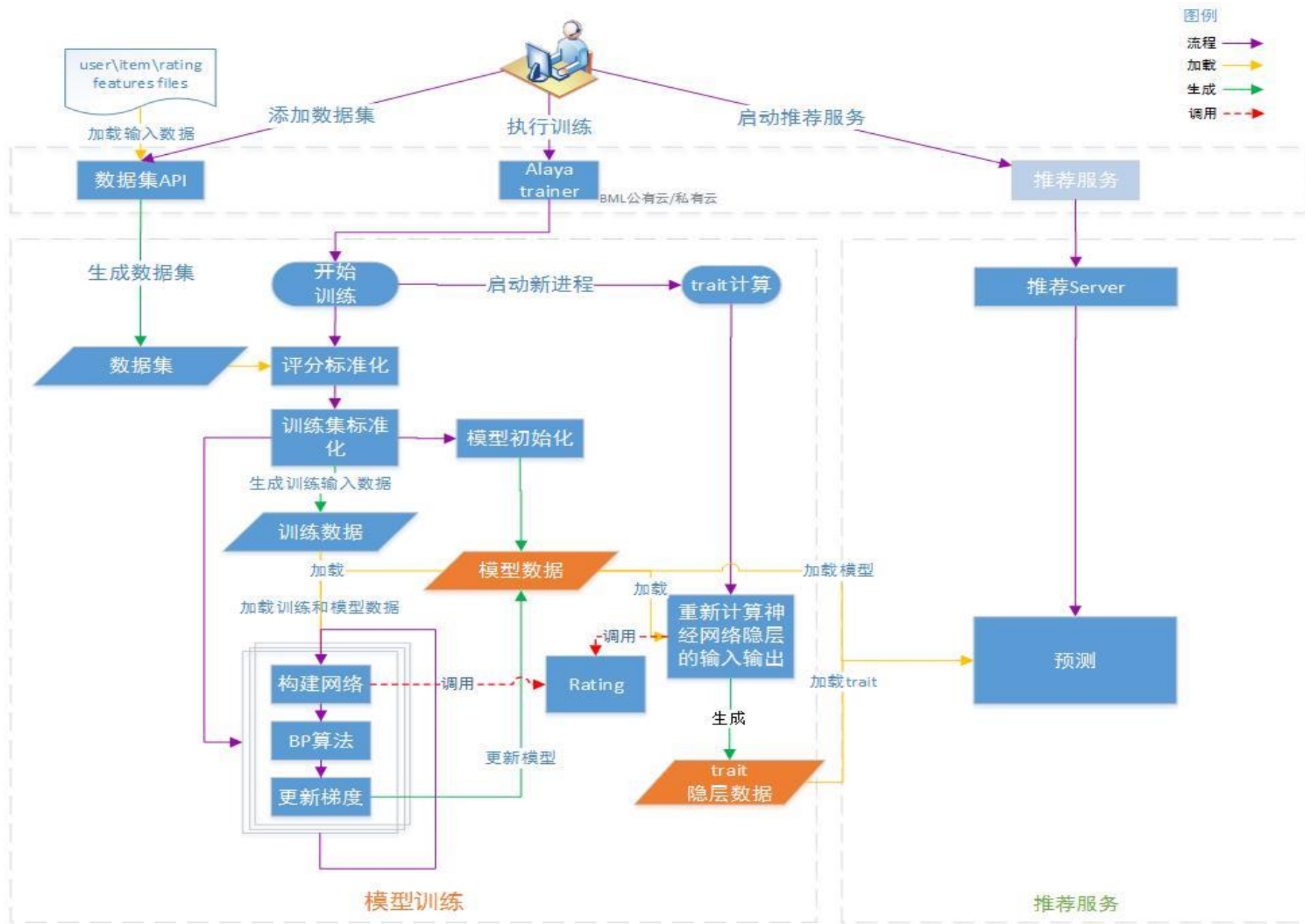
$$\cos(\vec{i}, \vec{j}) = \frac{\text{rating}_{u_1i} \times \text{rating}_{u_1j} + \dots + \text{rating}_{u_ni} \times \text{rating}_{u_nj}}{\sqrt{\text{rating}_{u_1i}^2 + \dots + \text{rating}_{u_ni}^2} \sqrt{\text{rating}_{u_1j}^2 + \dots + \text{rating}_{u_nj}^2}}$$

