



淘宝数据仓库环境介绍

陈吉平

淘宝网首席DBA，资深技术专家

Oracle ACE

《搭建Oracle高可用环境》一书作者



淘宝网 Taobao.com
阿里巴巴旗下网站

淘宝介绍

- 淘宝网（www.taobao.com）是国内领先的个人交易网上平台，由全球最佳B2B公司阿里巴巴公司投资4.5亿创办，致力于成就全球最大的个人交易网站。自2003年5月10日成立以来，淘宝网基于诚信为本的准则，从零做起，在短短的半年时间，迅速占领了国内个人交易市场的领先地位...

设淘宝为首页 您好, virpure! [请登录] [免费注册] 站内信 收藏夹 打听 交易安全 帮助中心

我要买 | 我要卖 | 我的淘宝

淘宝网

 消费者社区 | 支付宝 | 阿里旺旺

宝贝 | 店铺 | 商城 | 拍卖 | 全球购 | 打听

搜索 高级搜索 使用帮助

首页 | 商城 | 全球购 | 拍卖 | 礼物 | 促销 | 机票 | 彩票 | 资讯 | 男人 | 女人 | 亲子 | 服饰 | 美容 | 时尚 | 居家 | 手机 | 运动 | 爱心

公告栏

火热报名 淘宝“生”了，你呢？
淘宝创57万就业机会
淘宝获客户体验大奖

最新服务

- 手机淘宝
- i淘宝
- 试衣间
- 宝贝传奇
- 美容心得
- vip卡

·丰胸真实体验 · 半年瘦瘦瘦！
· 09春装发布 · 甜心冬服3折

Spring 流行春季风

as 热卖开襟短外套 (8色)
价格: **¥25.00**

安都超人气主打款 1 2 3 4 5

春装全新发布
· 春装流行发布会
· 全家老小都爱!
· Jackjones 3.4折
· 一百元日本购物

裙裤专场

春装全新发布

品牌鞋包流行特辑
幸福恋爱季
vol. 30

游戏快充 手机充值 机票 彩票

游戏点卡 游戏币

游戏: 请选择游戏
面值: 请选择面值

查看淘宝折扣价

消费者保障 假一赔三 七天无理由退货 限时抢购 抵价券

宝贝类目 店铺 淘1站 网游 手机 数码 化妆品 服装 美容 母婴

拼音索引 a b c[ch] d e f g h i j k l m n o p q r s[sh] t u w x y z[zh]

精彩活动 新奇宝贝 特价宝贝

574M 缩放

议程

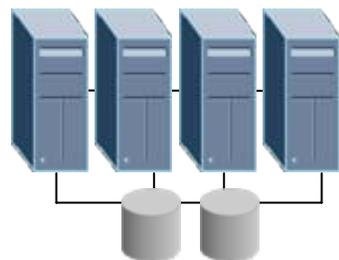
- 淘宝数据仓库的发展历史
- 淘宝数据仓库的现状
- 淘宝数据仓库的维护
- 淘宝数据仓库的未来发展

淘宝数据仓库的发展历史

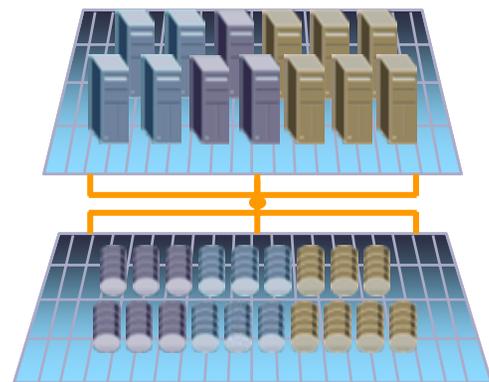
单机-
DELL6650



IBM-P550



4节点RAC环境



12节点RAC环境
计划扩展到20节点

2004

2006

2007

2008

为什么选择Oracle RAC

➤ 当时的可选方案

- Teradata, 在ebay有成功案例
- Oracle Rac, 对淘宝来说, 技术比较成熟

➤ 主要考虑的问题:

- 1. 效率, 是否能快速的满足日常需求
- 2. 开放性, 开放系统还是选择封闭系统
- 3. 可扩展性, 包括线性扩展能力, 对外接口等
- 4. 成本因素

