

白豚 ERP 系统技术白皮书

Version 1.10



版权所有©：2012-2022，上海重玄科技，保留所有权利

版本控制

版本号	修改时间	修改内容	修改人	审阅者
1.0	2013-01-12	创建	白杨	
1.1	2013-02-03	调整系统总体设计图和业务模块描述	白杨	
1.2	2013-12-21	增加系统实例截屏	白杨	
1.3	2015-06-21	加入 2.0 平台新特性描述	白杨	
1.4	2015-12-10	一些次要更新	白杨	
1.5	2017-05-07	新增 3.0 平台新特性描述	白杨	
1.6	2017-11-14	新增“1.概述”小节	白杨	
1.7	2017-11-15	新增“4.安全性和可靠性设计”小节	白杨	
1.8	2018-01-13	一些次要更新	白杨	
1.9	2020-03-21	新增 BYDMQ 和 BYST 等组件相关描述	白杨	
1.10	2022-02-20	对 BYPSS、BYDMQ 以及 BYST 等组件的个别细节进行澄清	白杨	

目录

版本控制.....	I
目录.....	II
1. 概述.....	4
1.1 应用支撑平台.....	9
1.1.1 单点支撑千万量级并发的高效 IO 服务器组件.....	10
1.1.2 强一致多活 IDC 高可用（HAC）和高性能（HPC）服务器集群.....	10
1.1.3 高效、高强度的密码编码学组件.....	16
1.1.4 数据查询分析引擎.....	16
1.1.5 更多.....	17
1.2 设计指标.....	17
2. 系统总体架构.....	20
2.1 通用后端服务.....	20
2.1.1 结构化存储（RDB、NewSQL、NoSQL）服务.....	21
2.1.2 分布式协调和消息分发服务（可选）.....	21
2.1.3 触发消息通知（可选）.....	21
2.2 App 服务器集群.....	22
2.3 反向代理服务（可选）.....	22
2.4 客户端.....	24
3. 系统总体设计.....	25
3.1 ERP 基础平台.....	25
3.1.1 应用支撑平台.....	26
3.1.2 基础服务群.....	29
3.2 可拔插业务模块.....	34
3.2.1 可拔插业务模块功能描述.....	34
3.2.2 可拔插业务模块设计框架.....	35
4. 安全性和可靠性设计.....	36
4.1 服务可用性和数据可靠性.....	36
4.1.1 在线数据可靠性.....	36
4.1.2 备份计划及离线数据可靠性.....	37
4.1.3 服务可用性.....	37
4.2 安全性设计.....	38
4.2.1 数据安全性.....	38
4.2.2 鉴权和授权安全性.....	39
4.2.3 传输安全性.....	39

4.2.4 系统安全性	40
5. 技术特点	43
5.1 REST 风格的 Web API 接口	43
5.2 nano-SOA 架构	44
5.2.1 SOA vs. AIO	44
5.2.2 nano-SOA 架构	46
5.2.3 白杨消息端口交换服务 (BYSS)	47
5.2.4 白杨分布式消息队列 (BYDMQ)	63
5.3 安全隧道服务 (BYST)	69

1. 概述

白豚 ERP 系统专门针对现代信息化环境设计。通过基于 HTML5 的新一代 B/S 架构，可以良好地兼容包括 IE、Firefox、Chrome、Opera、Safari、iOS (iPhone/iPad)、Android 以及 BlackBerry 等各类主流平台，能够方便地实现跨浏览器、跨平台的移动办公。

ERP 系统服务器端基于跨平台的 C/C++ 应用支撑框架而构建，支持 Windows、Linux、BSD、Solaris 等绝大部分主流操作系统，以及 x86/x64、IA64、ARM、RISC-V、MIPS、POWER、SPARC 等各类硬件平台。可根据客户要求来适应各类计算平台和云计算环境。

白豚 ERP 基于 SOA 思想进行设计，整个系统对外表现为一个独立的 ERP Service，除了 B/S 界面之外，还提供了一套清晰、完整的 API 接口。该接口基于 HTTP 协议，是满足 REST 风格的 WebAPI 实现，使用方便。非常适合用来进行二次开发。

ERP 系统支持多租户架构，可在一个服务器节点或同一分布式服务器集群中，动态承载和负载均衡多个租户，以此实现更高的资源（计算、存储、带宽等）利用率，大大降低客户总体拥有成本（TOC）。

通过使用灵活的业务模块插件（BMOD）机制，ERP 平台可以以可插拔的方式为用户提供人事管理、财务管理、客户关系管理、项目管理、档案管理、知识库、信息发布、办公协同等各种功能组件，并可按需对其进行组合、替换以及定制。BMOD 的安装和卸载均十分方便，只需几下鼠标轻点即可完成。

依托强大的统计分析插件机制，白豚 ERP 可以根据用户实际需要灵活地进行商务智能分析，并自动生成符合要求的各类统计报表。报表既可以方便地导出为 Excel、HTML 等便于阅读和打印的格式，也可以保存为 CSV、JSON 等利于其它工具进行二次处理的格式。

基于我方独有的分布式复杂事件处理（DCEP）技术，以及我方自主研发的先进工作流引擎、分布式数据查询分析引擎等优质组件，白豚 ERP 可以极其智能的方式为用户自动处理大量原本需要繁琐人工处理的业务逻辑。例如：在计算员工出差补贴时，系统可根据企业规则，综合关联员工所在部门、成本中心、职位、政治面貌、籍贯、性别、出发地、目的地、出差期间绩效表现，以及出差期间包含的工作日、节假日、公休日天数等各方面信息，自动、准确地完成相关计算工作。

不仅如此，得益于上述先进技术，对于考勤规则（休假规则、差旅规则、加班规则、打卡规则等）、绩效规则、工资计算、用工管理等系统所囊括的各方面业务逻辑及其中的每一处细节，白豚 ERP 均做到了智能完成相关的复杂业务关联计算，**高度智能的自动化数据处理为用户节省了大量人工成本，显著提高的用户工作效率。**

除此之外，白豚 ERP 平台还支持对员工、部门、站点、岗位、休假信息、薪资帐套等**所有关键数据的多版本无限回溯**——公司主管等管理人员可以随时审计这些数据的每一次修改，查看其

变更时间、变更人、相关工作流和变更数据摘要等关键信息，并将数据回溯到任意指定的修订版。除了可用于历史回顾和审计外，回溯功能还可有效应对用户对数据的恶意修改或意外破坏。

带有类似功能的其它产品大多仅能保留少量字段的变更记录（例如：可以查看组织架构和人员分布的历史信息，但在历史视图中查看一位员工的详细信息时，其所在站点、手机号、银行卡号、护照编号、联系地址、岗位、薪资结构、休假信息等所有相关数据仍旧是当前最新值，而非历史视图所在时间节点时的内容）。得益于我方自主研发的高效增量算法和高性能版本控制功能，白豚 ERP 平台可完整保留用户进行的每一次修订细节，不会丢失任何历史数据。因此可提供真实、完整、一致的精准历史视图，帮助用户更有效地完成企业规划和组织管理等工作。

在后台数据支持方面，得益于《白杨应用支撑平台》中提供的先进数据库连接器插件机制，以及由支撑平台提供的查询分析引擎等大数据处理组件。因此白豚 ERP 系统对底层 DBMS 产品的数据模型和功能模型均无重度依赖。通过简单地替换针对不同数据库产品的 DBC 插件，系统不但可快速支持 ORACLE、DB2、MySQL、MS SQL Server、PostgreSQL、SQLite 等多种主流企业级数据库产品，也能够方便地迁移到 MongoDB、Cassandra 等 NoSQL 产品，或 VoltDB、Clustrix、MySQL Cluster 等 NewSQL 产品上。同时，通过对逻辑并行（流水线）算法、数据并行算法、内存消息解析以及异步 IO+线程池等技术的综合利用，系统能够瞬间处理上百万并发请求。

在安全性方面，白豚 ERP 系统使用网银级的 TLS 通信强加密（HTTPS）以及可靠的 CHAP 鉴权协议保证通信安全，所有数据整库存储于我方自主研发的，支持实时数据强加密的分布式文件系统（VFS）保障用户数据安全。

得益于上述分布式数据处理和分布式强加密系统等优势技术，白豚云平台为用户提供了我方独有的，支持整库强加密的多活 IDC 6 副本技术的分布式数据库系统，其数据可靠性高达“17 个 9”（99.99999999999999%）。与此同时，白豚云平台还为用户保存了最近 90 天内的全量数据快照备份（7+12），在进一步增强数据可靠性的同时，还确保了用户可有效恢复被高权限员工恶意删除或破坏的数据。

白豚 ERP 系统基于由我方自主研发的，拥有多项国家和国际发明专利保护的“BYSS 分布式协调服务”而构建，支持抗脑裂（Split Brain）的强一致高可用（HAC）+ 高性能分布式集群（HPC），可实现极高的服务可用性和线性的横向扩展能力。与此同时，得益于我方汇编优化，基于 DMA+硬件中断的零拷贝异步 IO 伺服器模型，白豚 ERP 系统可提供单点千万量级并发连接和百万量级并发请求的强大处理能力。彻底摆脱了以往 ERP 软件低效、稳定性差、对数据库和操作系统等平台依赖性强、以及只支持 IE 浏览器平台等问题。

白豚云平台即为基于上述高效模型的分布式集群。依托于跨越多座 IDC 的强大分布式集群架构以及秒级的故障转移能力，白豚云平台可为用户提供高达“6 个 9”（99.9999%）的服务可用性保障，承诺平均每年的不可用（故障）时间低于 31.5 秒（详见可靠性说明和 SLA 文档）。

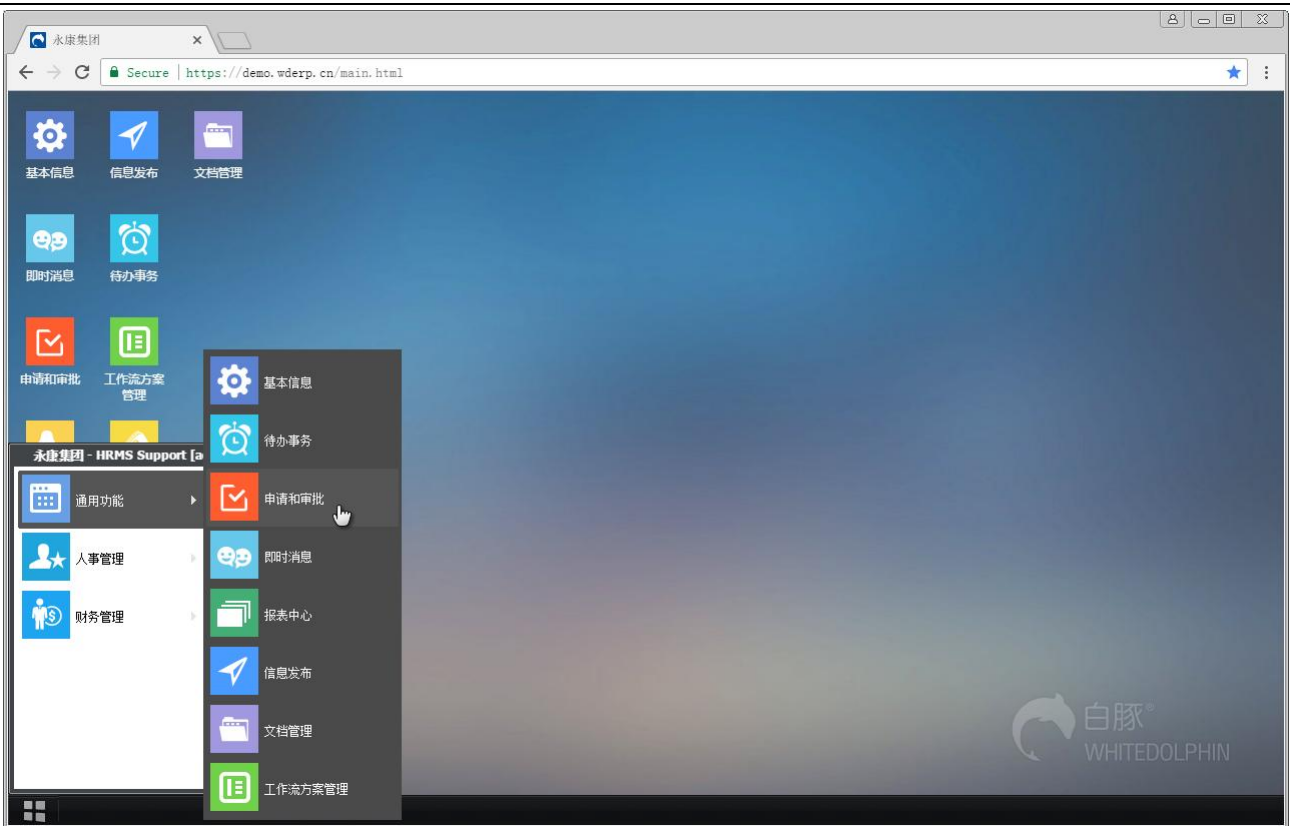


图 1 白豚 ERP 操作界面实例

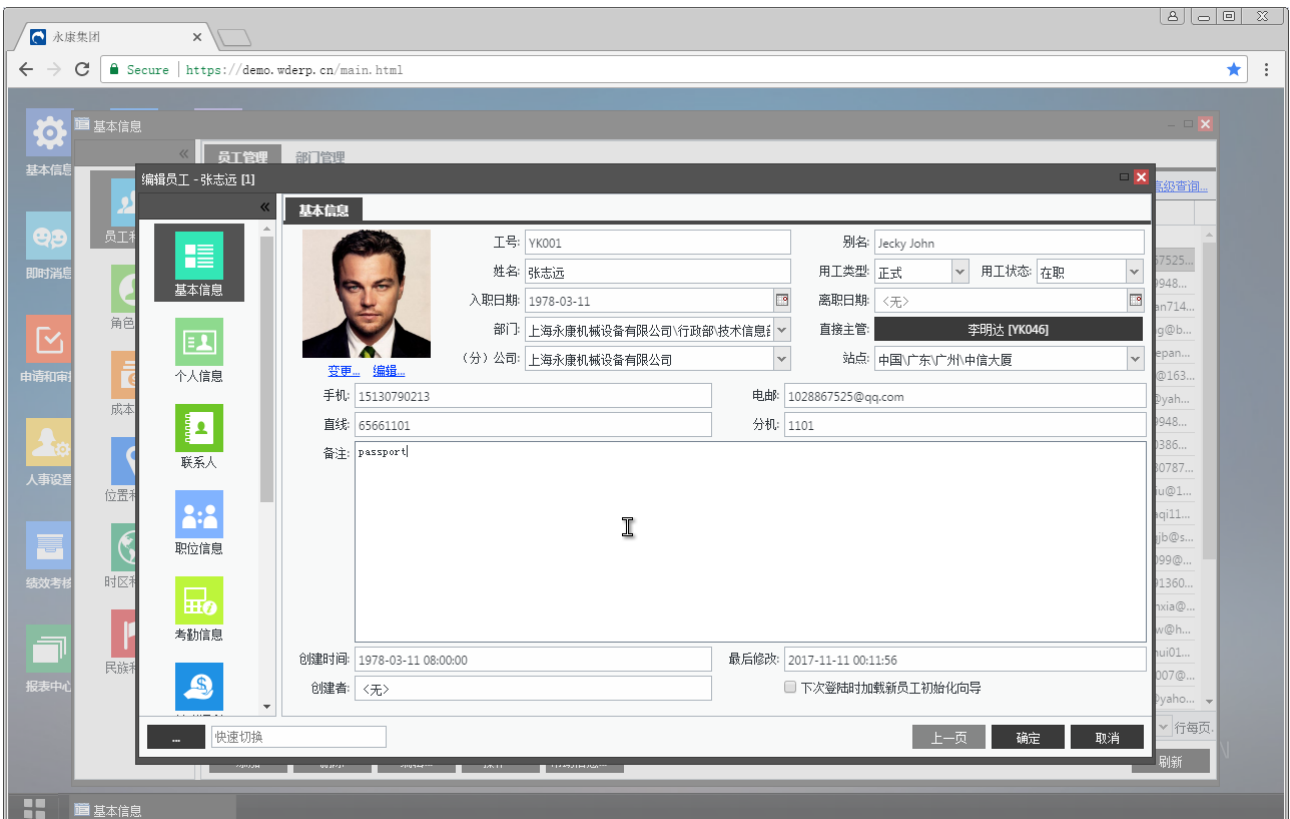


图 2 白豚 ERP 操作界面实例

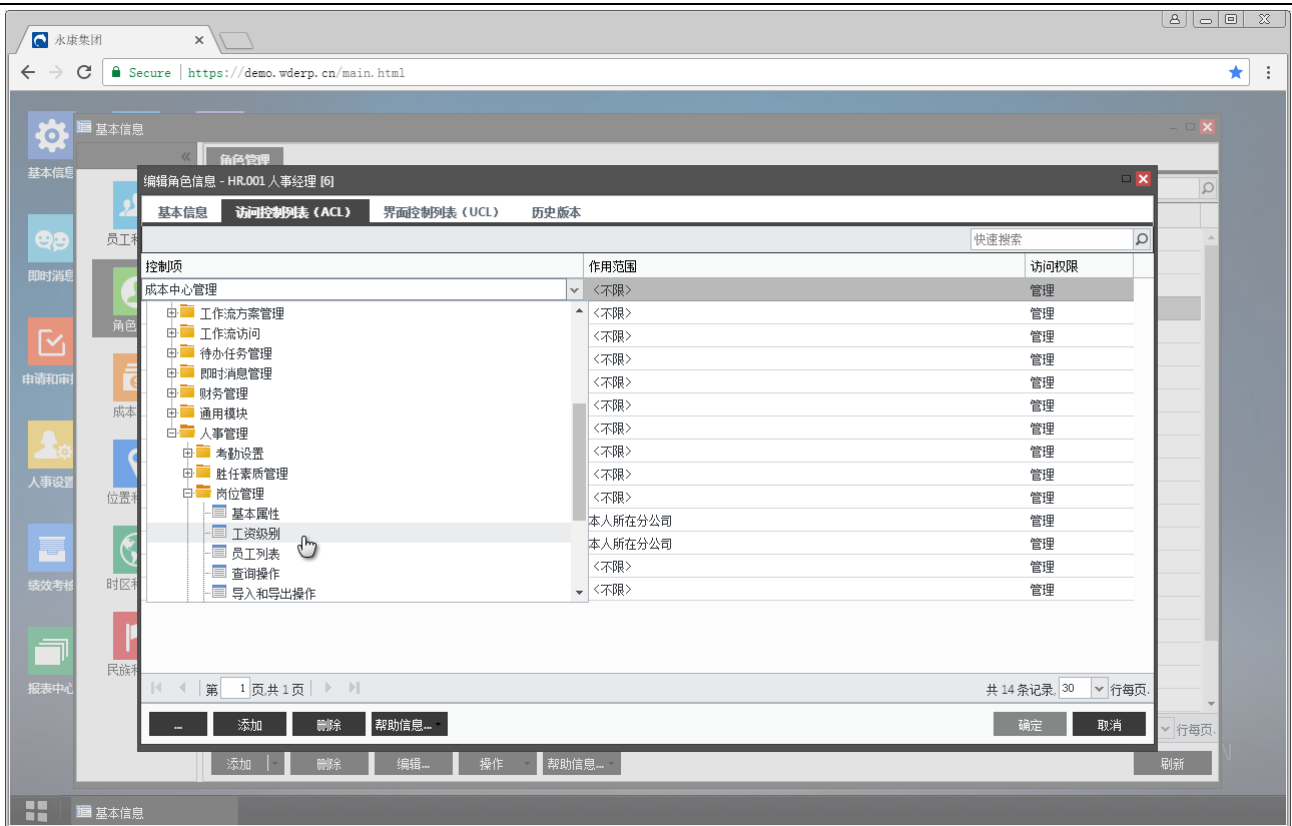


图 3 白豚 ERP 操作界面实例



图 4 白豚 ERP 操作界面实例

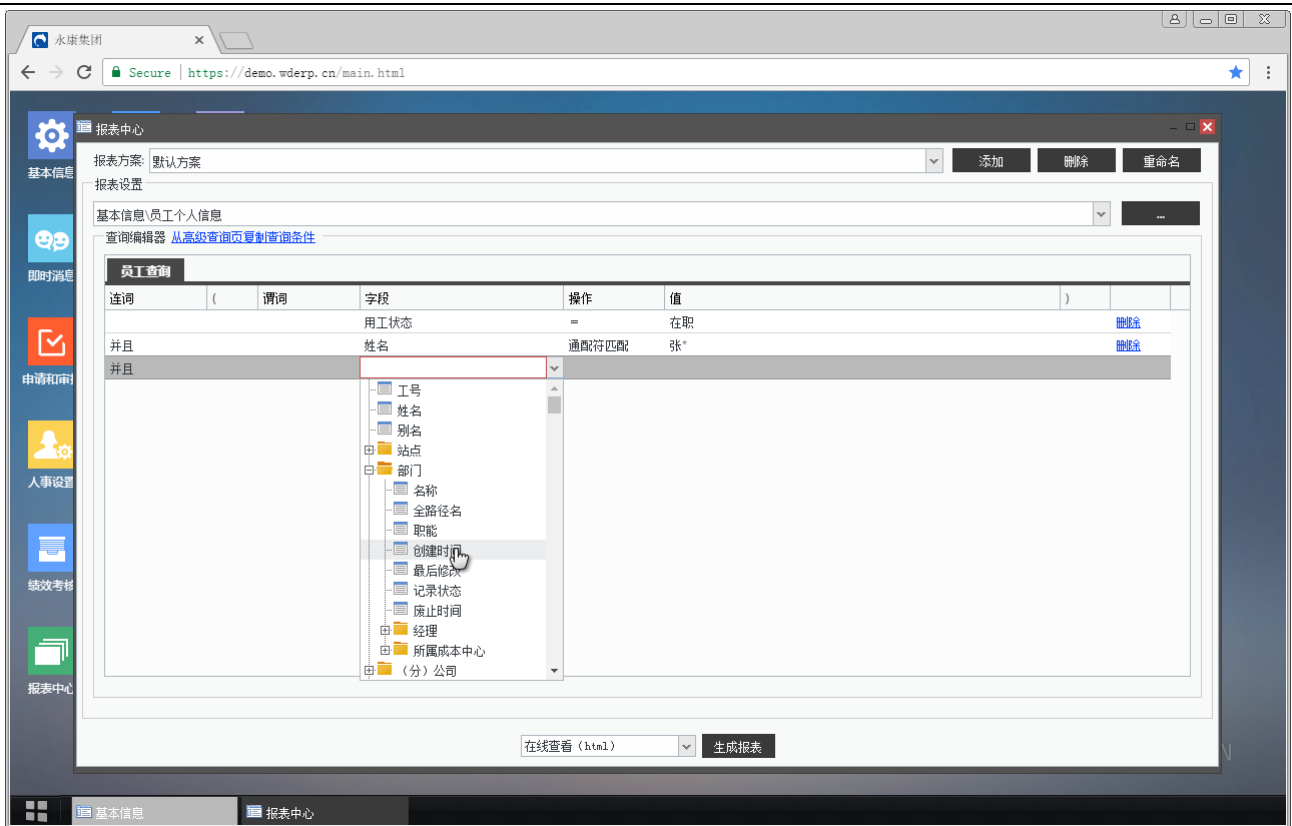


图 5 白豚 ERP 操作界面实例

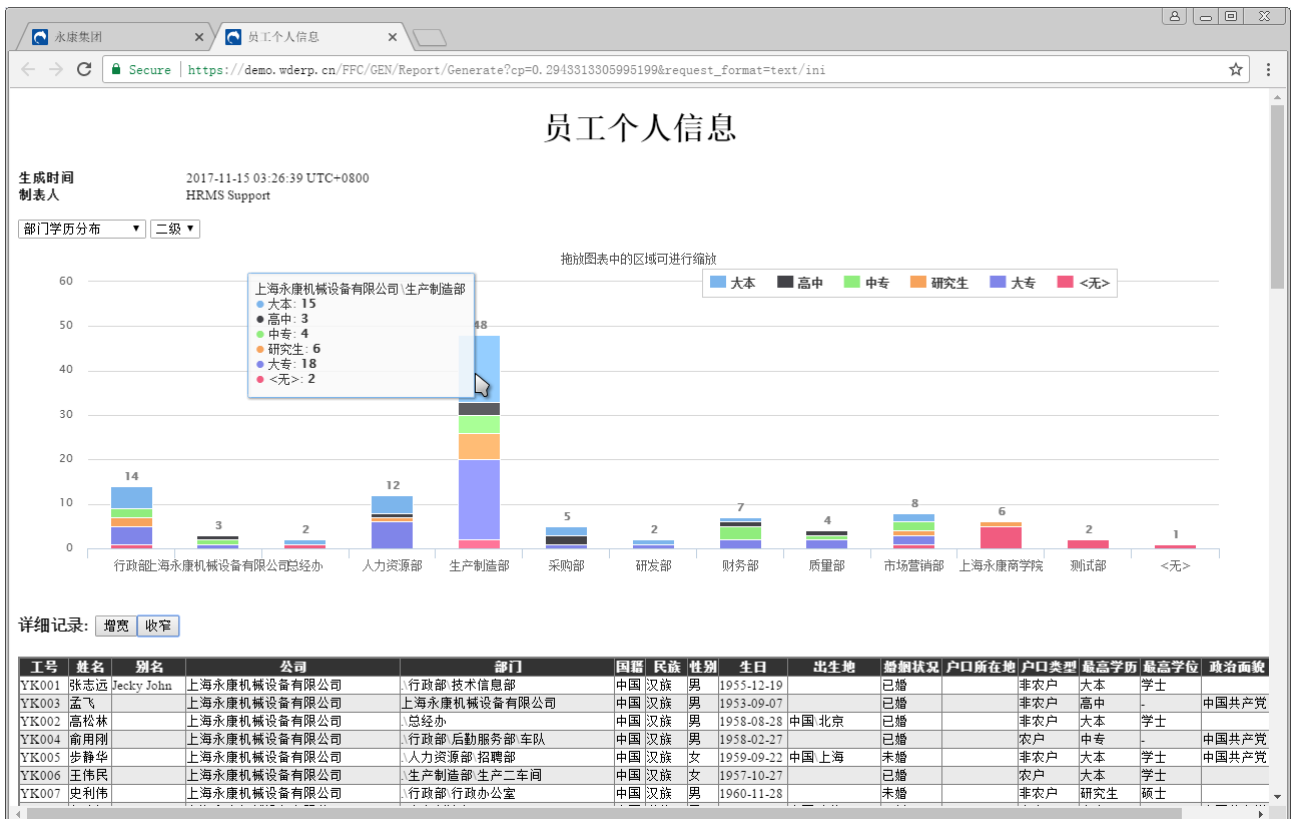


图 6 白豚 ERP 操作界面实例

1.1 应用支撑平台

《[白杨应用支撑平台](#)》作为我方的核心竞争力，使用汇编、C/C++ 构建，包含数百万行代码和上千个成熟的通用组件。上述单点高并发和多活 IDC 分布式集群等核心技术均由白杨应用支撑平台提供。多年来，基于支撑平台的各种产品已被广泛部署于包括兴业银行（China CIB）、中石油（CNPC）、华安保险（Sinosafe Insurance）、淘宝网（taobao.com）、易果网（yiguo.com）、烟台万华集团、法国兴业银行（SOCIETE GENERALE）、中国农业银行（ABC）、宝马集团（BMW）、陕汽集团（陕汽重卡）、德尔福汽车（Delphi）、美联航（United Airlines）、GE（美国通用电气）、贝塔斯曼（Bertelsmann）、埃森哲（Accenture）、光大银行（CEB）、壹基金（One Foundation）、中国宋庆龄基金会、中国移动（China Mobile）、中国联通（China Unicom）、国家电网（SGCC）等各大企业在内的不同生产环境中：



