

邮箱: dolphin0520@163.com

昵称: Matrix海子

园龄: 9年

粉丝: 8161

关注: 7

+加关注

最新随笔

1.深入理解Java中的不可变对象

2.再议Java中的static关键字

3.工程实践：如何给变量取一个好的名字

4.工程实践：给函数取一个"好"的名字

5.工程实践：如何规范地打印程序日志？

6.金山快盘+ TortoiseSVN构建版本控制仓库

7.在windows下安装配置Ulipad

8.Java并发编程：线程间协作的两种方式：wait、notify、notifyAll和Condition

9.JAVA多线程和并发基础面试题答（转载）

10.Java线程面试题 Top 50 （转载）

积分与排名

积分 - 470029

排名 - 560

随笔分类 (177)

Android开发(12)

C/C++(29)

Java NIO-入门篇(2)

Java并发编程-入门篇(19)

Java基础-进阶篇(2)

Java基础-入门篇(18)

Linux/Shell(3)

OnlineJudge(14)

Python(16)

工程实践系列(3)

计算机基础&网络(1)

Java并发编程：volatile关键字解析

Java并发编程：volatile关键字解析

volatile这个关键字可能很多朋友都听说过，或许也都用过。在Java 5之前，它是一个备受争议的关键字，因为在程序中使用它往往会导致出人意料的结果。在Java 5之后，volatile关键字才得以重获生机。

volatile关键字虽然从字面上理解起来比较简单，但是要用好不是一件容易的事情。由于volatile关键字是与Java的内存模型有关的，因此在讲述volatile关键字之前，我们先来了解一下与内存模型相关的概念和知识，然后分析了volatile关键字的实现原理，最后给出了几个使用volatile关键字的场景。

以下是本文的目录大纲：

一.内存模型的相关概念

二.并发编程中的三个概念

三.Java内存模型

四..深入剖析volatile关键字

五.使用volatile关键字的场景

若有不正之处请多多谅解，并欢迎批评指正。

请尊重作者劳动成果，转载请标明原文链接：

http://www.cnblogs.com/dolphin0520/p/3920373.html

一.内存模型的相关概念

大家都知道，计算机在执行程序时，每条指令都是在CPU中执行的，而执行指令过程中，势必涉及到数据的读取和写入。由于程序运行过程中的临时数据是存放在主存（物理内存）当中的，这时就存在一个问题，由于CPU执行速度很快，而从内存读取数据和向内存写入数据的过程跟CPU执行指令的速度比起来要慢的多，因此如果任何时候对数据的操作都要通过和内存的交互来进行，会大大降低指令执行的速度。因此在CPU里面就有了高速缓存。

也就是，当程序在运行过程中，会将运算需要的数据从主存复制一份到CPU的高速缓存当中，那么CPU进行计算时就可以直接从它的高速缓存读取数据和向其中写入数据，当运算结束之后，再将高速缓存中的数据刷新到主存当中。举个简单的例子，比如下面的这段代码：

```
1 | i = i + 1;
```

当线程执行这个语句时，会先从主存当中读取i的值，然后复制一份到高速缓存当中，然后CPU执行指令对i进行加1操作，然后将数据写入高速缓存，最后将高速缓存中i最新的值刷新到主存当中。

这个代码在单线程中运行是没有任何问题的，但是在多线程中运行就会有问题了。在多核CPU中，每条线程可能运行于不同的CPU中，因此每个线程运行时有自己的高速缓存（对单核CPU来说，其实也会出现这种问题，只不过是线程调度的形式来分别执行的）。本文我们以多核CPU为例。

比如同时有2个线程执行这段代码，假如初始时i的值为0，那么我们希望两个线程执行完之后i的值变为2。但是事实会是这样吗？

赞助

开发工具(7)
设计模式(1)
生活有感(1)
数据库(2)
算法与数据结构(36)
业余娱乐(1)
杂项(2)
装机数码系列(8)

随笔档案 (174)

2019年4月(3)
2019年3月(2)
2014年10月(2)
2014年9月(7)
2014年8月(15)
2014年7月(4)
2014年6月(10)
2014年5月(9)
2014年4月(5)
2014年1月(1)
2013年11月(4)
2013年10月(3)
2013年8月(2)
2013年6月(1)
2013年5月(3)
2013年4月(3)
2013年3月(11)
2013年1月(4)
2012年11月(3)
2012年10月(4)
2012年9月(3)
2012年8月(2)
2012年4月(1)
2012年3月(1)
2011年12月(2)
2011年11月(5)
2011年10月(20)
2011年9月(2)
2011年8月(3)

可能存在下面一种情况：初始时，两个线程分别读取i的值存入各自所在的CPU的高速缓存当中，然后线程1进行加1操作，然后把i的最新值1写入到内存。此时线程2的高速缓存当中i的值还是0，进行加1操作之后，i的值为1，然后线程2把i的值写入内存。