为什么选择Oracle RAC

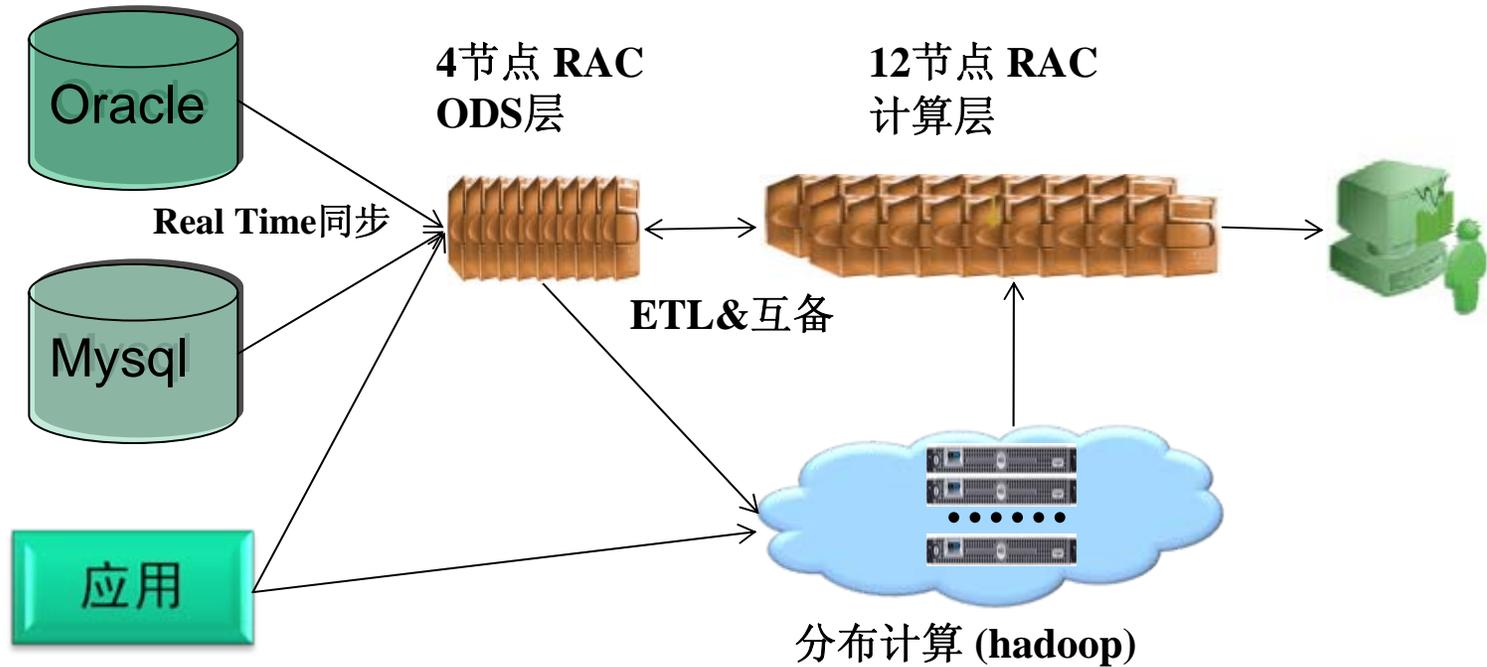
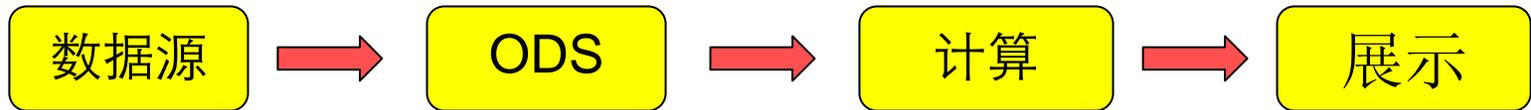
➤ Oracle RAC优点

- Oracle是一个相对开放的系统，文档很多
- Taobao本身对Oracle的技术比较了解
- 在实际的测试中，Oracle表现比较优越
- 可以做到线性的扩展能力
- 成本不会高于Teradata

