



DTCC

2016中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2016

数据定义未来

SequeMedia
盛拓传媒

IT168.com

ChinaUnix

ITPUB

简介：张粤磊 (Jackson)

邮箱：vzyuelei@126.com

@me:

微信：
vzyuelei

QQ：
416988515



- 飞谷云 (www.feiguyun.com) 创始人
- (2014-2016.3) 平安付大数据平台架构师
- (2012-2014) 外汇交易中心ETL项目开发经理
- (2010-2012) HP TRAM项目 ETL开发组长
- (2005-2010) DBA
- 10余年一线数据业务（制造，咨询服务，互联网金融）及数据处理技术实践经验

关注的技术产品及工具：



DTCC

2016年中国数据库技术大会
DATABASE TECHNOLOGY CONFERENCE CHINA 2015

SequeMedia

IT168

ChinaUnix

ITPUB

数据处理的哪些事

1. 传统数据仓库的数据处理技术及思考
2. 大数据环境下对于公共数据及行为数据的数据处理技术
3. 由传统数据仓库到大数据数据仓库的数据处理实践思考及建议



传统数据仓库的数据处理技术

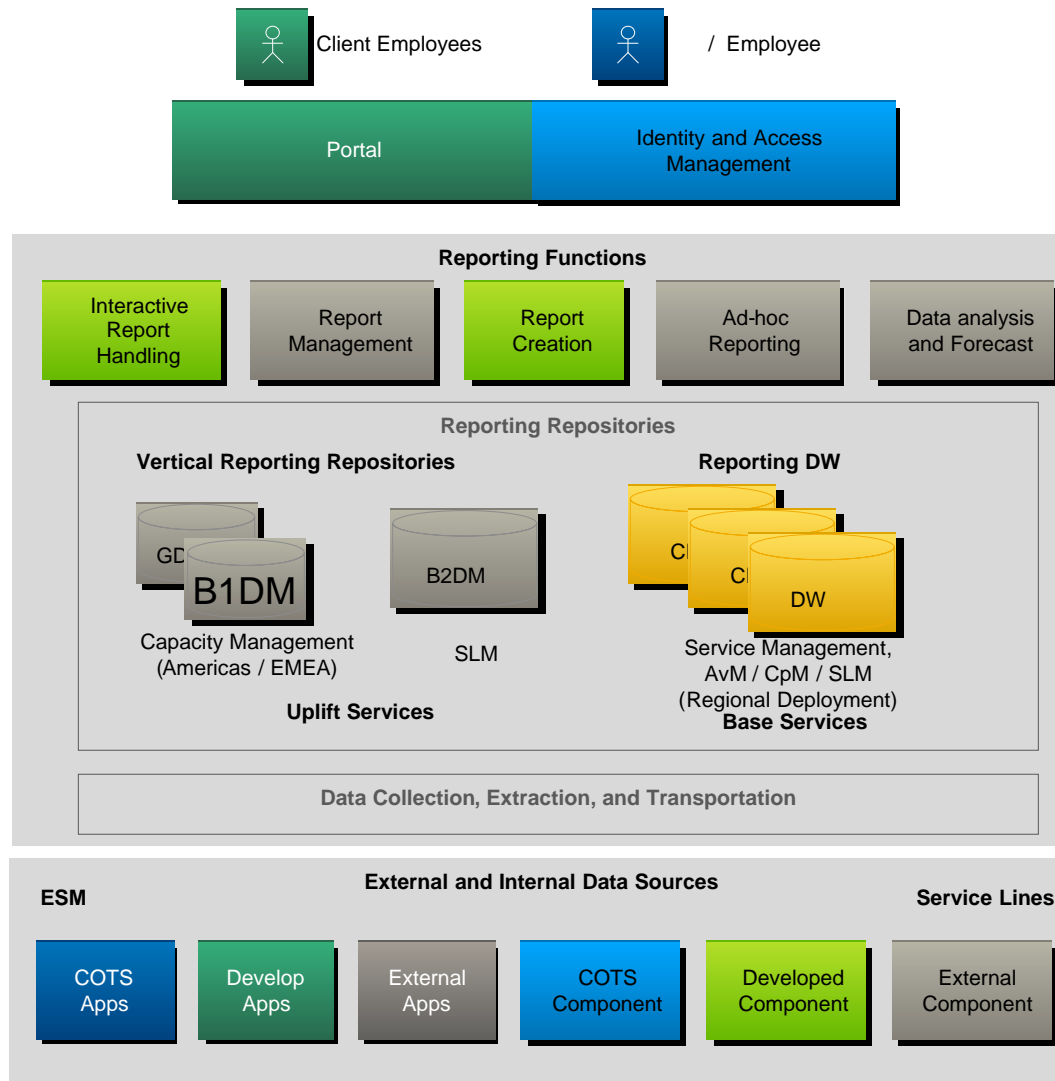
传统数据仓库的数据处理技术是什么？

从我参与某大型数据仓库项目经历为例来分享：

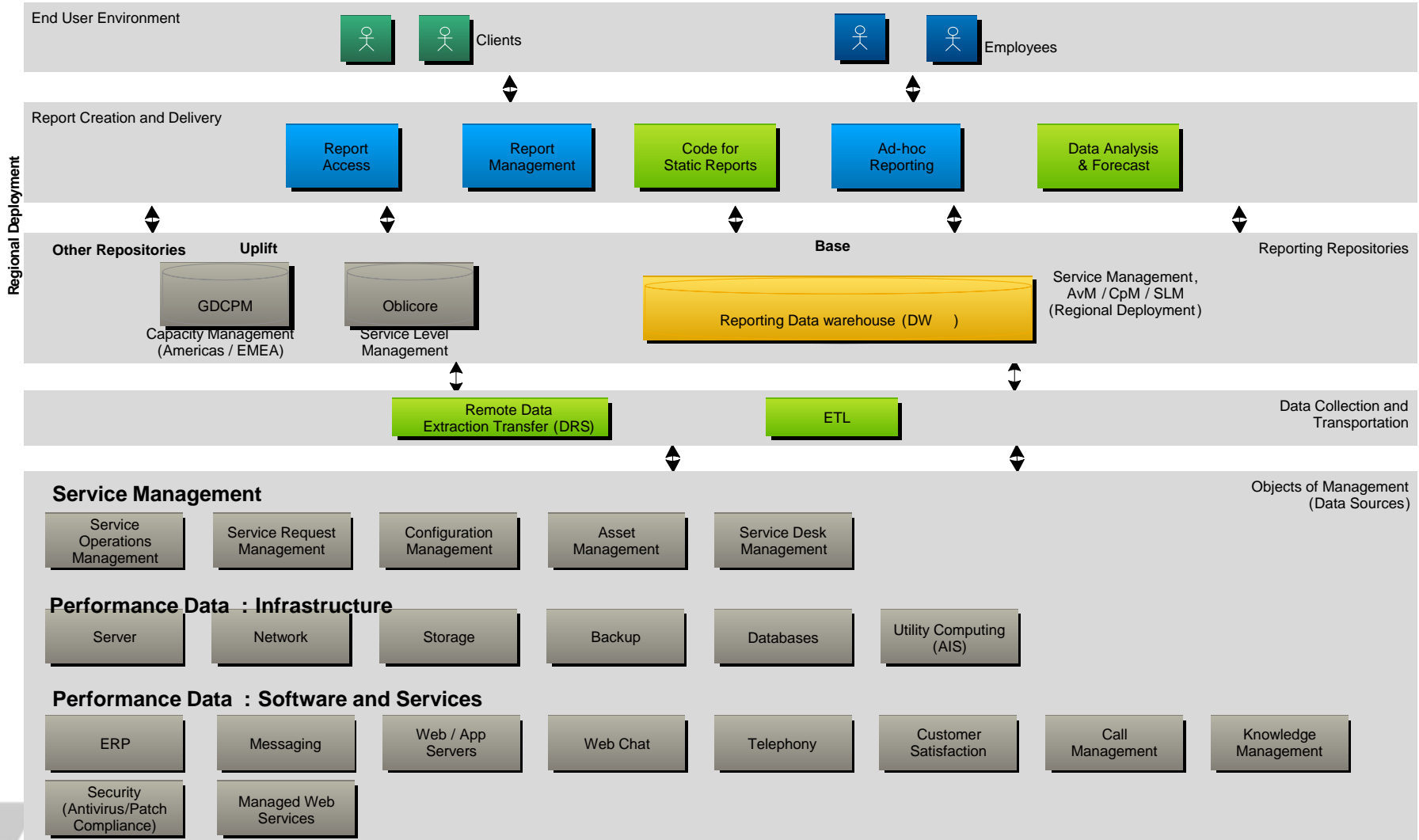
- 该过程涵盖传统数据仓库的标准流程和数据处理规范
- 数据处理方法和实践同样适用于数据平台数据处理