最终结果i的值是1，而不是2。这就是著名的缓存一致性问题。通常称这种被多个线程访问的变量为共享变量。

也就是说，如果一个变量在多个CPU中都存在缓存（一般在多线程编程时才会出现），那么就可能存在缓存不一致的问题。

为了解决缓存不一致性问题，通常来说有以下2种解决方法：

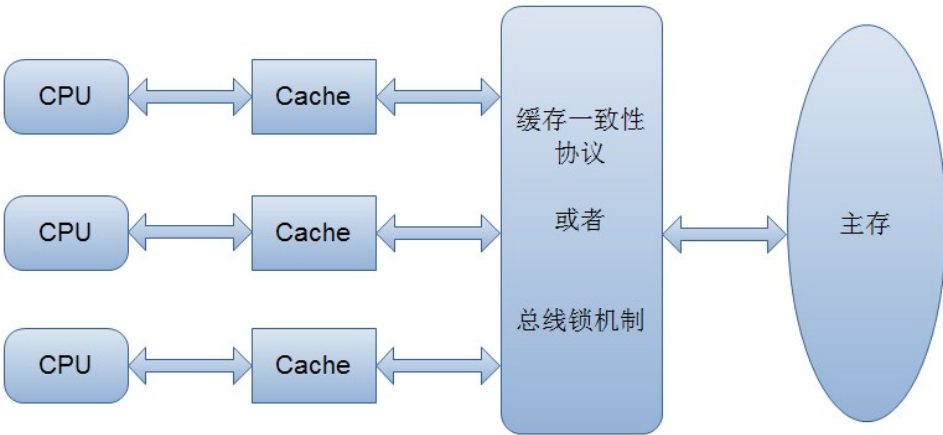
- 1) 通过在总线加LOCK#锁的方式
- 2) 通过缓存一致性协议

这2种方式都是硬件层面上提供的方式。

在早期的CPU当中，是通过在总线上加LOCK#锁的形式来解决缓存不一致的问题。因为CPU和其他部件进行通信都是通过总线来进行的，如果对总线加LOCK#锁的话，也就是说阻塞了其他CPU对其他部件访问（如内存），从而使得只能有一个CPU能使用这个变量的内存。比如上面例子中 如果一个线程在执行 $i = i + 1$ ，如果在执行这段代码的过程中，在总线上发出了LOCK#锁的信号，那么只有等待这段代码完全执行完毕之后，其他CPU才能从变量i所在的内存读取变量，然后进行相应的操作。这样就解决了缓存不一致的问题。

但是上面的方式会有一个问题，由于在锁住总线期间，其他CPU无法访问内存，导致效率低下。

所以就出现了缓存一致性协议。最出名的就是Intel 的MESI协议，MESI协议保证了每个缓存中使用的共享变量的副本是一致的。它核心的思想是：当CPU写数据时，如果发现操作的变量是共享变量，即在其他CPU中也存在该变量的副本，会发出信号通知其他CPU将该变量的缓存行置为无效状态，因此当其他CPU需要读取这个变量时，发现自己缓存中缓存该变量的缓存行是无效的，那么它就会从内存重新读取。



二.并发编程中的三个概念

在并发编程中，我们通常会遇到以下三个问题：原子性问题，可见性问题，有序性问题。我们先看具体看一下这三个概念：

1.原子性

原子性：即一个操作或者多个操作 要么全部执行并且执行的过程不会被任何因素打断，要么就都不执行。

一个很经典的例子就是银行账户转账问题：

比如从账户A向账户B转1000元，那么必然包括2个操作：从账户A减去1000元，往账户B加上1000元。

赞
助

2011年7月(10)
2011年4月(29)

最新评论
1. Re:Java并发编程：CountDownLatch、CyclicBarrier和Semaphore
@清风吹斜阳 getNumberWaiting方法里面，用了parties-count，而count是从parties一直减到0的，所以反过来说，getNumberWaiting就是加计数...
--aukocharlie
2. Re:Java并发编程：volatile关键字解析
牛逼！写mark，再研究！
--海纳分享
3. Re:Java ConcurrentModificationException异常原因和解决方法
@ 春冬秋...夏 这是我的理解，供你参考。 add方法添加两个元素，modCount(0) + 2 = 2；获得迭代器的时候expectedModCount = modCount = 2；remove一...
--xiaoBai3

阅读排行榜
1. Java并发编程：线程池的使用(727898)
2. Java并发编程：volatile关键字解析(492324)
3. Java中的static关键字解析(462452)
4. Java ConcurrentModificationException异常原因和解决方法(391552)
5. 深入理解Java的接口和抽象类(383684)

试想一下，如果这2个操作不具备原子性，会造成什么样的后果。假如从账户A减去1000元之后，操作突然中止。然后又从B取出了500元，取出500元之后，再执行 往账户B加上1000元 的操作。这样就会导致账户A虽然减去了1000元，但是账户B没有收到这个转过来的1000元。

所以这2个操作必须要具备原子性才能保证不出现一些意外的问题。

同样地反映到并发编程中会出现什么结果呢？

举个最简单的例子，大家想一下假如为一个32位的变量赋值过程不具备原子性的话，会发生什么后果？

```
1 | i = 9;
```

假若一个线程执行到这个语句时，我暂且假设为一个32位的变量赋值包括两个过程：为低16位赋值，为高16位赋值。

那么就可能发生一种情况：当将低16位数值写入之后，突然被中断，而此时又有一个线程去读取i的值，那么读取到的就是错误的数值。

2.可见性

可见性是指当多个线程访问同一个变量时，一个线程修改了这个变量的值，其他线程能够立即看到修改的值。

举个简单的例子，看下面这段代码：

```
1 | //线程1执行的代码
2 | int i = 0;
3 | i = 10;
4 |
5 | //线程2执行的代码
6 | j = i;
```

假若执行线程1的是CPU1，执行线程2的是CPU2。由上面的分析可知，当线程1执行 i = 10这句话时，会先把i的初始值加载到CPU1的高速缓存中，然后赋值为10，那么在CPU1的高速缓存当中i的值为10了，却没有立即写入到主存当中。

