

阿里分布式数据库双11实战

- 王晶昱（花名：沈询）
- 阿里巴巴 中间件&稳定性平台 – 资深技术专家
- 8年
- 阿里企业级应用平台技术总监
- 淘宝分布式数据库(TDDL/DRDS)
- 淘宝分布式消息系统(Notify/ONS)
- Weibo: @淘宝沈询_WhisperXD



- 历史长河中数据存储发展
- 数据库的未来方向展望
- 双11中的阿里分布式数据库DRDS

可以存尽可能多的数据
可以很方便的取出数据



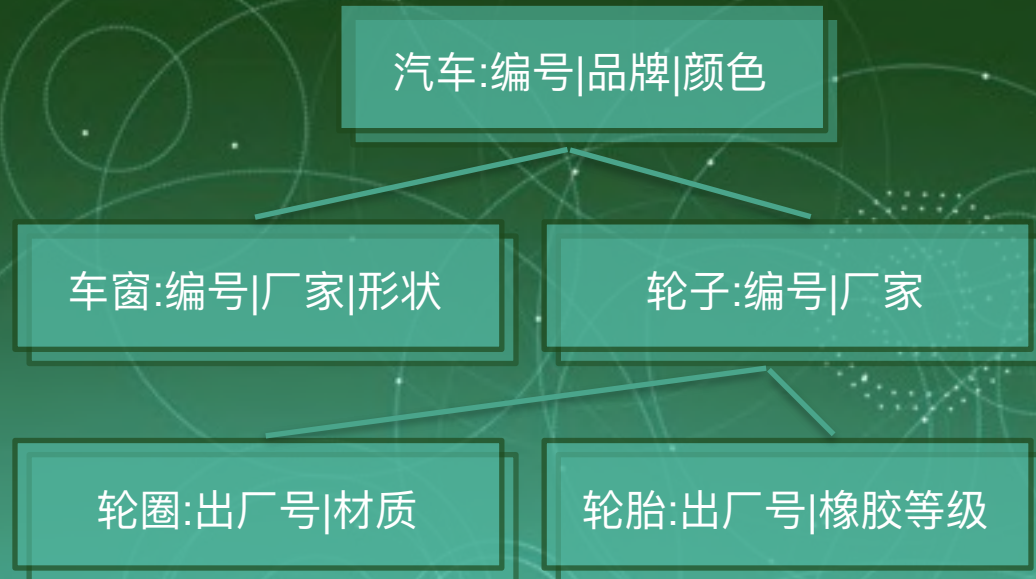
方便:

找到所有轮子生产厂家
为DRDS的汽车编号

```
select * from 汽车表  
a, 轮子表 b, on  
a.id=b.id where b.厂家  
= 'DRDS'
```

代价:

性能损耗



性能调优:

```
select * from 汽车表  
a, 轮子表 b, on  
a.id=b.id where b.厂家  
= 'DRDS'
```

轮子:编号|厂家

$O(n)$

轮子表:
编号>轮子编号,厂家名

汽车:编号|品牌|颜色

$O(\log_2 n)$

汽车表:
编号>品牌名, 颜色

性能调优:

```
select * from 汽车表  
a, 轮子表 b, on  
a.id=b.id where b.厂家  
= 'DRDS'
```

