

# B站的数据治理运营框架 实践

高隆 bilibili 数据仓库工程师

# 讲师简介



高隆

//

表哥

DAMA

B站

数据成本治理

数据质量

//

# 目录

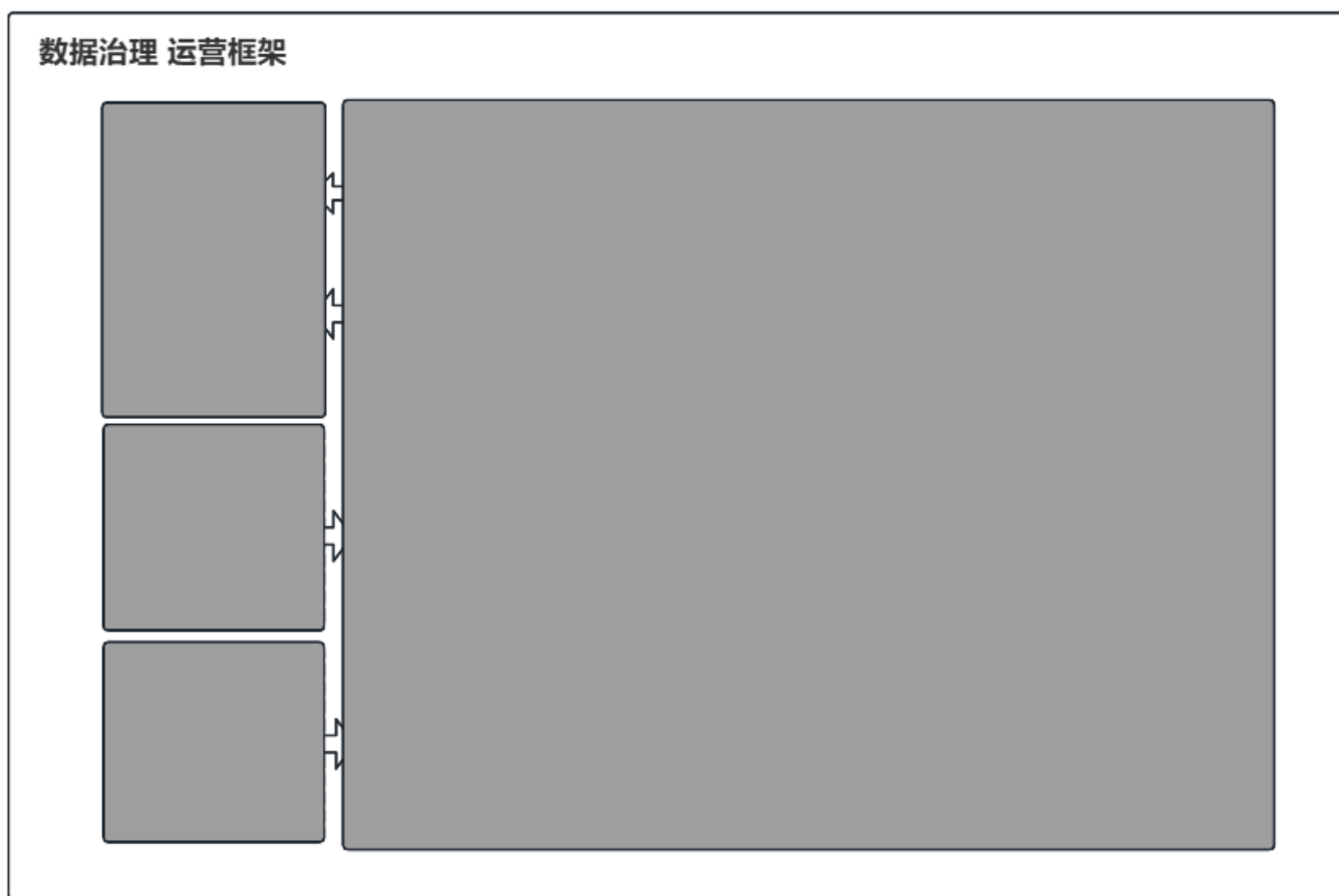
## 分析工具：DAMA-Bok

### 案例1--2022-05-11 存储水位风险

- 虚拟组织
- 嵌入治理
- 元数据管理

### 案例2--2023-10-30 数据丢失复盘

- 质量运营
- 质量的需求与满足
- 数据治理中的风险

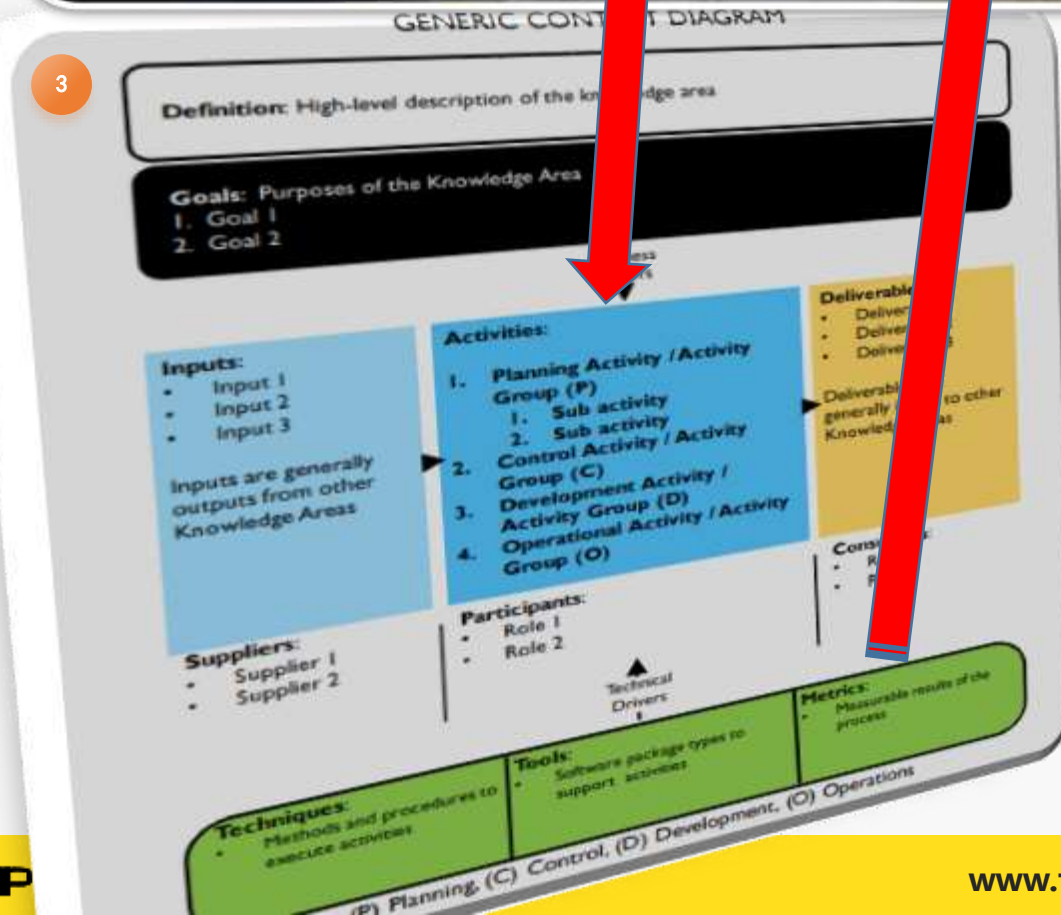




# 数据管理知识体系 DAMA-DMBOK

DAMA 成立于 1980年, 是一个全球性数据管理和业务专业志愿人士组成的非营利协会, 致力于数据管理的研究和实践

\* 左上图是买书的时候送的鼠标垫



1 车轮图: 11个知识领域

2 6边形图: 每个知识领域的7件事情

3 语境关系图: 每个知识领域都可以展开成具体的活动、方法、目标、指标

4 12原则: 采取行动和判断时的依据

# 案例1 -- 2022-05-11 存储水位风险

## 案例背景

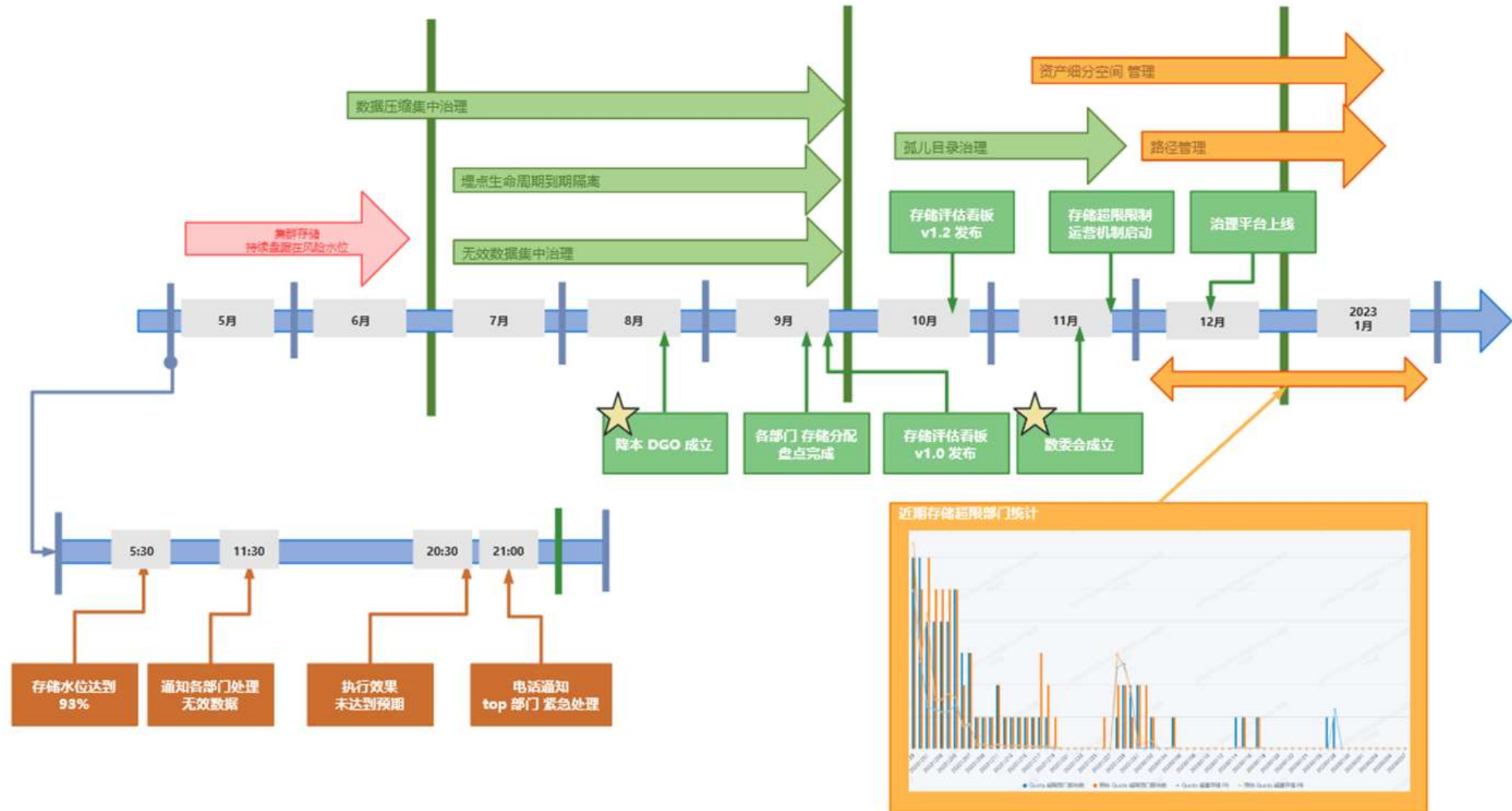


### 名词解释:

- 存储水位: HDFS集群存储
- 部门数管: 部门的“CDO”
- Quota: 部门预算资源分配
- A级数据: 多为跨部门使用数据
- Trash 数据: 一般存7日



# 案例背景





# 问题与挑战

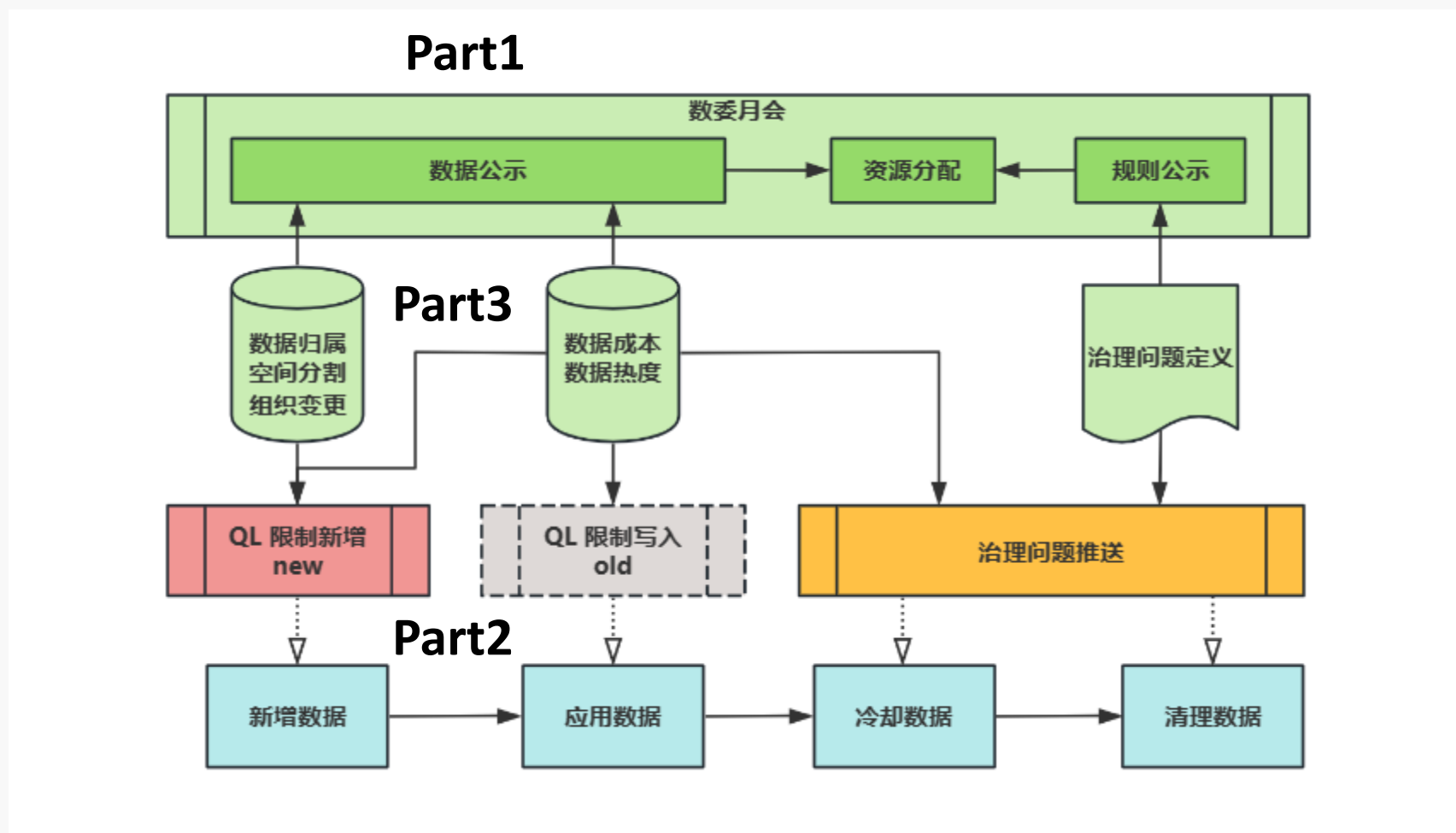
预案	执行项	问题	DAMA Bok
• 4级 (90%)	通知数管执行治理	<ul style="list-style-type: none"> <li>• 组织变更</li> <li>• 数管权责不明确</li> </ul>	<ul style="list-style-type: none"> <li>• 原则：数据管理需要领导力承诺</li> <li>• 领域：数据治理</li> </ul>
• 3级 (93%)	删除长期无访问数据	<ul style="list-style-type: none"> <li>• 没有执行驱动力</li> <li>• 删除数据存在风险</li> </ul>	<ul style="list-style-type: none"> <li>• 原则：数据价值使用经济术语表达</li> <li>• 领域：元数据</li> </ul>