真实生产环境下的大范围部署不但为上层应用提供了可靠的、平台无关的底层环境，也进一步检验了支撑平台的可靠性、稳定性、可移植性、高效性等各方面指标。

支撑平台使用汇编、C/C++ 构建，支持 Windows、Linux、BSD、IBM AIX、HP-UX、Solaris、MAC OS X、vxWorks、QNX、DOS、WinCE (Windows Mobile)、NanoGUI、eCos、RTEMS 以及 Android、iOS 等绝大多数主流操作系统。支持 x86/x64、ARM、RISC-V、IA64、MIPS、POWER、SPARC 等主流硬件平台。在提供了大量高品质可重用组件的同时，也确保了良好的可移植性。

全球专利证书一览



上千成熟可靠的高品质功能组件可在性能、功能和稳定性等方面大幅提升软件产品的品质，并为产品的开发带来了难以想象的便利，例如：

1.1.1 单点支撑千万量级并发的高效 IO 服务器组件

支撑平台使用汇编和异步 IO 对网络服务组件进行了优化，通过 DMA + 硬件中断实现内存零拷贝的高效异步网络服务。性能、可靠性和可伸缩性都很强。可在 2011 年出厂的，当时售价不足 2 万元人民币的单台至强 5600 系入门级 1U PC Server 上支撑上千万 TCP/HTTP 并发连接。相对应地，一般由 Java / .NET 开发的服务器端，在相同配置的机器上，单点最高仅可支撑 3000 到 5000 并发，PHP 则更低（具体可参考《白杨应用支撑平台技术白皮书》：3.2.1、3.3.1 以及 3.3.2 小节）。

1.1.2 强一致多活 IDC 高可用（HAC）和高性能（HPC）服务器集群

强一致、抗脑裂（Split Brain）的多活 IDC 分布式高可用（HAC）和高性能（HPC）计算集群支持：独有的 nano-SOA 大规模分布式架构在保持高内聚、低耦合设计的前提下，将单点性能提升到了远超传统 SOA 架构的水平，同时简化了集群部署，提高了集群的可维护性。

白杨消息端口交换服务（BYPSS）：一种基于多数派算法的，强一致（抗脑裂）、高可用的分布式协调组件，可用于向集群提供服务发现、故障检测、服务选举、分布式锁等传统分布式协调服务，同时还支持消息分发与路由等消息中间件功能。由于通过专利算法消除了传统 Paxos/Raft

中的网络广播和磁盘 IO 等主要开销，再加上批量模式支持、并发散列表、高并发服务组件等大量其它优化，使得 BYPASS 可在延迟和吞吐均受限的跨 IDC 网络环境中支持百万节点、万亿端口量级的超大规模计算集群。

带强一致保证的多活 IDC 技术是现代高性能和高可用集群的关键技术，也是业界公认的主要难点。作为实例：2018 年 9 月 4 日微软美国中南区某数据中心空调故障导致 Office、Active Directory、Visual Studio 等服务下线近 10 小时不可用；2015 年 8 月 20 日 Google GCE 服务中断 12 小时并永久丢失部分数据；2015 年 5 月 27 日、2016 年 7 月 22 及 2019 年 12 月 5 日支付宝多次中断数小时；以及 2013 年 7 月 22 日、2023 年 3 月 29 日微信服务中断数小时等重大事故均属于产品未能实现多活 IDC 架构，单个 IDC 故障导致服务全面下线的惨痛案例。

在上述方面，我方拥有超过十年的积累，[掌握多项受国家和国际发明专利保护的分布式架构和算法](#)。得益于这些领先的强一致、高可用、高性能分布式集群算法和架构，我们在蓝鲸、白豚、职业精等全线产品上，均实现了真正的多活 IDC 架构，为客户提供了无以伦比的数据可靠性和服务可用性保证。

1.1.2.1 分布式协调服务

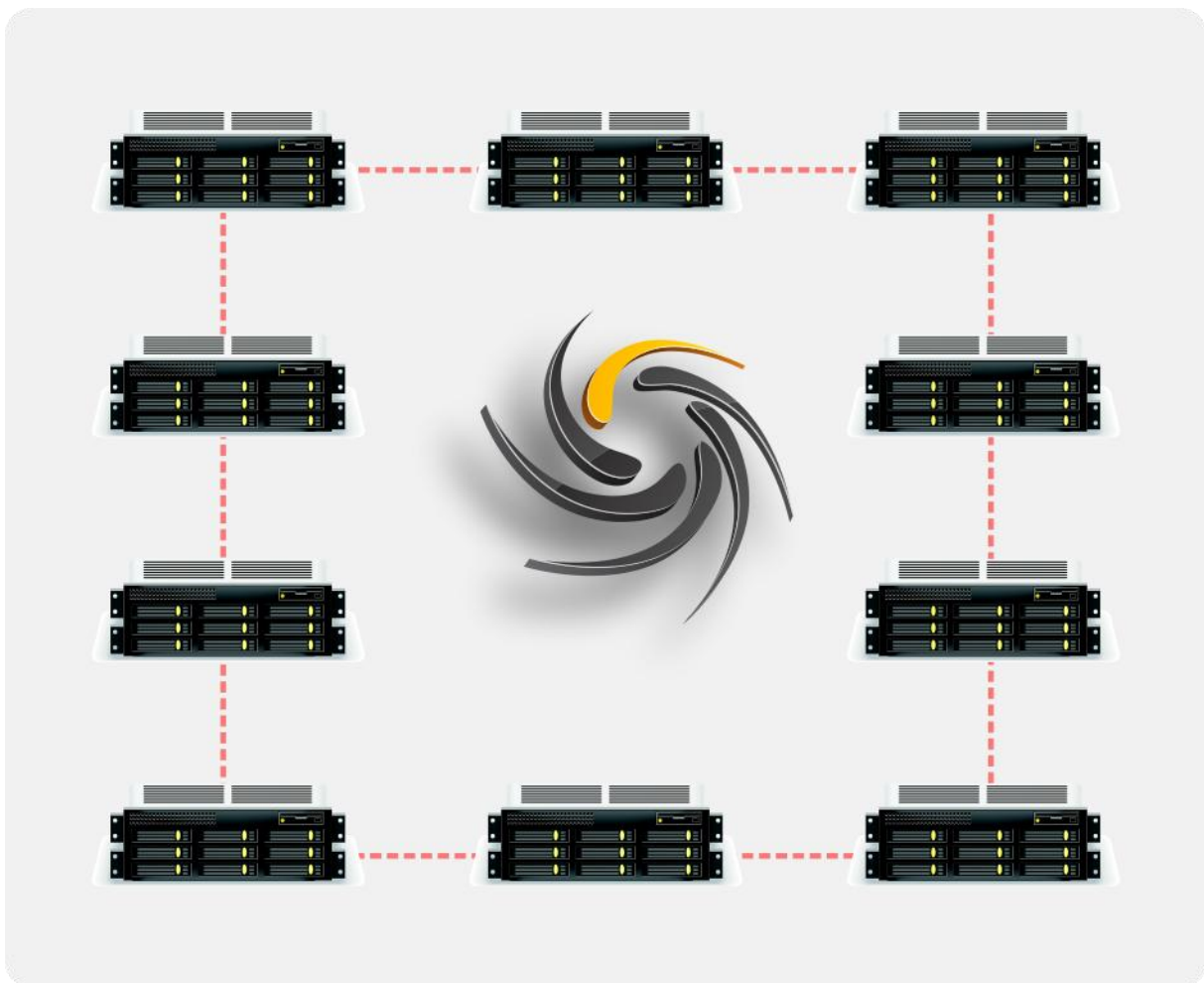


图 7

分布式协调服务为集群提供服务发现、服务选举、故障检测、故障转移、故障恢复、分布式锁、任务调度，以及消息路由和消息分发等功能。

分布式协调服务是分布式集群的大脑，负责指挥集群中的所有服务器节点协同工作。将分布式集群协调为一个有机整体，使其有效且一致地运转。实现可线性横向扩展的高性能（HPC）和高可用（HAC）分布式集群系统。

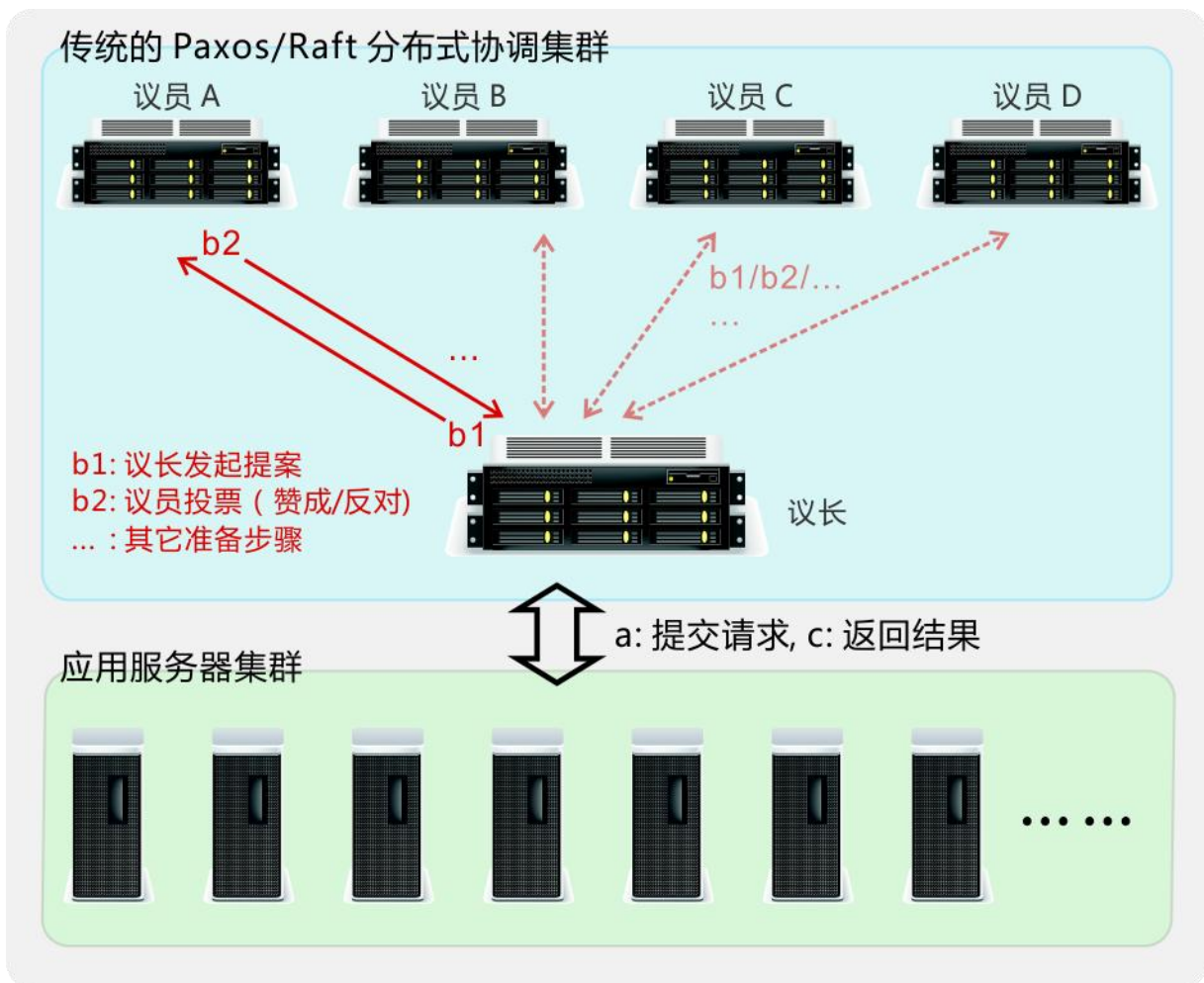


图 8

传统的 Paxos/Raft 分布式协调算法为每个请求发起投票，产生至少 2 到 4 次网络广播（b1、b2、...）和多次磁盘 IO。使其对网络吞吐和通信时延要求很高，无法部署在跨 IDC（城域网或广域网）环境。

我们的专利算法则完全消除了此类开销。因此大大降低了网络负载，显著提升整体效率。并使得集群跨 IDC 部署（多活 IDC）变得简单可行。

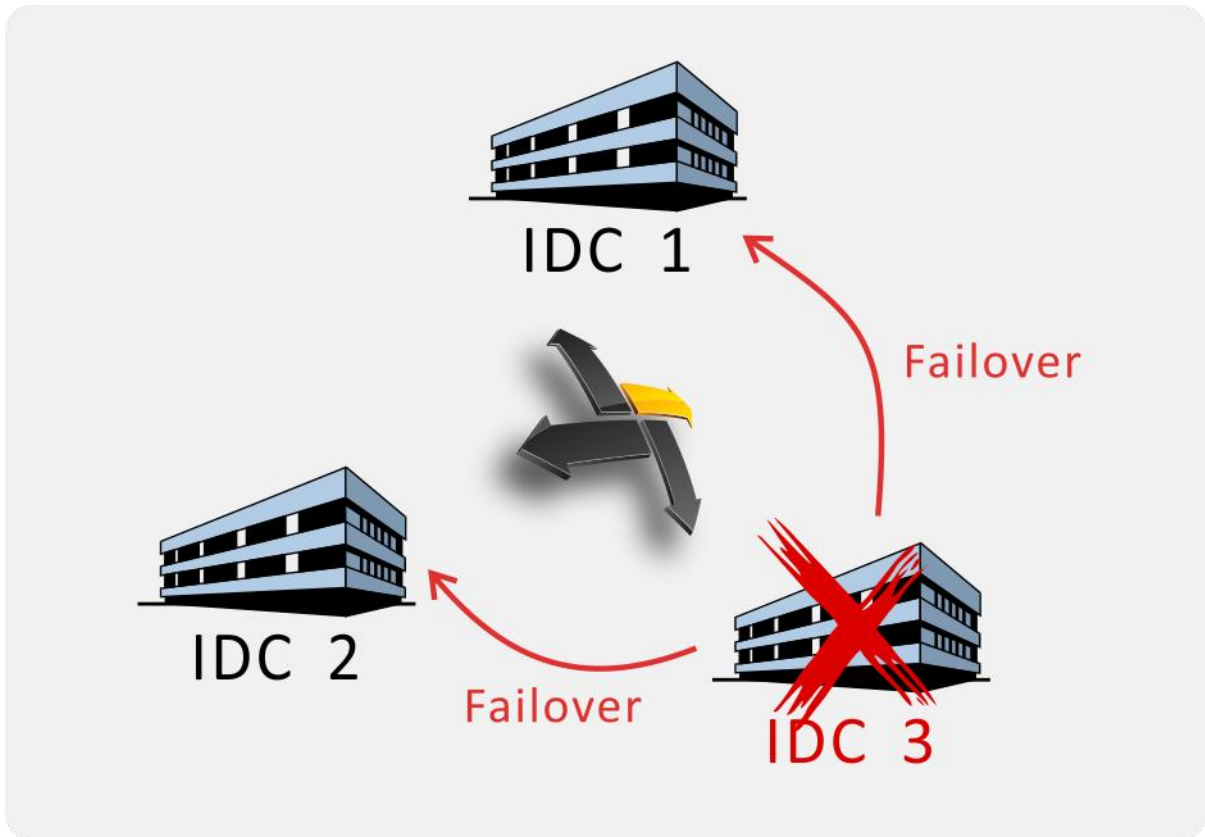


图 9

基于我方独有的分布式协调技术，可实现高性能、强一致的多活 IDC 机制。可在毫秒级完成故障检测和故障转移，即使整座 IDC 机房下线，也不会导致系统不可用。同时提供强一致性保证：即使发生了网络分区也不会出现脑裂（Split Brain）等数据不一致的情形。例如：

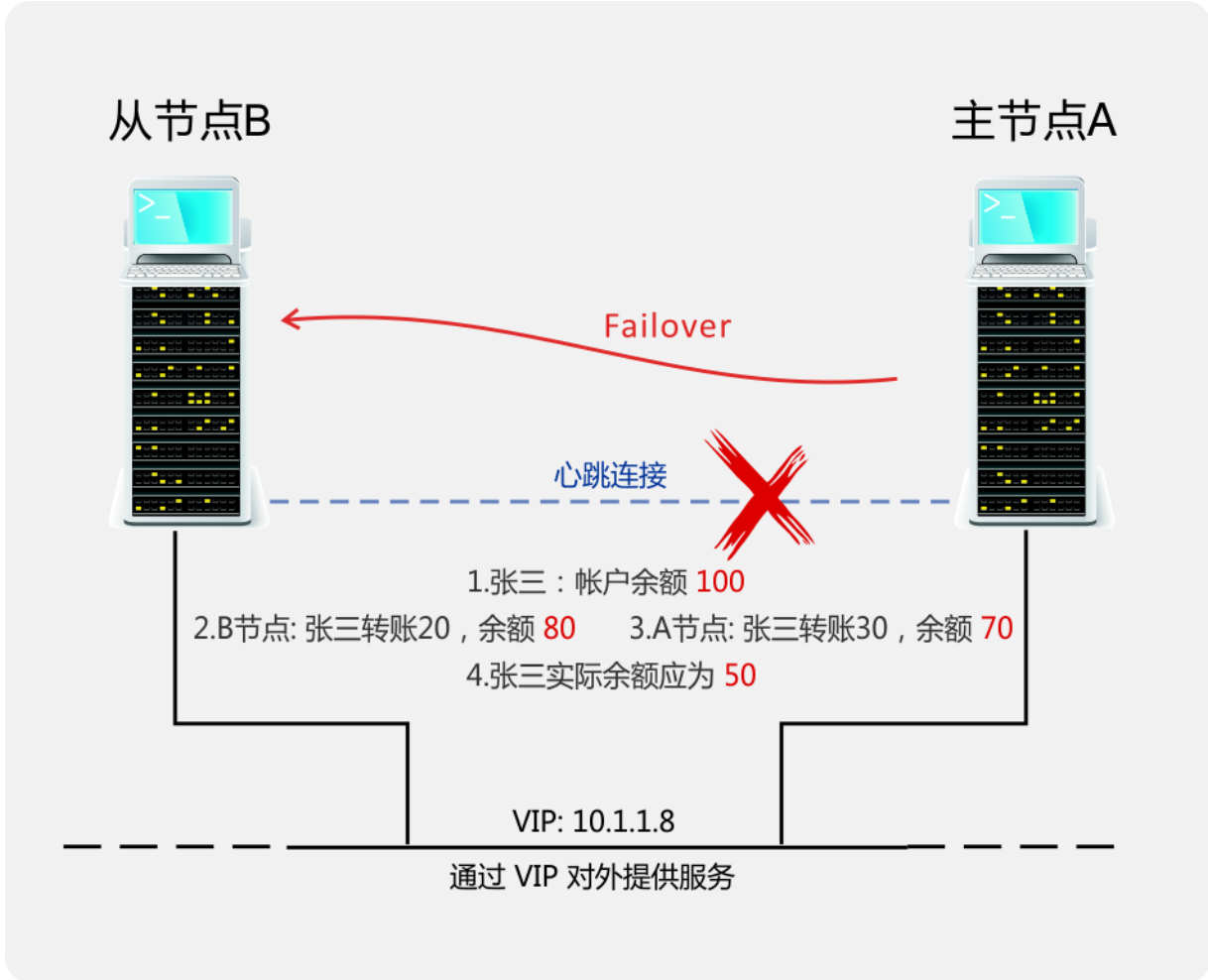


图 10

在传统的双机容错方案中，从节点在丢失主节点心跳信号后，会自动将自身提升为主节点，并继续对外提供服务，以实现高可用。在此种情形中，当主从节点均正常，但心跳连接意外断开时（网络分区），就会发生脑裂（Split Brain）问题，如图 10 所示：此时 A、B 均认为对方已下线，故将自己提升为主节点并分别对外提供服务，产生难以恢复的数据不一致。

我方 BYPSS 服务可提供与传统 Paxos/Raft 分布式算法相同水平的强一致性保证，从根本上杜绝脑裂问题的发生。

类似地：工行、支付宝等服务也有异地容灾方案（支付宝：杭州 → 深圳、工行：上海 → 北京）。但在其异地容灾方案中，两座 IDC 之间并无 Paxos 等分布式协调算法保护，因此无法实现强一致，也无法避免脑裂。

举例来说，一个在支付宝成功完成的转账交易，可能要数分钟甚至数小时后才会从杭州主 IDC 被异步地同步到深圳的灾备中心。杭州主 IDC 发生故障后，若切换到灾备中心，意味着这些未同步的交易全部丢失，并伴随大量的不一致。比如：商家明明收到支付宝已收款提示，并且在淘宝交易系统看到买家已付款，并因此发货。但由于灾备中心切换带来的支付宝交易记录丢失，导致在支付宝中丢失了相应的收入，但淘宝仍然提示买家已付款。因此，工行、支付宝等机构在主 IDC

发生重大事故时，宁可停止服务几个小时甚至更久，也不愿意将服务切换到灾备中心。只有在主 IDC 发生大火等毁灭级事故后，运营商才会考虑将业务切换到灾备中心（这也是灾备中心建立的意义所在）。

因此，异地容灾与我方的强一致、高可用、抗脑裂多活 IDC 方案具有本质区别。

此外，Paxos / Raft 在经历过半节点同时故障下线并维修恢复的过程中，无法保证数据的强一致性，**可能产生幻读等不一致问题**（例如：在一个三节点集群中，节点 A 因为电力故障下线，一小时后节点 B 和 C 则因为磁盘故障下线。此时节点 A 恢复电力供应重新上线，紧接着管理员更换了节点 B 和 C 的磁盘并让它们分别恢复上线。此时整个集群 1 小时内的修改将全部丢失，集群退回到了 1 小时前 A 节点下线时的状态）。而 BYPSS 则从根本上避免了此类问题的发生，因此 BYPSS 拥有比 Paxos / Raft 更强的一致性保证。

由于消除了 Paxos/Raft 算法中的大量广播和分布式磁盘 IO 等高开销环节，配合支撑平台中的高并发网络服务器、以及并发散列表等组件。使得 BYPSS 分布式协调组件除了上述优势外，还提供了更多优秀特性：

批量操作：允许在每个网络包中，同时包含大量分布式协调请求。网络利用率极大提高，从之前的不足 5% 提升到超过 99%。类似于一趟高铁每次只运送一位乘客，与每班次均坐满乘客之间的区别。实际测试中，在单千兆网卡上，可实现 400 万次请求每秒的性能。在当前 IDC 主流的双口万兆网卡配置上，可实现 8000 万次请求每秒的吞吐。比起受到大量磁盘 IO 和网络广播限制，性能通常不到 200 次请求每秒的 Paxos/Raft 集群，有巨大提升。

