

# 分布式操作系统终极完整版

## 目录

分布式操作系统终极完整版 .....	1
第一章 分布式系统概论 .....	6
1、中间件在分布式系统中扮演什么角色? .....	6
2、解释(分布)透明性的含义,并且给出各种类型透明性的例子。 .....	6
3、在分布式系统中,为什么有时难以隐藏故障的发生以及故障恢复过程? .....	6
4、为什么有时候要求最大程度地实现透明性并不好? .....	6
5、什么是开放的分布式系统?开放性带来哪些好处? .....	6
6、请对可扩展系统的含义做出准确描述 .....	6
7、可以通过应用多种技术来取得可扩展性。请说出这些技术。 .....	6
8、多处理器系统与多计算机系统有什么不同? .....	6
9、在多计算机系统上的 256 个 CPU 组成了一个 16 X 16 的网格方阵。在最坏的情况下,消息的延迟时间有多长(以跳(hop)的形式给出,跳是结点之间的逻辑距离)? .....	6
10、现在考虑包含 256 个 CPU 的超立方体,最坏情况下消息的延迟有多长? .....	7
11、分布式操作系统与网络操作系统有什么不同? .....	7
12、请解释如何使用微内核将操作系统组织成客户-服务器的方式。 .....	7
13、请解释基于分页的分布式共享存储器系统主要有哪些操作。 .....	7
14、为什么要开发分布式共享存储器系统?你认为是什么问题给这种系统的高效率实现造成了困难? .....	7
15、请解释什么是分布式共享存储器系统中的伪共享。你有没有针对这个问题的解决方案? .....	7
16、由于存在错误,某个实验性的文件服务器有 3 / 4 的时间能够正常工作,而另外 1 / 4 的时间无法工作。如果要确保服务至少在 99%的时间可用,需要将该文件服务器复制多少次? .....	7
17、什么是三层客户-服务器体系结构? .....	7
18、纵向分布与横向分布有什么不同? .....	7
19、考虑一个进程链,该进程链由进程 P1, P2, ..., Pn 构成,实现了一个多层客户-服务器体系结构。进程 Pi 是进程 Pi+1 的客户, Pi 只有得到 Pi+1 的应答之后才能向 Pi-1 发出应答。如果考虑到进程 P1 的请求-应答性能,这种组织结构主要存在什么问题? .....	7
20、比较分布式操作系统、多处理机操作系统、网络操作系统的特点和应用范围。 .....	7
21、比较单体型系统和微内核型系统的优缺点。 .....	8
22、解释透明性问题对系统和用户的 .....	8
23、给出不同类型透明性的例子。 .....	8
24、举例说明最高程度的透明性对于某些应用不一定是必须的。 .....	8
25、说明实现可伸缩性的主要技术。 .....	9
26、给出一个虚拟组织的实际例子。 .....	9
第二章 通信 .....	9
1、在许多分层协议中,每一层都有自己的报头。如果每个消息首部都只有单个报头,其中包含了所有控制信息,无疑会比使用单独的多个报头具有更高的效率。为什么不这么做? .....	9
2、为什么传输层通信服务常常不适于构建分布式应用程序? .....	9
3、一种可靠的多播服务允许发送者向一组接收者可靠地传递消息。这种服务是属于中间件层还是更低层的一部分? .....	9
4、考虑一个带有两个整型参数的过程 incr。该过程将两个参数的值分别增加 1。现在假定调用它时使用的两个参数是同一个变量,比如 incr(i, i)。如果 i 的初始值是 0,在执行引用调用之后 i 将变为什么值?如果使用复制——还原调用呢? .....	9
5、C 语言中有一种称为联合(union)的构造,其中的记录(在 C 语言中称作结构)的字段可以用来保存几种	

可能值中的一个。在运行时，没有可靠的办法来分辨其中保存的是那一个值。C 的这种特性是否与远程过程调用有某些相似之处？请说明理由。 .....	9
6、处理 RPC 系统中参数转换的一种方法是，每台机器以自己系统使用的表示方式来发送参数，由另一方在必要的情况下进行转换。可以通过首字节中的代码来表示发送消息机器所用的系统。然而，由于要在首个字中找到开头的字节这本身也是一个问题，这种方法能行得通吗？ .....	10
答：首先，当一台机器发送字节 0 时，消息肯定已经送到。因此目标机器可以访问字节 0，而代码就在消息里面。这种方法不考虑字节是高位优先还是低位优先的字节。另一个方法是将代码放在第一个单词的所有字节中。因此不管检查的是哪一个字节，代码都能被找到。 .....	10
7、假定客户通过异步 RPC 对服务器进行调用，随后等待服务器使用另一异步 RPC 返回结果。这种方法与客户端执行常规的 RPC 有没有什么不同？如果使用的是同步 RPC 而不是异步 RPC，情况又如何呢？ .....	10
8、在 DCE 中，服务器在守护程序中注册自身。如果换一种方法，也可以总是为它分配一个固定的端点，然后再指向服务器地址空间中对象的引用中就可以使用该端点。这种方法的缺陷在哪里？ .....	10
答：这种方法的主要缺陷是向服务器分配对象变得很难。另外，许多端点而不止一个需要被修复。如果机器中很多都是服务器，分配固定端点不是一个好办法。 .....	10
9、给出一种用来让客户端绑定到暂时远程对象的对象应用的实现示例。 .....	10
10、Java 和其他语言支持异常处理，当错误发生时会引起异常。如何在 RPC 和 RMI 中实现异常处理？ .....	10
11、将静态和动态 RPC 区分开来有用吗？ .....	10
12、某些分布式中间件系统的实现是完全基于动态方法调用的。甚至连静态的调用也被编译成动态的调用。这种方法的优点在哪里？ .....	10
13、描述一下客户端和服务端之间使用套接字的无连接通信是如何进行的。 .....	10
14、说明 MPI 中 <code>mp_based</code> 原语和 <code>mp_isend</code> 原语之间的区别。 .....	10
15、假定只能使用暂时异步通信原语，再加上异步 <code>receive</code> 原语，如何实现用于暂时同步通信的原语？ .....	11
16、假定只能使用暂时同步通信原语，如何实现用于暂时异步通信的原语？ .....	11
17、通过 RPC 实现持久化异步通信有意义吗？ .....	11
18、在本章中我们讲过，为了自动启动一个进程以从输入队列中获取消息，常常要使用守护程序来监视输入队列。请给出一种不使用守护程序的实现方法。 .....	11
19、IBM MQSeries 以及许多其他消息队列系统中的路由表是人工配置的。描述一种自动完成配置工作的简单方法。 .....	11
20、如何将持久异步通信加入到基于远程对象 RMI 的通信模型中去？ .....	11
21、在持久通信中，接收者一般拥有自己的本地缓冲区，如果接收者不在运行状态，可以将消息放入该缓冲区中去。为了创建这种缓冲区，必须指定它们的大小。分成两方进行辩论：一方认为这种指定缓冲区大小的行为是可取的，而另一方反对这种指定大小的行为。 .....	11
22、请说明为什么暂时同步通信在可扩展性方面存在固有的问题，以及如何解决这些问题。 .....	11
23、给出一个将多播应用于离散数据流的例子。 .....	11
24、当一组计算机组成一个逻辑上或物理上的环时，如何确保传输延迟不超过允许的最大端到端延迟时间？ .....	11
25、当一组计算机组成一个（逻辑上或物理上的）环时，如何确保传输延迟不小于允许的最小端到端延迟时间？ .....	11
26、想象一下，某个令牌存储桶规范说明中的最大数据单元大小是 1000B，令牌存储桶速率是 10MB/s，令牌存储桶的大小是 1MB，而最大传输速率是 50MB/s。突发的传输可以以最大传输速率持续传输多长时间？ .....	12
27、在这道练习题中，要求实现一个使用 RPC 简单客户——服务器系统。服务器提供一个名为 <code>next</code> 的过程，该过程接受一个整型数作为输入，并且返回紧接着该整型数之后的一个数。编写一个客户使用的名为 <code>next</code> 的存根，它的任务是使用 UDP 将参数发送给服务器，随后等待服务器响应，如果过长时间未收到响应就认为超时了。服务器过程必须在一个公开的端口监听并接受请求，完成这些请求所要求的操作，然后返回结果。 .....	12
答：无。 .....	12
28、举出一个分布式系统实例，该系统采用以数据为中心的软件体系结构。 .....	12
29、什么是三层客户-服务器体系结构，举出一个实际应用案例。 .....	12
30、在结构化覆盖网络中，如果消息是根据覆盖的拓扑结构来路由的，讨论其优缺点。 .....	12
31、在非结构化覆盖 P 网络中，比较洪泛式搜索算法与随机游走搜索算法的性能优缺点。 .....	12

32、设一个 P2P 网络中, N 为总的节点个数, r 为每个数据项被随机保存的节点的个数, m 为每个节点的平均邻居个数。现需搜索一个数据 d, 分别采用洪泛式搜索算法和随机游走算法, 在下列条件下估算出这两种算法的搜索规模 S 和搜索步数 k: .....	13
33、请从技术角度来解释为什么 BitTorrent 中的 tit-for-tat 策略比因特网中的文件共享好的多? .....	13
34、简述 p2p 网络的特点, 并分析如果采用集中型拓扑结构构建 p2p 网络, 其主要的优缺点是什么? .....	13
第三章 分布式进程管理 .....	14
1. 比较使用单线程文件服务器读取文件和使用多线程服务器读取文件有什么不同。花费 15ms 来接收请求、调度该请求并且完成其它必须的处理工作, 假定需要的数据存放在主存储器的缓存中。如果需要磁盘操作, 就需要额外多花 75ms 在磁盘操作的过程中线程处于睡眠状态。如果服务器采用单线程的话, 它每秒能处理多少个请求? 如果采用单线程呢? .....	14
2. 对服务器进程中的线程数目进行限制有意义吗? .....	14
3、在文中我们描述了一个多线程的文件服务器, 说明了为什么它比单线程服务器和有限状态服务器更好。有没有这样的环境, 在其中使用单线程服务器会更好? 给出这种环境的例子。 .....	14
4、将轻量级进程与单个线程静态关联起来并不好, 为什么? .....	14
5、如果每个进程只使用单个轻量级进程也不好, 为什么? .....	14
6、描述一种使用与可运行线程数目相等的轻量级进程的方法。 .....	14
7、本章解释过, 代理可以通过调用所有副本来支持引用透明。能否对(服务器端的)对象的副本进行调用? .....	14
8、通过生成进程来构建并发服务器与使用多线程服务器相比有优点也有缺点。给出部分优点和缺点。 .....	14
9、粗略地设想一种多线程服务器的设计, 该服务器必须使用套接字作为面对底层操作系统的传输级接口, 以支持多种协议。 .....	14
10、如何防止应用程序绕过窗口管理器破坏屏幕显示。 .....	15
11、解释对象适配器的概念。 .....	15
12、举出几个用于支持持久性对象的对象适配器设计方面的问题。 .....	15
13、改变对象适配器示例中的 thread_per_object 过程, 使用单个线程来处理受该适配器控制的所有对象。 .....	15
14、维护到客户的 TCP/IP 连接的服务器是状态相关的还是状态无关的。 .....	15
15、想象一下, 某个 web 服务器维护一个列表, 该列表中的内容是 IP 地址与该地址最近访问过的 web 页的映射关系。当一个客户连接到该服务器的时候, 该服务器在列表中查找该客户, 如果找到的话就返回注册过的页面。这个服务器是状态相关还是状态无关的? .....	16
16、Java RMI 对代码迁移依赖到何种程度? .....	16
17、在 UNIX 系统中, 可以通过让进程在远程机器上派生出一个子进程来支持强可移植性。请说明这种机制的工作机理。 .....	16
18、图 3.13 指出, 强可移植性不能与在目标进程中已迁移代码的执行结合在一起。请举出一个反例。 .....	16
19、考虑某个进程 P, 它请求访问与自己位于同一台机器上的本地文件 F。在 P 迁移到另一台机器上以后, 它还需要访问 F。如果文件对机器的绑定是紧固的, 如何实现对 F 的系统级引用? .....	16
20、D' Agents 系统中的每个代理都由单独的进程实现。代理通信的主要方式是共享文件以及消息传递, 而文件无法在机器之间传递。如果按照 3.4 节中给出的移动框架结构, 代理的那一部分状态包含了资源段? .....	16
21、将 D' Agents 的体系结构与 FIPA 模型所实现的代理平台的体系结构进行比较。 .....	16
22、代理通信语言(ACL)在哪些方面与 OSI 模型相符? .....	16
23、在系统顶层实现代理通信语言以进行电子邮件处理时, 这中代理通信语言又有哪些方面与 OSI 模型相符? 这种方案有什么优点? .....	16
24、为什么 ACL 消息中通常必须指定实体(ontology)? .....	16
25、有一个文件服务器, 当所需的数据放在缓存中时, 接受、分派或处理一个请求需要花费 15ms; 但在 1/3 的情形下, 还要做磁盘操作, 需要 75ms, 线程在这时处于睡眠状态。如果使用单线程, 每秒能处理多少请求? 而如果使用多线程呢? .....	16
26、假设完全在用户空间实现线程, 运行系统每隔 1 秒得到 1 个时钟中断。如果发生时钟中断时, 某个线程正在运行系统中执行。这时会出现什么问题? 给出一种解决方案。 .....	17
第四章 分布式通信管理 .....	17

