



# Enterprise-Ready Cluster Solution

Peter Cheng | 程勇  
Open Solution Alliance Manager

# 议程

- ➔ **企业面临的高可靠性挑战**
- ➔ 负载均衡解决方案
- ➔ JVM 集群解决方案
- ➔ 数据库集群解决方案

# 当前企业高可用性集群的挑战

## 最终使用者

“... 希望更专注于设计的本职工作少做或者不做 IT 工作”

- 高性能计算应用基于开放平台的提供与集成
- 简单易用的计算作业提交与监控

## IT 系统管理人员

“... 部署管理一个高性能计算机群实在太困难了”

- 更加简单的部署与建立计算机集群
- 简单统一的集群管理环境并与现有的 IT 架构一致

## 应用开发者

“... 编程很困难  
... 没有足够好用的开发工具与环境”

- 零编程的集群配置模式
- 无需额外的二次开发

# 全面的集群解决方案

服务器负载均衡

应用服务器集群

数据库集群



Load Director



Terracotta



Continuent

系统

用户群

政府、研究机构、大型商业企业、互联网企业

典型客户

[www.abc.com](http://www.abc.com) , [www.real.com](http://www.real.com) , [www.eds.com](http://www.eds.com) , [www.ctrip.com](http://www.ctrip.com)

