

阿里巴巴在线技术峰会
Alibaba Online Technology Summit

企业大数据平台下数仓架构

阿里云-飞天一部
介然

 Alibaba Group
阿里巴巴集团

 阿里云
aliyun.com

总体思路

模型设计

数加架构

数据治理

新环境下的数据应用特征



关键词



大数据平台特征

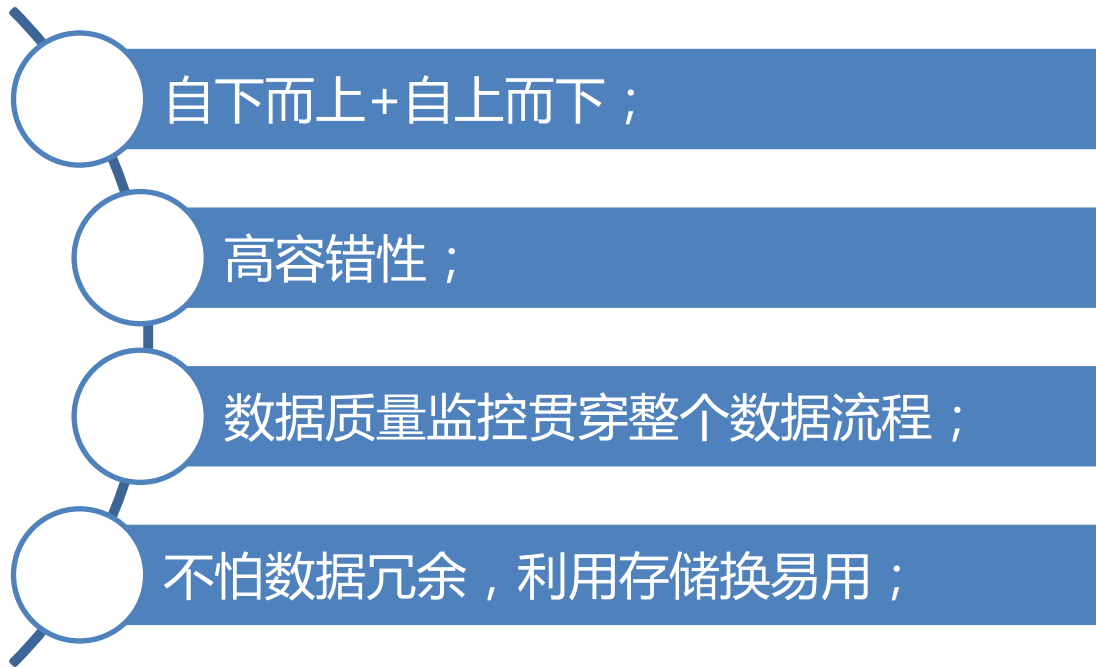
强大的计算和存储能力；

多样的编程接口和框架；

丰富的数据采集通道；

各种安全和管理措施；

仓库架构设计原则



总体思路

模型设计

数加架构

数据治理

维度建模 OR 实体关系建模

- 维度建模
 - 实施简单
 - 便于事实数据分析
 - 适合业务分析报表和BI

- 实体关系建模
 - 实施复杂
 - 便于主体数据打通
 - 适合复杂数据内容的深度挖掘

星型模型 AND 雪花模型

两种模型是并存的

星型是雪花的一种，理论上真实数据的模型都为雪花模型，实际数据仓库中两种模型会并存。

中间层将雪花转变成星型

星型模型相对结构简单，在数据中间层利用数据冗余将雪花转变成星型模型有利于数据应用和减少计算资源消耗。



数据分层

上下三层结构