1. 为什么传输层通信对建立分布式应用来讲不合适? .....	17
2. 编写并运行一个实现文件操作的 RPC 程序。 .....	17
3. 过程 incr 有两个整形参数。该过程将每个参数分别加 1。假设在调用 incr 时使用同一个变量, 如 incr(i, i)。若 i 的初值为 0, 如果按 call-by-reference 方式调用, i 值为多少?.....	17
4. 假设完成一个空 RPC 的时间是 1.0ms, 此外, 传输每 1K 数据需要 1.5ms。如果用 1 个 32K 的 RPC, 从一个文件服务器中读取 32K 数据需要多少时间? 如果用 32 个 1K 的 RPC, 又需要多少时间? .....	17
5. 假定对服务器做异步式 RPC 调用, 并等待直到服务器用另一个异步式 RPC 返回。这种方法与客户执行一个正常的 RPC 有何区别? .....	18
6. 解释为什么瞬时的同步通信天生具有可伸缩性问题, 以及如何解决。 .....	18
7. 在持久性通信中, 接受程序需要有保存消息的局部缓冲区。在创建该缓冲区时, 需要设定它的大小。给出这样做的理由。在给出不该这样做的理由。 .....	18
8. 在 IBM MQSeries 的路由表中是手工配置的, 设计一种自动配置方法。 .....	18
第五章 同步 .....	18
1、说出至少三种可在 WWV 广播时间和在分布式系统中处理机设置内部时钟之间引入的延迟源。 .....	18
2、考虑分布式系统中的两台机器的行为。这两台机器的时钟假设每毫秒滴答 1000 次。但实际上只有一台是这样的, 另一台每毫秒滴答 990 次。如果 UTC 每分钟更新一次, 那么时钟的最大偏移量是多大? .....	18
3、向图 5.7 中加入一个与消息 A 并发的新消息, 即它既不发生在 A 的前面也不发生在 A 的后面。 .....	18
4、要使用 Lamport 时间戳实现全序多播, 是不是每个消息都必须被严格地确认? .....	18
5、考虑一个消息只是以它们被发送的顺序被传递的通信层。给出一个不需要这种限制的例子。 .....	18
6、假设两个进程同时检测到协调者崩溃了, 并且它们都使用欺负算法主持一个选举。这时将发生什么。 ....	18
7、在图 5.12 中, 我们有两个 ELECTION 消息同时在循环, 当对它们都没有什么不利影响时, 杀掉一个将更好。设计这样一个算法, 该算法在不影响基本的选举操作算法的情况下完成此项任务。 .....	18
8、许多分布式算法需要使用协调进程。讨论一下, 这样的算法实际上可以在什么程度上被看作是分布式的? .....	19
9、在集中式互斥方法中 (图 5.13), 协调者在收到一个进程释放它独占访问的临界区的消息后, 通常给等待队列中的第一个进程授权, 允许其访问临界区。请给出协调者另一个可能的算法。 .....	19
10、请再考虑图 5.13, 假设协调者崩溃了, 这会必然导致整个系统瘫痪吗? 如果不会, 那么在什么情况下会呢? 有什么方法避免了这个问题的发生, 使系统能够忍受协调者的崩溃? .....	19
11、Ricart 和 Agrawala 算法会有这样的问题: 如果一个进程崩溃, 并且没有对另一个要求进入临界区的进程的请求作应答, 没有应答意味着拒绝请求, 我们建议所有的请求应用立即被应答, 以便很容易地检测到崩溃的进程, 是否存在一些情况, 即使使用这种方法也还不足够? 请讨论。 .....	19
12、如果我们假设该算法可以在支持硬件广播的 LAN 中实现, 那么图 5.16 中的何时将如何改变? .....	19
13、分布式系统可能有多个互相独立的临界区, 假设进程 0 想进入临界区 A 而进程 1 想进入临界区 B, Ricart 和 Agrawala 的算法会导致死锁吗? 请解释原因。 .....	19
14、在图 5.17 中, 我们看到一种用磁带实现对存货表自动更新的方法。因为磁带可以很容易地用磁盘上的一个文件来模拟, 你认为为什么不使用这种方法 (用磁盘文件代替) 呢? .....	19
15、在图 5.25 (D) 中有三个调度策略, 两个合法, 一个非法的。对于同一个事务处理, 给出临终 x 所有的可能值和合法的与非法的状态的一个完整表。 .....	19
16、当一个私有工作空间用于实现事务处理时, 可能需要将大量的文件索引复制到父辈工作区, 怎样在不引入竞争条件下实现这种操作? .....	19
17、给出一个完整的算法判定一个试图锁住一个文件的操作是否成功。要同时考虑到读锁和写锁, 以及文件被解锁, 读锁写或写锁写的可能性。 .....	19
18、用锁定的方法实现并发控制的系统通常区分读锁和写锁, 如果一个进程已经得到了一个读锁但现在又想将它换成写锁, 将会怎么样呢? .....	19
19、在分布式事务中使用时间戳排序, 假设写操作 write(T1, x) 被传给数据管理器, 因为仅有的可能发生冲突的操作 write(T2, x) 的时间戳较早, 为什么让调度管理器推迟传递 write(T1, x) 直到 T2 完成呢? .....	20
20、乐观并发控制与使用时间戳相比是更严格还是更不严格呢? 为什么? .....	20
21、使用时间戳来进行并发控制能保证串行性吗? 请讨论。 .....	20
22、我们常说事务被中止时, 一切都恢复到它以前的状态, 就好像事务从来没有发生一样, 实现上我们在说谎。	

请给出一个重新恢复一切是不可能的例子。 .....	20
23 举例说明, 需要进一步解析实体 E 的地址, 才能真正访问 E。 .....	20
1、对共享的 JAVA 对象的访问可以通过将其方法声明为同步的而被串行化。当这种对象被复制时, 这种方法足以保证访问的串行化吗? .....	21
2、对于第一章所讨论的监视器, 如果允许在一个复制的监视器中阻塞多个线程, 那么给条件变量发信号时, 需要保证什么? .....	21
3、请用自己的语言解释实际考虑一致性模型的主要原因。 .....	21
4、请解释 DNS 如何进行复制, 以及它实际运行很好的原因。 .....	21
5、讨论一致性模型时, 我们经常提及软件和数据存储间的合约。为什么需要这一合约? .....	21
6、线性化假设存在一个全局时钟。但是, 我们已经在严格一致性中指出, 这种假设对于大多数分布式系统都是不现实的。线性化可以应用于物理分布的数据存储吗? .....	22
7、如果多处理器使用单一总路线, 那么可以实现严格一致的存储器吗? .....	22
8、对于图 6.5 (B), 为什么 $W_1(x) \wedge R_2(x) \wedge \neg R_3(x) \wedge a$ 是非法的? .....	22
9、在图 6.7 中, 000000 是仅满足 FIFO 一致的分布式共享存储器的合法输出吗? 请解释你的答案。 .....	22
10、在图 6.8 中, 001110 是顺序一致的存储器的合法输出吗? 请解释你的答案。 .....	22
11、我们在 6.2.2 节的最后讨论了一个形式模型, 该模型规定顺序一致的数据存储上的每个操作集合可以用一个字符串 H 表示, 从这个字符串 H 可以衍生出所有的单独进程的顺序序列。对于图 6.9 中的进程 P1 和 P2, 给出 H 所有可能的值。忽略进程 P3 和 P4, 也不包括它们在 H 的操作。 .....	22
12、在图 6.13 中, 一个顺序一致的存储器允许 6 种可能的语句交叉。请列出这 6 种所有可能。 .....	22
13、通常认为弱一致性模型给程序开发人员强加了负担。这一命题在哪个程度上确实是正确的? .....	22
14、在分布式共享存储系统中的释放一致性的大多数实现中, 共享变量是在执行释放操作时被同步的, 而不是在执行获取操作时被同步的。那么, 为什么还需要获取操作呢? .....	22
15、ORCA 提供顺序一致性还是入口一致性? 请解释你的答案。 .....	22
16、使用序列器并为维持主动复制的一致性而进行的全序的广播是否违背系统设计中的端到端理论? .....	22
17、您会选择哪种类型的一致性来实现电子股票市场? 请解释您的答案。 .....	22
18、如果要将一个移动用户的个人邮箱作为广域范围的分布式数据库的一部分来实现, 哪种类型的以客户为中心的一致性最为合适? .....	22
19、请描述出一个用于显示刚被更新的 Web 页面的写后读一致性的简单实现。 .....	23
20、请给出以客户为中心的一致性易于导致写操作冲突的案例。 .....	23
21、使用租用时, 客户和服务器的时钟分别紧密同步是否是必要的? .....	23
22、考虑一个用于保证分布式数据存储上顺序一致性的非阻塞的主机备份协议。这样的数据存储总能提供写后读一致性吗? .....	23
23、为了使主动复制正常工作, 所有操作必须以相同的顺序在每个副本处被执行。这个顺序总是必需的吗? .....	23
24、为了使用序列器实现全序的广播, 一种方法是先将操作转发给序列器, 然后该序列器为这个操作分配一个惟一的号码, 随后广播这个操作。请提出另外两种可选择的方法, 并比较这三种方法。 .....	23
25、一个文件被复制在 10 个服务器上。请列出表决算法允许的所有读团体和写团体。 .....	23
26、为防止复制的调用, 我们在书中讨论了基于发送方的方法。在基于接收方的方法中, 是由接收副本辨认属于同一调用的输入消息的多个拷贝。请描述基于接收方的方法如何防止复制的调用。 .....	23
27、考虑在因果一致的懒惰复制中, 一个操作被从一个写队列中删除的确切时间。 .....	23
28、在实现一个支持广播 RPC 的简单系统时, 我们假设系统有多个复制的服务器, 每个客户可以通过 RPC 与一个服务器通信。但是, 处理复制时, 客户需要向每个副本发送一个 RPC 请求。设计客户程序, 以便客户好像只向应用程序发送单一的 RPC。假设复制的目的是为了提高性能, 而那些服务器可能是易于出故障的。 .....	23
29. 在下图中加入一条与 A 并发的新消息。它既不在 A 之前, 也不在 A 之后。 .....	23
30. 假定有两台机器 A 和 B, 它们的时钟都是每毫秒滴答 1000 次。但实际上 B 每毫秒滴答 900 次。如果每分钟根据 UTC 时间校正一次时钟, 那么 A 和 B 之间的最大时钟偏差是多少? .....	24
31. 对所介绍的 bully 选举算法, 进行优化。 .....	24
32. 假定 A 和 B 是相互独立的两个临界区, 进程 0 要进入 A, 进程 1 要进入 B, R-A 分布式互斥算法会导致死锁吗? 说明理由。 .....	24

33. 假设有标识符分别为 1-8 的 8 个进程, 其中进程 8 崩溃, 节点 5 发起 bully 算法进行选举, 完成选举需要发送消息的个数为多少? .....	24
34. 无线网络系统的选举算法中, 完成选举的延迟为多少? .....	24
第七章 分布式一致性与复制管理 .....	24
1. 在内存一致性模型的讨论中, 经常提到软件和内存的约定。为什么需要这样的约定, 举例说明? .....	24
2. 下图为顺序一致性内存的一个例子。对 P2 做少量改动, 使它违反顺序一致性。 .....	24
3. 在电子股票市场的交易过程中, 需要提供哪种一致性, 说明理由. ....	25
4. 释放一致性的大多数实现方法是在 release 时同步共享变量, 而不是在 acquire 时同步, 但为什么还需要 acquire 操作? .....	25
5. 假设两个变量 a 和 b, 恰好位于基于分页的 DSM 系统的同一页上。然而, 它们都不是共享变量。是否会发生错误共享? .....	25
6. 在如下并行执行的进程 P1 和 P2, 列出顺序一致性所允许的 6 种语句交叉执行情况 。 ....	25
7. 浏览网页时, 为了保证读自己写一致性, 在系统中如何实现?.....	25

# 第一章 分布式系统概论

## 1、中间件在分布式系统中扮演什么角色？

答：中间件主要是为了增强分布式系统的透明性（这正是网络操作系统所缺乏的），换言之，中间件的目标是分布式系统的单系统视图。

## 2、解释（分布）透明性的含义，并且给出各种类型透明性的例子。

答：分布透明性是一种现象，即一个系统的分布情况对于用户和应用来说是隐藏的。包括：访问透明、位置透明、移植透明、重定位透明、复制透明、并发透明、故障透明和持久性透明。

## 3、在分布式系统中，为什么有时难以隐藏故障的发生以及故障恢复过程？

答：通常，要探测一个服务器是停止服务还是该服务器的反应变慢这些情况是不可能的。因此，一个系统可能在服务 响应变慢的时候报告该服务已经停止了。

## 4、为什么有时候要求最大程度地实现透明性并不好？

答：最大程度地实现透明性可能导致相当大的性能损失，从而导致用户无法接受。

## 5、什么是开放的分布式系统？开放性带来哪些好处？

答：开放的分布式系统根据明确定义的规则来提供服务。开放系统能够很容易地与其它系统协作，同时也允许应用移植到同一个系统的不同实现中。

## 6、请对可扩展系统的含义做出准确描述

答：一个系统的可扩展包含下面几个方面：组件的数量、几何尺寸、管理域的数量与尺寸，前提是这个系统可以在上面几个方面进行增加而不会导致不可接受的性能损失。

## 7、可以通过应用多种技术来取得可扩展性。请说出这些技术。

答：可扩展性可以通过分布式、复制和缓存来获得。

## 8、多处理器系统与多计算机系统有什么不同？

答：在多处理器系统中，多个 CPU 访问共享的主存储器。在多计算机系统中没有共享存储器，CPU 之间只能通过消息传递来进行通信。

## 9、在多计算机系统 256 个 CPU 组成了一个 16 X 16 的网格方阵。在最坏的情况下，消息的延迟时间有多长（以跳(hop)的形式给出，跳是结点之间的逻辑距离）？

答：假设路由是最优的，最长的路由是从网格方阵的一个角落到对角的角落。那么这个路由的长度是 30 跳。如果一行或一列中的处理器彼此相连，则路由长度为 15 跳。



### 10、现在考虑包含 256 个 CPU 的超立方体，最坏情况下消息的延迟有多长？

答：在具有 256 个 CPUs 的超立方体中，每个结点可以用一个二进制地址，范围从 00000000~1111, 1111。一个 hop 表示地址改变了一位。因此从 0000, 0000~0000, 0001 表示一个 hop，而从 0000, 0001~0000, 0011 也是表示一个 hop。所以最长的路由有 8 个 hop。