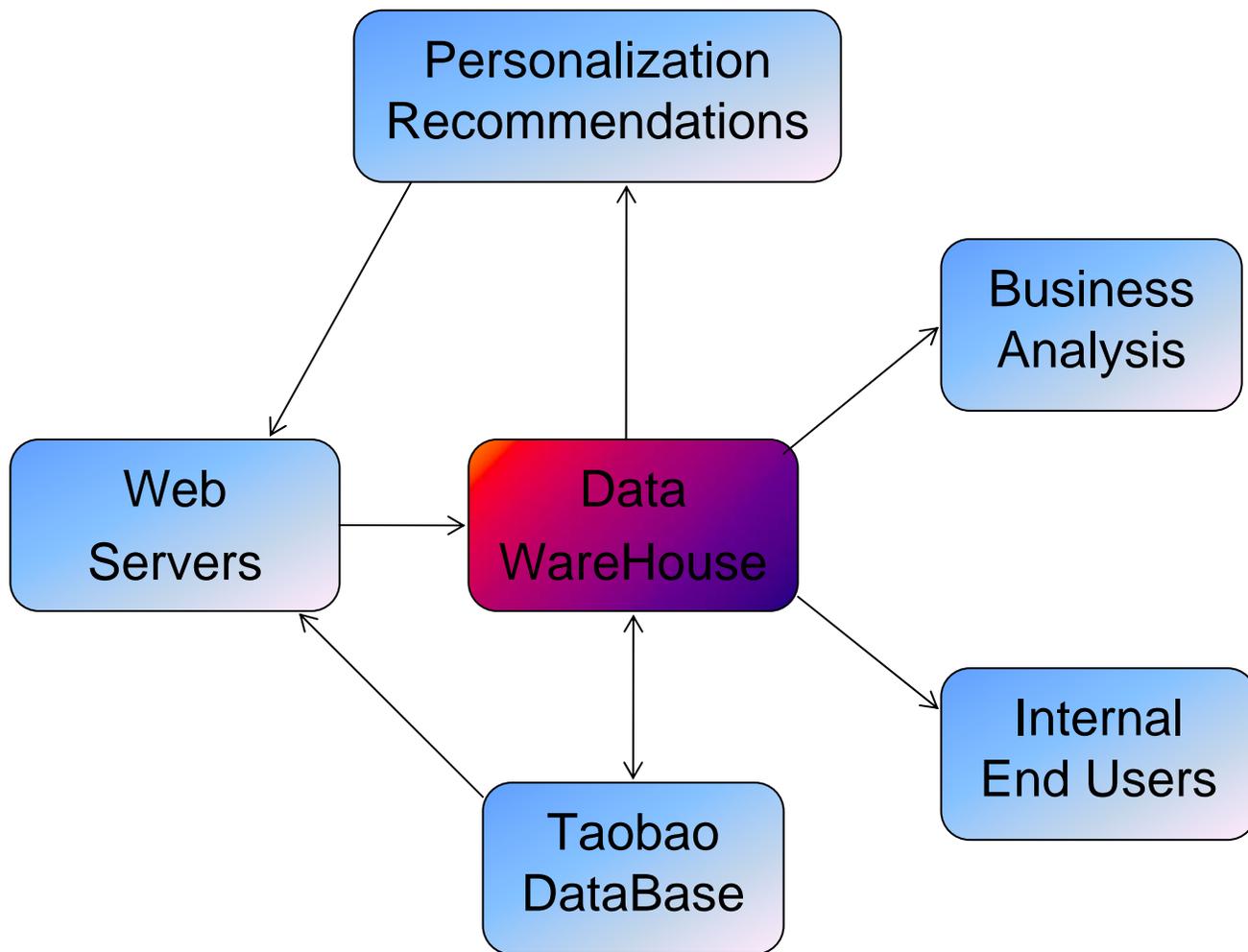
➤ Oracle RAC缺点

- 对技术的要求相对比较高，需要DBA深入的优化

淘宝数据仓库现状



淘宝数据仓库的应用场景



淘宝数据仓库应用案例（一）——淘我喜欢

首页

i 淘宝

Beta

我的收藏

我的空间

提建议

帮助

人气宝贝



移动手机充值卡



毛衣



牛仔裤

欢迎来这儿，快快 **登录** 吧！

让淘宝与你一起，淘宝贝，淘店铺，淘伙伴！唯一的页面只给唯一的你，一切只为，淘你喜欢！

收藏本页

设为首页

淘我喜欢的宝贝

全部

彩妆香水

ZIPPO眼镜

女装

服饰配件



phiten法力藤/法藤钛项圈硅胶防水运动型颈环

¥ 48.00



冲皇冠Hello Kitty钛项圈(Phiten)法力藤法藤钛

¥ 30.00



BVLGARI 宝格丽红茶女用香水 50ml

¥ 110.00



四冠店1001韩版外套OL手感超柔软珍珠糖果色V

¥ 19.00



热门活动

- 1 新春开运时尚穿搭新主张
- 2 施华洛手表新品发布5折上市！
- 3 佩饰达人让你变韩剧佳人！
- 4 情人节-以爱情的名义减肥
- 5 2008淘宝最热卖化妆品终评榜
- 6 百元淘遍日本
- 7 明天要送礼？不知道送啥就点这里！
- 8 09新装全场包邮