减少层次结构的目的是为了压缩整体数据处理流程的长度，扁平化的数据处理流程有助于数据质量控制和数据运维

流式数据作为数据体系的一部分

当前的数据应用方向会越来越关注数据的时效性，越实时的数据价值度越高。



基础数据层

数据采集

把不同数据源的数据统一采集到一个平台

数据清洗

清洗不符合质量要求的数据，避免脏数据参与后续数据计算

数据归类

建立数据目录，在基础层一般按照来源系统和业务域进行分类

数据结构化

对于半结构化或非结构化的数据，进行结构化

数据规范化

规范维度标识、统一计量单位

数据中间层

围绕实体打通行为

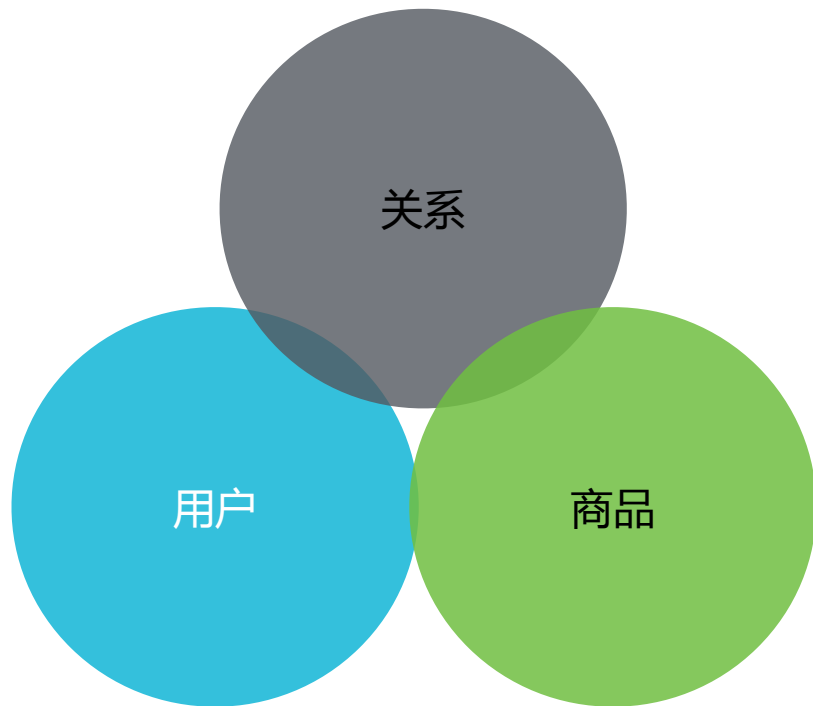
当前业务形态下，同一实体的数据可能分散在不同的系统和来源，且这些数据对同一实体的标示符可能不同。中间层最重要的目标是把同一实体不同来源数据打通起来。

从行为抽象关系

从行为中抽象出来的基础关系，会是未来上层应用一个很重要的数据依赖。如兴趣、偏好、习惯等关系数据是推荐、个性化的基础生产资料。

冗余是个好手段

在中间层，为了保证主题的完整性或提高数据的易用性，经常会进行适当的数据冗余。比如某一事实数据和两个主题相关但自身又没有成为独立主题，会放在两个主题库中。为了提高单数据表的复用性和减少计算关联，通常会在事实表中冗余部分维度信息。