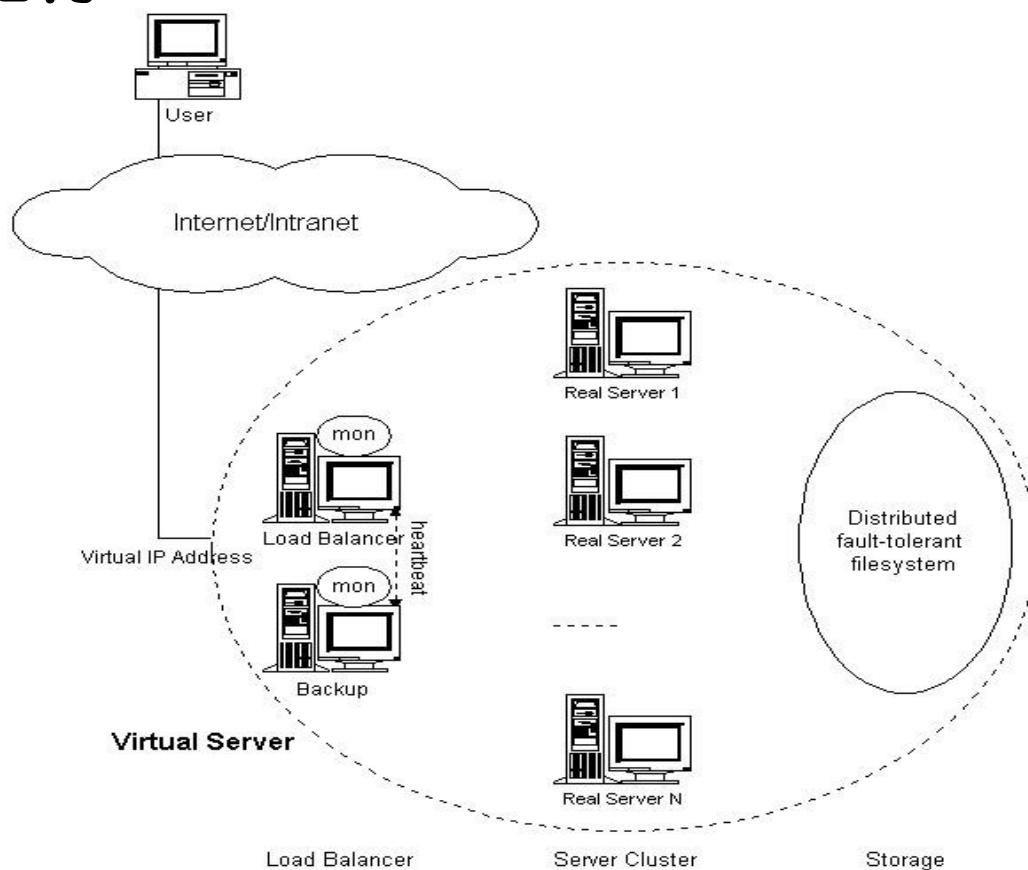
## 议程

- 企业面临的高可靠性挑战
- **负载均衡解决方案**
- JVM 集群解决方案
- 数据库集群解决方案
- Linux 操作系统高可靠性方案

## 核心技术：集群技术

### 虚拟服务器体系结构

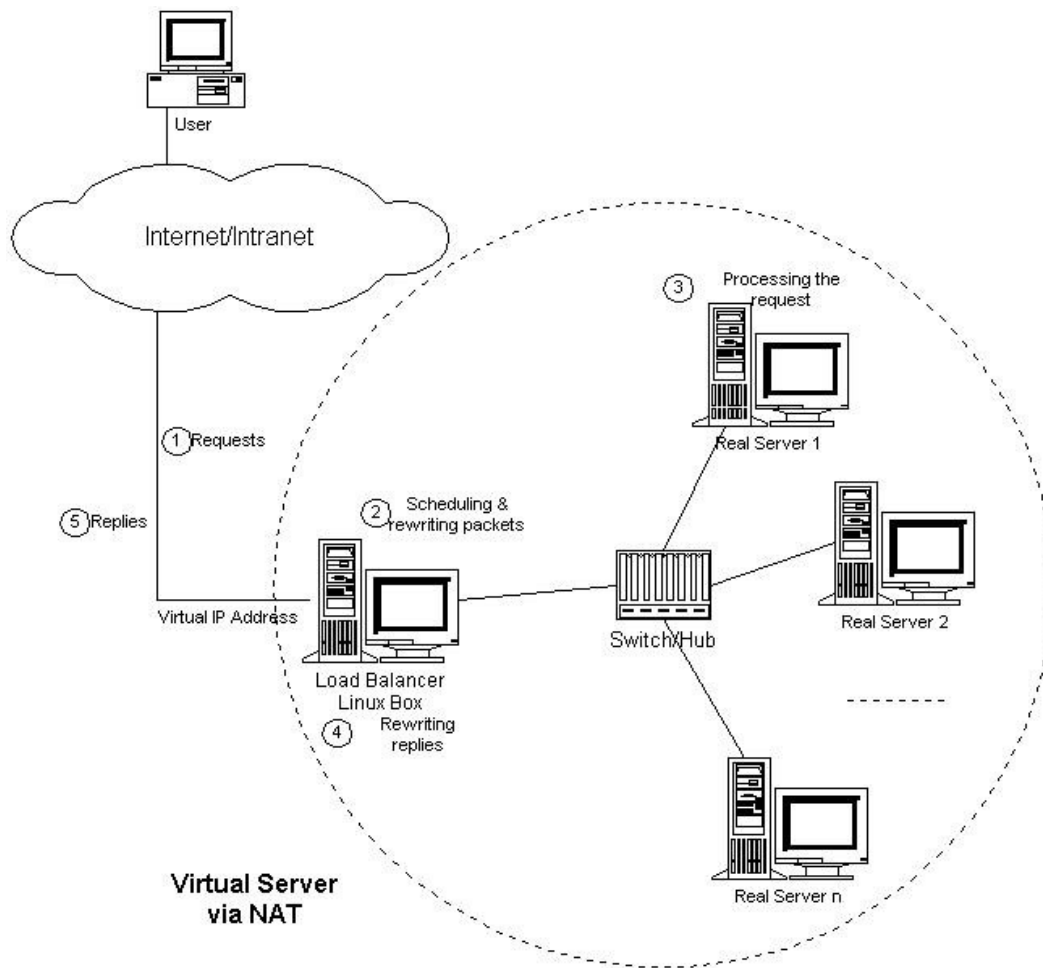
- ➔ 负载调度器
- ➔ 服务器池
- ➔ 后端存储



## LVS 集群技术-技术简介 (1/3)

### IP 负载均衡技术 —Virtual Server via NAT

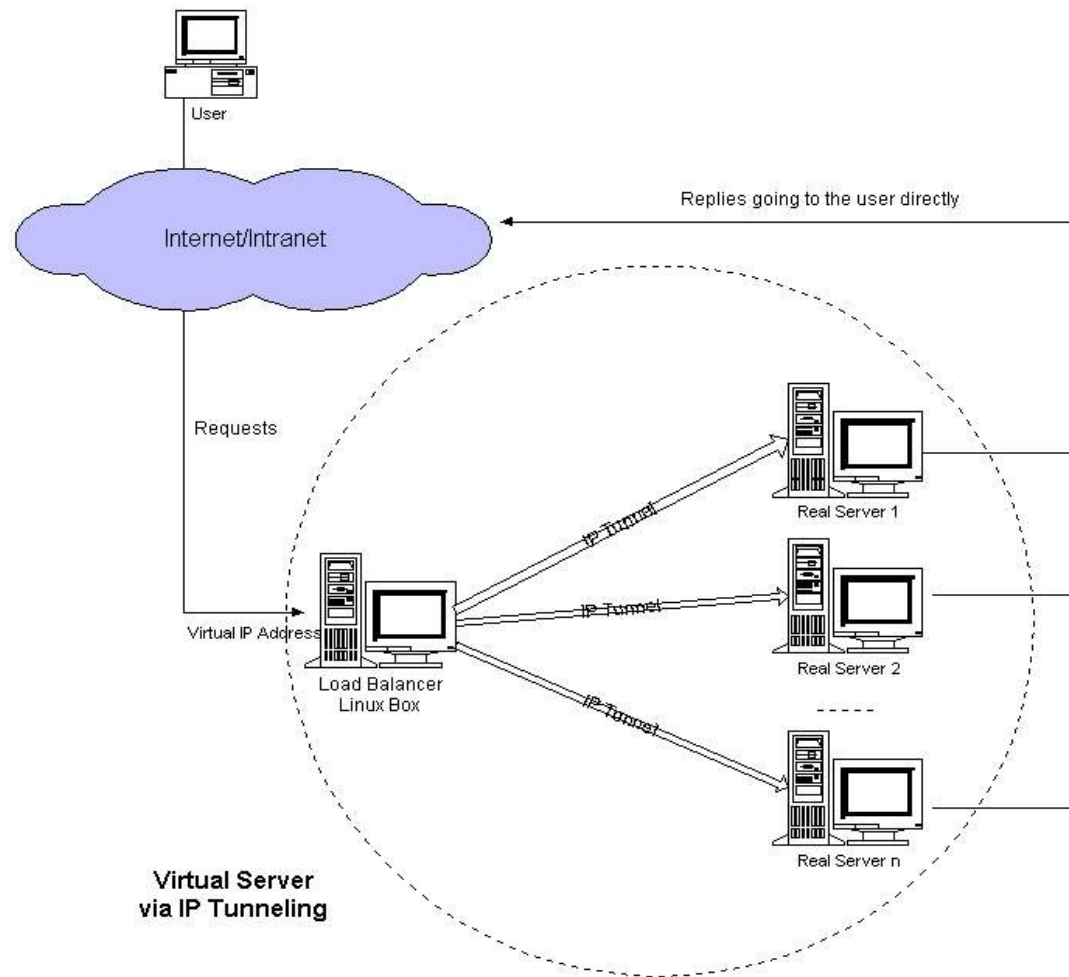
- ➔ 体系结构
- ➔ IP 隧道
- ➔ 工作流程
- ➔ 特 点



## LVS 集群技术-技术简介 (2/3)

### IP 负载均衡技术 —Virtual Server via IP Tunneling

- 体系结构
- IP 隧道
- 工作流程
- 特 点

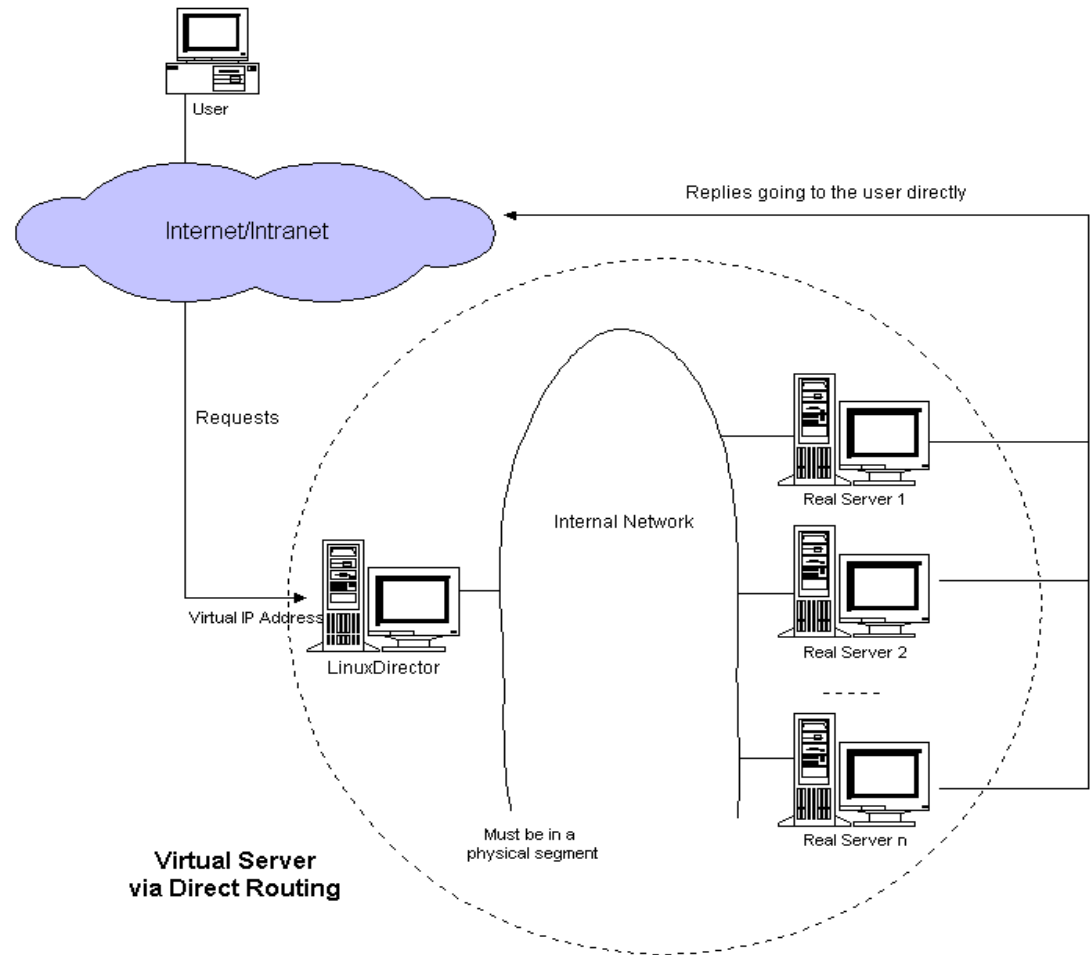




## LVS 集群技术-技术简介 (3/3)

### IP 负载均衡技术 —Virtual Server via Direct Routing

- ➔ 体系结构
- ➔ 原理
- ➔ 工作流程
- ➔ 特 点



## LVS 集群技术-应用实例

- ABC.com (American Broadcasting Corporation)
- www.real.com
- empas.com
- 英国国家 JANET Cache 网
- Dell
- EDS.com
- sourceforge.net & linux.com
- Sina.com
- QQ.com

### 集群技术 – 用户评价

“We tried virtually all of the commercial load balancers, LVS beats them all for reliability, cost, manageability, you-name-it”  
Jerry Glomph Black, Director, Internet & Technical Operations, Real Networks, Seattle Washington, USA

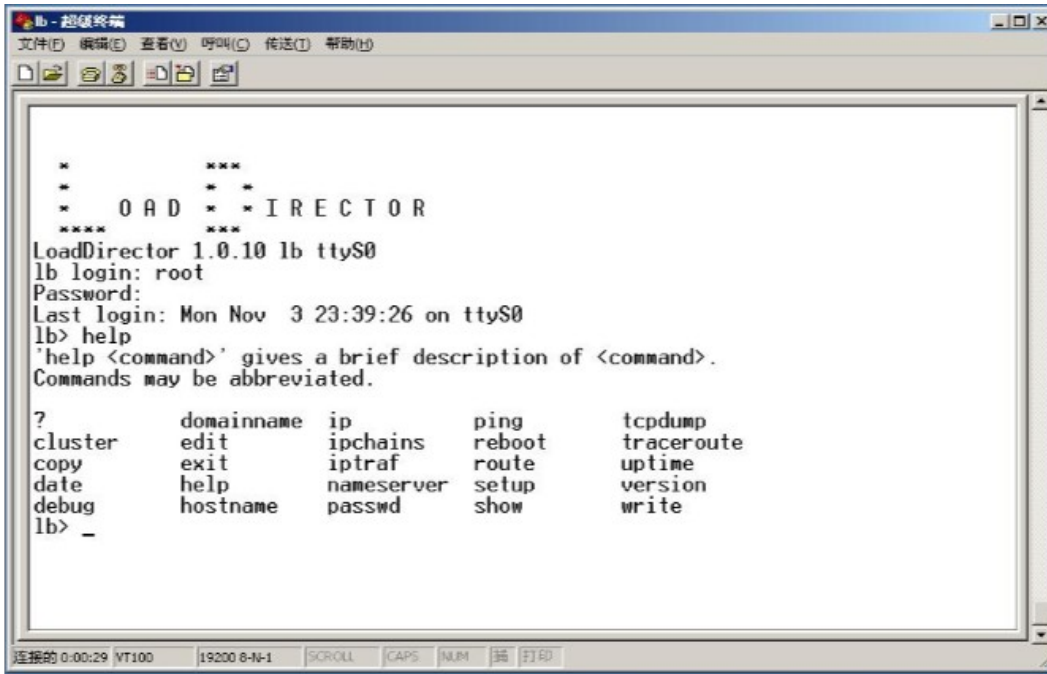
<http://marc.theaimsgroup.com/?1=linux-virtual-server&m=95385809030794&w=2>

“I can say without a doubt that lvs toasts F5/BigIP solutions, at least in our real world implementations. I wouldn't trade a good lvs box for a Cisco LocalDirector either”

Drew Streib, Information Architect, VA Linux Systems, USA

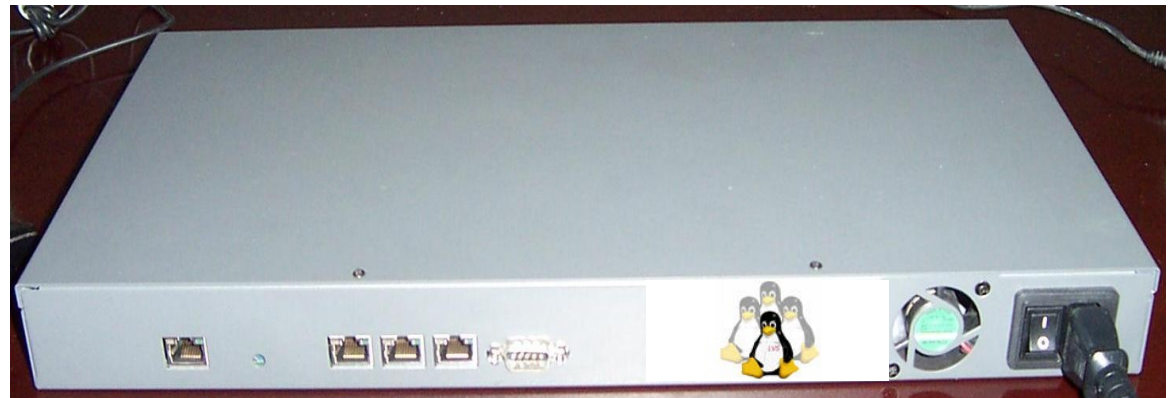
<http://marc.theaimsgroup.com/?1=linux-virtual-server&m=95385694529750&w=2>

