

推荐引文的高层次架构:

1. 由来源的收集与规范化
2. 主题的收集与规范化
3. 高维度的模型.

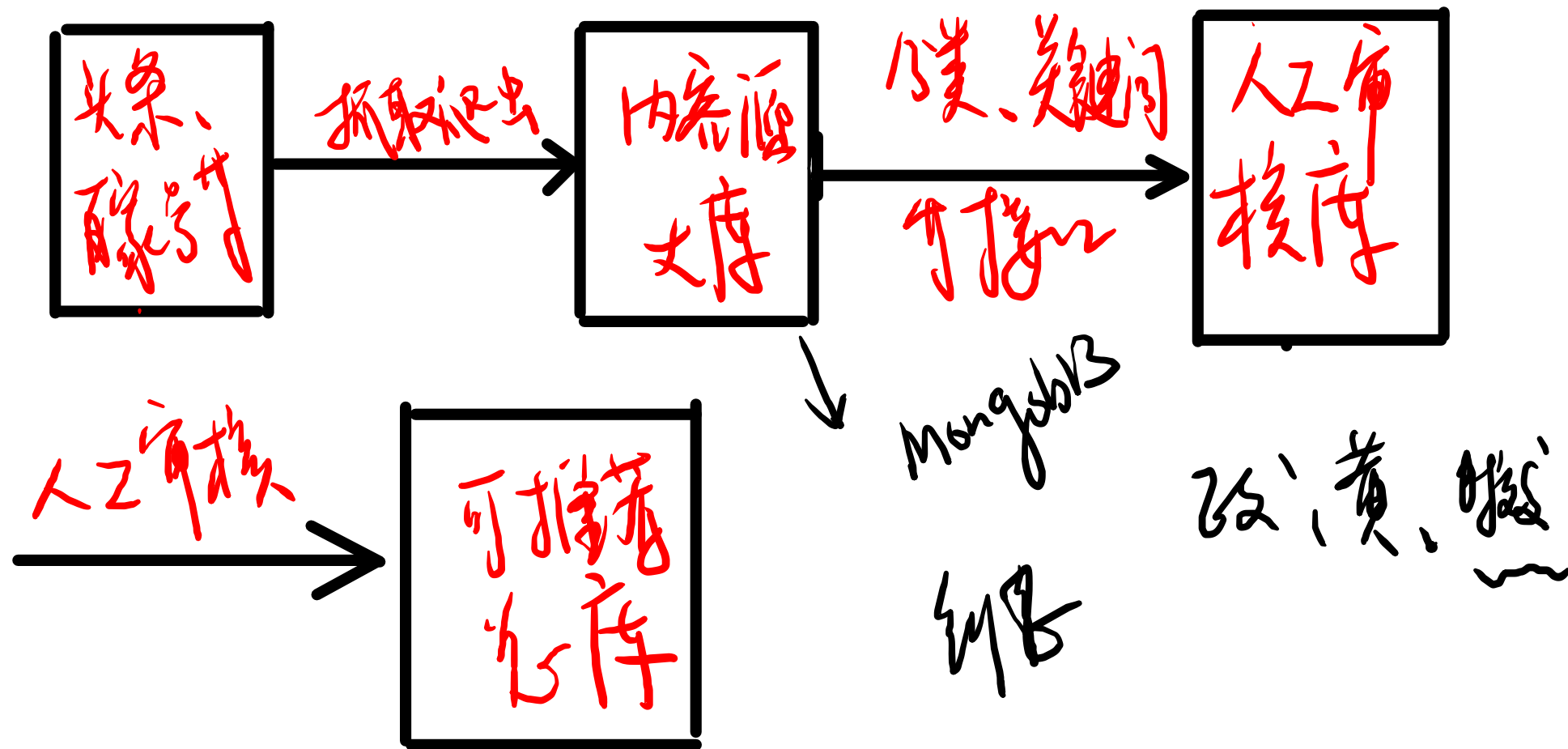
总体架构：

需求 → 设计 → 精排 → 展示

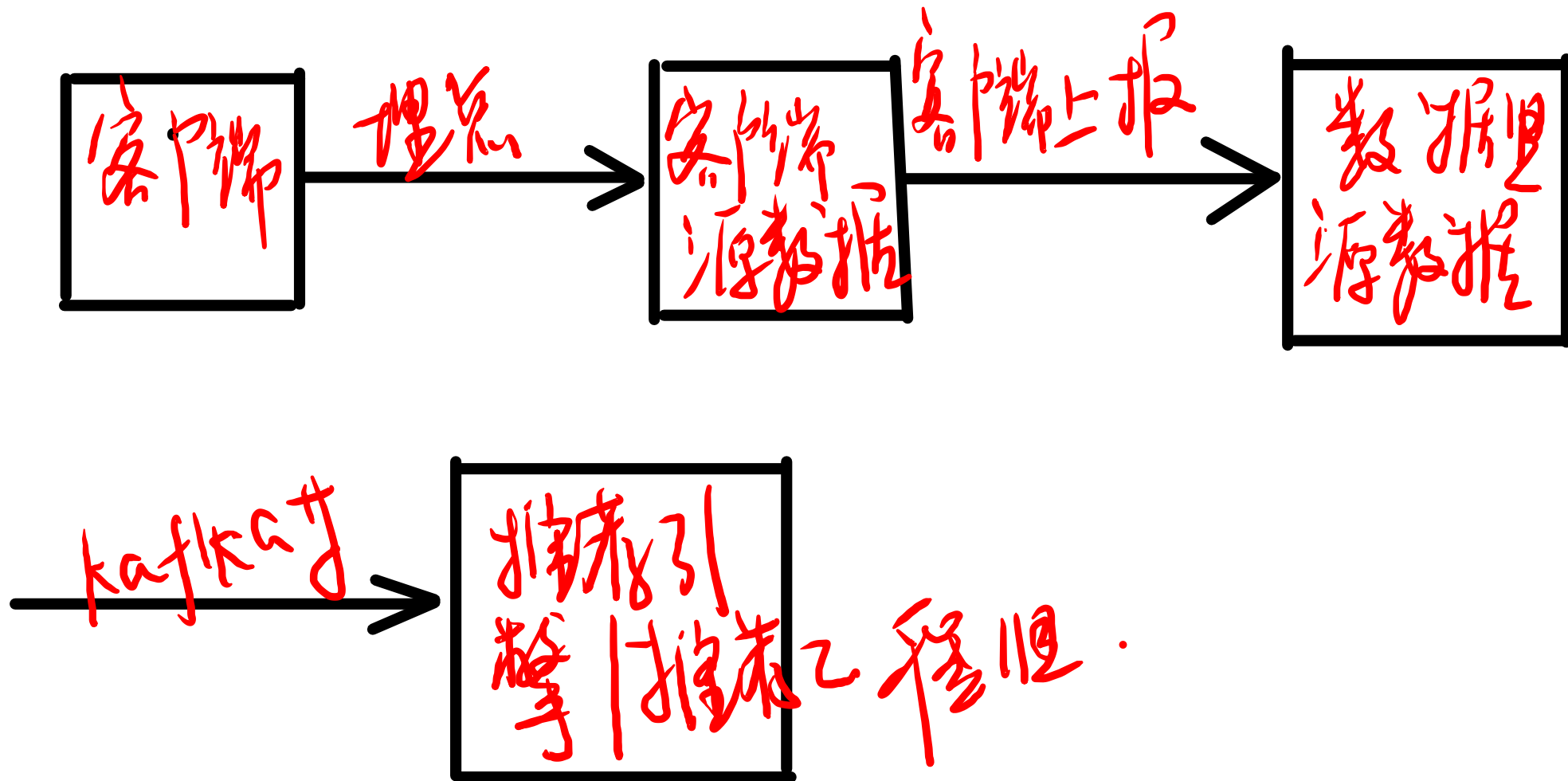
各个部分的目的：

- (一) 召回：从百万甚至千万级物品中
快速、全面的选出百万级别
的候选物品
- (二) 精排：对召回的结果进行精细化
打分

1. 内容源的收集与规范化



2. 日志接收与规范化



3. 高维训练模型.

3.1. 数据的收集.