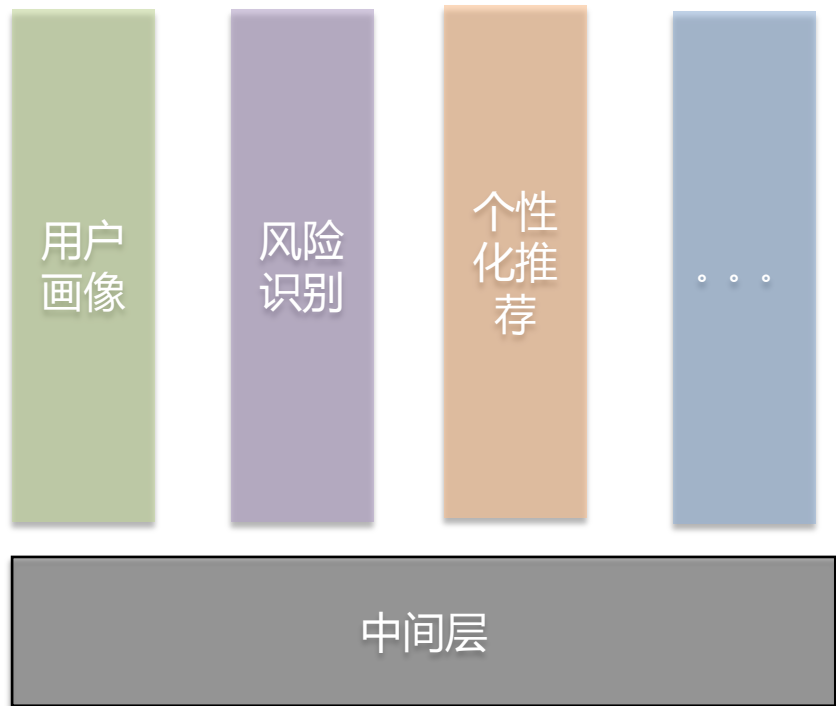


数据集市层

需求场景驱动的集市层建设，各集市之间垂直构建

集市层深度挖掘数据价值

集市层需要能够快速试错



流式数据集

需求驱动

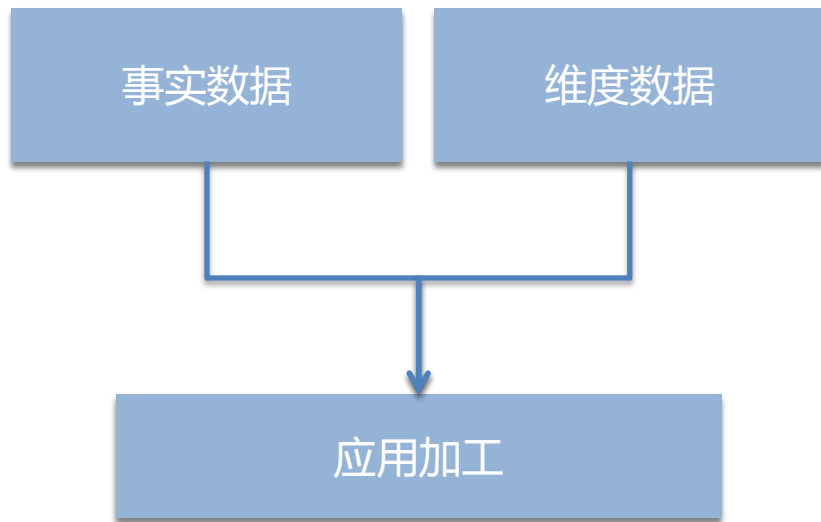
流式数据集的采集、加工和管理的成本较高，一般都会按照需求驱动的方式建设。

包含事实和维度

未来保障数据统计的准确度，流式数据集同样包含事实和维度。

结构更扁平

介于成本较高，流式数据体系的结构更扁平，通常不会设计中间层。

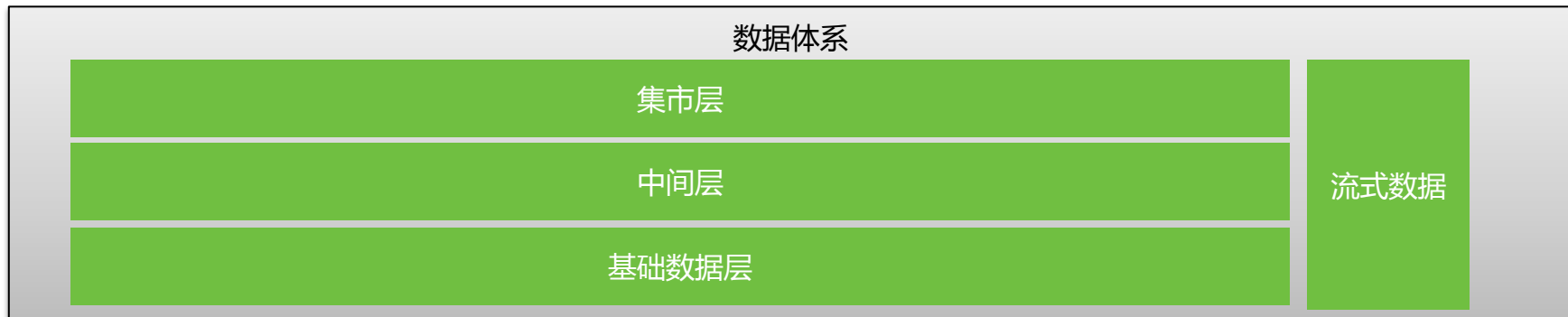


总体思路

模型设计

数加架构

数据治理



结构化数据采集

全量采集

- 每日采集数据库表的快照。
- 适合数据量较小的数据集。
- 前端库压力较小，不会影响前端应用，不会占用较大带宽，同步时间不会较长。
- 采集方式最简单，对库表没有特殊要求，后续使用较简单。

