

# 大数据集群 运维



- hadoop集群节点6000+
- 数据容量100P+
- 日处理数据量370T+

- 最高文件数2亿+
- 日处理条数3.7万亿+
- 日作业数30万+

HDFS

YARN

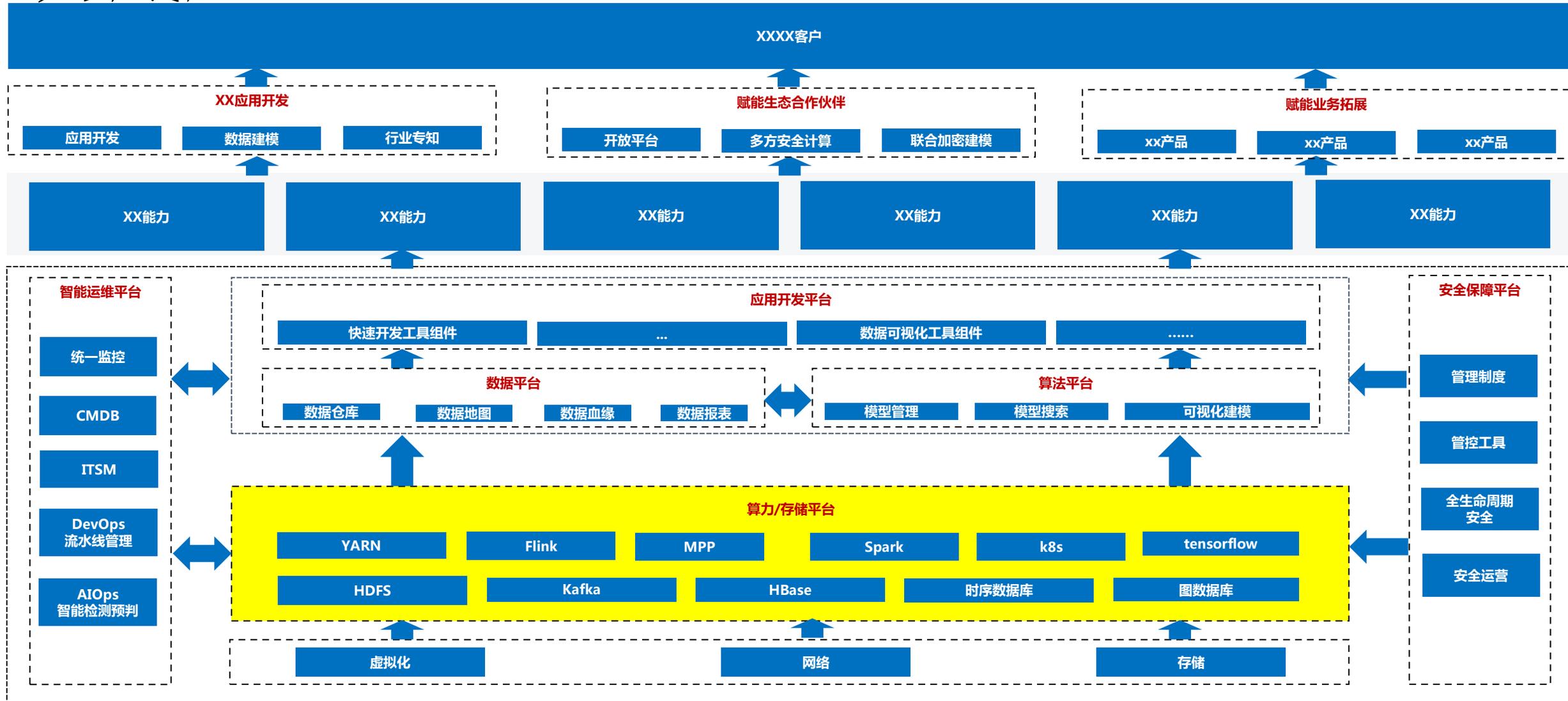
Kafka

HBase

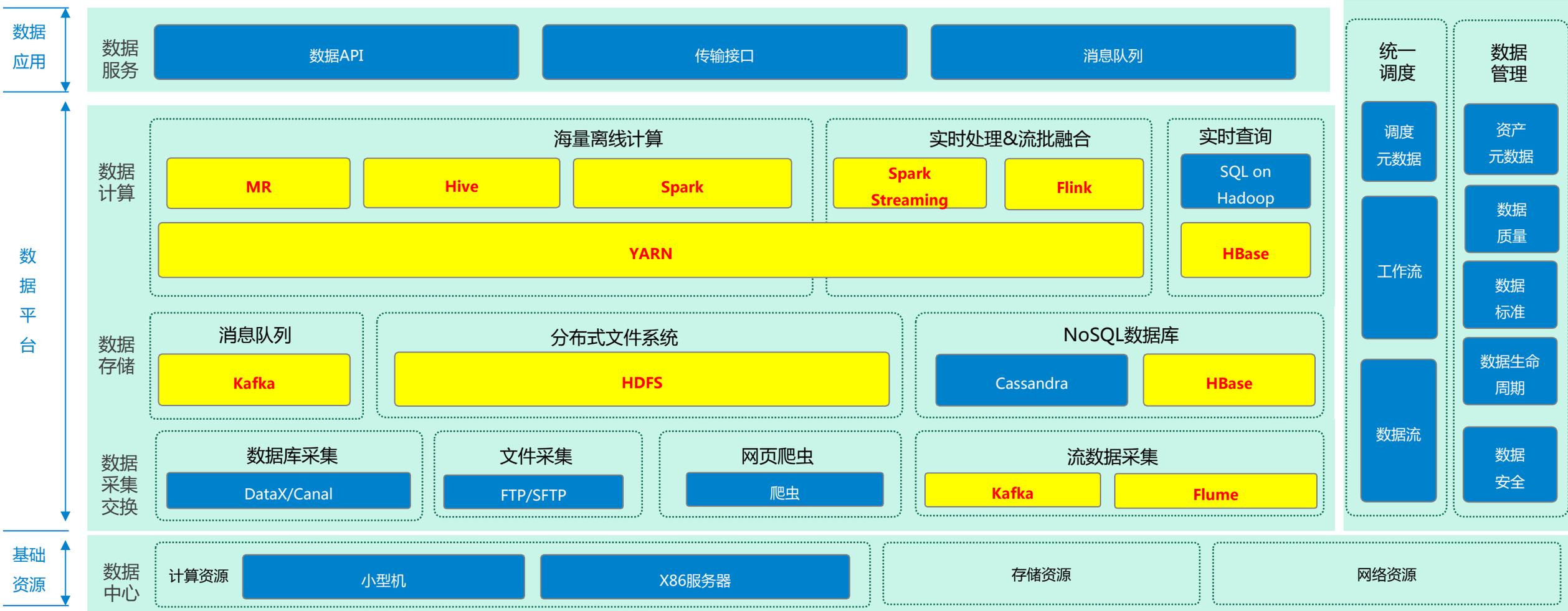