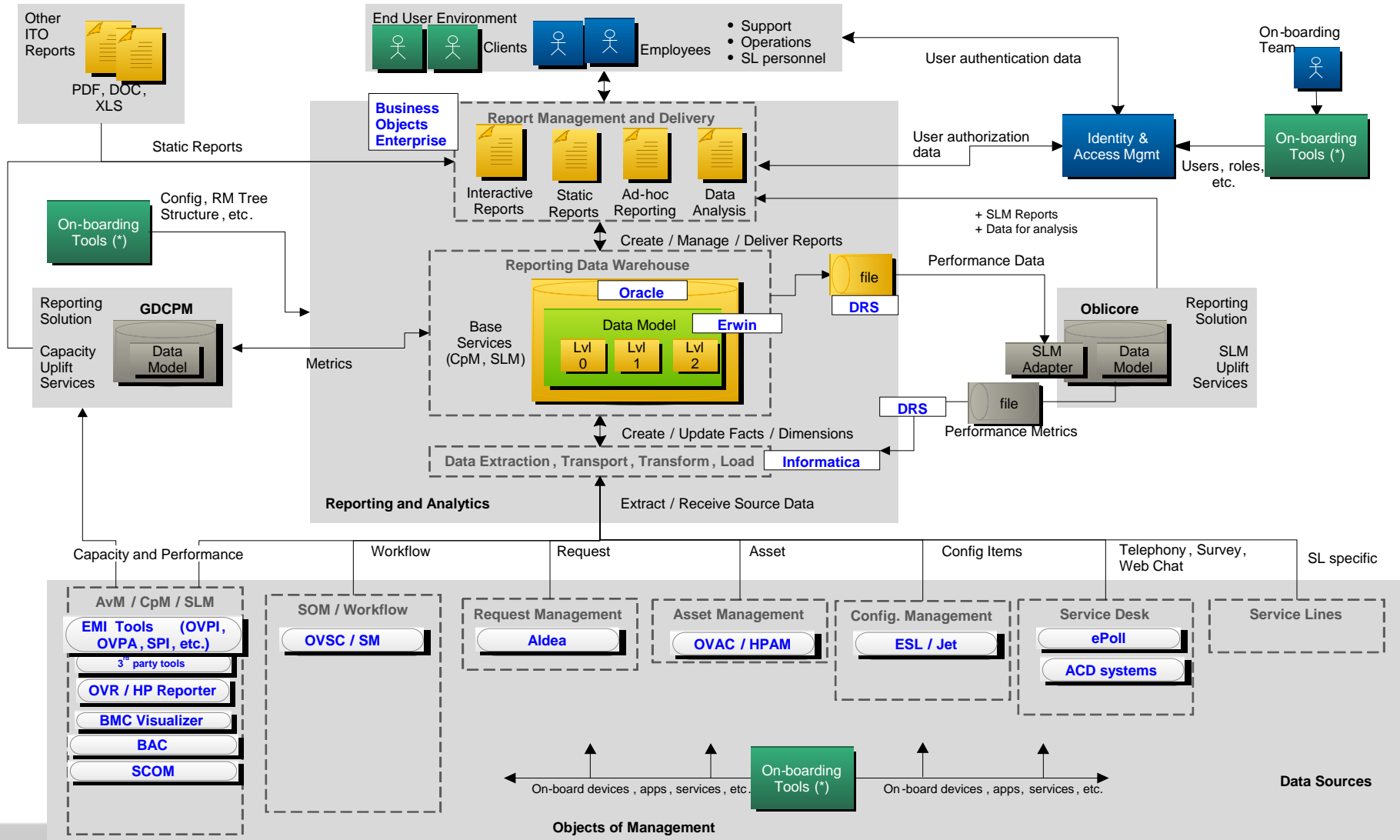
传统数据仓库的数据处理技术-概念定义



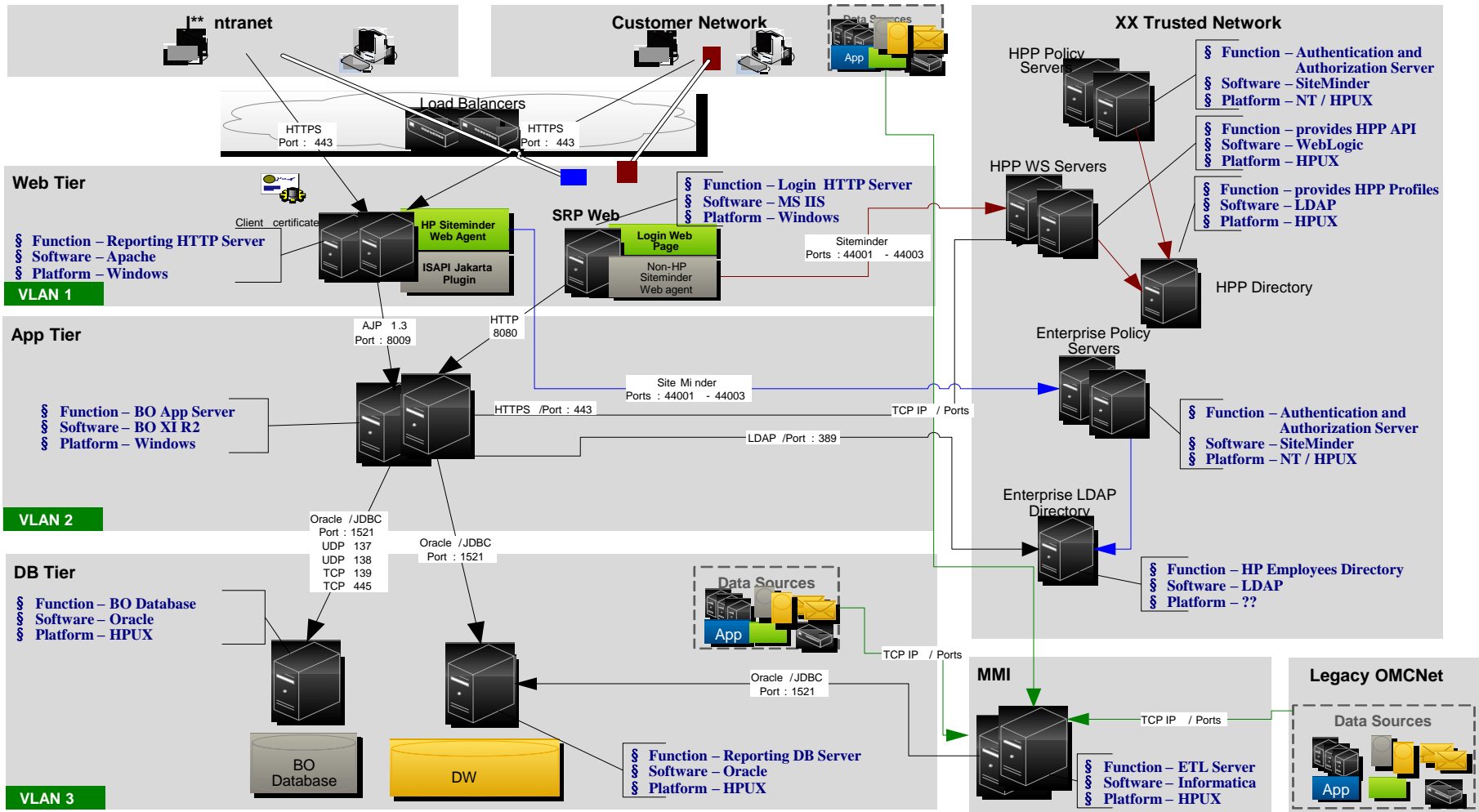
传统数据仓库的数据处理技术-业务定义



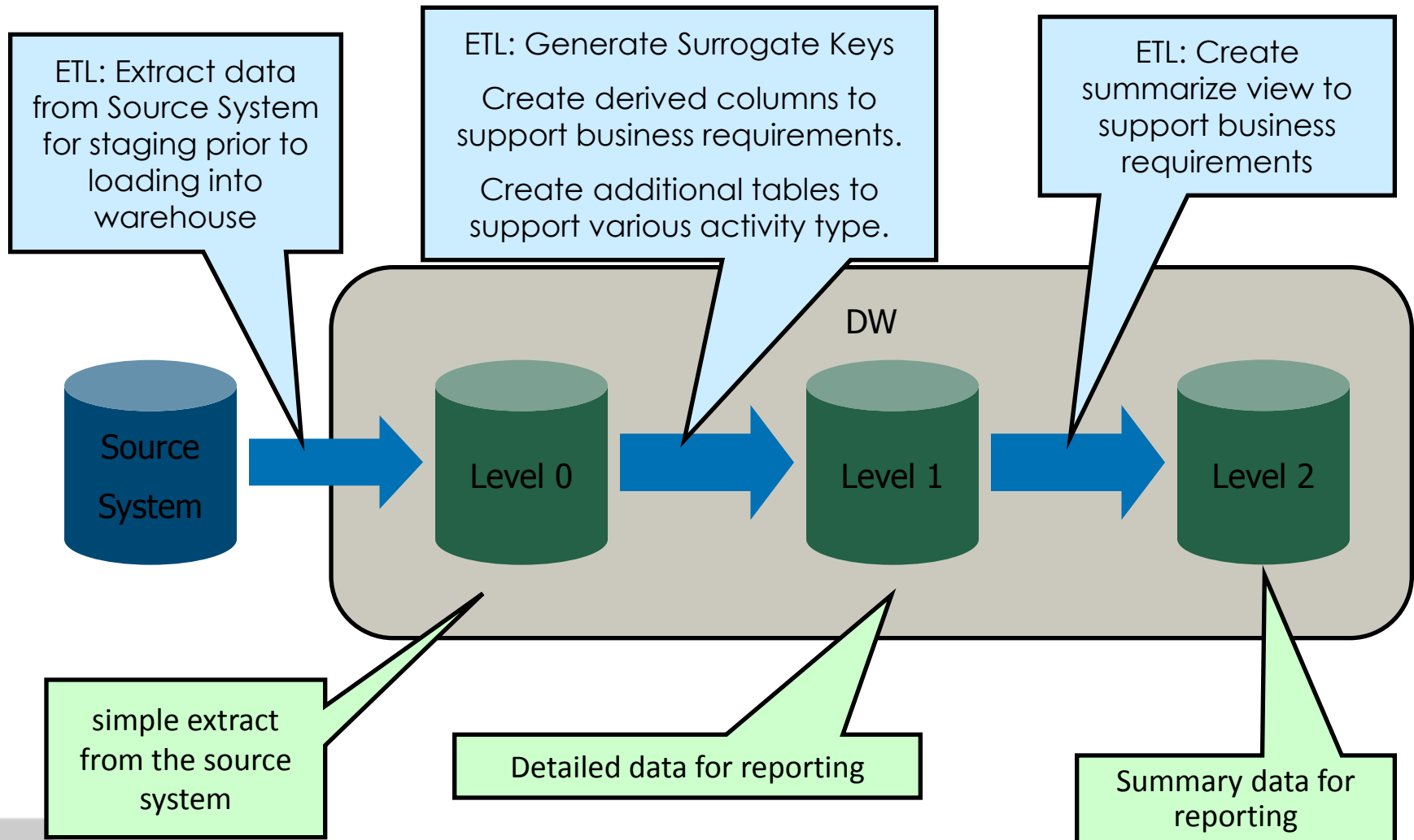
传统数据仓库的数据处理技术-逻辑定义



传统数据仓库的数据处理技术-物理定义



传统数据仓库的数据处理技术-规则定义

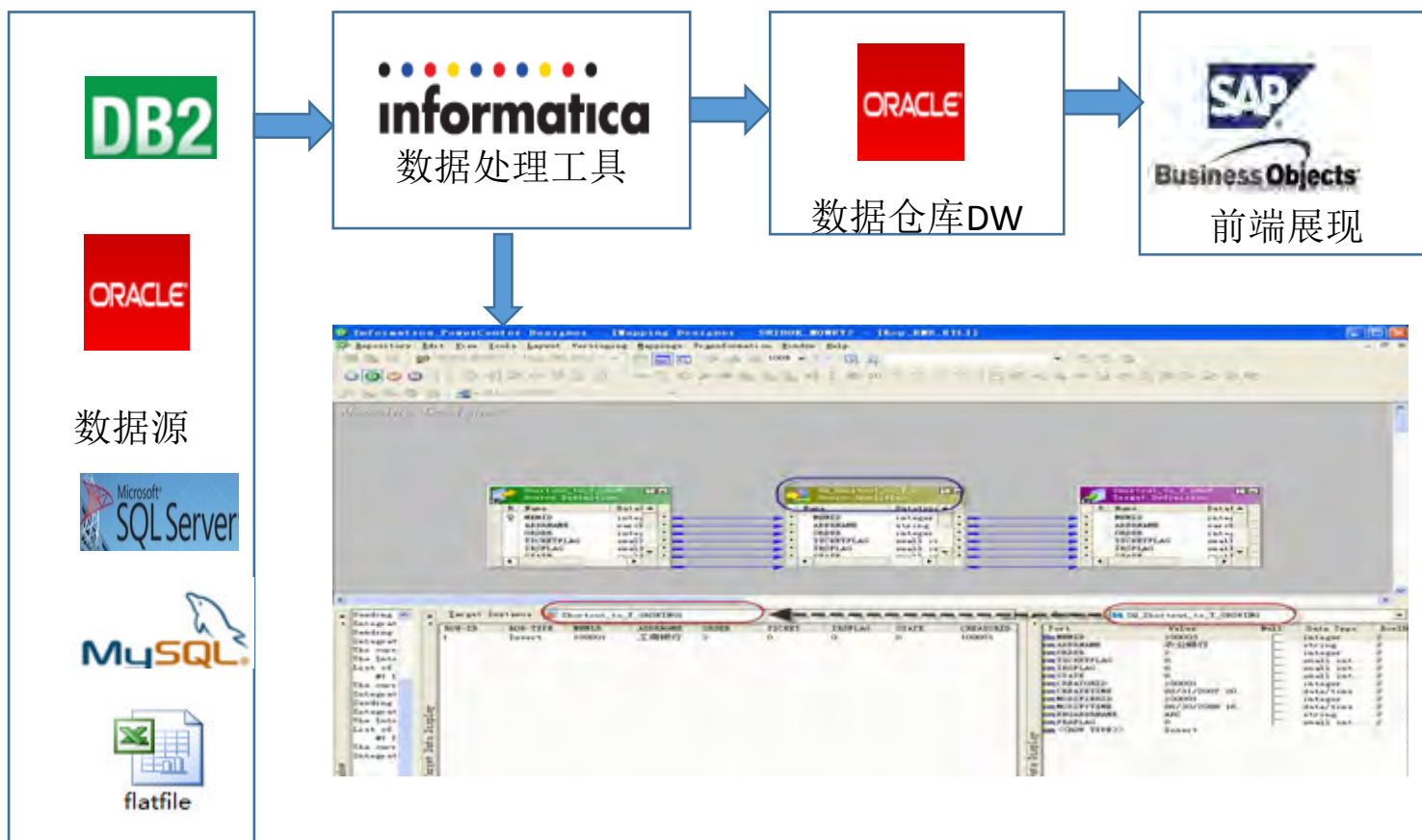


传统数据仓库的数据处理技术-设计定义

序号 (ID)	源系统信息 (Source System Information)						目标系统信息 (Target System Information)				ETL抽取 (PowerC)					
	源Schema (Source Schema)	表名 (Table Name)	标志位读取方式 (标志位数据生成 伴有标志位时, 需要填写: 标志 为表, 标志位名 称, 标志位表名)	表当前 的数据 量 (Row Count)	表每 月增 量 (Month Inc)	标识增量的 时间戳字段 (Record Increment Field)	该表是否 有物理删 除操作 (Is delete Action)	数据生 成时间 窗 (Data Creation Time)	目标 Schema (Target Schema)	目标表名 (Target Table Name)	加载 类型 (Load Type)	数据到达时 间需求 (Required Data Arrival Time)	ETL抽 取时间	ETL抽 取频率	标志 位等 待时 间窗 口 (天)	标志 位等 待循 环时 间隔