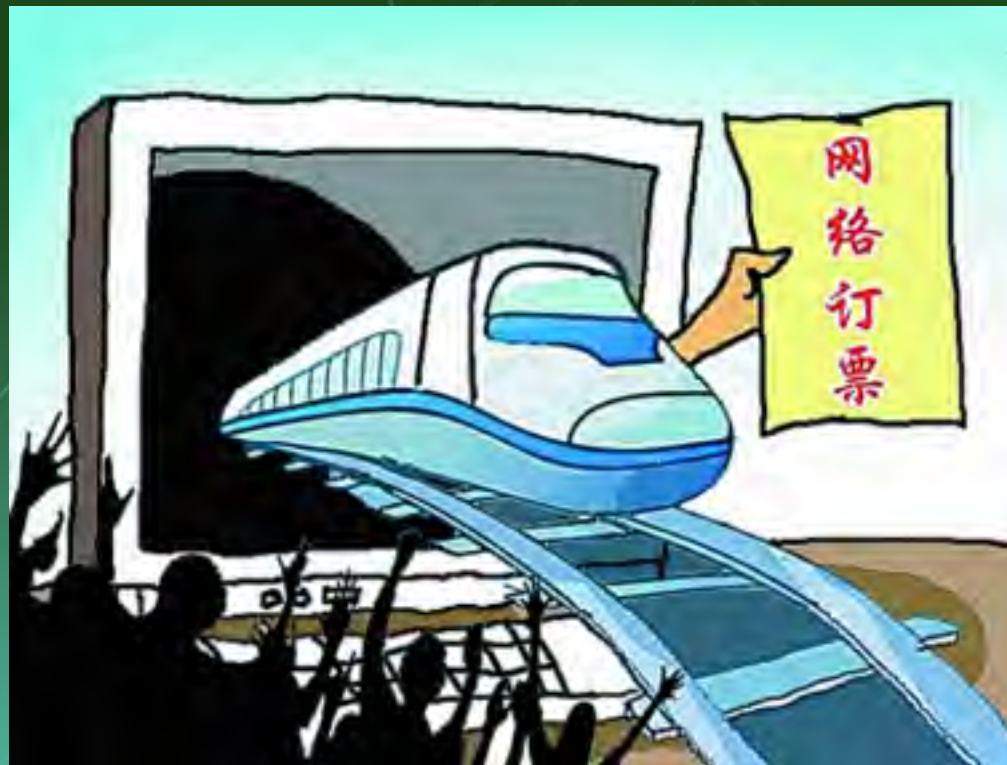


移动/互联网直接服务60
亿人口。

传统关系数据库性能瓶
颈

数据量

访问量



数据必须分散到更多机
器

读写能力要水平伸缩



更多的网络访问



L1 cache reference	0.5 ns
Branch mispredict	5 ns
L2 cache reference	7 ns
Mutex lock/unlock	25 ns
Main memory reference	100 ns
Compress 1K bytes with Zippy	3,000 ns
Send 2K bytes over 1 Gbps network	20,000 ns
Read 1 MB sequentially from memory	250,000 ns
Round trip within same datacenter	500,000 ns
Disk seek	10,000,000 ns
Read 1 MB sequentially from disk	20,000,000 ns
Send packet CA->Netherlands->CA	150,000,000 ns