Hive

Flink

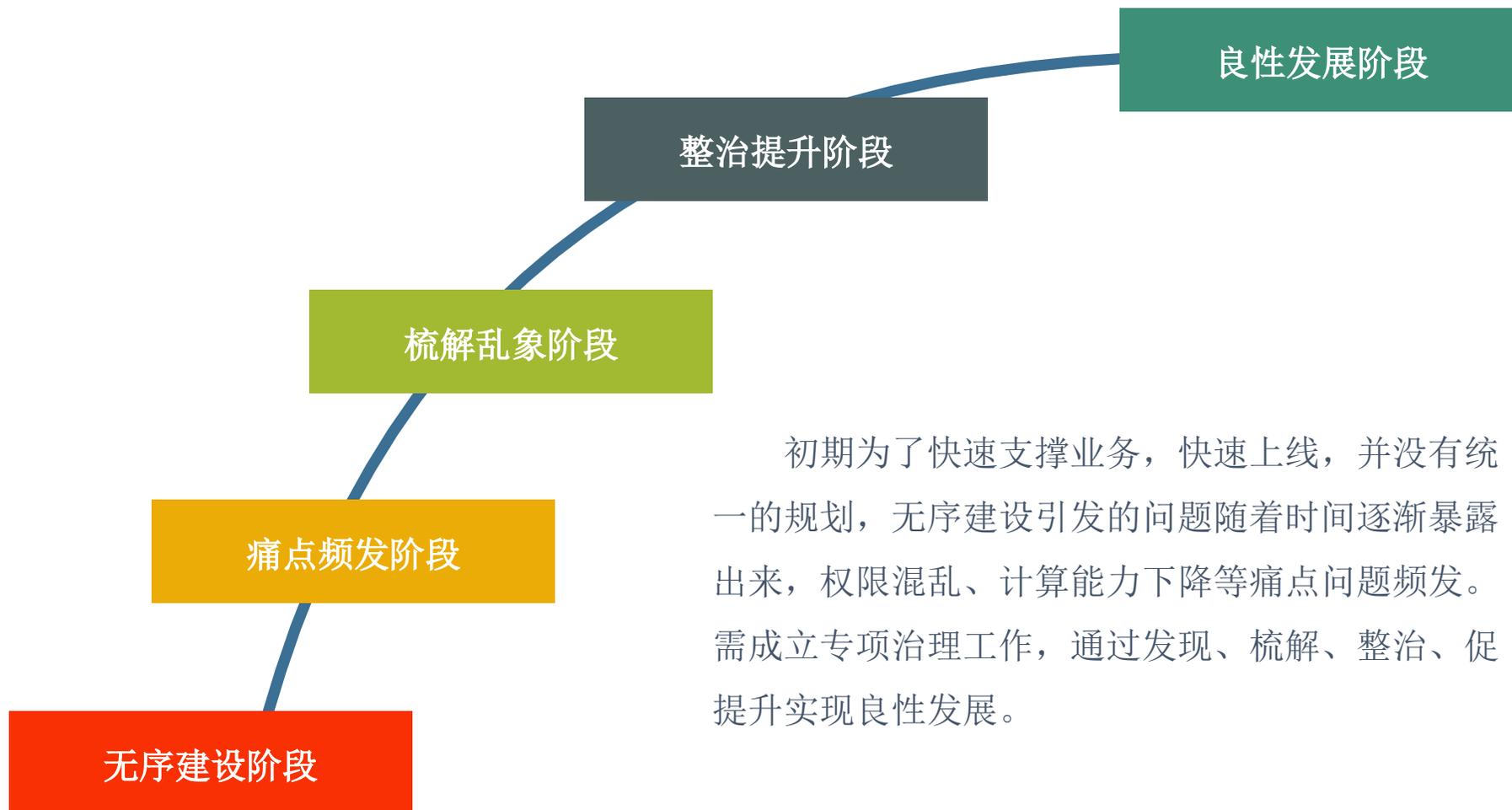
# 大数据运维的重要性-大数据平台能力底座



# 大数据常用组件在大数据项目中的位置



# 大数据项目经历的几个阶段



# 集群治理面临背景、挑战与效果

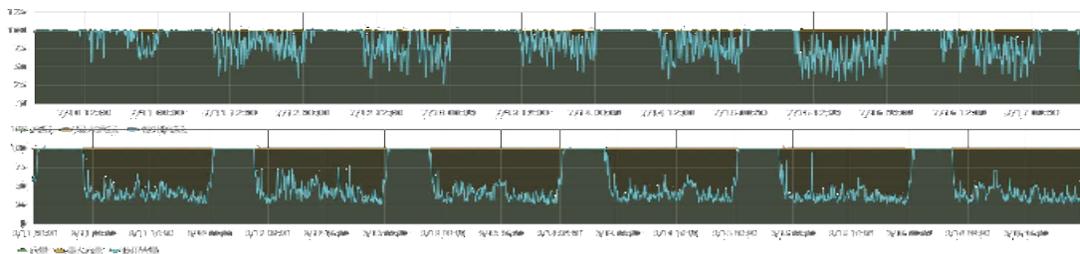
大数据公司业务高速发展过程中数据业务需求越来越复杂，所需要的算力也越来越大，进一步导致集群的规模越来越大，承担的产品也越来越多，集群面临资源负载过高、资源抢占严重、RPC请求负载过高等问题，存储系统也面临空文件过多、垃圾文件过多、小文件过多、平均文件大小过小、文件数持续增长等一系列问题，存储系统稳定性面临很大隐患，作业又面临执行耗时过长、耗资源大、数据倾斜严重等问题，直接导致数据加工异常率过高、数据具备时间有延迟风险、产品交付面临很多风险。