送礼,该送啥? **礼品频道**

淘宝数据仓库应用案例（二）—猜你喜欢

收藏该店铺的用户最喜欢的其它宝贝:



¥ 6.99



h61116【特】Fruit of the Earth 100%纯芦荟胶56g

收藏人气: 5552

加入我的推荐



¥ 1.90



J00905脚丫双面吸盘~颜色随机发送

收藏人气: 738

加入我的推荐



¥ 9.49



h60321【特】最爱瘦身教馆冲凉瘦★消脂收紧去橘皮

收藏人气: 1092

加入我的推荐



¥ 1.29



h61119【折】女人我最大韩伊橄榄特润美白保湿护手霜

收藏人气: 2472

加入我的推荐

淘宝猜你还喜欢什么店铺:



双生儿四周年双金冠 来

掌柜: 双生儿香港平价店

信用: 1013848 🏆🏆



心蓝t透双金冠冲击! 全

掌柜: 心蓝t透0

信用: 980570 🏆



『女人天生是妖精』店内

掌柜: 妖精de国度

信用: 567326 🏆

淘宝数据仓库应用案例（三）—内部系统

淘宝网 淘宝网数据仓库 Taobao.com

首页

会员登陆 数据仓库数据一览

用户名:

密码:

 登入数据仓库 SIGN IN

taobao数据仓库常用文档