### 11、分布式操作系统与网络操作系统有什么不同？

答：分布式操作系统管理多处理器和同种类的多计算机。网络操作系统连接不同的、独立的计算机，这些计算机有自己的操作系统以使用户可以容易地使用每台计算机所提供的服务。

### 12、请解释如何使用微内核将操作系统组织成客户-服务器的方式。

答：微内核可以把客户端应用从操作系统服务中分离出来，方法是通过强迫每个请求都通过内核来传递。因此，操作系统服务能够通过用户级的服务器来实现，这些用户级的服务器是作为普通的进程来运行的。如果微内核有网络功能，这些服务器也能作为远端机器。

### 13、请解释基于分页的分布式共享存储器系统主要有哪些操作。

答：基于分页的 DSM 利用了操作系统的虚拟存储器的功能。当一个应用程序对内存寻址时，如果该内存位置没有映射到当前物理存储器的存储器空间的时候，页错误就会发生，并将控制权转交给操作系统。操作系统定位到适当的页，通过网络传递它的内容，并映射到物理存储器中。从而，应用可以继续运行。

### 14、为什么要开发分布式共享存储器系统？你认为是什么问题给这种系统的高效率实现造成了困难？

答：主要的原因是，写基于消息传递的并行和分布式程序来进行通信要比使用共享存储器难得多。不管进行什么操作，网络中的页面传输都会发生，这导致了 DSM 系统效率的下降。如果页面被不同的处理器共享，在虚拟存储系统中很容易进入和“陷入”类似的状态。最后，DSM 系统比消息传递的解决方案要慢，而且由于需要跟踪页面而变得更慢。

### 15、请解释什么是分布式共享存储器系统中的伪共享。你有没有针对这个问题的解决方案？

答：当数据属于两个不同且独立的处理器（可能在不同的机器上）时会发生伪共享，这时数据被映射到相同的逻辑页上。效果是这个页面会在两个处理器之间发生交换，从而导致不明显且不必要的依赖。解决方案是让页面更小并禁止独立的处理器共享一个页面。

### 16、由于存在错误，某个实验性的文件服务器有 3 / 4 的时间能够正常工作，而另外 1 / 4 的时间无法工作。如果要确保服务至少在 99%的时间可用，需要将该文件服务器复制多少次？

答：令  $k$  为服务器的数量，则有  $(1/4)^k < 0.01$ 。即在最坏的情况下，这时所有的服务器都已关闭，发生这种情况的概率是  $1/100$ 。因此  $k=4$ 。

### 17、什么是三层客户-服务器体系结构？

答：三层客户——服务器体系结构包括三个逻辑层，每一层在理论上来说都在一台单独的机器上实现。最高层包括了客户的用户界面，中间层包括实际的应用程序，最底层包含了被使用的数据。

### 18、纵向分布与横向分布有什么不同？

答：纵向分布指的是多台机器组成的多层架构中各个层的分布。从理论上说，每一层都在一台不同的机器上实现。横向分布则处理多台机器上的一个层的分布，例如一个数据库的分布。

### 19、考虑一个进程链，该进程链由进程 $P_1, P_2, \dots, P_n$ 构成，实现了一个多层客户-服务器体系结构。进程 $P_i$ 是进程 $P_{i+1}$ 的客户， $P_i$ 只有得到 $P_{i+1}$ 的应答之后才能向 $P_{i-1}$ 发出应答。如果考虑到进程 $P_1$ 的请求-应答性能，这种组织结构主要存在什么问题？

答：如果  $n$  很大的话性能会很差。从理论上来说，两个邻接层之间的通信应该在两台不同的机器之间发生。因此， $P_1$  和  $P_2$  之间的性能由  $n-2$  次其它层之间的请求——应答过程决定。另一个问题是如果链中的一台机器性能很差甚至临时不可达，这会立刻使最高层的性能降低。

### 20、比较分布式操作系统、多处理机操作系统、网络操作系统的特点和应用范围。

答：分布式系统是由一组具有自治功能的独立计算机构成的系统，在用户看来好像是一个计算机系统一样。物理上分布，逻辑上是一个整体。

硬件方面：每台计算机都是独立、自主的计算机

软件方面：用户感觉在独占系统

分布式系统发展的前提条件有：计算机性能价格比在迅速提高；网络技术的发展与普及：有线网络、移动计算、无处不在的计算；计算量增大；任务本身需要分布式处理例；

例子：工作流处理系统：订单自动处理、办公自动化、电网调度等系统；

银行“一卡通”系统、电信“神州行”系统、飞机订票系统等

项目	网络操作系统	分布式操作系统	多处理机分时操作系统
表现是否象一台单机	否	是	是
各台机器是否运行同样的操作系统	否	是	是
操作系统有多少拷贝	多	多	1
通信如何完成	共享文件	消息	共享存储器
是否要求一致的网络协议	是	是	否
是否具有唯一的运行队列	否	否	是
文件共享是否具有统一的语义	通常没有	有	有

## 21、比较单体型系统和微内核型系统的优缺点。

单核操作系统简单理解就是单一内核，宏内核，编写难，修改难，效率高，效率比其他系统高很多，每一个模块又都独立服务于操作，内核代码高度集成，但是所有的模块都在同一内核空间运行，一个很小的 BUG 都回导致整个系统崩溃。

微核就是把单核改了下，把系统服务的实现和系统的基本操作规则分离开来，这样不至于因为一个小 BUG 而整个系统崩溃，便于维护，但是性能上要比单核低一些。微内核优点：(1)微内核系统具有高度的模块化，对于每一个服务都有一个定义好的借口，每一个服务程序对所有客户来说都是可以访问的，且和位置无关。

(2)微内核系统具有很高的灵活性，易于实现、安装和调试新的服务程序，因为增加或者改变一个服务程序不需要像有一个单内核那样停止系统和启动一个新的内核。

## 22、解释透明性问题对系统和用户的

透明性	描述
存取透明性	隐藏了数据表示和获取资源的具体实现
位置透明性	用户不必知道资源位于何处
迁移透明性	资源可以不改名随意移动
重定位透明性	用户不必知道资源是位置是否改变
复制透明性	用户不必知道有多少拷贝存在
并发透明性	多个用户可以自动的共享资源
容错透明性	用户不必知道系统出现错误

对于分布式系统而言，透明性是指它呈现给用户或应用程序时，就好像是一个单独是计算机系统。具体说来，就是隐藏了多个计算机的处理过程，资源的物理分布。

## 23、给出不同类型透明性的例子。

访问透明:分布式系统中的多个计算机系统运行可能是不同的操作系统，这些操作系统的文件命名方式不同,命名方式的差异以及由此引发的文件操作方式的差异应该对用户和应用程序隐藏起来。

位置透明:从 <http://www.prenhal.com/index.htm> 这个 url 看不出 parentice hall 的主 web 服务器所在的位置，同时也看不出 index.html 的位置情况。

迁移透明:分布式系统中的资源移动不会影响该资源的访问方式。

重定位透明:资源可以在接受访问的同时进行重新定位，而不引起用户和应用系统的注意。移动通信用户从一个点到另一个点，可以一直使用移动设备，而无中断连接。

复制透明:对同一个资源存在多个副本这样一个事实的隐藏。所有的副本同名。并发透明:访问位于同一个共享数据库中的一批表。

故障透明:用户不会注意到某个资源无法正常工作，以及系统随后的恢复过程。持久性透明:指对转位于易失性的存储器还是在磁盘上的隐藏。许多面向对象的数据库提供直接调用存储对象的方法的功能。

## 24、举例说明最高程度的透明性对于某些应用不一定是必须的。

答:最大程度地实现透明性可能导致相当大的性能损失，从而导致用户无法接受。例如,许多 web 程序会不断尝试连接某台服务器，多次失败后才放弃。这种在用户转向另一台服务器之前竭力隐藏服务器短暂故障的企图会导致整个系统变慢。



## 25、说明实现可伸缩性的主要技术。

答:实现分布式可伸缩性,基本的三种技术为:

- 减少通信延迟,即使用异步通信方式,使得发送方发送请求后不必阻塞以等待答复,而是处理其他本地任务。
- 分布化,即将一个组件分解为几个小层。一个好的例子是 DNS 域名系统,它将域名分为三层,均衡了系统负载。
- 复制和高速缓存,它能使得资源更容易就近获取,并且它能使资源分布于整个系统,均衡了负载。

## 26、给出一个虚拟组织的实际例子。

答:属于同一虚拟组织的人或应用,具有访问提供给该组织的资源的权限。这些资源可能分布在不同的计算机、服务器、数据工具和数据库中。虚拟组织定义了谁能访问什么。资源应该保持一个账户用户和访问权限,这通常使用了标准的访问控制机制(类似 UNIX 中的 `rwX`)。

27、给出与智能家庭、电子保健、传感器网络不同的其它普适应用系统的例子

b) 另一类普适系统面向于改造和增强常见事物的能力。

文献(PUC2006,Erez Kikin-Gil)设计了以光墙为界面的系统性思维辅助学习系统。系统采用光容器来覆盖墙面,通过对光容器的操作来模拟计算机系统的数据传输和存储。高年级学生可在光墙上使用 STELLA 这种模拟计算机功能的软件,低年级的学生使用故事性的介绍来了解系统元素之间的相互关系。在设计这个光墙学习系统时,作者采用了同其它类似系统一样的基于玩的方法,让学生用户在玩的过程中了解计算机系统的原理。

## 第二章 通信

1、在许多分层协议中,每一层都有自己的报头。如果每个消息前部都只有单个报头,其中包含了所有控制信息,无疑会比使用单独的多个报头具有更高的效率。为什么不这么做?

答:协议的每一层都必须和其它层相独立。从第  $k+1$  层传送至第  $k$  层的数据同时包含了报头和数据,但是第  $k$  层协议不能对它们进行辨别。如果使用单个大的报头来包含所有信息的话将会破坏透明性,使得一个协议层的变动会影响到其它层,这显然不是我们所希望的

2、为什么传输层通信服务常常不适于构建分布式应用程序?

答:它们通常不提供分布透明性,这意味着应用程序开发人员需要注意通信的实现,从而导致解决方案的可扩展性很差。分布式应用程序,例如基于套接字构建的分布式应用程序,将很难移植或者和其它应用程序交互。

3、一种可靠的多播服务允许发送者向一组接收者可靠地传递消息。这种服务是属于中间件层还是更低层的一部分?

答:从理论上来说,一种可靠的多播服务可以很容易的成为传输层,甚至是网络层的一部分。例如,不可靠的 IP 多播服务是在网络层实现的。但是,由于这些服务目前尚无法应用,它们通常使用传输层的服务来实现,传输层的服务将它们放在中间件中。如果把可扩展性加以考虑的话,只有充分考虑应用程序的需求时可靠性才能得到保证。用更高、更特殊的网络层来实现这些服务存在一个很大的争议。

4、考虑一个带有两个整型参数的过程 `incr`。该过程将两个参数的值分别增加 1。现在假定调用它时使用的两个参数是同一个变量,比如 `incr(i, i)`。如果  $i$  的初始值是 0,在执行引用调用之后  $i$  将变为什么值?如果使用复制——还原调用呢?

答:如果执行引用调用,指向  $i$  的指针被传入 `incr`。 $i$  将会被增加两次,因此最终的结果是 2。而使用复制——还原调用时, $i$  会被两次传值,每次的初始值均为 0。两次都会增加 1,因此结果均为 1。最后都复制到  $i$ ,第二次的复制会覆盖第一次的,因此最终  $i$  的值为 1,而不是 2。

5、C 语言中有一种称为联合(union)的构造,其中的记录(在 C 语言中称作结构)的字段可以用来保存几种可能值中的一个。在运行时,没有可靠的办法来分辨其中保存的是那一个值。C 的这种特性是否与远程过程调用有某些相似之处?请说明理由。

答:如果运行时系统不能分辨一个字段的值类型,它就不能对该字段进行正确的封送处理。除非有一个标签字段用来清楚的表明一个字段的值类型,联合不能在远程过程调用中使用。这个标签字段不能被用户所控制

6、处理 RPC 系统中参数转换的一种方法是，每台机器以自己系统使用的表示方式来发送参数，由另一方在必要的情况下进行转换。可以通过首字节中的代码来表示发送消息机器所用的系统。然而，由于要在首个字中找到开头的字节这本身也是一个问题，这种方法能行得通吗？

答：首先，当一台机器发送字节 0 时，消息肯定已经送到。因此目标机器可以访问字节 0，而代码就在消息里面。这种方法不考虑字节是高位优先还是低位优先的字节。另一个方法是将代码放在第一个单词的所有字节中。因此不管检查的是哪一个字节，代码都能被找到。

7、假定客户通过异步 RPC 对服务器进行调用，随后等待服务器使用另一异步 RPC 返回结果。这种方法与客户端执行常规的 RPC 有没有什么不同？如果使用的是同步 RPC 而不是异步 RPC，情况又如何呢？

答：二者并不相同。异步 RPC 向调用者返回一个通知，这意味着客户第一次调用之后，有一个额外的消息会被发送。类似地，服务端接收到它的响应消息已经发送到客户端的通知。如果保证通信可靠的话，两次异步 RPC 调用是一样的。

8、在 DCE 中，服务器在守护程序中注册自身。如果换一种方法，也可以总是为它分配一个固定的端点，然后再指向服务器地址空间中对象的引用中就可以使用该端点。这种方法的缺陷在哪里？