超大容量：通常每 10GB 内存可支持至少 1 亿端口。在一台插满 64 根 DIMM 槽的 1U 尺寸入门级 PC Server 上（8TB），可同时支撑至少 800 亿对象的协调工作；在一台 32U 大型 PC Server 上（96TB），可同时支撑约 1 万亿对象的分布式协调工作。相对地，传统 Paxos/Raft 算法由于其各方面限制，通常只能有效管理和调度数十万对象。

问题的本质在于 Paxos / Raft 等算法中，超过 99.99% 的代价都消耗在了网络广播（投票）和落盘等行为上。而这些行为的目的是要保证数据的可靠性（数据要同时存储在多数节点的持久化设备上）。而服务发现、服务选举、故障检测、故障转移、故障恢复、分布式锁、任务调度等分布式协调功能所涉及到的恰恰又都是没有长期保存价值的临时性数据。因此花费超过 99.99% 的精力来持久化地保存它们的多个副本是毫无意义的——就算真的发生主节点下线等罕见灾难，我们也可以极高的效率，在瞬间就重新生成这些数据。

就好像张三买了一辆车，这辆车有个附加保险服务，其条款为：在张三万一发生了致命交通事故时，它能提供一种时光倒流机制，将其带回到意外发生之前的一瞬间来避免这场意外的发生。当然，这么牛的服务肯定也很贵，它大概需要预付张三家在接下来三生三世里能获得的所有财富。而且即使张三在驾驶这辆车的过程中，始终未发生过致命交通事故，那这些预先支付的服务费也是一分钱都不能减免的。这么昂贵的服务，且不说一般人一生中大概率都不会发生致命交通事故（更别提还要指定具体的某辆车）。即使真发生了，这个三代赤贫的代价也难说就值得吧？

而我们则为自己的汽车产品提供了另一种不同的附加服务：虽然没有时光倒流功能，但我们

的服务可以在张三发生致命事故后，将所有受害方全体连车带人瞬间原地满血复活（是一根头发丝都不会少、一块漆皮都不会掉那种满血）。最关键的是，该服务无需预先收取任何费用。张三只需要在每次这样的灾难发生以后，支付相当于其半个月的工资的再生技术服务费就可以了。

综上，我方专利的分布式协调算法，在提供与传统 Paxos/Raft 算法相同等级的强一致性和高可用性保证之同时，极大地降低了系统对网络和磁盘 IO 的依赖，并显著提升了系统整体性能和容量。对于大规模、强一致分布式集群的可用性（HAC）和性能（HPC）等指标均有显著提升。

关于 BYPASS 服务的进一步描述，详见：5.2.3 白杨消息端口交换服务（BYPASS）。

1.1.3 高效、高强度的密码编码学组件

包含公钥算法、对称加密算法、数据编解码、散列和消息验证算法、数据压缩算法等基础功能组件。除此之外，支撑平台还提供了多个经过高度抽象、可开箱即用的高级密码编码学功能组件，例如：

支持实时压缩和强加密的虚拟文件系统（VFS），VFS 支持包括 AES（128/256）、TwoFish 等在内的数十种强加密算法，使用 AES-NI、SSE4 等汇编指令集优化，效率高。在蓝鲸、白豚、职业精等全线产品中，我们均使用该组件为产品的数据库和配置类数据提供整库级的实时压缩和强加密保护。另外还包括基于公钥体系架构（PKI）的强加密通信保护组件等。

近年来安全问题频发，亚马逊、沃尔玛、Yahoo、Linkedin（领英）、索尼、摩根大通、OpenAI（ChatGPT）、UPS、eBay、京东、支付宝、一号店、协程、12306、网易、CSDN、中国人寿、以及各大酒店集团（如家、汉庭、锦江、洲际、喜来登、万豪等）等国内外知名企业均频频报出大量用户信息泄露的严重安全事件，安全保障已经刻不容缓。

我方所有数据库（整库）和本地配置数据均存放在我方自主研发的，支持实时（on-the-fly）数据压缩和强加密的虚拟文件系统（VFS）中进行全方位保护。支持数十种业界公认的强加密安全算法，即使系统管理员也无法窥视企业数据。

基于业界标准的强加密算法保证了即使未来出现了每秒钟可完成一千万亿次密钥破解尝试的超级计算机，也需要平均五千四百万亿年才能破解一个密钥。安全性得到了极大保证（具体可参考《[白杨应用支撑平台技术白皮书](#)》中的第 4 节）。

1.1.4 数据查询分析引擎

支撑平台中还包含了表达力优于 SQL 的查询分析引擎，自有查询引擎除了可以摆脱对特定 DBMS 的依赖，使我们的产品可以自由地在 MySQL、MS SQL Server、Oracle、DB2、SQLite 等 RDBMS 以及 MongoDB、Cassandra 等 NoSQL 数据库间灵活切换。更增加了基于 UNICODE 字符集的 ARE 高级正则查询、表中套表的关联查询、复杂虚拟字段等 SQL 没有的高级特性。

查询引擎通过汇编优化的 C/C++ 代码实现词法分析、语法分析、语义分析/中间代码生成、优化等步骤，并可在 2010 出厂的 Thinkpad W510（4 核 8 线程，主频 1.6G）笔记本上，仅用其中一核一线程即可达到每秒 1300 万次以上的表达式求值效率。

1.1.5 更多...

我们并不依靠“商业秘密”来保护核心竞争力。相反，我们使用更公开透明的商标、认证、版权、专利和公证保管等手段来保护自己的合法权益。因此，我们的技术细节均公开于相应的文档中，详情可见《白杨应用支撑平台技术白皮书》：

http://baiy.cn/doc/asp_whitepaper.pdf

或者：http://baiy.cn/doc/asp_whitepaper_en.pdf（英文版）

等文档，包括 Hacker News（全球最大的计算机科学新闻网站）、Google Blogger、CSDN 以及博客园在内的多家国内外媒体均转载或报道了这些论文。相较于“严守秘密”，我们相信公开透明下的大量同行审评，加上实际生产环境下的严酷考验，更有利于产品品质的提升。

1.2 设计指标

白豚 ERP 产品可达到或满足如下设计标准：

- ★ **用户数：**单点单库支持用户数不低于 20 万，单点在线并发用户数支持不低于 100 万⁽¹⁾。并可通过向现有集群中增加节点、升级到分布式 DBMS（分库分表）等方式来线性增加用户容量。
- ★ **前台并发数：**前台伺服界面在主流 PC Server⁽¹⁾上可支持不少于 500 万个 HTTP/TCP 并发连接。
- ★ **后台并发数：**后台伺服界面在主流 PC Server⁽¹⁾上可支持不少于 500 个 HTTP/TCP 并发连接。
- ★ **稳定性：**系统能够长时间稳定连续地工作，期间无需人为干预。系统可以强一致、抗脑裂（Split Brain）的多活 IDC 高可用+高性能集群（HAC+HPC）方式来进行部署，因此可在保证高性能的同时，为用户提供高达“6 个 9”（99.9999%）的服务高可用性保障。
- ★ **可靠性：**针对所有支持的任务，系统均能够以可靠的方式完成。并确保任务执行的原子性、持久性、隔离性和一致性。如果在执行中发生不可避免的错误（如电源故障、磁盘损坏等），系统必须能够恢复到任务执行之前的状态。

系统中的所有数据均以强一致多副本的形式被同时保存在至少 2 座数据中心内，数据中

心之间应至少相隔十公里以上，每座数据中心中至少存储了数据的三个副本。上述分布式 6 副本技术可以提供高达“17 个 9”（99.99999999999999%）的数据可靠性，确保数据不会意外丢失。

除此之外，系统每日自动完成快照备份，保留最近 7 天和近 90 天内每周末的磁盘数据备份（7+12）。快照数据基于可靠的双重 EC 算法，可有效找回由员工误操作或恶意操作引起的数据丢失。同时也为上述分布式 6 副本技术额外增加了一层双重保险，进一步降低了本就微乎其微的数据丢失风险。

- ★ **弹性：**随着网络规模不断增长，系统所承受的负载将会相应上升。不同的负载量和健壮性要求将直接导致不同的软 / 硬件平台选择和体系结构的变化（例如：从初期运行 Windows 系统的单台 PC 服务器+嵌入式数据库到后期运行 UNIX 系统的小型机集群+分布式数据库）。系统的设计能够保证足够的伸缩性，使其能够以较高的性价比同时适应不同规模的实际需求。
- ★ **容量：**系统提供的横向扩展（HAC+HPC）能力可确保稳定支持单表万亿记录（行）、PB 级别的海量数据。
- ★ **可扩展性：**系统能够很好地支持数据并行、任务并发以及流水线并行等各种并发模式，以结合可配置的内存及处理器管理策略来实现较好的纵向扩展能力。系统中的各主要部件应支持松耦合的分布式架构，以便于便捷地实现较高的横向扩展能力。
- ★ **可移植性：**系统中的所有平台相关性均已被位于底层的《白杨应用支撑平台》透明地封装。使得系统具备极高的可移植性，可以同时运行在不同的软硬件环境中。以期在简化演示、测试、扩容、联调等任务的同时，降低特定硬件平台对产品的粘度，为用户提供更丰富的平台选择。
- ★ **安全性：**由于 ERP 系统中记录着大量敏感数据，保证这些数据不被非法取阅和篡改对企业具有重大意义。为此，系统能够对用户进行安全的鉴权和基于 ACL 的细粒度访问控制：系统能够有效地对用户的身份进行识别，管理和行为审查。数据通讯的可靠性、保密性和完整性也以业界标准的强加密方式加以保证。除此之外，包含结构化数据、配置项和附件在内的所有用户数据均以整库的形式，以密文存储在我方自主研发的分布式虚拟文件系统中，全面保障数据安全。
- ★ **易用性：**以 App、桌面版 Web UI、移动版 Web UI、微信企业号集成等形式为用户提供简单易用、清晰明了、兼容各主流平台的图形化管理界面，同时为用户提供了桌面快捷方式、UI 主题等丰富的自定义机制。并通过语言包机制支持多种语言文字的国际化界面。
- ★ **易维护性：**在系统正常运行的过程中，不需要或仅需要很少人工干预和维护工作（fire and forget）。做到了故障发生时，在秒级自动检测并执行故障转移动作；故障修复后，自动感知并完成故障恢复操作。
- ★ **灾难恢复与异地容错：**系统提供了在线备份和恢复能力，使用户免于由自然灾害、硬件

故障、系统崩溃等灾难事故带来的数据丢失。并为对可用性和可靠性要求较高的组织提供提供了异地容灾解决方案，以及完备的强一致多活 IDC HAC 集群支持。

- ★ **信息交互能力：**根据用户要求为财务软件、办公软件等第三方应用提供特定格式的数据输出能力。提供完整易用的 WebAPI 接口，用以帮助客户完成各种二次开发项目。

注 1：以上单点性能指标均可在当前主流入门级（2-3 万人民币/台）商用 PC Server 上实现。通过升级硬件配置（纵向扩展），还可进一步线性地提升系统的单点性能表现。

2. 系统总体架构

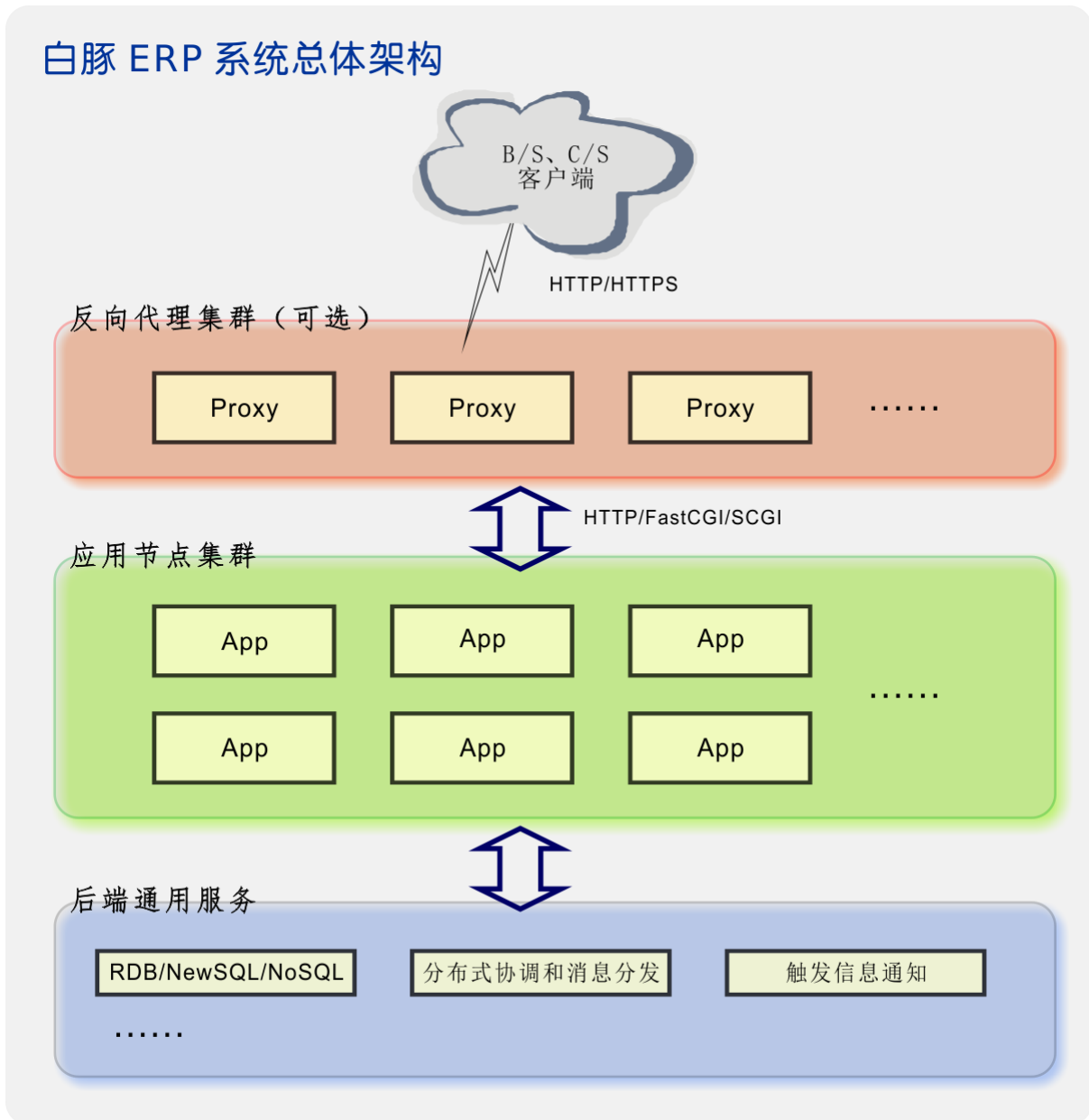


图 11

如图 11 所示，白豚 ERP 系统主要分为后端服务、App 集群以及 Web 客户端等部分。以下依循自底向上的顺序依次进行讨论。

2.1 通用后端服务

为应用提供支持的后端服务通用性较强，也多存在成熟的现有产品。短期内基本可使用现成

的第三方产品和服务来实现这些功能:

2.1.1 结构化存储（RDB、NewSQL、NoSQL）服务

数据库作为结构化数据的存储中心，存储所有与业务相关的数据，为所有数据访问提供 ACID 保证，并提供事务（Transaction）、索引和查询优化和等高级特性。系统通过《[白杨应用支撑平台](#)》中提供的 DBC、查询对象以及查询分析引擎等组件实现了跨数据模型（SQL / NoSQL / NewSQL）透明支持不同种类底层 DBMS 的目标，允许服务随着业务的发展从单点、主备到分布式集群逐步完成横向扩展的需求。

此外，还可基于由《[白杨应用支撑平台](#)》提供的，支持实时（on-the-fly）数据压缩和强加密的虚拟文件系统（VFS）来为数据提供整库级别的强加密保护。

白豚云 SaaS 平台基于由我方自主研发的，支持整库强加密的多活 IDC 6 副本分布式云端数据库系统，用户的每笔数据均在不同 IDC 中拥有多达 6 重同步副本。故其数据可靠性高达“17 个 9”（99.99999999999999%），同时还提供最近 90 天内的全量数据快照备份（7+12），可方便实现数据回溯。

2.1.2 分布式协调和消息分发服务（可选）

分布式协调服务负责实现服务发现、服务选举、故障检测以及分布式锁等职能。消息分发服务则实现了消息路由和分发算法。这些服务对于任何分布式系统而言，都属于核心基础组件，算法相对复杂，实现难度高。分布式白豚 ERP 集群系统使用《[白杨应用支撑平台](#)》中提供的现有 BYPSS 服务来实现上述功能，详见：5.2.3 白杨消息端口交换服务（BYPSS）。

受益于我方自主研发的 BYPSS 分布式协调服务，白豚云 SaaS 平台实现了强一致、高可用、高性能的多活 IDC 分布式集群，可为用户提供高达“6 个 9”（99.9999%）的可用性保障，承诺平均每年的不可用（故障）时间不超过 31.5 秒。专利的先进多活 IDC 集群架构和秒级的故障检测和故障转移能力使我们为用户提供了超越电信级“5 个 9”（99.999%，每年不可用时间 5.25 分钟）的极高可用性保证。

白豚 SaaS 平台中已包含此服务，但在用户独立部署时，此服务可选。若用户未部署此服务，则白豚 ERP 系统将无法部署为分布式的高可用集群（HAC），意即：存在单点故障。

2.1.3 触发消息通知（可选）

触发类信息通知主要包括触发邮件和触发短信。通过选用优质的电信运营商，以会员身份加入国内和国际上，各 ESP 组织的白名单。同时实现 SPF、DKIM、DMARC 等技术规范，并加入相应的管理组织等方式，白豚 ERP 系统实现了极高的邮件、短信送达率。

与此同时，我们的邮件、短信发送信令均经由 **SSL/TLS** 连接实现安全强加密保护，有效防止用户信息泄露。

白豚 SaaS 平台中已包含此服务，但在用户独立部署时，此服务可选。若用户未部署此服务，则白豚 ERP 系统将无法发送必要的触发通知（如： workflow 待审批等）给相应员工。

2.2 App 服务器集群

App 服务节点是系统的核心部分，依托底层提供各通用服务，实现了白豚 ERP 系统中的所有业务逻辑。由于使用了先进的 **nano-SOA** 架构（详见：[5.2 nano-SOA 架构](#)），在同时保留了 **SOA** 架构高内聚、低耦合优点的前提下，App 集群中的所有节点得以以完全对等的形式来实现。

得益于 **nano-SOA** 架构，我们：

- ★ 只要维护一套 App 节点集群即可实现负载均衡和强一致、多活 IDC 级别的高可用，不必为应用中的每个服务部署和维护单独的高可用集群。运维工作量大大减轻，服务器、存储等资源的购买或租用成本也大大降低。
- ★ 避免了大量节点间通信，时间和空间效率均大大提升。可以使用更少的硬件资源和运维投入来服务更多用户。并实现更短的页面加载时间，更快速的 UI 响应，以提供更优质的客户体验。
- ★ 完全对等的集群模式不仅消除了大部分节点间通信的需求，而且避免了对高可用、高可靠的消息中间件之依赖。甚至可以取消中心化的消息分发服务，而使用点对点的去中心化消息分发模式。这从根本上解决了传统 **SOA** 模式中，由消息中间件引入的各种可靠性、可用性和性能瓶颈。

基于 **nano-SOA** 的设计思想，我们将一个 App 节点划分为基础平台和可插拔业务模块两大部分。其中基础平台提供底层功能封装、**Web** 伺服框架、数据库连接、存储管理、服务管理、配置管理、插件管理、角色管理、用户管理、会话管理、部门管理、时区和时段管理、成本中心管理等各种基本服务。而诸如财务报销、考勤管理、工资管理、绩效管理、招聘管理、培训管理、文档管理、信息发布等具体业务逻辑则由不同的业务模块分别实现。

有关 App 服务节点设计的详细信息，请参考：[3. 系统总体设计](#)。

2.3 反向代理服务（可选）

位于三层构架中最外层的反向代理服务器负责接受用户的接入请求，在实际应用中，反向代理服务通常至少还要完成以下列表中的一部分任务：

连接管理：分别维护客户端和应用服务器的连接池，管理并关闭已超时的长连接。

攻击检测和安全隔离：由于反向代理服务无需完成任何动态页面生成任务，所有与业务逻辑相关的请求都转发至后端应用服务器处理。因此反向代理服务几乎不会被应用程序设计或后端数据漏洞所影响。反向代理的安全性和可靠性通常仅取决于产品本身。在应用服务的前端部署反向代理服务器可以有效地在后端应用和远程用户间建立起一套可靠的安全隔离和攻击检测机制。

如果需要的话，还可以通过在外网、反向代理、后端应用和数据库等边界位置添加额外的硬件防火墙等网络隔离设备来实现更高的安全性保证。

负载均衡：通常使用轮转（Round Robin）或最少连接数优先等策略完成基于客户请求的负载均衡；也可以使用 SSI 将一个客户请求拆分成若干并行计算部分分别提交到多个应用服务器。

分布式的 cache 加速（CDN）：可以将反向代理分组部署在距离热点地区地理位置较近的网络边界上。通过在位于客户较近的位置提供缓冲服务来加速网络应用。

静态文件伺服：当收到静态文件请求时，直接返回该文件而无需将该请求提交至后端应用服务器。

动态响应缓存：对一段时间内不会发生改变的动态生成响应进行缓存，避免后端应用服务器频繁执行重复查询和计算。

数据压缩传输：为返回的数据启用 GZIP/ZLIB 压缩算法以节约带宽。

数据加密保护（SSL Offloading）：为与客户端的通信启用 SSL/TLS 加密保护。

容错：跟踪后端应用服务器的健康状况，避免将请求调度到发生故障的服务器。

用户鉴权：完成用户登陆和会话建立等工作。

URL 别名：对外建立统一的 url 别名信息，屏蔽真实位置。

应用混搭：通过 SSI 和 URL 映射技术将不同的 web 应用混搭在一起。

协议转换：为使用 SCGI 和 FastCGI 的后端应用提供协议转换服务。

目前比较有名的反向代理服务包括：Apache httpd+mod_proxy / IIS+ARR / Squid / Apache Traffic Server / Nginx / Cherokee / Lighttpd / HAProxy 以及 Varnish 等等。其中 HAProxy 无论在效率、性能、功能还是安全性方面，均为其中的佼佼者。

白豚 ERP 系统支持上述常见反向代理产品，其 SaaS 平台中已包含此服务，但在用户独立部署时，此服务可选。若用户未部署此服务，则白豚 ERP 系统将可能无法支持上述部分功能。

2.4 客户端

白豚 ERP 系统提供基于 HTML5 技术的新一代 B/S 架构 Web 客户端，精心设计的新型客户端支持响应式布局的国际化多主题界面。可兼容包括 IE、Firefox、Chrome、Opera、Safari、iOS（iPhone/iPad）、Android 以及 BlackBerry 在内的各种浏览器和移动平台。可轻松实现跨平台的移动办公。

客户端是遵循功能与界面完全分离、以及 one page one application 的原则来进行设计的。这不但为用户提供了更快速的响应，以及更接近 C/S 架构的用户体验，开放式的 RESTful WebAPI 接口也使得二次开发工作变得简单易行。

与此同时，为了适应更多复杂的办公需求，我们还为用户提供了专门为手机等小屏幕设备优化，并且方便嵌入微信企业号的 Mobile Web 端，以及 Android 和 iOS 版本的专用 App 客户端。

