

QCon 全球软件开发大会 【北京站】2016

社会化数据的混合存储和高效处理

明略技术合伙人 任鑫琦

SPEAKER

任鑫琦，2009年硕士毕业于北京大学计算机学院，先后在百度、斯伦贝谢担任研发工作，后加入秒针系统负责大数据计算和平台管理；2014年正式加入明略，先后负责了NoahArk、LogM等多款产品研发，目前任大数据关联分析产品SCOPA负责人。

邮箱：renxinqi@mininglamp.com

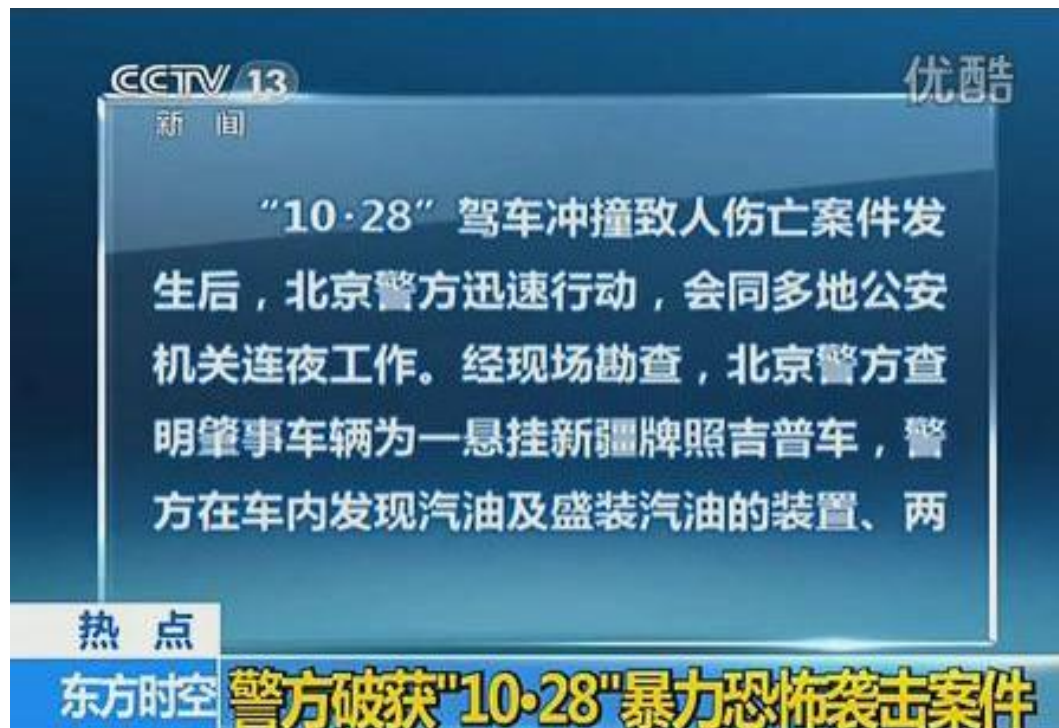
微信：lurkerpku



目录

- 1 社会化数据特点
- 2 社会关系网络的存储架构
- 3 混合存储体系的落地实践

数据到底“大”在了哪儿？



何为“社会化”数据？

“社会化数据” \neq “社会化媒体数据”

