

第十二章 多核与多处理器

在第 11 章，我们介绍了如何在多核硬件中使用合适的同步原语来保证程序的正确性。然而，多核硬件被简单地抽象为数个单核处理器。然而在实际中，多核处理器并非数个单核的简单叠加。多核硬件暴露给软件除了多个可用的计算单元之外，还有一系列硬件特性。这些硬件特性将对软件的性能造成巨大的影响。此外，在实际使用中，除了单一处理器中可以拥有多个核心，有的系统还拥有多个处理器，因此展现出更加复杂的硬件特性。

本章将首先介绍在操作系统视角下，多核多处理器硬件的部分特性，包括：缓存一致性、内存一致性模型以及非一致性内存访问。然后，我们将讨论操作系统开发者应该如何利用这些硬件特性，以获得最大的性能提升。

12.1 缓存一致性

本节主要知识点

- q 多核多处理器中高速缓存是什么样的？
- q 如何保证这些高速缓存中数据的一致性？
- q 缓存一致性对于操作系统开发者来说意味着什么？

12.1.1 多核缓存一致性结构

当前主流的多核处理器均采用了共享内存：不同的核心共享相同的内存资源，核心间可以通过访问同一个地址来共享数据。然而由于访问内存耗时较长，共享内存不是直接将多核处理器连接到同一个物理内存，而是添加多级缓存（如 L1、L2 和 L3 level Cache）来缓存高频访问的数据。访问缓存的延迟度表如下：

于访问内存的频度,因此使用缓存可以降低访问内存的速率从而减少访问开销。

我们首先回顾一下单核处理器中高速缓存的架构。缓存使用 (Cacheline) 作为最小的操作粒度,其大小与字长相等。在多级缓存中需要读一个地址的值时,处理器将逐级查找高速缓存中是否保存了该地址对应的缓存行。如果在任意一级缓存中找到,处理器将直接读取该值,从而避免耗时的内存访问操作。而当需要写一个地址的值时,处理器有多种策略可选。直写策略 (Write Through) 将修改的值刷入内存 (缓存值同时保留在缓存中); 写回策略 (Write Back) 则将修改的值暂时存在缓存中,避免频繁的内存写操作。当缓存行被替换 (Eviction) 时,或是 CPU 核心调用写回指令时,修改才会被更新至物理内存。

图 12.1: 多核下的缓存结构示意

而在多核处理器中,一级缓存型的多级缓存架构如图 12.1 所示。每个处理器核心均有自己的私有一级缓存 (Level 1 Cache, L1 Cache)¹, 多核心之间共享一个二级缓存 (Level 2 Cache, L2 Cache), 而系统级核心共享最后一级缓存 (Last Level Cache, LLC)。这种设计能够保证每个核心访问缓存时达到更好的性能。在这种架构下,不同核心访问时依据缓存行在缓存中的位置呈现别称。图 12.1 中的核心访问本级二级缓存时的缓存行名称快于访问系统级核心共享二级缓存时的缓存行。这种缓存架构称为非统一缓存访问 (Non-Uniform Cache Access, NUCA) 架构。系统的缓存设计除了指导非统一缓存访问以外,还会带来数据一致性问题。由于不同核心均拥有一级缓存 (含一级缓存), 某一地址上的数据可能同时存在于

¹通常 L1 Cache 进一步分为指令缓存 (Instruction Cache) 和数据缓存 (Data Cache)。

核心的一级缓存中。当这些核心同时使用写回策略修改该地址的数据时，会导致不同核心上一级缓存存储的地址数据不一致，违反了共享内存的抽象。为了保证私有缓存之间也能就某一地址的值达成共识，多核硬件提供了缓存一致性协议（Cache Coherence Protocol）。

12.1.2 目录式缓存一致性

缓存一致性协议有多种实现方案，目录式缓存一致性（Directory-based Cache Coherence）和 Snoop-based Cache Coherence。缓存一致性是由硬件保证的，其对于上层系统软件是透明的。这里通过介绍其中一种：目录式缓存一致性，来展现缓存一致性的基本硬件原理。

图 12.2: 目录式缓存一致性协议的简化硬件结构示意图

图12.2展示了使用目录式缓存一致性协议的简化硬件结构示意图。该图展示了多个核心拥有自己的私有缓存，且每个核心可以通过缓存力总线访问内存。系统将通过讨论这组简化模型向读者展示缓存一致性的基本实现原理。多级私有缓存，由于缓存之间关系更为复杂，因此需要遵守更多规则，但其背后的实现原理是一致的。每个私有缓存存储的某条缓存行，除了其所属缓存单元存储的区域存储缓存行的状态。这组缓存行状态可以概括为独占（Modified）、共享（Shared）以及失效（Invalid）等。这些状态的含义以及在不同状态之间迁移的条件将在下文详细介绍。除此之外，全局缓存单元共享的目录，用于记录每个缓存行所处的状态。某条缓存行对应了目录项，其包含两项内容：用于记录

²在简化模型中，单级访问操作导致的缓存一致性操作是原子的。实际硬件需要其机制保证并发情况下缓存一致性的正确性。此外，该简化模型也不考虑缓存写回操作。

是否处理器已纠正修改过这缓存行的 1 表示被修改, 0 表示未被修改), 以及用于记录缓存行的拥塞者的 Vector (拥塞者对应缓存行被逐出)。全局目录、高速缓存与内存三者均通过高速内部互联总线相连。目录式缓存一致性协议纠正是通过在全局共享目录记录缓存行在不同核上的状态, 确保对于缓存行的修改互斥且能够及时对全局可见, 从而保证了缓存行的一致性。下节将通过分析缓存行纠错行读操作在状态之间的迁移, 来详细介绍目录式缓存一致性协议的具体原理。

图 12.3: MSI 状态迁移

图12.3展示了 MSI 三种不同的状态以及当发生读写操作时状态的迁移。当一处理器要访问的缓存行在本地缓存中时, 缓存行可能处于以下状态:

1. **独占修改**。这状态代表着当前缓存行在全局只有本处理器一份拷贝。因此当前的核心独占缓存行, 可以直接进行操作, 不触发缓存行的状态变化。
2. **共享**。这状态代表当前缓存行在全局可能存在多份拷贝, 且本处理器的拷贝是有效的。因此当前核心可以直接读该缓存行。如果处理器要写缓存行, 则当前核心要查找全局共享目录, 找到拥有该缓存行拷贝的核心, 并通知这些核心将缓存行状态转换为独占。同时, 目录项的 Bit 位, 更新拥塞者的 Bit Vector。最后, 将本处理器的缓存行状态转换为独占修改, 方能对缓存行进行写操作。
3. **失效**。这状态代表当前缓存行本处理器的拷贝失效, 不能直接读/写缓存行。如果处理器要写, 则应当在目录找到拥有该缓存行的核心, 向其索要缓存行的数据, 同时通知核心将缓存行状态改为独占修改。最后, 在本处理器更新目录项的 Bit 位, 并更新拥塞者的 Bit Vector。最后, 在本

核心 P_0 将拿到的缓存行设置为脏状态，方能读取缓存行。啥果 P_1 要写缓存行，则 P_1 要通过目录找到脏缓存行的核心，通知脏缓存行的核心，脏缓存行状态才能拿到缓存行的数据，并更新目录脏的状态。最后，将本缓存行的脏状态设置为脏，方能写缓存行。

上藐啥逐情况均唔 P_1 要访唔的缓存行已纠在脏缓存行。脏缓存行不在脏缓存行逐时，脏核心将查看全局共享目录，检查缓存行是否在其缓存核心上。啥果其缓存核心拥脏缓存行的拷贝，则处力流呈同缓存行唔一致。啥果其缓存核心上也藐脏，则脏缓存行被标记为脏(Write Back)。此时 P_1 要从内存逐获取脏缓存行，放到本缓存的脏缓存行，同时更新全局的共享目录。而该缓存行状态将根据操作类型为读或写，设脏缓存行状态为脏或脏。占修改现代硬件为了更好的性能，会采用更加复杂的协议，拥脏更多的状态，同时也需要保证并发正确性并支持缓存写回操作，但其基本思路与我们介绍的简化模型一致。

下藐通过一缓存具体例子，展示目录式缓存一致性的具体操作流程。图12.4逐袁啥缓存核心，藐缓存核心脏袁自己脏袁的缓存行，分0、1、2高速缓存。这些核心共享一个全局目录项。图逐的力子逐关逐脏址脏在的缓存行，表项包含脏中保存的值以及该缓存行状态。图逐共展示了 $T_1 \sim T_6$ 时刻藐缓存核心的缓存行内容袁状态，以及对应时刻全局共享目录项中关于该缓存行的元数据。 T_1 时刻唔初始状态 T_1 而 T_6 时刻，不同核心将逐行缓存核心上方标逐的指令。 T_1 时刻，核心0逐行指令 $\& * \& *$ ，脏指令代表脏向脏逐写啥值，硬件将依照袁圈哑号依次完成操作，保证缓存一致性。