• 2级 (95%)	删除 trash 调整冷数据容量	<ul style="list-style-type: none"> <li>• trash类 数据如何归属</li> <li>• 用户没有直接控制trash大小的能力</li> </ul>	<ul style="list-style-type: none"> <li>• 原则：数据管理需求驱动技术决策</li> <li>• 领域：数据存储</li> </ul>
• 1级 (97%)	根据分配限制部门写入	<ul style="list-style-type: none"> <li>• 组织变更预算归属变更</li> <li>• 限制数据写入风险极高</li> </ul>	<ul style="list-style-type: none"> <li>• 原则：数据管理是数据生命周期的管理</li> <li>• 领域：数据仓库与商务智能</li> </ul>



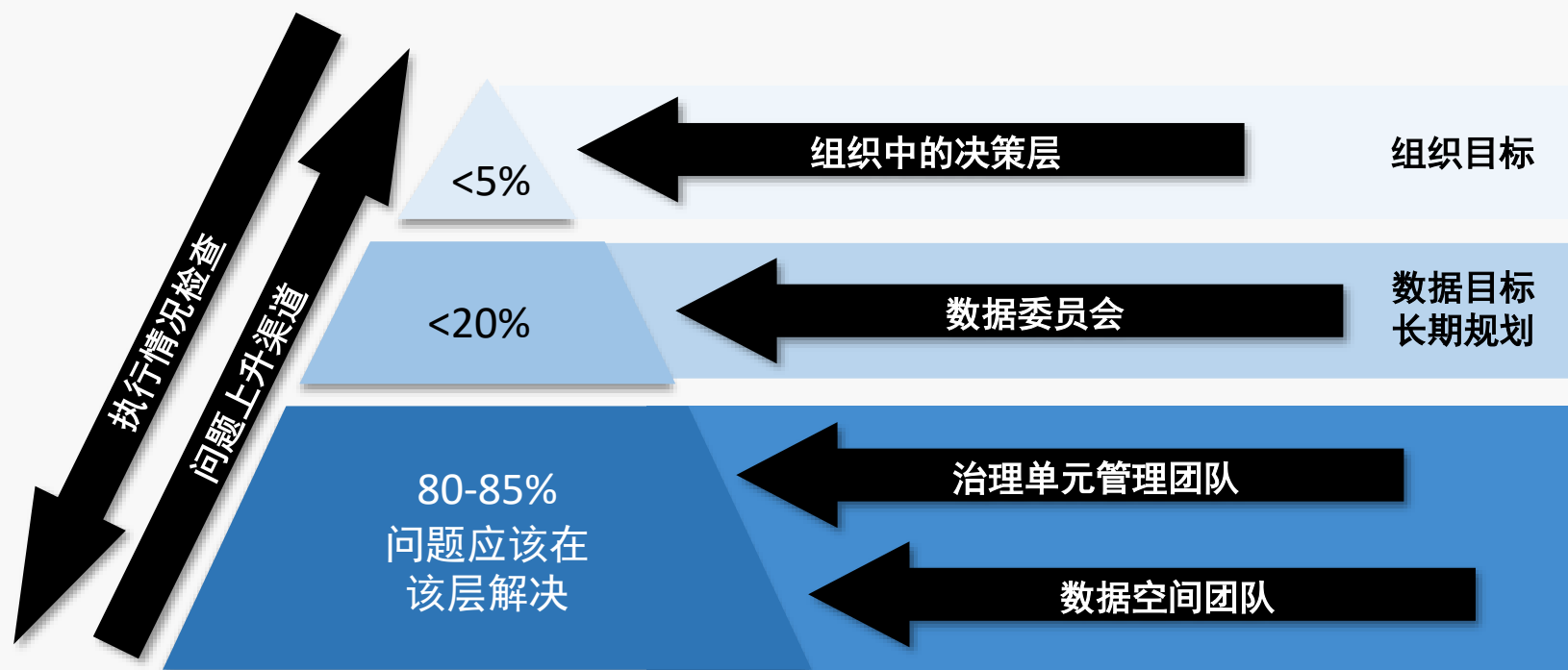
## 破题思路

	问题	DAMA Bok	方案
Part1	组织变更	<ul style="list-style-type: none"><li>原则：数据管理需要领导力承诺</li><li>领域：数据治理</li></ul>	<ul style="list-style-type: none"><li>虚拟组织<ul style="list-style-type: none"><li>数委会（人的虚拟组织）</li><li>资产空间（数据的虚拟组织）</li></ul></li></ul>
Part2	限制数据写入风险极高	<ul style="list-style-type: none"><li>原则：数据管理是数据生命周期的管理</li><li>领域：数据仓库与商务智能</li></ul>	<ul style="list-style-type: none"><li>嵌入治理<ul style="list-style-type: none"><li>预算分配</li><li>QuotaLimit（限制“新增”）</li></ul></li></ul>
Part3	没有执行驱动力	<ul style="list-style-type: none"><li>原则：数据价值使用经济术语表达</li><li>领域：元数据</li></ul>	<ul style="list-style-type: none"><li>元数据的管理与应用<ul style="list-style-type: none"><li>元数据数仓</li><li>资产账单</li><li>数据治理平台</li></ul></li></ul>