1 互联网社会数据

新闻媒体数据

社交网络数据

消费行为数据

物联网数据

。 。 。

2 现实社会数据

实名制轨迹类数据

公共安全类基础数据

基础设施建设监控数据

。 。 。

为何“社会化”数据？



公安情报分析



反恐



风险控制



企业内审

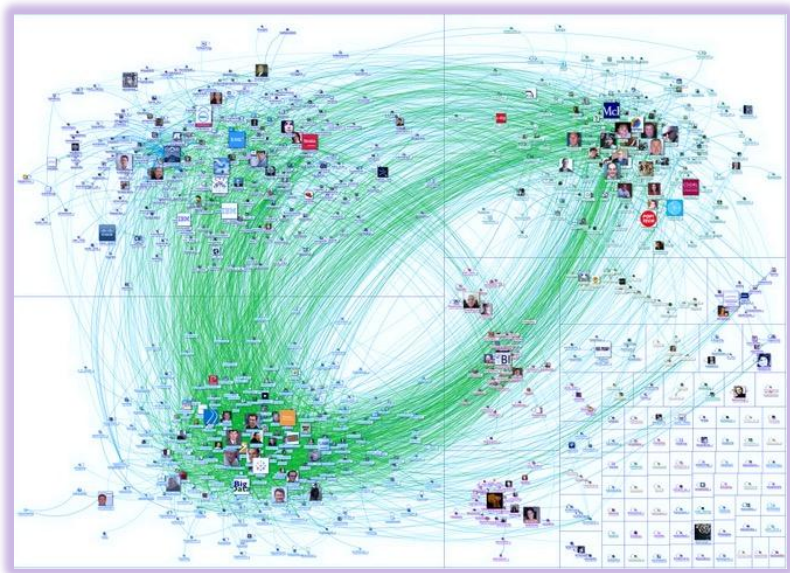


金融反洗钱



反偷税漏税

企业面临的数据难题



如何做数据关联分析？

如何透过多层次、多维度的数据分析实现对于某一个人、某一件事或某一种社会状态的现实态势的聚焦，在时间序列上离散的、貌似各不相同的数据集合中，找到一种或多种与人的活动、事件的发展以及社会的运作有机联系的连续性数据的分析逻辑。

从业务视角看数据

数据离散，价值低，挖掘和变现能力难度大

真正的大数据分析，全面系统包括挖掘，分析，关联等

数据特征与局限

变更困难

某个数据集，某个属性发生变化时，接入-治理-表结构-服务程序-业务一系列功能都需要修改。

数据质量

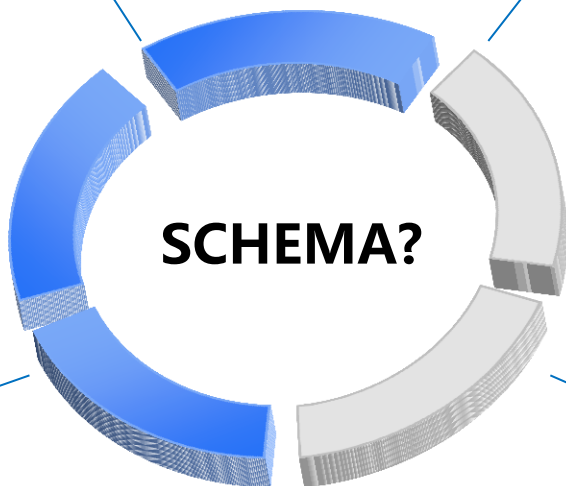
清洗过程中会将不符合规则的数据进行删除、修改，虽然符合了目标数据库，然而这种操作的正确性却很难保证。