3. 系统总体设计

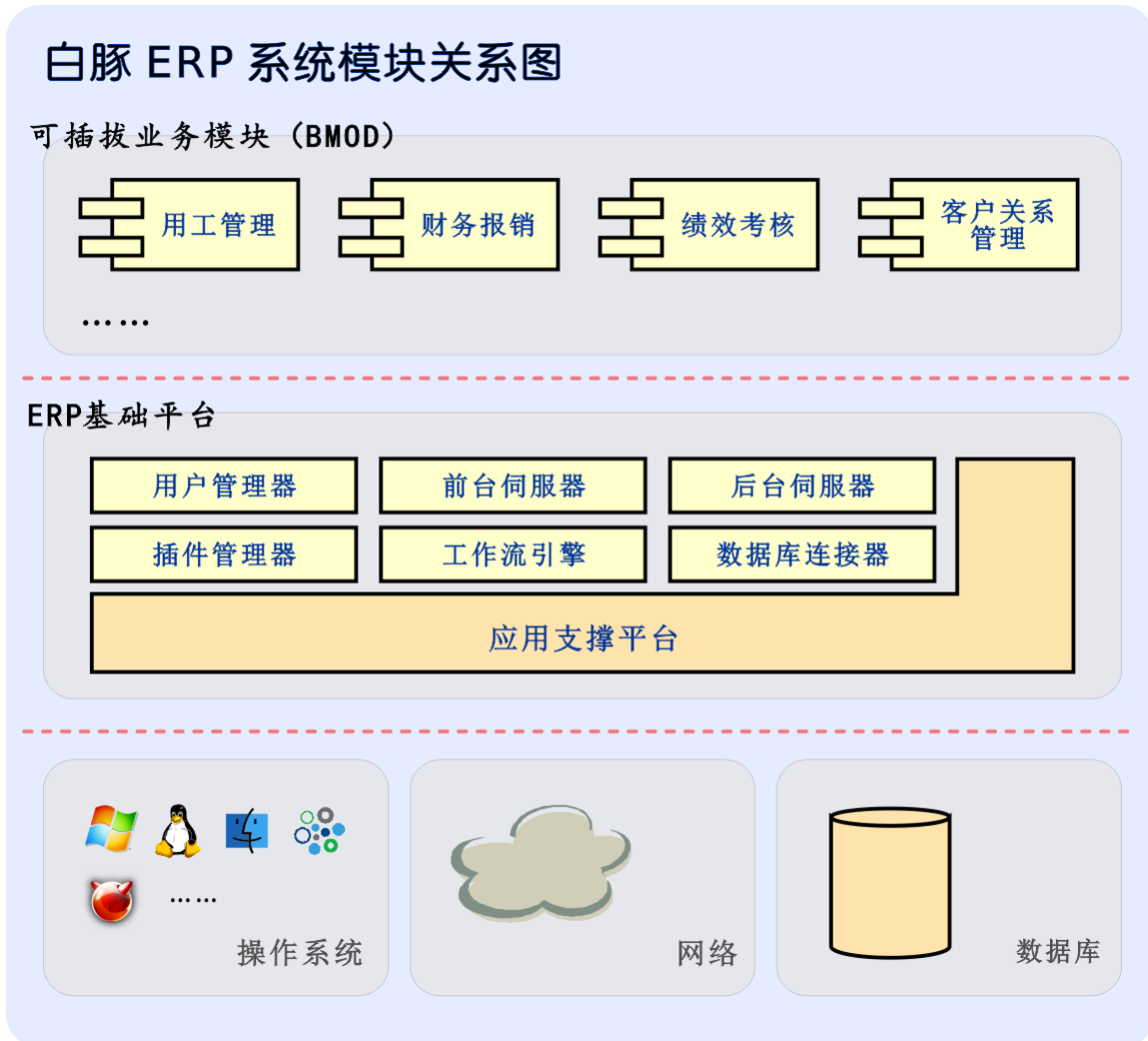


图 12

如图 12 所示，白豚 ERP 系统主要可分为 ERP 基础平台和可插拔业务模块两大层面，以下逐一讨论：

3.1 ERP 基础平台

基础平台向下封装所有与操作系统、网络和数据访问相关的底层操作；向上则作为各个业务模块的运行支撑环境，并为这些业务模块提供用户管理、配置管理等基础服务。

3.1.1 应用支撑平台

应用支撑平台是整个 ERP 系统的基石，是白豚 ERP 与底层操作系统及硬件平台间的重要接口，应用支撑平台在完整封装操作系统功能的同时，还提供各种通用工具。应用支撑平台是高质量、低成本快速开发及支持跨平台应用的关键组件。

3.1.1.1 应用支撑平台功能描述

应用支撑平台为 ERP 系统提供了众多基础功能，包括：

- ★ 跨平台底层支持：封装所有与操作系统相关的操作，如：信号量、原子操作、共享内存 / 文件映射、进程、线程、网络操作（Socket）、文件管理、服务控制、注册表访问、进程间通讯、服务器框架、异步 IO 框架等等。是 ERP 系统实现跨平台、多平台的关键组件。
- ★ 通用功能：包括用户权限管理、基于 PKI 体系结构的强加密算法、常用网络协议、二进制和字符集编码转换、自动化脚本引擎、表单处理、数据压缩、任务管理、日志记录、LRU Cache 容器等等。
- ★ 平台无关的高效 Web 扩展框架。
- ★ 平台无关的复杂数据结构表示以及支持实时压缩和强加密的虚拟文件系统。
- ★ nano-SOA 架构基础设置，包括：平台无关的功能插件框架；灵活，支持实时（on-the-fly）数据压缩和强加密保护的数据库连接器（DBC）插件；通用 API 注册和分发机制（API Nexus）；ETL 辅助工具；基于平台查询分析引擎的自定义高级查询辅助工具等等；以及 BYPSS 分布式协作与消息分发服务等等。
- ★ 平台无关的国际化支持：为用户提供一个平台无关的多语言、国际化工作环境。

等等，共上千种功能组件和通用服务。应用支撑平台完整封装了几乎所有与底层系统相关的功能，以及很多常用的算法、框架和功能组件。

3.1.1.2 应用支撑平台设计框架

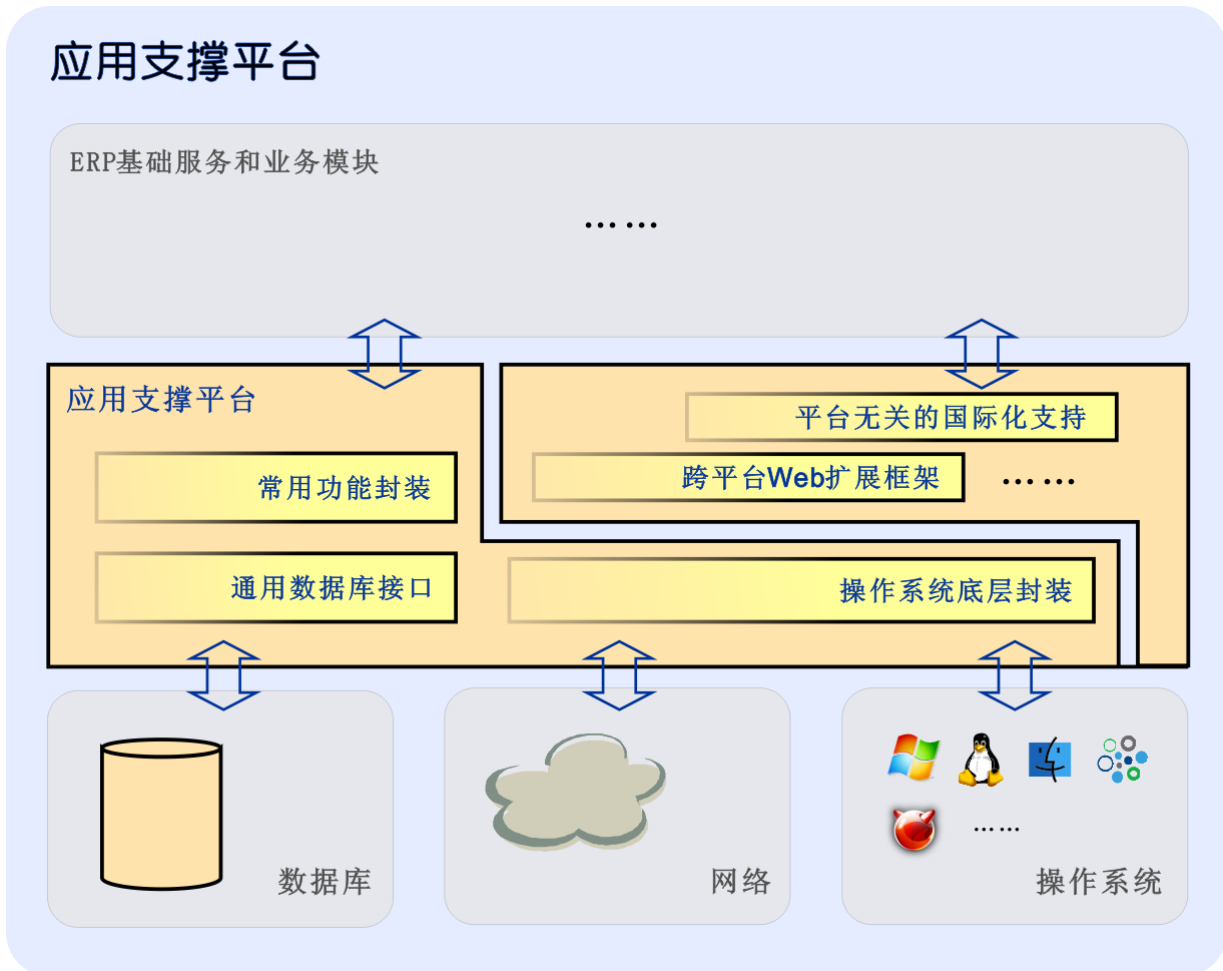


图 13

如图 13 所示，应用支撑平台位于操作系统、网络、数据库系统等底层环境及其它 ERP 基础服务和业务模块等高层应用之间，是白豚 ERP 的构造基础。向下封装各类底层的、操作系统和硬件平台相关的通用功能，向上为白豚 ERP 系统提供稳定强大、平台无关的通用服务和常用组件。

应用支撑平台为我方自主开发的跨平台应用支撑框架：跨平台框架统一封装了进程；线程；线程本地存储；信号量；共享内存；命名管道；文件映射；原子操作；消息机制；网络访问；文件系统管理；系统时钟；任务计划；服务控制；动态库装载等各类操作系统功能。

同时提供了基于 PKI 体系架构的高强度数据加密；基于数字证书的 License 生成和验证；数据压缩算法；通用句柄操作；数据编码转换；字体访问；字符集转换；数据库访问；报表处理；图像生成；用户鉴权和访问控制机制；音频、视频输入输出及编解码；嵌入式自动化脚本引擎等各类通用功能。

支撑平台从一开始就是专为对负载和可靠性要求很高的企业及互联网行业而设计的。其中的每个组件都会在保证可靠性的前提下，使用对当前底层平台来说，最高效的方式予以实现。

例如：支撑平台中的网络 IO 和 Web 扩展框架总是会尽可能地利用底层系统所支持的最优模式来完成网络访问：在 Windows 上使用 IOCP；在 Linux 上使用 epoll；在 FreeBSD/NetBSD/OpenBSD 等 BSD 环境中使用 kqueue；在 Solaris 上使用 event ports；在 HP-UX 上使用 /dev/poll、在 IBM AIX 上使用 pollset、而在其它 POSIX 环境下使用 POSIX AIO + Realtime Signal 等等。

再比如：根据实际情况直接使用内联汇编或编译器 Intrinsic 对原子量操作、内存屏障以及字节序交换等热点功能进行优化。并直接使用这些已被极度优化的组件构建快速互斥量、自旋锁、以及快速信号量等组件。

应用支撑平台完全使用高效强大的 C++ 语言和少量汇编代码进行开发，由一组跨平台的功能组件构成，全面支持当今流行的各种软、硬件环境^[1]。目前已被广泛应用于多个大型项目^[2]，其高效性、可靠性和成熟度已被广泛验证。

使用统一的应用支撑平台和用户界面组件构建软件可为用户带来很多显而易见的好处，例如：

- ★ 更高的产品质量：大量重用经过岁月积累和时间考验的组件显著提高了产品的稳定性。与此同时，重用经过反复调优的组件也有助于提高性能。
- ★ 跨平台能力和国际化支持：统一应用支撑平台和统一用户界面组件库封装了几乎所有系统相关的操作，并提供了很强的跨平台支持能力，以它们为基础开发的应用因此具有平台无关特性。典型情况下，跨平台移植的全部工作仅仅是对应用的简单重编译。同时，国际化作为统一用户界面组件的一个重要附加功能，为应用提供了平台无关的国际化支持能力。
- ★ 极高的运行时效率：100% 高效 C++/汇编编码，无论是时间效率还是空间效率上都没有任何妥协。

应用支撑平台包含繁多的功能和复杂的设计。其复杂度甚至还要远超白豚 ERP 系统本身。作为一个完善、独立的产品，应用支撑平台的具体细节已经超出了本文的讨论范围，关于和应用支撑平台相关的进一步信息，请参考：《[应用支撑平台技术白皮书](http://baiy.cn/doc/asp_whitepaper.pdf)》（http://baiy.cn/doc/asp_whitepaper.pdf）。

注 1：我们当前支持的操作系统主要包括：

- Windows 98/ME, Windows NT4/2000/XP/2k3/Vista/2k8/Win7/2k8r2。
- Linux、FreeBSD、NetBSD、IBM AIX、HP-UX、Solaris、MAC OS X 系统及众多 Un*x/POSIX 系统。
- vxWorks, Nuclear, QNX, SMX, DOS, Windows CE (Windows Mobile), NanoGUI、eCos、RTEMS、Android、iOS 等嵌入式系统。

当前支持的主要硬件平台包括：x86/x64、ARM、RISC-V、IA64、MIPS、POWER、SPARC 等。

注 2：应用支撑平台的典型客户中，包括了兴业银行（China CIB）、中石油（CNPC）、华安保险（Sinosafe Insurance）、淘宝网（taobao.com）、烟台万华集团、法国兴业银行（SOCIETE GENERALE）、德尔福汽车（Delphi）、美联航（United Airlines）、GE（美国通用电气）、贝塔斯曼（Bertelsmann）、埃森哲（Accenture）等各大企业。

3.1.2 基础服务群

ERP 基础平台中的基础服务群构建于应用支撑平台之上。以应用支撑平台为依托，实现了诸如：用户管理、配置管理、插件管理、数据库连接管理以及前台/后台 Web 服务器等各类基础服务。此外，基础服务群在与应用支撑平台一起为各个业务模块提供运行环境的同时，也负责对正在运行的各业务模块进行管理。

3.1.2.1 基础服务群功能描述

基础服务群实现以下功能：

- ★ 鉴于配置数据与业务数据相分离的设计思想，需要专门提供可靠的配置信息存储服务。该服务通常由带有自动备份和自动版本回退机制、使用强加密和实时数据压缩的本地 VFS 来实现。但在需要实现横向扩展时，也可以基于远程分布式存储。
- ★ 实现用户管理、组管理、权限管理、插件管理等各种管理组件。
- ★ 实现通用工作流引擎、附件服务器、待办事务服务器等通用功能组件。
- ★ 在插件机制的基础上，实现业务模块管理器。完成对各业务模块的启动、停止、加载、卸载、启用、禁用等各种调度和管理工作。与此同时，由于不同业务模块间可能存在依赖关系，因此业务模块管理器内还提供了一套可靠的模块间通信机制。
- ★ 为业务数据实现后端 DBMS 的访问和连接管理机制。并为其提供对用户透明的实时数据压缩和强加密保护、基于 Revision 的分布式乐观锁、高级自定义查询等功能。
- ★ 分别实现用于前台访问和后台管理界面的 Web 服务器框架。Web 框架应满足设计要求中定义的各项指标，并允许各业务模块灵活地进行功能扩展。
- ★ 提供非结构化大对象存储服务，以及支持多语种分词的全文搜索服务等通用服务。
- ★ 基于 BYPSS 服务实现节点间消息分发和集群协作功能。
- ★ 实现自动更新服务器，依赖后端 CMDb 服务实现配置和应用程序自动更新功能。
- ★ 实现分布式、可扩展的用户管理和会话管理机制。
- ★ 实现分布式的即时消息分发和推送机制。
- ★ 实现待办事务服务器，为用户提供待办事务计划和提醒。

- ✦ 实现计划任务伺服器，提供任务预约执行功能。
- ✦ 与第三方平台对接，实现触发消息通知功能。

3.1.2.2 基础服务群设计框架

由于需要支持高并发环境，前台 Web 伺服器框架使用高效的异步 IO 架构来实现，其工作模型如下图所示：

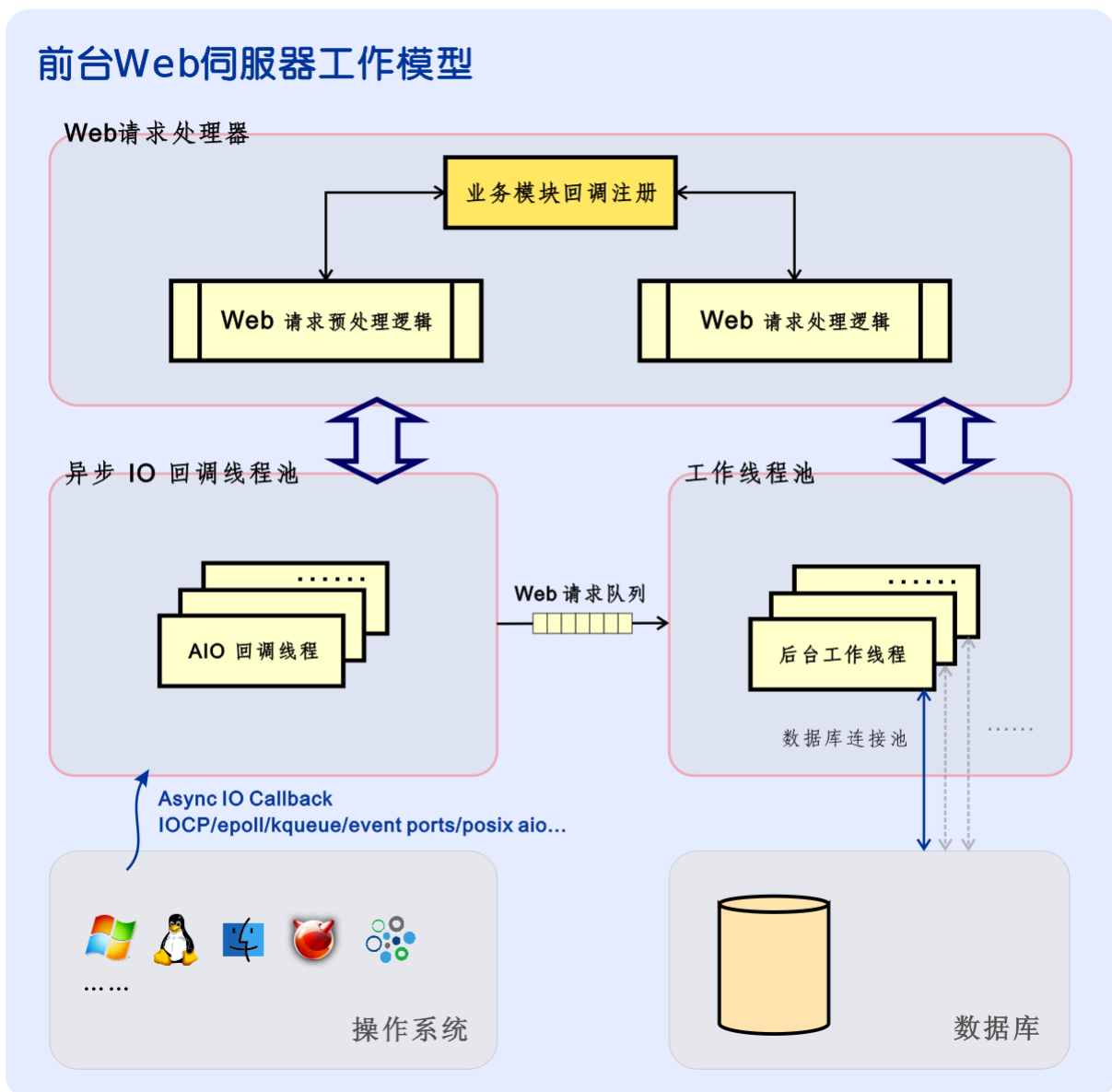


图 14

如图 14 所示，前台 Web 伺服器工作流程如下：

- ★ 当一个 Web 请求到达后，底层操作系统通过 IOCP、epoll、kqueue、event ports、real time signal (posix aio)、/dev/poll、pollset 等各类与具体平台紧密相关的 IO 完成（或 IO 就绪）回调机制通知异步 Web 扩展框架的 AIO 回调线程，对这个已到达的 Web 请求进行处理。
- ★ 在 AIO 回调池中的工作线程接收到一个已到达的请求后，首先尝试通过 Web 请求处理器对该请求进行预处理。在预处理过程中，会先尝试根据 URI 匹配结果来回调相应的业务模块，若无匹配或业务模块放弃处理，则会继续尝试使用位于本地的高速缓存来避免成本较高的磁盘访问及数据库查询。如果业务模块的预处理请求成功，或本地缓存命中，则直接将缓存中的结果（仍然以异步 IO 的方式）返回客户端，并结束本次请求。
- ★ 如果指定的 Web 请求要求查询的数据无法被本地缓存命中，或者这个请求需要数据库写入操作，则该请求将被 AIO 回调线程追加到指定的队列中，等待工作线程池中的某个空闲线程对其进行进一步处理。
- ★ 工作线程池中的每个线程都分别维护着两条长连接：一条与底层到数据库服务相连，另一条则连接到分布式缓存（memcached）网络（可选，图中省略）。通过让每个工作线程维护属于自己的长连接，后台工作线程池实现了数据库和分布式缓存服务的长连接池机制。长连接（Keep-Alive）机制通过为不同的请求重复使用同一条网络连接大大提高了应用程序和网络 IO 的处理效率。
- ★ 工作线程在 Web 请求队列上等待新的请求到达。在从队列中取出一个新的请求后，工作线程将该请求和可用的数据库连接提交给 Web 请求处理器进行处理。
- ★ 请求处理器首先使用工作线程指定的数据库和分布式缓存服务长连接处理指定的 Web 请求（如果存在匹配的 URI 规则，则将其送交相应的业务模块进行处理）。然后将处理结果提交给底层 Web 扩展框架，将结果通过异步 IO 的方式返回给客户端。

应注意到：前台 Web 伺服器中的 Web 请求处理器及 Web 请求预处理均允许业务模块进行注册。所有已正确注册的业务模块都可以通过 URI 匹配等方式参与 Web 请求的处理和预处理机制。这就使得业务模块可以对 B/S 界面和 Web API 进行扩展，从而向用户暴露自己的外部接口。

需要进一步说明的是，与 epoll/kqueue/event ports 等相位触发的通知机制不同，对于 Windows IOCP 和 POSIX AIO Realtime Signal 这类边缘触发的 AIO 完成事件通知机制，为了避免操作系统底层 IO 完成队列（或实时信号队列）过长或溢出导致的内存缓冲区被长时间锁定在非分页内存池，在上述系统内的 AIO 回调方式实际上是由两个独立的线程池和一个 AIO 完成事件队列组成的：一个线程池专门负责不间断地等待系统 AIO 完成队列中到达的事件，并将其提交到一个内部的 AIO 完成队列中（该队列工作在用户模式，具有用户可控的弹性尺寸，并且不会锁定内存）；与此同时另一个线程池等待在这个内部 AIO 完成队列上，并且处理不断到达该队列的 AIO 完成事件。这样的设计降低了操作系统的工作负担，避免了在极端情况下可能出现的消息丢失和内存泄露，同时也可以帮助操作系统更好地使用和管理非分页内存池。

在这一模型中，应用程序可将需要传输的内存地址直接提交给底层硬件，硬件通过 DMA 直接在这个内存位置完成 IO 操作，这就实现了内存零拷贝。在完成 IO 操作后，硬件通过触发中断

通知操作系统，并由操作系统回调应用程序。不但如此，由于可以向底层驱动并发地提交多个 IO 请求，这就使操作系统和底层硬件有机会实现操作合并（例如：将多个消息合并到一个网络帧或磁盘 IO 请求中完成读写）以及最优化请求顺序（例如：调度磁头从最近的磁道开始读写请求）等等。

得益于上述高效异步 IO 模式以及使用汇编代码的精细手工优化，《[白杨应用支撑平台](#)》中的高并发 Web 伺服组件可提供极高的 IO 吞吐和并发能力。即使在一台 2011 出厂的双路 Intel 至强 5600 入门级 1U PC Server 上（当时售价 RMB 2 万以内），亦可达到单点支持千万量级超高并发的水平。与 IIS+asp.net / Apache+php / Nginx+php 以及相应的 Java / Python / RoR 等方案相比，基于汇编/C/C++的 Web 应用框架具备几何级数的性能优势。

相反，由于不需要考虑高并发和高负载等使用场景，后台管理 Web 伺服器使用较简单的多线程，每连接一线程的阻塞 IO 模式：

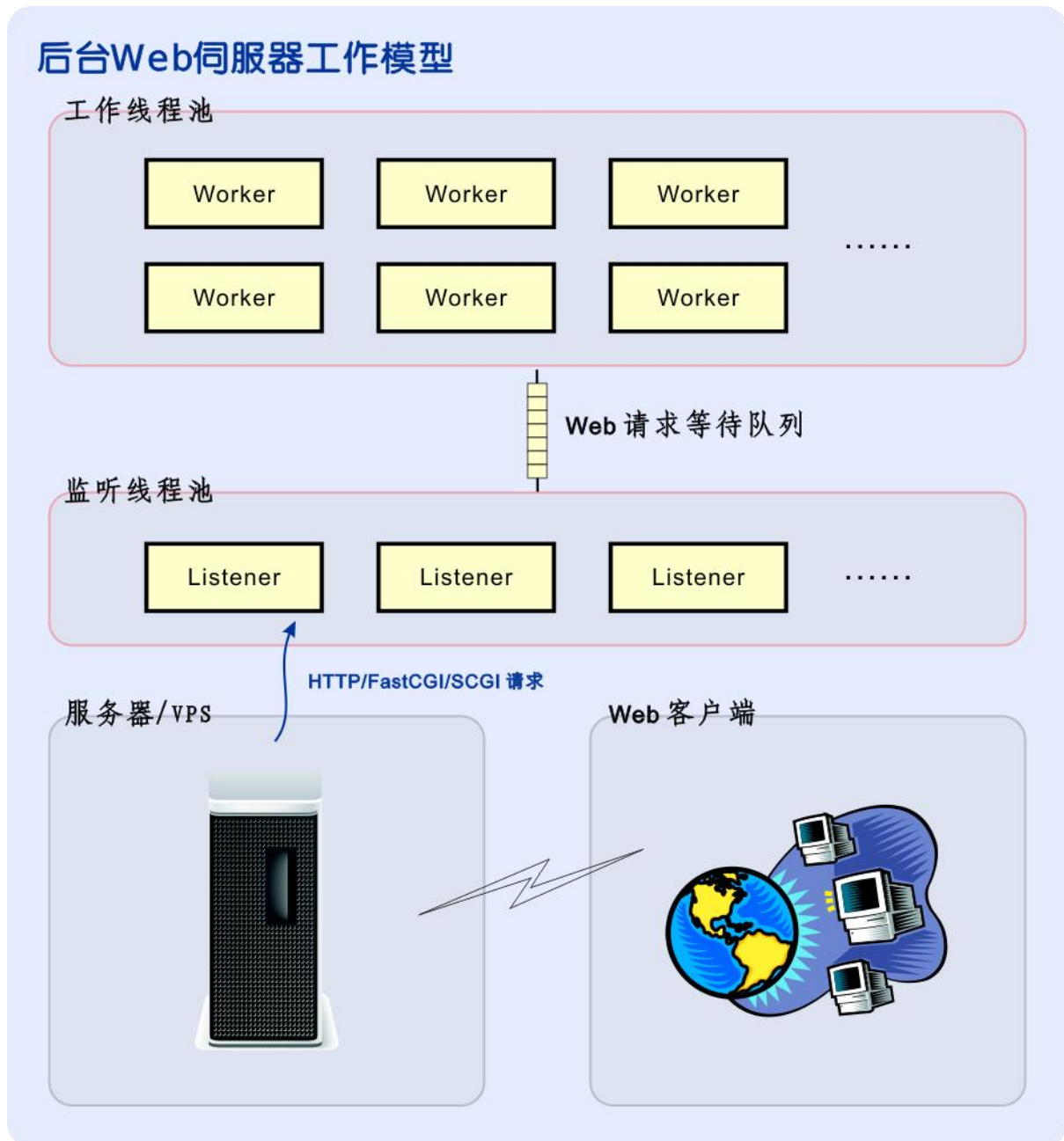


图 15

如图 15 所示，后台 Web 伺服器工作流程如下：

- ★ 监听线程池等待到达用户和会话管理服务器的 Web 请求，并完成与该请求相关的初始化工作。然后将其压入 Web 请求队列的尾端，并继续监听下一个请求。为了增强高负载模式下的并发性，这些工作由单独的线程池来完成。监听线程池中的线程数量可以由用户配置，也可以根据负载情况动态调整。
- ★ 工作线程池不断从等待队列的首部取出待处理的请求，将其提交给 Web 请求处理器。Web 请求处理器可视情况将其转交给相应的业务模块，最后返回处理结果。工作线程池