## Load Director CLI 配置界面

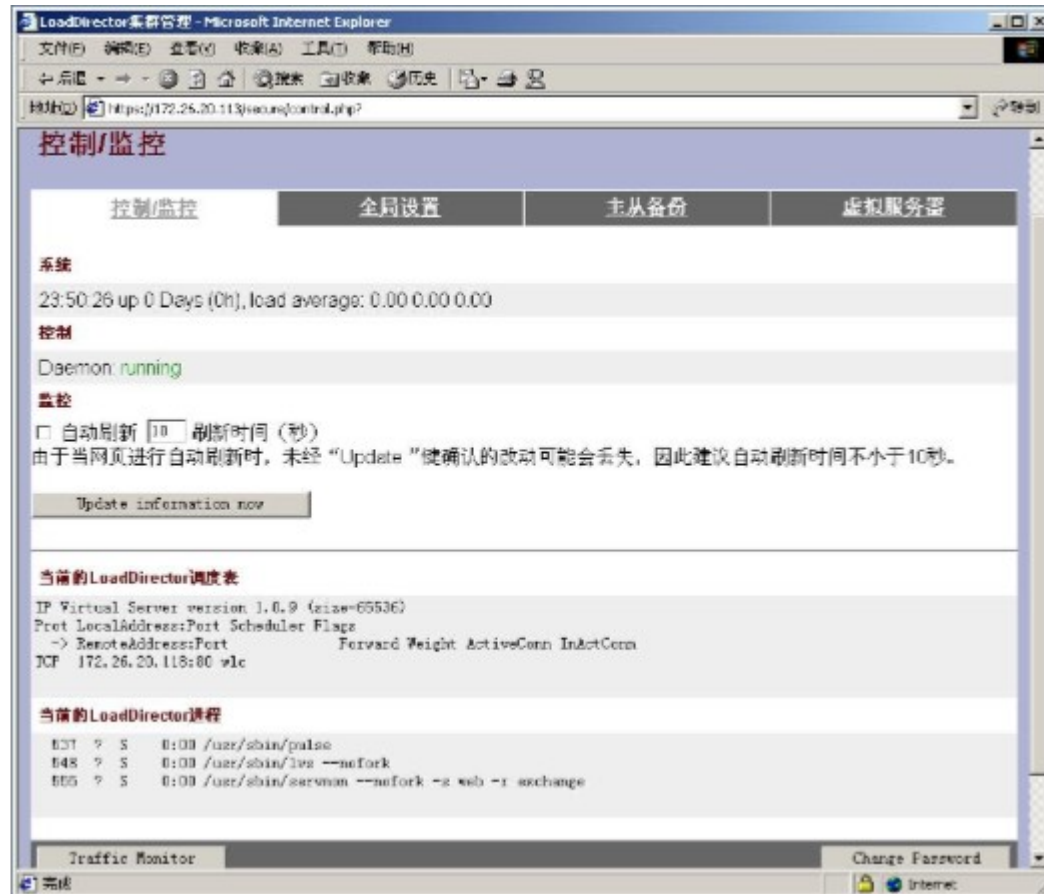


```
lb - 超级终端
文件(F) 编辑(E) 查看(V) 呼叫(O) 传送(T) 帮助(H)
*      *
*      *
*  O A D  *  *  I R E C T O R
*      *
*      *
LoadDirector 1.0.10 lb ttyS0
lb login: root
Password:
Last login: Mon Nov  3 23:39:26 on ttyS0
lb> help
'help <command>' gives a brief description of <command>.
Commands may be abbreviated.

?      domainname  ip          ping        tcpdump
cluster edit        ipchains   reboot      traceroute
copy   exit        iptraf     route       uptime
date   help       nameserver setup        version
debug  hostname   passwd     show        write
lb> _
```



## Load Director Web 管理监控界面



## Load Director 集群技术-应用实例

- Load Director 用户
  - ➔ 携程网 ctrip.com
  - ➔ 浩方在线
  - ➔ Ku6.com
  - ➔ 世纪互联
  - ➔ 中国移动湖南分公司
  - ➔ 中国联通广东分公司
  - ➔ 广州白云机场
  - ➔ 湖南省信息中心
  - ➔ 晓通网络 www.xiaotong.com.cn
  - ➔ 万维易化 www.ezcross.com
  - ➔ 凯思昊鹏 www.hopen.com.cn
  - ➔ 鸿联九五 www.hl95.com
  - ➔ 河北邯郸信息网 [www.hd.gov.cn](http://www.hd.gov.cn)
  - ➔ 轻点万维

## Load Director 竞争性分析

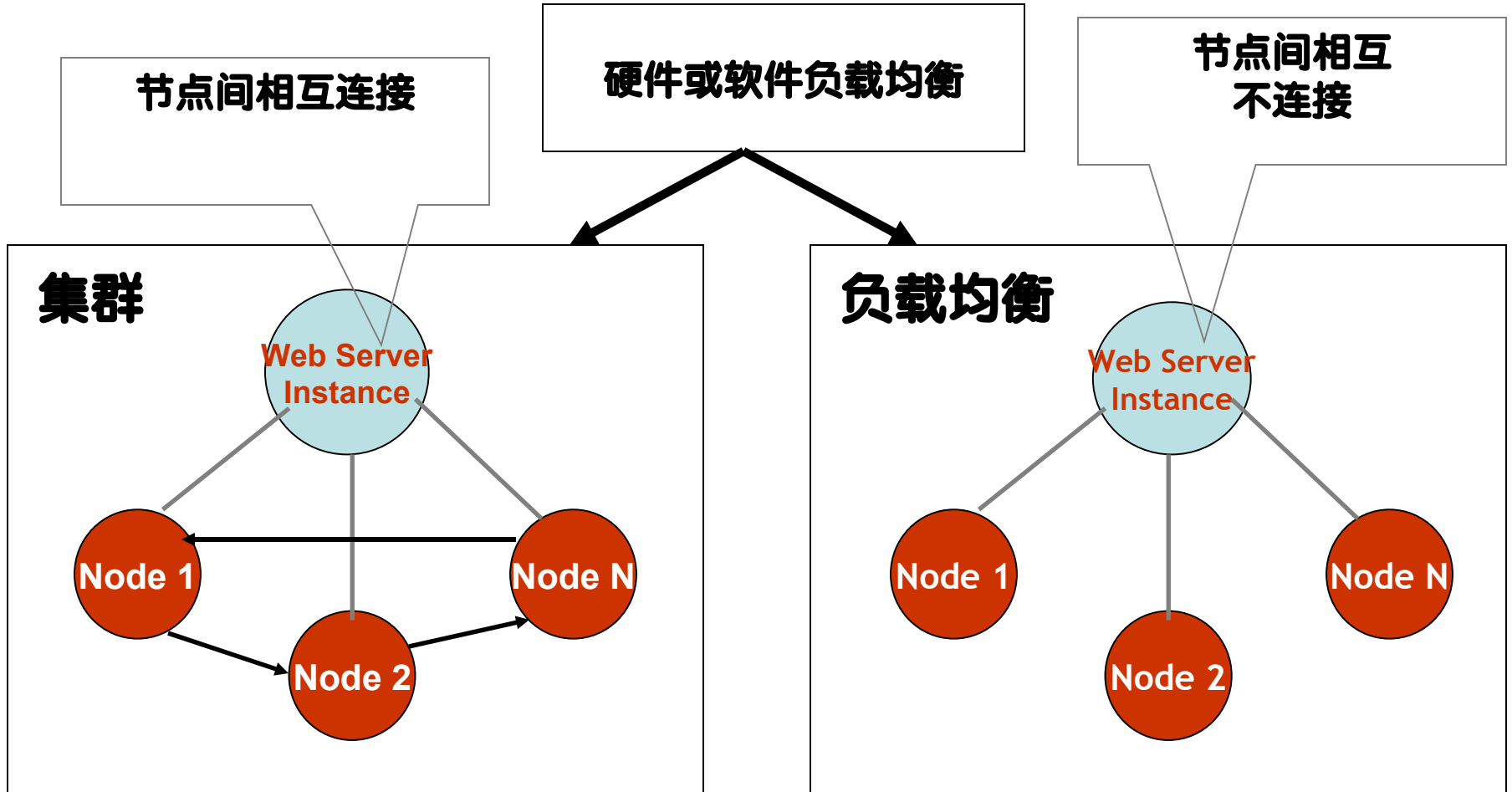
	Load Director 2000	F5 BIG/IP 2U
<b>操作系统</b>	<ul style="list-style-type: none"><li>支持任何 TCP/IP 操作系统包括 :Linux, Solaris, AIX, HP UX, Mac/OS, Windows</li></ul>	<ul style="list-style-type: none"><li>支持 WIN NT 各种 Unix 平台和 Mac/OS</li></ul>
<b>内存</b>	<ul style="list-style-type: none"><li>512M, 可扩展到 2G</li></ul>	<ul style="list-style-type: none"><li>512M 可扩展到 1G</li></ul>
<b>网络端口 负载均衡技术</b>	<ul style="list-style-type: none"><li>Intel 10/100/1000 以太网口 3~5 个</li><li>网络地址转换 (NAT)</li><li>直接路由 (Direct Routing)</li><li>IP 隧道 (IP Tunneling)</li></ul>	<ul style="list-style-type: none"><li>Intel 10/100/1000 以太网口 3~5 个</li><li>网络地址转换 (NAT)</li></ul>
<b>调度算法</b>	<ul style="list-style-type: none"><li>轮叫和加权轮叫</li><li>最小连接和加权最小连接</li><li>基于局部性最小连接调度</li><li>带复制的基于局部性最小连接调度</li></ul>	<ul style="list-style-type: none"><li>轮叫</li><li>最小连接</li></ul>
<b>最大系统 吞吐量</b>	<ul style="list-style-type: none"><li>700Mbps(NAT)</li><li>9Gbps(DR, TUN)</li></ul>	<ul style="list-style-type: none"><li>700Mbps(NAT)</li></ul>
<b>最大可同 时连接数</b>	<ul style="list-style-type: none"><li>500 万</li></ul>	<ul style="list-style-type: none"><li>500 万</li></ul>
<b>价格</b>	<ul style="list-style-type: none"><li>LD2000      <b>¥ 119,500</b></li><li>LD2000S    <b>¥ 189,500</b></li></ul>	<ul style="list-style-type: none"><li>F5 BIG-IP-3400    <b>¥ 370000</b></li><li>F5 BIG-IP-1500    <b>¥ 230000</b></li></ul>