同一台机器

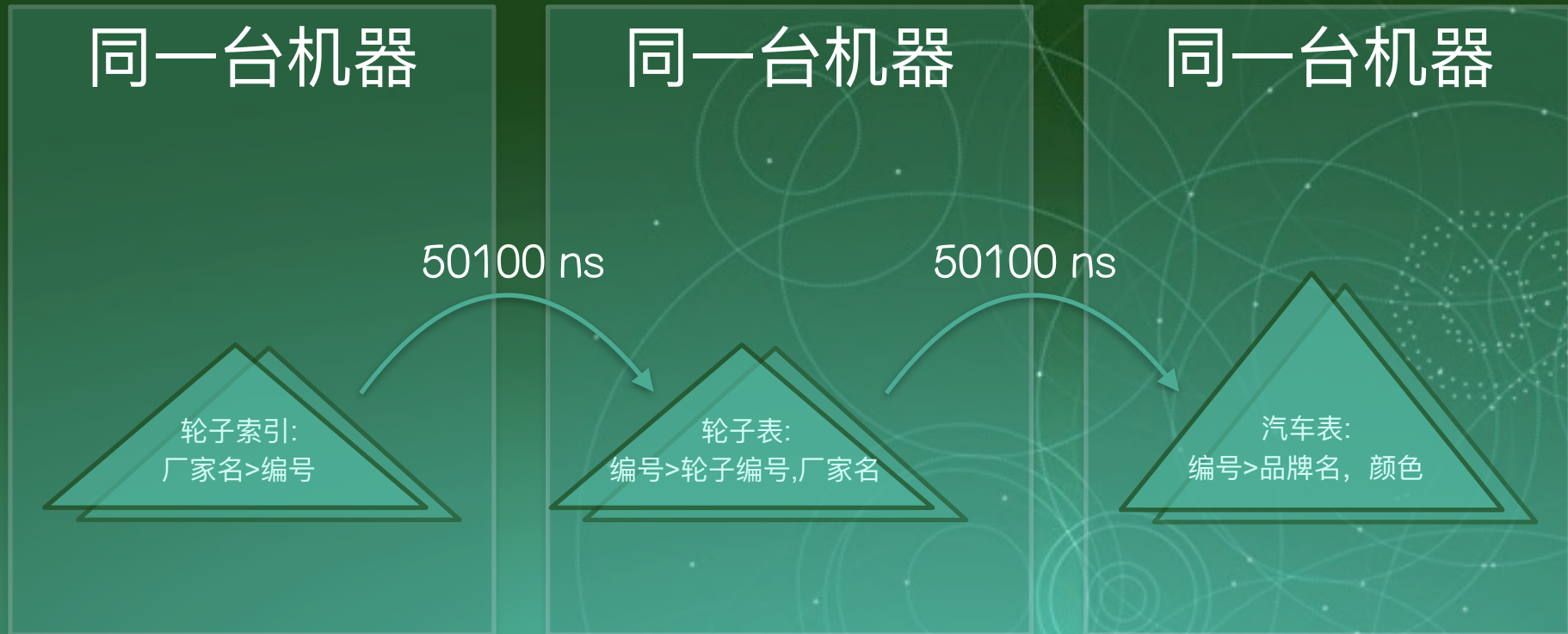
100 ns

100 ns

轮子索引:
厂家名>编号

轮子表:
编号>轮子编号,厂家名

汽车表:
编号>品牌名, 颜色



No SQL!

sdcc.csdn.net

No SQL!

xxx is a Web Scale db
doesn't use SQL or
Joins

xxx is web scale.

xxx is web scale.

xxx is web scale.

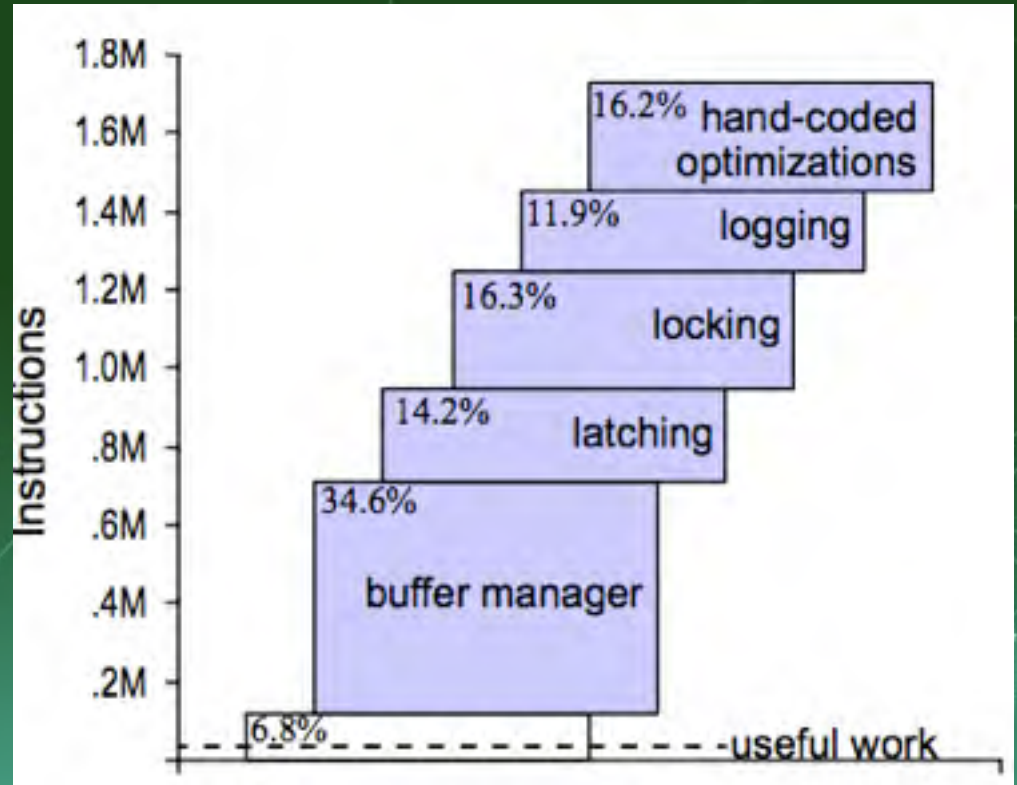


如何到达高性能?