3.2. 模型的训练与更新

3-1. 数据的收集

关系型 MySQL

源数据解析 → 存入数仓/数据库

(hdfs / hive / HBase)

→ 算法工程师 → 读取数据.

↓
(一般情况下是 hive)

3-1 数据的收集.

算法包: hive spark, mapReduce, HLL

形成数据 搬运至 hdfs 或 本地服务器

pytorch / tf / keras 于从 hdfs 或 本地
服务器读取并训练。

3.2. 模型的训练与更新

业界主流的几种方式:

- 1) 模型增量更新 (天/小时/分钟)
- 2) 模型每隔一段时间全量手动train

3.2. 模型的训练与更新

模型的全量与增量训练：比较常见。

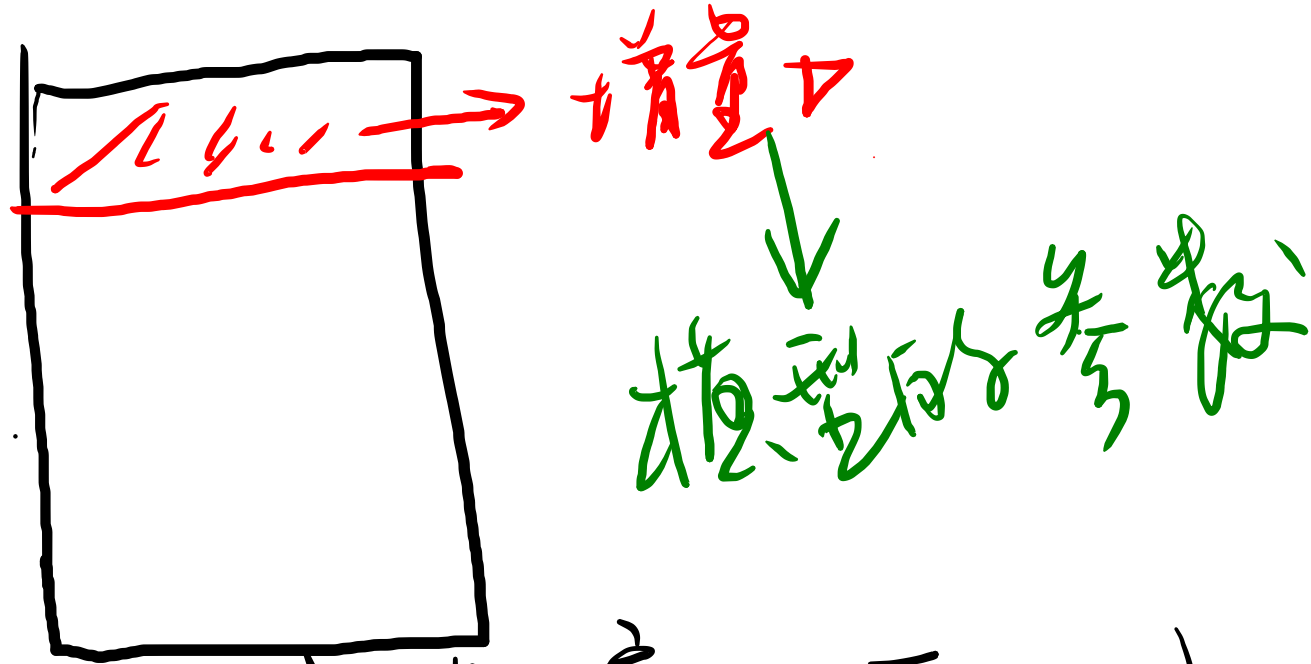
每天、每周、每天或间隔时间。

全量训练一次。

全量：重新训练

模型的增量更新.

1. 首先明确什么是增量更新



2. 更新频率: 天, 小时.

模型的增量更新.

2. 如何增量更新



后面讲模型的时候再详细讲

异步召回：

即，召回模型A 离线一直计算

基于模型A的召回结果并放入

缓存db (luredi's 等)

推荐系统架构的组成部分:

展示:

一个user请求展示推荐列表的
行为都发生了什么?

uid → 推荐引擎 → 解析user/物品
→ 请求召回模块 → 返回ves 子入db

→ 引擎 / 召回模块 →
调用精排服务 → 精排打分 结果 → 按规则
展示列表