## 议程

- 企业面临的高可靠性挑战
- 负载均衡解决方案
- **应用服务器集群解决方案**
- 数据库集群解决方案

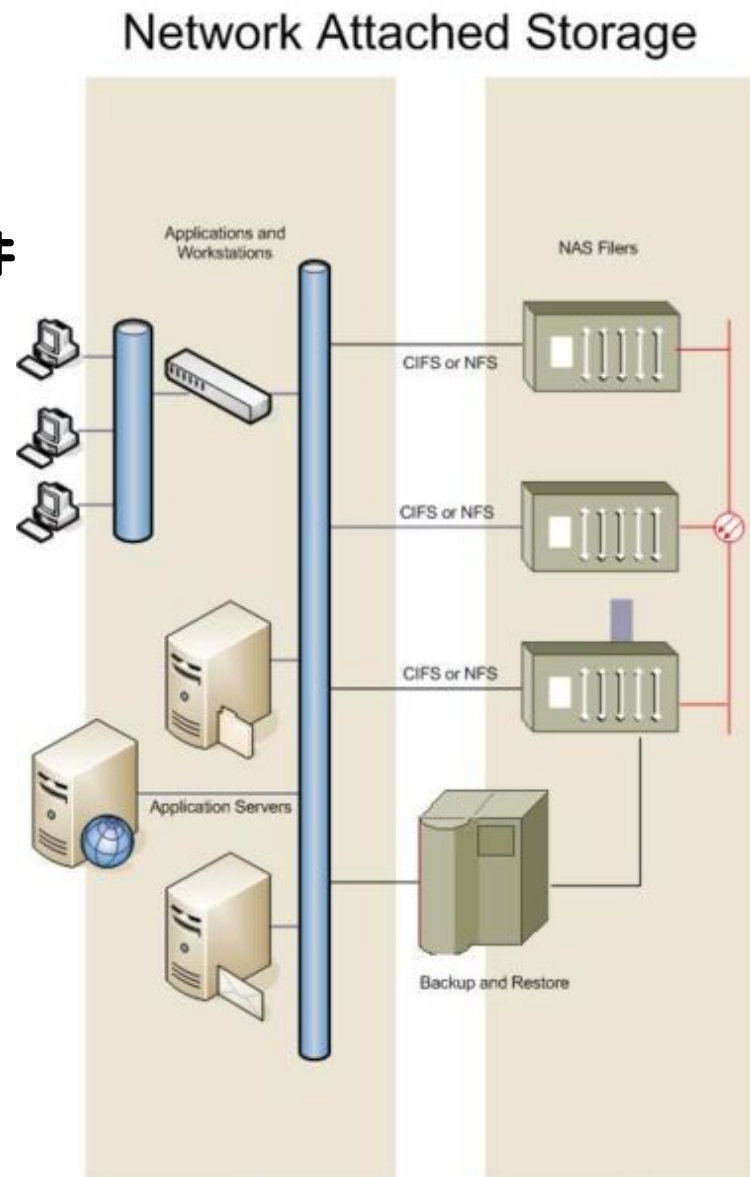


## 负载均衡与集群的节点

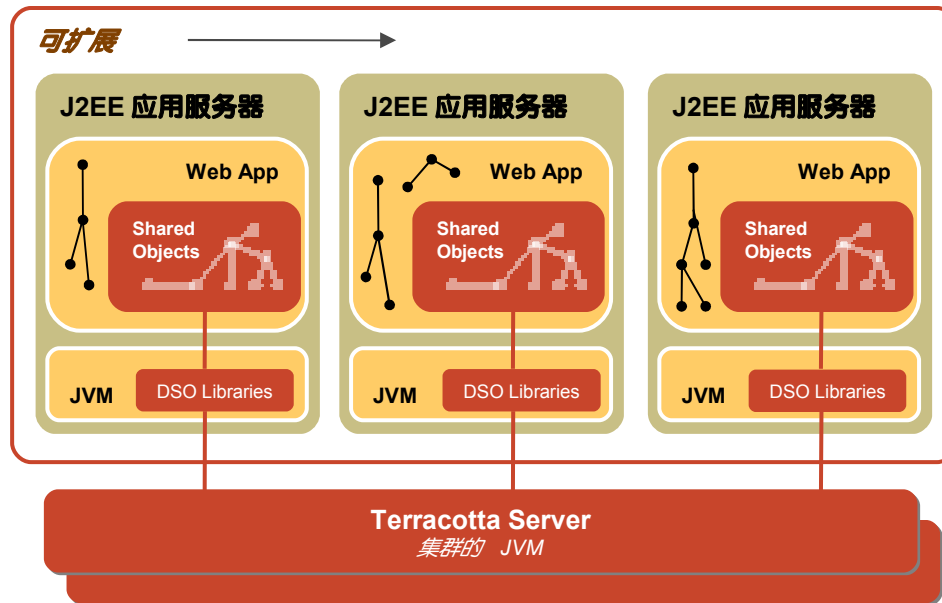


## 什么是集群的 JVM?

- 多个 JVM 如何解决内存共享的问题?
- JVM 内存共享比较满足一下两个基本条件
  - 1. 对于应用来说应该像 RAM 一样
  - 2. 必须作为一个基础设施服务



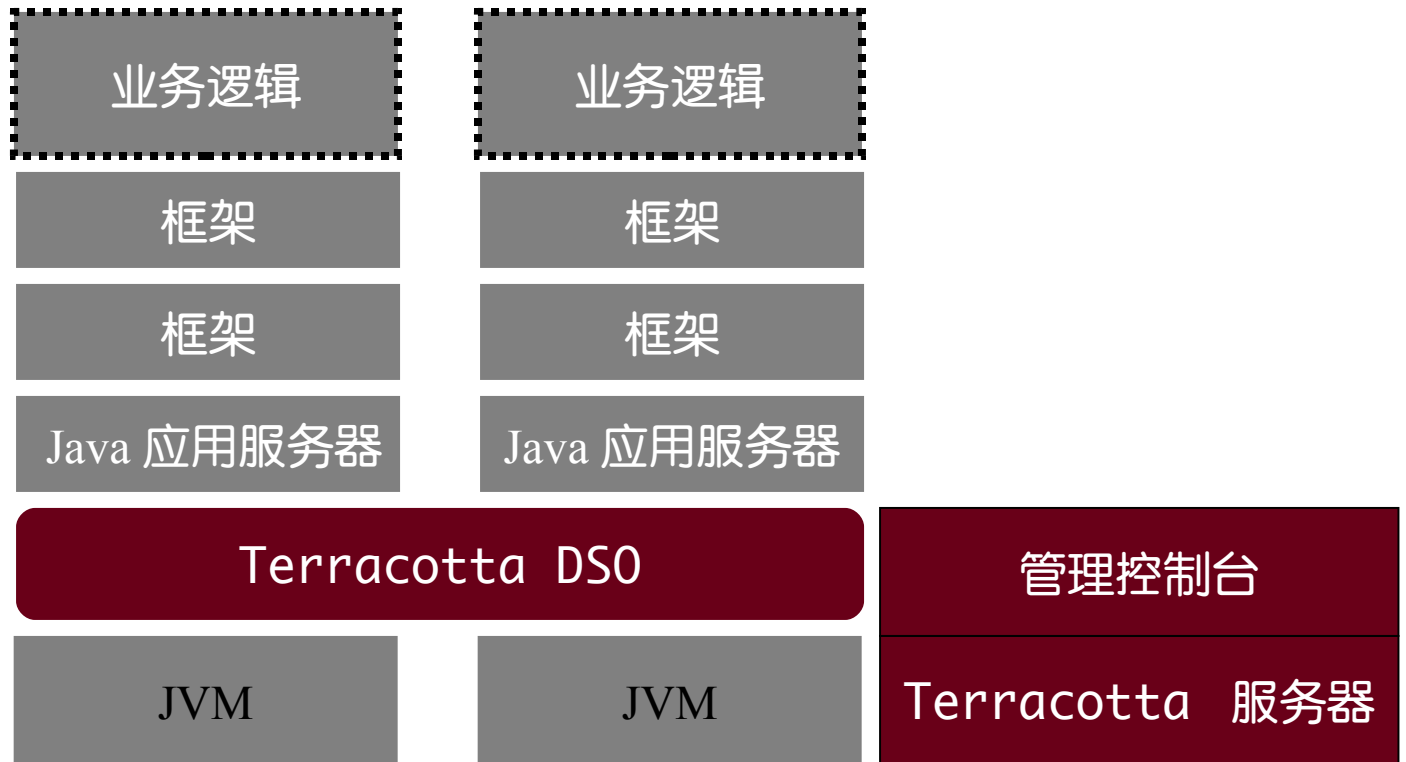
## Terracotta JVM 集群架构：集群你的 JavaEE 应用资源



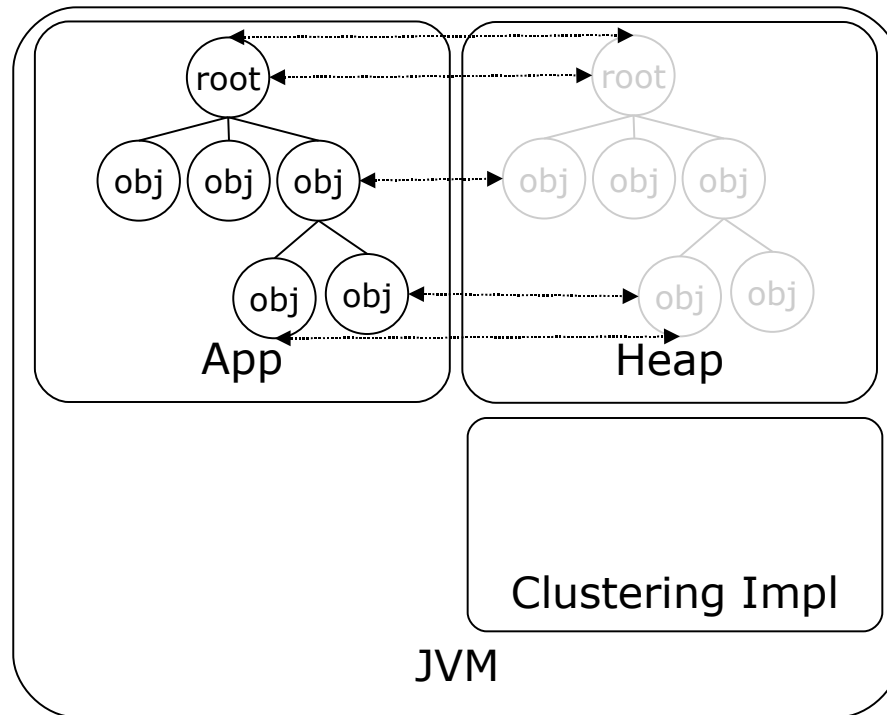
### 特性

- 内存堆级别的复制
- 中心化存储
- 虚拟内存
- 可监控和管理

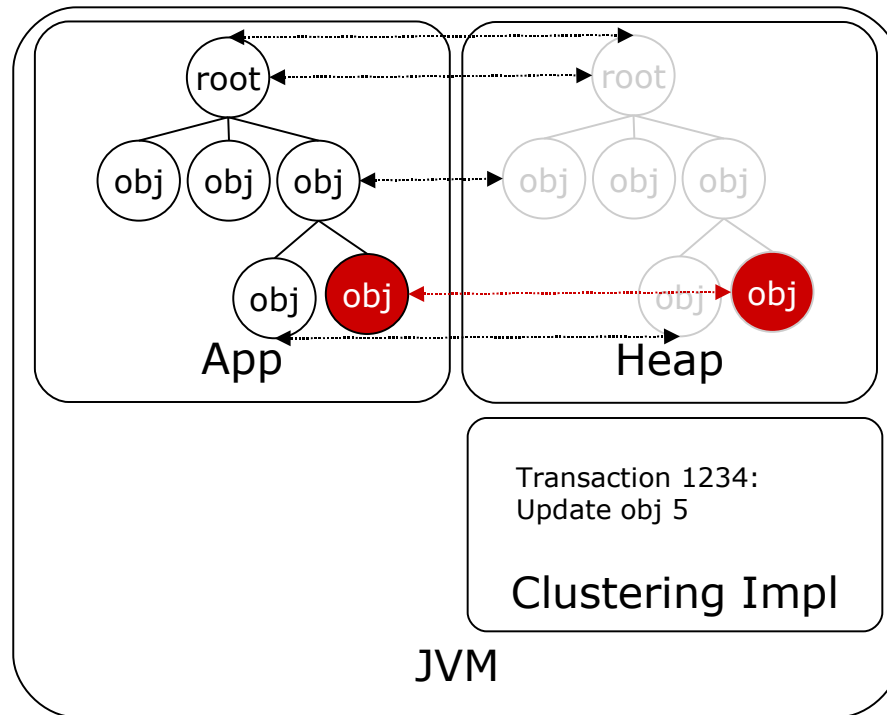
## Terracotta 如何实现?