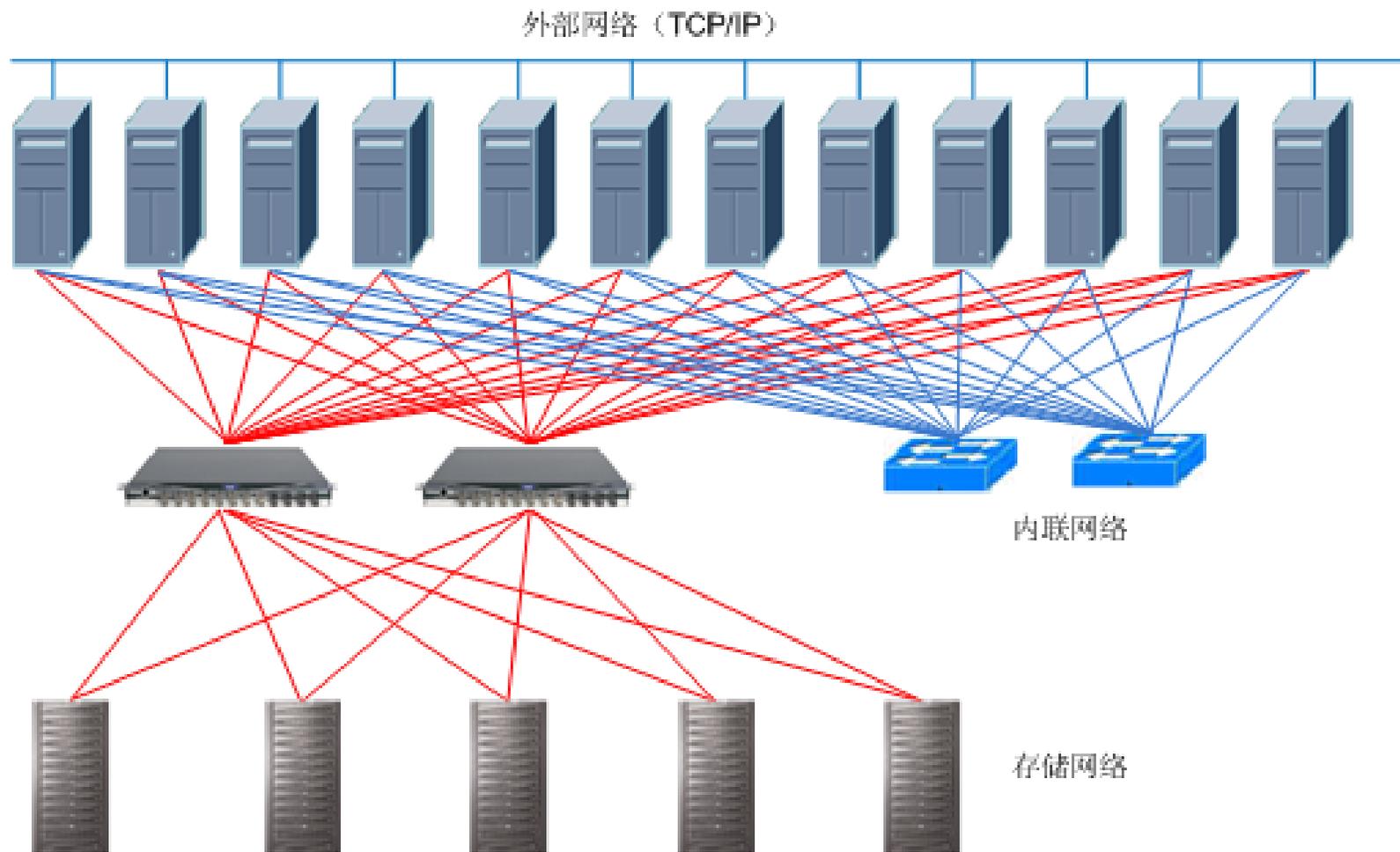
- 淘宝之家
- CRM系统入口



1 2 3

- [导航地图 - 我要看的报表在哪里?](#)
- [报表内容 - 报表中都有一些什么数据?](#)
- [动画课堂 - 怎样查看数据更便利?](#)
- [日常需求 - 数据提取需求入口](#)
- [关于我们 - 数据仓库团队介绍](#)

淘宝数据仓库12节点物理拓扑图

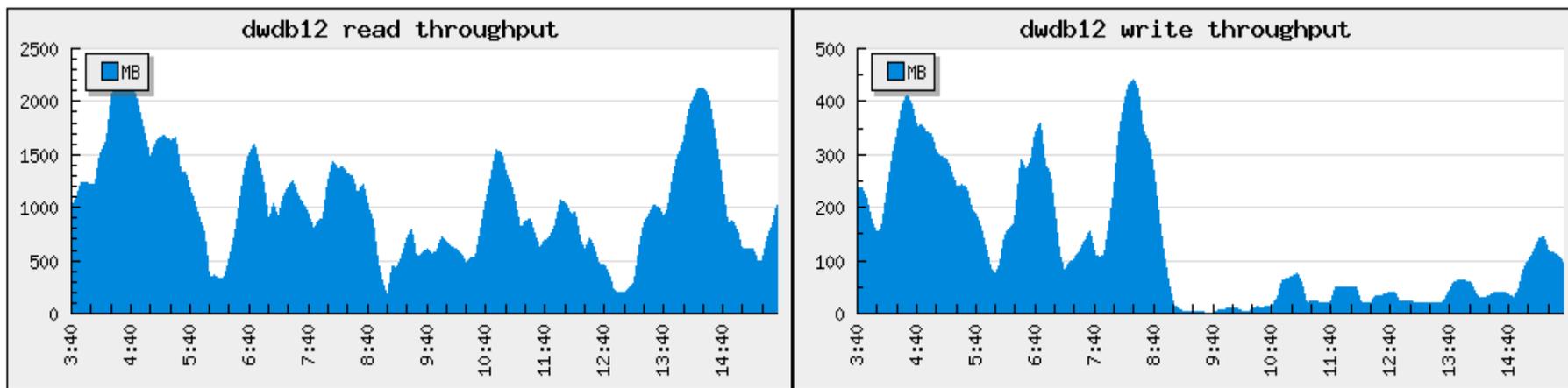


淘宝数据仓库12节点的承载量

➤ 上线前测试情况

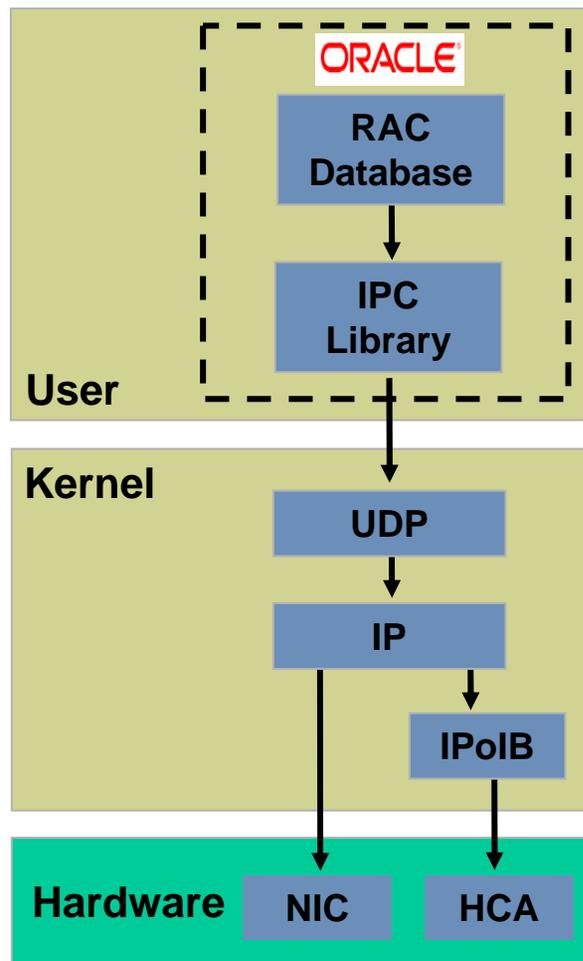
- 5台CX3-80，每台120块146G 15K光纤硬盘，12个节点同时跑Orion测试：
- `orion -run advanced -testname mbps_all_point_20 -num_disks 500 -size_small 1024 -size_large 1024 -matrix point -num_small 0 -num_large 20 -type seq -verbose`
- 读总吞吐量达到7721.26MB/s