小思考

在图12.4的力子逐时刻，核心1需要写地址*脏在的缓存行，更新其值88以覆盖旧的值，唔何唔哑脏值传脏自己？能否省略脏步骤？

当多缓存核心对脏逐行写操作时，也不能由于缓存发生数据覆盖而省略传输旧缓存行的步骤。这是由于缓存一致性是以缓存行为粒度进行的，而一缓存行的大小通常呈字节，写操作唔唔逐作用于其逐的一小部分。新的修改不会覆盖整个缓存行，依旧需要缓存行其他部分的旧值。因此 P_1 要等待核心0将整缓存行发脏脏自己缓存(图12.3)，才能修改。

图 12.4: 目录式缓存一致性协议示例

12.1.3 系统软件视角下的缓存一致性

缓存一致性是硬件提供的、对上层软件透明的硬件特性。一些读者可能会疑惑：“既然如此，为何系统软件设计者要知道这些硬件知识？”这是因为系统软件的性能依赖于硬件实现，系统软件要通过了解真实硬件实现，总结出其特性，才能针对这些特性设计出高性能系统软件。这些特性如下：

1. 多核心对于同一缓存行的修改将导致严重的性能开销。当多核心要修改同一缓存行时，需要缓存一致性协议来保证一致性。由于缓存一致性协议同一时刻只允许一核心独占修改缓存行，制造串行流行为，无法充分发挥出多核的性能优势；此外，多核心对于同一缓存行的修改会导致总线流量，从而造成性能瓶颈。因此，系统软件开发者要避免可能对单一缓存行的竞争。系统软件将在多线程中以自锁避免互斥锁，介绍由于对共享变量在缓存行的竞争导致的可扩展性问题，并介绍如何通过可扩展互斥锁来避免这个问题。
2. 伪共享(False Sharing)在多核造成严重的性能开销。伪共享是指本身在不同核心之间共享的内容被错误地划分到同一缓存行，并引起了多核缓存行下的对于同一缓存行的竞争，从而导致伪共享的性能开销。例如在软件中直接使用整型数组（如`static int S[1024]`）伪共享核心提供一核心独占的计数器，线程按照顺序在核心更新这些计数器。由于不同核心将更新不同的计数器，这些计数器本身不是共享的。但如果直接分配一核心整型数组，这些计数器很可能落到同一缓存行中。当运行在不同核心的多线程同时更新这些计数器时，就会导致额外的缓存一致性开销。因此，系统软件开发者要避免伪共享的发生，比如可将不需要共享的每核心本地数据分配到不同缓存行。
3. 多核环境下局部性同样重要。局部性包含时间局部性和空间局部性。时间局部性是指访同一缓存地址的数据呈现在一段时间内被反复访问相同的地址；空间局部性则是指相邻的内存地址很可能被访问。良好的局部性能够保证较高的缓存命中率，减少访问开销。而局部性较差的应用不制造大量的缓存命中率，导致巨大的访问开销，缓存命中率出其不意地降低，进一步影响性能。在多核多线程任意核心的局部性问题也会导致共享的缓存命中率下降（如共享的LLC）受到影响，从而影响整个系统的性能，可牵一发而动全身。软件开发者要注意让软件具备良好的局部性，高效使用高速缓存。

更多内存一致性对于系统的性能以及可扩展性的影响，将在第12.2节讨论。

12.2 内存一致性与硬件内存屏障

本节主要知识点

- q 有哪些常见的内存一致性模型？它们有哪些特性？
- q 常见的架构中使用的是什么样的内存一致性模型？
- q 在弱序一致性模型上编程，怎么保证访存操作可见顺序？
- q 不同架构采用不同的内存一致性模型对于操作系统开发者来意味着什么？
- q 访存操作到底为何会出现乱序可见？x86架构和ARM架构上又有什么区别？

12.2.1 多核的访存乱序

现代处理器往往允许指令乱序执行。对于间隔时间的访存指令，处理器可以选择任意顺序执行其指令，从而隐藏访存操作的时延。然而，乱序执行意味着指令将不再按照呈顺序执行，访存操作的乱序可能导致“全局可见”（即被其他核心观测到）最终导致执行结果不符合预期。

³注意，即使处理器不乱序执行访存操作，访存操作的乱序也可能被其他核心观测到。将在第12.2.5节进行详细介绍。

图 12.5: 菩璫 特啥璫 塌璫 法前身One 塌璫 法

图12.5展示了1; / 7! : 1塌璫 法, 其是璫璫介绍的皮特森算法的前身。这力璫使用1; / 7! : 1算法是为了方便读者理解, 访存乱哑璫 在菩璫 特啥璫 塌璫 法逐璫 也悔璫 导逐璫 错误, 其袁璫 因袁璫 7! : 1算法完全一致。1; / 7! : 1塌璫 法使用一组标记位3来表示两个线程是否希望进入临界区。因此, 在申请纠璫 啥璫 临界区时, 线呈璫 将先设逐璫 自己的标记梧璫 , 然后检查对方的标记位。啥璫 果对方同时想要纠璫 啥璫 临界区, 则哑璫 要等待。在藐璫 袁璫 访存操作乱哑璫 的前提下7! : 1塌璫 法保证了悔璫 斥访梧璫 , 也即两个线程不会同时进入临界区。哑璫 要逐璫 意的是, 不同于菩璫 特啥璫 塌璫 法塌璫 法逐璫 保证了悔璫 斥访梧璫 , 袁璫 等袁璫 空闲让纠璫 啥璫 果两隔璫 线呈璫 同时希望纠璫 啥璫 临界区, 且在互相读对方的标记位之前, 都已经设置了自己的标记位。此时, 两个线程都不能进入临界区, 陷入了无限等待。

在现代处理器上, 由于悔璫 出现访存乱哑璫 的情况, 因此7! : 1塌璫 法的悔璫 斥访梧璫 梧璫 法保证。假设在线袁璫 线呈逐璫 , 第行的写操作袁璫 第行的读操作发生了乱哑璫 , 导致检查标记位的操作在自己的写标记全局可见(即对方的读操作可以读到其纠璫 果)之前发生。此时啥璫 果线袁璫 线呈璫 同时申请纠璫 啥璫 临界区, 它们乱序后的读操作可能读到对方的标记位均为2- 8? 1, 导逐璫 两隔璫 线呈璫 同时纠璫 啥璫 临界区, 打菩璫 7! : 1的互斥访问的保证。

12.2.2 内存一致性模型

内存一致性模型(Memory Consistency Model, 简呈璫 梧璫 内存模型)藐璫 确栋璫 义了不同核心对于共享内存操作需要遵循的顺序。读写操作之间共有四类先后顺哑璫 哑璫 要保证, 即读操作与读操作的顺序、读操作与写操作的顺序、写操作与读操作的顺序、写操作与写操作的顺序(梧璫 了简便, 下文分别简写为读读、读

写、写读、读写)。针对同一地址或者数据依赖关系访问操作,处理器可以保证其顺序。因此在不同的内存模型下,要讨论的是不同地址及数据依赖关系的访问操作之间的顺序。下面本节将按照从强到弱的顺序,依次介绍四种不同的内存模型。

严格一致性模型 (Strict Consistency)

严格一致性模型是最强的内存模型,也是最符合开发者直觉的内存模型。在严格一致性模型下,每个处理器访问操作都是按照呈现实际编写的顺序。此外,其要求所有核心对一个地址的任意读操作都能读到这个地址最近一次写的数据。因此,每个核心看到的访问操作顺序与其发生的时间顺序完全一致。在这模型下实现的互斥算法一定能保证互斥访问。但是,实现严格一致性模型要求使用全局统一的时钟,不同核心上执行的访问指令的时间先后顺序,增加了系统的实现难度。

顺序一致性模型 (Sequential Consistency)

顺序一致性模型弱于严格一致性模型,其不要求操作按照其真实发生的时间顺序(即依据全局时钟定义的顺序)全局可见。顺序一致性模型提供了以下保证:首先,不同核心看到的访问操作顺序完全一致,这保证了顺序一致性;其次,在这模型全局顺序下,每个核心自己的读写操作可见顺序必须与其程序顺序保持一致。不同于严格一致性模型,顺序一致性模型放松了对于时间顺序的要求,因此其中的读操作不一定能读到其他核上最新的修改。

图12.6展示了顺序一致性模型下运行互斥算法可能的结果。这力使用了局部变量 x 来分别存储读取的对方标记的值。图中纵向箭头代表最慢的全局顺序,横向箭头则表示访问操作在全局顺序中发生的时机。注意,这力黑的箭头并非代表实际发生的具体时间,纠缠代表在全局顺序下的先后顺序。图中列出来的a、b、c三种运行顺序都满足顺序一致性要求,因此对于任意一次运行,其最慢全局可见顺序可能是c中的任意一种。此时, x 的值可能: , 表。而对于图的最慢情况,由于其要求线号的操作乱序到操作之可见,违背了呈现实顺序,违背了顺序一致性模型的要求,因此不可能出现。所以,顺序一致性模型能够保证互斥算法的互斥访问。而如果使用严格一致性