1	CRM	V_CNY_ENTY_MASTER_LOG						etluser	TRDX_ENTY_MASTER_LOG	UPD	每5分钟	每5分钟	每5分钟			任何时候不删
2	CRM	V_CNY_ENTY_TYPE_MASTER_LOG						etluser	TRDX_ENTY_TYPE_MASTER_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
3	CRM	V_CNY_ENTY_TYPE_DTLS_LOG						etluser	TRDX_ENTY_TYPE_DTLS_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
4	CRM	V_CNY_ENTY_ALT_CODE_MASTE_LOG						etluser	TRDX_ENTY_ALT_CODE_MASTE_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
5	CRM	V_CNY_RGN_CNFG_MASTER_LOG						etluser	TRDX_RGN_CNFG_MASTER_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
6	CRM	V_CNY_MKT_PRMTTR_DTLS_LOG						etluser	TRDX_MKT_PRMTTR_DTLS_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
7	CRM	V_CNY_ENTY_TRMNL_CNT_HSTR_LOG						etluser	TRDX_ENTY_TRMNL_CNT_HSTR_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
8	CRM	V_CNY_ENTY_ELGBLT_DTLS_LOG						etluser	TRDX_ENTY_ELGBLT_DTLS_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
9	CRM	V_CNY_ENTY_MKT_MAKING_DTLS_LOG						etluser	TRDX_ENTY_MKT_MAKING_DTLS_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
10	CRM	V_CNY_ENTY_MKT_STATUS_DTLS_LOG						etluser	TRDX_ENTY_MKT_STATUS_DTLS_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
11	CRM	V_CNY_ENTY_TYPE_RLTN_MAST_LOG						etluser	TRDX_ENTY_TYPE_RLTN_MAST_LOG		每5分钟	每5分钟	每5分钟			任何时候不删
12	CRM	V_CNY_ENTY_RLTN_DTLS_LOG						etluser	TRDX_ENTY_RLTN_DTLS_LOG		每5分钟	每5分钟	每5分钟			任何时候不删