事务只占了2x%

优化只占1x%

不提供Join,Join的需求
也是客观存在的

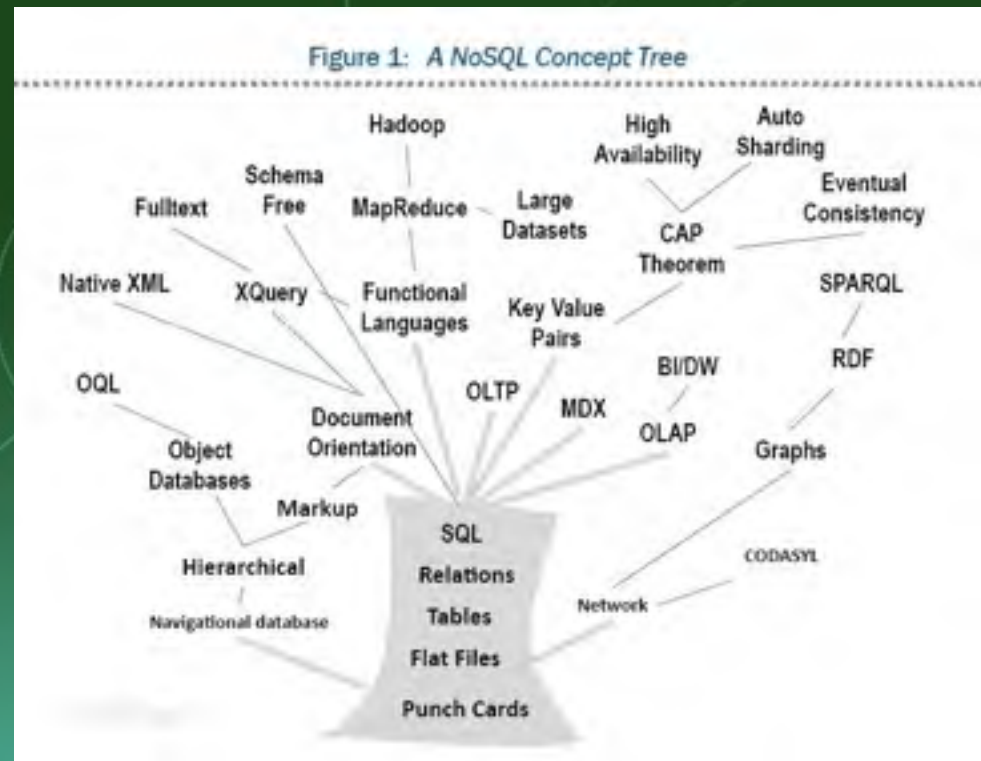


OLTP Through the Looking Glass, and What We Found There

Not Only SQL!

sdcc.csdn.net

合久必分
战国时代



New SQL!