⁴包括数据依赖、地址依赖等。我们将在第12.2.5详细介绍。

⁵严格一致性模型认为两个操作不可能在同一时刻发生。

图 12.6: 顺序一致性模型下的 LockOne 崩溃法

致性模型，则访存操作的可见顺序必须与执行顺序完全一致。因此对于任意一次逐行， r 的值都只有一种可能，其由访存操作的执行顺序决定。

TSO 一致性模型 (Total Store Ordering)

为了达到更好的性能，TSO 一致性模型进一步弱化了访存一致性保证。在 TSO 一致性模型中，其保证对不同地址且梧依赖读、读写、写写操作之间的全局可见顺序，逐行袁写读的全局可见顺序得不到保证。TSO 一致性模型通过加入一个写缓冲区达成优化性能的目的，该写缓冲区能够保证写操作按照顺序全局可见。TSO 一致性模型 (Total Store Ordering) 也因此得名。梧藐将在第 5 节进一步详细介绍。因此在图 12.6 崩溃法逐行 TSO 一致性模型允许图 12.6 逐行最悔一逐情况发生，即写操作乱序到达，操作之后可见，破坏了 r 的互斥访问保证。

弱序一致性模型 (Weak-ordering Consistency)

啥哑一致性模型 (梧简呈梧啥哑一致性模型) 提供了较致性保证。在一隔核心上，弱序一致性模型不保证任何不同地址且梧依赖存操作之间的顺序，也即读读，读写，写读与写写操作之间都可以乱序全局可见。这个特性导致在使用弱序一致性模型的设备上运行代码很容易出现违背开发者意图的情况。在啥哑一致性模型逐行，崩溃法同样不能保证互斥访问。除此之外，很多其他使用共享内存进行同步的程序也存在正确性问题。

代码第 1 展示了一逐基于共享内存的消息传栋机制。两隔哑要通过哑的线呈分别运行 r 袁 r 代码段。发塌者先填充数据，再通过标记 28-3 来通知接收者数据准备就绪 (第 6 行)。梧了保证消息被正确的传栋，发送者需要保证写数据与写 28-3 之间的顺序，而接收者需要保证读 28-3 与读


```
5: @ 0- @
5: @ 28- 3    ž! & $, ~ ~ +
B; 50 <>; / ~ B; 50

0- @
28- 3    $, ~ ~ +

B; 50 <>; / ~ B; 50

C4581 28- 3    $, ~ ~ +
      哑 悔 忙等
4- : 081 0- @
```

代码苦 段2.1: 基于共享内存的消息传栋

数据之间的顺序。在 TSO 一致性模型中，写写与读读的顺序均可以得到保证，因此接收者能够读到正确的数据。而在弱序一致性模型中，由于写写与读读操作之间均允许乱序全局可见，在这种情况下，接收者可能读到错误的数

表12.1总结了本小节中介绍的四种内存模型对于不同地址、梧 依赖的访存操作可见顺序的保证。

表 12.1: 不同内存模型对于不同地址、无依赖访存操作可见顺序保证

	读读	读写	写读	写写
严格一致性模型	X	X	X	X
顺序一致性模型	X	X	X	X
TSO 一致性模型	X	X	×	X
弱序一致性模型	×	×	×	×

12.2.3 内存苦 障

不同于缓存一致性，内存模型对于上呈 啥 不是透明的。在较啥 的内存一致性模型中，梧 了保证特栋 访存操作的全局可见顺哑 ，开发者必须手 添加件内存苦 障Barrier/Fence，简呈 内存苦 障)。硬件内存苦 障可以要求硬件保证访存操作之间的顺序。硬件往往提供多种不同的内存屏障指令来保证不同类型的访存操作的顺序。我们将在12.2.5节详细介绍内存屏障的工作原理。

代码苦 段2.2展示了通过添加内存苦 障来保证数据 @ 能够正确栋 从

```

5: @ 0- @
5: @ 28- 3      ž! & $, ~ ~ +
B; 50 <>; / ~ B; 50

0- @
.- >>51>
28- 3      $, ~ ~ +

B; 50 <>; / ~ B; 50

C4581 28- 3      $, ~ ~ +
      哑 悔 忙等
.- >>51>
4: 081 0- @

```

代码苦 段2.2: 内存苦 障示力

逐 行; / ~ 的线呈 传到逐 行/ ~ 的线程。这力 假设使用 啥 哑 一逐 性模型，即不保证 梧 依赖的读写操作之间的顺序。因此在设 3来隔 知数据纠 绪之前，哑 要保 此时已经全局可见。这力 使用一隔 内存苦 障行第 保证写0- @ 表 28- 3之间的顺哑 ，来完呈 这隔 目标。同样的; 在 逐 ，读28- 3表 读的顺哑 也哑 要通过添加一隔 内存苦 障来保证。而啥 果使用 TSO 一致性模型，读读和写写乱序都是不允许的，因此不哑 要加啥 任何内存苦 障。总的来说，不同的架构悔 隔 据自身使用的内存模型，提供不同的内存苦 障指令。软件通过调用对应指令告诉硬件来保证特定类型访存操作之间的顺序。

12.2.4 呈 见架构使用的内存模型

哑 隔 的内存模型对开发者更加友好，处力 器的行梧 更容易被开发者力 解。但其会造成处理器设计复杂，导致制造成本高、处理器能效低。而较啥 的内存模型硬件设计简单，能够挖掘更多的并行潜能。但开发者必须更加小心栋 添加硬件内存屏障来保证程序的正确性。除此之外，对于同步哑 求大的并行应用呈 哑 ，在弱序一致性模型下频繁使用内存屏障会带来显著的性能开销。因此在多维度的权衡之下，不同的处理器制造商在其架构中使用了不同的内存模型。啥 Intel 袁 AMD 在 x86 架构下栋 使用了较强的 TSO 一逐 性模型而 ARM 架构处力 器则使用 啥 哑 一逐 性模型。哑 择是隔 据 ARM 应对的呈 纠 (桌藐 袁 移栋)不同，考量处理器性能、功耗与成本等因素后确定的。2展 示了几种常见的体系结构使用的内存模型。

表 12.2: 几逐 呈 见架构使用的内存模型

	弱序一致性模型	TSO 一致性模型	顺序一致性模型
体系纠 构	ARM PowerPC	x86	Dual 386 MIPS R10000
使用呈 纠	嵌啥 式, 手 机, 桌 面 高能效服务器	嵌啥 式, 手 机, 桌 面 高性能服务器	被淘汰

在 x86 使用的 TSO 模型下, 纠 有写读操作悔 出现乱序。因此逐 哑 要在特 栋 情况下添加内存菩 障即可。ARM 处力 器上的呈 哑 纠 藐 袁 这藐 幸运了。由于其使用的是啥 哑 一逐 性模型, 必须逐 意塌 袁 的梧 依赖的访存之间是否哑 添加内存屏障保证可见顺序。在实际呈 纠 逐 , 逐 袁 涉及到多核协同隔 作时, 才 哑 要考虑访存操作的可见顺哑 , 使用合适的内存屏障。而哑 要同步的应用梧 梧 使用同步原语来保证正确性, 其隐含了内存屏障的语义。在实现这些同步袁 袁 时需要根据不同的内存模型添加合适的内存屏障。此外, 对于一些梧 塌 的设计 (啥 梧 塌 队列), 我们也需要考虑添加内存屏障来保证其正确性。

12.2.5 硬件视角下的内存模型袁 内存菩 障

为了达到更好的性能, 现代处理器往往允许访存操作乱序执行。在单隔 核 心逐 , 处力 器使用了多逐 技术保证呈 哑 逐 行纠 果袁 按照呈 哑 顺哑 依 致。

首先, 处理器保证了有依赖关系的访存指令之间的顺序。比啥 , 如果写操作的数据依赖于前序读操作的结果, 则写操作必须等待读操作结束后才能开始执行。啥 9< <@> <@> @9< 这两隔 操作逐 , 写 入 栋 址逐 的值依赖于读<@> 的值的纠 果, 则呈 这两隔 操作 数据依赖关系。除了袁 数据依赖, 悔 藐 址依赖如写操作需要写的地址依赖于读操作读出来的值 \ 控逐 依赖如写操作只有在读操作结束分支满足后才能执行⁶)。

除此之外, 现代处力 器悔 设重排哑 悔 缓冲区-Order Buffer, ROB), 让指令按照呈 哑 顺假 (Retire)。这力 , 退役对应顺哑 逐 行逐 的逐 行纠 束, 其意味着该条指令对系统的影响终将全局可见。啥 果袁 的指令由于分支袁 测得到提前逐 行, 但最逐 分支袁 测错误时, 由于分支判断袁 句悔 梧 退役, 这些指令不悔 提前退役, 处力 器纠 可以通过追踪到这些被错误逐 行的指令, 舍弃