答：这种方法的主要缺陷是向服务器分配对象变得很难。另外，许多端点而不止一个需要被修复。如果机器中很多都是服务器，分配固定端点不是一个好办法。

9、给出一种用来让客户端绑定到暂时远程对象的对象应用的实现示例。

答：使用 Java 实现的类如下：

```
public class Object_reference {
    InetAddress server3address; // network address of object's server
    int server3endpoint; // endpoint to which server is listening
    int object3identifier; // identifier for this object
    URL client3code; // (remote) file containing client-side stub
    byte[] init3data; // possible additional initialization data
}
```

Object\_reference 类至少需要包含对象所属的服务器的传输层地址。在具体实现中，使用了一个 URL 来标识包含了所有必需的客户端代码文件，用一个字节数组来保存进一步初始化后的代码。另一种实现可以直接保存客户端代码而不是一个 URL。这种方法将代理对象作为引用传递，Java RMI 采用了这种做法。

10、Java 和其他语言支持异常处理，当错误发生时会引起异常。如何在 RPC 和 RMI 中实现异常处理？

答：由于异常通常在服务端发生，服务器存根只能捕获这个异常并把它作为一个特殊的错误响应传送给客户端。另一方面，客户存接收这个消息并抛出这个异常以保持对服务器访问的透明性。因此，接口定义语言中也需要有对异常的描述。

11、将静态和动态 RPC 区分开来有用吗？

答：有用。和远程对象调用有用的理由一样：它会带来更多的弹性。缺点是很多的分布透明性会被牺牲掉，其程度视哪种 RPC 先被使用而定。

12、某些分布式中间件系统的实现是完全基于动态方法调用的。甚至连静态的调用也被编译成动态的调用。这种方法的优点在哪里？

答：动态调用的实现能够处理所有的调用，而静态调用仅仅是其中的一种特殊情况。优点是仅仅需要实现一种机制。一个可能的缺点是性能并不是总比使用静态调用的情况要好。

13、描述一下客户端和服务端之间使用套接字的无连接通信是如何进行的。

答：客户端和服务端都需要创建一个套接字，但是只有服务器把套接字绑定到本地的端点上。然后，服务器可以执行一个阻塞的 read 调用以等待从客户端发送的数据。类似地，在创建套接字之后，客户端仅仅执行一个阻塞调用以向服务端写数据。关闭连接是没有必要的。

14、说明 MPI 中 mp\_based 原语和 mp\_isend 原语之间的区别。

答：mpi\_based 原语使用有缓冲的通信，调用者将包含了信息的整个缓冲传送到本地的 MPI 运行时系统。当调用完

成时，这些信息要么被已发送，要么被拷贝到一个本地的缓冲区。如果使用 `mpi_isend`，调用者仅仅将指向信息的指针传送给本地的 MPI 运行时系统，然后继续往下执行。调用者需要保证在消息被拷贝或被传送之前不能覆盖它。

**15、假定只能使用暂时异步通信原语，再加上异步 receive 原语，如何实现用于暂时同步通信的原语？**

答：考虑一个同步的 `send` 原语。一个简单的实现是使用异步通信向服务器发送一个消息，然后让调用者不停地查询接收到的来自服务器的通知或响应。另一种实现方案是，如果假设本地操作系统将接收的消息保存在一个本地缓冲区中，那么阻塞调用程序直到接收到系统的消息到达信号，之后调用程序执行异步 `receive`。

**16、假定只能使用暂时同步通信原语，如何实现用于暂时异步通信的原语？**

答：异步 `send` 可以通过如下方式实现：调用者将它的消息拷贝到一个缓冲区，实际处理消息发送的进程共享该缓冲区。每当客户端将消息拷贝到缓冲区，消息发送线程被唤醒，它将该消息从缓冲区中删除并使用一个阻塞的 `send` 原语将其发送到目标机器。接收方的实现与此类似，它提供缓冲区，一个应用程序可以检查该缓冲区以确定是否有消息。

**17、通过 RPC 实现持久化异步通信有意义吗？**

答：仅仅在管理了一个队列的进程通过 RPC 将消息发送给下一个队列管理器时有意义。一个队列管理器为另一个管理器提供的服务是保存消息，调用的队列管理器可以获得远程队列的一个代理对象，并可能接收到每一个操作成功或失败的状态。使用这种方法时，队列管理器看到的仅仅是队列，而不会发生通信。

**18、在本章中我们讲过，为了自动启动一个进程以从输入队列中获取消息，常常要使用守护程序来监视输入队列。请给出一种不使用守护程序的实现方法。**

答：一个简单的实现方案是，每当接收端的进程将一条消息放进它的一个队列时，同时检查一下是否接收到了消息。

**19、IBM MQSeries 以及许多其他消息队列系统中的路由表是人工配置的。描述一种自动完成配置工作的简单方法。**

答：最简单的是现实使用一个集中的组件，该组件维护消息队列系统的拓扑结构。它使用一种已知的路由算法来计算各个队列管理器之间的最佳路由，然后为每一个队列管理器生成路由表，这些表可以由各个管理器分别下载。这种方法适合于队列管理器相对较少但是特别分散的消息队列系统。

**20、如何将持久异步通信加入到基于远程对象 RMI 的通信模型中去？**

答：RMI 应该是异步的，也就是说，在调用的时候不能即时得到调用结果。RMI 应该被保存在一个特殊的服务端，它在对象服务器中的对象启动并运行时立即将 RMI 转发到该对象。

**21、在持久通信中，接收者一般拥有自己的本地缓冲区，如果接收者不在运行状态，可以将消息放入该缓冲区中去。为了创建这种缓冲区，必须指定它们的大小。分成两方进行辩论：一方认为这种指定缓冲区大小的行为是可取的，而另一方反对这种指定大小的行为。**

答：用户指定缓冲区大小使得实现更加容易。系统创建一个固定大小的缓冲区，之后缓冲区管理会十分容易。但是，如果缓冲区已满，消息可能会丢失。另一种方案是由通信系统来管理缓冲区大小，首先先制定一个默认大小，然后视情况需要扩充或者压缩缓冲区。这种方法减少了由于缺少空间而导致丢失消息的机率，但是需要系统做更多的工作。

**22、请说明为什么暂时同步通信在可扩展性方面存在固有的问题，以及如何解决这些问题。**

答：问题在于地理上的可扩展性受到了限制。由于同步通信需要调用者被阻塞直到消息被接收，因此，如果接收方相隔较远的话调用者可能需要阻塞很长的时间才能继续运行。解决这个问题的惟一办法是设计调用程序，使得它在发生通信时可以做其它有用的工作，以有效地实现一种异步通信。

**23、给出一个将多播应用于离散数据流的例子。**

答：将一个大文件传送给许多用户就是一个例子。例如，更新 Web 服务或软件发布的镜像站点。

**24、当一组计算机组成一个逻辑上或物理上的环时，如何确保传输延迟不超过允许的最大端到端延迟时间？**

答：用一个令牌在环上循环，只有当计算机获取到该令牌的时候，它才能通过环发送数据。另外，计算机持有环的时间不能超过  $T$  秒。因此，如果假设两台邻接电脑之间的通信是有限制的，那么可以得到令牌循环的最大时间，这个时间就是发送一个包的最大端到端延迟时间。

**25、当一组计算机组成一个（逻辑上或物理上的）环时，如何确保传输延迟不小于允许的最小端到端延迟时间？**

答：从理论上来说，接收端计算机不应该在特定的时间之前接收数据。唯一的解决方案是将数据包存储在缓冲区中，缓冲可能发生在发送端、接收端或它们之间的某一个环节。缓冲数据最好的地方是在接收端，因为这里没有不可预

见的、可能延迟数据传输的障碍。接收端仅需要将数据从缓冲区中删除并将它传送至使用了一个简单的计时机制的应用程序。缺点是必须提供足够的缓冲能力。

**26、想象一下，某个令牌存储桶规范说明中的最大数据单元大小是 1000B，令牌存储桶速率是 10MB/s，令牌存储桶的大小是 1MB，而最大传输速率是 50MB/s。突发的传输可以以最大传输速率持续传输多长时间？**

答：令最大的突发时间间隔为  $\Delta t$  秒。在极端的情况下，存储桶在开始时就已经满了，还有  $10 \Delta t \text{MB}$  的数据在该间隔内被传进来。在突发的传输过程中，输出的数据为  $50 \Delta t \text{MB}$ ，等于  $(1+10 \Delta t)$ 。因此， $\Delta t$  等于 25 毫秒。

**27、在这道练习题中，要求实现一个使用 RPC 简单客户——服务器系统。服务器提供一个名为 next 的过程，该过程接受一个整型数作为输入，并且返回紧接着该整型数之后的一个数。编写一个客户使用的名为 next 的存根，它的任务是使用 UDP 将参数发送给服务器，随后等待服务器响应，如果过长时间未收到响应就认为超时了。服务器过程必须在一个公开的端口监听并接受请求，完成这些请求所要求的操作，然后返回结果。**

答：无。

**28、举出一个分布式系统实例，该系统采用以数据为中心的软件体系结构。**

组件间的通信，通过基于一个公用的存储来实现，共享的分布式文件系统，基于 web 的分布式系统，组件使用共享的基于 web 的数据服务。Hadoop 实现分布式计算 Hadoop 是一个开发和运行处理大规模数据的软件平台，是 Apache 的一个用 java 语言实现开源软件框架，实现在大量计算机组成的集群中对海量数据进行分布式计算。Hadoop 框架中最核心设计就是：MapReduce 和 HDFS。MapReduce 提供了对数据的计算，HDFS 提供了海量数据的存储。

29.0 举出一个分布式系统实例，该系统采用基于事件的软件体系结构。

发布/订阅 (publish/subscribe) 系统，发布-订阅模式其实是一种对象间一对多的依赖关系，当一个对象的状态发送改变时，所有依赖于它的对象都将得到状态改变的通知。订阅者 (Subscriber) 把自己想订阅的事件注册 (Subscribe) 到调度中心 (Event Channel)，当发布者 (Publisher) 发布该事件 (Publish Event) 到调度中心，也就是该事件触发时，由调度中心统一调度 (Fire Event) 订阅者注册到调度中心的处理代码。

**29、什么是三层客户-服务器体系结构，举出一个实际应用案例。**

适用于网络环境的事件服务是发布/订阅系统<sup>[9]</sup>。发布/订阅系统通过分布式环境将事件提供者和事件消费者连接。发布/订阅系统通常设计为覆盖在物理网络之上的“事件代理”网络。这些事件代理将事件从发布者路由至订阅者。发布/订阅系统是一种主动、实时的信息传递方式。在科学数据网络中，用户

答：三层客户——服务器体系结构包括三个逻辑层：用户接口层、处理层和数据层，每一层在理论上来说都在一台单独的机器上实现。最高层包括了客户的用户界面，中间层包括实际的应用程序，最底层包含了被使用的数据。

例如，大型 M I S（管理信息系统）中具有良好的应用前景。或者 ERP（企业资源计划）系统

**30、在结构化覆盖网络中，如果消息是根据覆盖的拓扑结构来路由的，讨论其优缺点。**

答：1、结构清晰，便于管理维护 2、灵活性高，适应各种不同需求 4、便于扩充，提高了系统的可靠性问题是我们处理的是逻辑路径。很有可能出现的情况是，在覆盖网络中的两个邻居节点 A 和 B，在物理距离上是相距很远的。因此，A 和 B 之间的逻辑最短路径实际上在底层物理网络上是很长的一条路径。

**31、在非结构化覆盖 P 网络中，比较洪泛式搜索算法与随机游走搜索算法的性能优缺点。**

洪泛式搜索算法：不能保证性能，网络带宽消耗非常大，可伸缩性差。优点是完全分布无单点故障的问题，不用维护索引结构。

随机游走搜索算法：维度大的时候概率低。优点是实现起来简单，节点之间定期交换表项，节点加入时通知任一已知节点进行视图交换，删除时可自行离开，无需通知其他节点，当其他节点发现其不在时，将它从表中删除。

32、设一个 P2P 网络中， $N$  为总的节点个数， $r$  为每个数据项被随机保存的节点的个数， $m$  为每个节点的平均邻居个数。现需搜索一个数据  $d$ ，分别采用洪泛式搜索算法和随机游走算法，在下列条件下估算出这两种算法的搜索规模  $S$  和搜索步数  $k$ ：

1) 已知： $N=100,000$ ， $r=100$ ， $m=10$ ；

2) 已知： $N=100,000$ ， $r=100$ ， $m=11$ 。

答：（1）洪泛式： $s=100000$ ， $k=10$ ；随机游走： $s=100$ ， $k=1000$

（2）洪泛式： $s=100000$ ， $k=11$ ；随机游走： $s=100$ ， $k=1100$

33、请从技术角度来解释为什么 BitTorrent 中的 tit-for-tat 策略比因特网中的文件共享好的多？

答：tit-for-tat 选择的是最适合的邻居节点的数据，采用分块选择策略（最少优先）保证了共享文件最快的分发到网络中。

大多数的 BitTorrent 客户端是被例如 ADSL 或调制解调器所提供的不对称链接所处理的。一般来说，BitTorrent 会提供给客户端比较高的进入带宽，但是并不希望客户端可以提供服务。BitTorrent 不会使用这种假设，把客户端变成协作服务器。对称连接会比 tit-for-tat 策略提供更好的匹配。

34、简述 p2p 网络的特点，并分析如果采用集中型拓扑结构构建 p2p 网络，其主要的优缺点是什么？

P2P 打破了传统的 Client/Server (C/S) 模式，在网络中的每个结点的地位都是对等的。每个结点既充当服务器，为其他结点提供服务，同时也享用其他结点提供的服务。