➤ 目前使用量：读2.5GB/s，写500MB/s



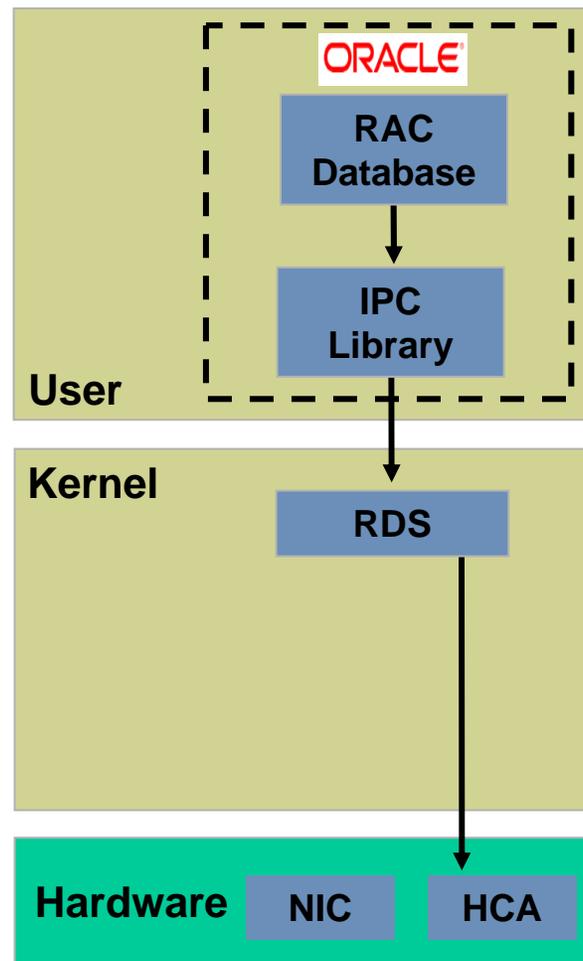
实施细节（一）——内联网络

- 普通千兆网卡
- Infiniband 高速交换



IPOIB

300-500MB/s, 延迟几十us



RDS (certified by oracle 10.2.0.2)

速度1GB/s, 延迟10us

关于infiniband

➤ 对于IPOIB

- oracle端不需要做任何调整
- 检测
 - oradebug setmypid
 - oradebug ipc
 - oradebug tracefile_name
 - netperf -H 10.0.100.104 -l 60(Linux tool)

➤ 对于RDS

- 安装
 - cd \$ORACLE_HOME/rdbms/lib
 - make -f ins_rdbms.mk ipc_rds ioracle
 - 详细见“Oracle Real Application Clusters 10g Release 2: Installation and Configuration of Linux Clusters Using RDS over InfiniBand Interconnect”
http://www.oracle.com/technology/products/database/clustering/pdf/rds_installation_on_oracle_rac_10g_release_2_linux_clusters.pdf
- 对于RDS速度的检测
 - crload/crtest(Oracle Developer team's tool)
 - Oradebug中无法看到RDS信息, oracle BUG 5610067
 - paul.tsien@oracle.com is the PM of RAC & RDS team