传统数据仓库的数据处理技术-开发实现



传统数据仓库的数据处理技术-开发实现

粤磊 informatica

为您找到相关结果14个

[粤磊informatica powercenter学习笔记\(九\) - Oracle数据库管理...](#)
INFORMATICA 的部署实施之一 INFORMATICA 的UNIX安装实施 INFORMATICA 一般为了保证其高可用性大多在UNIX环境安装实施,以下是我INFORMATICA在HP UNIX环境下的安装实施经历...
www.itpub.net/thread-1392140-1...html 2015-12-18

[粤磊informatica powercenter学习笔记\(九\) - Oracle数据库管理...](#)
INFORMATICA 的部署实施之一 INFORMATICA 的UNIX安装实施 INFORMATICA 一般为了保证其高可用性大多在UNIX环境安装实施,以下是我INFORMATICA在HP UNIX环境下的安装实施经历...
www.itpub.net/forum.php?mod=viewthread 2016-3-7

[粤磊informatica powercenter学习笔记\(九\) - Oracle数据库管理...](#)
INFORMATICA 的部署实施之一 INFORMATICA 的UNIX安装实施 INFORMATICA 一般为了保证其高可用性大多在UNIX环境安装实施,以下是我INFORMATICA在HP UNIX环境下的安装实施经历...
www.itpub.net/thread-1392140-1... 2016-1-15

[粤磊informatica powercenter学习笔记\(九\) - Oracle数据库管理...](#)
INFORMATICA 的部署实施之一 INFORMATICA 的UNIX安装实施 INFORMATICA 一般为了保证其高可用性大多在UNIX环境安装实施,以下是我INFORMATICA在HP UNIX环境下的安装实施经历...
www.itpub.net/forum.php?mod=viewthread 2015-7-18

[粤磊informatica powercenter学习笔记\(五\) - 数据仓库与数据挖掘...](#)
 粤磊informatica powercenter学习笔记(五) [复制链接] vzyuele9 注册会员 ...这两天做了一下测试用INFORMATICA来实现行列互换的功能。 列转行的SQL 实现 ENV...
www.itpub.net/forum.php?mod=viewthread 2016-4-16



传统数据仓库的数据处理技术-数据治理思考

完整性

哪些数据丢失了或者哪些数据不可用?

无论选择任何一种RDBMS，都无法涵盖大量的非结构化业务数据

准确性

哪些数据和信息是不正确的，或者数据是超期的?

不同RDBMS对数据类型的定义精度各有区别

规范性

哪些数据未按统一格式存储?

基于RDBMS的数据存储并不能真实反映业务数据本源格式，文本视频，邮件在DB中的存储

唯一性

哪些数据是重复数据或者数据的哪些属性是重复的

一致性

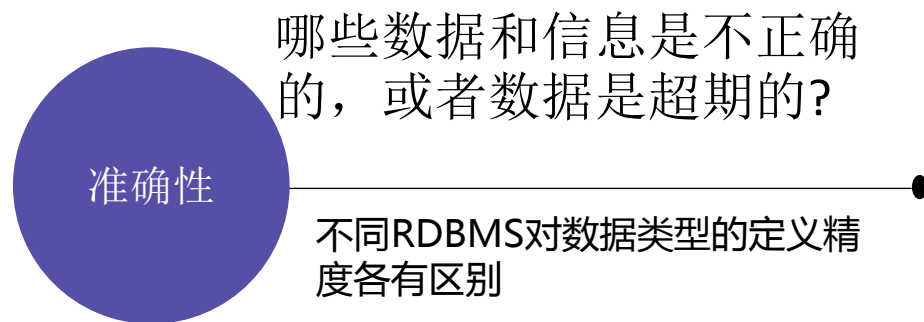
哪些数据的值在信息含义上是冲突的?

关联性

哪些关联的数据缺失或者未建立索引



传统数据仓库的数据处理技术-数据治理思考



当源系统与目标系统属于不同RDBMS或字符集等情况，可能存在字符类型不兼容问题，如：Oracle 的 date 数据类型有时分秒而db2 的date数据类型不含时分秒；oracle的Integer数据类型是8字节38位精度，db2 的Integer数据类型是4字节10位精度等等。