P2P 技术的特点体现在以下几个方面：

- **非中心化 (Decentralization)**：网络中的资源和服务分散在所有结点上，信息的传输和服务的实现都直接在结点之间进行，可以无需中间环节和服务器的介入，避免了可能的瓶颈。P2P 的非中心化基本特点，带来了其在可扩展性、健壮性等方面的优势。
- **可扩展性**：在 P2P 网络中，随着用户的加入，不仅服务的需求增加了，系统整体的资源和服务能力也在同步地扩充，始终能较容易地满足用户的需要。整个体系是全分布的，不存在瓶颈。理论上其可扩展性几乎可以认为是无限的。
- **健壮性**：P2P 架构天生具有耐攻击、高容错的优点。由于服务是分散在各个结点之间进行的，部分结点或网络遭到破坏对其它部分的影响很小。P2P 网络一般在部分结点失效时能够自动调整整体拓扑，保持其它结点的连通性。P2P 网络通常都是以自组织的方式建立起来的，并允许结点自由地加入和离开。P2P 网络还能够根据网络带宽、结点数、负载等变化不断地做自适应式的调整。
- **高性能/价格比**：性能优势是 P2P 被广泛关注的一个重要原因。随着硬件技术的发展，个人计算机的计算和存储能力以及网络带宽等性能依照摩尔定理高速增长。采用 P2P 架构可以有效地利用互联网中散布的大量普通结点，将计算任务或存储资料分布到所有结点上。利用其中闲置的计算能力或存储空间，达到高性能计算和海量存储的目的。通过利用网络中的大量空闲资源，可以用更低的成本提供更高的计算和存储能力。
- **隐私保护**：在 P2P 网络中，由于信息的传输分散在各节点之间进行而无需经过某个集中环节，用户的隐私信息被窃听和泄漏的可能性大大缩小。此外，目前解决 Internet 隐私问题主要采用中继转发的技术方法，从而将通信的参与者隐藏在众多的网络实体之中。在传统的一些匿名通信系统中，实现这一机制依赖于某些中继服务器节点。而在 P2P 中，所有参与者都可以提供中继转发的功能，因而大大提高了匿名通讯的灵活性和可靠性，能够为用户提供更好的隐私保护。
- **负载均衡**：P2P 网络环境下由于每个节点既是服务器又是客户机，减少了对传统 C/S 结构服务器计算能力、存储能力的要求，同时因为资源分布在多个节点，更好的实现了整个网络的负载均衡。

最大的优点是维护简单发现效率高，

缺点：中央服务器的瘫痪容易导致整个网络的崩溃，可靠性和安全性较低。



随着网络规模的扩大，对中央索引服务器进行维护和更新的费用将急剧增加，所需成本过高。中央服务器的存在引起共享资源在版权问题上的纠纷，并因此被攻击为非纯粹意义上的 P2P 网络模型。对小型网络而言，集中目录式模型在管理和控制方面占一定优势。但鉴于其存在的种种缺陷，该模型并不适合大型网络应用。

## 第三章 分布式进程管理

1. 比较使用单线程文件服务器读取文件和使用多线程服务器读取文件有什么不同。花费 15ms 来接收请求、调度该请求并且完成其它必须的处理工作，假定需要的数据存放在主存储器的缓存中。如果需要磁盘操作，就需要额外多花 75ms 在磁盘操作的过程中线程处于睡眠状态。如果服务器采用单线程的话，它每秒能处理多少个请求？如果采用多线程呢？

答：在单线程情况下，命中 cache 花了 15ms，未命中 cache 花了 90ms。加权平均值为  $\frac{2}{3} \times 15 + \frac{1}{3} \times 90$ 。这意味着请求花了 40ms，服务器每秒可以完成 25 次。对于多线程，所有磁盘等待都是交迭的，所以每个请求花了 15ms，服务器每秒可以处理  $66\frac{2}{3}$  个请求。

2. 对服务器进程中的线程数目进行限制有意义吗？

答：有。原因有两个：

线程需要内存来设置他们的私有堆栈。因此，线程太多可能导致消耗过多的存储器。

更严重的情况是，对于一个操作系统，独立的线程是以无序的方式在运行。在虚拟存储器系统中，构建一个相对稳定的工作环境可能比较困难，从而导致许多的页错误和过多的 I/O 操作，结果可能导致系统性能的下降。

3、在文中我们描述了一个多线程的文件服务器，说明了为什么它比单线程服务器和有限状态服务器更好。有没有这样的环境，在其中使用单线程服务器会更好？给出这种环境的例子。

答：有。如果服务器完全是 CPU 绑定的，就没有必要使用多线程了。多线程可能只是增加了不必要的复杂性。例如，某个地区拥有 1 百万人口，现要建立一个数据库来保存每个人的一些信息，如（名字，电话号码），假定每个人的信息大小为 64 字节，则数据库的总大小为 64M 字节，为了快速查找，应该把这些数据保存在服务器的存储器中。

4、将轻量级进程与单个线程静态关联起来并不好，为什么？

答：这样的关联将在很大程度上迫使只有一个内核级的线程，这就意味着多线程的性能优势将会被损失掉。

5、如果每个进程只使用单个轻量级进程也不好，为什么？

答：在这种模式下，我们只能拥有用户级的线程，这意味着任何阻塞系统调用都将阻塞整个进程。

6、描述一种使用与可运行线程数目相等的轻量级进程的方法。

答：开始只有单个轻量级进程并让它选择一个可运行的线程。当发现一个可运行的线程后，轻量级进程创建另一个轻量级进程来寻找下一个线程来执行。如果没有找到可运行的线程，轻量级进程就销毁它本身。

7、本章解释过，代理可以通过调用所有副本来支持引用透明。能否对（服务器端的）对象的副本进行调用？

答：可以。假设一个副本对象 A 调用另一个对象 B（非副本），如果 A 包含 K 个副本，一个 B 的调用将会被 A 的每个副本执行，然而，B 应该只被调用一次。这种复制调用应该采取一些特殊的措施。

8、通过生成进程来构建并发服务器与使用多线程服务器相比有优点也有缺点。给出部分优点和缺点。

答：一个重要的优点是每个进程都被保护。在超级服务器处理完全独立的服务的时候是非常有必要的。另一方面，创建进程的代价也较高，同时，如果进程需要通信，则使用线程在很多情况下将更为简单，因为它避免了用内核来实现通信。

9、粗略地设想一种多线程服务器的设计，该服务器必须使用套接字作为面对底层操作系统的传输级接口，以支持多种协议。

答：一个相对简单的设计是，使一个单一的线程 T 等待接收传输层消息 (TPDU)。如果我们假定每个 TPDU 的报头包含一个数字来识别更高层的协议，这个线程可以把它传递给处理特定协议的模块。每个这样的模块都有一个专门的线程来处理消息，这些模块把消息看作是一个输入请求。当处理完这个请求后，一个应答消息被传递给 T，这时 T 将把这个应答封装成一个传输层消息并把它传递给适当的目的地。



## 10、如何防止应用程序绕过窗口管理器破坏屏幕显示。

答：使用微内核方法，一个包含窗口管理器的窗口系统通过这个方法运行，因此所有的窗口操作都需要通过这个微内核。从效果上看，这正是第一章中所讲的把客户端——服务器端模式转换为单计算机的本质。

## 11、解释对象适配器的概念。

答：对象适配器是一种普通的程序，它可以接收输入的调用请求并把它们传递给对象的服务器存根。对象适配器主要负责实现调用策略，这些策略确定是否、怎样以及多少多线程被用来调用一个对象。

## 12、举出几个用于支持持久性对象的对象适配器设计方面的问题。

答：最重要的问题可能是产生一个对象引用，这个引用可以被当前服务器和适配器独立地使用。这样的引用应该能被传递给一个新的服务器并可能为相应的对象指定特定的活动策略。其它的问题包括，持久性对象什么时候、怎样被写到磁盘上，以及对象保存在主存储器中的内容和保存到磁盘上的内容有什么区别。

## 13、改变对象适配器示例中的 thread\_per\_object 过程，使用单个线程来处理受该适配器控制的所有对象。

答：代码基本保持不变，除了我们需要创建一个 single 线程，该线程用修改过的版本 thread\_per\_object 来运行。这个线程作为 adapter\_thread 来被引用。当多路输出选择器调用这个适配器的时候，它向 adapter\_thread 的缓冲区中存放一条消息，然后再读取它。适配器线程调用适当的 stub 并处理响应消息。

```
#include <header.h>
#include <thread.h>
#define MAX3OBJECTS 100
#define NULL 0
#define ANY -1
METHOD3CALL invoke[MAX3OBJECTS]; /* array of pointers to stubs */
THREAD *root; /* demultiplexer thread */
THREAD *adapter3thread /* thread that runs single3thread*/
void single3thread(long object3id) {
    message *req, *res; /* request/response message*/
    unsigned size; /* size of messages */
    char **results; /* array with all results*/
    while(TRUE) {
        get3msg(&size, (char*) &req); /* block for invocation request */
        /* Pass request to the appropriate stub. The stub is assumed to */
        /* allocate memory for storing the results. */
        (invoke[req->object 3id])(req->size, req->data, &size, results);
        res = malloc(sizeof(message)+size); /* create response message */
        res->object3id = object3id; /* identify object */
        res->method3id = req.method3id; /* identify method */
        res->size = size; /* set size of invocation results */
        memcpy(res->data, results, size); /* copy results into response */
        put3msg(root, sizeof(res), res); /* append response to buffer */
        free(req); /* free memory of request */
        free(*results); /* free memory of results*/
    }
}

void invoke3adapter(long oid, message *request) {
    put3msg(adapter3thread, sizeof(request), request);
}
```

## 14、维护到客户的 TCP/IP 连接的服务器是状态相关的还是状态无关的。

答：假设服务器没有在客户端上保存其它消息，就可能认为服务器是状态无关的。问题在于不是服务器，而是服务器的传输层在客户端上保存了状态。本地操作系统所跟踪的与服务器无关。

15、想象一下，某个 web 服务器维护一个列表，该列表中的内容是 IP 地址与该地址最近访问过的 web 页的映射关系。当一个客户连接到该服务器的时候，该服务器在列表中查找该客户，如果找到的话就返回注册过的页面。这个服务器是状态相关还是状态无关的？

答：是状态无关的。状态无关设计的主要问题不是服务器是否保存了客户端的任何信息，而是所保存信息的准确度。在本例中，如果表丢失了，客户端与服务器还能正常的进行交互，就象什么事都没有发生一样。在一个状态相关设计中，这种交互只有在服务器从错误中恢复过来以后才能进行。

16、Java RMI 对代码迁移依赖到何种程度？

答：把对象引用看作是可移动的代理，对象引用每次被传递的时候，我们实际上通过网络来迁移代码。幸运的是，代理没有执行状态，因此支持简单的移动性。

17、在 UNIX 系统中，可以通过让进程在远程机器上派生出一个子进程来支持强可移植性。请说明这种机制的工作机理。

答：UNIX 中派生的意思是把父进程的完整镜像拷贝给子进程，这意味着在调用完 fork 后，子进程继续运行。一个相似的方法可以用来做远端克隆，这里假定了目标平台与父进程的平台是一样的。第一步是让目标操作系统为新建的子进程保存资源并创建相应的进程和存储器映射。这一步完成之后，父进程的镜像可以被拷贝了，子进程也可以被激活了。

18、图 3.13 指出，强可移植性不能与在目标进程中已迁移代码的执行结合在一起。请举出一个反例。

答：如果强可移植通过线程迁移发生的话，它就应该能使一个迁移的线程在目标进程的环境中执行。

19、考虑某个进程 P，它请求访问与自己位于同一台机器上的本地文件 F。在 P 迁移到另一台机器上以后，它还需要访问 F。如果文件对机器的绑定是紧固的，如何实现对 F 的系统级引用？

答：一个简单的解决方案是，创建一个单独的进程 Q，用 Q 来处理对 F 的远端请求。提供给进程 P 与以前一样的接口。例如以代理的形式。从效果上讲，进程 Q 是作为文件服务器来工作的。

20、D' Agents 系统中的每个代理都由单独的进程实现。代理通信的主要方式是共享文件以及消息传递，而文件无法在机器之间传递。如果按照 3.4 节中给出的移动框架结构，代理的那一部分状态包含了资源段？

答：资源段包含了所有本地以及全局的资源引用。同样地，它由那些指向其它代理、本地文件等等的变量组成。在 D' Agents 系统中，这些变量主要包含在全局程序变量部分中。事实上，在 D' Agents 系统中，所有的资源都是不可转移的，这就使得事情变得简单了。只有代理能在机器之间传递。因为代理已经被全局引用命名，即：一个 (address, local-id) 对，在提供迁移的情况下转换引用到资源，在 D' Agents 中就相对比较简单了。

21、将 D' Agents 的体系结构与 FIPA 模型所实现的代理平台的体系结构进行比较。

答：它们之间的主要区别是，D' Agents 没有单独的目录服务，事实上，它只提供了低级的命名服务，通过这个服务，代理能被全局引用。FIPA 体系结构中的管理组件对应 D' Agents 中的服务器，然而 ACC 是由通信层实现的。与 D' Agents 相比，FIPA 没有提供代理的体系结构的更进一步的信息。

22、代理通信语言 (ACL) 在哪些方面与 OSI 模型相符？

答：这些语言是应用层的一部分。

23、在系统顶层实现代理通信语言以进行电子邮件处理时，这中代理通信语言又有哪些方面与 OSI 模型相符？这种方案有什么优点？

答：它应该还是应用层的一部分。在 e-mail 的顶层实现 ACL 的一个重要原因是简单。一个完整，全世界范围内的通信基础设施对于处理代理之间的异步消息传递来说是可使用的。实际上，它与第二章讨论的消息队列系统很相近。