增量采集

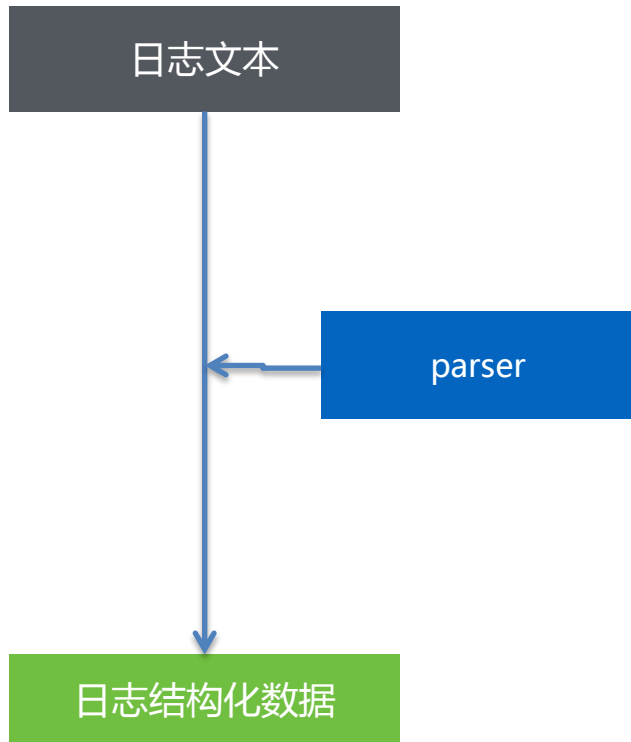
- 采集数据集每日变化的数据。
- 适合数据已经有一定规模，全量采集时间较长且对前端应用产生压力的状况。
- 前端业务表的设计需包含时间戳字段，且任何对数据的操作都会更新时间戳。
- 增量采集后一般会和昨天的全量数据merge一个今天的全量数据。

实时采集

- 实时采集数据集的变化数据。
- 比较适合数据量巨大，增量数据同步资源也消耗严重的情况。
- 或者后续的数据应用需要用到准实时数据。
- 实时采集对采集端系统有一定的要求。
- 采集质量最难控制。

日志结构化

- 日志采集到平台之前不做结构化
- 通过换行符分割每条日志，整条日志存储在一个数据表字段
- 通过UDF或MR计算框架实现日志结构化
- 日志原始结构越规范，解析的成本越低
- 并不一定需要完全平铺数据内容，结构化出重要常用字段，为了保障扩展性，利用数据冗余保存原始符合字段，如useragent字段



非结构化数据特征提取

语音转文本

图片识别

自然语言处理

图片打标

视频识别

...

数据服务化

统计服务

偏传统的报表服务，利用大数据平台将数据加工后的结果放入关系型数据库中，由前端的报表系统或业务系统查询。

分析服务

提供明细的事实数据，利用大数据平台的实时计算能力，允许操作人员自主灵活的进行各种维度的交叉组合查询。能力类似于传统cube提供的内容，但是在大数据平台下不需要预先建好cube，更灵活，更节省成本。