⁶需要注意的是控制依赖不保证所有访存操作之间的顺序, 需要结合具体硬件进行分析。

图 12.7: 现代处理器简化微体系结构示例

其执行结果并拿取正确的指令重新执行。因此，在单核上运行的应用并不会受到影响，这些程序展现出的行为与严格按照程序顺序执行的行为相同。

然而，在多核系统下，由于系统核心也可以观测到当前核心的运行结果，上述这些措施无法保证被其他核心观测到的行为与本地核心行为一致。其次，访存指令要完成在 12.1 节中介绍的缓存一致性流程后，才能顺利地与其它核心观测到。因此，如果访存指令等待缓存一致性流程结束后再退役，则会阻塞后续指令进入重排序缓冲区，导致性能受损。

现代处理器设计了另一套缓冲区来解决这问题。图 12.7 显示，处理器在每个核心旁设置单表 (Load/Store Unit, 简称为 LSU) 中预留了读缓冲区与写缓冲区。这缓冲区用于暂存待满足缓存一致性的访存指令。在这设计下，访存指令不再需要等待耗时的缓存一致性流程，而是放入对应的读缓冲区或写缓冲区后，指令可以先行完成。此时，访存指令等待前一条指令退役，且保证当前指令一栋可以逐行 (分支预测正确)，便可以退役。但要注意的是，此时访存指令退役并不代表其可以被其他核心观测到，其还在读写缓冲区中等待完成缓存一致性流程。退役指令一栋可以逐行，并一定在未来可以被其他核心观测到。这意味着，能否让其核心按照顺序观测到另一核心的逐行结果

重任落到了读缓冲区与写缓冲区身上。我们将一个访存操作完成缓存一致性流程、真正变得全局可见的过呈（提交 Commit）。提交与退役并不相等，一个访存操作需要等到其从重排序缓冲区退役后才会提交。

本节将详细介绍 x86 架构下的内存模型的微体系结构实现，并介绍内存屏障指令如何保证访存操作可见顺序。注意，由于 x86 架构与 ARM 架构的手册描述商业硬件的真正实现，因此本节所述的实现逐条是符合手册行规的一条实现可能，真实硬件逐条可能有不同的实现。

x86 架构下的 TSO 一致性模型

x86 架构采用了较强的 TSO 一致性模型。处理器逐条的读缓冲区与写缓冲区被用于维护不同访存操作全局可见的时机，因此 x86 架构的处理器逐条这两隔缓冲区也被命名为内存顺序缓冲区（Memory Ordering Buffer）。

处理器首先维护 x86 架构下的写缓冲区。写缓冲区用于缓存已逐条行结束但尚未变得全局可见的写操作。这缓冲区的逐条行结束包括隔条将指令提交给 LSU，处理器通过重排缓冲区退役的指令，以及已退役但是处理器能变得全局可见的写操作。处理器这规则，如果一个写操作满足了缓存一致性条件，真正力开写缓冲区到达 cache，则代表隔条写操作的因果全局可见，呈隔条写操作成功提交。对于任意写操作，离开写缓冲区需要满足以下几个前提条件：

1. 首先，隔条写操作必须要退役。一旦退役，当前核心将缓存啥处理器这隔条写指令已经完成，该指令不可撤销。
2. 其次，隔条写操作缓冲区的缓存一致性流呈必须结束。处理器在第 1 节介绍的缓存一致性协议中，逐条当前缓存行处于“独占修改”状态时，该写操作才能执行。
3. 最后，x86 架构的写缓冲区按照先啥先出的顺序提交写操作。也即一隔条写操作必须等到前条写操作力开写缓冲区之缓存才能力开。因此，处理器使用的 TSO 模型逐条写操作之间的顺序能够得到保证。

我们再来看读缓冲区。x86 处理器允许任意的读操作乱序执行。读操作并不存在提交的观念，但仍需要按照先入先出的顺序离开读缓冲区（即退役）。但是，如果某个在读缓冲区中的读操作还未退役，且此时隔条操作目标缓存行被其他核心修改，其会受到缓存一致性中的“失效”命令影响，舍弃当前读操作读到的结果，重新逐条行隔条读操作。处理器因此不会出现读读乱序（乱序逐条行的读操作读到的结果顺序逐条一行一逐条，否则缓存被要求重做）。因此，一隔条读操作的退役前提是：

图 12.8: TSO 一致性模型中四类不同的操作组合行为

1. 首先，该读操作需要读取的值必须已经被读取到该核心。
2. 其次，程序顺序中该操作之前的所有操作必须已经退役。
3. 最后，读操作从读到值开始到退役之间必须收到目标缓存行“失效”的命令。若在此期间收到了目标缓存行“失效”的命令，则需要重新执行读操作。

在图 12.8 展示了 x86 架构的读缓存冲突和写缓存冲突，下图展示了不同力度的访存操作组合在 TSO 一致性模型下的行为。其逐行或分别表示先后的两个访存操作。

写写。正如在介绍写缓冲区时讨论，由于写缓冲区保证了写操作会按照程序顺序提交，因此读操作必须按顺序全局可见。全局顺序一致性模型下，逐行或分别表示先后的两个访存操作。

读读。以代码片段 12.1 为例，假设乱序执行的读操作（第 13 行，图 12.8）在读操作（第 11 行，图 12.8）之前发生，且读到的是初始数据而非修改的。而此时读操作读到

了 $\langle \rangle$; / 修改的 S_{i+1} , 那藐 由于 $\langle \rangle$; / 可以保证写 O_{i+1} 到写 S_{i+1} 的提交顺哑 (即保写的可见顺哑), 因此此时 O_{i+1} 的写操作一定也已经提交完成。而根据缓存一致性协议, 此时对于 O_{i+1} 悔 存行悔 存失通知一栋 已纠 到达了 $\langle \rangle$; / 所在核心。由于读缓冲区按程序顺序退役, 此时 $\langle \rangle$; / 读28-3操作隔 隔 完呈, 悔 梧 退役 @ 有读存行操作也藐 袁 完呈 退役。塌 以, 对于 O_{i+1} 悔 存行悔 存失通知悔 使得 $\langle \rangle$; / 重做隔 读操作 (图逐时刻), 并读到修改值 。通过以上策略, x86 处力 器纠 实现了读操作的顺哑保证。

读写。在 x86 处理器中, 保证读写操作的顺序较为简单。由于重排哑 悔 冲区已纠 保证了塌 袁 指令悔 按照呈 哑 顺哑 退役, 而写操作的提交操作一定发生在退役之悔, 因此当一隔 写操作全局可见时 (即提交), 塌 之前的读操作一栋 已纠 完呈 并退役了 (通过保证, 图逐时刻)。塌 以 TSO 架构能保证对于读写操作之间的顺序。

写读。最悔 对写读操作, 由于当写操作退役时, 其值可能悔 藐 袁 变得全局可见 (图逐时刻, $\%S$ 退役但悔 在写悔 冲区逐), 但啥 其悔 哑 的读操作已纠 满足了退役的条件, 造成读操作在写操作真正变得全局可见之前退役 (图逐时刻), 最终破坏了写读之间的顺序。如果此时有使得读操作目标缓存行失效的命令到达, 写读乱序将会被其他核心所观测到。

梧 了避藐 乱哑 x86 处力 器提供 921: / 1 指令用于避藐 写读乱序。在梧 藐 上述的这逐 实现逐, 可以简单啥 梧 悔 阻啥 悔 哑 指令发塌 到 LSU, 直到前哑 塌 袁 访存操作完成。12.5.2 节的 7!: 1 塌 法梧 力, 啥 果在两隔 线呈 的 行之间添加了 921: / 1 指令, 读对方 28-3 的操作必须要等到写自己的 28-3 操作提交之悔 才能发塌 到 LSU 并逐 行。因此能够保证 1: / 7!: 1 算法的互斥访问。

ARM 架构下的弱序一致性模型

ARM 架构处力 器哑 用了啥 哑 一逐 性模型。啥 哑 一逐 性模型不保证任何依赖且针对不同栋 址的读写操作之间的顺序。这力 哑 要逐 意架构袁 一栋 的特殊性。塌 然市藐 上很多处力 器栋 使架构, 但是塌 藐 处力 器内部实现可能悔 袁 藐 显呈 别。这 对于处力 器提供了很多可以由呈 商决栋 啥 何实现的栋 方, 其手册逐 栋 义了最逐 行为。因此, 对于特栋 的处力 其内部遵循的内存模型可能更加严格。同样, 下藐 梧 藐 纠 绍的一逐 实现也逐 是遵哑 ARM 手册的一种简化实现。

ARM 架构处力 器塌 然也袁 写悔 冲区袁 读悔 冲区的 塌 架构

处理器不同，这两类缓存不再提供顺序的一致性保证，写缓冲区逐条的写操作可以不按照“先入先出”的顺序离开写缓冲区，当对应的写指令退役且目标缓存行缓存一致性流程结束后，该写操作便可以离开写缓冲区，变得全局可见。因此，对于任意的无依赖且针对不同地址的两个写操作，其全局可见的顺序不再受到约束。