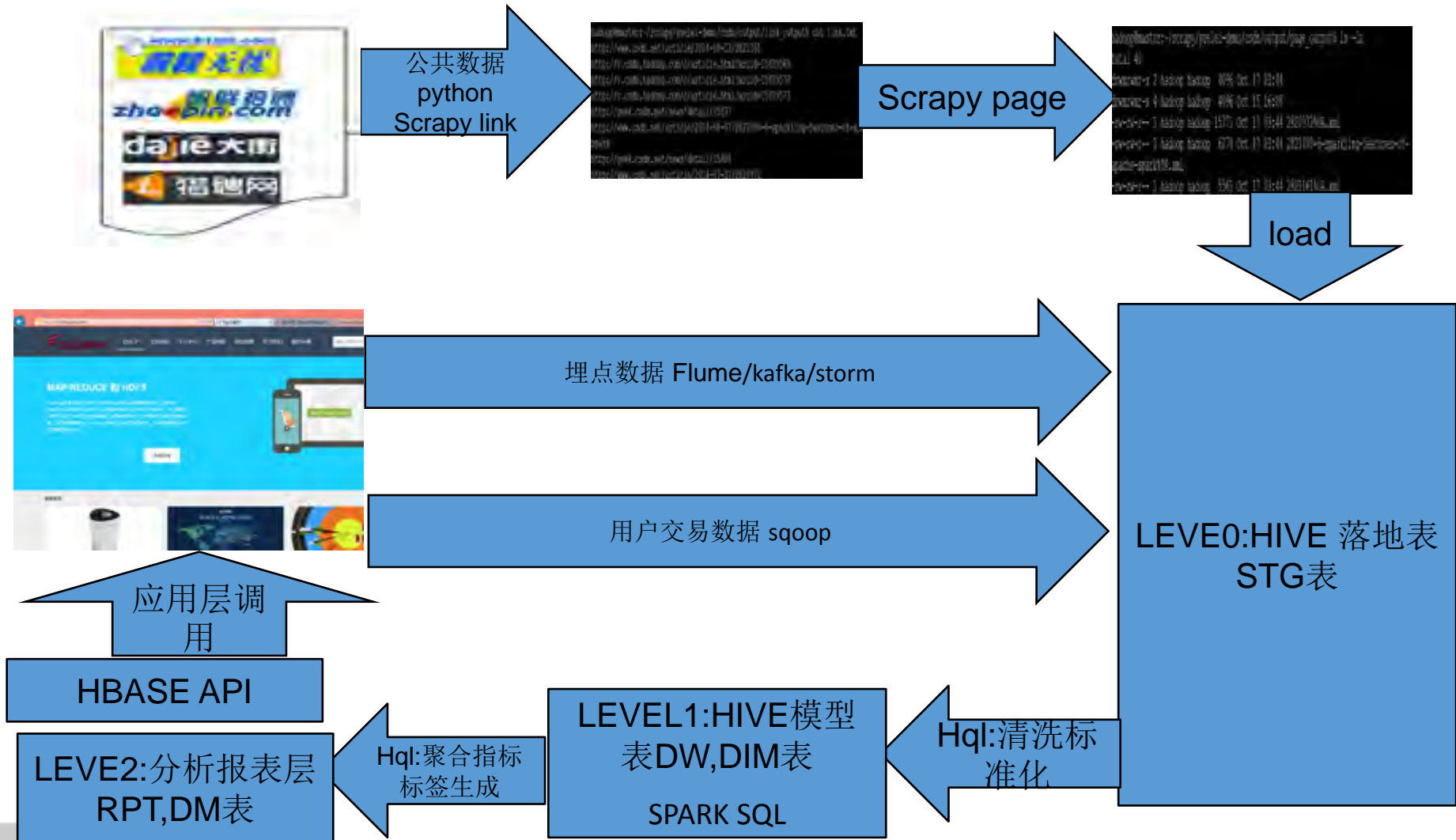


数据处理的哪些事

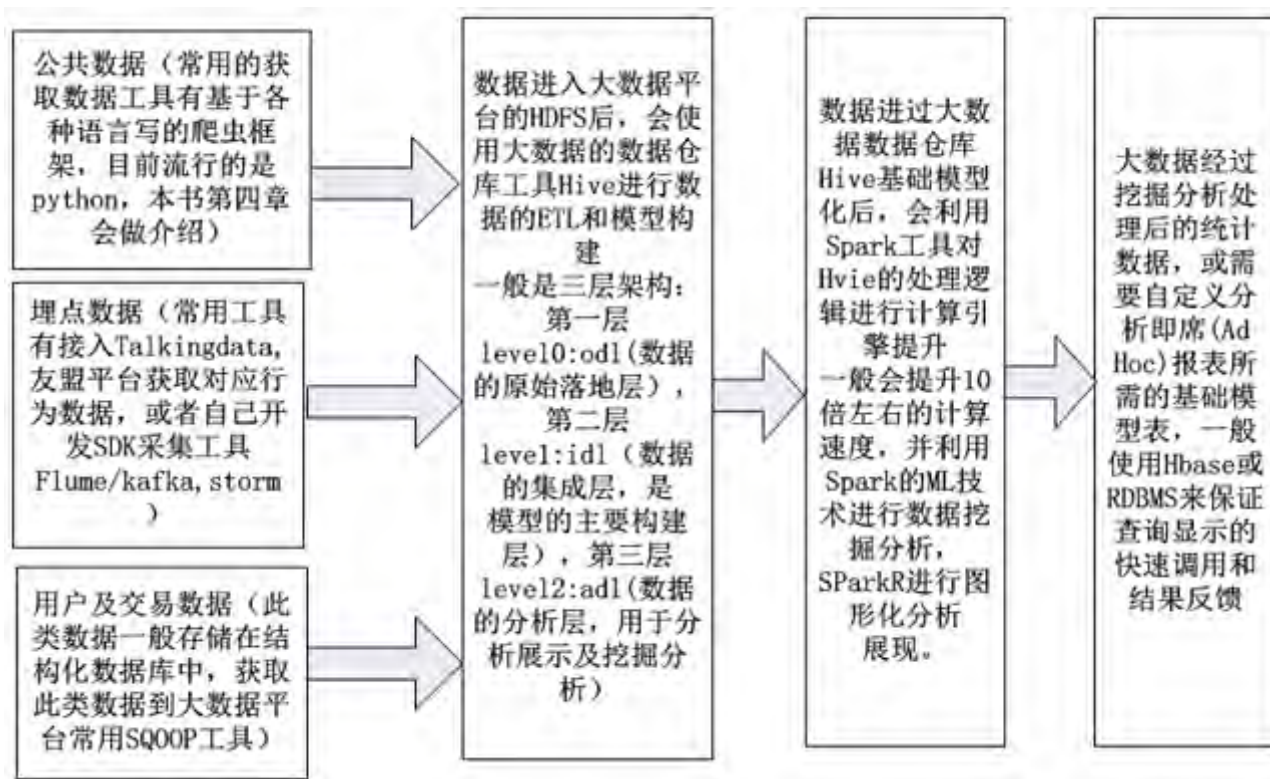
1. 传统数据仓库的数据处理技术及思考
2. 大数据环境下对于公共数据及行为数据的数据处理技术
3. 由传统数据仓库到大数据数据仓库的数据处理实践思考及建议