标签服务

大数据的应用场景下，经常会对主体进行特征刻画，比如客户的消费能力、兴趣习惯、物理特征等等，这些数据会转换成KV的数据服务，提供前端应用查询。



架构设计中一些实用的点

● 巧用虚拟节点

- 多系统数据源同步
- 跨系统间数据传输
- 多应用间数据交互

● 强制分区

- 所有数据表都应该加上时间分区
- 保障每个任务都能够独立重跑
不产生数据质量问题
- 所有数据处理过程都需要增加
分区裁剪

架构设计中一些实用的点

- 计算框架应用

- 日志结构化
- 同类数据计算过程
- 减少数据扫描次数

- 优化关键路径

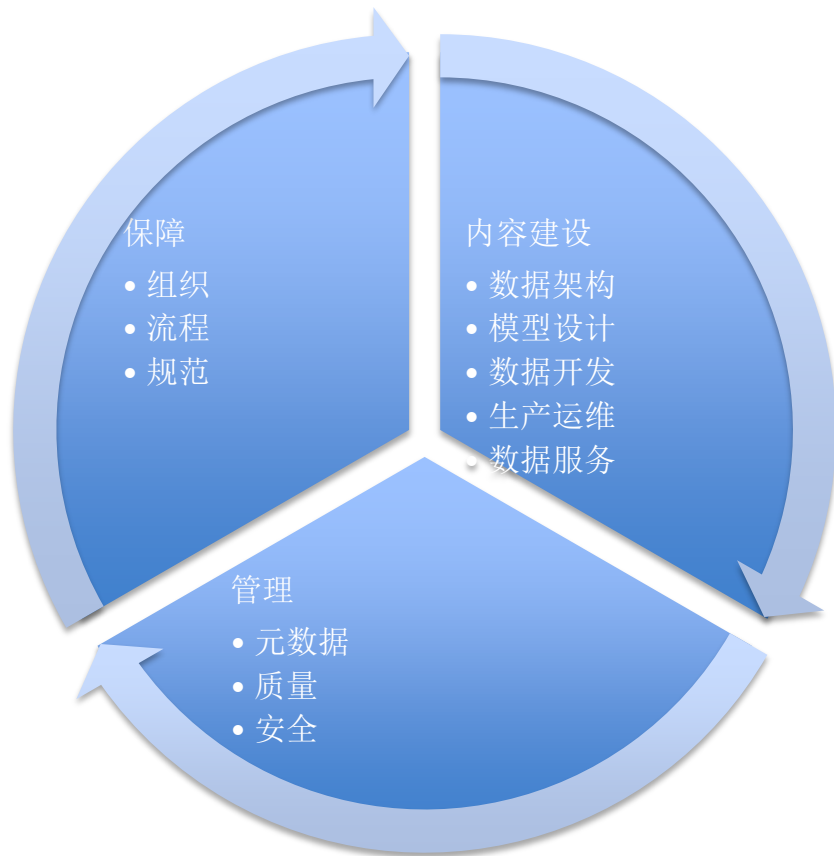
- 每份数据的产出都有一个关键数据加工路径
- 优化关键路径中耗时最长的任务是最有效的保障数据产出时间的手段
- 对重要数据产出增加基线监控