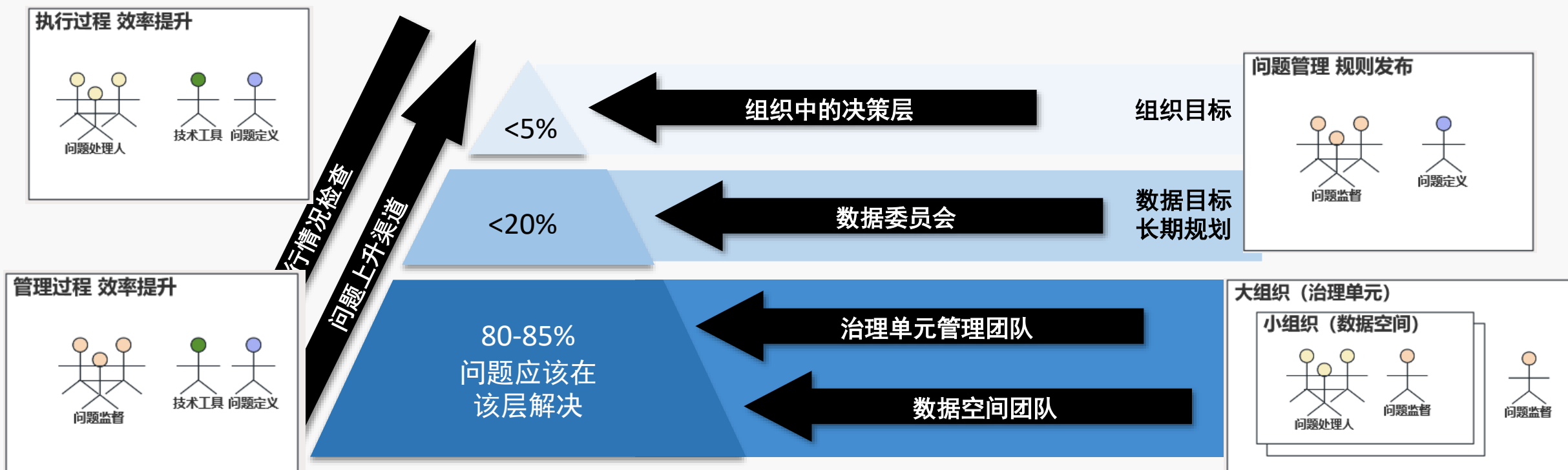
# 破题思路 – 变更方案



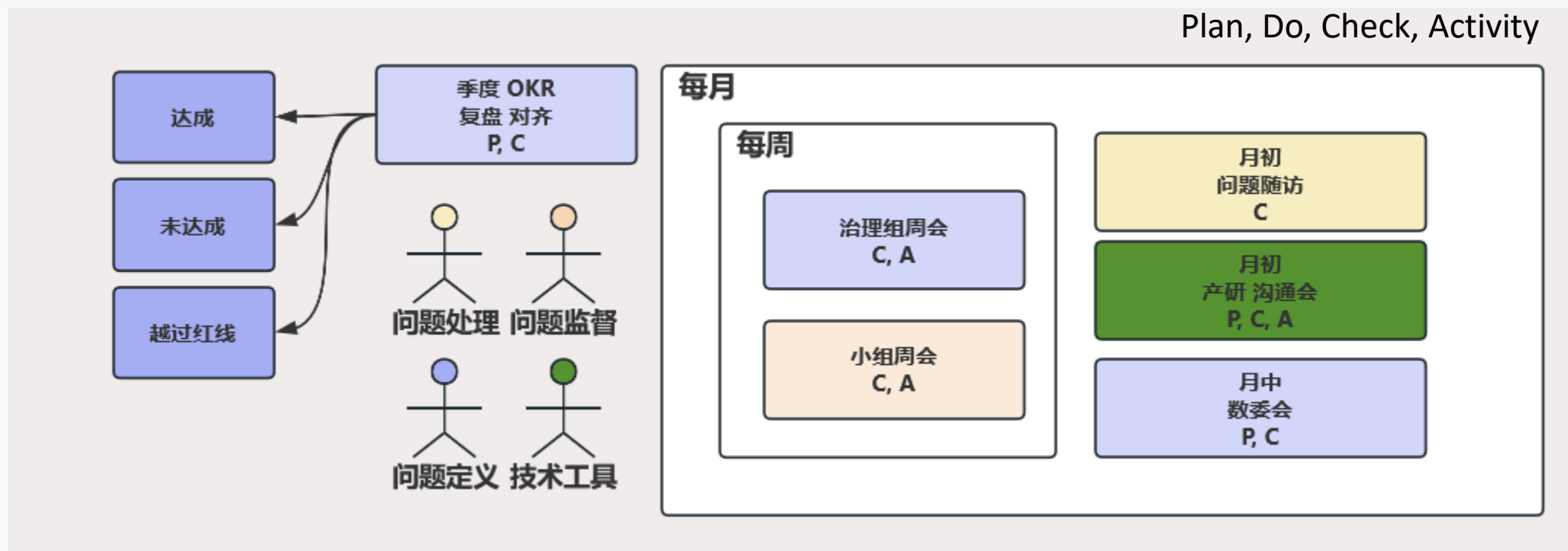
## Part1：数委会 -- 组织



# Part1：数委会 -- 角色



# Part1：数委会 -- 活动



## 季度目标

[illegible]

## 月度审计

[illegible]

## 资源分配

<b>9.7 Quota 配额审查与review</b>	<ul style="list-style-type: none"> <li>• 在员工入职时提供: A77 关于 <a href="#">配额与Quota 使用指南</a> 给员工</li> <li>• <b>web 端:</b> A-Care-7789 18871871, KPI-Care-7478 18101871, + 888 可访问, Quota <a href="#">web 端操作手册</a> 1855</li> <li>• <b>plan 端:</b> A-Care-7789 18842301, PU-Care-8478 1711871, + 888 可访问, Quota <a href="#">plan 端操作手册</a> 1855</li> <li>• 另外提供 <a href="#">MSP-Care-8478 1711871</a> 帮助解决 前端的内部 Quota 数据问题, 解决策略 给过文档和培训</li> <li>• <b>for 8 区:</b> A-Care-8478 1887187, K-Care-828 1711871, 已上线, 给过一些培训帮助解决部门内部问题, 另外给过 <a href="#">for 8 区操作手册</a></li> </ul>
------------------------------	--

## 概念说明

<p>4.0.0.0 软件 (v4.0) 说明</p>	<p>说明</p> <ul style="list-style-type: none"> <li>• 模块: 系统控制</li> <li>• 模块: 1.1.1.1 工具: 系统控制 (System Control)</li> </ul> 	<p>说明</p> <p>模块名称: 系统控制</p> <ul style="list-style-type: none"> <li>• 工具名称: 系统控制 (System Control)</li> <li>• 工具名称: 系统控制 (System Control)</li> <li>• 工具名称: 系统控制 (System Control)</li> </ul>															
<p>系统控制模块</p>		<p>系统控制模块</p> <table border="1"> <thead> <tr> <th>系统控制模块</th> <th>系统控制模块</th> <th>系统控制模块</th> </tr> </thead> <tbody> <tr> <td>System Control</td> <td>System Control</td> <td>System Control</td> </tr> <tr> <td>System Control</td> <td>System Control</td> <td>System Control</td> </tr> <tr> <td>System Control</td> <td>System Control</td> <td>System Control</td> </tr> <tr> <td>System Control</td> <td>System Control</td> <td>System Control</td> </tr> </tbody> </table>	系统控制模块	系统控制模块	系统控制模块	System Control	System Control	System Control	System Control	System Control	System Control	System Control	System Control	System Control	System Control	System Control	System Control
系统控制模块	系统控制模块	系统控制模块															
System Control	System Control	System Control															
System Control	System Control	System Control															
System Control	System Control	System Control															
System Control	System Control	System Control															
<p>系统控制模块</p>	<ul style="list-style-type: none"> <li>• 系统控制模块</li> <li>• 系统控制模块</li> </ul>																

## 组织调整

[illegible]