ARM 架构下的读操作的退役机制。ARM 架构也存在很大差异，x86 架构下必须等到真正的值读啥处理器内部（寄存器）且前写指令退役后才能退役。如果在等待过程中被缓存一致性协议标记为失效，则需要重新执行。因此，x86 架构能保证读写的顺序。而在 ARM 架构下，处理器可以在确保该读操作一栋缓存发生（啥分支预测呈隔）且前写指令退役时，就将该读操作退役。此时，这隔读操作要读到的值可能已经被读到处理器中。以读操作当然按照顺序退役，但读操作发生的时机可能会违背程序顺序读呈。

读写读这两情况，其顺序更无法得到保证。这是由于读操作与写操作退役之后都不要求其真正执行完毕（需要读的值到达处理器或写的值全局可见），因此它们之间真实发生的顺序均没有严格要求，很可能会发生乱序。

如果要在 ARM 架构下访存操作之间的顺序，则几条不同的选择。首先，像 x86 架构，ARM 提供了一系列硬件内存屏障指令来满足不同的顺序要求。其最常用的一类内存屏障指令是全屏障（Full Memory Barrier）。O9 指令能够保证其之前的访存指令及其之后的访存指令之间的顺序。O9 指令要隔一隔，缀用于表示其影响范围影响的访存操作力型。对于影响范围，ARM 定义运行操作系统的处理器均属于，即 Inner Shareable Domain，因此在操作系统一颁使用后缀。对于影响的访存操作力型 ARM 提供了啥逐力型：一逐是影响读的访存操作，这逐缀增加任何新的关键字，啥 O9. 5?4；一逐是保 O9 之前的写操作到其之后的写操作之间的顺序，这逐缀 O9. 5?4?@；最一逐是保证 O9 之前的读操作到其之后的写操作之间的顺序，这逐缀要添加 80 关键字，啥 O9. 5?480。处理器执行 O9 指令时，根据其缀阻啥对应的访存指令发到 SU，直到前缀对应的访存操作能够保证在指栋的范围可见的相对顺序。这力的相对顺序是指前缀指令并非已全局可见，而是下层硬件能够保证同一核心后续的访存操作一定会在当前指令之后全局可见。除了 O9 外，ARM 还提供 O7.、80-、?@8 等指令，本书不做过多介绍，感兴趣的读者可以参阅 ARM 手册 [4]。

除了硬件内存屏障，在 ARM 架构逐条还可以通过构造依赖关系来保证访存

操作的顺序。由于处理器逐行指令时，要先等待前指令依赖的访存操作完成，所以存在数据依赖、地址依赖或控制流依赖时，访存操作之间的顺序能够得到保证。

表 12.3: ARMv8 架构保序方案

	读读	读写	写读	写写
顺序一致 (Seq. Cons.)	X	X	X	X
顺序一致 (Seq. Cons.) 480	X	X	×	×
顺序一致 (Seq. Cons.) 480 @	×	×	×	X
顺序一致 (Seq. Cons.) >	X	X	×	×
顺序一致 (Seq. Cons.) @8>	×	X	×	X
数据依赖	×	X	×	×
地址依赖	X	X	×	×
控制流依赖	×	X	×	×
控制流依赖 57.	X	X	×	×

表12.3展示了在 ARMv8 架构下不同的保序方案以及其适用的场景。其中 X 表示能够保证，而 × 则代表不能保证。开发者需要根据具体使用场景选择合适的保序方案来保证程序的正确性。

12.3 非统一内存访问

本节主要知识点

- 什么是非统一内存访问 (NUMA) 架构？它与统一内存访问 (UMA) 架构有何区别？
- NUMA 架构有哪些特性？

随着单处理器中核心数量增多以及多处理器系统的出现，单一的内存控制器逐渐呈现出了性能瓶颈。因此多核及多处理器系统将多个内存控制器分布在不同的核心或处理器上。这种设计在多处理器系统中非常常见，每个处理器往往被分配一个单独的内存控制器。因此，不同的核心可以更快地访问其本地的内存，其也可以访问其他处理器上的远程内存，其访问时间要长于访问本地内存的时延。由于访问不同物理内存的时延不同，这种架构被称为非统一内存访问 (Non-Uniform Memory Access, NUMA)。

图 12.9: 非一致性内存访问示意图

图 12.10: Intel 和 AMD 示例服务器拓扑结构图

图12.9是一组 NUMA 系统的示意图，该系统包含两个 NUMA 节点。在这组示意图系统中，一个处理器上的多个核心（每个核心）访问内存的特性一致，因此被分到同一 NUMA 节点（NUMA 节点）上。任一节点上的任意核心能够快速访问本地的内存。一旦其需要访问远端节点的内存时，其需要通过系统总线（interconnect）或远端节点通信，造成了较长时间的延迟。NUMA 架构有多种组织方式，任一节点可以是一组物理处理器，也可以是处理器中一部分核心。

现代单台服务器为了应对更高的计算需求，可以插入多个处理器，而任一处理器中同时拥有多个核心，因此其 NUMA 架构也更为复杂。图12.10展示了两个分别基于 AMD 架构和 Intel 架构的示例服务器拓扑结构图 [10]。这两个服务器都拥有 8 个处理器插槽，分别对应着 8 个 NUMA 节点。由于节点数量较多，某些节点之间可能直接相连，请求需要通过两跳才能到达目标节点（如 Xeon 和 AMD Opteron 服务器中的节点之间均需要两跳）。因此访问时间差异也更加复杂。

表 12.4: 表 呈 内存访问 (ns) 时延特性

fi: ?@	4; <	4; <	4; <
核心 fi: @18 * 1; : 服务器			
l; - 0 %@, >1			
核心 ~Ž" !<@1>; : 服务器			
l; - 0 %@, >1			

表12.4 [10] 展示了这两塌 服务器上访梧 本栋 以及表 端内存的时延。表逐 数据的单梧 处理器周期 (Cycles), 数字越大代表访问内存的时延越高。表隔 共表 啥 列, 分别代表访梧 本栋 的内存(访梧 一跳的表 呈 内存(, 啥 Intel 服务器中节点 0 上的处理器访问节点 1、2、5 上的内存), 表 访梧 两跳的表 呈 内存hop, 啥 Intel 服务器逐 节栋 上的处力 器访梧 节栋 6、7 上的内存)。可以从表逐 看到, 访梧 一跳的表 呈 内存时哑 呈 较本栋 的内存访梧 慢倍以上, 而访梧 两跳的表 呈 内存访梧 更是达到了。可见, 啥 果应用塌 用内存塌 机分布在不同节栋 上, 表 逐 访梧 本栋 内存的应用相比, 隔 应 将面临严重的远程内存访问开销。

表 12.5: 内存访梧 带宽(MB/s)

~ // 1??	4; <	4; <	4; <	fi: @1>81- B10
核 fi: @18 * 1; : 服务器				
%l=A1: @5- 8 S- : 0; 9				
核 ~Ž" !<@1>; : 服务器				
%l=A1: @5- 8 S- : 0; 9				

除了访梧 时哑 表 塌 不同, 表 呈 内存访梧 的带宽上限也悔 受到内部悔 力的限制。表12.5 [10] 展示了这两台服务器本地及远端内存访问带宽的区别。由于内部悔 力 总线的带宽限逐 , 表 端的内存访梧 的带宽较本栋 塌 失应用哑 要访梧 大量表 呈 内存时, 访存带宽较访梧 本栋 内存更容易呈 梧 性能颈。

NUMA 架构对于操作系统不是透藐 的, 因此梧 了应对下呈 硬件的非一逐

内存访问，操作系统应当分配、管理好硬件资源，以避免应用出现复杂的呈现内存访问造成的性能消耗。目前已解决很多NUMA问题提出了设计方法，通过分配应用本地的内存，避免使用远端的内存来降低对于应用性能的影响。更多细节我们将在第12.4.2节进行讨论。

12.4 操作系统性能可扩展性

本节主要知识点

Q 为何软件会有可扩展性问题？其背后的原因是什么？

Q 如何设计拥有良好可扩展性的软件？

上一节我们介绍了操作系统视角下的三种硬件特性，包括缓存一致性、内存一致性模型、非一致性内存访问。本节将结合系统领域一些前沿问题，介绍系统软件的研究人员如何针对这些特性优化系统软件的性能，提升操作系统的性能可扩展性。什么是性能可扩展性？理想的可扩展性是随着核数增加 N 倍，软件的性能也能提升 N 倍。不过这明显是不切实际的。比如，在现实世界中，一个人不可能一天造起一栋房子，如果将人数增加100倍、1000倍，也无法将造房子的速度提升1000倍。在并行计算领域，阿姆达尔定律（Amdahl's Law）[3] 来描述并行计算的加速比。

$$S = \frac{1}{(1-p) + \frac{p}{s}}$$

公式描述了加速比与可以并行的部分占比（ p ）和 s 可以并行部分的加速比。在理想情况下，如果 $p=1$ ，那么加速比可以达到 N 。如果 $p < 1$ ，那么加速比会小于 N 。而当 $p = 1$ ，应用的理想加速比会在足够的核数下趋近于 $S = \frac{1}{(1-p)}$ 。因此对于多核运行的应用，如果本身无法并行部分 $1-p$ 占比较大，添加核心数能带来的加速比便十分有限。此时需要应用本身逻辑想办法，以提升应用的可扩展性。除了算法本身导致的可扩展性问题，上一节介绍的多核多处理器硬件特性同样也会影响应用的可扩展性，甚至出现应用在多个核心上运行的性能不及单个核心的情况。下面将分别介绍在多核环境下由于对单一内存竞争造成的可扩展性问题，并从同步原语角度介绍如何改善这两种场景下的性能可扩展性。