总体思路

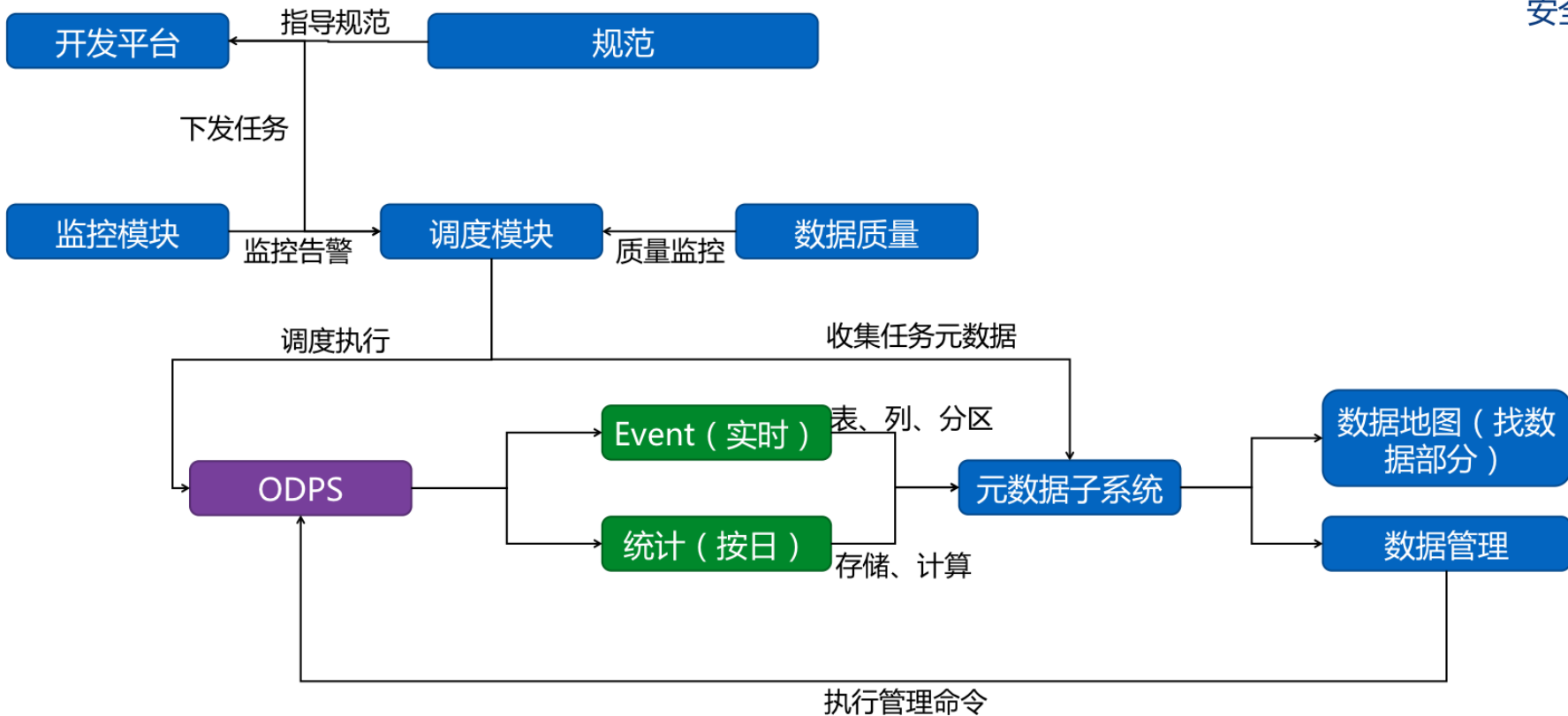
模型设计

数加架构

数据治理

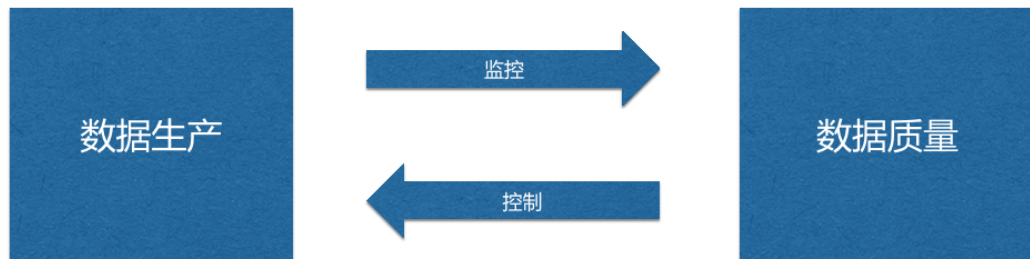


安全管控



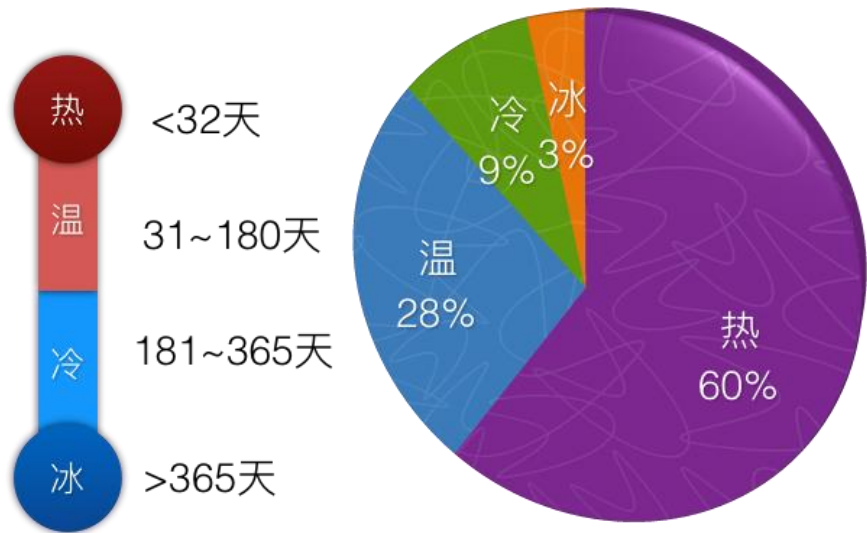
• 数据质量

- 事前：制定每份数据的数据质量监控规则
- 事中：监控和影响数据生产过程，不符合质量要求的数据不算产出数据；
- 事后：数据质量情况分析和打分，推动数据质量提升；