## 红线规则

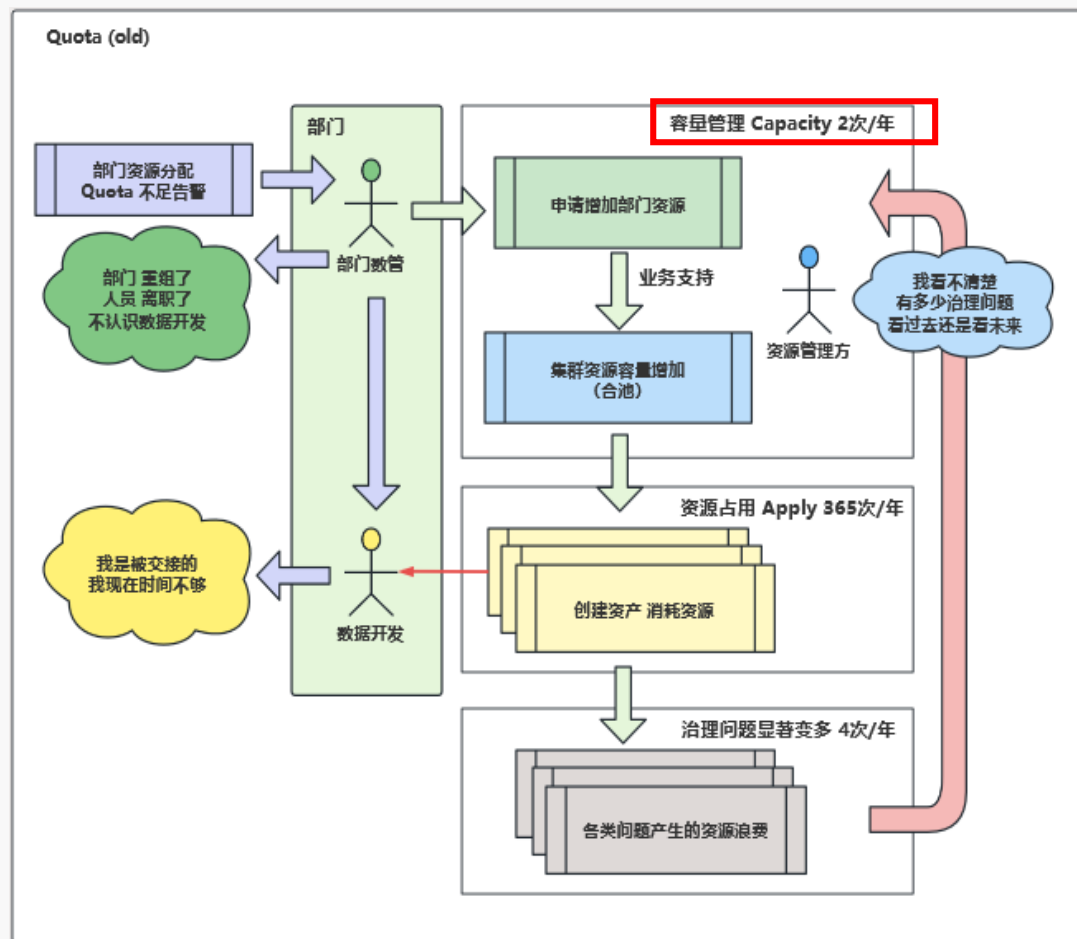
**3.1 数据源**

**3.2 数据源**

# Part2：嵌入治理

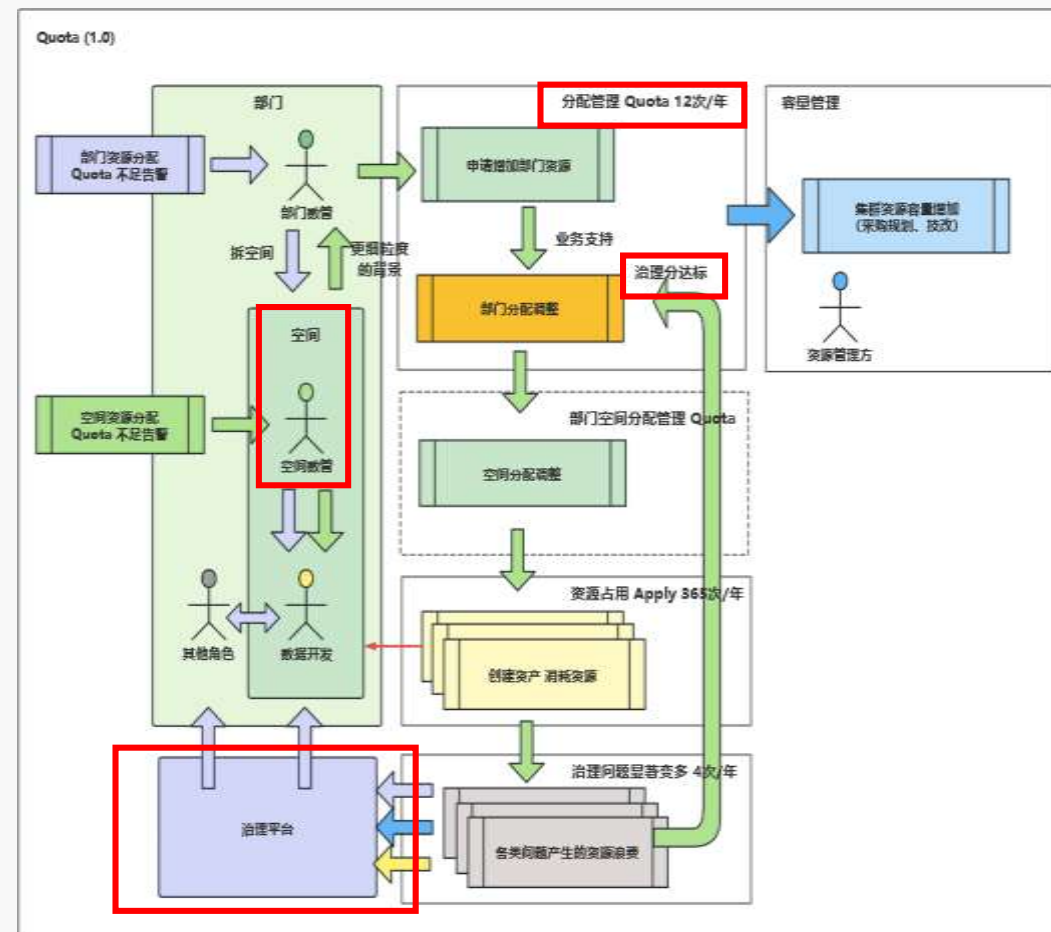
Quota – old (till 2022.H1)