图 12.11: 华为泰山服务器（鲲鹏 916）互斥锁可扩展性测试

C4581

```
8; /7 38; . - 8 8; /7
38; . - 8 /: @
>- : 0; 9 - //1?? /- /4185: 1 38; . - 8 /- /4185: 1
A: 8; /7 38; . - 8 8; /7
:: <?
```

代码片段 2.3: 微基准测试

12.4.1 单一缓存行高度竞争导致的可扩展性问题

在多核环境中，应用不能并行的部分需要使用互斥锁来保证程序的正确性。但是，这些同步原语在引入可扩展性问题。图12.11展示了在华为鲲鹏ARM架构服务器上使用互斥锁面临的可扩展性问题。在这项测试中，互斥锁在每颗核心上运行一个测试程序。代码段显示了微基准测试每一线程的逐行流程。线程将不断执行一个互斥锁循环。线程进入，线程将竞争一个全局的互斥锁；一旦获得互斥锁，线程在临界区中将更新全局的计数器；然后读取并修改全局互斥锁行38; . - 8 /- /4185: 1用于模拟应用临界区对于全局共享变量的修改。而两个临界区之间，线程将执行一段数量的空操作，用来模拟非临界区代码。图逐行的逐行从左逐行右增加并行逐行的核心数逐行则代表对应核心数下测试程序的吞吐率。从图逐行可以发现，当增加核数时，使用自旋锁（Spin Lock，在8.1.4节中介绍）的测试程序性能达到峰值。一旦增加更多的核数，其性能不仅不会上升，反而会出现断崖式下跌。

系统领域的研究者通过对互斥锁性能建模，详细分析了这个问题。研究者发现，当执行核心数超过一定阈值后，线程会出现同一时刻多个核心在等

待获取锁的情况，且等待的线程数量随着核心数的增加而增多。但是在理想情况下，即使更多核心同时竞争锁，也应该导致性能逐步提升，并不损害性能，更不会出现如图 12-15 所示的断崖式下跌。性能断崖式下跌的现象实际上是由于在第 12.1 节介绍的缓存一致性导致的。由于自旋锁（代码清单 8-7）是通过修改全局单一变量 `lock` 来代表获取以及释放锁，因此自旋锁的获取与释放操作均会造成对单一内存行拥权的竞争（同一时刻，至多一个核心能够“独占修改”这一内存行）。而当多核心对同一内存行竞争时，缓存一致性开销十分巨大。这是由于自旋锁在等待期间不断读取并尝试修改锁变量，当同一时刻竞争核数增多时，内存行的状态拥塞者也不断改变，最终消耗大量时间在内存一致性协议上导致互斥锁无法快速有效地在不同的核心之间传递。此时，锁在竞争者之间传递的开销（主要是缓存一致性开销）伴随着竞争者数量增加而增加。

既然断崖式下跌是由于同一时刻有太多竞争者在抢锁，减少这些竞争者数量就成为了解决可扩展性问题的关键。一种直观的、无需大量修改的策略：退避策略（Back-Off）。

退避策略的核心思想非常简单，当竞争者拿不到锁时，就不再继续尝试修改内存行，而是选择等一段时间再去拿锁。为了避免多个竞争者的等待时间相同导致再次修改时依然出现竞争，退避策略让竞争者设定了不同的等待时间。比如等待时机时呈指数分布或依照一定比例依次加乘等待时间。图 12-16 显示了使用指数退避策略的退避锁（Back-Off Lock）在华为服务器（鲲鹏 916）上微基准测试的性能。退避锁在竞争获取失败时，以指数时间退避。可以发现在少于 14 个核心时，退避锁的性能略高于自旋锁。这是由于退避锁在竞争程度不高时，就已经开始采用回退策略。此时会出现一些竞争者睡觉时被释放，且其他竞争者获取锁的情况，从而导致一定的性能开销。但当更多核心时，其性能超过自旋锁，且随着核心数增加，退避锁的吞吐率十分稳定，不再出现可扩展性断崖。

退避策略虽然解决了对于内存行的竞争，但并未能从根本上避免性能问题。一定程度减少了问题的出现。而良好的可扩展锁要保证其竞争开销（内存行失效的次数）不应随着竞争者数量增多而加大。下面将通过介绍可扩展的队列锁 MCS 锁，分析如何设计具有良好性能可扩展性的系统软件。

MCS 锁

MCS 锁[8]是由 John M. Mellor-Crummey 和 Michael L. Scott 在 1991 年提出，MCS 锁也因此得名（MC 分别指两位作者姓氏首字母）。同排号

```
B; 50  *^/1 B; 50  - 00>  B; 50  : 1C B- 8A1
```

```
B; 50  @9<      - 00>
- 00>      : 1C B- 8A1
>1@A>:  @9<
```

代码苦图 段2.4: XCHG 操作示意

塌 (读者可以悔 顾第1节) 力 似MCS 塌 拥袁 一隔 等待队列。梧 了避藐 对单一悔 存行的纠 拿CS 梧 藐 一隔 纠 争者栋 准备了一隔 节栋 , 并插啥 到一链表中。这样锁的持有者可以通过链表找到下一任竞争者并将锁传递。因此纠 争者逐 哑 等待在自己的节栋 上, 由前任塌 的呈 袁 者通过修改自己节栋 上的标记来完成锁的传递, 无需像排号锁一样通过竞争全局缓存行的方式来检查是否轮到自己。下面详细介绍 MCS 的加锁和放锁操作流程。

梧 了实现CS 队列塌 , 梧 藐 这力 引啥 一隔 新的袁/子操作: *^/1 的操作啥 代码 苦图 2.4塌 示, 塌 悔 0>栋 址上的值修改梧 新值: 1C B- 8A1 并返悔 隔 栋 址上袁 来的值95而 *^/1 则悔 保证交悔 操作的原子性。

代码苦图 2.5展示了 MCS 塌 的实现。隔 实现省去了内存苦 障, 要求访存操作哑 隔 按照呈 哑 顺 顺 顺 的袁 数据逐 记录着等待队列的队梧 指针@ 58。当塌 梧 释放状塌 时, 该指针为空。一旦袁 纠 争者出现, 隔 指针将指向等待队列的队尾。藐 一隔 塌 的纠 争者栋 哑 要创建自建的% : ; 01。隔 节栋 包括一隔 极逐, 梧 可以表示两逐 状塌 1&fz1, 表示当前节栋 对应的纠 争者应当等待; ž&~, 表示当前节栋 对应的纠 争者被授权可以纠 啥 临界区。此外, 藐 隔 节栋 悔 袁 一隔 指向等待队列的指针。逐 意, 梧 了保证藐 隔 节栋 栋 在不同悔 存行上 1.1 避藐 提到 在共享的梧 题, 这力 要求 : ; 01依照缓存行大小对齐。不同于之前的自旋锁, MCS 塌 的 /7袁 8; /7操作除了哑 要传啥 塌 的纠 构体, 悔 哑 要传啥 当前线呈 的CS 节栋 1。

当一隔 线呈 栋 /7 试获取MCS 塌 时, 其先初始悔 自己的等待队列节栋 (第3、14 行)。其悔 , 隔 线呈 力 5用 *^/1 操作将队梧 指针8; /7 @ 58交悔 梧 指向自己的指针, 并将袁 来的队梧 指针值8存啥 (第15 行)。啥 袁 来的指针梧 空, 代表隔 塌 梧 空闲状塌 , 此时隔 线呈 可以纠 啥 临界区逐 行。否则, 隔 线呈 将力 啥 自己的节栋等待在自己的 MCS 节栋 上 (第8 行)。

```
?@>A/@ Ž°%::; 01
B; 8- @581 5: @ 28- 3
B; 8- @581 ?@>A/@ Ž°%::; 01 : 1D@
- @>5. A@1 - 853: 10 °~°/ ċ ħž ċ %
```

```
?@>A/@ Ž°% 8; /7
?@>A/@ Ž°%::; 01 @ 58
```

```
B; 50 8; /7 ?@>A/@ Ž°% 8; /7 8; /7 ?@>A/@ Ž°%::; 01 91
```

```
?@>A/@ Ž°%::; 01 @ 58
91 : 1D@ ž' ħħ
91 28- 3 ) ~ ħ&ħž1
@ 58 - @, 95/ *°/1 8; /7 @ 58 91
52 @ 58
@ 58 : 1D@ 91
C4581 91 28- 3 1 $~ž&~
```

```
B; 50 A; 8; /7 ?@>A/@ Ž°% 8; /7 8; /7
?@>A/@ Ž°%::; 01 91
```

```
52 91 : 1D@
52 - @, 95/ °~°% 8; /7 @ 58 91 91
>1@A>:
C4581 91 : 1D@
```

```
91 : 1D@ 28- 3 1 $~ž&~
```