- 基于深度学习的推荐系统
- 支持基于用户和商品的协同过滤
- 可以使用上更多的用户信息
- 并行分布式训练

# Alaya训练



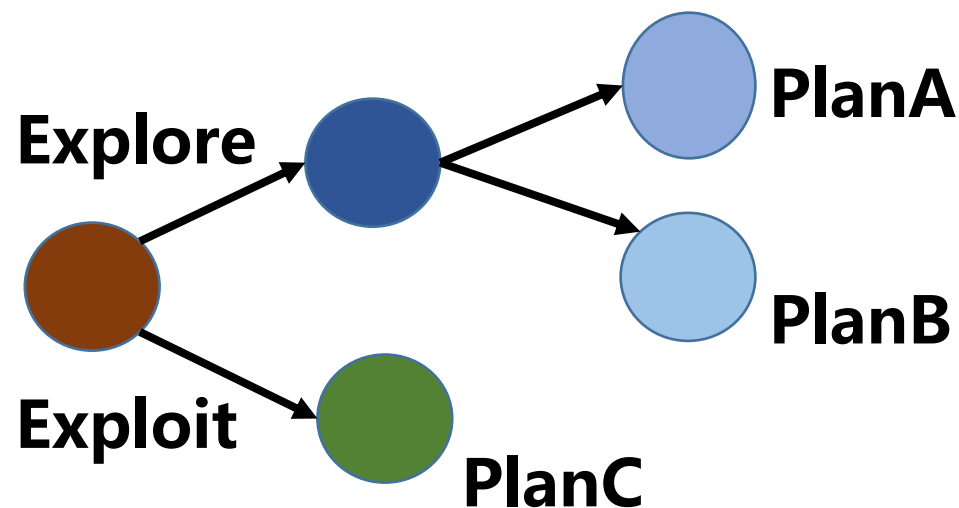
# Alaya整体架构



# Exploration and Exploitation

假设有一个广告资源池，每个广告的收益都服从一个未知的分布。怎么知道该给每个用户展示哪个广告，从而获得最大的点击收益？

- **Epsilon-Greedy方法**



- **UCB方法**

- 先在前K轮每个Plan都执行一次
- 在后面的N轮，选值最大的那个

$$\bar{x}_j + \sqrt{\frac{2 \ln n}{n_j}}$$

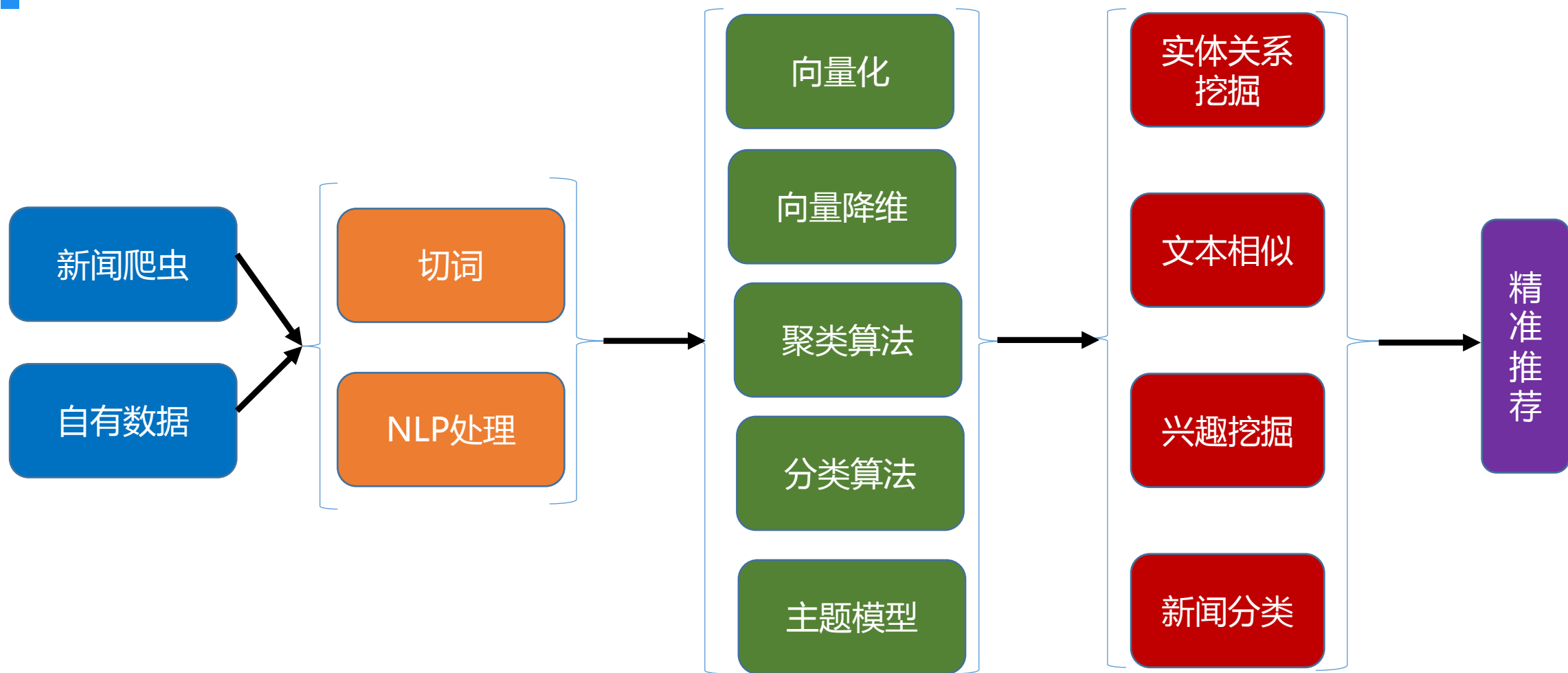
# Exploration and Exploitation

- **LinUCB**

$$a_t \stackrel{\text{def}}{=} \arg \max_{a \in \mathcal{A}_t} \left( \mathbf{x}_{t,a}^\top \hat{\boldsymbol{\theta}}_a + \alpha \sqrt{\mathbf{x}_{t,a}^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a}} \right)$$

- **Thompson sampling**

# 新闻推荐





# THANK YOU

cloud.baidu.com

email: qinduohao@baidu.com