2次/年，组织架构变更



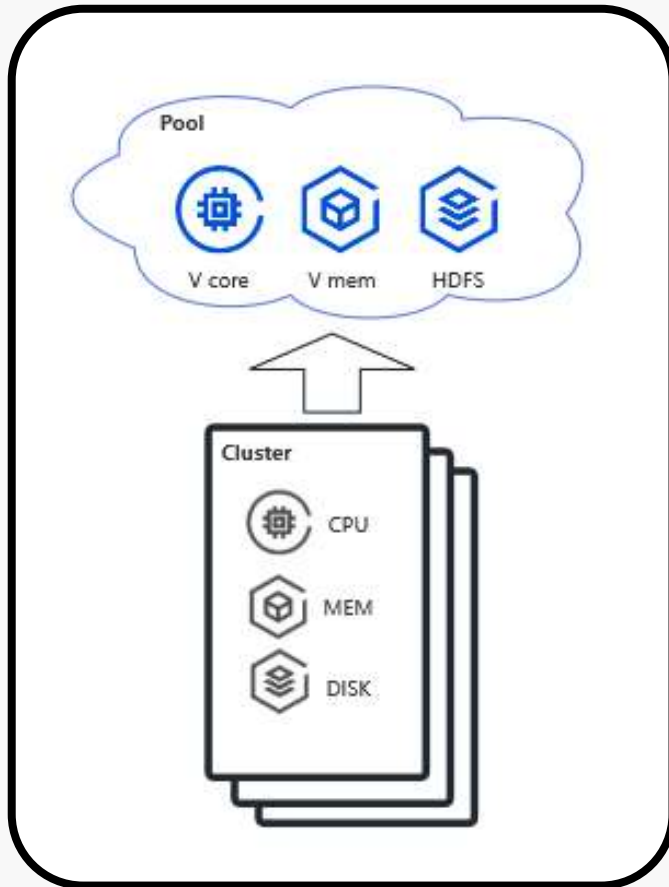
Quota – v1.0

数委会=12次/年，空间=中间组织，待治理问题=备用资产

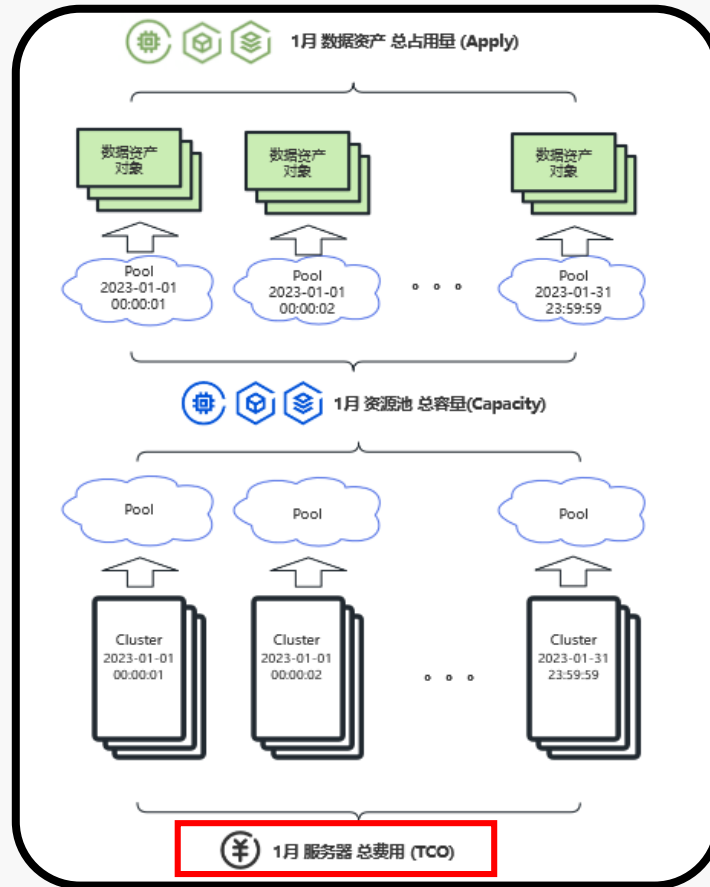




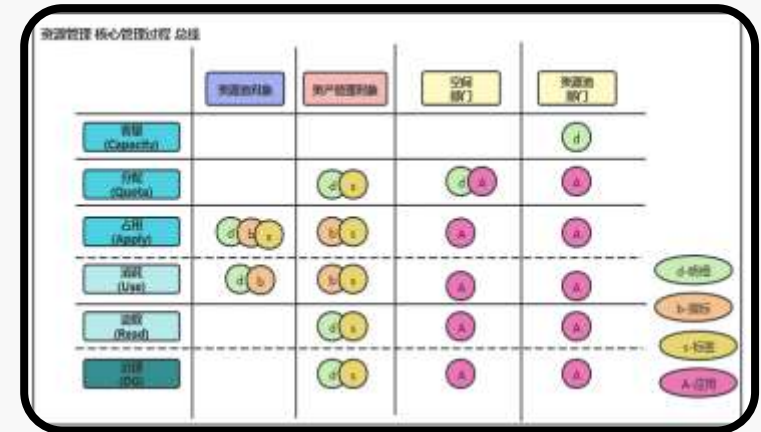
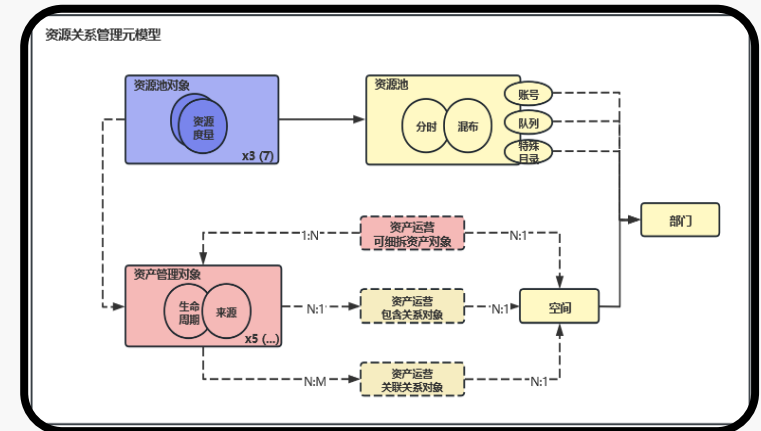
## Part3：元数据的应用 -- 资产账单



01 物理资源→虚拟资源



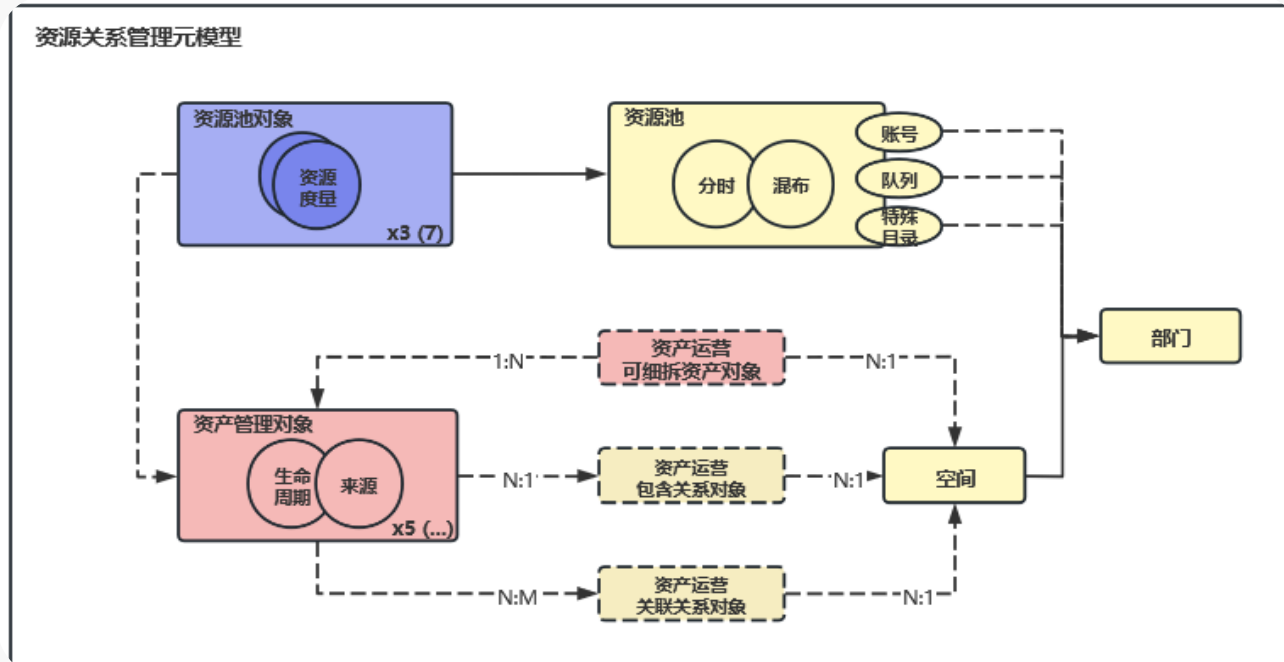
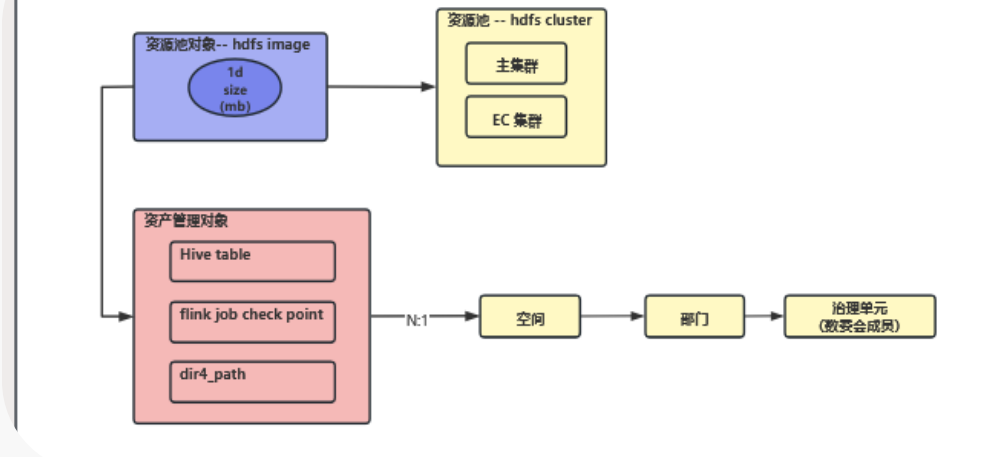
02 资源分配 + 运营成本