## 问题挑战

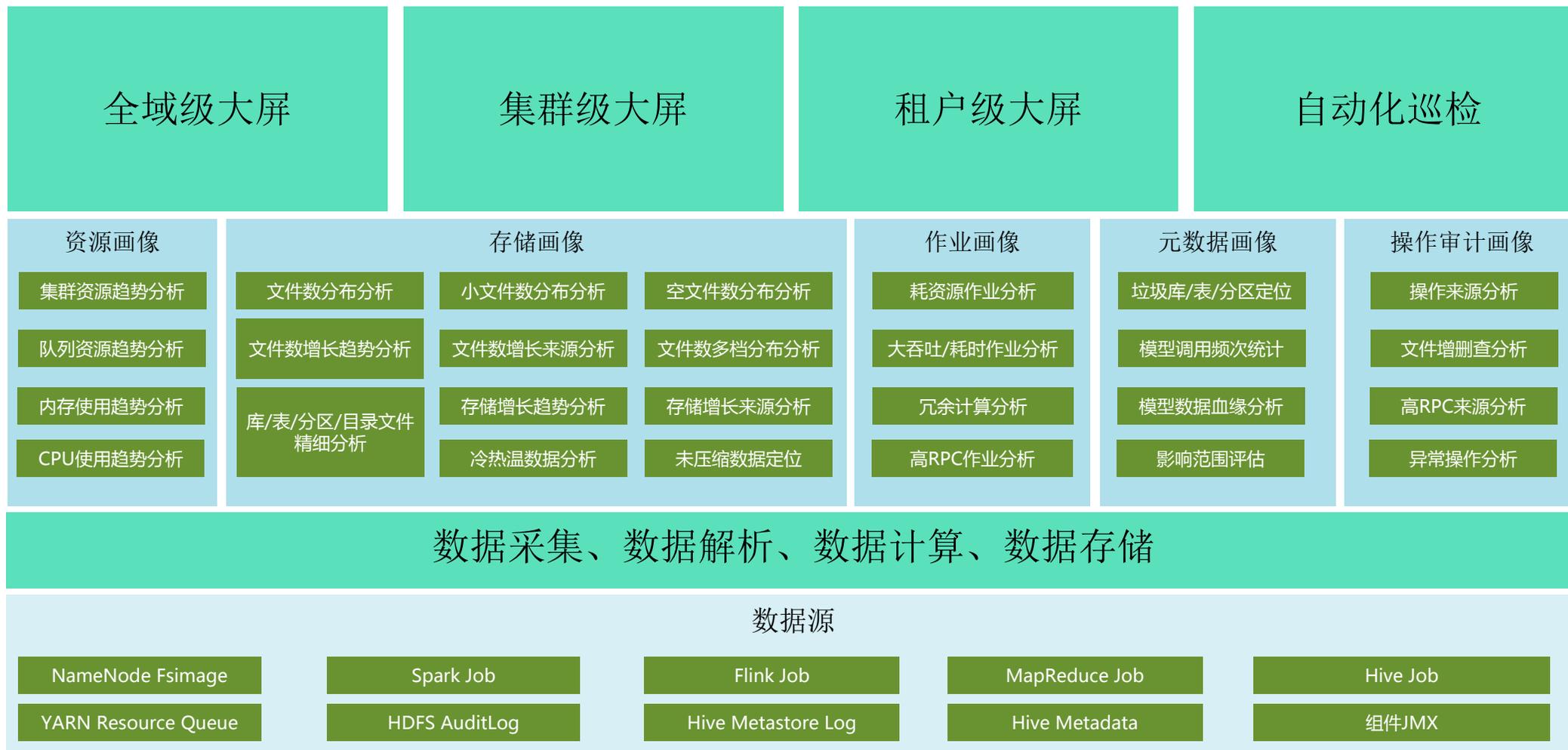
| 集群问题    | 资源问题   | 计算问题   | 存储问题   | 文件问题   | 元数据问题   |
|---------|--------|--------|--------|--------|---------|
| 集群稳定性差  | 资源负载过高 | 计算耗时间长 | 存储负载过高 | 文件数太多  | 数据库太多   |
| RPC负载过高 | 资源抢占严重 | 计算耗资源大 | 冷数据占比高 | 平均大小太小 | 数据表太多   |
| 请求响应时间长 | 打满时间过长 | 作业数据倾斜 | 未压缩数据多 | 空文件太多  | 垃圾表多    |
| 大小故障频发  | 资源浪费严重 | 冗余计算严重 | 存储持续增长 | 小文件太多  | 数据字段太多  |
|         |        |        |        | 文件夹太多  | 垃圾分区多   |
|         |        |        |        | 文件持续增长 | 数据价值不明确 |
|         |        |        |        | 冷文件太多  | 缺少数据血缘  |
|         |        |        |        |        | 重复加工多   |