此时线程2执行 j = i，它会先去主存读取i的值并加载到CPU2的缓存当中，注意此时内存当中i的值还是0，那么就会使得j的值为0，而不是10。

这就是可见性问题，线程1对变量i修改了之后，线程2没有立即看到线程1修改的值。

3.有序性

有序性：即程序执行的顺序按照代码的先后顺序执行。举个简单的例子，看下面这段代码：

```
1 | int i = 0;
2 | boolean flag = false;
3 | i = 1;           //语句1
4 | flag = true;     //语句2
```

上面代码定义了一个int型变量，定义了一个boolean类型变量，然后分别对两个变量进行赋值操作。从代码顺序上看，语句1是在语句2前面的，那么JVM在真正执行这段代码的时候会保证语句1一定在语句2前面执行吗？不一定，为什么呢？这里可能会发生指令重排序（Instruction Reorder）。

下面解释一下什么是指令重排序，一般来说，处理器为了提高程序运行效率，可能会对输入代码进行优化，它不保证程序中各个语句的执行先后顺序同代码中的顺序一致，但是它会保证程序最终执行结果和代码顺序执行的结果是一致的。

比如上面的代码中，语句1和语句2谁先执行对最终的程序结果并没有影响，那么就有可能在执行过程中，语句2先执行而语句1后执行。

但是要注意，虽然处理器会对指令进行重排序，但是它会保证程序最终结果会和代码顺序执行结果相同，那么它靠什么保证的呢？再看下面一个例子：

赞
助

```
1 | int a = 10;    //语句1
2 | int r = 2;    //语句2
3 | a = a + 3;    //语句3
4 | r = a*a;      //语句4
```

这段代码有4个语句，那么可能的一个执行顺序是：



那么可不可能是这个执行顺序呢：语句2 语句1 语句4 语句3

不可能，因为处理器在进行重排序时会考虑指令之间的数据依赖性，如果一个指令Instruction 2必须用到Instruction 1的结果，那么处理器会保证Instruction 1会在Instruction 2之前执行。

虽然重排序不会影响单个线程内程序执行的结果，但是多线程呢？下面看一个例子：

```
1 | //线程1:
2 | context = loadContext(); //语句1
3 | initied = true;         //语句2
4 |
5 | //线程2:
6 | while(!initied ){
7 |     sleep()
8 | }
9 | doSomethingwithconfig(context);
```

上面代码中，由于语句1和语句2没有数据依赖性，因此可能会被重排序。假如发生了重排序，在线程1执行过程中先执行语句2，而此是线程2会以为初始化工作已经完成，那么就会跳出while循环，去执行doSomethingwithconfig(context)方法，而此时context并没有被初始化，就会导致程序出错。

从上面可以看出，指令重排序不会影响单个线程的执行，但是会影响到线程并发执行的正确性。

也就是说，要想并发程序正确地执行，必须要保证原子性、可见性以及有序性。只要有一个没有被保证，就有可能导致程序运行不正确。

三.Java内存模型

在前面谈到了一些关于内存模型以及并发编程中可能会出现的一些问题。下面我们来看一下Java内存模型，研究一下Java内存模型为我们提供了哪些保证以及在java中提供了哪些方法和机制来让我们在进行多线程编程时能够保证程序执行的正确性。

在Java虚拟机规范中试图定义一种Java内存模型（Java Memory Model, JMM）来屏蔽各个硬件平台和操作系统的内存访问差异，以实现让Java程序在各种平台下都能达到一致的内存访问效果。那么Java内存模型规定了哪些东西呢，它定义了程序中变量的访问规则，往大一点说是定义了程序执行的次序。注意，为了获得较好的执行性能，Java内存模型并没有限制执行引擎使用处理器的寄存器或者高速缓存来提升指令执行速度，也没有限制编译器对指令进行重排序。也就是说，在java内存模型中，也会存在缓存一致性问题 and 指令重排序的问题。

Java内存模型规定所有的变量都是存在主存当中（类似于前面说的物理内存），每个线程都有自己的工作内存（类似于前面的高速缓存）。线程对变量的所有操作都必须在工作内存中进行，而不能直接对主存进行操作。并且每个线程不能访问其他线程的工作内存。

举个简单的例子：在java中，执行下面这个语句：

```
1 | i = 10;
```

执行线程必须先在自己的工作线程中对变量i所在的缓存行进行赋值操作，然后再写入主存当中。而不是直接将数值10写入主存当中。

那么Java语言 本身对 原子性、可见性以及有序性提供了哪些保证呢？

1.原子性

在Java中，对基本数据类型的变量的读取和赋值操作是原子性操作，即这些操作是不可被中断的，要么执行，要么不执行。

上面一句话虽然看起来简单，但是理解起来并不是那么容易。看下面一个例子i：

请分析以下哪些操作是原子性操作：

```
1  x = 10;           //语句1
2  y = x;            //语句2
3  x++;              //语句3
4  x = x + 1;        //语句4
```

乍一看，有些朋友可能会说上面的4个语句中的操作都是原子性操作。其实只有语句1是原子性操作，其他三个语句都不是原子性操作。

语句1是直接将数值10赋值给x，也就是说线程执行这个语句的会直接将数值10写入到工作内存中。

语句2实际上包含2个操作，它先要去读取x的值，再将x的值写入工作内存，虽然读取x的值以及将x的值写入工作内存 这2个操作都是原子性操作，但是合起来就不是原子性操作了。