非结构化数据

传统数据整理，通常无法有效处理文本型或日志型数据，造成大量有价值数据的流失。

工作量大，性能差

数据抽取、数据修改、数据入库、数据统计分析等过程独立到不同的工具，当增量数据量大时，数据延迟会非常严重。



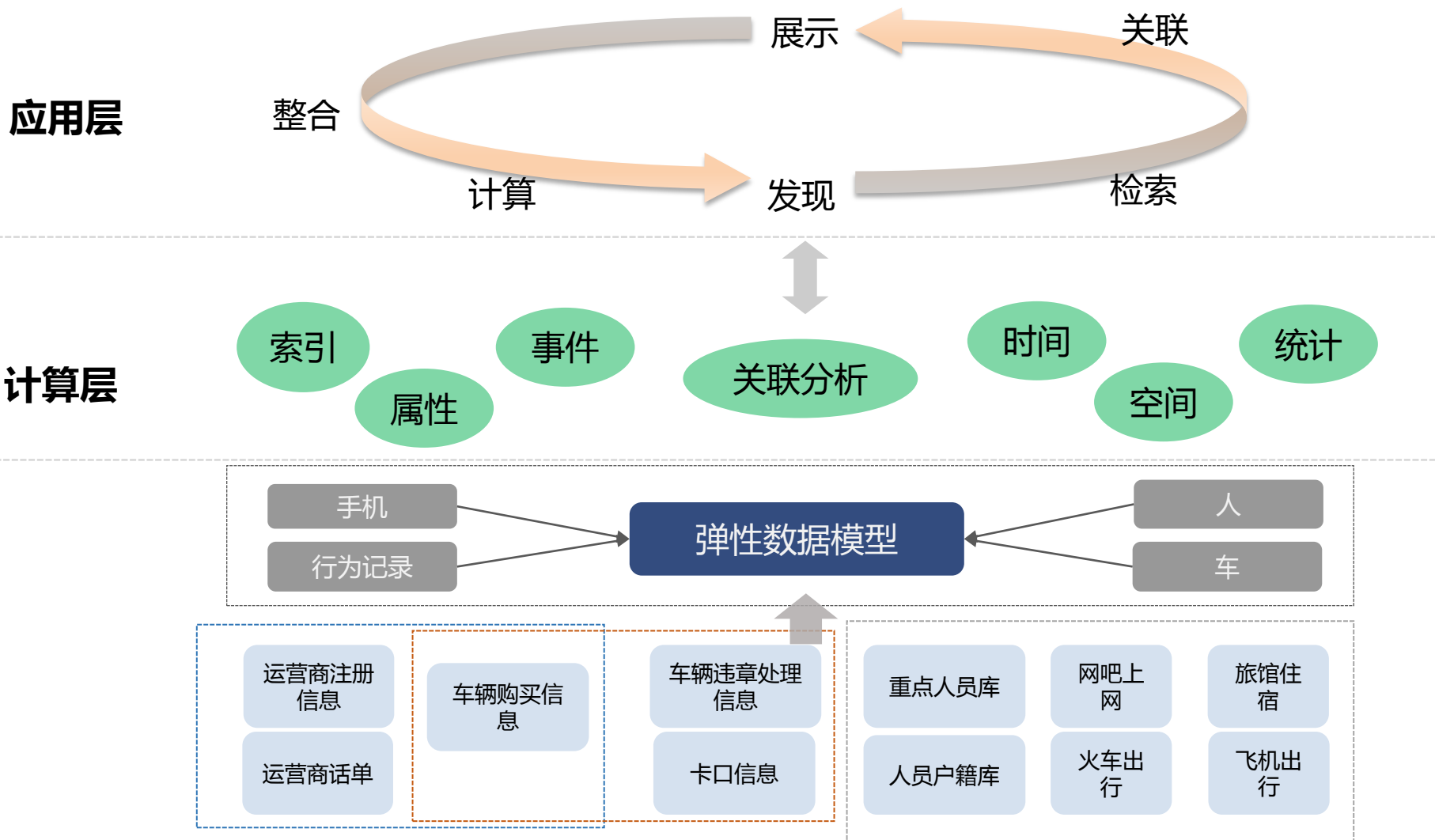
目录

1 社会化数据特点

2 社会化关系网络的存储架构

3 混合存储体系的落地实践

数据应用的过程



数据模型 — 对象 — “本体”