24、为什么 ACL 消息中通常必须指定实体 (ontology)？

答：在这个特定的环境下，实体最好被解释为对 ACL 消息中包含的实际数据的标准解释的一个引用。通常，消息传递系统中的数据总是被假定能被消息的发送者和接收者正确地解释。代理经常被认为是高度独立的。因此，不能总是假定接收方能正确地解释接收到的数据。当然，对实体域的解释形成一些约定是必须的。

25、有一个文件服务器，当所需的数据放在缓存中时，接受、分派或处理一个请求需要花费 15ms；但在 1/3 的情形下，还要做磁盘操作，需要 75ms，线程在这时处于睡眠状态。如果使用单线程，每秒能处理多少请求？而如果使用多线程呢？

答：在单线程情况下，命中 cache 花了 15ms，未命中 cache 花了 90ms。加权平均值为  $23 \times 15 + 31 \times 90$ ，这意味着请求花了 40ms，服务器每秒可以完成 25 次。对于多线程，所有磁盘等待都是交迭的，所以每个请求花了 15ms，服务器每秒可以处理 66 个请求。/

**26、假设完全在用户空间实现线程，运行系统每隔 1 秒得到 1 个时钟中断。如果发生时钟中断时，某个线程正在运行系统中执行。这时会出现什么问题？给出一种解决方案。**

答：A：runtime 系统可以正好在这一时刻阻塞或者解除阻塞某个线程，并且忙于处理调度队列。此时并不适合于时钟中断处理程序开始检查该队列是否应该进行线程切换，因为它们可能处于不一致的状态。

解决方法可以是：当进入 runtime 系统后，设置一个标志。时钟处理程序将看到该标志，并且设置其自己的标志，然后返回。当 runtime 系统完成时，它将检测时钟标志，看是否有时钟中断发生，并且现在运行时钟处理程序。

**27、设计一种简单的线程实现方式，其中的轻量进程的个数与可运行的线程一样多。**

答：开始只有单个轻量级进程并让它选择一个可运行的线程。当发现一个可运行的线程后，轻量级进程创建另一个轻量级进程来寻找下一个线程来执行。如果没有找到可运行的线程，轻量级进程就销毁它本身。

**28 设计一种支持多通信协议的多线程服务器。它与底层操作系统的传输层接口采用 socket。**

**29 设计支持持久对象的对象适配器时，需考虑哪些问题？**

答：最重要的问题可能是产生一个对象引用，这个引用可以被当前服务器和适配器独立地使用。这样的引用应该能被传递给一个新的服务器并可能为相应的对象指定特定的活动策略。其它的问题包括，持久性对象什么时候、怎样被写到磁盘上，以及对象保存在主存储器中的内容和保存到磁盘上的内容有什么区别。

**30 修改示例对象适配器中的 thread\_per\_object 过程，使得该适配器控制的所有对象都用一个单一的线程处理。**

答：代码基本保持不变，除了我们需要创建一个 single 线程，该线程用修改过的版本 thread\_per\_object 来运行。这个线程作为 adapter\_thread 来被引用。当多路输出选择器调用这个适配器的时候，它向 adapter\_thread 的缓冲区中存放一条消息，然后再读取它。适配器线程调用适当的 stub 并处理响应消息。

**31 在一个 Web 服务器中，维护有一个映射表，描述客户的 IP 地址与最近访问的网页的映射关系。当一个客户连接该服务器时，服务器在表中查找该客户，如果找到，则返回登记的网页。这种服务器，是有状态的，还是无状态的？**

答：是状态无关的。状态无关设计的主要问题不是服务器是否保存了客户端的任何信息，而是所保存信息的准确度。在本例中，如果表丢失了，客户端与服务器还能正常的进行交互，就象什么事都没有发生一样。在一个状态相关设计中，这种交互只有在服务器从错误中恢复过来以后才能进行。

## 第四章 分布式通信管理

### 1. 为什么传输层通信对建立分布式应用来讲不合适？

答：它们通常不提供分布透明性，这意味着应用程序开发人员需要注意通信的实现，从而导致解决方案的可扩展性很差。分布式应用程序，例如基于套接字构建的分布式应用程序，将很难移植或者和其它应用程序交互。

### 2. 编写并运行一个实现文件操作的 RPC 程序。

**3. 过程 incr 有两个整形参数。该过程将每个参数分别加 1。假设在调用 incr 时使用同一个变量，如 incr(i, i)。若 i 的初值为 0，如果按 call-by-reference 方式调用，i 值为多少？**

答：如果执行引用调用，指向 i 的指针被传入 incr。i 将会被增加两次，因此最终的结果是 2。而使用复制——还原调用时，i 会被两次传值，每次的初始值均为 0。两次都会增加 1，因此结果均为 1。最后都复制到 i，第二次的复制会覆盖第一次的，因此最终 i 的值为 1，而不是 2。

**4. 假设完成一个空 RPC 的时间是 1.0ms，此外，传输每 1K 数据需要 1.5ms。如果用 1 个 32K 的 RPC，从一个文件服务器中读取 32K 数据需要多少时间？如果用 32 个 1K 的 RPC，又需要多少时间？**

答：(1) 由题知——我们记一次读取 32k 数据的 RPC 所需的时间记为 T，则  $T = 1.0 + 32 * 1.5 = 49\text{ms}$ 。故一次读取 32k 数据的 RPC 所需的时间是 49ms。

(2) 32 次读取 1k 数据所需的时间记为 T1 则  $T1 = 32 * (1.0 + 1.5) = 80\text{ms}$ ，故 32 次读取 k 所需的时间为 80ms。

5. 假定对服务器做异步式 RPC 调用，并等待直到服务器用另一个异步式 RPC 返回。这种方法与客户执行一个正常的 RPC 有何区别？

答：二者并不相同。异步 RPC 向调用者返回一个通知，这意味着客户第一次调用之后，有一个额外的消息会被发送。类似地，服务端接收到它的响应消息已经发送到客户端的通知。如果保证通信可靠的话，两次异步 RPC 调用是一样的。

6. 解释为什么瞬时的同步通信天生具有可伸缩性问题，以及如何解决。

7. 在持久性通信中，接受程序需要有保存消息的局部缓冲区。在创建该缓冲区时，需要设定它的大小。给出这样做的理由。在给出不该这样做的理由。

答：用户指定缓冲区大小使得实现更加容易。系统创建一个固定大小的缓冲区，之后缓冲区管理会十分容易。但是，如果缓冲区已满，消息可能会丢失。另一种方案是由通信系统来管理缓冲区大小，首先先制定一个默认大小，然后视情况需要扩充或者压缩缓冲区。这种方法减少了由于缺少空间而导致丢失消息的机率，但是需要系统做更多的工作。

8. 在 IBM MQSeries 的路由表中是手工配置的，设计一种自动配置方法。

## 第五章 同步

1、说出至少三种可在 WWV 广播时间和在分布式系统中处理机设置内部时钟之间引入的延迟源。

答：信号在大气中的传播延迟，当机器在协调 WWV 接受者与以太网时的碰撞延迟，包在局域网上的传播延迟，各个处理器由于中断处理延迟和内部队列延迟所产生的延迟。

2、考虑分布式系统中的两台机器的行为。这两台机器的时钟假设每毫秒滴答 1000 次。但实际上只有一台是这样的，另一台每毫秒滴答 990 次。如果 UTC 每分钟更新一次，那么时钟的最大偏移量是多大？

答：第二台时钟每秒滴答 990,000 次，每秒的偏移量为 10ms，所以一分钟偏移量为 600ms。

3、向图 5.7 中加入一个与消息 A 并发的新消息，即它既不发生在 A 的前面也不发生在 A 的后面。

答：在 0 时刻从进程 2 发送消息，并且在 8 时刻到达进程 1，或者是在 0 时刻从进程 1 发送消息并在 10 时刻到达进程 2。这两种方法都可以满足要求。

4、要使用 Lamport 时间戳实现全序多播，是不是每个消息都必须要被严格地确认？

答：不是。只要消息的时间戳比所接收的消息的时间戳大，就有足够的时间多播其他类型的消息。将消息 m 传送到应用程序的条件是：以一个很大的时间戳从其它所有进程那里接受到其他消息。这保证了进行中的消息都使用更小的时间戳。

5、考虑一个消息只是以它们被发送的顺序被传递的通信层。给出一个不需要这种限制的例子。

答：想象传输一个大的图像，图像被分成连续的块。每块以它在原始图像中的位置标识，可能还有高度和宽度。如果是那样的话，FIFO 的顺序就不是必要的，因为接受者可以很容易地把到来的块拼接到正确的位置。

6、假设两个进程同时检测到协调者崩溃了，并且它们都使用欺负算法主持一个选举。这时将发生什么。

答：每个更高位置的处理器都会收到两个 ELECTION 消息，但是会忽略掉第二个。选举会象往常一样进行下去。

7、在图 5.12 中，我们有两个 ELECTION 消息同时在循环，当对它们都没有什么不利影响时，杀掉一个将更好。设计这样一个算法，该算法在不影响基本的选举操作算法的情况下完成此项任务。

答：当一个进程接受到 ELECTION 消息，它会检查消息是谁开始发送的，如果是它自己开始的（例如它的位置在列表的首位），它会把消息变成协调者消息并在正文里描述它。如果消息不是它开始的，它会加入其进程号并沿着环向前发送。然而，如果它更早地发送了自己的选举消息并发现了竞争者，它会将创始人的进程号与自己的相比较。如果其他进程拥有较小的号码，它会丢弃那个消息而不是传递消息。如果竞争者更大，那么消息将以平常的方式发送。这样，如果多个选举消息被开始发送，那么入口最大的消息将会幸存，其他消息将沿着路由被忽略掉。

8、许多分布式算法需要使用协调进程。讨论一下，这样的算法实际上可以在什么程度上被看作是分布式的？

答：在集中式的算法中，常常是固定的进程充当协调者。分布来源于其他进程在不同的机器上运行的事实。在分布式算法中，没有固定的协调者，协调者从组成部分算法的进程中选出。事实是协调者能使算法更具分布性。

9、在集中式互斥方法中（图 5.13），协调者在收到一个进程释放它独占访问的临界区的消息后，通常给等待队列中的第一个进程授权，允许其访问临界区。请给出协调者另一个可能的算法。

答：请求会和优先权联系在一起，优先权取决于他们的重要性。协调者应该保证最高优先级的最先请求。

10、请再考虑图 5.13，假设协调者崩溃了，这会必然导致整个系统瘫痪吗？如果不会，那么在什么情况下会呢？有什么方法避免了这个问题的发生，使系统能够忍受协调者的崩溃？

答：在允许和拒绝的条件下，假定算法是只要有请求，就立即回应。如果没有进程在临界区并且没有进行在排队，那么这个崩溃不是毁灭性的。下一个请求准许的进程将不会获得任何回应并且初始一个协调者的选举。在发送回应之前，使用协调者存储每个到来的请求将会使系统更加健壮。这样的话，在崩溃事件中，新的协调者将会重建一个活动临界区域列表并将从磁盘读文件的行为排队。

11、Ricart 和 Agrawala 算法会有这样的问题：如果一个进程崩溃，并且没有对另一个要求进入临界区的进程的请求作应答，没有应答意味着拒绝请求，我们建议所有的请求应用立即被应答，以便很容易地检测到崩溃的进程，是否存在一些情况，即使使用这种方法也还不足够？请讨论。

答：假定一个进程在拒绝许可并在那时崩溃。请求进程会认为它是活动的，但许可永远不会到来。一种方法为是请求者并不真正阻塞，而是休眠一定的时间。在休眠后，请求者将会测试所有拒绝许可的进程是否还是活动的。

12、如果我们假设该算法可以在支持硬件广播的 LAN 中实现，那么图 5.16 中的实何时将如何改变？

答：仅仅只有分布式案例的实体改变。因为发送一个点对点的消息的花费与广播消息一样多，所以我们只需要发送一个广播消息到所有的进程去请求实体到临界区。同样地，仅仅一个出口的广播消息是需要的。延迟变成  $1+(n-1)$ ：一个延迟来源于广播请求，再加上  $n-1$ ，因为在被允许进入临界区之前，我们仍然需要从其他每个进程接受到消息。

13、分布式系统可能有多个互相独立的临界区，假设进程 0 想进入临界区 A 而进程 1 想进入临界区 B，Ricart 和 Agrawala 的算法会导致死锁吗？请解释原因。

答：这取决于程序的基本规则。如果进程严格地依次进入临界区域，也就是说一个临界区内的进程不会试图进入其他进程的临界区，那么当它拥有其他进程想要获取的资源（比如临界区）的时候，它是不会阻塞的。系统于是死锁释放。另一方面，如果进程 0 进入临界区 A 并且试图进入临界区 B，如果另外的进程以相反的顺序获取它们时，死锁就会发生。因为每个临界区被所有的进程独立处理，所以 Ricart 和 Agrawala 算法不会造成死锁。

14、在图 5.17 中，我们看到一种用磁带实现对存货表自动更新的方法。因为磁带可以很容易地用磁盘上的一个文件来模拟，你认为为什么不使用这种方法（用磁盘文件代替）呢？

答：第一个原因大概是人类变得更加贪婪和想要做更多他们过去习惯的东西。如果用户满足于每天运行一次简单的保存详细目录，可以用磁盘。问题在于每个人都想即时访问数据库，这会使得将详细目录保存在磁带是不可能的。

15、在图 5.25 (D) 中有三个调度策略，两个合法，一个非法的。对于同一个事务处理，给出临界 x 所有的可能值和合法的与非法的状态的一个完整表。

答：合法的值为 1, 2, 3。非法的值为 4, 5, 6。4 可以通过中间的事务处理来实现，其他两个的交错是不正确的。5 在策略 3 中实现。6 可以通过首先三次设置 x 为 0，然后三次增加它实现。