同样的，x++和 x = x + 1包括3个操作：读取x的值，进行加1操作，写入新的值。

所以上面4个语句只有语句1的操作具备原子性。

也就是说，只有简单的读取、赋值（而且必须是将数字赋值给某个变量，变量之间的相互赋值不是原子操作）才是原子操作。

不过这里有一点需要注意：在32位平台下，对64位数据的读取和赋值是需要通过两个操作来完成的，不能保证其原子性。但是好像在最新的JDK中，JVM已经保证对64位数据的读取和赋值也是原子性操作了。

从上面可以看出，Java内存模型只保证了基本读取和赋值是原子性操作，如果要实现更大范围操作的原子性，可以通过synchronized和Lock来实现。由于synchronized和Lock能够保证任一时刻只有一个线程执行该代码块，那么自然就不存在原子性问题了，从而保证了原子性。

2.可见性

对于可见性，Java提供了volatile关键字来保证可见性。

当一个共享变量被volatile修饰时，它会保证修改的值会立即被更新到主存，当有其他线程需要读取时，它会去内存中读取新值。

而普通的共享变量不能保证可见性，因为普通共享变量被修改之后，什么时候被写入主存是不确定的，当其他线程去读取时，此时内存中可能还是原来的旧值，因此无法保证可见性。

另外，通过synchronized和Lock也能够保证可见性，synchronized和Lock能保证同一时刻只有一个线程获取锁然后执行同步代码，并且在释放锁之前会将对变量的修改刷新到主存当中。因此可以保证可见性。

3.有序性

在Java内存模型中，允许编译器和处理器对指令进行重排序，但是重排序过程不会影响到单线程程序的执行，却会影响到多线程并发执行的正确性。

在Java里面，可以通过volatile关键字来保证一定的“有序性”（具体原理在下一节讲述）。另外可以通过synchronized和Lock来保证有序性，很显然，synchronized和Lock保证每个时刻是有一个线

赞
助

程执行同步代码，相当于是让线程顺序执行同步代码，自然就保证了有序性。

另外，Java内存模型具备一些先天的“有序性”，即不需要通过任何手段就能够得到保证的有序性，这个通常也称为 happens-before 原则。如果两个操作的执行次序无法从happens-before原则推导出来，那么它们就不能保证它们的有序性，虚拟机可以随意地对它们进行重排序。

下面就来具体介绍下happens-before原则（先行发生原则）：

- 程序次序规则：一个线程内，按照代码顺序，书写在前面的操作先行发生于书写在后面的操作
- 锁定规则：一个unLock操作先行发生于后面对同一个锁锁lock操作
- volatile变量规则：对一个变量的写操作先行发生于后面对这个变量的读操作
- 传递规则：如果操作A先行发生于操作B，而操作B又先行发生于操作C，则可以得出操作A先行发生于操作C
- 线程启动规则：Thread对象的start()方法先行发生于此线程的每一个动作
- 线程中断规则：对线程interrupt()方法的调用先行发生于被中断线程的代码检测到中断事件的发生
- 线程终结规则：线程中所有的操作都先行发生于线程的终止检测，我们可以通过Thread.join()方法结束、Thread.isAlive()的返回值手段检测到线程已经终止执行
- 对象终结规则：一个对象的初始化完成先行发生于他的finalize()方法的开始

这8条原则摘自《深入理解Java虚拟机》。

这8条规则中，前4条规则是比较重要的，后4条规则都是显而易见的。

下面我们来解释一下前4条规则：

对于程序次序规则来说，我的理解就是一段程序代码的执行在单个线程中看起来是有序的。注意，虽然这条规则中提到“书写在前面的操作先行发生于书写在后面的操作”，这个应该是程序看起来执行的顺序是按照代码顺序执行的，因为虚拟机可能会对程序代码进行指令重排序。虽然进行重排序，但是最终执行的结果是与程序顺序执行的结果一致的，它只会对不存在数据依赖性的指令进行重排序。因此，在单个线程中，程序执行看起来是有序执行的，这一点要注意理解。事实上，这个规则是用来保证程序在单线程中执行结果的正确性，但无法保证程序在多线程中执行的正确性。

第二条规则也比较容易理解，也就是说无论在单线程中还是多线程中，同一个锁如果出于被锁定的状态，那么必须先对锁进行了释放操作，后面才能继续进行lock操作。

第三条规则是一条比较重要的规则，也是后文将要重点讲述的内容。直观地解释就是，如果一个线程先去写一个变量，然后一个线程去进行读取，那么写入操作肯定会先行发生于读操作。

第四条规则实际上就是体现happens-before原则具备传递性。

四.深入剖析volatile关键字

在前面讲述了很多东西，其实都是为讲述volatile关键字作铺垫，那么接下来我们就进入主题。

1.volatile关键字的两层语义

一旦一个共享变量（类的成员变量、类的静态成员变量）被volatile修饰之后，那么就具备了两层语义：

1) 保证了不同线程对这个变量进行操作时的可见性，即一个线程修改了某个变量的值，这新值对其他线程来说是立即可见的。

2) 禁止进行指令重排序。

先看一段代码，假如线程1先执行，线程2后执行：

```
1 //线程1
2 boolean stop = false;
3 while(!stop){
```

```

4      doSomething();
5  }
6
7  //线程2
8  stop = true;