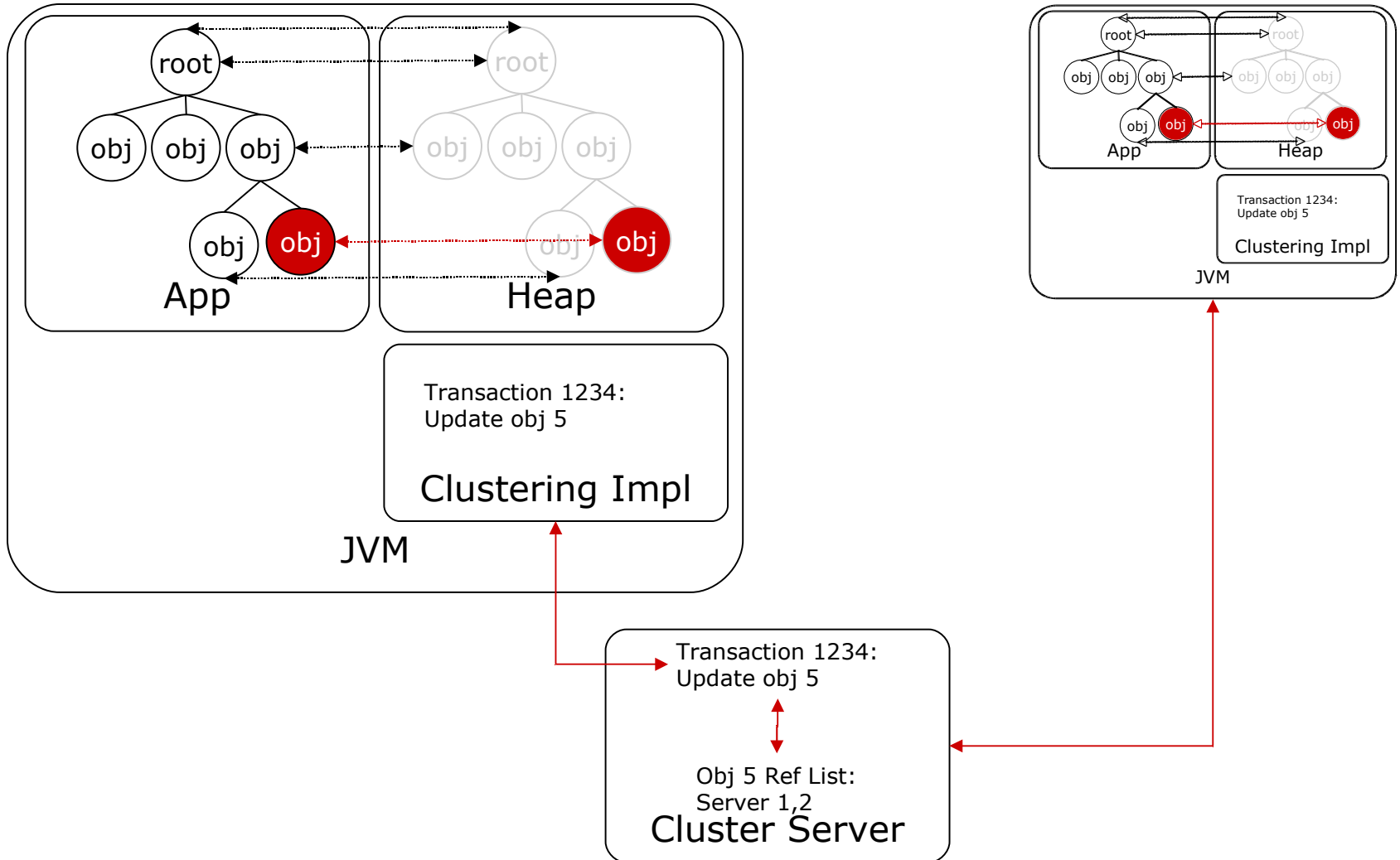
## Terracotta 工作流程 (1/3)



## Terracotta 工作流程 (2/3)



## Terracotta 工作流程 (3/3)



## 配置级管理，零编程 DEMO

### Shared JTable (spreadsheet)

```
<terracotta-config>
  <dso>
    <server-host>localhost</server-host>
    <server-port>9510</server-port>
    <dso-client>

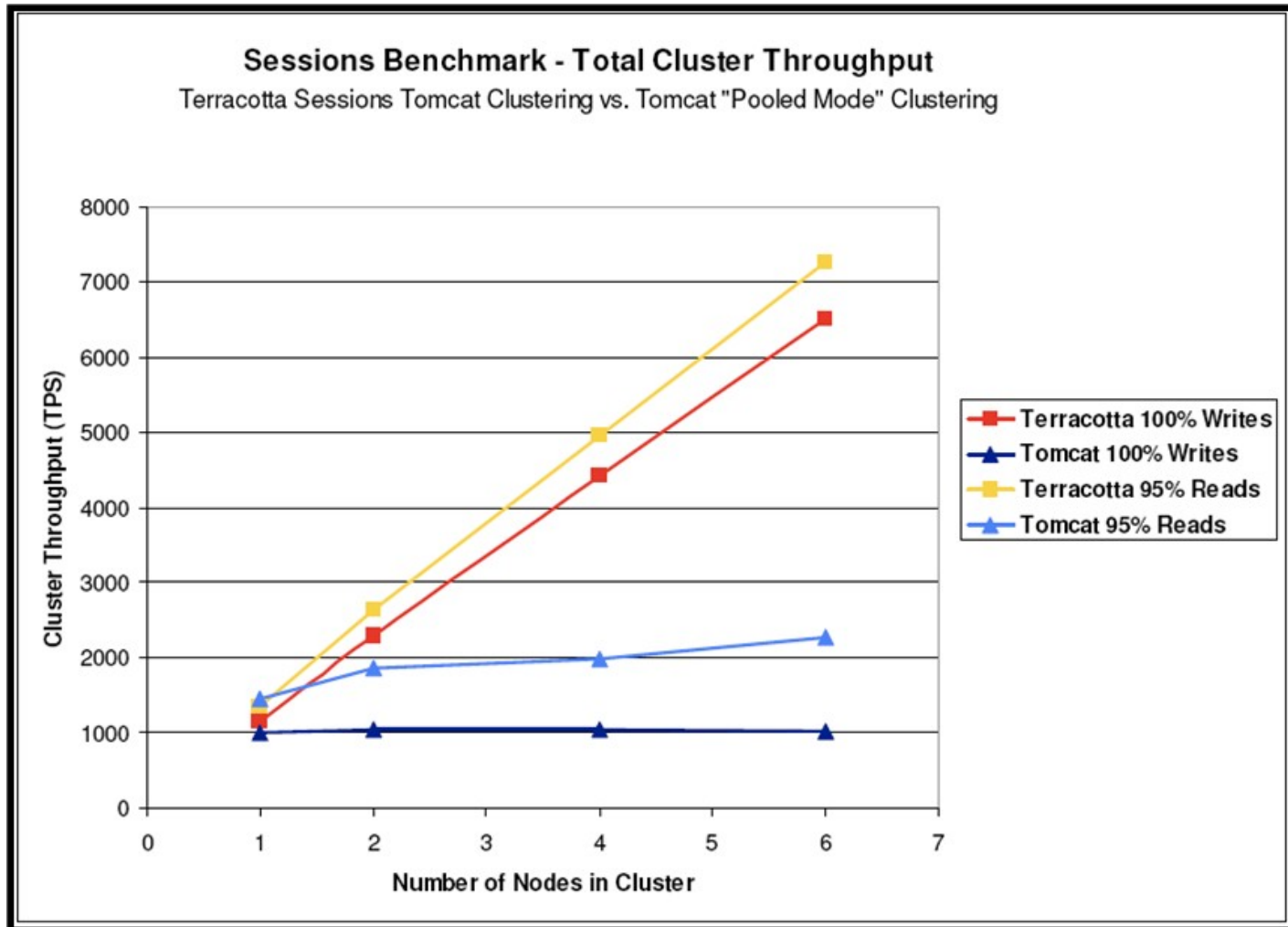
      <roots>
        <root>
          <field-name>demo.jtable.TableDemo.model</field-name>
        </root>
      </roots>

      <included-classes>
        <include><class-expression>demo..*</class-expression></include>
      </included-classes>

    </dso-client>
  </dso>
</terracotta-config>
```