可以根据当前负载状况和用户配置的参数进行动态调整。在处于轻载或空闲时，工作线程的数量保持在用户指定的最小值。当负载持续加重时，为提供额外的并发处理能力，线程数量会逐渐增加，直到达到用户指定的最大值为止。当业务高峰结束，负载量持续下降时，应用框架会逐渐回收已空闲的工作线程，直到达到用户指定的最小值为止。工作线程的回收策略可以由用户指定。

除前台和后台 Web 伺服器以外，基础服务群中的其它组件均容易理解和实现，本文不再赘述。

3.2 可拔插业务模块

正如前文所述，为了提供灵活的定制能力和良好的功能扩展性，白豚 ERP 系统被特别划分成了基础平台和可插拔业务模块两个部分。

从平台的角度讲，每个可插拔业务模块都是一个运行在 ERP 基础平台上的插件。业务模块依托于基础平台提供的运行时环境，享受基础平台提供的服务，并服从基础平台的管理和调度命令。在必要时，多个业务模块间也可通过基础平台提供的服务实现相互通信。

从用户的角度讲，诸如人事管理、财务管理、库存管理等所有实际业务功能其实都是由各个业务模块分别实现的。一个完整的 ERP 系统就是由 ERP 基础平台加上用户选择的数个业务模块组合而成的。从这个角度看，基础平台好似一套空房间，而业务模块则是各种可以摆放在房间中的家居物件。

实际上，诸如“人事管理”、“财务管理”等业务大类也是由多个业务模块组合实现的。例如：“人事管理”相关的业务中，包含了“绩效考核”、“考勤管理”、“用工管理”、“薪资管理”、“招聘管理”、“岗位管理”等多个具体的业务模块。这样的细化粒度进一步增强了系统组合的灵活性，并且降低了专用模块定制化的成本。

3.2.1 可拔插业务模块功能描述

可拔插业务模块完成各类实际业务逻辑，可能包括人事管理、财务管理、CRM 管理、库存管理等各类通用业务。也可以为有特殊要求的客户定义实现一些特殊业务。

可插拔式的业务模块是基于 ERP 应用所固有的需求多样化本质而设计的。为说明可拔插业务模块的特点，以下列举几种 ERP 系统在与用户的实际需求对接时，可能碰到的典型场景：

- ★ 现有业务模块足以满足用户全部需求：这是最理想的情况，无需任何额外工作。用户仅需从现有业务模块集合中，挑选出需要的子集即可。
- ★ 现有业务模块覆盖了用户全部需求，但用户对某个（或某些）模块有特殊需求，即：用户对某个现有模块中的部分功能不满意。例如：用户需要使用“绩效考核”、“考勤管理”、

和“财务报销”三个已有模块，但希望考勤管理模块可连接某特定厂商的指纹考勤机，以便实时获取考勤数据。这时仅需为用户定制考勤模块即可。当然，用户也可以通过 WebAPI 接口自行添加相关的功能扩展。

- ★ 现有业务模块满足用户的部分需求。例如：白豚 ERP 现有“绩效考核”、“考勤管理”、和“财务总账”三个业务模块，这三个模块均可满足用户在这些方面的需求。但是该用户还需要一个尚未被实现的“品质管理”业务。这时可有以下几种解决方案：
 - 保持现有部分不变，为用户新增一个“品质管理”业务模块。在这种情形下，重用现有的基础平台和已有业务模块既可为用户节约大量定制成本，又能最大程度地保证产品品质。
 - 保持现有部分不变，用户通过 WebAPI 接口自行完成二次开发，以满足自己的特殊需求。这种方式适用于需求频繁变更，并且自己拥有技术团队的客户。
- ★ 其它所有可能的场景均可由以上几种案例相互组合而得到。

通过以上对各种经典案例的讨论，可以明显地看出：可插拔式业务模块配合 SOA 风格的开放式 WebAPI 所带来极其优异的灵活性。这也是可插拔业务模块的主要功能及其价值所在。

3.2.2 可插拔业务模块设计框架

可插拔业务模块与数据库连接器插件、统计分析插件一样，遵循标准的插件接口，并通过专用的管理器进行调度和管理。

4. 安全性和可靠性设计

白豚 ERP 系统通过在线复制（分布式文件系统）、磁盘快照、自动备份和自动恢复、数据强加密、CHAP 鉴权、TLS（HTTPS）安全传输、以及高可用集群（HAC）等多种手段，从数据存储、网络传输、节点容错等各个方面来保证服务的安全性、可靠性、稳定性和健壮性。

4.1 服务可用性和数据可靠性

正如前文所述，基于我方独有的专利技术，白豚云服务以强一致，抗脑裂（Split Brain）的高可用（HAC）+ 高性能（HPC）分布式集群的方式，**同时部署于相互之间至少相隔数十公里的多座 IDC 中。**

其中每座 IDC 均为国际顶级的多线 BGP 机房，电力供应、网络接入、制冷设备、服务器节点等均有相互隔离的冗余配置，并配有大规模 UPS 和备份发电机组等设施。机房建筑为高防灾等级的混凝土结构，并配有完备的自动灭火设施等防灾装置。

除此之外，机房还配备了无死角的监控系统和高安全等级的门禁设施，并有专业保安团队和运维工程师团队 7x24 小时值守。完全符合**国际电信联盟制定的最高等级机房标准(ANSI/TIA-942, Tire IV)**，机房本身的高等级为系统整体可靠性保障打下了坚实的基础。

4.1.1 在线数据可靠性

得益于我方独有的**跨 IDC 的分布式 6 副本技术**，白豚 ERP 系统确保包括业务数据、配置数据、附件数据、版本控制数据等在内的所有在线数据均以强一致多副本的形式，被同时存储在多座至少相互间隔数十公里的高等级 IDC 机房中。**其数据可靠性高达 17 个 9 (99.99999999999999%)**，即使整座 IDC 被灾难性地完全摧毁，也不会丢失任何用户数据！

在线数据可靠性按服务周期统计，一个服务周期为一个自然月，如不满一个月不计算为一个服务周期。在线数据可靠性计算公式为： $\text{每服务周期完好数据的云租户} / (\text{每服务周期完好数据的云租户} + \text{每服务周期丢失数据的云租户})$ 。

在线数据的可靠性保证为：不低于 99.99999999999999%。即，每 10000000000000000（十亿亿）个白豚云租户中，每月数据不丢失的概率为 99.99999999999999%，即其中每月最多只有 1 个租户可能发生数据丢失。

4.1.2 备份计划及离线数据可靠性

白豚云平台每日对其下所有用户数据自动完成**全量快照备份**，平台保留最近 7 天的快照备份，以及最近 90 天内，每周末的磁盘快照（7+12）。

快照数据基于**可靠的双重 EC 算法**，可有效防止单点故障带来的数据损失。与此同时，每次快照备份都会按照备份计划中事先设置好的规则，**交替地保存在不同的高等级 IDC 机房中**，有效降低单个 IDC 发生毁灭级灾难时的损失。

离线备份数据可靠性按服务周期统计，一个服务周期为一个自然月，如不满一个月不计算为一个服务周期。离线备份数据可靠性计算公式为：每服务周期完好数据的云租户/(每服务周期完好数据的云租户+每服务周期丢失数据的云租户)。

单一快照备份的可靠性保证为：不低于 **99.99999999%（10 个 9）**。即，每 10000000000（百亿）个白豚云租户中，每月数据不丢失的概率为 99.99999999%，即其中每月最多只有 1 个租户可能发生数据丢失。

分别存储于两座不同 IDC 内的任意两份快照数据同时发生丢失的可靠性保证为：不低于 **99.9999999999999999%（19 个 9）**。即，每 100000000000000000000（千亿亿）个白豚云租户中，每月最多只有 1 个租户可能发生快照备份数据完全丢失。

除此之外，白豚 ERP 平台还提供了**无限量的多版本控制技术**：平台为系统中的员工、部门、站点、岗位、休假信息、薪资帐套等在内的所有关键数据提供**历史版本无限回溯**支持。例如：自某一数据对象创建后，平台在其每一次修改后，都会永久保留相关增量数据。管理人员可以随时审计这些数据的修改记录，查看其修改时间、操作者、相关工作流和变更操作摘要等关键信息，并将数据回溯到任意指定的修订版。除了可用于历史回顾和审计外，回溯功能还可有效应对用户对数据的恶意修改或意外破坏。

上述无限多版本技术的增量数据同样保存在我方独有的，**跨 IDC 的分布式 6 副本存储**中，其数据可靠性承诺详见上一小节：4.1.1 在线数据可靠性。与此同时，这些增量数据也会与其它数据一起包含在每日的全量快照备份中，被保存于前文所述的“7+12”**离线备份集**内。

4.1.3 服务可用性

得益于我方自主研发的强一致、抗脑裂（Split Brain）多活 IDC **高可用+高性能分布式集群（HAC+HPC）专利架构**，加之集群被同时部署于相互之间至少相隔数十公里的，至少三座高等级 IDC 中。因此**白豚云的整体服务可用性可达 6 个 9（99.9999%）**。确保即使整座 IDC 由于市政施工、自然灾害、人为失误等原因下线，也可在秒级完成故障转移，自动切换到其它可用 IDC，不间断地为用户提供服务。

精心设计的故障检测、服务发现、服务选举、分布式锁，以及消息分发、消息路由等分布式

协调和数据同步算法提供的强一致性保证确保了跨 IDC 的故障转移（Failove）和故障恢复（Failback）也不会导致数据一致性问题。

白豚云可用性计算公式为：每服务周期单个云租户所有可用时间/（每服务周期单个云租户所有可用时间+每服务周期单个云租户所有不可用时间）。即：每年不可用时间不超过 31.536 秒。其中：

1. 白豚云可用性按服务周期统计，一个服务周期为一个自然年，如不满一年不计算为一个服务周期，统计的业务单元为单个白豚云实例（集群），时间单位为秒钟。

2. 不可用时间：白豚云所提供的服务在连续的 15 秒钟或更长时间不可使用方计为不可用时间，不可使用的服务时间低于 15 秒钟的，不计入不可用时间，白豚云实例不可用时间不包括日常系统维护时间，以及由用户原因、第三方原因或不可抗力导致的不可用时间。

带强一致保证的多活 IDC 技术是现代高性能和高可用集群的关键技术，也是业界公认的主要难点。作为实例：2015 年 8 月 20 日 Google GCE 服务中断 12 小时并永久丢失数据；2015 年 5 月 27 日、2016 年 7 月 22 及 2019 年 12 月 5 日支付宝多次中断数小时；以及 2013 年 7 月 22 日微信服务中断数小时等重大事故均属于产品未实现多活 IDC 架构，单个 IDC 故障导致服务全面下线的惨痛案例。

在上述方面，我方拥有超过十年的丰富积累，受多项国家和国际发明专利保护的分布式架构和算法。得益于领先的强一致、高可用、高性能分布式集群算法和架构，我们在白豚 ERP 平台上实现了真正的多活 IDC 架构，为客户提供了无以伦比的数据可靠性和服务可用性保证。

4.2 安全性设计

ERP 平台中包含大量企业内部敏感数据。为此，除了做好机房安保等管理、配置网络边界安全规则等基本工作外，系统自身也应该包含完备的安全防护机制。白豚 ERP 平台从设计之初即将安全性放在首位，下面从数据安全性、鉴权和授权、传输安全性、系统安全性等各个方面逐一介绍白豚 ERP 产品的安全性相关设计。

4.2.1 数据安全性

白豚 ERP 使用基于我方自主研发，支持实时（on-the-fly）数据压缩和强加密的虚拟文件系统（VFS），对数据库和所有配置信息提供整库级别的强加密保护，支持包括 AES、BlowFish、SERPENT 在内的数十种业界认可的强加密算法。

普通的，以数据库记录为单位的加密通常会将相同的明文转换为同样的密文。例如：若数据库内的用户表中，有一个名为 Name 的字段用来存放员工姓名，并且明文“张三”被加密后的密文为“ABCD”，则对于逐行（逐字段）加密的算法来说，表中每次出现名为“张三”的员工，都

将被加密为“ABCD”（类 ECB 模式）。

所有类 ECB 模式对静态分析、差分分析、已知部分明文攻击等破解手段防护能力都很弱，无法确保数据安全。

白豚云使用基于 CTR 模式的 AES 强加密算法，对包含数据库、配置信息以及用户上传附件等在内的所有用户数据进行整库强加密，提供数学上可证明的强安全性（可靠性、一致性、保密性）防护。即使是相同的明文，也会被变换为不同且完全无关的密文，从而杜绝了上述攻击方式。

整库强加密支持基于《白杨应用支撑平台》中提供的，支持实时压缩和强加密的分布式虚拟文件系统（VFS）组件。该组件经过多家 500 强企业生产环境、以及互联网、电信高负载复杂环境下，多年、多产品的广泛考验，拥有优异的性能表现和极强的稳定性。

基于业界标准的强加密算法保证了即使未来出现了每秒钟可完成一千万亿次密钥破解尝试的超级计算机，也需要平均五千四百万亿年才能破解一个密钥。安全性得到了极大保证，可以确保即使 IDC 运维人员或系统管理员将服务器存储设备打包拆走，也无法破解其中保存的用户数据（具体可参考《[白杨应用支撑平台技术白皮书](#)》中的第 4 节）。

4.2.2 鉴权和授权安全性

白豚 ERP 使用加盐的安全散列值存储用户密钥。这些信息将经过压缩和强加密处理后才保存到数据库存储中。正如前文所述，数据库存储本身又以实时强加密的形式保存于我方自主研发的分布式文件系统中，为用户密钥安全提供了二次保障。

白豚 ERP 使用基于 SHA1 和 HMAC 安全算法的 CHAP 鉴权协议完成鉴权流程。即使不依赖 TLS 加密保护（全明文传输），也不会发生密钥泄露等安全问题。与此同时，CHAP 安全鉴权协议可很好地低于中间人、重放攻击等目前已知的各类攻击手段。

防御暴力破解尝试：同一用户连续 5 次登陆失败后，强制冷却 5 分钟，配合入侵检测和恶意请求过滤技术，防止穷举破解或字典攻击。此外，管理员还可自定义最小长度，大小写字母和数字组合等密码复杂度要求，防止用户使用过于简单的登陆密码。

白豚 ERP 使用基于角色的严格授权机制和访问控制列表（ACL）来统一管理用户权限。拒绝所有未经授权的数据访问请求。

4.2.3 传输安全性

白豚云使用公认权威机构发布的数字证书，以及 STS（Strict Transport Security）指令等技术，实现强制的网络传输全程加密保护，提供**网银级别的 TLS（HTTPS）网络数据传输保护**。

4.2.4 系统安全性

Java (JRE)、C# (.NET) 以及 Tomcat 等环境存在大量严重漏洞，有极高安全隐患。通过国际安全组织公开的 CVE 数据库可发现，上述产品几乎每年都会产生多个安全等级为“严重”的新漏洞。例如：

- ★ Java: https://www.cvedetails.com/vulnerability-list/vendor_id-93/product_id-19117/Oracle-JRE.html
- ★ Tomcat: https://www.cvedetails.com/vulnerability-list/vendor_id-45/product_id-887/Apache-Tomcat.html
- ★ .NET: https://www.cvedetails.com/vulnerability-list/vendor_id-26/product_id-2002/Microsoft-.net-Framework.html

等等。

其它 ERP/HRMS 产品组成

某某 ERP / 某某 HRMS.....

应用程序业务逻辑

第三方框架

JBoss、Struts、Spring.....

Web 容器

Tomcat、Jetty、IIS、Apache.....

动态运行时

JRE (Java) 、CLR (.NET) 、ZEND/HHVM (PHP).....

底层基础设施

操作系统     

硬件平台 (x86/x64、ARM、MIPS、POWER.....)

图 16

如图 16 所示：对 Java、.NET 等动态运行时环境、Tomcat 等容器，以及 JBoss 等第三方框架的依赖，导致系统中被引入了大量额外的，不受开发商控制的第三方组件。

可以看到，使用 Java、.NET、PHP 等动态语言开发的产品，其应用程序与底层基础设施间包含了第三方框架、Web 容器、动态运行时等多个额外层面。这些不受开发商控制的，由第三方软件组成的额外层面包含大量广为人知，可被轻易利用的安全漏洞。

甚至任何稍有常识的普通用户都可以通过搜索引擎简单地找到可成功利用这些漏洞进行攻击的免费开源工具，并以此入侵企业系统，完成数据窃取、数据篡改、数据销毁等非法操作。

由此可见，引入不受控制的额外层面大大增加了系统的受攻击面，严重增加了系统的安全风

险。同时大大增加了运维人员的工作负担，提升了系统的总体拥有成本（CTO）。

例如：近年来发生的 Yahoo 泄露 30 亿用户数据、京东泄露包括信用卡、登陆密码、身份证号码等在内的数千万用户敏感数据（Java Struts）、信用巨头 Equifax 泄露 1.43 亿（三分之二美国成年人）的敏感征信记录（包括社保、驾照、住址、信用卡等，Java Struts）、优步（Uber）外泄 5700 万用户资料、支付宝泄露 2500 万用户敏感数据、网易泄露过亿用户数据等惨痛的安全事故案例，均是由于系统中引入的上述第三方软件层而引起。

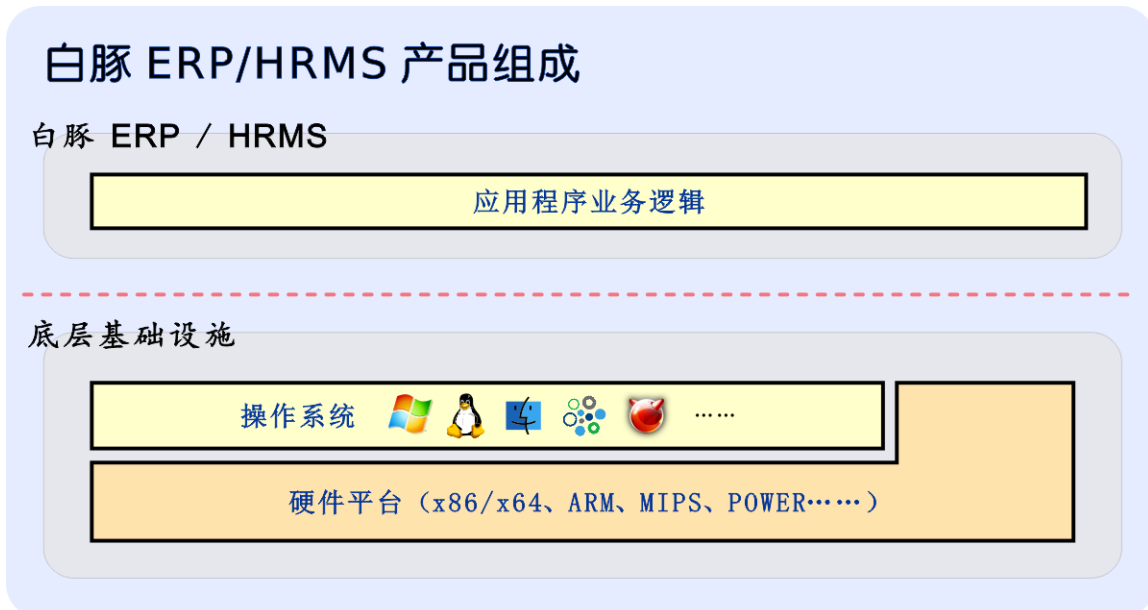


图 17

如图 17 所示：白豚 ERP 系统使用汇编、C/C++ 等工具开发，摒弃了第三方框架、容器和语言动态运行等多余的第三方软件层，直接构建在底层操作系统和硬件平台之上。

这种以高效的原生机器码（Native Code）直接运行在裸硬件和操作系统构成的底层基础设施上的设计，为 ERP 系统带来了无与伦比的时空效率和安全性保障。

由于无需依赖动态运行时、Web 容器等不可控的第三方组件，因此大大减少了软件的受攻击面，极大地降低了系统安全隐患。同时也在很大程度上减轻了运维人员的工作负担，降低了客户总体成本。

5. 技术特点

本节集中讨论白豚 ERP 系统的主要技术特点和关键设计考量。

5.1 REST 风格的 Web API 接口

在综合考虑性能、网关穿透性、安全性、易用性和学习曲线等多方面因素后，白豚 ERP 系统使用 RESTful 风格的 HTTP 接口（也常被称为 Web API）为其它子系统提供功能调用。并通过 SSL/TLS 层为其提供通信安全保护服务。HTTP/HTTPS 作为当前应用最广泛的网络协议，已被几乎所有开发人员认可并熟知，并已有大量成熟的封装组件实现了完整的 HTTP 协议。这些因素使得其学习成本和易用性都非常理想。另一方面，HTTP/1.1 对 KeepAlive 和 Pipelining、以及 HTTP/2 对全双工多路复用通信的支持使其能够很好地工作在高并发、高负载环境。

在今天的网络环境中，Web API 已被广泛应用于当今各大站点中，例如：amazon.com 的 EC2/S3、淘宝的 TOP（即：淘宝提供给其各个 ISV 的开发接口）、以及 Google、Facebook、Twitter、Flicker、百度、腾讯等网站都以 Web API 的形式开放了自己的二次开发接口。

Web API 的优越性主要体现在以下几方面：

- ★ **实现方便：**B/S 架构的应用本身就使用 HTTP 实现，只需要进行较少的扩展即可作为 API 提供给用户使用。
- ★ **连通性好：**能够通过任何 HTTP Proxy 透传。
- ★ **基础设施完备：**可基于大量现有的，经过高度优化和长时间测试的 Web Server、Web 扩展框架和反向代理来实现。可直接利用 HTTP 协议中已实现的 SSL/TLS 加密和 gzip/zlib 压缩算法完成数据传输的加密和压缩等功能。
- ★ **使用 HTTP 语义：**不同于 SOAP 等仅将 HTTP 作为传输层协议的技术，RESTful HTTP 接口要求使用 HTTP 语义来进行通信，意即：请求中的 GET/PUT/POST/DELETE/PUT 等方法、URI 以及 HTTP 消息头和返回中的状态码及消息头均遵循标准的 HTTP 规范而代表着对应的行为和参数等含义，以能够真实反应调用的实际语义为佳。
- ★ **无状态：**HTTP 本身为无状态协议，便于实现负载均衡、容错和提高实现效率。

5.2 nano-SOA 架构

5.2.1 SOA vs. AIO

长久以来,服务器端的高层架构大体被区分为对立的两类:SOA(Service-oriented architecture)以及AIO(All in one)。SOA 将一个完整的应用分割为相互独立的服务,每个服务提供一个单一标准功能(如:会话管理、交易评价、用户积分等等)。服务间通过RPC、WebAPI等IPC机制暴露功能接口,并以此相互通信,最终组合成一个完整的应用。

而AIO则相反,它将一个应用规约在一个独立的整体中,SOA中的不同服务在AIO架构下呈现为不同的功能组件和模块。AIO应用的所有组件通常都运行在一个地址空间(同一进程)内,所有组件的代码也常常放在同一个产品项目中一起维护。

AIO的优势是部署简单,不需要分别部署多个服务,并为每个服务实现一套高可用集群。与此同时,由于可避免网络传输、内存拷贝等IPC通信所带来的大量开销,因此AIO架构的单点效率通常远高于SOA。