```

这段代码是很典型的一段代码，很多人在中断线程时可能都会采用这种标记办法。但是事实上，这段代码会完全运行正确么？即一定会将线程中断么？不一定，也许在大多数时候，这个代码能够把线程中断，但是也有可能无法中断线程（虽然这个可能性很小，但是只要一旦发生这种情况就会造成死循环了）。

下面解释一下这段代码为何有可能导致无法中断线程。在前面已经解释过，每个线程在运行过程中都有自己的工作内存，那么线程1在运行的时候，会将stop变量的值拷贝一份放在自己的工作内存当中。

那么当线程2更改了stop变量的值之后，但是还没来得及写入主存当中，线程2转去做其他事情了，那么线程1由于不知道线程2对stop变量的更改，因此还会一直循环下去。

但是用volatile修饰之后就变得不一样了：

第一：使用volatile关键字会强制将修改的值立即写入主存；

第二：使用volatile关键字的话，当线程2进行修改时，会导致线程1的工作内存中缓存变量stop的缓存行无效（反映到硬件层的话，就是CPU的L1或者L2缓存中对应的缓存行无效）；

第三：由于线程1的工作内存中缓存变量stop的缓存行无效，所以线程1再次读取变量stop的值时会去主存读取。

那么在线程2修改stop值时（当然这里包括2个操作，修改线程2工作内存中的值，然后将修改后的值写入内存），会使得线程1的工作内存中缓存变量stop的缓存行无效，然后线程1读取时，发现自己的缓存行无效，它会等待缓存行对应的主存地址被更新之后，然后去对应的主存读取最新的值。

那么线程1读取到的就是最新的正确的值。

2.volatile保证原子性吗？

从上面知道volatile关键字保证了操作的可见性，但是volatile能保证对变量的操作是原子性吗？

下面看一个例子：

```

1  public class Test {
2      public volatile int inc = 0;
3
4      public void increase() {
5          inc++;
6      }
7
8      public static void main(String[] args) {
9          final Test test = new Test();
10         for(int i=0;i<10;i++){
11             new Thread(){
12                 public void run() {
13                     for(int j=0;j<1000;j++)
14                         test.increase();
15                 };
16             }.start();
17         }
18
19         while(Thread.activeCount()>1) //保证前面的线程都执行完
20             Thread.yield();
21         System.out.println(test.inc);
22     }
23 }

```

大家想一下这段程序的输出结果是多少？也许有些朋友认为是10000。但是事实上运行它会发现每次运行结果都不一致，都是一个小于10000的数字。

可能有的朋友就会有疑问，不对啊，上面是对变量inc进行自增操作，由于volatile保证了可见性，那么在每个线程中对inc自增完之后，在其他线程中都能看到修改后的值啊，所以有10个线程分别进行了1000次操作，那么最终inc的值应该是 $1000 \times 10 = 10000$ 。

这里面就有一个误区了，volatile关键字能保证可见性没有错，但是上面的程序错在没能保证原子性。可见性只能保证每次读取的是最新的值，但是volatile没办法保证对变量的操作的原子性。

在前面已经提到过，自增操作是不具备原子性的，它包括读取变量的原始值、进行加1操作、写入工作内存。那么就是说自增操作的三个子操作可能会分割开执行，就有可能导致下面这种情况出现：

假如某个时刻变量inc的值为10，

线程1对变量进行自增操作，线程1先读取了变量inc的原始值，然后线程1被阻塞了；

然后线程2对变量进行自增操作，线程2也去读取变量inc的原始值，由于线程1只是对变量inc进行读取操作，而没有对变量进行修改操作，所以不会导致线程2的工作内存中缓存变量inc的缓存行无效，所以线程2会直接去主存读取inc的值，发现inc的值是10，然后进行加1操作，并把11写入工作内存，最后写入主存。

然后线程1接着进行加1操作，由于已经读取了inc的值，注意此时在线程1的工作内存中inc的值仍然为10，所以线程1对inc进行加1操作后inc的值为11，然后将11写入工作内存，最后写入主存。


那么两个线程分别进行了一次自增操作后，inc只增加了1。

解释到这里，可能有朋友会有疑问，不对啊，前面不是保证一个变量在修改volatile变量时，会让缓存行无效吗？然后其他线程去读就会读到新的值，对，这个没错。这个就是上面的happens-before规则中的volatile变量规则，但是要注意，线程1对变量进行读取操作之后，被阻塞了的话，并没有对inc值进行修改。然后虽然volatile能保证线程2对变量inc的值读取是从内存中读取的，但是线程1没有进行修改，所以线程2根本就不会看到修改的值。


根源就在这里，自增操作不是原子性操作，而且volatile也无法保证对变量的任何操作都是原子性的。

把上面的代码改成以下任何一种都可以达到效果：


采用synchronized：

 [View Code](#)

采用Lock：

 [View Code](#)

采用AtomicInteger：

 [View Code](#)

在java 1.5的java.util.concurrent.atomic包下提供了一些原子操作类，即对基本数据类型的自增（加1操作），自减（减1操作）、以及加法操作（加一个数），减法操作（减一个数）进行了封装，保证这些操作是原子性操作。atomic是利用CAS来实现原子性操作的（Compare And Swap），CAS实际上是利用处理器提供的CMPXCHG指令实现的，而处理器执行CMPXCHG指令是一个原子性操作。