16、当一个私有工作空间用于实现事务处理时，可能需要将大量的文件索引复制到父辈工作区，怎样在不引入竞争条件下实现这种操作？

答：一种办法为在系统顶层设置锁以阻止所有的活动直到所有的索引都被重写。在预防崩溃以前，创建一个目的表也许是明智的。

17、给出一个完整的算法判定一个试图锁住一个文件的操作是否成功。要同时考虑到读锁和写锁，以及文件被解锁，读锁写或写锁写的可能性。

答：算法如下：如果文件没被锁，那么操作通常是成功的。如果文件被读锁定并且操作是另一个读锁定，那么操作也是成功的。在其他情况，操作是失败的（例如，写锁定一个被锁文件会失败，当一个文件被写锁定时，读锁定会失败）

18、用锁定的方法实现并发控制的系统通常区分读锁和写锁，如果一个进程已经得到了一个读锁但现在又想将它换成写锁，将会怎么样呢？

答：一个读锁只有在没有其它读进程的情况下才可能被转换成一个写锁。一种有效的方法就是先释放这个读锁，然后马上再获得一个写锁。如果有其它的读进程的话这个操作将会失败。如果有其它进程正在等待获得一个写锁时，

这个操作是否成功取决于具体的策略，而不是技术上的问题。将一个写锁降低优先级为一个读锁问题可行的，并且不会出错。

**19、在分布式事务中使用时间戳排序，假设写操作 `write(T1, x)` 被传给数据管理器，因为仅有的可能发生冲突的操作 `write(T2, x)` 的时间戳较早，为什么让调度管理器推迟传递 `write(T1, x)` 直到 T2 完成呢？**

答：通过推迟到 T2 完成，调度器可以避免由于 T2 运行的失败而造成的其它一系列的运行失败。

**20、乐观并发控制与使用时间戳相比是更严格还是更不严格呢？为什么？**

答：它更加严格。因为通过使用乐观并发控制的机制，如果一个事务正在交付并且发现另外一个事务也在更改当前事务使用的文件，它问题会退出。如果使用时间戳的话，如果另一个事务的时间戳更低的话，这个事务有时候可能会成功。

**21、使用时间戳来进行并发控制能保证串行性吗？请讨论。**

答：可以的。它要么按照事务的时间戳全部串行地完成，要么只完成一部分。

**22、我们常说事务被中止时，一切都恢复到它以前的状态，就好像事务从来没有发生一样，实现上我们在说谎。请给出一个重新恢复一切是不可能的例子。**

答：物理 I/O 发生的事务不能够被重置。比如说，一个进程打印出了一些数据，这些墨不能够从纸上移掉。同样，一个控制所有进程的系统也不大可能会恢复一切事务到它们以前的状态。

**23 举例说明，需要进一步解析实体 E 的地址，才能真正访问 E。**

答：IP 地址被用来定位主机。不过访问一个主机，它的 IP 地址需要被解析成 MAC 地址，比如说，一个以太网地址。

24 在基于 DHT 的系统中，使用递归查询方式解析一个键值时，有何缺点？

答：问题在于，如果没有结果返回，那么请求客户端将无法发现问题出在哪里。在解析键值的过程中，可能发生消息丢失或者节点失败的情况。因此动态查询有时更受欢迎，这样客户端将知道查询的哪部分出了问题，并可能求助另一个可用的节点。

25 在举例的 chord 系统中，假设加入了节点 7，为它建立指示表，并修改其他相关的指示表。

答：节点 7: **【9, 9, 11, 18, 28】**。比如，2 号入口是这样计算的： $\text{succ}(7+21)=\text{succ}(9)=9$ 。更多的指状表需要进行改变，节点 4: **【7, 7, 9, 14, 28】**，节点 21: **【28, 28, 28, 1, 7】**，节点 1: **【4, 4, 7, 9, 18】**。

26 考虑 DNS。设一个子域（subdomain）实现在与当前域不同的区域（zone）中。为了引用该子域中的节点 N，需要为该子域定义一个名字服务器。是否总是必须为该服务器的地址设置一个资源记录，或者只需提供它的域名就足够了。

答：当代表名称服务器的节点 NS 和节点 N 不在一个子域中，仅提供域名就可以了。在这种情况下，可以通过分离的 DNS 查询找到名称。当 NS 和 N 在同一个子域中时，需要从名称服务器获取它的地址。

27 在大多数 Unix 系统中，如何查找安装点（mounting point）？

答：使用挂载表，它包含了一个指向挂载点的入口。这意味着当查找一个挂载点的时候，我们需要遍历挂载表来找到与挂载点相匹配的入口。

（存储节点标识符的目录结点称为挂结点，外部名称空间中的目录结点称为挂载点，通常挂载点是名称空间的根，在名称解析的过程中，挂载点被查询，并通过访问目录表来完成解析。）

28 举例说明，URL 的闭合机制是如何工作的。

答：假设已知一个过程是处理 URL 的，那么它会首先从 URL 中提取协议标识符，比如 `ftp:`。它随后在表中查找这个字符串，进而找到在本地执行 FTP 协议的接口。此闭合机制的另一工作就是从 URL 中提取主机名，比如 `www.cs.vu.nl`，并且将主机名传给本地的 DNS 名字服务器。了解从何处以及如何与 DNS 服务器通信也是此闭合机制的重要部分。这通常硬编码进 URL 名字解析器，解析器执行整个解析过程。最后，URL 的最后一部分，指向着查找的文件，被传给标识符主机。后者通过自己的本地闭合机制来开始文件名的名字解析。

29 解释 Unix 中硬链接和软链接的区别。

答：硬链接是目录文件的一个条目，指向了同样的文件描述器（共享 inode），并成为另一个条目（可能是在不同的目录上）。符号链接（软链接）是包含了另一个文件名的文件（不共享 inode）。而通过软链接，则可以链接到不同的



磁盘分区甚至不同的机器。硬链接:新建的文件是已经存在的文件的一个别名,当原文件删除时,新建的文件仍然可以使用(因为即使原文件删除了,只要硬链接个数不为0,文件在物理位置就不会被删掉,这样通过硬链接仍然能够访问到物理号 inode 对应的内容,除非硬链接数为0了,物理位置上的内容才被删除)。软链接:新建的文件以“路径”的形式来表示另一个文件,和 Windows 的快捷方式十分相似,新建的软链接可以指向不存在的文件。主要区别如下:

1. 硬链接和原来的文件没有什么区别,而且共享一个 inode 号(文件在文件系统上的唯一标识);而软链接不共享 inode,也可以说是个特殊的 inode,所以和原来的 inode 有区别。
2. 若原文件删除了,则该软连接则不可以访问,而硬连接则是可以的。
3. 由于软链接的特性,导致其可以跨越磁盘分区,但硬链接不具备这个特性。

30 解释 DNS 如何用于实现基于 home 的方法,以定位移动主机。

答:移动主机的 DNS 名字将被用作标识符(虽然不是很好)。每当解析名称的时候,它将返回主机的当前 IP 地址。这意味着负责提供 IP 地址的 DNS 服务器将作为主机的名称服务器。每当主机移动,它就与主服务器通信,将自己的当前地址发送出去。注意需要设置一个避免抓取地址的机制,换句话说,其他的名称服务器不应该获取查找到的地址。

31. 在一个深度为 k 的分层定位服务中,当一个移动实体改变位置后,最多需要更新多少个定位纪录。

答:改变位置其实就是插入操作以及删除操作的结合。一次插入操作最坏情况下需要改变 k+1 条位置记录。同样的,一次删除操作也需要改变 K+1 条记录,并且根记录被两种操作共享。这导致 2K+1 条记录的改变。(因为叶节点中存储的是当前页节点在页域中的位置,所以深度为 k 时,需要执行 K+1 次更新)。

32 说明为什么权重引用计数比简单计数更有效(假设通信是可靠的)

答:初始化引用权重为一个很大的值(例如:  $2^{16}$ )。从这个引用创建新引用的时候将这个权重分为两份,一半留给原引用,一半分给新引用。这个地方并不需要更新对象的引用计数。

如果当要分裂的引用权重为 1 时,重新向对象申请比较大的计数权重。

当引用销毁的时候引用计数减去对应的权重,当引用计数为 0 则说明对象已经没有引用了,可以安全的释放内存。

## 第六章 一致性和复制

**1、对共享的 JAVA 对象的访问可以通过将其方法声明为同步的而被串行化.当这种对象被复制时,这种方法足以保证访问的串行化吗?**

答:不能.问题是对每一个复制对象的访问是串行化的.但是对不同的复制对象可以在同一时间进行不同的操作,使得复制的实例变量不一致.

**2、对于第一章所讨论的监视器,如果允许在一个复制的监视器中阻塞多个线程,那么给条件变量发信号时,需要保证什么?**

答:在每一组复制线程中,相同的线程会被唤醒,所以所有的复制对象都发生了完全相同的流控制.实际上,这意味着运行时系统应该完全控制着每一组复制对象的线程的运行过程.

**3、请用自己的语言解释实际考虑一致性模型的主要原因.**

答:弱一致性模型的出现是为了满足性能复制的需要.但是,只有当我们能够阻止全局同步,复制才会有效.而要达到这样的目的,就得放松一致性的限制。

**4、请解释 DNS 如何进行复制,以及它实际运行很好的原因。**

答:基本思想是域名服务器预先查询结果。结果可以存在高速缓存中很长一段时间,因为 DNS 认定主机名——IP 地址的映射不会常常改变。

**5、讨论一致性模型时,我们经常提及软件和数据存储间的合约。为什么需要这一合约?**

答:如果程序期望连续的一致性的数据存储且不能缺少这些任何数据,存储必须保证连续一致。但是,为了提高性能,一些系统提供弱一点的模式。事实上软件能遵守模式强加的规则。总的说来,这意味着遵守规则的程序可以感知到什么看起来像是连续一致的数据存储。

6、线性化假设存在一个全局时钟。但是，我们已经在严格一致性中指出，这种假设对于大多数分布式系统都是不现实的。线性化可以应用于物理分布的数据存储吗？

答：可以。线性化假设放松了同步时钟，就是说，它假设几个事件发生在同一时间段。这些事件需要根据时间分类。

7、如果多处理器使用单一总路线，那么可以实现严格一致的存储器吗？

答：可以。总线连续请求因此它们以绝对的时间顺序出现在内存中。

8、对于图 6.5 (B)，为什么  $W1(x) \ a \ R2(x) \ NIL \ R3(x) \ a$  是非法的？

答：它违背了数据一致性。

9、在图 6.7 中，000000 是仅满足 FIFO 一致的分布式共享存储器的合法输出吗？请解释你的答案。

答：是的。假设 (a) 第一个运行。它打印 00。现在 (B) 运行。如果存入 (A) 的还没有到达，它也打印 00。现在 (C) 运行。如果前面的两个在存储的都没到达，它也打印 00。

10、在图 6.8 中，001110 是顺序一致的存储器的合法输出吗？请解释你的答案。

答：是的。如果程序以 (a) (c) (b) 的顺序运行，就能得到这个结果。

11、我们在 6.2.2 节的最后讨论了一个形式模型，该模型规定顺序一致的数据存储上的每个操作集合可以用一个字符串 H 表示，从这个字符串 H 可以衍生出所有的单独进程的顺序序列。对于图 6.9 中的进程 P1 和 P2，给出 H 所有可能的值。忽略进程 P3 和 P4，也不包括它们在 H 的操作。

答：很明显 P2 在被写之前不能读 1，因此 H 的所有值都将在读 1 前写入 1。程序的顺序是有关的，因此写 2 一定是在写 1 之后。唯一可能出现的结果是：

$H = W(x) \ a \ R(x) \ a \ W(x) \ b \ W(x) \ c$

$H = W(x) \ a \ R(x) \ a \ W(x) \ c \ W(x) \ b$

12、在图 6.13 中，一个顺序一致的存储器允许 6 种可能的语句交叉。请列出这 6 种所有可能。

答：(1)  $a=1; \text{ if } (b==0); \ b=1; \text{ if } (a==0);$

(2)  $a=1; \ b=1; \text{ if } (a==0); \text{ if } (b==0);$

(3)  $a=1; \ b=1; \text{ if } (b==0); \text{ if } (a==0);$

(4)  $b=1; \text{ if } (a==0); \ a=1; \text{ if } (b==0)$

(5)  $b=1; \ a=1; \text{ if } (b==0); \text{ if } (a==0);$

(6)  $b=1; \ a=1; \text{ if } (a==0); \text{ if } (b==0);$

13、通常认为弱一致性模型给程序开发人员强加了负担。这一命题在哪个程度上确实是正确的？

答：这要看情况。许多程序员习惯于使用同步机制例如锁或通信的方式来保护他们的共享数据。主要意图是他们需要一种能提供只进行读或写操作的方式而不是两者同时进行。但是，程序员非常期望对同步变量的操作一致连续。

14、在分布式共享存储系统中的释放一致性的大多数实现中，共享变量是在执行释放操作时被同步的，而不是在执行获取操作时被同步的。那么，为什么还需要获取操作呢？

答：获取操作需要用于延迟进程当它试图访问共享变量时而另一个进程正在对共享变量进行此操作。

15、ORCA 提供顺序一致性还是入口一致性？请解释你的答案。

答：形式上，ORCA 只提供入口一致性，因对相同对象的并行操作是完全序列化的。但是，当使用完全顺序广播用为操作排序的实施方法时，所有操作顺序都是全局安排的，并取决于它们操作的对象。在这种情况下，它提供顺序一致性。