公共数据及行为数据的数据处理技术



公共数据及行为数据的数据处理技术



公共数据及行为数据的数据处理技术

按数据特征分类

■ 结构化数据

定义：目前其实专指的是关系模型数据，即以关系型数据库表形式管理的数据。绝大多数的企业业务数据都以此格式进行存放。

简析：虽然从专业角度讲，结构化就是关系模型的说法并不准确。但针对目前业内现状，还是将其定义为关系模型数据为最为妥当，因为它清晰而准确地代表了我們传统上最熟悉的企业业务数据，基本没有歧义。

■ 半结构化数据

定义：半结构化与非结构化常常一同被提及，两者其实专指所有其他“非”结构化数据。但如果想更加清晰地描述，可以将“半结构化数据”定义为：那些非关系模型的、有基本固定结构模式的数据，例如应用日志文件、XML文档、JSON文档和电子邮件等。

简析：此部分数据可以用程序化格式解析处理，公共数据，行为数据多以此种格式

■ 非结构化数据

定义：除去结构化与半结构化的所有数据，即没有固定结构模式的数据，例如WORD、PDF、PPT、EXL文档，以及各种格式的图片 and 视频等。

简析：区分半结构化与非结构化数据的意义在于，目前在企业内对两者的处理方法（包括存储、访问与分析）是不同的。非结构化数据大多采用内容管理的方法，展示上需要采用对应的组件工具。