实施细节（二）——网格存储ASM

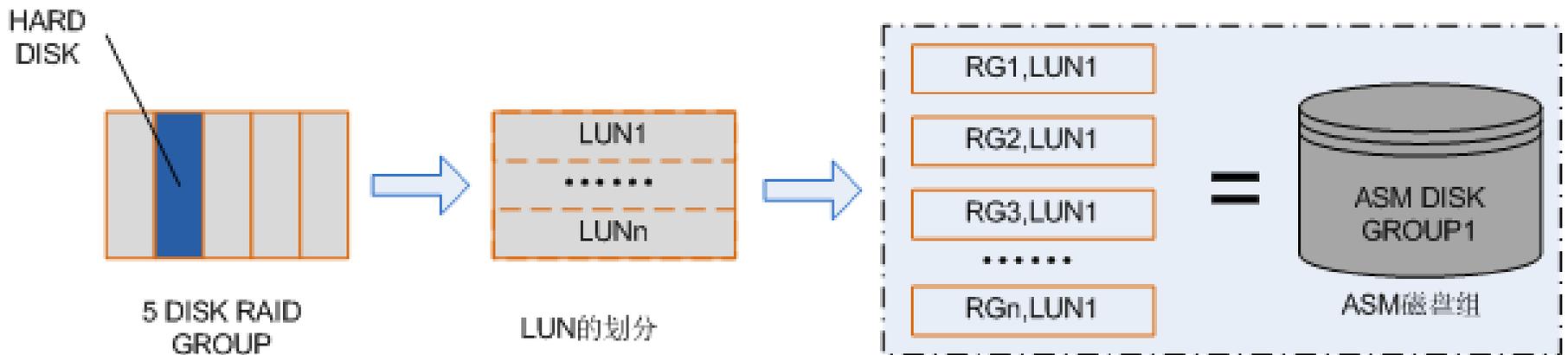
➤ 关于ASM

- ASM的性能十分好和裸设备基本没有差别，在测试中可以把存储性能压榨干
- ASM的strip做的也很好，通过存储监控可以看到I/O分布很均匀

➤ 对于rebalance

- 在RAC中可以利用ASM来rebalance
- `alter diskgroup taobao rebalance power 11;`
- Rebalance速度很快，测试过程中，新增加3TB磁盘的同步用时1.5个小时

ASM磁盘组的划分



- 4+1 的RAID5
- 每个RAID组划分2-4个LUN
- 每个RAID组的第一个LUN组成ASM磁盘组1，第二个LUN组成ASM磁盘组2，依次类推
- 数据库建立在ASM磁盘组上
- 可以确保数据分布在所有的物理磁盘上