另一方面,由于AIO架构中组件依赖性强,组件间经常知晓并相互依赖对方的实现细节,因此组件的可重用性及可替换性差,维护和扩展也较困难。特别是对于刚加入团队的新人来说,面对包含了大量互相深度耦合之组件和模块的“巨型项目”,常常需要花费大量努力、经历很多挫折并且犯很多错误才能真正接手。而即使对于老手来说,由于模块间各自对对方实现细节错综复杂的依赖关系,也容易发生在修改了一个模块的功能后,莫名其妙地影响到其它看起来毫不相干功能的情况。

与此相反,SOA模型部署和配置复杂——现实中,一个大型应用常常被拆分为数百个相互独立的服务,《程序员》期刊中的一份公开发表的论文显示,某个国内“彻底拥抱”SOA的著名(中国排名前3)电商网站将他们的Web应用拆分成了一千多个服务。可以想象,在多活数据中心的高可用环境内部署成百上千个服务器集群,并且配置他们彼此间的协作关系是多大的工作量。最近的携程(ctrip.com)网络瘫痪事件也是因为上千个服务组成的庞大SOA架构导致故障恢复缓慢。

除了部署复杂以外,SOA的另一个主要缺点就是低效——从逻辑流的角度看,几乎每次来自客户端的完整请求都需要依次流经多个服务后,才能产生最终结果并返回用户端。而请求(通过消息中间件)每“流经”一个服务都需要伴随多次网络IO和磁盘访问,多个请求可累计产生较高的网络时延,使用户请求的响应时间变得不可确定,用户体验变差,并额外消耗大量资源。

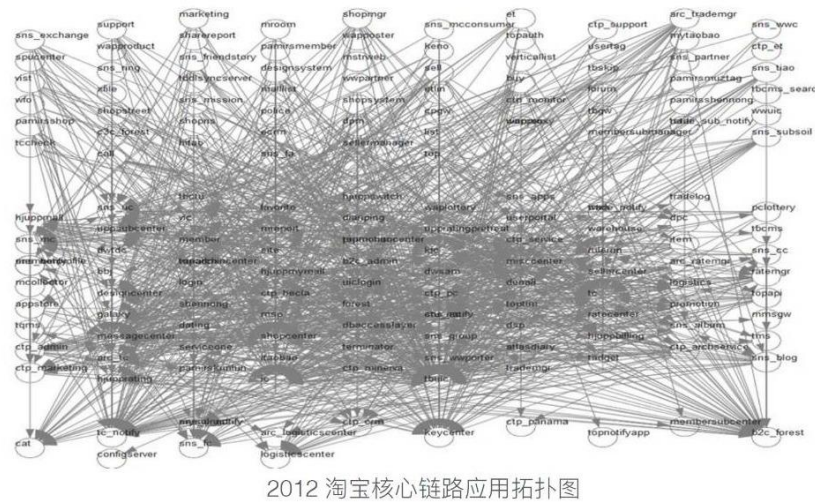


图 18 混乱的 SOA 依赖关系（图片来自互联网）

此外,无论是每个 Service 各自连接不同的 DBMS 还是它们分别接入同一个后端分布式 DBMS 系统,实现跨服务的分布式事务支持工作都要落到应用层开发者手中。而分布式事务 (XA) 本身的实现复杂度恐怕就已超过大部分普通应用了,更何况还需要为分布式事务加上高可靠和高可用保证——需要在单个数据切片上使用 Paxos/Raft 或主从+Arbiter 之类的高可用、强一直性算法,同时在涉及多个数据切片的事务上使用 2PC/3PC 等算法来保证事务的原子性。因此 SOA 应用中的跨 Service 事务基本都只能退而求其次,做到最终一致性保证,即便如此,也需要增加大量的额外工作——在稍微复杂点的系统里,高可用,并能在指定时间内可靠收敛的最终一致性算法实现起来也不是那么容易。

与此同时,大部分 SOA 系统还经常需要使用消息中间件来实现消息分发服务。如果对消息中间件的可用性(部分节点故障不会影响正常使用)、可靠性(即使在部分节点故障时,也确保消息不丢失、不重复、并严格有序)、功能性(如:发布/订阅模型、基于轮转的任务分发等)等方面有所要求的话,那么消息中间件本身也容易成为系统的瓶颈。

SOA 架构的优点在于其高内聚、低耦合的天然特性。仅通过事先约定的 IPC 接口对外提供服务,再配合服务间隔离(通常是在独立节点中)运行的特质,SOA 架构划分出了清晰的接口和功能边界,因此可以被非常容易地重用和替换(任何实现了兼容 IPC 接口的新服务都可替换已有的老服务)。

从软件工程和项目管理的视角来看,由于每个服务本身通常有足够高的内聚性,并且单个服务实现的功能也较独立,因此相对于 AIO 意大利面式的,相互交织的结构来说,SOA 的服务非常便于维护——负责某一服务的开发人员只需要看好自己这一亩三分地即可,只要保持服务对外提供的 API 没有发生不兼容的变化,就不需要担心修改代码、替换组件等工作会影响到其它“消费者”。

同时,由多个独立服务所组成的应用也更容易通过加入新服务和重新组合现有服务来进行功能变更和扩展。

5.2.2 nano-SOA 架构

在经历了大量实际项目中的权衡、思索和实践后，我逐步定义、实现和完善了能够兼两者之长的“nano-SOA”专利架构。并将功能插件 (IPlugin)、数据库连接器 (DBC) API Nexus 等关键性的基础功能组件封装在了《[白杨应用支撑平台](#)》中，以提供便于重用的高品质基础架构。

在 nano-SOA 架构中，独立运行的服务被替换成了支持动态插拔的跨平台功能插件 (IPlugin)；而插件则通过（并仅可通过）API Nexus 来动态地暴露（注册）和隐藏（注销）自身所提供的功能接口，同时也使用 API Nexus 来消费其它插件提供服务。

nano-SOA 完全继承了 SOA 架构高内聚、低耦合的优点，每个插件如独立的服务一样，有清晰的接口和边界，可容易地被替换和重用。在可维护性上，nano-SOA 也与 SOA 完全一致，每个插件都可以被单独地开发和维护，开发人员只需要管好自己维护的功能插件即可。通过加入新插件以及对现有功能插件的重新组合，甚至可比 SOA 模式更容易地对现有功能进行变更和扩展。

而在性能方面，由于所有功能插件都运行在同一个进程内，因此通过 API Nexus 的相互调用不需要任何网络 IO、磁盘访问和内存拷贝，也没有任何形式的其它 IPC 开销，因此其性能和效率均可与 AIO 架构保持在相同量级。

与此同时，nano-SOA 的部署与 AIO 同样简单——部署在单个节点即可使用，只需部署一个集群即可实现高可用和横向扩展。在配置方面也远比 SOA 简单，仅需要比 AIO 应用多配置一个待加载模块列表而已，并且这些配置也可通过各种配置管理产品来实现批量维护。简单的部署和配置过程不但简化了运营和维护工作，也大大方便了开发和测试环境的构建。

此外，nano-SOA 也在极大程度上避免了对消息中间件的依赖，取而代之的是通过 API Nexus 的直接 API 调用；或是在需要削峰填谷的场合中，使用由内存零拷贝和无锁算法高度优化的线程间消息队列。这一方面大大增加了吞吐，避免了延迟，另一方面也避免了部署和维护一个高可用的消息分发服务集群所带来的巨大工作量——nano-SOA 集群内的节点间协作和协调通信需求已被将至最低，对消息分发的可靠性、可用性和功能性都没有太高要求。在多数情况下，使用 Gossip Protocol 等去中心化的 P2P 协议即足以满足需要，有时甚至可以完全避免这种集群内的节点间通信。

从 nano-SOA 的角度看，也可以将 DBC 视作一种几乎所有服务器端应用都需要使用的基础功能插件，由于其常用性，因此他们被事先实现并加进了 libapidbc 中。由此，通过提供 IPlugin、API Nexus 以及 DBC 等几个关键组件，libapidbc 为 nano-SOA 架构奠定了良好的基础设施。

当然，nano-SOA 与 SOA 和 AIO 三者间并不是互斥的选择。在实际应用场景中，可以通过三者间的有机组合来达成最合理的设计。例如：对于视频转码等非常耗时并且不需要同步等待其完成并返回结果的异步操作来说，由于其绝大部分开销都耗费在了视频编解码计算上，因此将其作为插件加入其它 App Server 就完全没有必要，将它作为独立的服务，部署在配置了专用加速硬件的服务器集群上应该是更好的选择。

5.2.3 白杨消息端口交换服务 (BYPSS)

白杨消息端口交换服务 (BYPSS) 设计用于单点支撑万亿量级端口、百万量级节点规模，每秒处理千万至十亿量级消息的高可用、强一致、高性能分布式协调和消息交换服务。其中关键概念包括：

- ★ **连接 (Connection)**：每个客户端 (应用集群中的服务器) 节点至少与端口交换服务保持一个 TCP 长连接。
- ★ **端口 (Port)**：每个连接上可以注册任意多个消息端口，消息端口由一个 UTF-8 字符串描述，必须在全局范围内唯一，若其它客户端节点已注册了相同的消息端口，则端口注册失败。

端口交换服务对外提供的 API 原语包括：

- ★ **等待消息 (WaitMsg)**：客户端集群中的每个节点均应保持一个到端口交换服务的 TCP 长连接，并调用此方法等待消息推送。此方法将当前客户端连接由消息发送连接升级为消息接收连接。

每个节点号只能对应一个消息接收连接，若一个节点尝试同时发起两个消息接收连接，则较早的那个接收连接将被关闭，并且绑定到当前节点上的所有端口都将被注销。

- ★ **续租 (Relet)**：若端口交换服务在指定的间隔内未收到来自某个消息接收连接的续租请求，则判定该节点已经下线，并释放所有属于该节点的端口。续租操作用来周期性地向端口交换服务提供心跳信号。
- ★ **注册端口 (RegPort)**：连接成功建立后，客户端应向端口交换服务注册所有属于当前节点的消息端口。可以在一个端口注册请求中包含任意多个待注册端口，端口交换服务会返回所有注册失败 (已被占用) 的端口列表。调用者可以选择是否需要为注册失败的端口订阅端口注销通知。

需要注意的是，每次调用 **WaitMsg** 重建消息接收连接后，都需要重新注册当前节点上的所有端口。

- ★ **注销端口 (UnRegPort)**：注销数据当前节点的端口，可一次提交多个端口执行批量注销。
- ★ **消息发送 (SendMsg)**：向指定的端口发送消息 (BLOB)，消息格式对交换服务透明。若指定的端口为空串，则向端口交换服务上的所有节点广播此消息；亦可同时指定多个接收端口，实现消息组播。若指定的端口不存在，则安静地丢弃该消息。客户端可在一次请求中包含多个消息发送命令，主动执行批量发送，服务器端也会将发往同一节点的消息自动打包，实现消息批量推送。

- ★ **端口查询 (QueryPort)**：查询当前占用着指定端口的节点号，及其 IP 地址。此操作主要用于实现带故障检测的服务发现，消息投递时已自动执行了相应操作，故无需使用此方法。可在同一请求中包含多个端口查询命令，执行批量查询。
- ★ **节点查询 (QueryNode)**：查询指定节点的 IP 地址等信息。此操作主要用于实现带故障检测的节点解析。可以一次提交多个节点，实现批量查询。

端口交换服务的客户端连接分为以下两类：

- ★ **消息接收连接 (1:1)**：接收连接使用 WaitMsg 方法完成节点注册并等待消息推送，同时通过 Relet 接口保持属于该节点的所有端口被持续占用。客户端集群中的每个节点应当并且仅能够保持一个消息接收连接。此连接为长连接，由于连接中断重连后需要重新进行服务选举（注册端口），因此应尽可能一直保持该连接有效并及时完成续租。
- ★ **消息发送连接 (1:N)**：所有未使用 WaitMsg API 升级的客户端连接均被视为发送连接，发送连接无需通过 Relet 保持心跳，仅使用 RegPort、UnRegPort、SendMsg 以及 QueryPort 等原语完成非推送类的客户端请求。客户端集群中的每个节点通常都会维护一个消息发送连接池，以方便各工作线程高效地与端口转发服务保持通信。

与传统的分布式协调服务以及消息中间件产品相比，端口转发服务主要有以下特点：

- ★ **功能性**：端口转发服务将标准的消息路由功能集成到了服务选举（注册端口）、服务发现（发送消息和查询端口信息）、故障检测（续租超时）以及分布式锁（端口注册和注销通知）等分布式协调服务中。是带有分布式协调能力的高性能消息转发服务。通过 QueryPort 等接口，也可以将其单纯地当作带故障检测的服务选举和发现服务来使用。
- ★ **高并发、高性能**：由 C/C++/汇编实现；为每个连接维护一个消息缓冲队列，将所有端口定义及待转发消息均保存在内存中（Full in-memory）；主从节点间无任何数据复制和状态同步开销；信息的发送和接收均使用纯异步 IO 实现，因而可提供高并发和高吞吐的消息转发性能。
- ★ **可伸缩性**：在单点性能遭遇瓶颈后，可通过级联上级端口交换服务来进行扩展（类似 IDC 接入、汇聚、核心等多层交换体系）。
- ★ **可用性**：最低 5 毫秒内完成故障检测和主备切换的高可用保证，基于多数派的选举算法，避免由网络分区引起的脑裂问题。
- ★ **一致性**：保证任意给定时间内，最多只有一个客户端节点可持有某一特定端口。不可能出现多个客户端节点同时成功注册和持有相同端口的情况。
- ★ **可靠性**：所有发往未注册（不存在、已注销或已过期）端口的消息都将被安静地丢弃。系统保证所有发往已注册端口消息有序且不重复，但在极端情况下，可能发生消息丢失：

- 端口交换服务宕机引起主从切换：此时所有在消息队列中排队的待转发消息均会丢失；所有已注册的客户端节点均需要重新注册；所有已注册的端口（服务和锁）均需要重新进行选举/获取（注册）。
- 客户端节点接收连接断开重连：消息接收连接断开或重连后，所有该客户端节点之前注册的端口均会失效并需重新注册。在接收连接断开到重连的时间窗口内，所有发往之前与该客户端节点绑定的，且尚未被其它节点重新注册的端口之消息均被丢弃。

可见，白杨消息端口转发服务本身是一个集成了故障检测、服务选举、服务发现和分布式锁等分布式协调功能的消息路由服务。它通过牺牲极端条件下的可靠性，在保证强一致、高可用、可伸缩（横向扩展）的前提下，实现了极高的性能和并发能力。

可以认为消息端口交换服务就是为 nano-SOA 架构量身定做的集群协调和消息分发服务。nano-SOA 的主要改进即：将在 SOA 中，每个用户请求均需要牵扯网络中的多个服务节点参与处理的模型改进为大部分用户请求仅需要同一个进程空间内的不同 BMOD 参与处理。

这样的改进除了便于部署和维护，以及大大降低请求处理延迟外，还有两个主要的优点：

- ★ 将 SOA 中，需要多个服务节点参与的分布式事务或分布式最终一致性问题简化成为了本地 ACID Transaction 问题（从应用视角来看是如此，对于分布式 DBS 来说，以 DB 视角看来，事务仍然可以是分布式的），这不仅极大地简化了分布式应用的复杂度，增强了分布式应用的一致性，也大大减少了节点间通信（由服务间的 IPC 通信变成了进程内的指针传递），提高了分布式应用的整体效率。
- ★ 全对等节点不仅便于部署和维护，还大大简化了分布式协作算法。同时由于对一致性要求较高的任务都已在同一个进程空间内完成，因此节点间通信不但大大减少，而且对消息中间件的可靠性也不再有过高的要求（通常消息丢失引起的不一致可简单地通过缓存超时或手动刷新来解决，可确保可靠收敛的最终一致性）。

在此前提下，消息端口交换服务以允许在极端情况下丢失少量未来得及转发的消息为代价，来避免磁盘写入、主从复制等低效模式，以提供极高效率。这对 nano-SOA 来说是一种非常合理的选择。

5.2.3.1 极端条件下的可靠性

传统的分布式协调服务通常使用 Paxos 或 Raft 之类基于多数派的强一致分布式算法实现，主要负责为应用提供一个高可用、强一致的分布式元数据 KV 访问服务。并以此为基础，提供分布式锁、消息分发、配置共享、角色选举、服务发现、故障检测等分布式协调服务。常见的分布式协调服务实现包括 Google Chubby (Paxos)、Apache ZooKeeper (Fast Paxos)、etcd (Raft)、Consul (Raft+Gossip) 等。

Paxos、Raft 等分布式一致性算法的最大问题在于其极低的访问性能和极高的网络开销：对这些服务的每次访问，无论读写，都会产生至少 2 到 4 次网络广播——以投票的方式确定本次访问经过多数派确认（读也需要如此，因为主节点需要确认本次操作发生时，自己仍拥有多数票支持，仍是集群的合法主节点）。

在实践中，虽可通过降低系统整体一致性或加入租期机制来优化读操作的效率，但其总体性能仍十分低下，并且对网络 IO 有很高的冲击：Google、Facebook、Twitter 等公司的历次重大事故中，很多都是由于发生网络分区或人为配置错误导致 Paxos、Raft 等算法疯狂广播消息，致使整个网络陷入广播风暴而瘫痪。

此外，由于 Paxos、Raft 等分布式一致性算法对网络 IO 的吞吐和延迟等方面均有较高要求，而连接多座数据中心机房（IDC）的互联网络通常又很难满足这些要求，因此导致依赖分布式协调算法的强一致（抗脑裂）多活 IDC 高可用集群架构难以以合理成本实现。作为实例：2018 年 9 月 4 日微软美国中南区某数据中心空调故障导致 Office、Active Directory、Visual Studio 等服务下线近 10 小时不可用；2015 年 8 月 20 日 Google GCE 服务中断 12 小时并永久丢失部分数据；2015 年 5 月 27 日、2016 年 7 月 22 及 2019 年 12 月 5 日支付宝多次中断数小时；2013 年 7 月 22 日微信服务中断数小时；以及 2017 年 5 月英国航空瘫痪数日等重大事故均是由于单个 IDC 因市政施工（挖断光纤）等原因下线，同时未能成功构建多活 IDC 架构，因此造成 IDC 单点依赖所导致的。

前文也已提到过：由于大部分采用 SOA 架构的产品需要依赖消息中间件来确保系统的最终一致性。因此对其可用性（部分节点故障不会影响正常使用）、可靠性（即使在部分节点故障时，也确保消息不丢失、不重复、并严格有序）、功能性（如：发布/订阅模型、基于轮转的任务分发等）等方面均有较严格的要求。这就必然要用到高可用集群、节点间同步复制、数据持久化等低效率、高维护成本的技术手段。因此消息分发服务也常常成为分布式系统中的一大主要瓶颈。

与 Paxos、Raft 等算法相比，BYPSS 同样提供了故障检测、服务选举、服务发现和分布式锁等分布式协调功能，以及相同等级的强一直性、高可用性和抗脑裂（Split Brain）能力。在消除了几乎全部网络广播和磁盘 IO 等高开销操作的同时，提供了数千、甚至上万倍于前者的访问性能和并发处理能力。可在对网络吞吐和延迟等方面无附加要求的前提下，构建跨多个 IDC 的大规模分布式集群系统。

与各个常见的消息中间件相比，BYPSS 提供了一骑绝尘的单点百万至千万条消息每秒的吞吐和路由能力——同样达到千百倍的性能提升，同时保证消息不重复和严格有序。

然而天下没有免费的午餐，特别是在分布式算法已经非常成熟的今天。在性能上拥有绝对优势的同时，BYPSS 必然也有其妥协及取舍——BYPSS 选择放弃极端（平均每年 2 次，并且大多由维护引起，控制在低谷时段，基于实际生产环境多年统计数据）情形下的可靠性，对分布式系统的具体影响包括以下两方面：

- ★ 对于分布式协调服务来说，这意味着每次发生 BYPSS 主节点故障掉线后，所有的已注册端口都会被强制失效，所有活动的端口都需要重新注册。

例如：若分布式 Web 服务器集群以用户为最小调度单位，为每位已登陆用户注册一个消息端口。则当 BYPSS 主节点因故障掉线后，每个服务器节点都会得知自己持有的所有端口均已失效，并需要重新注册当前自己持有的所有活动（在线）用户。

幸运的是，该操作可以被批量化地完成——通过批量端口注册接口，可在一次请求中同时提交多达数百万端口的注册和注销操作，从而大大提升了请求处理效率和网络利用率：在 2013 年出厂的至强处理器上（Haswell 2.0GHz），BYPSS 服务可实现每核（每线程）100 万端口/秒的处理速度。同时，得益于我方自主实现的并发散列表（每个 arena 都拥有专属的汇编优化用户态读者/写者高速锁），因此可通过简单地增加处理器核数来实现处理能力的线性扩展。

具体来说，在 4 核处理器+千兆网卡环境下，BYPSS 可达成约每秒 400 万端口注册的处理能力；而在 48 核处理器+万兆网卡环境下，则可实现约每秒 4000 万端口注册的处理能力（测试时每个端口名称的长度均为 16 字节），网卡吞吐量和载荷比都接近饱和。再加上其发生概率极低，并且恢复时只需要随着对象的加载来逐步完成重新注册，因此对系统整体性能几乎不会产生什么波动。

为了说明这个问题，考虑 10 亿用户同时在线的极端情形，即使应用程序为每个用户分别注册一个专用端口（例如：用来确定用户属主、完成消息分发等），那么在故障恢复后的第一秒内，也不可能出现“全球 10 亿用户心有灵犀地同时按下刷新按钮”的情况。相反，基于 Web 等网络应用的固有特性，这些在线用户通常要经过几分钟、几小时甚至更久才会逐步返回服务器（同时在线用户数=每秒并发请求数 x 用户平均思考时间）。即使按照比较严苛的“1 分钟内全部返回”（平均思考时间 1 分钟）来计算，BYPSS 服务每秒也仅需处理约 1600 万条端口注册请求。也就是说，一台配备了 16 核至强处理器和万兆网卡的入门级 1U PC Server 即可满足上述需求。

作为对比实例：官方数据显示，淘宝网 2015 年双十一当天的日活用户数（DAU）为 1.8 亿，同时在线用户数峰值为 4500 万。由此可见，目前超大型站点瞬时并发用户数的最高峰值仍远低于前文描述的极端情况。即使再提高数十倍，BYPSS 也足可轻松支持。

- ★ 另一方面，对于消息路由和分发服务来说，这意味着每次发生 BYPSS 主节点故障掉线后，所有暂存在 BYPSS 服务器消息队列中，未及发出的待转发消息都将永久丢失。可喜的是，nano-SOA 架构不需要依赖消息中间件来实现跨服务的事务一致性。因此对消息投递的可靠性并无严格要求。

意即：nano-SOA 架构中的消息丢失最多导致对应的用户请求完全失败，此时仍可保证数据的全局强一致性，绝不会出现“成功一半”之类的不一致问题。在绝大多数应用场景中，这样的保证已经足够——即使支付宝和四大行的网银应用也会偶尔发生操作失败的问题，这时只要资金等帐户数据未出现错误，那么稍候重试即可。

此外，BYPSS 服务也通过高度优化的异步 IO，以及消息批量打包等技术有效降低了消息在服务器队列中的等待时间。具体来说，这种消息批量打包机制由消息推送和消息发送机制两方面组成：

BYPSS 提供了消息批量发送接口，可在一次请求中同时提交数以百万计的消息发送操作，从而大大提升了消息处理效率和网络利用率。另一方面，BYPSS 服务器也实现了消息批量打包推送机制：若某节点发生消息浪涌，针对该节点的消息大量到达并堆积在服务器端消息队列中。则 BYPSS 服务器会自动开启批量消息推送模式——将大量消息打包成一次网络请求，批量推送至目的节点。

通过上述的批量处理机制，BYPSS 服务可大大提升消息处理和网络利用效率，确保在大部分情况下，其服务器端消息队列基本为空，因此就进一步降低了其主服务器节点掉线时，发生消息丢失的概率。

然而，虽然消息丢失的概率极低，并且 nano-SOA 架构先天就不怎么需要依赖消息中间件提供的可靠性保证。但仍然可能存在极少数对消息传递要求很高的情况。对于此类情况，可选择使用下列解决方案：

- 自行实现回执和超时重传机制：消息发送方对指定端口发送消息，并等待接收该消息处理回执。若在指定时段内未收到回执，则重新发送请求。
- 直接向消息端口的属主节点发起 RPC 请求：消息发送方通过端口查询命令获取该端口属主节点的 IP 地址等信息，并直接与该属主节点建立连接、提交请求并等待其返回处理结果。BYPSS 在此过程中仅担当服务选举和发现的角色，并不直接路由消息。对于视频推流和转码、深度学习等有大量数据流交换的节点间通信，也建议使用此方式，以免 BYPSS 成为 IO 瓶颈。
- 使用第三方的可靠消息中间件产品：若需要保证可靠性的消息投递请求较多，规则也较复杂，也可单独搭建第三方的可靠消息分发集群来处理这部分请求。

当然，实际上不存在完全可靠的消息队列服务（能保证消息不丢失、不重复、不乱序）。因此在确实需要实现跨应用服务器节点的分布式事务时，建议通过 BYPSS 以及 BYDMQ 配合 SAGA 等模式来实现，详见：5.2.4 白杨分布式消息队列(BYDMQ)。