## Tomcat Http Session Cluster



## 竞争性分析 : Terracotta vs. Coherence

- ➔ **Terracotta: 分布式的 Java 运行环境**
  - 易于集成
  - 可视的中心化操作管理
  - 企业级的开放源码
  
- ➔ **Coherence: 开发者框架**
  - 入侵式的 API
  - 有限的可操作控制 – 非运行环境
  - Oracle 式的许可证

## Terracotta 性能测试

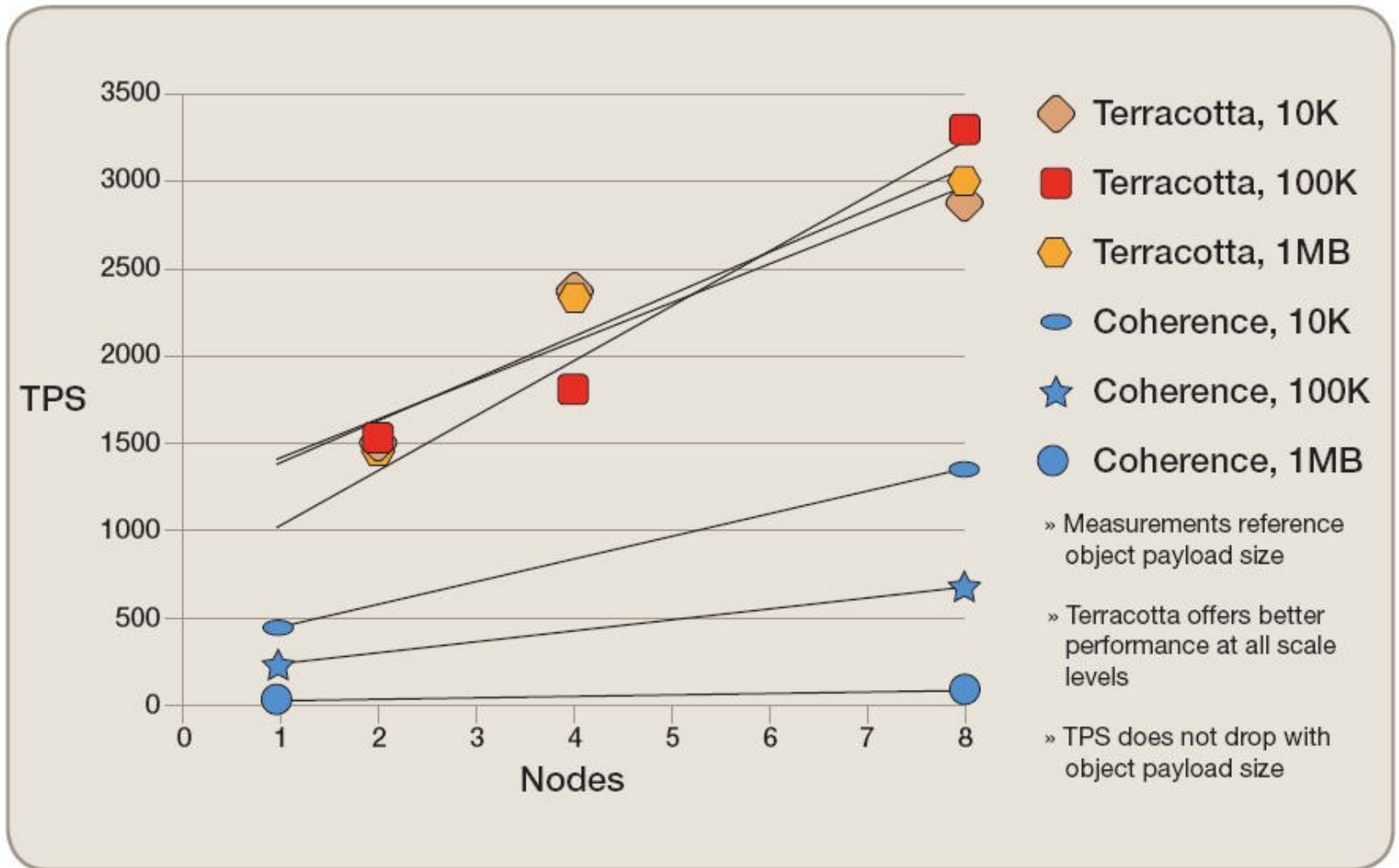
- 测试环境
  - 文件系统： SAN/Linux ext3
  - JVM 内存堆大小： 4GB
  - 硬件： DL380
  - 存储类型： 内部磁盘 40G/HDB 335G Ext3
  - 内存： 12GB RAM
  - 网络： GigE
- 测试方法
  - 每秒交易数
  - CPU 利用率
  - 网络利用率
  - 内存特征

## Terracotta 性能测试

### → 评估标准

- 高性能：主机上服务能够处理交易的能力 (10,000TPS)
- 易于使用：技术应该使开发成本低廉并且保持简单。技术应该从应用中抽象分离出基础设施层所关注的部分 ( 集群, 可扩展, HA)
- 高可扩展：服务应该能够满足业务增长的需求而扩展, 并将运营和开发的成本降到最低。在相同的性能和 HA 下, 离线主机系统能够处理 20 倍主机 (1MB) 系统容量的负载
- 高可靠性：服务应该满足跨企业的 99.999% 的 SLA
- 运营监控与管理：技术应该能够提供丰富的调试工具集, 监控并管理这些集群服务

## Terracotta vs Coherence 性能



## 测试结果关键点

- Terracotta 性能是 Oracle Coherence 的 6 倍
- Terracotta 实现应用需要 3 天而不是 6 周
- Terracotta 展现出卓越的运营监控和管理能力

## Terracotta 竞争性分析

	Terracotta	Oracle Coherence
简单	<ul style="list-style-type: none"><li>• 选择适当的开发模式</li><li>• 选择最佳的框架和容器</li><li>• 自动判断工作负载</li></ul>	<ul style="list-style-type: none"><li>• 入侵式 API 必须显示编码</li><li>• 仅适用于数据缓存模式</li><li>• 缓存策略必须显示的配置</li></ul>
可扩展	<ul style="list-style-type: none"><li>• Delta 复制, 精密的数据路由 = 高性能</li><li>• 有效的网络复制</li><li>• 在应用恶化服务器级别扩展</li></ul>	<ul style="list-style-type: none"><li>• 交易随可靠性和数据一致性规模化</li><li>• 复制是网络和 CPU 敏感的</li></ul>
高可用	<ul style="list-style-type: none"><li>• 没有单点故障</li><li>• 数据在磁盘上是持久化的</li></ul>	<ul style="list-style-type: none"><li>• 数据只在内存中复制</li><li>• 没有数据在哪存活概念</li></ul>
数据一致性	<ul style="list-style-type: none"><li>• 直接的共享状态模型</li><li>• 有保障的连贯更新</li><li>• ACID 兼容</li></ul>	<ul style="list-style-type: none"><li>• 高风险的数据失误</li><li>• 高扩展的交易一致性更新</li><li>• Get/Put 模式易损坏和出错</li></ul>
操作管理	<ul style="list-style-type: none"><li>• 服务器基础设施</li><li>• 特征鲜明, 操作直接的控制台</li><li>• 开发与运营拥有公共的接口</li><li>• 持久化的数据意味着简单的提供 / 销毁</li></ul>	<ul style="list-style-type: none"><li>• 非运行时环境</li><li>• 有限的可视化与控制</li><li>• 没有数据在哪存活的概念</li><li>• 提供 / 销毁会带来数据损害</li></ul>
价格	<ul style="list-style-type: none"><li>• <b>\$ 10000/JVM</b></li></ul>	<ul style="list-style-type: none"><li>• <b>\$ 30000/CPU</b></li></ul>