并行执行相关参数

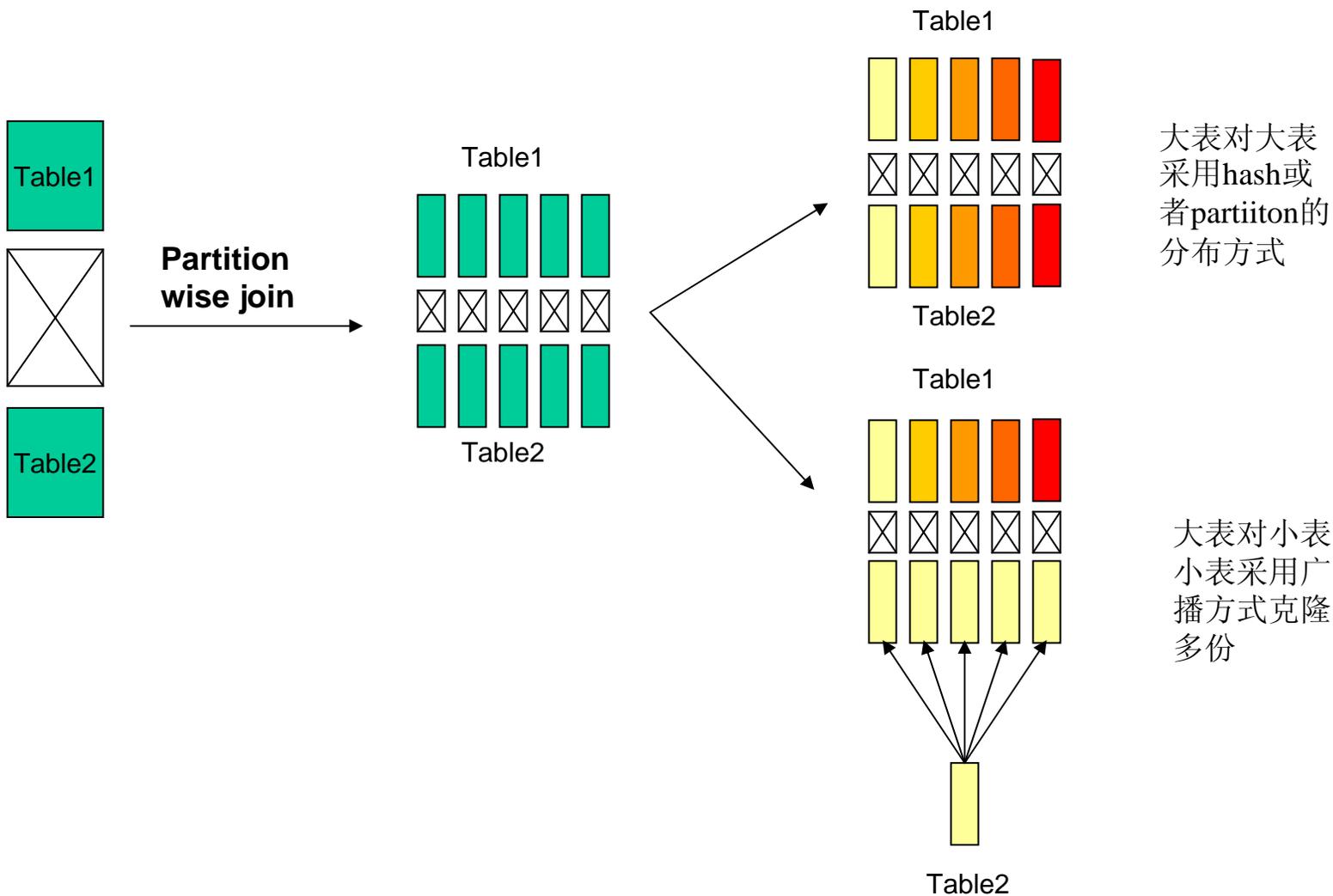
➤ 并行度参数，如

- `parallel_max_servers=256`
- `parallel_min_servers=16`
- `parallel_threads_per_cpu=2`
- `parallel_execution_message_size=16384`

➤ 跨节点并行的参数

- `instance_groups=crm,erp,oltp`
- `parallel_instance_group=crm`
- 关于跨节点并行的更多资料，可以参考我们团队的Blog：
：http://rdc.taobao.com/blog/dba/html/217_rac-span-paralle.html

并行执行的Join方式



案例分析（一）——升级中的问题

➤ 08年数据仓库从10.2.03升级到10.2.04

- 在升级之后碰到了一个新的BUG，数据库会报如下600错误：
 - ORA-00600: internal error code, arguments:
[kddummy_blkchk], [47], [935468], [18038], [], [], [], []
- 数据库部分查询报块损坏，无法执行

升级中的问题—解决办法

- 临时解决办法是将db_block_checksum这个参数改为FALSE
- 后面ORACLE新发布了一个补丁5386204，打上补丁后问题得以解决。

案例分析（二）——节点异常退出

➤ 08年12节点刚开始使用cisco的infiniband交换机，总是出现如下的情况：

- [CSSD]2008-09-02 17:48:15.706 [1241577824] >WARNING: clssnmPollingThread: node dw6 (6) at 50 3.118329e-317artbeat fatal, eviction in 29.490 seconds
.....
- [CSSD]2008-09-02 17:49:15.677 [1241577824] >WARNING: clssnmPollingThread: node dw5 (5) at 50 3.118519e-317artbeat fatal, eviction in 29.780 seconds
.....
- [CSSD]2008-09-02 17:49:46.460 [1262557536] >ERROR: clssnmCheckDskInfo: Aborting local node to avoid splitbrain.
[CSSD]2008-09-02 17:49:46.460 [1262557536] >ERROR: : my node(11), Leader(1), Size(10)
VS Node(6), Leader(1), Size(12)

- 节点6与其它节点内部通信超时，甚至不通，而自动重启。
 -
- 节点6重启之前，节点5也出现重启
- 在节点6重启之前，有更多的节点被踢出（重启），直到开始正式重启

节点异常退出—解决办法

➤ 原因分析

- 怀疑cisco的infiniband在高压下不稳定，网络出现延迟。
- 当一个节点负载压力太高，被迫重启的时候，与其它节点内部网络出现超时甚至不通
- 在该节点重启之前，其它节点认为该节点还正常（Voting Disk是正常的）
- 其它节点会出现脑裂，也会不同程度的把自己从RAC环境中踢出

➤ 解决办法

- 将infiniband网卡替换成普通网卡，内部互联就恢复正常，再也没有节点自动宕机的情况了

未来的计划

- 扩展到20节点
- Oracle 11g的新特性
- 继续尝试使用Infiniband交互技术
- Oracle新硬件Exadata的测试

谢谢

➤ Q&A