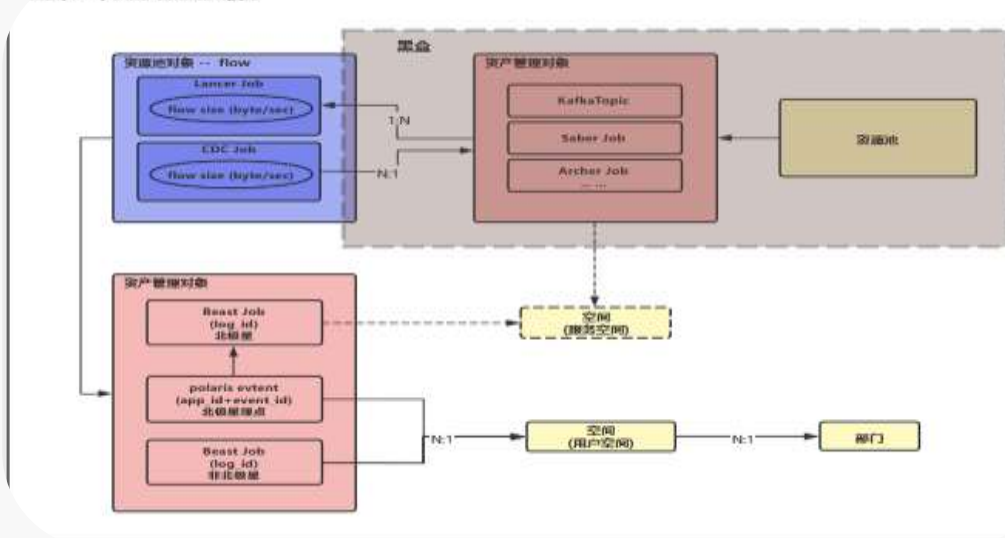
03 资源管理 元数据

# Part3：元数据的管理 -- 元模型

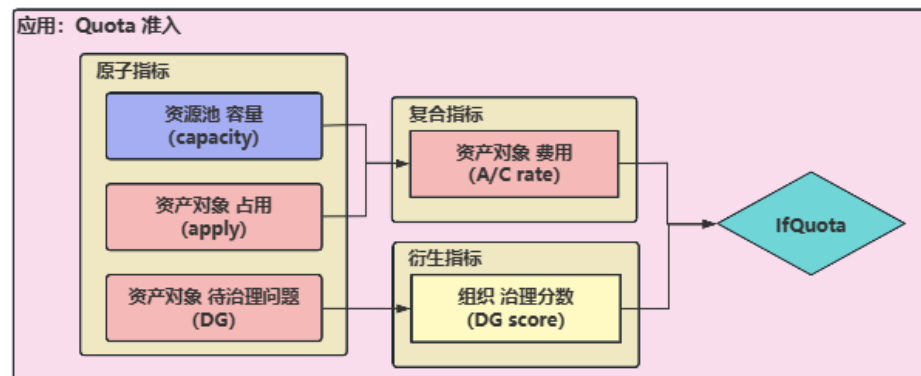
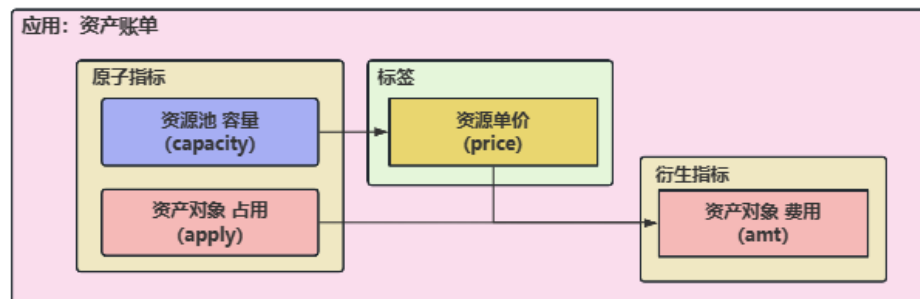
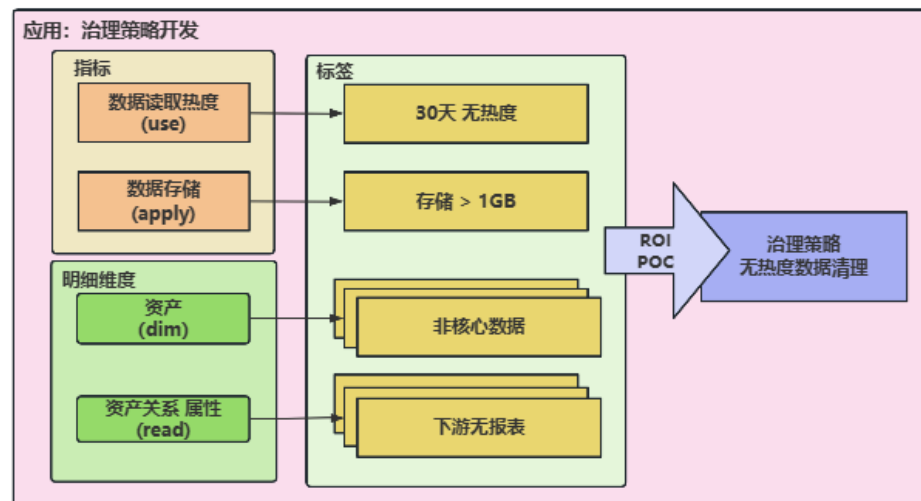
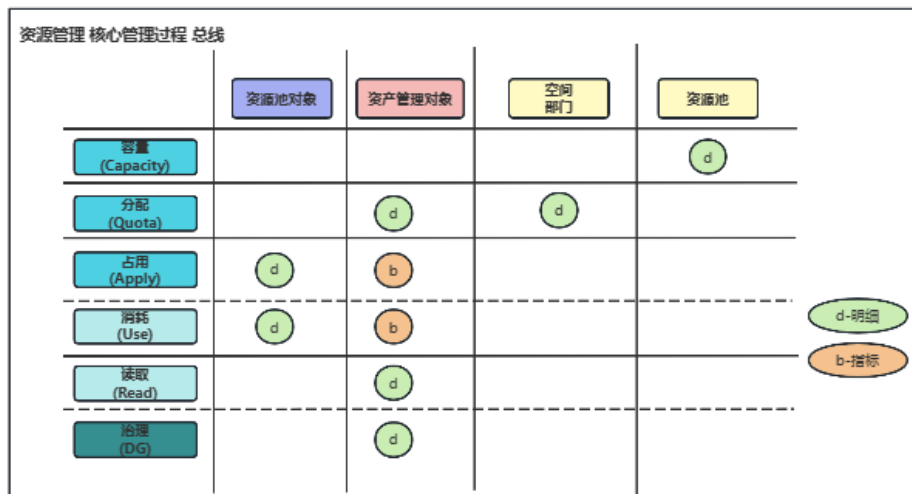
资源关系管理元模型