# Enterprise-Ready Cluster Solution

## 企业级用户



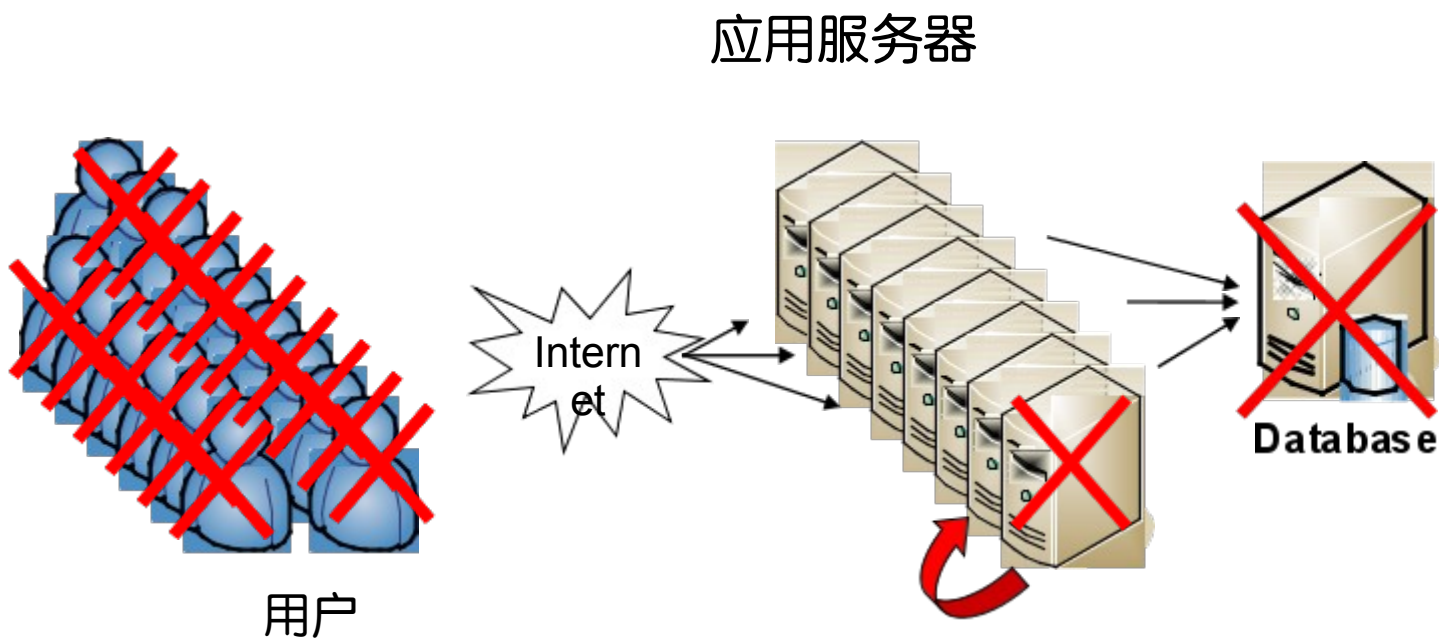


## 议程

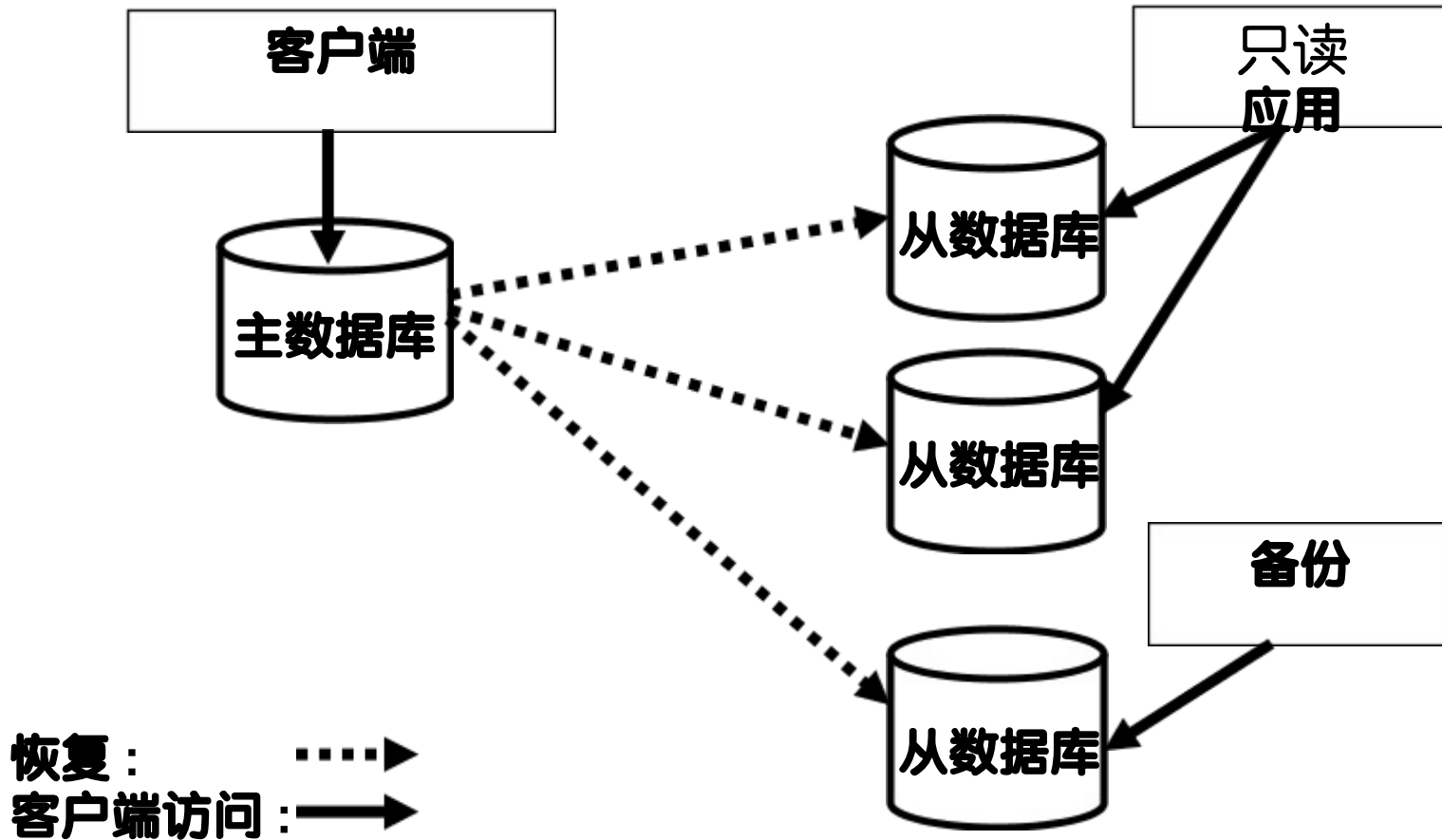
- 企业面临的高可靠性挑战
- 负载均衡解决方案
- 应用服务器集群解决方案
- **数据库集群解决方案**

## DB 的高可用性问题

- 客户端连接到应用服务器
- 应用都基于数据库来构建
- 大量的工具解决应用层的可用性问题
- **数据库仍然有单点故障**



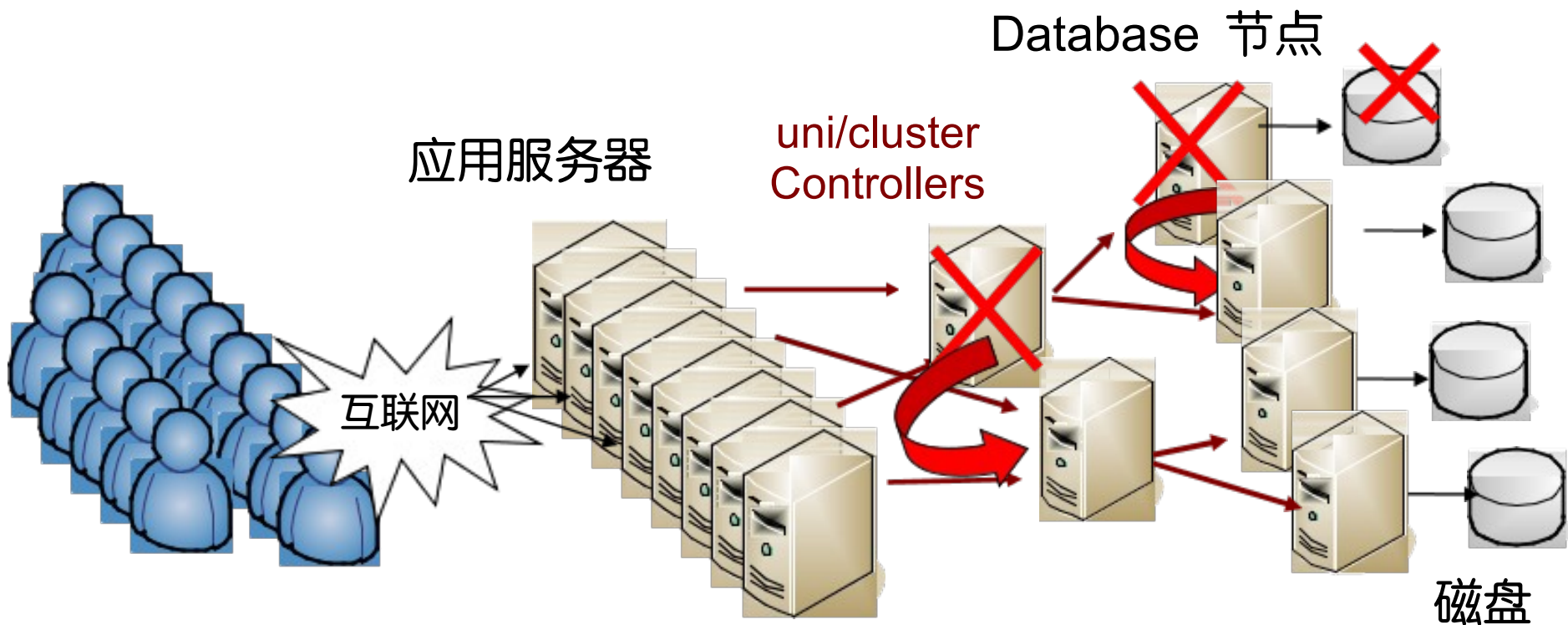
## DBMS 集群常见解决方案 例如：MySQL



## DBMS 集群解决方案 – 限制

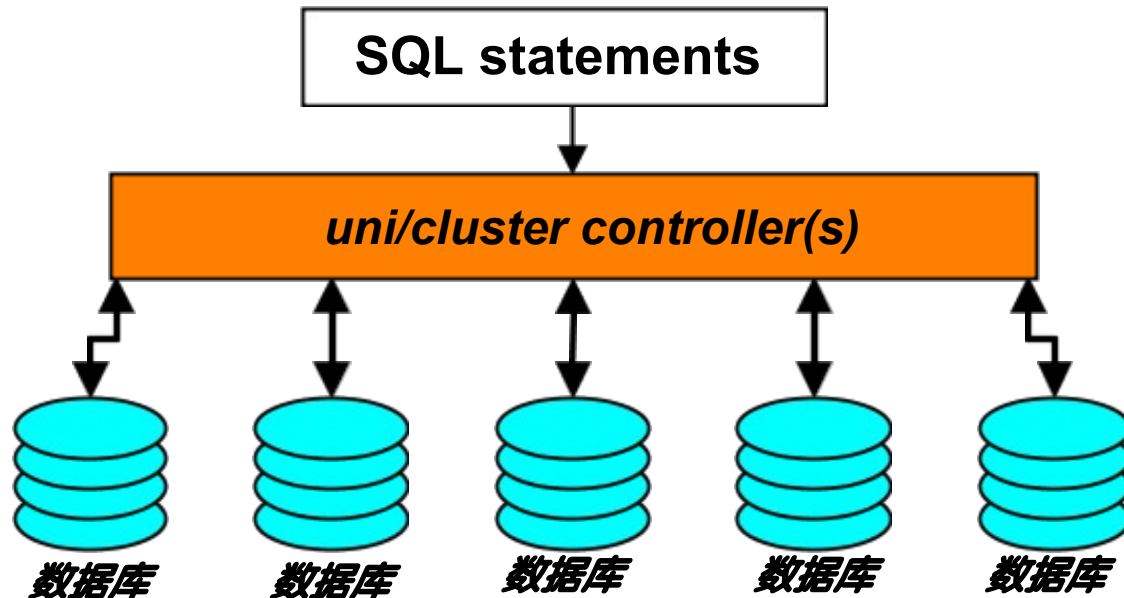
- 主数据仍然会是单点故障
- 故障恢复不是自动的，并不能透明的完成
- 异步更新 – 从数据库总是在后端
- 配置和管理非常复杂
- 没有一致的复制管理方式
  - MySQL vs. PostgreSQL