16、使用序列器并为维持主动复制的一致性而进行的全序的广播是否违背系统设计中的端到端理论？

答：是的。端到端理论表明了问题应该在它们发生的相同层次上去解决。既然这样，我们以全序广播的方式处理问题以达到主动复制的一致性。在基本协议里，一致性是通过趋向原始作为第一性来达到的。而如果使用一个序列器，我们可以做到相同的事但在更低的一抽象层次。在这种情况下，可能使用原始协议更好因为它的升级是通过广播发送操作。

17、您会选择哪种类型的一致性来实现电子股票市场？请解释您的答案。

答：因果一致性应该能满足需要。问题是对于股价的变化，反应应该是连续的。独立股票的变化应该能以不同的顺序显示出来。

18、如果要将一个移动用户的个人邮箱作为广域范围的分布式数据库的一部分来实现，哪种类型的以客户为中心的一致性最为合适？

答：所有的类型都可以。关键是不论是写还是更新，邮箱对于用户来说应该都是一样的。对于这种邮箱来说，最简单的实现应该是基于主备份的本地写协议，该协议中主备份总是位于用户的移动电脑上。

**19、请描述出一个用于显示刚被更新的 Web 页面的写后读一致性的简单实现。**

答：最简单的实现是让浏览器检查它所显示的是不是最新版本的页面，这需要发送一个请求到 Web 服务器。由于许多系统都已经实现了该方案，所有它非常简单。

**20、请给出以客户为中心的一致性易于导致写操作冲突的案例。**

答：在描述以客户为中心的一致性模型时，假设每一个数据项都只有一个拥有着，只有该拥有着具有写权限。如果放弃这个假设的话肯定将导致写操作冲突。例如，如果数据的两个独立用户最终都被绑定到了同一个副本上，他们各自的更新都需要被传播到这个副本上。这是，更新冲突就很明显了。

**21、使用租用时，客户和服务器的时钟分别紧密同步是否是必要的？**

答：不需要。如果客户端认为它的时钟和服务器的时钟同步有困难时，它会在当前服务器的租借到期之前获取一个新的租借。

**22、考虑一个用于保证分布式数据存储上顺序一致性的非阻塞的主机备份协议。这样的数据存储总能提供写后读一致性吗？**

答：不能。如果更新进程接收到它的更新请求已被处理的通知，它就从数据存储断开并重新连接

另一个副本。更新是否应用在副本上则没有通知。相比来说，如果使用一个阻塞协议，更新进程仅在他的更新请求被完全应用在其它副本上时才断开连接。

**23、为了使主动复制正常工作，所有操作必须以相同的顺序在每个副本处被执行。这个顺序总是必需的吗？**

答：不是。对于那些对只读数据的操作或者可交换的写操作，从理论上来说，这些情况允许不同的副本有不同的顺序。但是，要检查两个写操作是否为可交换的十分困难。

**24、为了使用序列器实现全序的广播，一种方法是先将操作转发给序列器，然后该序列器为这个操作分配一个惟一的号码，随后广播这个操作。请提出另外两种可选择的方法，并比较这三种方法。**

答：另一种方法是对操作进行广播，但是延迟发送直到序列器为该操作广播一个 *ixuliehao*。后者在该操作被序列器接收后发生。第三种方法是首先从序列器获取一个序列号，然后广播这个操作。

第一种方法（将操作发送给序列器），涉及到发送包含操作信息的点到点的消息和一个广播信息。第二种方法需要两个广播消息：一个包含操作信息，另一个则包含序列号。第三种方法需要一个包含序列号的点对点的消息，和一个包含操作信息的广播消息。

**25、一个文件被复制在 10 个服务器上。请列出表决算法允许的所有读团体和写团体。**

答：下面这些是合法的：

(1, 10), (2, 9), (3, 8), (4, 7), (5, 6), (6, 5), (7, 4), (8, 3), (9, 2), and (10, 1).

**26、为防止复制的调用，我们在书中讨论了基于发送方的方法。在基于接收方的方法中，是由接收副本辨认属于同一调用的输入消息的多个拷贝。请描述基于接收方的方法如何防止复制的调用。**

答：假设对象 A 调用对象 B，如以前，每一次调用指定的 A 的副本相同的唯一的标志。每一个副本多点调用并请求 B 的副本。当 B 的副本接收到调用请求时，它检查它是否已经接收到过 A 的副本的请求。如果没有，调用就会在本地副本上执行，然后对 A 的副本发一个回应的广播。如果调用已经执行，请求信息就会进一步被忽视。同样地，A 的一个副本只有在获得对象的本地副本的显著调用请求第一个引入回应时才执行，而此后相同的调用请求都会被简单地忽视。

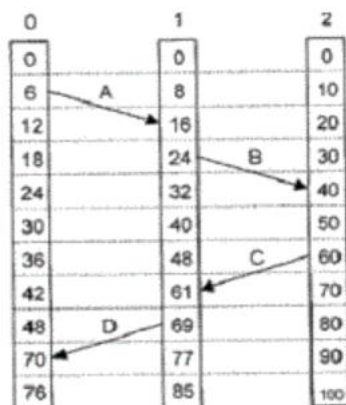
**27、考虑在因果一致的懒惰复制中，一个操作被从一个写队列中删除的确切时间。**

答：当知道此操作已经在各处都被执行，这可以通过追踪每一个  $L_j$  中的副本最新的一次升级来检测，如果每一个副本都执行过了来自  $L_i$  的  $W$  操作， $L_i$  可在更小的时间内先移去  $W_2$  的所有操作，（比如：每一个  $W'[k] \leq W[k]$ ）。因此，每个  $L_j$  发送一个回应给  $L_i$  其中包含最近执行的操作的时间。

**28、在实现一个支持广播 RPC 的简单系统时，我们假设系统有多个复制的服务器，每个客户可以通过 RPC 与一个服务器通信。但是，处理复制时，客户需要向每个副本发送一个 RPC 请求。设计客户程序，以便客户好像只向应用程序发送单一的 RPC。假设复制的目的是为了提高性能，而那些服务器可能是易于出故障的。**

答：无。

**29. 在下图中加入一条与 A 并发的新消息。它既不在 A 之前，也不在 A 之后。**



答:在 0 时刻从进程 2.发送消息, 并且在 8 时刻到达进程 1, 或者是在 0 时刻从进程 1 发送消息并在 10 时刻到达进程 2。这两种方法都可以满足要求。

30. 假定有两台机器 A 和 B, 它们的时钟都是每毫秒滴答 1000 次。但实际上 B 每毫秒滴答 900 次。如果每分钟根据 UTC 时间校正一次时钟, 那么 A 和 B 之间的最大时钟偏差是多少?

答:第 2 台计算机始终每秒钟发送 990000 次脉冲, 误差为 10 毫秒, 一分钟, 这个误差会增至 600 毫秒。

31. 对所介绍的 bully 选举算法, 进行优化。

答: 欺负算法存在的问题是网络通信代价和选举时间开销, 欺负算法时随机的选择比自己号码大的节点进行消息发送, 优化的思路是, 每个节点都存储其他节点的编号及对应的地址信息, 当发起者发送选举消息时, 首先向编号最大的节点发送选举消息, 如果

接收到应答消息, 则此次选举成功, 由编号最大的节点广播协调者消息给其他所有节点, 如果没有收到应答消息, 则重新选择一个编号第二大的节点, 进行上述操作。

32. 假定 A 和 B 是相互独立的两个临界区, 进程 0 要进入 A, 进程 1 要进入 B, R-A 分布式互斥算法会导致死锁吗? 说明理由。

答:这取决于基础规则, 如果说各个进程严格的按照线性时间顺序来访问资源, 也就是说一个已经获取某种资源的进程不会在尝试获取另外一种资源, 这样的活他也就没有办法在拥有另外一个进程打算获取的资源的时候去阻塞了。那么, 系统就不会产生死锁了。另一方面, 如果进程 0 在已经获取资源 A 的情况下又尝试获取资源 B, 那么如果另外一个进程此时按照相反的顺序进行操作就有可能产生死锁。因此, 算法本身并不会产生死锁, 因为每一种资源的处理过程相互之间都是相互独立的。

33. 假设有标识符分别为 1-8 的 8 个进程, 其中进程 8 崩溃, 节点 5 发起 bully 算法进行选举, 完成选举需要发送消息的个数为多少?

答: 15 个

34. 无线网络系统的选举算法中, 完成选举的延迟为多少?

答: 其生成树的最大高度\*2 或者该无线网络的直径\*2

## 第七章 分布式一致性与复制管理

1. 在内存一致性模型的讨论中, 经常提到软件和内存的约定。为什么需要这样的约定, 举例说明?

答:如果程序期望连续的一致性的数据存储且不能缺少这些任何数据, 存储必须保证连续一致。但是, 为了提高性能, 一些系统提供弱一点的模式。事实上软件能遵守模式强加的规则。总的说来, 这意味着遵守规则的程序可以感知到什么看起来像是连续一致的数据存储。如果进程同意遵守某些合约, 那么数据存储将正常运行。正常情况下, 一个进程执行读操作时, 它期待该操作返回的是该数据在其最后一次写操作之后的结果。在没有全局时钟情况系啊, 精确定义哪次写操作是最后一次很困难。因此产生了一些了一致性模型, 每种模型有效限制在一个数据项上执行一次读操作所应返回值。

2. 下图为顺序一致性内存的一个例子。对 P2 做少量改动, 使它违反顺序一致性。

**P<sub>1</sub>: W(x) 1**

**P<sub>2</sub>: R(x) 0 R(x) 1**

答: P<sub>2</sub>: R(X) 1 R(X) 0

### 3. 在电子股票市场的交易过程中, 需要提供哪种一致性, 说明理由.

答: 因果一致性可能是足够的。问题是, 对股票价值变动的反应应该是一致的。在不同的订单中独立的股票的变化。

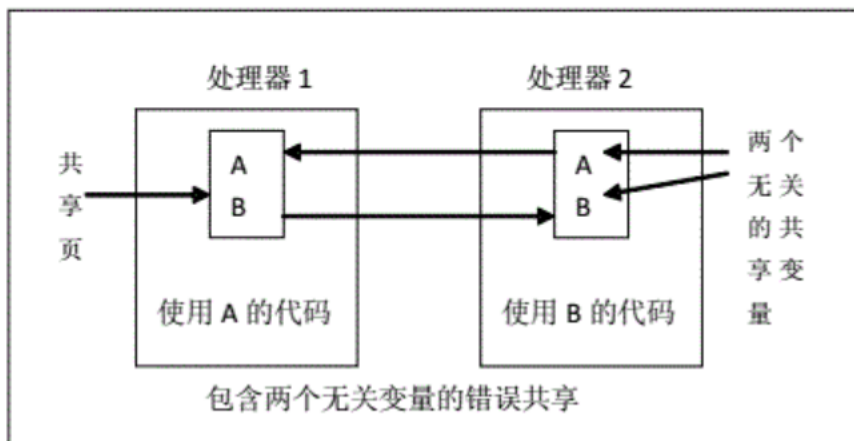
4. 释放一致性的大多数实现方法是在 release 时同步共享变量, 而不是在 acquire 时同步, 但为什么还需要 acquire 操作?

答: 获取操作 (acquire) 需要用于延迟进程当它试图访问共享变量时而另一个进程正在对共享变量进行此操作。获取操作时进程需要等待其他进程对共享变量完成释放操作, 这样保证该进程得到共享变量当前值。否则不能保证该进程会获得共享变量当前值。

5. 假设两个变量 a 和 b, 恰好位于基于分页的 DSM 系统的同一页上。然而, 它们都不是共享变量。是否会发生错误共享?

答: 错误共享是指无关的变量出现在同一页上, 当一进程使用它们之一时, 进程也得到了其他变量。有效页越大, 发生错误共享的可能性越大。假设有一页包含了两个无关的共享变量 A 和 B, 如下图所示, 进程 1 经常对 A 执行读写操作, 进程 2 经常使用 B, 这时包含两个变量的页将经常在两个进程之间传送。

“错误共享”问题是由于有效页太大产生的, 有效页越大, 发生错误共享的可能性越大, 共享页越小, 发生错误共享的可能性越小。



6. 在如下并行执行的进程 P1 和 P2, 列出顺序一致性所允许的 6 种语句交叉执行情况。

<b>a=1;</b>	<b>b=1;</b>
<b>If (b==0) kill (P2)</b>	<b>if(a==0) kill (P1)</b>
<b>(a) P1</b>	<b>(b) P2</b>

答: (1) a=1; if (b==0); b=1; if (a==0);

(2) a=1; b=1; if (a==0); if (b==0);

(3) a=1; b=1; if (b==0); if (a==0);

(4) b=1; if (a==0); a=1; if (b==0)

(5) b=1; a=1; if (b==0); if (a==0);

(6) b=1; a=1; if (a==0); if (b==0);

### 7. 浏览网页时, 为了保证读自己写一致性, 在系统中如何实现?

答: 它可以保证, 如果用户刷新页面, 他们总会看到自己刚提交的任何更新。它不会对其他用户的写入做出承诺,

其他用户的更新可能稍等才会看到，但它保证用户自己提交的数据能马上被自己看到。

最简单的方案，对于某些特定的内容，都从主库读。举个例子，知乎个人主页信息只能由用户本人编辑，而不能由其他人编辑。因此，永远从主库读取用户自己的个人主页，从从库读取其他用户的个人主页。

如果应用中的大部分内容都可能被用户编辑，那这种方法就没用了。在这种情况下可以使用其他标准来决定是否从主库读取，例如可以记录每个用户最后一次写入主库的时间，一分钟内都从主库读，同时监控从库的最后同步时间，任何超过一分钟没有更新的从库不响应查询。

还有一种更好的方法是，客户端可以在本地记住最近一次写入的时间戳，发起请求时带着此时间戳。从库提供任何查询服务前，需确保该时间戳前的变更都已经同步到了本从库中。如果当前从库不够新，则可以从另一个从库读，或者等待从库追赶上来。