制作：郎丰利1519 制作时间：2023年 睿利而行

## 治理效果

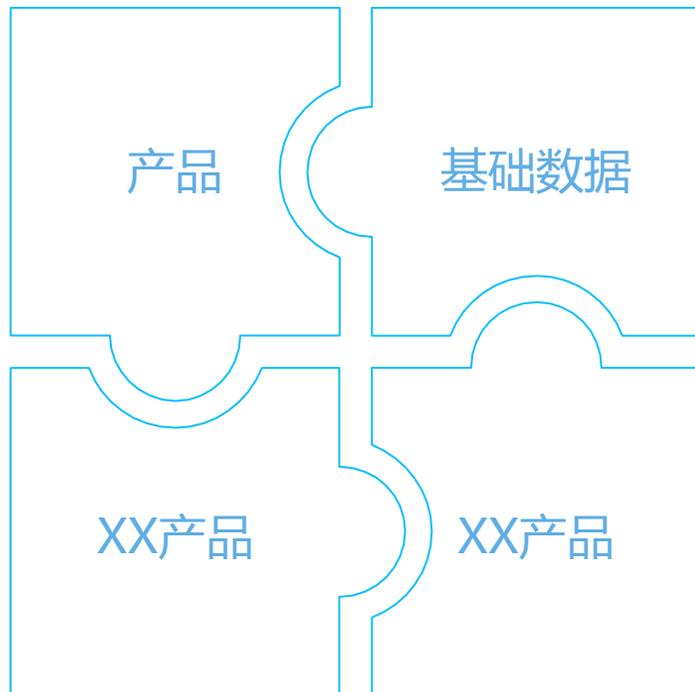


# 集群治理平台技术实现



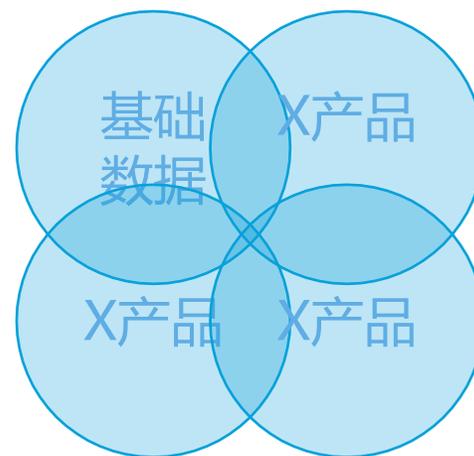
# 接口机复用严重

## 接口机应用

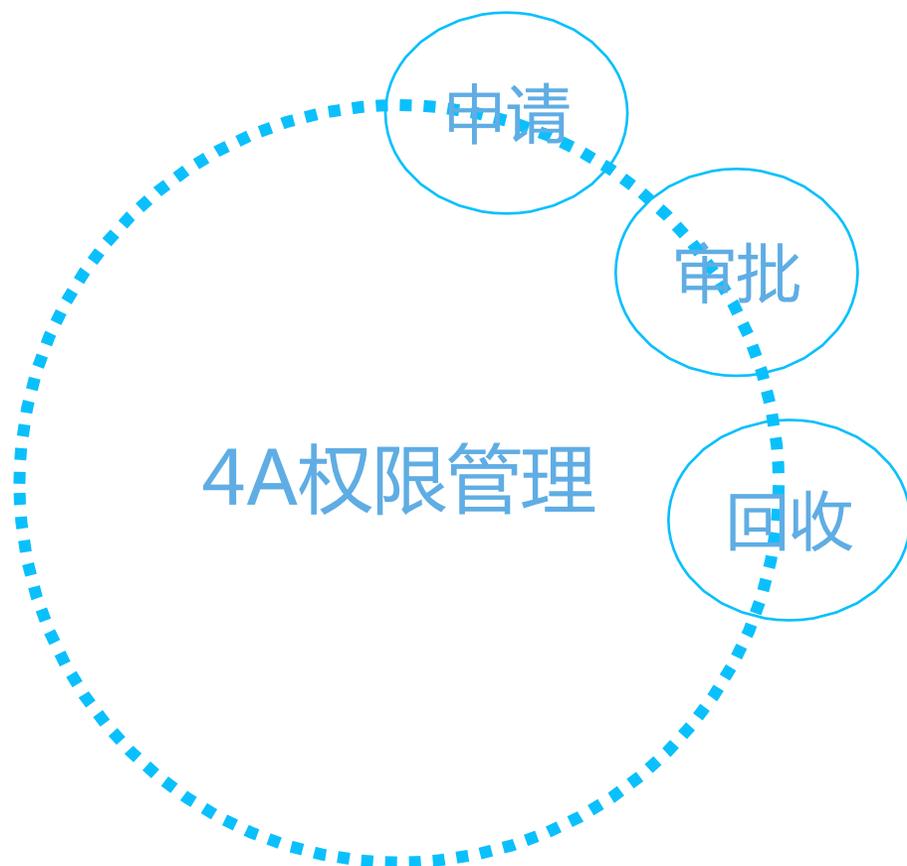


不同业务共用用户的现象非常严重，权限泛滥相当于没有权限管理，安全存在重大隐患，一旦出现问题很难排查。

- 接口机的应用主要分为四大类，基础数据的采集和加工、产品的数据加工、XX产品的数据加工以及XX产品的数据加工。现在很多接口机存在大量的交叉使用，甚至存在基础数据、产品共同使用一台接口机的情况。
- 接口机复用存在严重安全隐患，一方面服务器资源不能统筹规划，各应用会争抢资源，导致服务器不能健康运行；另一方面会导致用户权限泛滥，难以管理和监控。



# 4A权限管理宽松



## 申请

使用者申请权限时，不清楚具体资源细节或为了省事，申请时申报了过多接口机的权限。