公共数据处理的注意点



接口定义加入接口规范变更版本及内容到数据字段中

对于网站抓取或接口调用的变化版本记录有利于对数据准确和完整性的可追溯



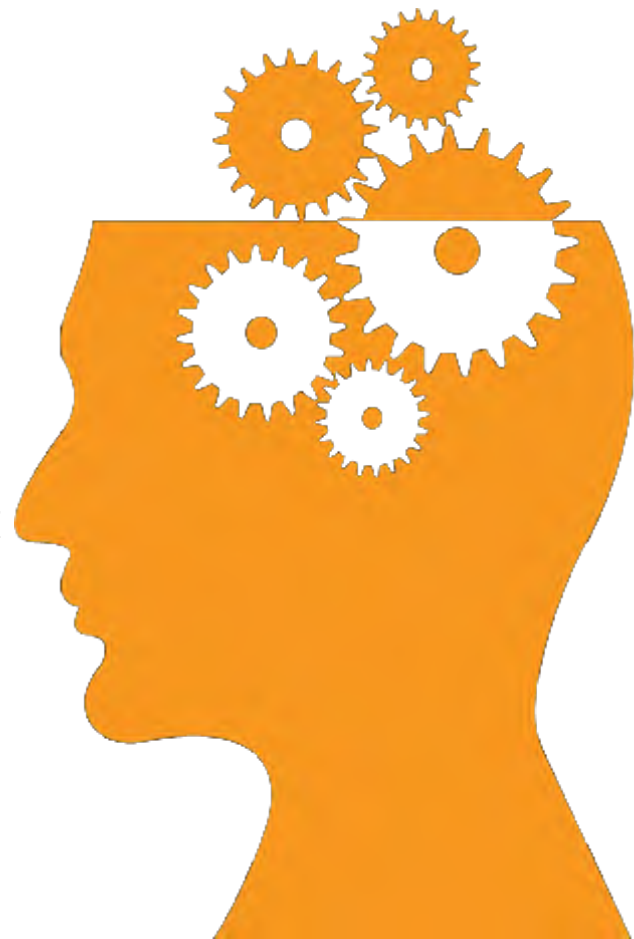
落地后的文件时间和成功标志信息同样参与数据处理

合并到数据落地层 (LEVEL0) 后数据的落地时间和数据大小行数记录到数据监控表中



在数据仓库处理和分析展示中添加数据处理的可追溯信息

对于核心指标及对应元数据显示和监控，确保对于数据的理解和定义全局一致



行为数据处理的注意点



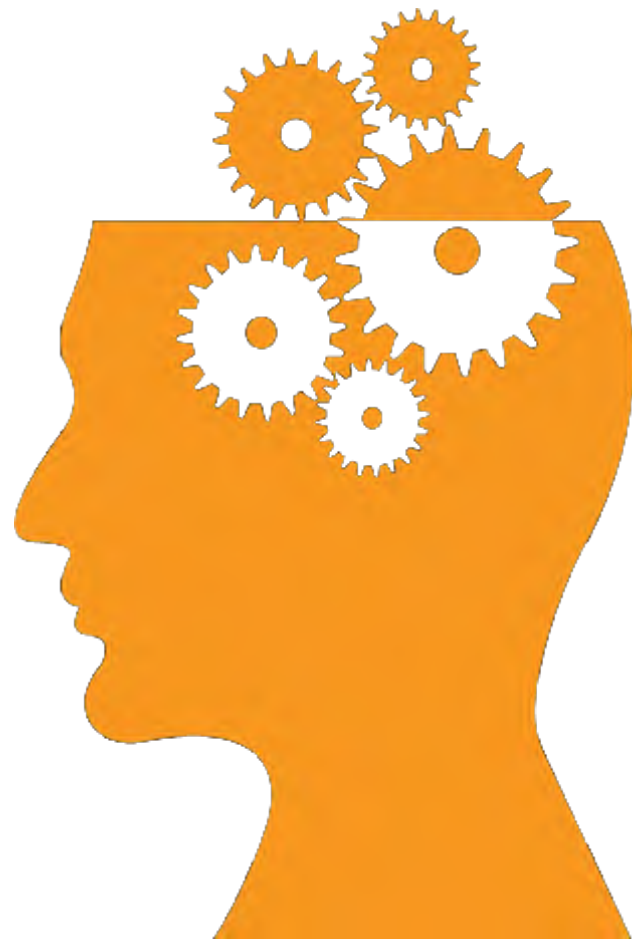
埋点数据一定要符合业务数据信息流才能保证数据处理的完整性和确保数据的业务可用性



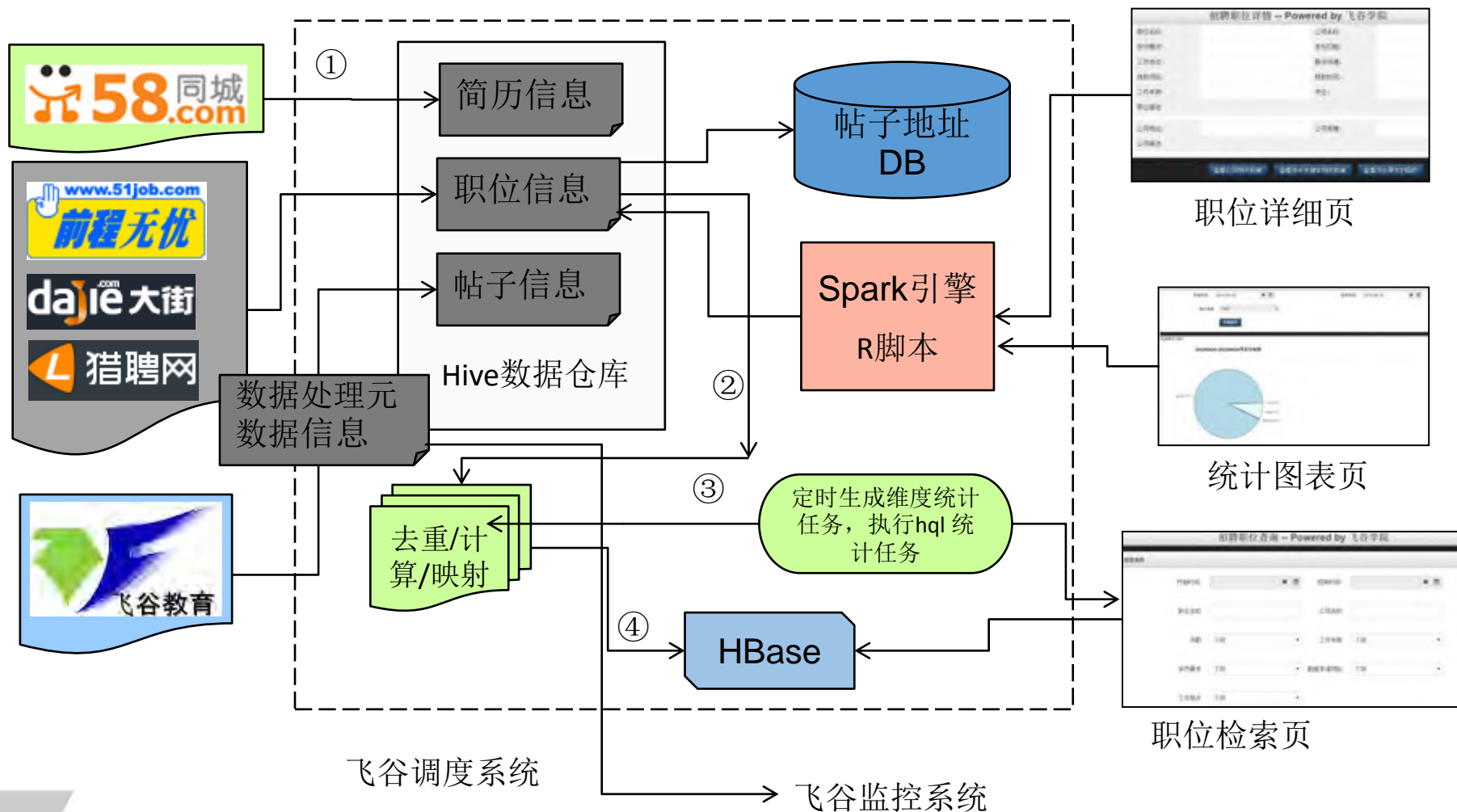
行为数据的标识键 (UID , DID) 要与其它数据源统一关联键和对应时间周期，确保数据的一致性和关联性。



行为数据的元数据信息尽可能从源头以字段化方式植入数据处理的数据文件中



公共数据及行为数据的数据处理技术案例图



数据处理的哪些事

1. 传统数据仓库的数据处理技术及思考
2. 大数据环境下对于公共数据及行为数据的数据处理技术
3. 由传统数据仓库到大数据数据仓库的数据处理实践思考及建议



传统数仓到大数据数仓的数据处理

大数据平台的迁移与构建

传统数据仓库

- 1 以RDBMS为主要的
数据处理存储层。
- 2 数据处理采用通用的
ETL产品工具
- 3 报表层是报表产品通
过标准的数据库连接驱
动连接到数据仓库DB
中。
- 4 数据库安全级别可以
通过RDBMS安全管理



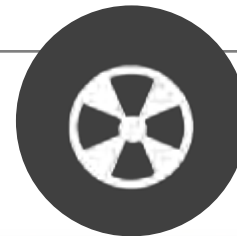
数据同步与脱敏

- 1 数据处理的重点是
全样本数据的基于业
务时间周期同步落地
- 2 基于RDBMS的敏
感数据在进入大数据
平台后进行脱敏处理，
确保数据安全



大数据平台

- 1 以HDFS为主要的数
据处理存储层
- 2 数据处理根据数据来
源采用不同工具，其中
同步RDBMS数据的
SQOOP，实时数据采
集的 Flume，
kafka,storm，及公共
数据的接口API
- 3 报表层产品采用自主
开发或支持大数据平台
的工具tablelu等
- 4 数据安全管理工作大，
需要全局设计控制



由传统数据仓库到大数据数据仓库的数据处理思考及建议

构建数据平台时的数据基因一定要准确完整，这是整个数据平台的根基

**数据基因
定义完整
准确**

01

**数据血缘
设计清晰
可溯**

02

数据平台的数据处理开始就需要同业务数据流一同设计数据的元数据血缘流，确保业务数据断点可查可控

03

**数据安全
机制原子
化**

对数据平台的分层数据做到基于存储机制的原子化安全控制，确保从底层实现数据的安全分层控制。主数据及业务权限数据等

04

**核心指标
及元数据
做到可视
化和监控
自动化**

可视化设计时除了正常的业务数据报表外，对于主线重要的业务元数据及技术元数据的信息同样要做可视化设计，并加入自动化监控内容中。



关于飞谷云

飞谷云是大数据爱好者的家园，是共同有着‘诚信进取协同分享’文化的码农们聚在一起共同打造的大数据学习实践云平台，旨在帮助大学生或需要职业技能提升的码农们通过飞谷云平台（老师，实战环境，大数据生产项目）达到企业大数据相关岗位的技能要求。我们正在进行各大学和中小企业的免费公益交流活动，欢迎各大学社团或院系组织及企业联系我们，一起交流合作，一起让大数据落地！



大数据公益班

大数据生产项目

大数据人才服务

@飞谷云邮箱: feigu@vip.163.com





THANKS
THANKS

SequeMedia
威拓传媒

IT168.com

ChinaUnix

ITPUB