代码苦龘 段2.5: MCS 队列塌龘 实现

当呈袁塌的线程通过释放塌时,将先检查等待队列逐是否悔袁其塌线呈等待在隔塌第啥果袁其塌线呈等待,则依照力表顺哑,通过修改悔哑等待者节栋逐的标记梧来通知隔节栋纠啥临界区(第32行)。啥当前线呈已纠是等待队列力的最悔一隔,隔线呈则哑要袁子等待队列的队梧指针逐梧空,表示隔塌已纠释放。啥果袁子操作失颁,塌藐此时袁新的线呈已纠交悔了队梧指针,隔线呈哑要等待新纠并将塌传栋隔新的线呈行第

图12.12纠合一隔实际力子,展示了获取塌袁释放S锁的流程。

- 在 T_1 时刻,等待队列中一个有四个节点。此时,处于队首的纠争者拥袁塌,其标逐梧来。
- 在 T_2 时刻,有一个新的竞争者出现。其首先初始悔自己的节栋:将标记梧填梧,并将:1D@指针置为空。接着,隔纠争者将袁子交悔袁数据逐的队梧指针指向自己。其发现,之前队梧指针不梧空,而是指向袁等待队列的队梧,塌藐隔塌正被其塌纠争者呈有。因此,隔纠改袁队梧的指针指向自己,从而力啥等待队列,并等待前哑节栋将锁传递给自己。
- 在 T_3 时刻,持有锁的第一个线程希望释放锁。其发现等待队列逐存在其塌的等待者,因此通过修改悔哑等待者的标记梧来通知其纠啥临界区,完成锁的传递。之后,隔线呈可以悔收自己的节栋,或者留到下次继续使用。
- 最悔在 t_n 时刻,逐于轮到等待队列最悔一隔线呈放塌时,不巧又袁一隔新的线呈希望得到塌(S节栋)。由于新线呈悔藐呈隔交悔梧指针袁队梧线呈并不知晓新纠争者的存在。作梧队梧节栋,其哑要通过袁的CAS操作,判断队梧指针是否指向自己,如果指向自己则将其换为空。由于新的纠争者的存在,其也在呈试交悔队梧指针,因此袁队梧线呈的CAS操作可能悔失败。当袁CAS失颁,意梧着袁新的线呈已纠呈隔交悔了队梧指针。此时逐哑要等待新的线呈通过更新袁队梧的力啥等待队列,即可按照袁时刻相同的流程将锁传递下去。而啥袁子的CAS呈隔,此时意味着锁已经被成功释放。新的纠争者呈隔交悔队梧指针时,会发现队尾指针为空。因此其可以直接进入临界区。

下藐分析梧何这样设计能够给塌带来更好的可扩展性。如果一隔塌要具有良好的可扩展性,塌纠争的开销(啥悔存行失效的次数)不应塌着同一

图 12.12: MCS 操作示力圖

时刻竞争者数量增多而加大。在自旋锁、阻塞锁、逐行锁，阻塞锁着竞争者数量增多，单一行锁存行的竞争加剧，最终导致关键路径上平均每次获取锁造成的缓存失效数量急剧上升。而对 MCS 阻塞锁而自旋锁，当同一时刻竞争者数量增多时，力锁想情况下关键路阻塞锁上阻塞锁传栋栋逐锁涉及两次行锁失效：分别行锁被继任者修改行锁在行锁存行袁继任者初始行锁的：1D@ 28-3所在缓存行。当然，对于@ 58的竞争也袁可能暴露在关键路阻塞锁上，但当参袁锁竞争的线呈较多时，等待队列逐一颁栋已行锁力啥了多隔锁竞争者，因此阻塞锁可以快阻塞栋在等待队列逐传栋，在关键不行锁涉及对争 58 的修改。图 12.11 的微基准测试逐 MCS 阻塞锁展现出良好的性能可扩展性。可以看到在核心数增多时，MCS 锁的性能不会出现大幅下降。

小思考

MCS 锁性能一定比自旋锁好？我们应该怎么样在这两个锁之间选择？

自旋锁、阻塞锁、获取阻塞锁、袁释放阻塞锁、阻塞锁而的操作均非呈行锁简，在逐袁少数核其性能较 MCS 阻塞锁更好。因此这两隔阻塞锁并不存在好坏之分，行锁要依据具体呈行锁下的竞争程度来判断到底选用哪一种锁。操作系统逐很多行锁题栋袁其力似，很难找到一隔能应对阻塞锁袁呈行锁的最袁行锁决方颁，很多情况下行锁要对具体呈具体分析。系统研究者也想办法结合多种方案之间的优点，设计自栋依据具体呈行锁特征切行锁策略的行锁决方案。Linux 的 3.15 颁本逐，行锁引啥了 qspinlock (queue spinlock) 这一逐新的锁。这逐阻塞锁行锁合了自旋锁阻塞锁袁队列阻塞锁的袁势，在不同行锁争呈度下使用不同的策略，阻塞锁拥袁快阻塞路行锁与慢速路径。当一隔线呈呈行锁试获取隔阻塞锁时，其首先呈行锁试快阻塞路行锁，通过袁CAS 获取该锁。一旦获取失败，隔线呈将采取力阻塞锁的机逐，将自己加入等待队列。这种设计能适应竞争程度变化的复杂场景。当行锁争呈度低时，线呈可以走快阻塞路行锁快阻塞栋拿到阻塞锁，而当行锁争呈度高时，设计似能提供良好的性能可扩展性。

12.4.2 NUMA 架构中频繁远程内存访问导致的可扩展性问题

本节逐要知识栋

- q NUMA 架构为可扩展性带来什么样的新挑战？
- q 行锁在 MCS 阻塞锁在 NUMA 架构下会出现可扩展性问题？
- q 行锁藐隔啥何设计才能在 NUMA 架构下拥有良好的可扩展性？

图 12.13: 华为泰山服务器 (鲲鹏 916) 内存一致性 NUMA 环境可扩展性测试

上一节我们分析了由于对于单一内存行的竞争导致的内存一致性可扩展性问题。然而，内存一致性问题是呈隐性的（共享资源保障全局变量）本身同样可能造成可扩展性问题。比如，两不同核上的线程在临界区内访问了同一个缓存行。当锁在这两个线程之间传递时，缓存行也要在这两不同的核之间传递，造成一定的通信开销。这在 NUMA 架构下会更加显著。假设这两线程分别运行在不同的 NUMA 节点上，缓存行涉及的跨 NUMA 节点的内存一致性通信。此外，一旦缓存行被移出缓存，需要重新从内存读取时，访问远程内存会带来更加严重的性能开销，造成严重的可扩展性问题。在 NUMA 环境下，无论是使用之前介绍的排号锁或是锁的队列锁，缓存行按照申请获取的顺序依次传递。如果来自不同 NUMA 节点的竞争者依次申请获取缓存行，那缓存行的数据以及临界区内共享数据所在的内存行将在不同的节点之间传递，这将会非常耗时。

我们进一步验证这隐性问题，下文使用华为泰山服务器进行了系列测试。该服务器一共 8 个处理器插槽，每个插槽上插了 3 颗基于 ARM 架构的处理器。每颗处理器中又进一步分成了 2 个 NUMA 节点，因此整个系统一共 48 个 NUMA 节点，且每个节点有 6 个核心。图 12.13 展示了在内存一致性服务器上使用内存一致性锁面临的隐性问题。相较于前一小节的微基准测试，这里稍有不同。该代码段 6 显示，这力临界区将访问并修改共享缓存行的数据，而非之前另外，测试使用的核心数进一步从 32 核扩展到 64 核。可以从图 12.13 看到 12.4.1 节所述的隐性问题。在核心数上升时，其性能开始逐步下降。对比图 12.11 的结果，这力出现的性能下降要是由于线程临界区访问了更多的缓存行，而当核心数达到 64 核心时，缓存的竞争者可能出现在不

⁷主要为拥有缓存一致性保证的 NUMA 架构 (Cache Coherence-NUMA, CC-NUMA)。

C4581

```

8; /7 38; . - 8 8; /7
38; . - 8 /: @
2; > 5 5 5
>-: 0; 9 - //1?? /- /4185: 1 38; . - 8 /- /4185: 1 5
A: 8; /7 38; . - 8 8; /7
:: <?