资源关系管理元模型  
hdfs

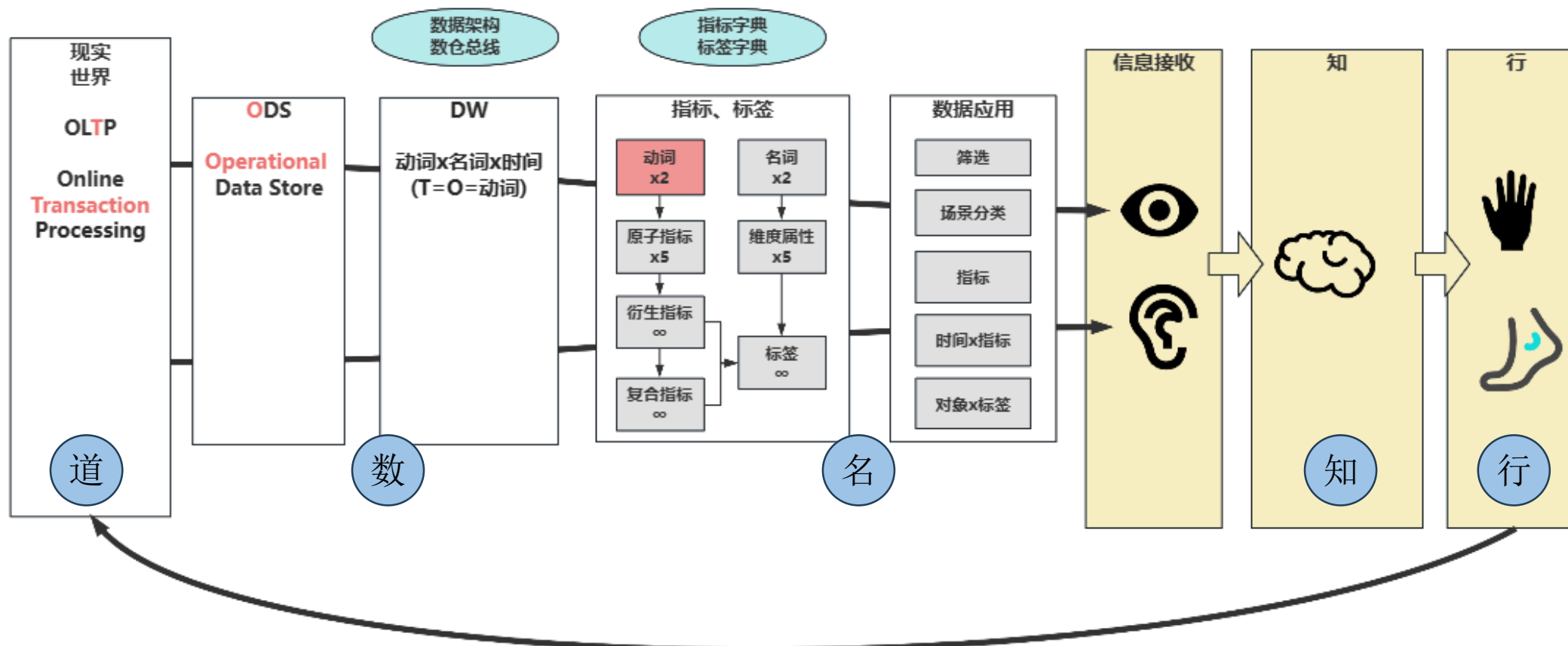
数据流 (flow) 资源关系管理元模型

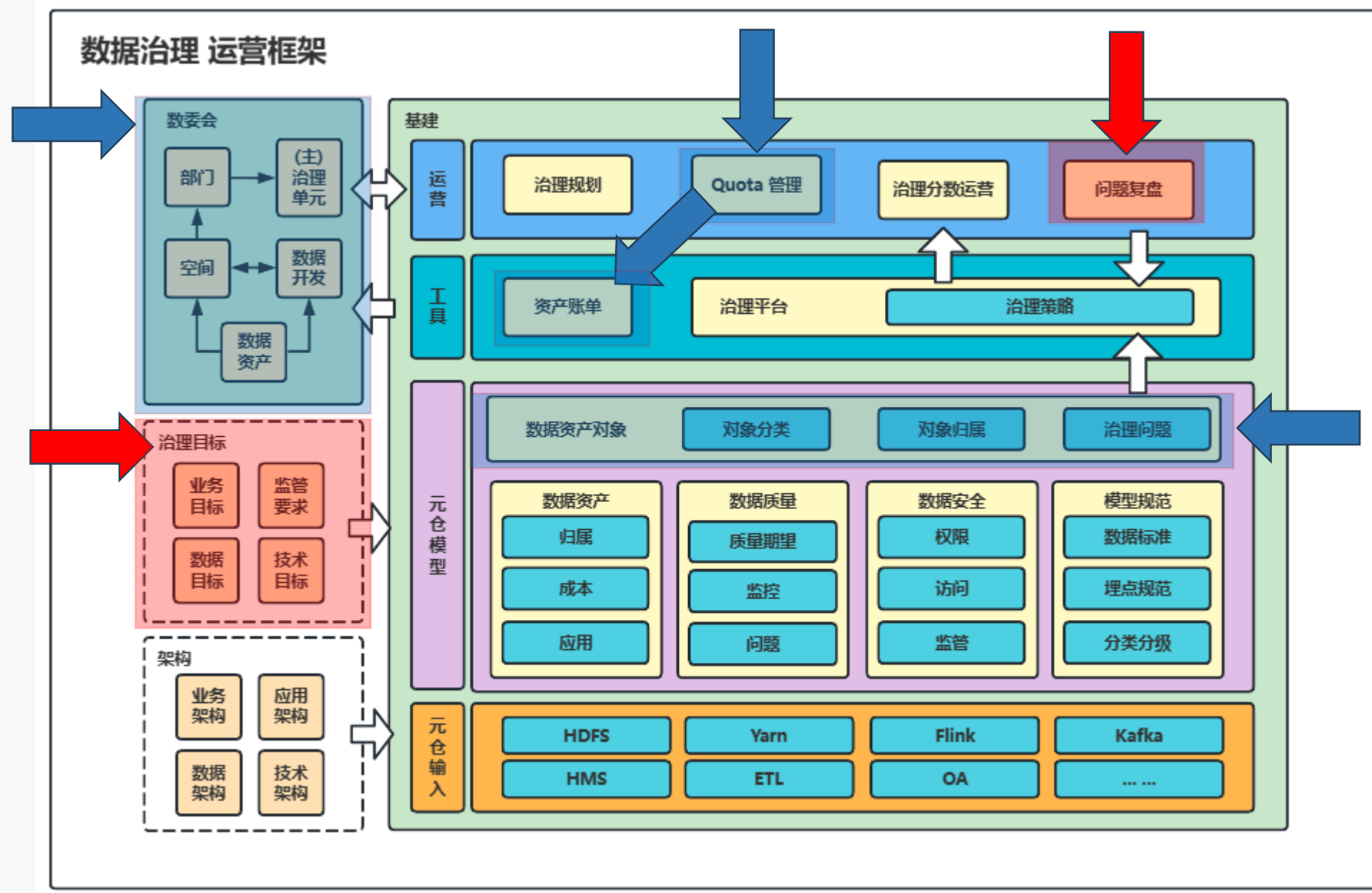


# Part3：元数据的管理 -- 元数据的指标与标签



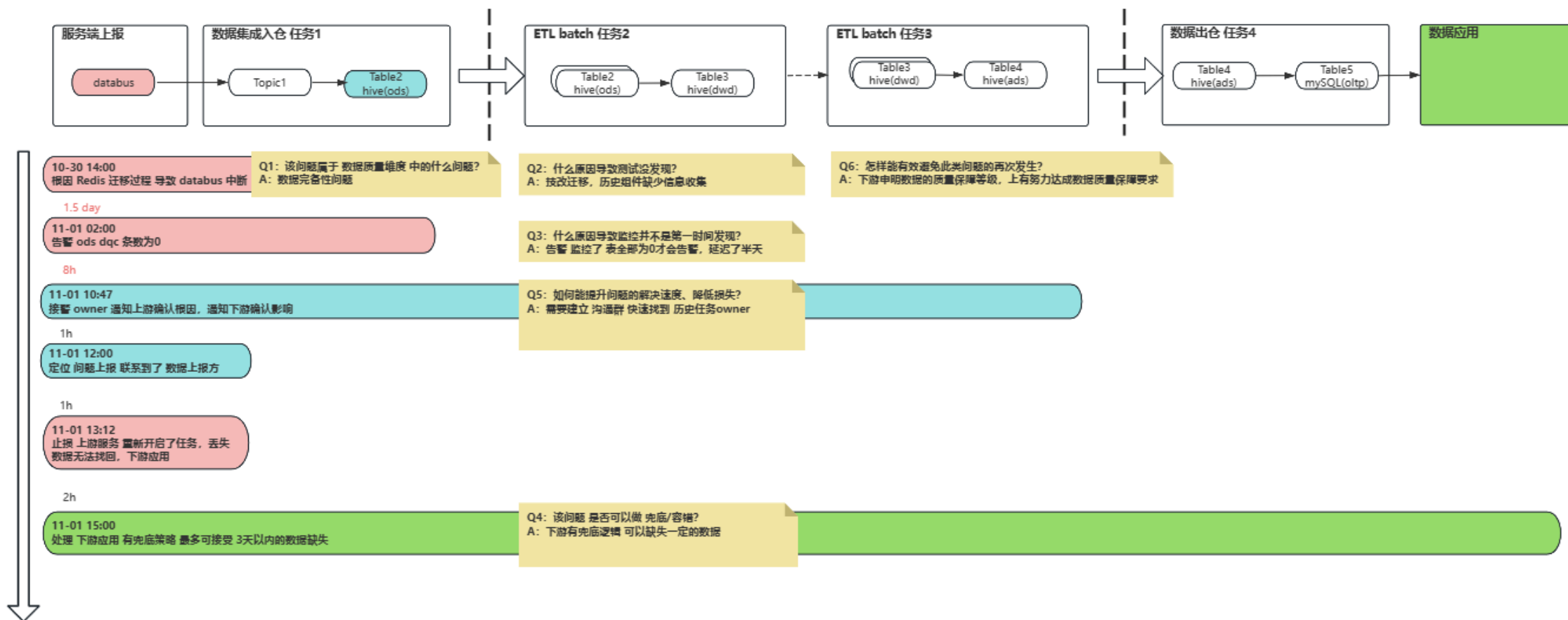
## \* 关于“指标 标签”与“道行数知名”





# 案例2 -- 2023-10-30 数据丢失复盘

# 案例背景（2线6问）





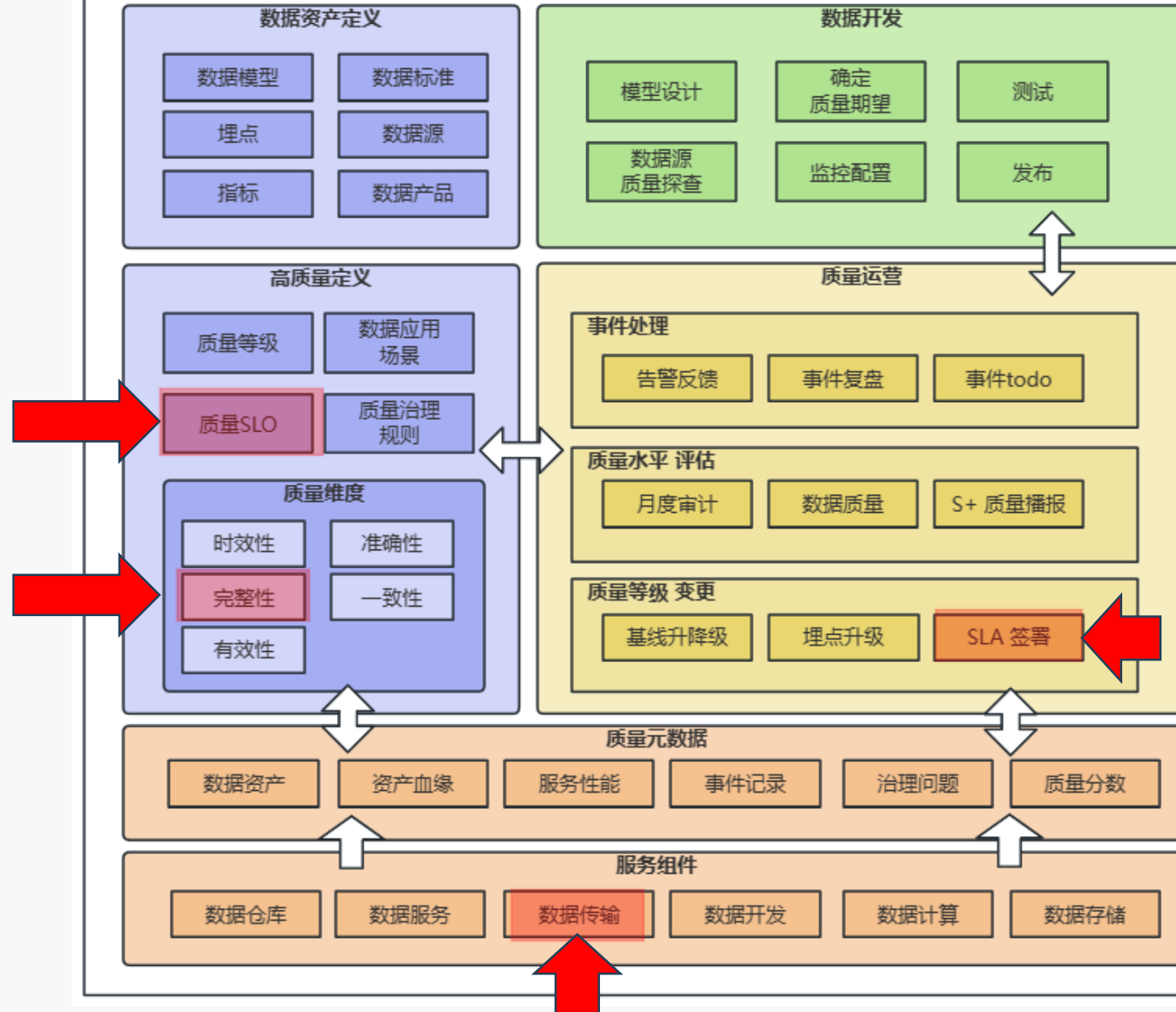
## 问题与挑战

	问	答	DAMA Bok
Q1	该问题属于 数据质量维度 中的什么问题？	数据完整性问题	<ul style="list-style-type: none"> <li>原则：数据管理即使对数据的质量进行管理</li> <li>领域：数据质量</li> </ul>
Q2	什么原因导致测试没发现？	技改迁移， <b>历史组件缺少信息收集</b>	<ul style="list-style-type: none"> <li>原则：数据管理需要元数据</li> <li>领域：数据架构</li> </ul>
Q3	什么原因导致监控并不是第一时间发现？	<b>告警 监控了，表全部为空才会告警，延迟了半天</b>	<ul style="list-style-type: none"> <li>原则：数据价值使用经济术语表达</li> <li>领域：数据质量</li> </ul>
Q4	该问题 是否可以做 兜底/容错？	<b>下游有兜底逻辑</b> 可以缺失一定的数据	<ul style="list-style-type: none"> <li>原则：数据管理需要全景视角</li> <li>领域：数据质量</li> </ul>
Q5	如何能提升问题的解决速度、降低损失？	<b>数据提供方</b> 应该感知下游数据应用的重要等级	<ul style="list-style-type: none"> <li>原则：数据管理需要全景视角</li> <li>领域：数据质量</li> </ul>
Q6	怎样才能有效避免此类问题的再次发生？	<b>数据服务方</b> 应该感知下游数据应用的重要等级	<ul style="list-style-type: none"> <li>原则：数据管理是跨功能的</li> <li>领域：数据质量</li> </ul>