## 审批

审批时，由于没有时间和精力对权限逐条核对，导致审批失去意义。



## 回收

目前权限回收机制依赖权限一年有效期和员工离职，导致无法及时回收人员无关的权限。

# 数据资源管理混乱

建表

使用

更新

清理

建表

所有人都有建表的权限，导致现在很多表用途不明，归属不明。

使用

基础数据表直接对外提供使用，没有采用中间层进行隔离。

更新

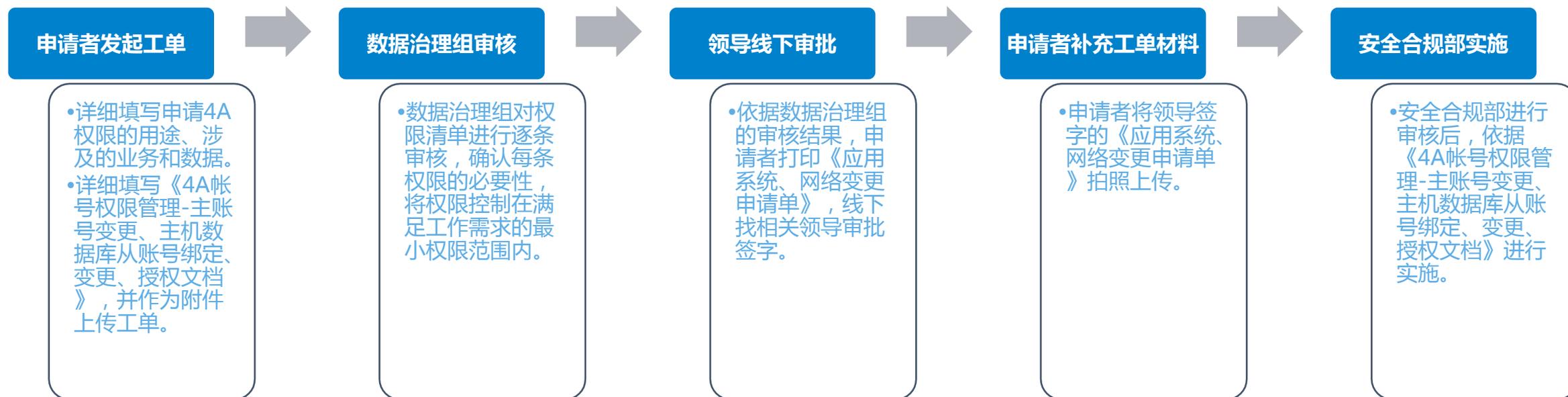
由于无法确认某张表被谁引用，所以有新需求要更新表结构时，确认的工作量非常大。

清理

尽管很多表不再被使用，但表用途不明，归属不明，不敢轻易进行清理。

# 4A权限管理规范

4A权限的申请审批流程采用ITSM工单系统



定期梳理并回收4A权限

# 库表管理规范

回收所有人建库表权限统一由数据资产管理专员统一进行操作并记录文档

DDL需求采用ITSM工单系统进行流程管理

以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：  
<https://d.book118.com/567152114142006102>