3.volatile能保证有序性吗？

在前面提到volatile关键字能禁止指令重排序，所以volatile能在一定程度上保证有序性。

volatile关键字禁止指令重排序有两层意思：

赞
助

- 1) 当程序执行到volatile变量的读操作或者写操作时，在其前面的操作的更改肯定全部已经进行，且结果已经对后面的操作可见；在其后面的操作肯定还没有进行；
- 2) 在进行指令优化时，不能将在对volatile变量访问的语句放在其后面执行，也不能把volatile变量后面的语句放到其前面执行。

可能上面说的比较绕，举个简单的例子：

```
1 //x、y为非volatile变量
2 //flag为volatile变量
3
4 x = 2;          //语句1
5 y = 0;          //语句2
6 flag = true;    //语句3
7 x = 4;          //语句4
8 y = -1;         //语句5
```

由于flag变量为volatile变量，那么在进行指令重排序的过程的时候，不会将语句3放到语句1、语句2前面，也不会讲语句3放到语句4、语句5后面。但是要注意语句1和语句2的顺序、语句4和语句5的顺序是不作任何保证的。

并且volatile关键字能保证，执行到语句3时，语句1和语句2必定是执行完毕了的，且语句1和语句2的执行结果对语句3、语句4、语句5是可见的。

那么我们回到前面举的一个例子：

```
1 //线程1:
2 context = loadContext(); //语句1
3 initied = true;          //语句2
4
5 //线程2:
6 while(!initied ){
7     sleep()
8 }
9 doSomethingwithconfig(context);
```

前面举这个例子的时候，提到有可能语句2会在语句1之前执行，那么久可能导致context还没被初始化，而线程2中就使用未初始化的context去进行操作，导致程序出错。

这里如果用volatile关键字对initied变量进行修饰，就不会出现这种问题了，因为当执行到语句2时，必定能保证context已经初始化完毕。

4.volatile的原理和实现机制

前面讲述了源于volatile关键字的一些使用，下面我们来探讨一下volatile到底如何保证可见性和禁止指令重排序的。

下面这段话摘自《深入理解Java虚拟机》：

“观察加入volatile关键字和没有加入volatile关键字时所生成的汇编代码发现，加入volatile关键字时，会多出一个lock前缀指令”

lock前缀指令实际上相当于一个内存屏障（也成内存栅栏），内存屏障会提供3个功能：

- 1) 它确保指令重排序时不会把其后面的指令排到内存屏障之前的位置，也不会把前面的指令排到内存屏障的后面；即在执行到内存屏障这句指令时，在它前面的操作已经全部完成；
- 2) 它会强制将对缓存的修改操作立即写入主存；
- 3) 如果是写操作，它会导致其他CPU中对应的缓存行无效。

五.使用volatile关键字的场景

synchronized关键字是防止多个线程同时执行一段代码，那么就会很影响程序执行效率，而volatile关键字在某些情况下性能要优于synchronized，但是要注意volatile关键字是无法替代

赞
助

synchronized关键字的, 因为volatile关键字无法保证操作的原子性。通常来说, 使用volatile必须具备以下2个条件:

- 1) 对变量的写操作不依赖于当前值
- 2) 该变量没有包含在具有其他变量的不变式中

实际上, 这些条件表明, 可以被写入 volatile 变量的这些有效值独立于任何程序的状态, 包括变量的当前状态。

事实上, 我的理解就是上面的2个条件需要保证操作是原子性操作, 才能保证使用volatile关键字的程序在并发时能够正确执行。

下面列举几个Java中使用volatile的几个场景。

1.状态标记量

```
1 volatile boolean flag = false;
2
3 while(!flag){
4     doSomething();
5 }
6
7 public void setFlag() {
8     flag = true;
9 }
```

```
1 volatile boolean initd = false;
2 //线程1:
3 context = loadContext();
4 initd = true;
5
6 //线程2:
7 while(!initd ){
8     sleep()
9 }
10 doSomethingwithconfig(context);
```

2.double check

```
1 class Singleton{
2     private volatile static Singleton instance = null;
3
4     private Singleton() {
5
6     }
7
8     public static Singleton getInstance() {
9         if(instance==null) {
10             synchronized (Singleton.class) {
11                 if(instance==null)
12                     instance = new Singleton();
13             }
14         }
15         return instance;
16     }
17 }
```

至于为何需要这么写请参考:

《Java 中的双重检查 (Double-Check) 》

<http://blog.csdn.net/dl88250/article/details/5439024>

赞助

和<http://www.iteye.com/topic/652440>

参考资料:

《Java编程思想》

《深入理解Java虚拟机》

<http://jiangzhengjun.iteye.com/blog/652532>

http://blog.sina.com.cn/s/blog_7bee8dd50101fu8n.html

<http://ifeve.com/volatile/>

<http://blog.csdn.net/ccit0519/article/details/11241403>

http://blog.csdn.net/ns_code/article/details/17101369

<http://www.cnblogs.com/kevinwu/archive/2012/05/02/2479464.html>

<http://www.cppblog.com/elva/archive/2011/01/21/139019.html>

<http://ifeve.com/volatile-array-visibility/>

<http://www.bdq.cn/news/201312/12579.shtml>

<http://exploer.blog.51cto.com/7123589/1193399>

<http://www.cnblogs.com/Mainz/p/3556430.html>

作者: [Matrix海子](#)

出处: <http://www.cnblogs.com/dolphin0520/>

本博客中未标明转载的文章归作者[Matrix海子](#)和博客园共有, 欢迎转载, 但未经作者同意必须保留此段声明, 且在文章页面明显位置给出原文连接, 否则保留追究法律责任的权利。

分类: [Java并发编程-入门篇](#)

好文要顶

关注我

收藏该文

[Matrix海子](#)
[关注 - 7](#)
[粉丝 - 8161](#)
[+加关注](#)

3791

« 上一篇: [Java并发编程: Lock](#)
» 下一篇: [Java并发编程: 深入剖析ThreadLocal](#)

posted @ 2014-08-22 22:07 Matrix海子 阅读(492324) 评论(235) 编辑 收藏

评论列表

#201楼 2019-06-11 16:40 qweenhool

看完了一部分的评论, 感觉越来越懵逼了, inc那个例子到底是怎么回事?