综上所述，可以认为 BYPSS 服务就是为 nano-SOA 架构量身定做的集群协调和消息分发服务。BYPSS 和 nano-SOA 架构之间形成了扬长避短的互补关系：BYPSS 以极端条件下系统整体性能的轻微波动为代价，极大提升了系统的总体性能表现。适合用来实现高效率、高可用、高可靠、强一致的 nano-SOA 架构分布式系统。

5.2.3.2 BYPSS 特性总结

BYPSS 和基于 Paxos、Raft 等传统分布式一致性算法的分布式协调产品特性对比如下：

项目	BYPSS	ZooKeeper、Consul、etcd...
----	-------	--------------------------

项目	BYPSS	ZooKeeper、Consul、etcd...
可用性	高可用，支持多活 IDC。	高可用，但难以支持多活 IDC。
一致性	强一致，主节点通过多数派选举。读写操作均提供强一致保证。	写入强一致，多副本复制；大部分实现为提升性能，读取时牺牲了一致性（仅 Consul 支持配置成强一致读取模式）。
并发性	千万量级并发连接，可支持数十万并发节点	不超过 5000 节点
容量	每 10GB 内存可支持至少 1 亿消息端口；每 1TB 内存可支持至少 100 亿消息端口；两级并发散列表结构确保容量可线性扩展至 PB 级。	通常最高支持数十万 KV 对。开启了变更通知时则更少。
延迟	相同 IDC 内每次请求延迟在亚毫秒级（阿里云中实测为 0.5ms）；相同区域内的不同 IDC 间每次请求延迟在毫秒级（阿里云环境实测 2ms）。	由于每次请求需要至少 2 到 4 次网络广播和多次磁盘 IO，因此相同 IDC 中的每操作延迟在十几毫秒左右；不同 IDC 间的延迟则更长（详见下文）。
性能	每 1Gbps 网络带宽可支持约 400 万次/秒的端口注册和注销操作。在 2013 年出厂的入门级至强处理器上，每核心可支持约 100 万次/秒的上述端口操作。性能可通过增加带宽和处理器核心数量线性扩展。在现代处理器和双口 4 万兆网卡上可达 3 亿次操作/秒的处理能力。	算法本身的特性决定了无法支持批量操作，相同测试条件下不到 200 次每秒的请求性能（由于每个原子操作都需要至少 2 到 4 次网络广播和多次磁盘 IO，因此支持批量操作毫无意义，详见下文）。
网络利用率	高：服务器端和客户端均具备端口注册、端口注销、消息发送、端口查询、节点查询等原语的批量打包能力，网络载荷比可接近 100%。	低：每请求一个独立包（TCP Segment、IP Packet、Network Frame），网络载荷比通常低于 5%。
可伸缩性	有：可通过级联的方式进行横向扩展。	无：集群中的节点越多（因为广播和磁盘 IO 的范围更大）性能反而越差。
分区容忍	无多数派分区时系统下线，但不会产生广播风暴。	无多数派分区时系统下线，有可能产生广播风暴引发进一步网络故障。
消息分发	有，高性能，客户端和服务端均包含了消息的批量自动打包支持。	无。
配置管理	无，BYPSS 认为配置类数据应交由 Redis、MySQL、MongoDB 等专门的产品来维护和管理。当然，这些 CMDb 的主从选举等分布式协调工作仍可由 BYPSS 来完成。	有，可当作简单的 CMDb 来使用，这种功能和职责上的混淆不清进一步劣化了产品的容量和性能。
故障恢复	需要重新生成状态机，但可以数千万至数亿端口/秒的性能完成。实际使用中几无波动。	不需要重新生成状态机。

上述比较中，延迟和性能两项主要针对写操作。这是因为在常见的分布式协调任务中，几乎全部有意义的操作都是写操作。例如：

操作	对服务协调来说	对分布式锁来说
端口注册	成功：服务选举成功，成为该服务的属主。 失败：成功查询到该服务的当前属主。	成功：上锁成功。 失败：上锁失败，同时返回锁的当前属主。
端口注销	放弃服务所有权。	释放锁。

操作	对服务协调来说	对分布式锁来说
注销通知	服务已下线，可更新本地查询缓存或参与服务竞选。	锁已释放，可重新开始尝试上锁。

上表中，BYPSS 的端口注册对应 ZooKeeper 等传统分布式产品中的“写/创建 KV 对”；端口注销对应“删除 KV 对”；注销通知则对应“变更通知”服务。

由此可见，为了发挥最高效率，在生产环境中通常不会使用单纯的查询等只读操作。而是将查询操作隐含在端口注册等写请求中，请求成功则当前节点自身成为属主；注册失败自然会返回请求服务的当前属主，因此变相完成了属主查询（服务发现/名称解析）等读操作。

需要注意的是，就算是端口注册等写操作失败，其实还是会伴随一个成功的写操作。因为仍然要将发起请求的当前节点加入到指定条目的变更通知列表中，以便在端口注销等变更事件发生时，向各个感兴趣的节点推送通知消息。因此写操作的性能差异极大地影响了现实产品的实际表现。

注：以上所述 nano-SOA 架构和 BYPSS 分布式协调算法均受到多项国家和国际发明专利保护。

5.2.3.3 基于 BYPSS 的高性能集群

从高性能集群（HPC）的视角来看，BYPSS 与前文所述的传统分布式协调产品之间，最大的区别主要体现在以下两个方面：

1. **高性能：**BYPSS 通过消除网络广播、磁盘 IO 等开销，以及增加批处理支持等多种优化手段使分布式协调服务的整体性能提升了上万倍。
2. **大容量：**每 10GB 内存至少 1 亿个消息端口的容量密度，由于合理使用了并发散列表等数据结构，使得容量和处理性能可随内存容量、处理器核心数量以及网卡速率等硬件升级而线性扩展。

由于传统分布式协调服务的性能和容量等限制，在经典的分布式集群中，多以服务或节点作为单位来进行分布式协调和调度，同时尽量要求集群中的节点工作在无状态模式。服务节点无状态的设计虽然对分布式协调服务的要求较低，但同时也带来了集群整体性能低下等问题。

与此相反，BYPSS 可轻松实现每秒千万次请求的处理性能和万亿量级的消息端口容量。这就给分布式集群的精细化协作构建了良好的基础。与传统的无状态集群相比，基于 BYPSS 的精细化协作集群能够带来巨大的整体性能提升。

我们首先以最常见的用户和会话管理功能来说明：在无状态的集群中，在线用户并无自己的属主服务器，用户的每次请求均被反向代理服务随机地路由至集群中的任意节点。虽然 LVS、Nginx、HAProxy、TS 等主流反向代理服务器均支持基于 Cookie 或 IP 等机制的节点粘滞选项，但由于集群中的节点都是无状态的，因此该机制仅仅是增加了相同客户端请求会被路由到某个确

定后台服务器节点的概率而已，仍无法提供所有权保证，也就无法实现进一步的相关优化措施。

而得益于 BYPASS 突出的性能和容量保证，基于 BYPASS 的集群可以用户为单位来进行协调和调度（即：为每个活动用户注册一个端口），以提供更优的整体性能。具体的实现方式为：

1. 与传统模式一样，在用户请求到达反向代理服务时，由反向代理通过 HTTP Cookie、IP 地址或自定义协议中的相关字段等方式来判定当前请求应该被转发至哪一台后端服务器节点。若请求中尚无粘滞标记，则选择当前负载最轻的一个后端节点来处理。
2. 服务器节点在收到用户请求后，在本地内存表中检查该用户的属主是否为当前节点。
 - a) 若当前节点已是该用户属主，则由此节点继续处理用户请求。
 - b) 若当前节点不是该用户的属主，则向 BYPASS 发起 RegPort 请求，尝试成为该用户的属主。此请求应使用批量方式发起，以进一步提高网络利用率和处理效率。
 - i. 若 RegPort 请求成功，说明当前节点已成功获取该用户的所有权，此时可将用户信息由后端数据库加载到当前节点的本地缓存中（应使用批量加载优化），并继续处理此用户相关请求。
 - ii. 若 RegPort 请求失败，说明指定用户正归于另一个节点管辖，此时应重新设置反向代理能够识别的 Cookie 等粘滞字段，将其指向正确的属主节点。并要求反向代理服务或客户端重试请求。

与传统架构相比，考虑到无状态服务也需要通过 MySQL、Memcached 或 Redis 等有状态技术来实现专门的用户和会话管理机制，因此以上实现并未增加多少复杂度，但是其带来的性能提升却非常巨大，对比如下：

项目	BYPASS HPC 集群	传统无状态集群
1 运维	省去用户和会话管理集群的部署和维护成本。	需要单独实施和维护用户管理集群，并为用户和会话管理服务提供专门的高可用保障，增加故障点、增加系统整体复杂性、增加运维成本。
2 网络	几乎所有客户请求的用户匹配和会话验证工作都得以在其属主节点的内存中直接完成。内存访问为纳秒级操作，对比毫秒级的网络查询延迟，性能提升十万倍以上。同时有效降低了服务器集群的内部网络负载。	每次需要验证用户身份和会话有效性时，均需要通过网络发送查询请求到用户和会话管理服务，并等待其返回结果，网络负载高、延迟大。 由于在一个典型的网络应用中，大部分用户请求都需要在完成用户识别和会话验证后才能继续处理，因此这对整体性能的影响很大。
3	因为拥有了稳定的属主服务器，而用户在某个时间	无专属服务器，用户请求随机到达集群

项目	BYSS HPC 集群	传统无状态集群
缓存	<p>段内总是倾向于重复访问相同或相似的数据（如自身属性，自己刚刚发布或查看的商品信息等）。因此服务器本地缓存的数据局部性强、命中率高。</p> <p>相较于分布式缓存而言，本地缓存的优势非常明显：</p> <ol style="list-style-type: none"> 省去了查询请求所需的网络延迟，降低了网络负载（详见本表“项目 2”中的描述）。 直接从内存中读取已展开的数据结构，省去了大量的数据序列化和反序列化工作。 仅由属主节点缓存对应数据也避免了分布式缓存与 DB 之间的不一致问题，提供了强一致保证。 <p>与此同时，如能尽量按照某些规律来分配用户属主，还可进一步地提升服务器本地缓存的命中率。例如：</p> <ol style="list-style-type: none"> 按租户（公司、部门、站点）来分组用户； 按区域（地理位置、游戏中的地图区域）来分组用户； 按兴趣特征（游戏战队、商品偏好）来分组用户。 <p>等等，然后尽量将属于相同分组的用户优先分配给同一个（或同一组）服务器节点。显而易见，选择合适的用户分组策略可极大提升服务器节点的本地缓存命中率。</p> <p>这使得绝大部分与用户或人群相关的数据均可在本地缓存命中，不但提升了集群整体性能，还消除了集群对分布式缓存的依赖，同时大大降低了后端数据库的读负载。</p>	<p>中的任意服务器节点；本地缓存命中率低；各节点重复缓存的内容多；需要以更高的成本为代价依赖分布式缓存。</p> <p>后端数据库服务器的读压力高，要对其进行分库分表、读写分离等额外优化。</p> <p>分布式缓存与 DB 之间存在无法避免的数据不一致问题（除非在分布式缓存与 DB 间使用 Paxos 等协议来保证一致性，但随之而来的高昂性能损失也将使分布式缓存失去意义——这比不用分布式缓存还慢）。</p>
4 更新	<p>由于所有权确定，能在集群全局确保任意用户在给定时间段内，均由特定的属主节点来提供服务。再加上现代服务器突发宕机故障的概率也较低。</p> <p>因此可以将用户属性中频繁变化但重要性或时效性较低的部分缓存在内存中，待积累一段时间后再批量更新至数据库。这可大大降低后端数据库服务器的写压力。</p>	<p>由于用户的每次请求都可能被转发到不同服务器节点来处理，因此无法实现累积写入优化和批量写入优化。后端数据库的写负担非常重。</p> <p>存在多个节点争抢更新同一条记录的问题，进一步加重了数据库负担。</p> <p>为此要对后端数据库进行额外的分库</p>

项目	BYPSS HPC 集群	传统无状态集群
	<p>例如：商城系统可能随着用户的浏览（比如每次查看商品）进程，随时收集并记录用户的偏好信息。若每次用户查看了新商品后，都需要即时更新数据库，则负载较高。再考虑到因为服务器偶发硬件故障导致丢失最后数小时商品浏览偏好数据完全可以接受，因此可由属主节点将这些数据临时保存在本地缓存中，每积累数小时再批量更新一次数据库。</p> <p>再比如：MMORPG 游戏中，用户的当前位置、状态、经验值等数据随时都在变化。属主服务器同样可以将这些数据变化积累在本地缓存中，并以适当的间隔（比如：每 5 分钟一次）批量更新到数据库中。</p> <p>这不但极大地降低了后端数据库要执行的请求数量，而且将多个用户的数据在一个批量事务中打包更新也大大减少数据库操作时的磁盘刷新动作，进一步提升了效率。</p> <p>此外，由专门的属主节点发起对用户属性的更新也避免了无状态集群中多个节点同时请求更新同一对象时的争抢问题，进一步提高了数据库性能。</p>	<p>分表等优化，还会引发“需要由应用层来自行处理分布式事务”等副作用。</p>
5 推送	<p>由于同一用户发起的所有会话均被集中在同一个属主节点内统一管理，因此可非常方便地向用户推送即时通知消息（Comet）。</p> <p>若发送消息的对象与消息接消息的收用户处于相同节点，则可直接将该消息推送给收件人麾下的所有活动会话。</p> <p>否则只需将消息定向投递到收件人的属主节点即可。消息投递可使用 BYPSS 实现（直接向收件人对端口发消息，应启用消息批量发送机制来优化），亦可通过 BYDMQ 等专用的高性能消息中间件来完成。</p> <p>若按照本表“项目 3”中描述的方法，优先将关联更紧密的用户分配到相同属主节点的话，则可大大提升消息推送在相同节点内完成的概率，此举可显著降低服务器间通信的压力。</p> <p>因此我们鼓励针对业务的实际情况来妥善定制用</p>	<p>由于同一用户的不同会话被随机分配到不同节点处理，因此需要开发、部署和维护专门的消息推送集群，同时专门确保该集群的高性能和高可用性。</p> <p>这不但增加了开发和运维成本，而且由于需要将每条消息先投递到消息推送服务后，再由该服务转发给客户端，因此也加重了服务器集群的内部网络负载，同时也加大了用户请求的处理延迟。</p>

项目	BYPSS HPC 集群	传统无状态集群
	<p>户分组策略，合理的分组策略可实现让绝大部分消息都在当前服务器节点内本地推送的理想效果。</p> <p>例如：对游戏类应用，可按地图对象分组，将处于相同地图副本内的玩家交由同一属主节点进行管理——传统 MMORPG 中的绝大部分消息推送都发生在同一地图副本内的玩家之间（AOI 范围）。</p> <p>再比如：对于 CRM、HCM、ERP 等 SaaS 应用来说，可按照公司来分组，将隶属于相同企业的用户集中到同一属主节点上——很显然，此类企业应用中，近 100%的通信都来自于企业内部成员之间。</p> <p>这样即可实现近乎 100%的本地消息推送率，达到几乎消除了服务器间消息投递的效果，极大地降低了服务器集群的内部网络负载。</p>	
6 平衡	<p>集群可使用主被动负载平衡相结合的手段进行调度。</p> <p>被动平衡：集群中的每个节点均会定期将其麾下不再活跃的用户和会话卸载掉，同时批量通知 BYPSS 服务释放这些用户所对应的端口。此算法实现了宏观上的负载平衡（以较长的时间周期来说，集群是平衡的）。</p> <p>主动平衡：集群会通过 BYPSS 服务推选出负载平衡协调节点，该节点连续监视集群中各个节点的负载情况，并主动发出指令进行负载调度（如：要求 A 节点将其麾下 5000 位用户的所有权转移给 B 节点）。不同于宏观层面的被动平衡，主动平衡机制可以在更短的时间片内，以迅捷的反应速度来达成集群的快速配平。</p> <p>主动平衡通常在集群中的部分节点刚刚从故障中恢复（因此处于空载状态）时效果明显，它比被动平衡反应更加迅速。如：在一个多活 IDC 集群中，某个 IDC 的光缆故障刚刚被修复而恢复上线时。</p>	<p>若启用了反向代理中的节点粘滞选项，则其负载平衡性与 BYPSS 集群的被动平衡算法相当。</p> <p>若未启用反向代理中的节点粘滞选项，则在从故障中恢复时，其平衡性低于 BYPSS 主动平衡集群。与此同时，为了保证本地缓存命中率等其它性能指标不被过分劣化，管理员通常不会禁用节点粘滞功能。</p> <p>另外，SOA 架构的多个服务间，容易产生负载不平衡，出现一些服务超载，另一些轻载的情况，nano-SOA 集群则无此弊端。</p>

值得一提的是，这样的精准协作算法并不会造成集群在可用性方面的任何损失。考虑集群中的某个节点因故障下线的情况：此时 BYPSS 服务会检测到节点已下线，并自动释放属于该节点的所有用户。待其用户向集群发起新请求时，该请求会被路由到当前集群中，负载最轻的节点。这个新节点将代替已下线的故障节点，成为此用户的属主，继续为该用户提供服务（见前文中的

步骤 2-b-i)。此过程对用户透明，不需要在客户端中加入额外的处理逻辑。

上述讨论以几乎所有网络应用中都会涉及的用户和会话管理功能为例，为大家展示了 BYPSS HPC 集群精细协调能力的优势。但在多数真实应用中，并不只有用户管理功能。除此之外，应用中通常还会包含可供其用户操作的其它对象。例如在优酷、土豆、youtube 等视频网站中，除了用户以外，至少还有“可供播放的视频”这种对象。

下面我们就以“视频对象”为例，探讨如何使用 BYPSS 的精细化调度能力来大幅提升集群性能。

在这个假想的视频点播类应用中，与前文描述的用户管理功能类似，我们首先通过 BYPSS 服务为每个**活动的视频对象**选取一个属主节点。其次，我们将视频对象的属性分为以下两大类：

1. **普通属性**：包含了那些较少更新，并且尺寸较小的属性。如：视频封面和视频流数据在 S3/OSS 等对象存储服务中的 ID、视频标题、视频简介、视频标签、视频作者 UID、视频发布时间等等。这些属性均符合读多写少的规律，其中大部分字段甚至在视频正式发布后就无法再做修改。

对于这类尺寸小、变化少的字段，可以将其分布在当前集群中，各个服务器节点的本地缓存内。本地缓存有高性能、低延迟、无需序列化等优点，加上缓存对象较小的尺寸，再配合用户分组等进一步提升缓存局部性的策略，可以合理的内存开销，有效地提升应用整体性能（详见下文）。

2. **动态属性**：包含了所有需要频繁变更，或尺寸较大的属性。如：视频的播放次数、点赞次数、差评次数、平均得分、收藏数、引用次数，以及视频讨论区内容等。

我们规定这类尺寸较大（讨论区内容）或者变化较快（播放次数等）的字段只能由该视频对象的属主节点来访问。其它非属主节点如需访问这些动态属性，则需要将相应请求提交给对应的属主节点来进行处理。

意即：通过 BYPSS 的所有权选举机制，我们将那些需要频繁变更（更新数据库和执行缓存失效），以及那些占用内存较多（重复缓存代价高）的属性都交给对应的属主节点来管理和维护。这就形成了一套高效的分布式计算和分布式缓存机制，大大提升了应用整体性能（详见下文）。

此外，我们还规定对视频对象的任何写操作（不管是普通属性还是动态属性）均必须交由其属主来完成，非属主节点只能读取和缓存视频对象的普通属性，不能读取动态属性，也不能执行任何更新操作。

由此，我们可以简单地推断出视频对象访问的大体业务逻辑如下：

1. 在普通属性的读取类用户请求到达服务器节点时，检查本地缓存，若命中则直接返回结果，否则从后端数据库读取视频对象的普通属性部分并将其加入到当前节点的本地缓存

中。

2. 在更新类请求或动态属性读取类请求到达服务器节点时，通过本地内存表检查当前节点是否为对应视频对象的属主。
 - a) 若当前节点已是该视频的属主，则由当前节点继续处理用户请求：读操作直接从当前节点的本地缓存中返回结果；写操作视情形累积在本地缓存中，或直接提交给后端数据库并更新本地缓存。
 - b) 若当前节点不是该视频的属主，但在当前节点的名称解析缓存表中找到了与该视频匹配的条目，则将当前请求转发给对应的属主节点。
 - c) 若当前节点不是该视频的属主，同时并未在当前节点的名称解析缓存表中查找到对应的条目，则向 BYPSS 发起 RegPort 请求，尝试成为该视频的属主。此请求应使用批量方式发起，以进一步提高网络利用率和处理效率。
 - i. 若 RegPort 请求成功，说明当前节点已成功获取该视频的所有权，此时可将视频信息由后端数据库加载到当前节点的本地缓存中（应使用批量加载优化），并继续处理此视频相关请求。
 - ii. 若 RegPort 请求失败，说明指定视频对象正归于另一个节点管辖，此时可将该视频及其对应的属主节点 ID 加入到本地名称解析缓存表中，并将请求转发给对应的属主节点来处理。

注意：由于 BYPSS 能够在端口注销时（无论是由于属主节点主动放弃所有权，还是该节点故障宕机），向所有对此事件感兴趣的节点推送通知。因此名称解析缓存表不需要类似 DNS 缓存的 TTL 超时淘汰机制，仅需在收到端口注销通知或 LRU 缓存满时删除对应条目即可。这不但能够大大增强查询表中条目的时效性和准确性，同时也有效地减少了 RegPort 请求的发送次数，提高了应用的整体性能。

与经典的无状态 SOA 集群相比，上述设计带来的好处如下：

项目	BYPSS HPC 集群	传统无状态集群
1 运维	基于所有权的分布式缓存架构，省去 Memcached、Redis 等分布式缓存集群的部署和维护成本。	需要单独实施和维护分布式缓存集群，增加系统整体复杂性。
2 缓存	普通属性的读操作在本地缓存命中，若使用“优先以用户偏好特征来分组”的用户属主节点分配策略，则可极大增强缓存局部性，增加本地缓存命中率，降低本地缓存在集群中各个节点上的重复率。 正如前文所述，相对于分布式缓存而言，本地缓存有消除网络延迟、降低网络负载、避免数据结构频	无专属服务器，用户请求随机到达集群中的任意服务器节点；本地缓存命中率低；各节点重复缓存的内容多；需要以更高的成本为代价依赖额外的分布式缓存服务。 后端数据库服务器的读压力高，要对其

项目	BYPSS HPC 集群	传统无状态集群
	<p>繁序列化和反序列化等优点。</p> <p>此外，动态属性使用基于所有权的分布式缓存来实现，避免了传统分布式缓存的频繁失效和数据不一致等问题。同时由于动态属性仅被缓存在属主节点上，因此也显著提升了系统整体的内存利用率。</p>	<p>实施分库分表、读写分离等额外优化。</p> <p>此外，即使为 Memcached、Redis 等产品加入了基于 CAS 原子操作的 Revision 字段等改进，这些独立的分布式缓存集群仍无法提供数据强一致保证（意即：缓存中的数据与后端数据库里的记录无法避免地可能发生不一致）。</p>
3 更新	<p>由于所有权确定，能在集群全局确保任意视频对象在给定时间段内，均由特定的属主节点来提供写操作和动态属性的读操作等相关服务，再加上现代服务器突发宕机故障的概率也较低。</p> <p>因此可以将动态属性中频繁变化但重要性或时效性较低的部分缓存在内存中，待积累一段时间后再批量更新至数据库。这可大大降低后端数据库服务器的写压力。</p> <p>例如：视频的播放次数、点赞次数、差评次数、平均得分、收藏数、引用次数等属性都会随着用户点击等操作密集地变化。若每次发生相关的用户点击事件后，都需要即时更新数据库，则负载较高。而在发生“属主节点由于硬件故障宕机”等极端情况时，丢失几分钟的上述统计数据完全可以接受。因此，我们可以将这些字段的变更积累在属主节点的缓存中，每隔数分钟再将其统一地批量写回后端数据库。</p> <p>这不但极大地降低了后端数据库收到的请求数量，而且将多个视频的数据在一个批量事务中打包更新，也大大减少数据库操作时的磁盘刷新动作，进一步提升了效率。</p> <p>此外，由专门的属主节点单独发起对视频记录的更新也避免了无状态集群中多个节点同时请求更新同一对象时的争抢问题，进一步提高了数据库性能。</p>	<p>由于每次请求都可能被路由到不同服务器节点来处理，因此无法实现累积写入优化和批量写入优化。后端数据库服务器的写负担非常重。存在多个节点争抢更新同一条记录的问题，这进一步加重了数据库负担。</p> <p>为此要对后端数据库进行额外的分库分表等优化，还会引发“需要由应用层来自行处理分布式事务”等副作用。</p>
4 平衡	<p>集群可使用主被动负载平衡相结合的手段进行调度。</p> <p>被动平衡：集群中的每个节点均会定期将其麾下不</p>	<p>在从故障中恢复时，其平衡性低于 BYPASS 主动平衡集群。正常情况下则相差不多。</p>