Ontology : 某一领域内的研究对象及其之间的联系

计算机领域

数据库、知识工程

生物学

门、纲、目、科、属、种

公安

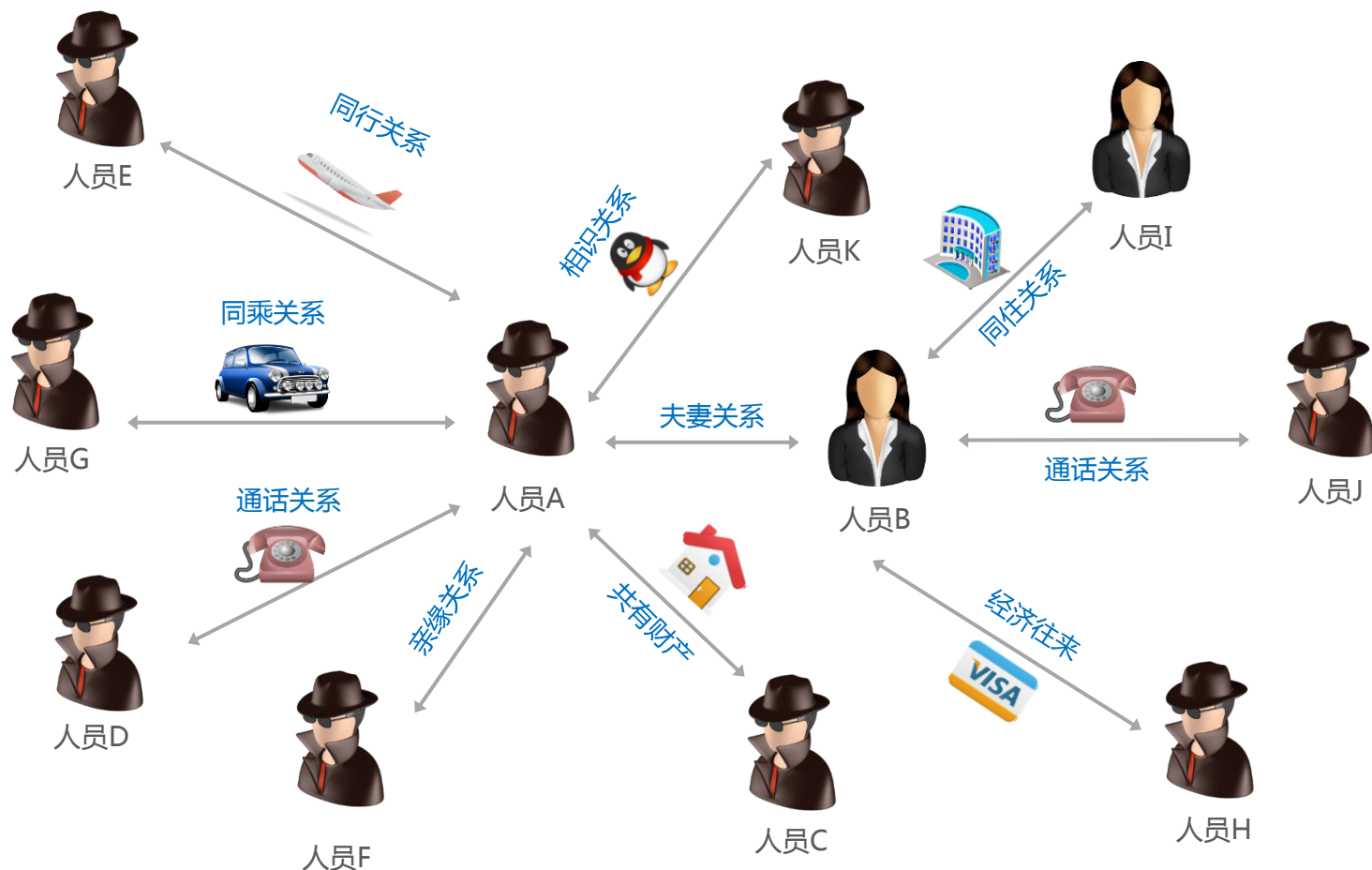
人、事、地、物、组织

弄清数据的本质：概念 — 关联 — 应用 — 表示

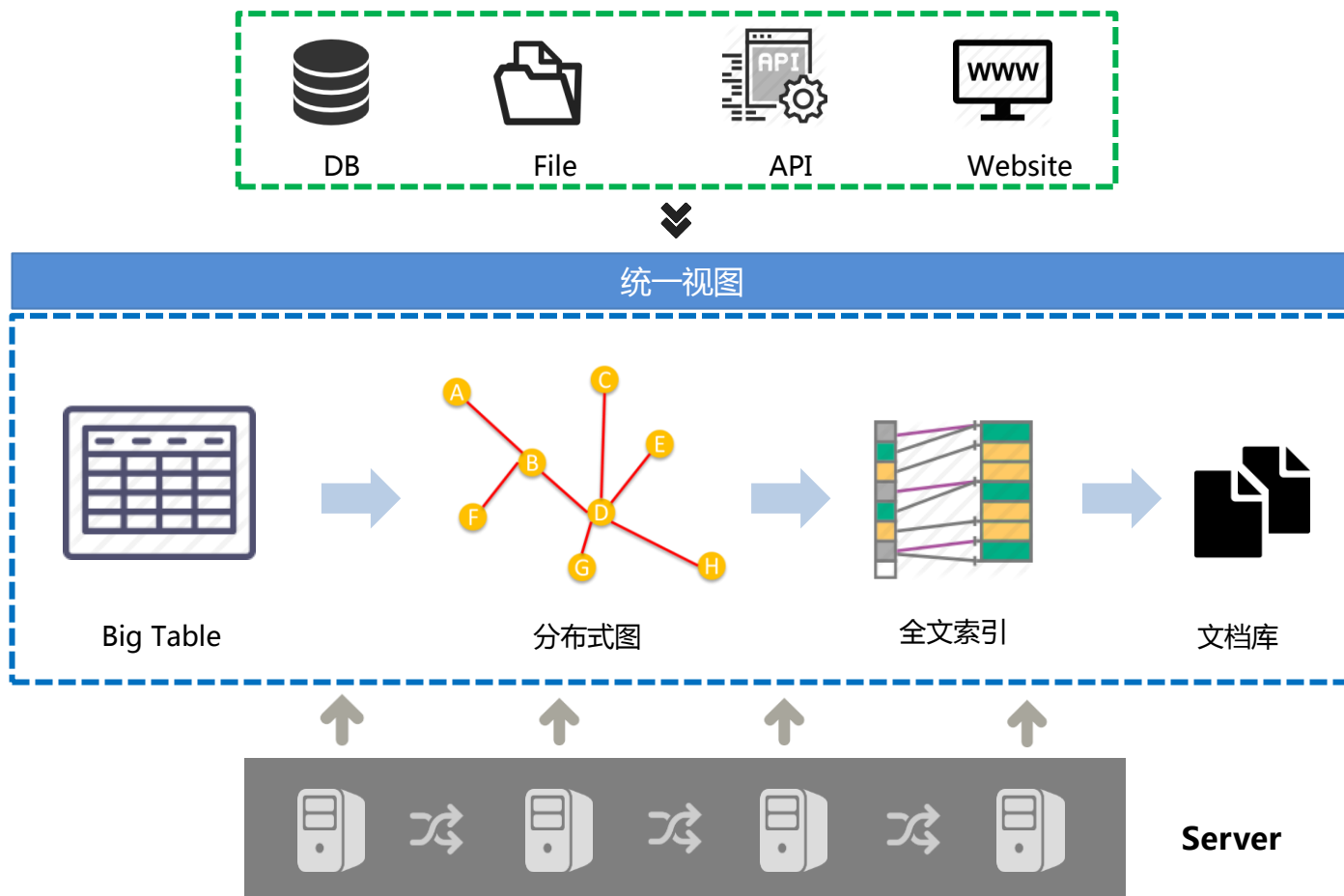
发现“对象”的联系



弹性的社会化关系网络



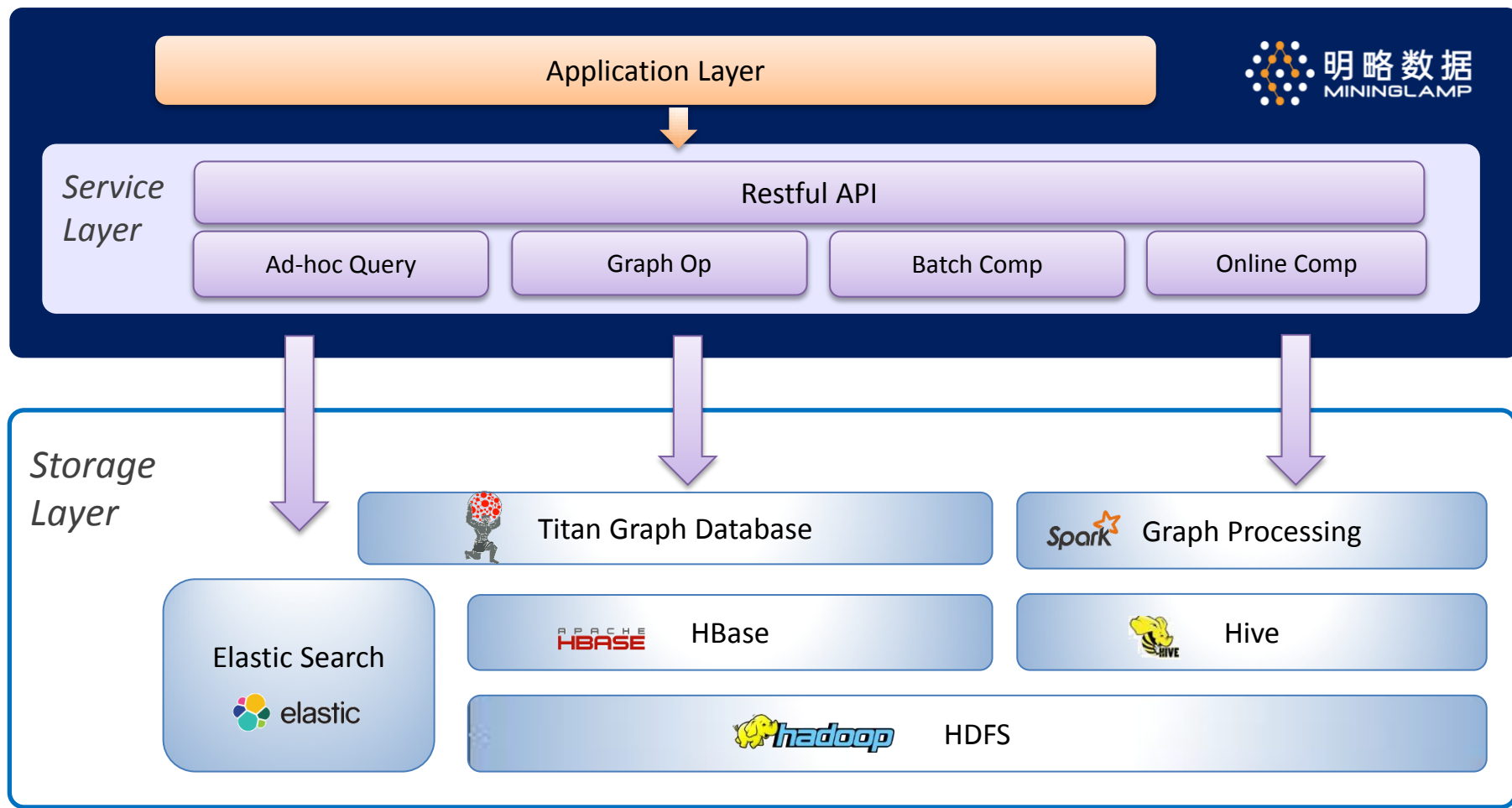
如何把整个网络存下来？



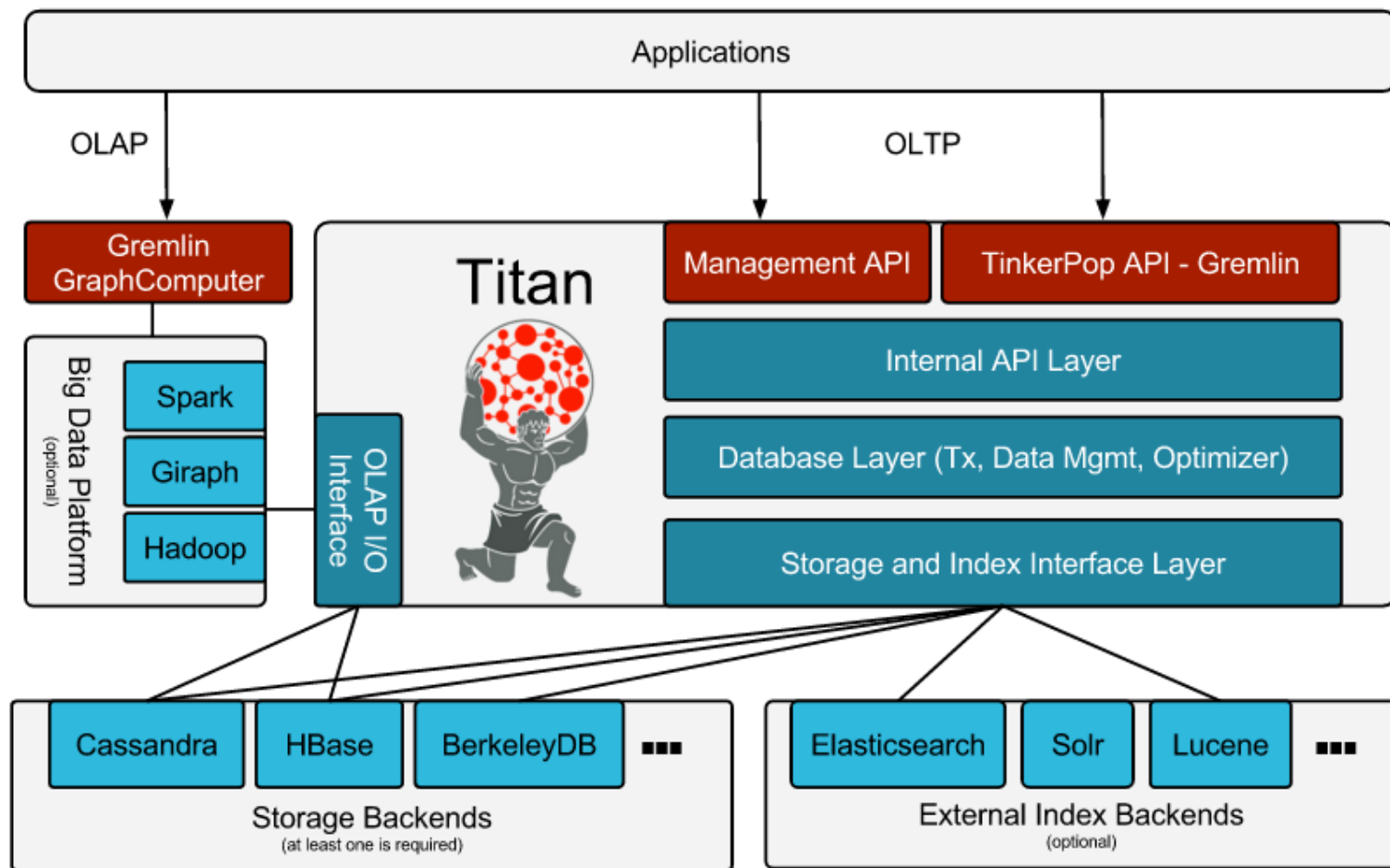
目录

- 1 社会化数据特点
- 2 社会化关系网络的存储架构
- 3 混合存储体系的落地实践

混合存储体系