支持(5) 反对(0)

#202楼 2019-06-12 12:13 naffan

可以再进行一次结构化调整, 感觉后来写的越来越乱了

支持(0) 反对(0)

赞助

#203楼 2019-06-14 16:01 那些年的代码

确实明白点了, 之前接触的少, 还得多看几遍巩固

支持(0) 反对(0)

#204楼 2019-06-22 11:34 李大鹏1991

写的很好, 给赞

支持(3) 反对(0)

#205楼 2019-07-09 08:07 陈俊成

@_ qweenhool

引用

看完了一部分的评论, 感觉越来越懵逼了, inc那个例子到底是怎么回事?

inc那个例子, 作者是想说明, inc++本来是原子性的, 用volatile也不能保证原子性, 所以每次运行结果都有可能小于1000

支持(3) 反对(0)

#206楼 2019-07-16 11:25 关外夜店

@_ qweenhool

引用

看完了一部分的评论, 感觉越来越懵逼了, inc那个例子到底是怎么回事?

我是这样理解的, volatile修饰的变量, 多个线程不去同时操作不是原子性的操作, 类似inr++就不是原子性, 但是可以去操作是原子性的操作的比如简单的赋值。赋值操作是一个线程 其他的线程去读取这个变量的时候 读取到的是修改过的值。这时候修改线程改的值 对其他线程是可见的。

支持(0) 反对(0)

#207楼 2019-07-17 15:09 zbmcz

感谢作者, 我花了3天时间才读完并消化里面的内容, 受益匪浅

支持(1) 反对(0)

#208楼 2019-07-29 13:57 dnliu

写的很好, 点赞

支持(1) 反对(0)

#209楼 2019-07-30 14:47 云起龙沙暗

@_ 深夜不说话

这里加volatile不是为了解决可见性, 而是解决对象实例化过程中可能发生重排序问题

支持(0) 反对(0)

#210楼 2019-08-20 21:27 刘知昊

@_ qweenhool

我来说一下吧, **作者对那个例子的解释就是错误的**, 线程2写完内存, 到线程1去加1之前, 线程1确实是会发现缓存行失效然后重新读取, **照作者说的这个顺序并不会造成不一致的问题**。真正造成不一致的原因, 也就是对这个程序执行现象正确的解释见**评论区15楼寒铁衣**的回答。

我来总结一下吧, 真正造成不一致的是这样的顺序: 线程1读取Inc的值, 然后对Inc+1, 然后线程1被阻塞。线程2读取Inc的值, 对Inc加1, 然后写回给Inc, 此时volatile起作用了! 这个值立即被写回主内存, 并导致其它线程的Inc变量的缓存行失效, 这时候线程1的Inc变量缓存失效了! 接下来线程1继续执行, 将Inc的值写回。

也就是说, **线程1中断的地方是在加1完成写回之前! 不是作者写的读取完成加1之前!** 这样才会出现只加了一次的的问题!

支持(7) 反对(1)

#211楼 2019-08-22 20:36 周李焕

赞
助

我就上面的说法有个疑问, 线程1读取了值, 堵塞了, 然后线程2读取值自增, 为什么线程2修改了值, 线程1看不到修改, 还继续自增, 不是说只要这个变量被任何线程修改都会使其它线程缓存失效吗, 为什么线程1读取了值之后就不失效了

支持(1) 反对(0)

#212楼 2019-08-30 16:57 lewis_gary

原子性的解释有点乱啊

支持(0) 反对(0)

#213楼 2019-09-25 17:20 hadeser

过段时间来看看评论, 果然又多一层体会!

支持(0) 反对(0)

#214楼 2019-11-17 20:57 胖兔向前冲

@ 刘知昊

说的更具体一点, 线程1执行inc+1之后是以temp=inc+1的形式存在, 这个temp是个非volatile变量, 不会受后面线程2写inc导致缓存行失效的影响

支持(0) 反对(0)

赞

#215楼 2019-11-26 15:00 mgwanli

这篇文章我看了三遍, 第三遍才彻底明白volatile

支持(0) 反对(0)

#216楼 2019-12-20 16:38 Silent_Walker

@ NeverAlone

问题是不是这个, 问题是文章前后矛盾了

支持(0) 反对(0)

#217楼 2019-12-20 16:45 Silent_Walker

@ 刘知昊

但是, 你有没有向想过, 如果+1后还能阻塞一段时间后再写回主存, 那可见性就是个笑话。要保持可见性, +1后必须立马写回主存, 期间不能切换线程。

支持(0) 反对(0)

#218楼 2019-12-23 21:09 vincent_glk

然后线程1接着进行加1操作, 由于已经读取了inc的值, 注意此时在线程1的工作内存中inc的值仍然为10, 所以线程1对inc进行加1操作后inc的值为11, 然后将11写入工作内存, 最后写入主存。

你好, 作者

这时候线程1在进行加1操作的话, 它的缓存变量的值应该失效了吧?