sdcc.csdn.net

OLTP

NoSQL + +

SQL - -

关系数据库

全SQL支持

可控事务支持

NoSQL数据库

高可扩展性

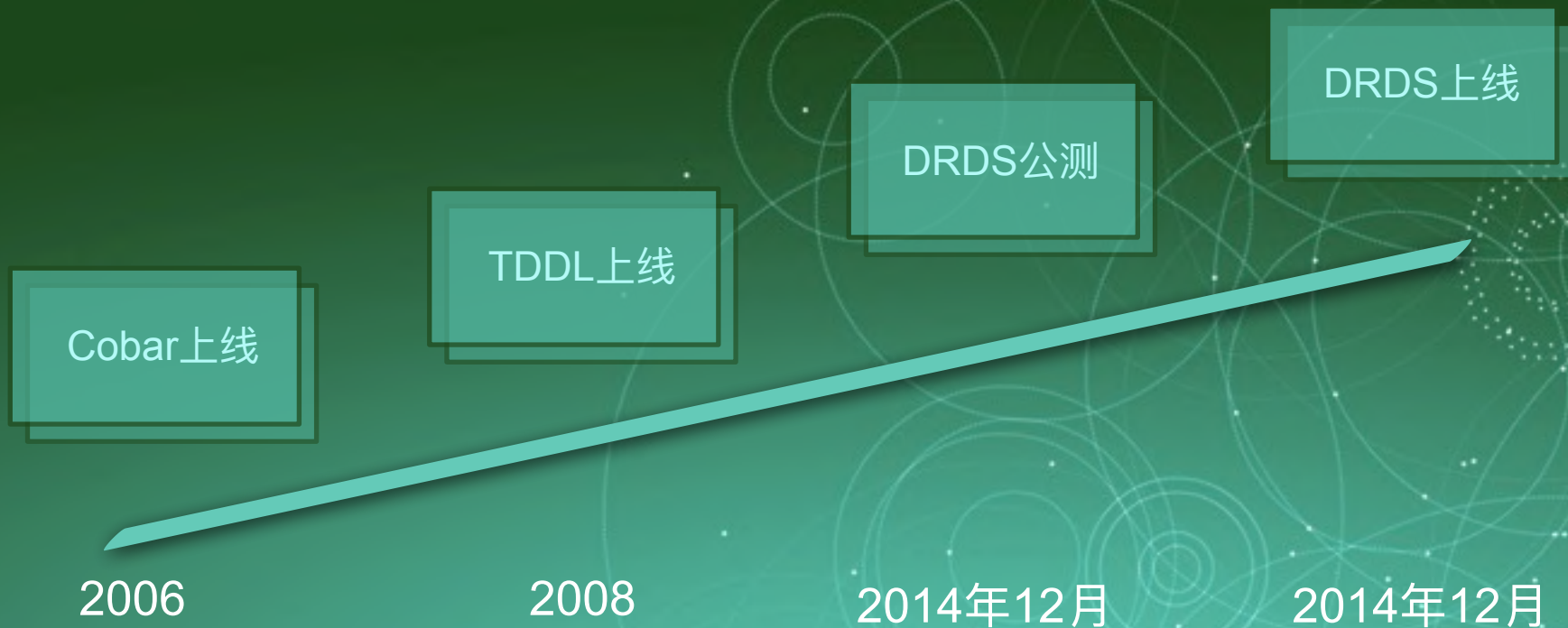
schema less

DRDS是New SQL

DRDS 是一个产品

集合众家之长

稳定运行多年



海量的数据库节点如何
运维？

我的proxy经常出现故
障怎么办？

多年线上稳定运维经验

无单点高可用设计

运维能不能更简单?

配个新的规则还需要
登录机器?

加个表需要重启?

忘记增加列?

DDL语句支持

平滑控制台体验

MySQL图形工具兼容



2pc性能太低了，有没有更好的?

阿里多年积累的 分布式事务体系



- 历史长河中数据存储发展
- 数据库的未来方向展望
- 双11中的阿里分布式数据库DRDS

DRDS(TDDL)从09年开始就一直在参与双11

在双11中一直保持稳定。



DRDS 的核心价值在双11
之前

平滑动态扩容

最大集群有几百台机
器

DRDS 的核心价值在双11
之后

机器立刻归还



线上运维保障

容量评估

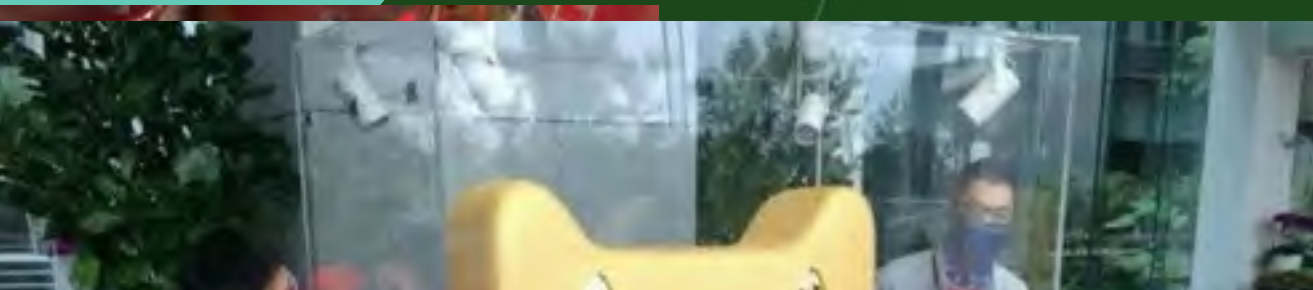
复杂SQL性能优化

高可用

集群切换

机房多活容灾演练





大家辛苦了



同志们辛苦了

谢谢聆听

shenxun@taobao.com