数据生命周期管理

- 合理的数据生命周期管理要保证温热数据占整个数据体系大部分
- 为了保障数据资产的完整性，对于重要的基础数据会长久保留
- 对于数据中间计算过程数据，在保障满足绝大部分应用访问历史数据需要的前提下，缩短数据保留周期，有助于降低存储成本
- 冷备已经成为历史，在大数据平台下不需要单独的冷备设备



数加平台



新手引导，帮你快速入门

- <https://data.aliyun.com>
- 一站式大数据开发、分析及应用平台

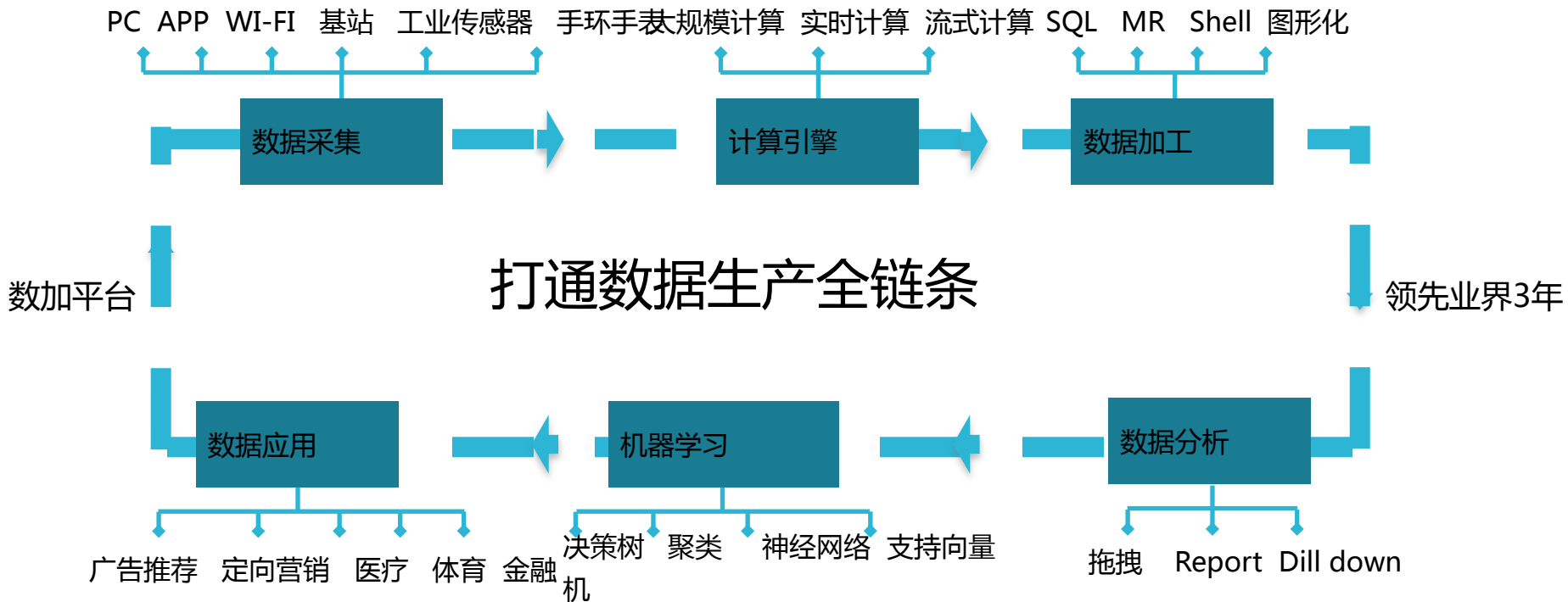
大数据体验馆



免费体验 + 教程 = 大数据零距离！



<https://data.aliyun.com/experience?spm=a2c0j.7906235.started.1.mIBzT8>



The background of the slide is a deep space scene. In the center, there is a large, circular nebula with a complex, filamentary structure. The colors of the nebula range from a deep, dark blue to a bright, fiery orange-red. The surrounding space is dark, filled with numerous small, distant stars of various colors, including blue, white, and yellow. Some stars appear as bright, multi-pointed patterns, possibly due to diffraction or the presence of interstellar dust. The overall atmosphere is one of vastness and cosmic beauty.

Thank you !