认识Titan



趟过的坑（1）

1 边爆炸问题

- Titan使用邻接表存储点边
- 边经过编码后存储和查找代价仍然很大
- 同类边合并
- 原始信息使用其他存储

2 Super Node问题

- Titan对Super Node有优化，但效果不好
- Super Node标记

趟过的坑（2）

3 多点查询效率

- TinkerPop
- TitanMultiVertexQuery接口
- 根据场景特定优化

4 索引性能和灵活度

- 优化相应组件
- 复杂索引统一管理

趟过的坑（3）

5 导入数据性能

- 并行导入 = 多进程+多线程
- 划分子图, 并行处理
- 打开batch-loading, 解决一致性问题
- 优化参数, 提高效率
- 调优底层存储, 如避免hotspotting等

整体写数据性能

批量数据插入效率 （单位 个数/秒）

	Single	MT	On Yarn
实体 (50M)	1020	6723	17202
关系 (100M)	540	3812	8710
事件 (300M)	-	13832	41230

测试集群环境：5台 Intel Xeon E52620 (24Core) 128G内存

多线程插入使用24线程

On Yarn 申请10个Container，每个10G内存

实时查询性能

几个典型的实时查询与计算场景

	Average	Median	Min	Max	Error %
根据主键值查询实体	32	22	8	102	0
图析中按标签过滤的1跳推演	61	26	16	310	0
图析中3跳推演	79	32	22	610	0
100个实体间的3条内路径查找	2101	1324	619	13120	0

测试数据120M实体，300M关系，100线程并发，循环10次



明略数据

MININGLAMP



明略数据 是专注于关联关系挖掘的大数据解决方案提供商

以自主研发的安全大数据平台MDP为基础，围绕数据关联分析挖掘产品SCOPA和分布式数据挖掘系统DI，凭借明略大数据科学家丰富的多领域知识积累，实现明略独特的挖掘复杂数据价值的能力，帮助政府、公安、税务、金融机构等客户，在安全可靠的环境下，整理、分析、利用不同来源的结构化和非结构化数据。核心理念在于，利用数据的连接性，挖掘数据间的关系，激发大数据的真正价值，从而创造一种人脑智慧和计算机智能“共生”的关系，发挥两者各自的特长，帮助中国及中国企业解决实际、困难的、最重要的发展问题。

明 万 象 · 筑 方 略

THINK BIG | DATA

START SMALL



THANKS!