项目	BYPSS HPC 集群	传统无状态集群
	<p>再活跃的视频对象卸载掉，同时批量通知 BYPSS 服务释放这些视频对应的端口。此算法实现了宏观上的负载平衡（以较长的时间周期来说，集群是平衡的）。</p> <p>主动平衡：集群会通过 BYPSS 服务推选出负载平衡协调节点，该节点连续监视集群中各个节点的负载情况，并主动发出指令进行负载调度（如：要求 A 节点将其麾下 10000 个视频对象的所有权转移给 B 节点）。不同于宏观层面的被动平衡，主动平衡机制可以在更短的时间片内，以迅捷的反应速度来达成集群的快速配平。</p> <p>主动平衡通常在集群中的部分节点刚刚从故障中恢复（因此处于空载状态）时效果明显，它比被动平衡反应更加迅速。如：在一个多活 IDC 集群中，某个 IDC 的光缆故障刚刚被修复而恢复上线时。</p>	<p>另外，SOA 架构的多个服务间，容易产生负载不平衡，出现一些服务超载，另一些轻载的情况，nano-SOA 集群则无此弊端。</p>

与前文提及的用户管理案例类似，上述精准协作算法不会为集群的服务可用性方面带来任何损失。考虑集群中的某个节点因故障下线的情况：此时 BYPSS 服务会检测到节点已下线，并自动释放属于该节点的所有视频对象。待用户下次访问这些视频对象时，收到该请求的服务器节点会从 BYPSS 获得此视频对象的所有权并完成对该请求的处理。至此，这个新节点将代替已下线的故障节点成为此视频对象的属主（见前文中的步骤 2-c-i）。此过程对用户透明，不需要在客户端中加入额外的处理逻辑。

通过上述两个案例可以看出，相对于传统的 Redis + MySQL 等现有分布式缓存 + DB 的集群架构，BYPSS 集群还存在如下额外优势：

1. **BYPSS 集群保证强一致性**：正如前文所述，BYPSS 提供强一致、多活 IDC 高可用的高可用和高性能集群计算(HAC+HPC)能力。而 Redi 等分布式缓存，即使加入了 Revision 字段并利用其实现原子(CAS)操作，仍然无法保证其与 DB 之间的强一致性，只能通过加入失效超时(TTL)等机制来缓解数据不一致带来的危害（但此举同时也引入了缓存击穿、缓存雪崩等其它问题）。
2. **BYPSS 集群可提供数万倍量级的网络和 CPU 性能提升**：BYPSS 集群以纳秒级的内存直接访问替代了毫秒级的 Redis 网络查询，性能提升几个数量级。同时避免了每次缓存访问时频繁的数据序列化和反序列化过程，显著降低了 CPU 和网络开销。
3. **BYPSS 集群可完全避免缓存击穿（失效）问题**：由于 BYPSS 集群中的每个对象均仅缓存在其属主节点上，仅该对象的唯一属主节点有权向 DB 发送访问请求，因此不可能发生缓存击穿（不会产生针对某个热点对象的并发 DB 查询）。

4. BYPSS 集群可**轻易避免缓存穿透**（无效）问题：由于 BYPSS 集群中的每个对象均有其唯一属主节点，因此很容易在属主节点上对无效（不存在等）的对象 ID 进行标记和过滤（例如：高性能 ID 并发散列集合、布隆过滤器等）。因此可以轻易避免缓存穿透问题。
5. BYPSS 集群可**完全避免缓存雪崩**（大面积失效）问题：首先，由于 BYPSS 集群中缓存与 DB 数据间是强一致的，因此无需设置一个缓存过期时间来缓和缓存数据不一致的危害。因此缓存雪崩的一大前提：大量缓存同时到期就不复存在了。其次，由于消除了专门的缓存集群，所有缓存数据均内嵌在 App 服务器内，自然也就不存在“缓存节点宕机”这种问题了。与此同时，BYPSS 提供的强一致多活 IDC 高可用集群能力可保证 App Server 集群不会出现大规模宕机事故（不可抗力除外，真发生不可抗事件导致集群整体下线的话，那跟缓存雪崩已不是一个级别的问题了）。

以上对“用户管理”和“视频服务”案例的剖析均属抛砖引玉。在实际应用中，BYPSS 通过其高性能、大容量等特征提供的资源精细化协调能力可适用于包括互联网、电信、物联网、大数据批处理、大数据流式计算等广泛领域。

5.2.4 白杨分布式消息队列（BYDMQ）

白杨分布式消息队列服务（BYDMQ，读作“by dark”）是一种强一致、高可用、高性能、高吞吐、低延迟、可线性横向扩展的分布式消息队列服务。可支持单点千万量级的并发连接以及单点每秒千万量级的消息转发性能，并支持集群的线性横向扩展。

BYDMQ 自身亦依赖 BYPSS 来完成其服务选举、服务发现、故障检测、分布式锁、消息分发等分布式协调工作。BYPSS 虽然也包含了高性能的消息路由和分发功能，但其主要设计目还是为了传递任务调度等分布式协调相关的控制类信令。而 BYDMQ 则专注于高吞吐、低延迟的大量业务类消息投递等工作。将业务类消息转移到 BYDMQ 后，可使 BYPSS 的工作压力显著降低。

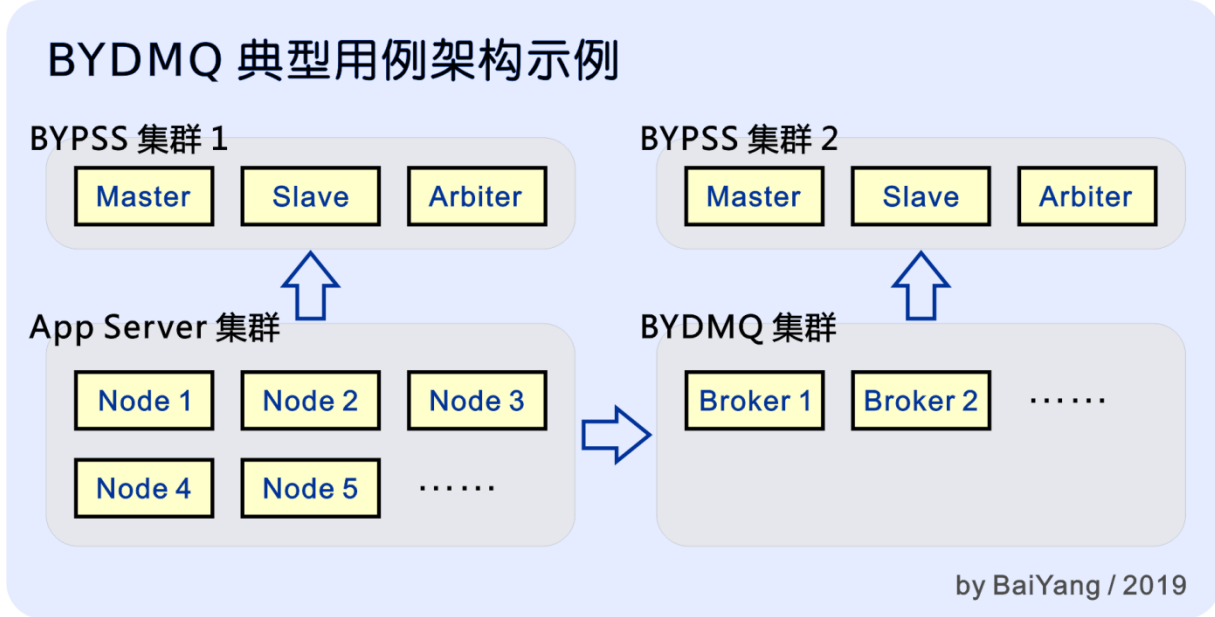


图 19

如图 19 所示,在典型用例中,BYDMQ 与 App Server 集群各自拥有一套独立的 BYPSS 集群,它们分别负责各自的分布式协调任务,App 集群依赖 BYPSS1 完成分布式协调,而其消息通信则依赖 BYDMQ 集群来完成。

不过在研发、测试环境,或业务量不大的生产环境中,也可以让 AppServer 和 BYDMQ 集群共享同一套 BYPSS 服务。另外需要指出的是,此处描述的“独立集群”仅是指逻辑上的独立。而从物理上说,即使是两个逻辑上独立的 BYPSS 集群,也可以共享物理资源。例如:一个 Arbiter 节点完全可以被多个 BYPSS 集群所共享;甚至两个 BYPSS 集群中的 Master、Slave 节点还可以互为主备,这样即简化了运维管理负担,又能够有效节约服务器硬件和能源消耗等资源。

在继续介绍 BYDMQ 的主要特性前,我们首先要澄清一个概念,即:消息队列(消息中间件, MQ)的可靠性问题。众所周知,“可靠的消息传递”包含三个要素——消息在投递过程中,要能够做到不丢失、不乱序和不重复才能称之为可靠。令人遗憾的是,这世间目前并不真正存在同时满足以上三个条件的消息队列产品。或者换句话说,我们目前尚无法在可接受的成本范围内实现同时满足以上三要素的消息队列产品。

要说明这个问题,请考虑以下案例:



图 20

如上图所示，在这一案例中，消息生产者由节点 A、B、C 组成，而消息消费者包含了节点 X、Y、Z，生产者与消费者之间通过一个消息队列连接。现在消息生产者已经生产了 5 条消息，并将它们依序成功提交到了消息队列。

在这样的情形下，我们来逐一讨论消息投递的可靠性问题：

★ **消息不丢失**：这点是三要素里最容易保证的。可拆分为两步来讨论：

- **存储可靠性**：可以通过将每条消息同步复制到消息队列服务中的其它节点（Broker）并确保落盘来保证；同时需要使用 Paxos、Raft 等分布式共识协议来确保多个副本间的一致性。但是显而易见，与无副本的纯内存方案相比，由于增加了磁盘 IO、网络复制以及共识投票等步骤，此方案会极大（数千甚至上万倍）地降低消息队列服务的性能。
- **ACK 机制**：在生产者向消息队列提交消息，以及消息队列向消费者投递消息时，均增加 ACK 机制来确认消息投递成功。发送方在指定时间内未收到 ACK 机制则重发（重新投递）消息。

以上两步措施虽可在很大程度保证消息不丢失（至少一次送达），但是可以看到其开销十分巨大，对性能的劣化非常显著。

与此同时，也应该注意到多副本间的故障检测和主从切换、以及消息收发时的超时重传等技术手段均会引入各自的延迟。而且其中每个步骤中引入的延迟通常都超过数秒钟。

在真正的用户场景中，这些额外的延迟使得“消息不丢失”的保证除了平白增加了巨大开销以外，大多没有任何实际用处：如今的用户在发起一个请求（如：打开一个链接、提交一个表单等）后，很少有耐心等待许久——若等待数秒后仍然没有响应，他们早就关闭页面离开或 F5 重新刷新了。

此时无论用户是关闭了页面还是重新发起了请求，已经延迟了几秒（甚至更久）才到达的那条消息（老请求）都已经没有了价值。不但如此，处理这些请求更是在白白消耗网络、运算和存储资源而已——因为其处理结果已经无人问津。

- ★ **消息不重复：**考虑前文提到的 ACK 机制：消费者在处理完一条消息后，需要向消息队列服务回复一条对应的 ACK 信令来确认该消息已被消费。例如：假设上图中的 1 号消息是一个转账请求，消息队列将该消息投递给节点 X 后，节点 X 必须在处理完该转账请求后，向消息队列服务发一条形如“ACK: Msg=1”的信令来告知消息队列服务，该消息已被处理。而 MQ 无法确保消息不重复的矛盾即在于此：

仍然按照上例中的假设，MQ 将消息 1 投递到了节点 X，但在规定的时间内却并未收到来自节点 X 的确认（ACK）信令，此时有很多种可能，例如：

- 消息 1 未被处理：由于网络故障，节点 X 并未收到该消息。
- 消息 1 未被处理：节点 X 收到了消息，但由于故障掉电而未能及时保存和处理该消息。
- 消息 1 已被处理：节点 X 收到了并处理了该消息，但由于故障掉电而未能及时向 MQ 服务返回对应的 ACK 信令。
- 消息 1 已被处理：节点 X 收到了并处理了该消息，但由于网络故障，对应的 ACK 信令未能成功返回至 MQ 服务。

等等。由此可见，在消息投递超时后，MQ 服务是无法得知该消息是否已被消费的。雪上加霜的是，由于前文所述的原因（不能让用户等太久），这个超时通常还要设置的尽量短，这就让 MQ 服务正确感知实际情况变得更加不可能。

此时为了保证消息不丢失，MQ 通常会假设消息未被处理，而重新分发该消息（例如在超时后，将消息 1 重新分发给节点 Y）。而这就势必无法再保证消息不重复了，反之亦然。

- ★ **消息不乱序：**从上例中可以看出，所谓“消息不乱序”是指 MQ 中的消息要按照先来后到，以“1、2、3、4、5”的顺序逐一被消费。要保证严格的不乱序，就要求 MQ 必须等待一条消息处理结束（收到 ACK）后，才能继续分发队列中的下一条消息，这至少带来了以下问题：

- 首先，MQ 中的多条消息无法被并行地消费。例如：MQ 无法将消息 1、2、3 同时分别派发给节点 X、Y、Z，这使得大量消费者节点长期处于饥饿（空闲）状态，甚至于即使在正在执行消息处理的节点上（比如节点 X）也会有大量处理器核心、SMT 单元等计算资源被浪费。
- 其次，在处理一条消息的过程中，所有其后续消息均只能处于等待状态。若一条消

息投递失败（超时），则在其“超时-重新投递”期间内，则会长时间阻塞其所有后续消息，使得它们无法被及时处理。

由此可见，保证消息严格有序会极大地影响系统整体的消息处理性能、增加硬件采购和运维成本，同时也会显著破坏用户体验。

由以上论述可知，现阶段尚无在合理成本下提供消息可靠传递的MQ产品问世。在此前提下，目前的解决方案主要是依赖 App Server 自身的业务逻辑（如：等幂操作、持久化状态机等）算法来克服这些问题。

反过来说：无论使用号称多“可靠”的MQ产品，现在的App业务逻辑中也均需要处理和克服上述种种消息投递不可靠的情形。既然MQ本质上做不到消息可靠，同时App也已经克服了这些不可靠性，那又何必再花费性能被劣化几千、甚至几万倍的代价来在MQ层实现支持“分布式存储 + ACK 机制”的方案呢？

基于上述思想，BYDMQ并不像RabbitMQ、RocketMQ等产品那样，提供所谓（实际无法达到）的“可靠性”保证。相反，BYDMQ采用“尽力送达”的模式，仅在确保不损失性能的前提下，尽可能地保证消息被可靠送达。

正如前文所述，由于App已经克服了消息传递过程中偶尔出现的不可靠。因此这样的设计抉择在极大提升了系统性能之余，并未实际增加业务逻辑的开发工作量。

基于上述设计理念的BYDMQ包含了以下特性：

- ★ 与BYPSS一样基于白杨应用支撑平台中的各优质跨平台组件实现，如：单点支持千万量级并发的网络服务器组件；支持多核线性扩展的并发散列表容器等等。这些优质高性能组件使得BYDMQ在可移植性、可扩展性、高容量、高并发处理能力等方面均拥有非常好的表现。
- ★ 与BYPSS一样，在客户端和服务端均拥有成熟的消息批量打包机制。支持连续消息的自动批量打包，大大提高网络利用率和消息吞吐量。
- ★ 与BYPSS一样支持pipelining机制：使得客户端无需等待一条指令的响应结果即可连续发送下一条指令，显著降低了命令处理延迟、提高了网络吞吐、有效增加了网络利用率。
- ★ 每个客户端（App Server）可注册一个专属MQ，并通过长连接 + 心跳的方式保活，对应的属主Broker（BYDMQ节点）也通过此长连接向客户端实时推送到达的消息（带有批量打包机制）。
- ★ 客户端通过一致性散列算法来推测一个MQ的属主（Broker），Broker在首次收到针对指定MQ的请求（如：注册、发消息等）时，将通过BYPSS服务来竞选成为该MQ的属主。若竞选成功则继续处理，竞选失败则引导客户端重新连接到正确的属主节点。

与此同时，BYDMQ 会通过 BYPSS 服务实时感知集群变化（如：现有 Broker 下线、新的 Broker 上线等），并将这些变化实时推送到每个客户端节点。这就保证了除非 BYDMQ 集群正在发生剧烈变化（大量 Broker 节点上线或下线等），否则通过一致性散列算法推定属主的准确性是非常高的，从而基本无需再次重定向请求。

此外，即使一致性散列算法推定错误，该 MQ 的实际属主也会被客户端节点自动记忆到本地快查表中，确保下次向这个 MQ 发送消息时能够直接投递到正确的 Broker。

这种由客户端直接向对应 MQ 之属主（Broker）投递消息的方法避免了服务器集群中的复杂路由和消息的多次中转，将消息投递的网络路由降至最简，极大地提升了消息投递的效率，有效降低了网络负载。

通过 BYPSS 来为 MQ 选举属主节点的做法则为集群提供了“每个 MQ 在全局范围内唯一”的一致性保证。与此同时，BYPSS 服务也负责在各个 Broker 节点之间分发一些控制类信令（如：节点上下线通知等等），使 BYDMQ 集群可以被更好地统一协调和调度。

- ★ 所有待分发的消息仅存储在对应 MQ 属主（Broker）节点的内存中，避免了写盘、复制、共识投票等大量无用开销。
- ★ 发送方可以为每条消息分别设置其生命期（TTL）和发送失败时的最大重试次数。可根据消息的类型和价值精确控制其消耗的资源。对时效性短或不重要的请求，可及时使其失效，避免在各个环节继续空耗资源，反之亦然。
- ★ 支持分散投递：当指定 MQ 所属客户端节点未上线，或其连接断开超过指定的时长后，此消息队列中的所有待投递消息将被随机发送到任意一个仍可正常工作的 MQ 中。

分散投递方案预期客户端（App Server）也是一个基于 BYPSS 的 nano-SOA 架构集群。此时若一个 App Server 节点因为运维、硬件故障或网络分区等原因下线，则该节点下辖的所有对象都会被管理该集群的 BYPSS 释放。

此时系统可随机将发往此节点的请求散布至其它仍正常工作的节点，可让这些目标节点通过 RegPort 重新获取该请求相关对象的所有权，并接替其已经下线的原属主节点继续完成处理。这就大大降低了在一个 App Server 节点异常下线的瞬间，与之相关请求的失败率，优化了客户体验。同时随机散布也使得下线节点麾下的对象被集群中仍然工作的其它节点均分，保证了集群的负载均衡。

综上，BYDMQ 通过在一定程度上牺牲了本就无法真正保证的消息可靠性，再配合消息打包、pipelining、属主直接投递等方式，极大地提升了消息队列服务的单点性能。同时得益于 BYPSS 为其引入的强一致、高可用、高性能的分布式集群计算能力，使其拥有了优异的线性横向扩展能力。加之其对每条消息的灵活控制、以及分散投递等特性，最终为用户提供了一款超高性能的高品质分布式消息队列产品。

5.3 安全隧道服务 (BYST)

白杨安全隧道服务 (BYST, 读作 “best”) 为用户提供了一种端到端的安全隧道服务, 主要支持以下特性:

1. 支持 PSK 和 PKI 鉴权: 可使用预分配密钥和公钥体系架构进行鉴权。
2. 支持数十种强加密算法: 支持数十种主流块加密和流式加密算法。
3. 支持实时数据压缩: 基于 lz4 算法的高性能实时数据压缩, 可通过配置选项开启。
4. 支持消息完整性校验: 使用散列算法计算校验和, 以确保消息的可靠性和完整性 (校验和与压缩后的数据均会被加密后再传输)。此功能可通过配置选项开启。
5. 数据混淆: 不同于 OpenVPN、SSL、SSH、IPSec 等现有方案, 本隧道协议无任何可被观测的特征。不知道密钥的第三方拦截者只能观察到随机二进制字节流, 难以通过任何有效手段侦测出通信双方正在使用本隧道协议。

不但自己无特征, BYST 还可以扮演 http 等合法白名单协议 (BYST over http), 确保在某些即使针对未知协议, 也仅依据流量限制进行阻断的极端环境下, 也可正常工作。

与此同时, BYST 还支持 HTTP chunked transfer encoding, 可以做到每个 package (几十 KB~几 MB) 仅增加 3.5 个 Bytes 左右的开销, 数据膨胀率低于 0.01%, 避免了笨重的 HTTP 包头导致数据膨胀的问题。远远低于 SSL/TLS、SSH、h2、WebSocket 等其它白名单协议所带来的额外通信开销, 在实现白名单通信的同时仍维持了极高的网络利用率。

6. 高效率: 得益于基于汇编优化、零内存拷贝以及 DMA + 硬件中断的纯异步 IO 通信组件, 即使在受限的硬件平台上, BYST 也可以满足高性能和高并发业务的严苛需求。

此外, 得益于我方成熟的 IO 批量打包、专利的分布式 N:M 动态连接池加速, 以及高性能数据实时压缩等算法, BYST 显著提升了网络利用率 (高载荷比) 和网络吞吐量。

与此同时, 相比于各类现有方案, 已将握手 (协商)、挥手、确认等环节进行了精心化简的 BYST 隧道协议亦明显具备更低的通信延迟。

7. 高性能和高可用集群: 支持基于 BYPSS (详见: 5.2.3 白杨消息端口交换服务 (BYPSS)) 的高性能和高可用多活 IDC 集群部署。

此外还支持多路径自动验活和路由自动切换等功能。可通过增加备用链路来显著提升服务整体可用性, 同时也支持多个链路之间的自动流量聚合及负载平衡等功能。

8. 线路级实时补偿: 独有的线路级实时跟踪和预测智能算法, 对各线路的丢包和抖动分别

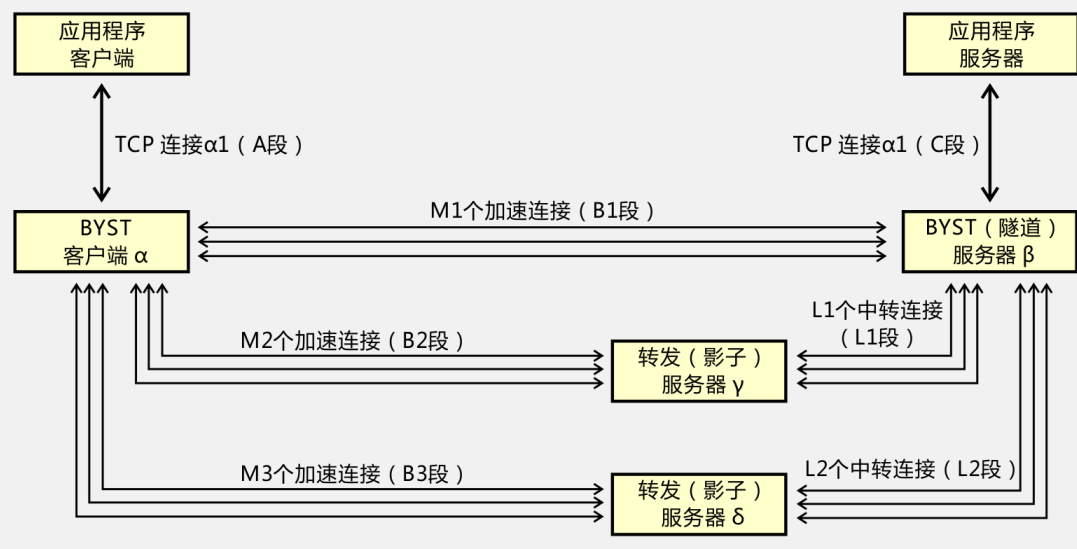
进行实时跟踪和补偿。

BYST 服务主要用于在 Internet、卫星、微波、SDH (MSTP)、以及跨区域光纤专线等广域或城域网环境中建立安全可靠的数据传输通道，通过强加密和强校验算法来保证数据的安全性、可靠性和完整性，并通过数据压缩降低带宽成本。同时还可通过难以被分析的混淆算法来保护其隧道通信不会被识别、拦截和屏蔽。

综上所述，BYST 主要带来如下技术优势：

1. 保障通信安全：提供带有强加密、强校验以及防回放攻击机制的安全通信隧道。
2. 打满带宽：得益于独有的 IO 自动批量打包、实时数据压缩，以及我方专利的分布式 N:M:N 动态连接映射加速算法，BYST 可做到打满用户带宽上限，显著提升 Site-to-Site Tunnel 通信性能。大量用户实测表明，在同城通信 (MAN) 环境中，BYST 仅启用单点 N:M 加速即可提升超过 6 倍的带宽吞吐；而在异地通信 (WAN) 环境中，BYST 的单点 N:M 加速则可实现高达 70 倍的吞吐性能提升。而分布式 N:M 加速则可在在此基础上随着分布式加速节点的增加而线性提升聚合吞吐率。

分布式 N:M 连接动态映射加速算法最简示例



3. 防止误杀：正如前文所述，为提高通信性能以及降低握手延迟，BYST 本身被设计为完全无特征的通信协议。此外，通过数据强加密、IO 自动批量打包、实时数据压缩、以及我方专利的分布式 N:M:N 动态连接映射加速算法等数据混淆机制，BYST 所承载的所有上层协议也将失去其可识别的特征。而 BYST over HTTP 的白名单协议通信支持则进一步保证了其优秀的防火墙通过性。
4. 代理鉴权：BYST 可通过 PSK、PKI 以及 EAL5+ 等级的安全智能卡硬件等技术提供安全可靠的双端鉴权访问，从而免除了上层应用间相互对接时还需要自行实现 CHAP、IKE、

LDAP 等复杂算法来完成鉴权的麻烦。