```

代码片段 2.6: NUMA 环境下的微基准测试

同的 NUMA 节点。当线程呈现者在不同节点上时，临界区内访问的缓存行也要在不同的 NUMA 节点之间迁移，导致巨大的开销。

为了解决这问题，线程提供者提供了一组（NUMA-aware）的设计，用于增强访问操作的局部性，减少发生跨 NUMA 节点的缓存一致性。这组设计的核心，是在保证正确性的同时可能保证一段时间内访问能在本地节点（包括本地的缓存行和内存），从而避免在关键路径上出现过多的内存访问（包括通过缓存一致性访问在其节点隔邻进行）。为了进一步说明这组设计方针，这组文档用作一组详细的性能子分析啥设计的应用。隔邻是由 Dice 在 2012 年提出，其展现了在 NUMA 环境下优化可扩展性的基本思想。

Cohort 的设计核心是在一段时间内限制线程斥逐在本地节点内部传播，从而在关键路径上剔除耗时的跨节点缓存一致性。为达到这目的 Cohort 采取了两呈现的设计，第一呈现是唯一的全局锁，第二呈现是对应隔邻 NUMA 节点的本锁。当某一 NUMA 节点上的核心呈 Cohort 时，其需要同时持有全局锁以及本地锁。而当隔邻核心逐行线程操作时，啥果本节点呈现其线程竞争者，则隔邻核心不释放全局锁，而是释放本节点锁并将全局锁的权转隔邻下一隔邻本节点的线程竞争者，从而保证线程在一段时间内在本地节点之间传递。

图 12.14 展示了在 3 隔邻 NUMA 节点的系统逐锁传递的示例。最初，NUMA 节点 0 在核心 0 上的线程竞争者（记作 0）呈现袁锁，即同时呈现袁节点的本地锁以及全局锁。当其释放锁时，其将逐释放本节点锁，将其传播隔邻下一隔邻线程竞争者。在一段时间内，逐袁 NUMA 节点 0 内的线程竞争者可以获取锁。当所有的本地竞争者都执行完自己的临界区之后，本节点最悔的呈现袁者将释放本节点袁全局锁，将全局锁传播隔邻节点悔线程竞争者线程竞争者要重新获取全局锁才能进入临界区。

图 12.14: Cohort 塌缩 传栋 示力

```
?@>A/@ 8; /7 38;. - 8 8; /7
?@>A/@ 8; /7 ž' Ž~ 8; /7 ž' Ž~ ž' Ž

B; 50 >; A@5: 1 B; 50

5: @ 9E ;; 01_50 31@ >A : 5: 3 :; 01_50
52 8; /7 ž' Ž~ 8; /7 9E ;; 01_50 ž! ) ~ fi& $
8; /7 38;. - 8 8; /7
临界区
52 A: 8; /7 ž' Ž~ 8; /7 9E :; 01_50 ž! ) ~ fi& $
A: 8; /7 38;. - 8 8; /7
```

代码菩 段2.7: 使用 Cohort 塌缩 示力

显式地指定分配内存所在位置。其次，对于藐视袁使用这些接口的应用，操作系统还要纠正可能将内存分配在本节点，避免远程内存访问。最悔，操作系统维度时还要纠正可能避免节点的线程迁移。

12.5 了解并分析linux内核逐节点的UMA隔离知设计

理解了纠正一步纠正绍操作系统啥提供隔离知的支呈，这力以内核逐节点UMA隔离知的内存管理能力以及 NUMA 隔离知的节点度纠正行详细分析。

12.5.1 NUMA 隔离知的内存管理能力

针对 NUMA 悔纠正Linux 内核提供了多逐节点内存分配策略以及内存分配模式，以便在 NUMA 环境中合理地分配内存。

其逐节点内存分配策略指纠正了内存分配时，应隔离谁逐节点的规则（即内存分配模式）。这些策略包含了：由任务纠正呈自己指纠正自己分配内存时采用的模式（即任务指定策略）；以及直接隔离纠正特纠正虚拟纠正址空间，在隔离空间内纠正行内存分配时采用的模式（即虚拟纠正址空间指纠正策略）等。用悔可以分别使用内核提供的?1@ 919<; 85/E袁. 5: 0这两个系统调用来实现上述功能。啥用悔藐视袁纠正选择内存分配策略，则操作系统内核将接过决纠正分配模式的大旗，了解逐纠正呈纠正选择合适的分配模式。

而逐节点于内存分配模式Linux 内核提供了绑定模式、优先模式与交错模式。了解纠正模式顾名思义，其将从指纠正节点（通过上一段纠正绍的系统纠正用指纠正）上分配内存。应用可以使用隔离模式纠正选择在自己运行的节点上分配本地内存。而袁先模式则在分配失败时，呈试从指纠正节点最纠正的节点上分配内存。交错模式则会从给定的节点中以页为粒度交错地分配内存。啥果应用藐视袁纠正用任何模式，操作系统则悔采用袁先模式了解逐应用在其运行的节点及邻纠正的节点上分配内存。

总而言之Linux 内核从两方藐视针对NUMA 悔纠正下的内存分配纠正行了袁化。首先，Linux 内核提供了系统纠正用，让应用能够隔离据实际纠正求指纠正特纠正节点纠正行内存分配；而对于藐视袁使用这些接口的应用，内核则纠正可能在应用运行的节点上分配本地的内存。

图 12.15: 调度域示意图

12.5.2 NUMA 系统中的调度域

Linux 内核中的调度策略也考虑 NUMA 环境。对于多核多处理器硬件，Linux 内核中的调度器需要权衡两个问题：一方面需要避免可能让任务分布在系统的各个核心上，达到负载均衡；另一方面避免，避免频繁地在不同的核心之间迁移任务意味着丧失了任务的本地性，特别是在 NUMA 环境下，跨 NUMA 节点的迁移将导致巨大的迁移开销。因此 Linux 内核引入了调度域（Scheduling Domain）的概念，将 CPU 分成了不同的调度域。每个调度域是具有相同属性的一组 CPU 的集合，并根据硬件特征被分成不同的层级，呈现出一种树状结构。

图 12.15 展示了一个 NUMA 系统的调度域。最下层级是逻辑核调度域，每个调度域包含了同一组逻辑核，这些逻辑核共享本地内存，因此在域之间迁移的开销最低；向上一级是核调度域，每个调度域包含了共享本地内存的多个核心；再向上分别是处理器调度域（包含了同一个处理器中的所有核心）、NUMA 节点调度域（包含了整个 NUMA 节点的所有核心）、全系统调度域（包含了整个系统的所有核心）。内核在启动时根据硬件拓扑结构归为不同的域。不同的域有不同的负载均衡策略，越下层的域，任务在核心之间迁移的开销越低，逐层进行负载均衡越频繁。NUMA 节点域这一层级，负载均衡就很少执行。通过划分调度域的方法，避免可以在实现一个域的负载均衡的同时，避免在 NUMA 节点之间频繁迁移任务。

参考书目 献

- [1] Linux memory policy. 4@@<? CCC.71>: 18.; >3 0; /
"; /A91: @ @5; : B9 : A9- 919; >E <; 85/E.@D@.
- [2] Linux scheduler. 4@@<? CCC.71>: 18.; >3 0; / "; /A91: @ @5; :
?/410A81> ?/410 0; 9- 5: ?.@D@.
- [3] Gene M Amdahl. Validity of the single processor approach to achieving large scale computing capabilities. In *Proceedings of the April 18-20, 1967, spring joint computer conference*, pages 483–485, 1967.
- [4] ARM. Programmer ' s guide for armv8-a. 4@@<?
?@ @5/.0; /?.- >9./; 9 01: - " ċ ž ~ B
- >/45@1/@A>1 " 1 .<02, 2015.
- [5] Silas Boyd-Wickizer, M Frans Kaashoek, Robert Morris, and Nickolai Zeldovich. Non-scalable locks are dangerous. In *Proceedings of the Linux Symposium*, pages 119–130, 2012.
- [6] David Dice, Virendra J Marathe, and Nir Shavit. Lock cohorting: a general technique for designing numa locks. *ACM SIGPLAN Notices*, 47(8):247–256, 2012.
- [7] Danny Hendler, Itai Incze, Nir Shavit, and Moran Tzafrir. Flat combining and the synchronization-parallelism tradeoff. In *Proceedings of the twenty-second annual ACM symposium on Parallelism in algorithms and architectures*, pages 355–364, 2010.
- [8] John M Mellor-Crummey and Michael L Scott. Algorithms for scalable synchronization on shared-memory multiprocessors. *ACM Transactions on Computer Systems (TOCS)*, 9(1):21–65, 1991.
- [9] Mark S Papamarcos and Janak H Patel. A low-overhead coherence solution for multiprocessors with private cache memories. In *Proceedings of the 11th annual international symposium on Computer architecture*, pages 348–354, 1984.

- [10] Kaiyuan Zhang, Rong Chen, and Haibo Chen. Numa-aware graph-structured analytics. In *Proceedings of the 20th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 183–193, 2015.

多核与多处理器：扫码反馈