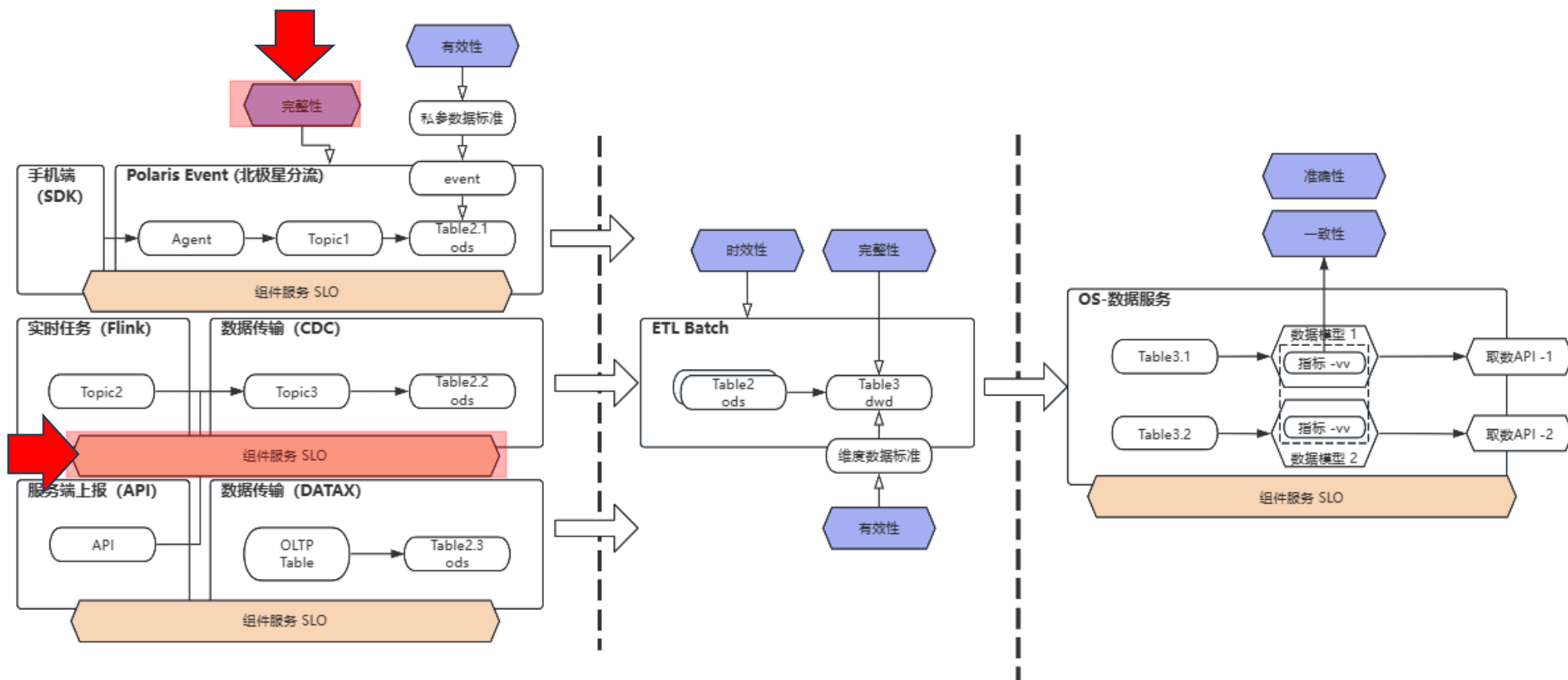
## 破题思路

问题	DAMA Bok	方案
历史组件缺少信息收集?	<ul style="list-style-type: none"><li>原则：数据管理需要元数据</li><li>领域：数据架构</li></ul>	<ul style="list-style-type: none"><li>各类组件都应具备自回收能力，基于该能力补充需要的元数据与功能更流程</li></ul>
能否不等表全部为空就可以发出告警	<ul style="list-style-type: none"><li>原则：数据价值使用经济术语表达</li><li>领域：数据质量</li></ul>	<ul style="list-style-type: none"><li>数据传输链路中需要 增加完整性监控方案</li></ul>
是否下游都有兜底逻辑?	<ul style="list-style-type: none"><li>原则：数据管理需要全景视角</li><li>领域：数据质量</li></ul>	<ul style="list-style-type: none"><li>重要数据下游使用方需要明确提供数据质量的问题识别与容忍区间</li></ul>
数据提供方如何感知下游数据的重要性?	<ul style="list-style-type: none"><li>原则：数据管理需要全景视角</li><li>领域：数据质量</li></ul>	<ul style="list-style-type: none"><li>数据下游的使用场景，质量容忍度需要传递给数据提供方</li></ul>
数据服务方是否提供了足够的质量保障?	<ul style="list-style-type: none"><li>原则：数据管理是跨功能的</li><li>领域：数据质量</li></ul>	<ul style="list-style-type: none"><li>数据服务方需要定期进行质量服务水平的审计</li></ul>

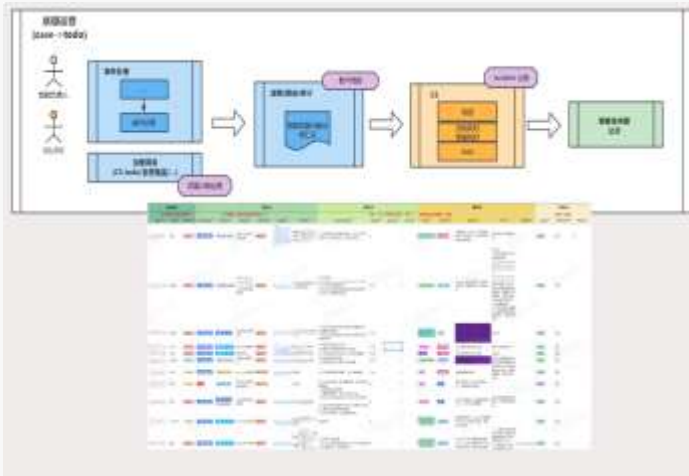
## 数据质量运营框架



# 数据质量监控



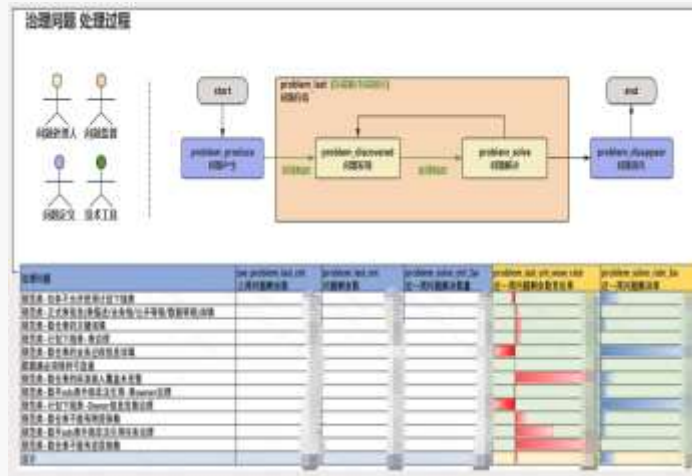
# 从 CS 到 TODO



## 事前 分析问题



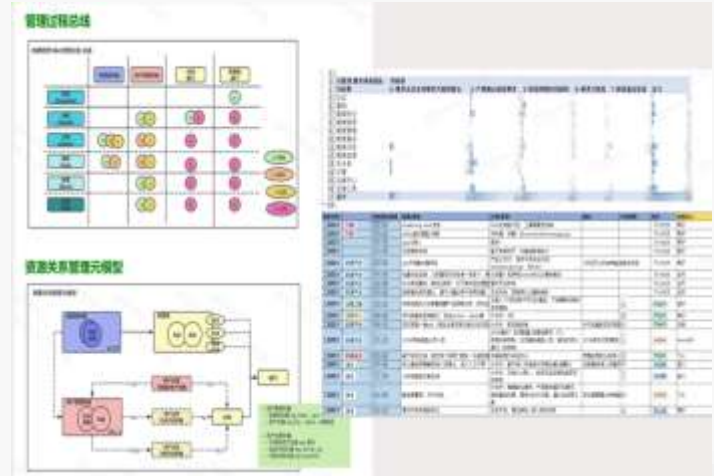
通过 各种渠道收集问题  
通过 复盘寻找问题处理方法  
通过 治理策略、开发基建 控制问题的影响



## 事中 处理问题



需要观察发布的治理项目是否有在被顺利执行  
需要观察问题新增速度是否远远大于处理速度  
需要评估使用人力参与数治理的工作是否合理



## 事后 沉淀基建



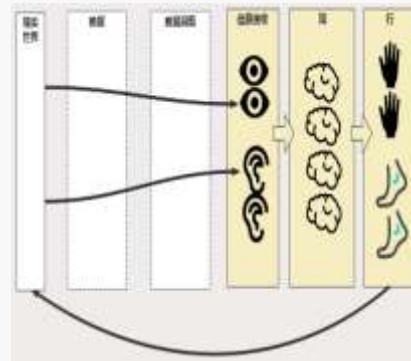
开发 元数据增强 数据管理能力  
开发 治理策略 提高问题处理 效率  
开发 基建降低问题产生的几率

# TakeAway



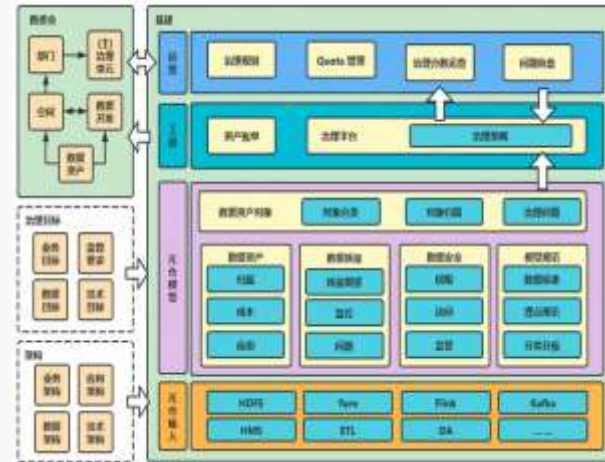
**DAMA 车轮图**

11 个数据管理知识领域



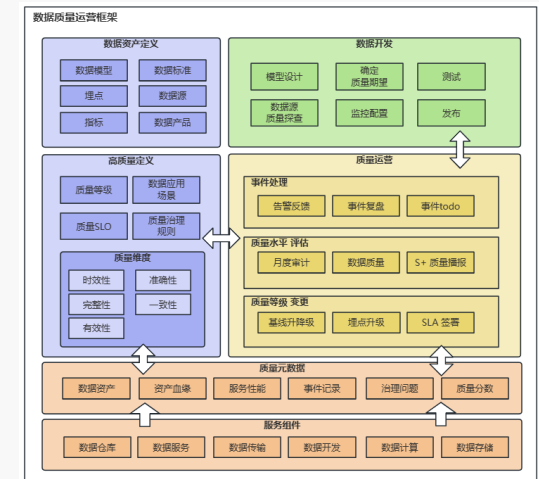
**道-行-数-知-名**

分析数据问题往往就是分析管道问题



**数据治理运营框架**

实施数据治理需要操作系统



**数据质量运营框架**

处理 复盘 执行 沉淀





**微信官方公众号：壹佰案例**  
**关注查看更多年度实践案例**



**微信官方公众号：哔哩哔哩技术**  
**关注查看更多B站技术分享**