# Uni/Cluster 数据库解决方案

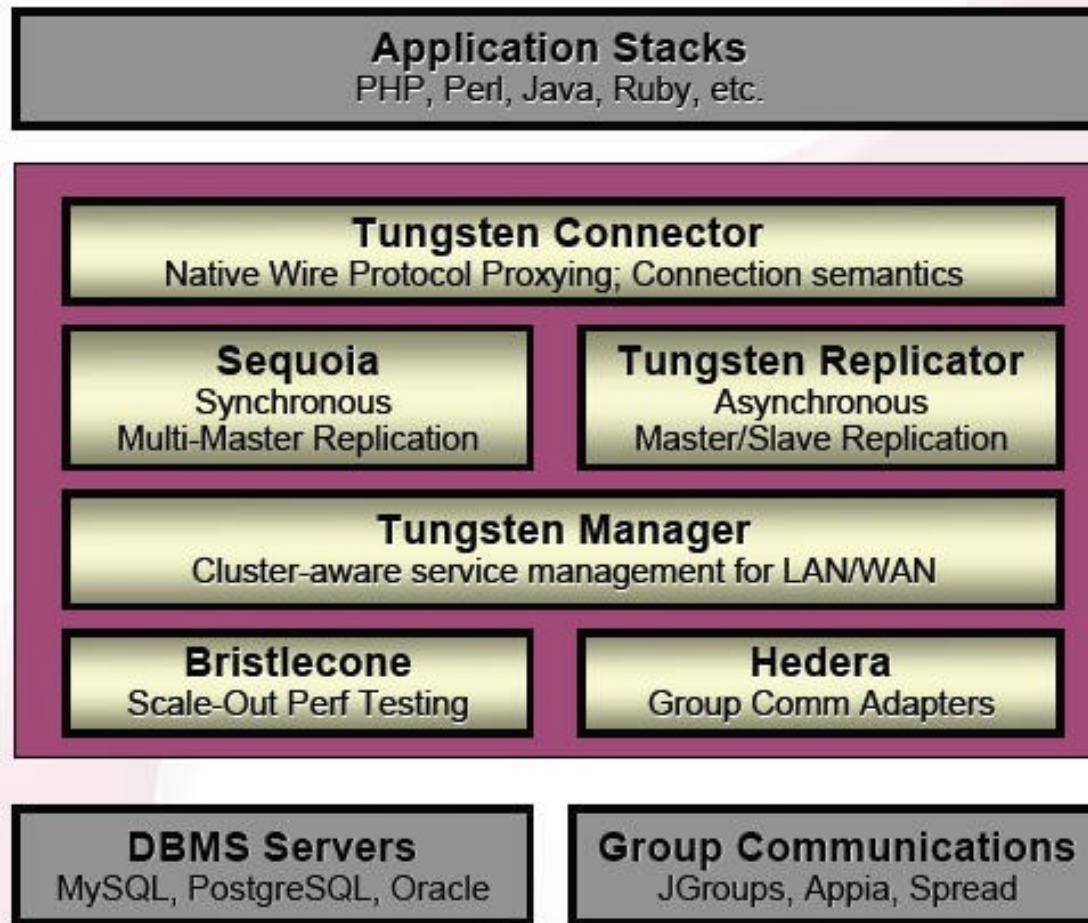


## 冗余的数据库集群架构

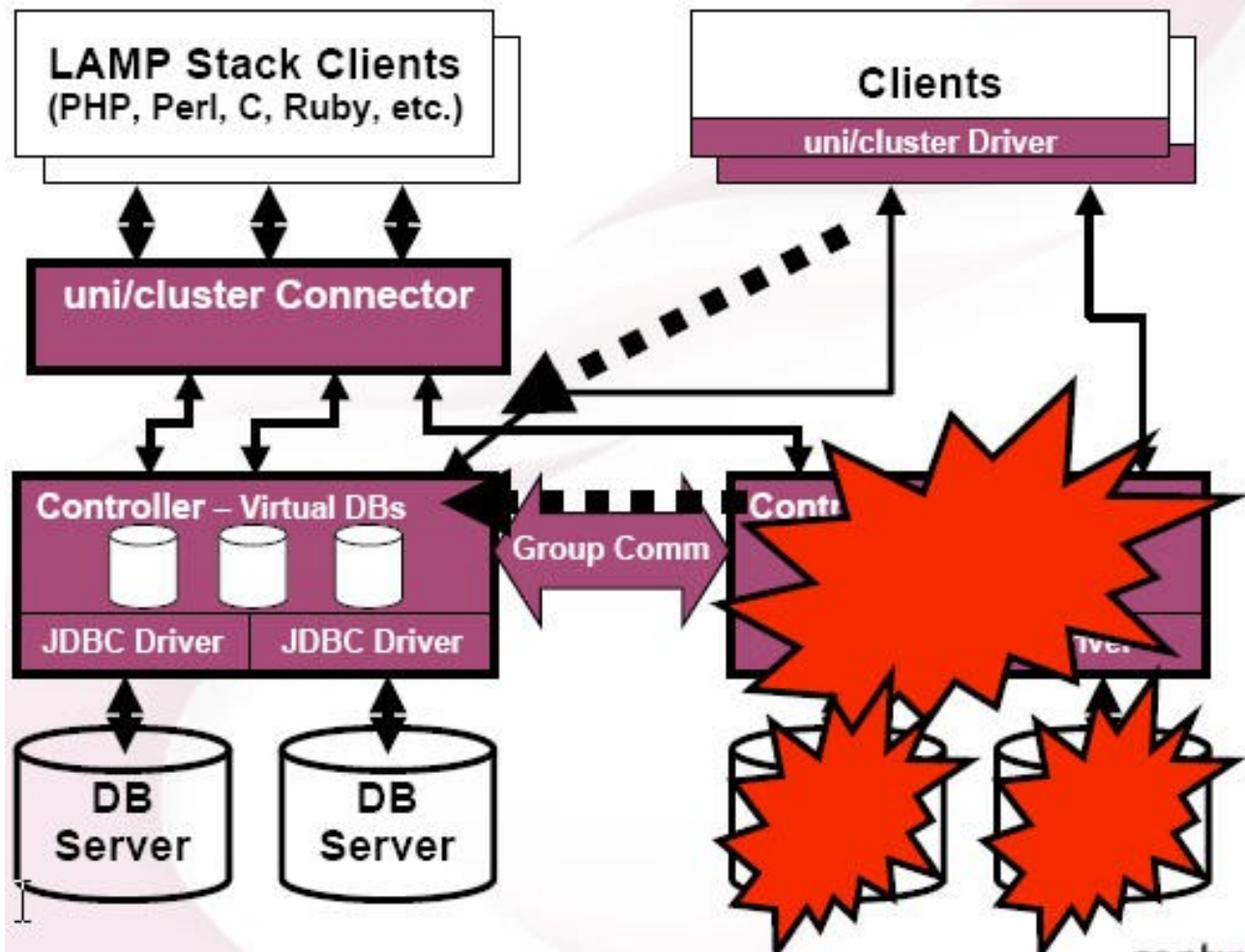
- **和磁盘 RAID 概念一样**
- **uni/cluster 数据库集群中间件**
  - 透明性：为客户端提供数据库的单一视图 client
  - 可扩展性：平衡后台数据库的负载
  - 高可靠性：管理实务，重新恢复和自动故障切换。



## Tungsten



## Failover/Recover





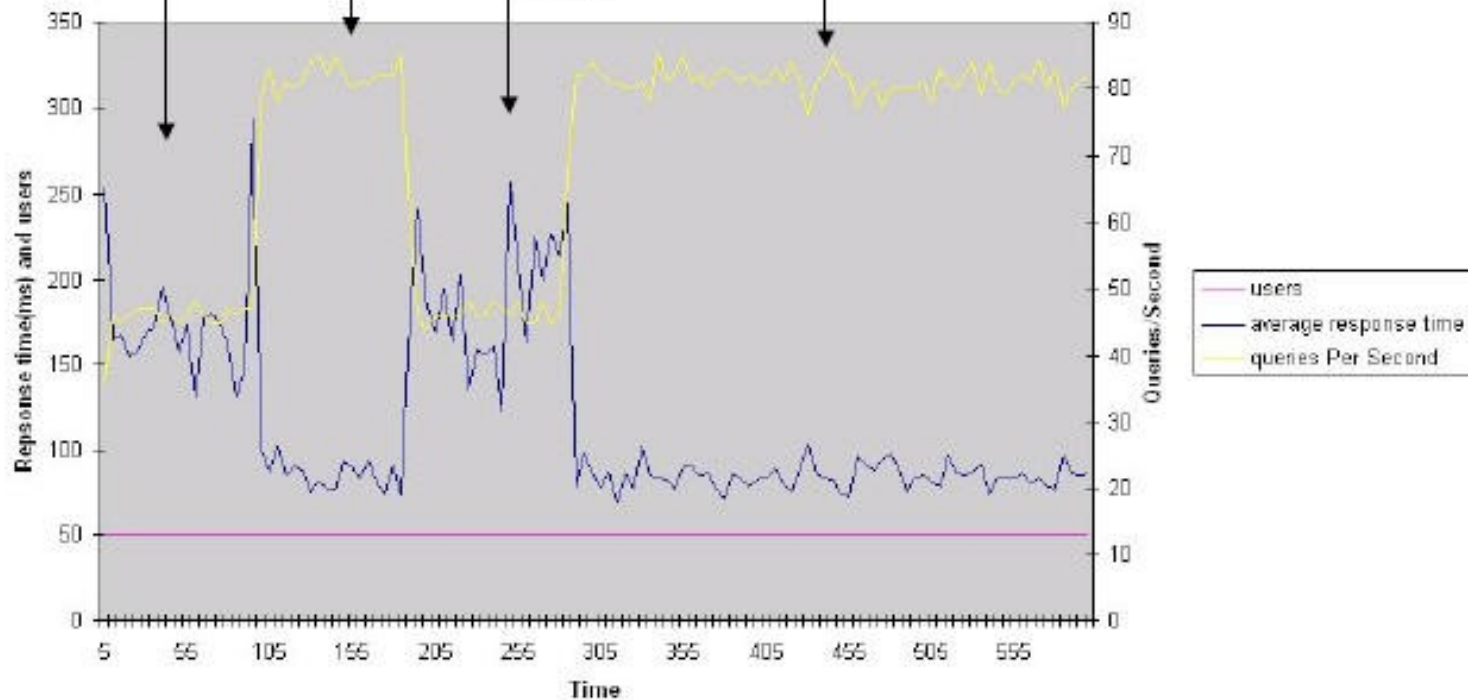
## uni/cluster - Non-stop Availability

Single Active Node

Two Active Nodes

Two Active Nodes

2<sup>nd</sup> Node  
Down



# Enterprise-Ready Cluster Solution

## Uni/Cluster GUI



## Continuent 竞争性分析

	Continuent	Oracle RAC
数据库	<ul style="list-style-type: none"><li>• 支持 MYSQL, PostgreSQL</li><li>• Oracle</li></ul>	<ul style="list-style-type: none"><li>• 仅适用于 Oracle</li></ul>
可扩展	<ul style="list-style-type: none"><li>• 多层次的扩展框架, 可解决不同层面的集群问题</li></ul>	<ul style="list-style-type: none"><li>• 独立的集群产品, 仅针对 Oracle 数据库</li></ul>
高可用	<ul style="list-style-type: none"><li>• 没有单点故障</li><li>• 错误故障自动恢复</li></ul>	<ul style="list-style-type: none"><li>• 数据请求负载均衡</li><li>• 错误故障自动恢复</li></ul>
数据一致性	<ul style="list-style-type: none"><li>• 直接的共享状态模型</li><li>• 有保障的连贯更新</li><li>• ACID 兼容</li></ul>	<ul style="list-style-type: none"><li>• 高风险的数据失误</li><li>• 高扩展的交易一致性更新</li><li>• Get/Put 模式易损坏和出错</li></ul>
操作管理	<ul style="list-style-type: none"><li>• 提供统一的管理监控客户端</li></ul>	<ul style="list-style-type: none"><li>• 依赖于 Oracle DB 的控制台</li></ul>
价格	<ul style="list-style-type: none"><li>• <b>¥ 112000/2CPU/5*8</b></li></ul>	<ul style="list-style-type: none"><li>• <b>¥ 330000/License</b></li></ul>