支持(0) 反对(0)

#219楼 2020-01-15 14:39 原来是超人啊

@ Silent_Walker

我理解的跟你差不多, 对于变量的操作过程分三步, 读-->运算-->写
也就是说volatile保证了运算操作与写回主存这两个操作是原子的, 这样才能保证所谓的可见性, 也就是说对volatile修饰的变量进行修改线程只会在读和运算之间切换, 运算和写之间是不会切换的, 不知道这样理解对不对, 也没办法证明啊!

支持(0) 反对(0)

#220楼 2020-01-15 15:17 Silent_Walker

@ 原来是超人啊

不对, 读, 运算和写都是独立的操作, 只有读和写是原子操作, 并且是互斥的。后来我看了一些文章, 现在再来看看作者写。其实本文的问题在于把硬件层面和语言层面混合讲解, 但是硬件层面又没有讲清楚, 并且不同的cpu有可能有不同的实现。问题就出在缓存行失效, 是怎么失效的。两个线程对volatile变量的自增在硬件层面是如何操作的。这个我建议保留疑问, 可以看看其他有关底层

书。volatile的语言层面的分析我建议看这篇文章 <https://www.infoq.cn/article/java-memory-model-4/> , 最好看看这一个系列的都可以看看。

支持(0) 反对(0)

#221楼 2020-01-15 16:04 原来是超人啊

@ Silent_Walker

按你推荐的这篇文章来说, 如果线程1修改完了, 但是线程2已经读完了, 那么线程2的值也在本次操作中也不再改变, 只有等下次读时才从内存里读, 跟前面部分博友的理解是一样的

支持(0) 反对(0)

#222楼 2020-01-15 16:06 Silent_Walker

@ 原来是超人啊
是不会改变的,

支持(0) 反对(0)

#223楼 2020-02-20 17:53 joininjoy

好详细啊, 谢谢大神分享。

支持(0) 反对(0)

赞
助

#224楼 2020-02-26 23:12 烟雨星空

@ 刘知昊

既然线程1进行加1操作了, 那么线程2应该去读取最新值啊, 也就是说必须线程1把加1后的操作写回主存, 才会可见? 这个时候, 线程2才会读取到最新的值2? 这样理解有毛病吗

支持(0) 反对(0)

#225楼 2020-02-26 23:18 烟雨星空

@ 刘知昊

文中提到volatile会“强制将对缓存的修改操作立即写入主存”。如果按照你的理论, 线程1进行了自加1操作才被阻塞, 那就会立即写入主存, 线程2应该读取到最新值2啊。。我越来越迷惑了。怎么解释原子性都解释不通

支持(0) 反对(0)

#226楼 2020-02-26 23:39 刘知昊

@ 烟雨星空

[读], [加1], [写回给变量], 是三个操作, 只有在**[写回给变量]**完成之后才会发生**[其它线程的该变量的缓存行失效]**, 最终导致其它线程如果要进行**[加1]**操作就需要重新从内存里读取, 而不是取用刚刚读到缓存里的值。

补充: 这里**[写回给变量]**指的是227楼伪代码的**step3+step4合在一起**, 也就是写回到主存里才算。

支持(0) 反对(0)

#227楼 2020-02-26 23:42 刘知昊

@ 烟雨星空

Inc = Inc + 1 请理解成

```
1  step1: 从内存中读取Inc的值给存到缓存中, 为c
2  step2: 临时值 t = c + 1 (在寄存器里完成加法操作, 不要认为是在缓存里把自己加1)
3  step3: 写回缓存的操作 c = t
4  *step4: 写回主存的操作 Inc = c
```

当Inc用volatile修饰时, step3发生会导致step4发生, 即“强制将对缓存的修改操作立即写入主存”。其它时候, step4不会立即发生(否则缓存也就起不到加速作用了)。

支持(0) 反对(0)

#228楼 2020-02-26 23:42 烟雨星空

@ 刘知昊

嗯, 你说这个我能理解。我意思是, 这句话“强制将对缓存的修改操作立即写入主存”发生在什么时候? 就是说, 什么时候会强制立即写入主存呢

支持(0) 反对(0)

#229楼 2020-02-26 23:47 烟雨星空

@ 刘知昊

非常感谢, 这个三步曲我明白了, “强制将对缓存的修改操作立即写入主存”, 这句话我非常迷惑, 对缓存的修改不就是第二步吗, 然后它又说这样会立即写入主存也就是第三步。好像意思是第二步和第三步是一起发生的。。。按你的意思, 也能说得通, 但是, 这是不是也就说明上边这句话会有歧义。。

支持(0) 反对(0)

#230楼 2020-02-26 23:48 刘知昊

@ 烟雨星空

你好, 我把step3详细展开了, 修改了227楼, 你看了就明白了。

支持(0) 反对(0)

#231楼 2020-03-10 00:22 "类库☹"

妙啊 妙

支持(0) 反对(0)

#232楼 2020-03-20 21:42 花生和馒头

volatile保证了assign、store、write的连续, 会把工作内存的值直接刷到了主存, 但是在use阶段, jvm把工作内存交给执行引擎了, 这里执行引擎压入操作栈以后, 你切换线程了, 修改了i以后, 刷到主内存了, 通知过来了, 执行引擎也不会和再去重读了。这是我的理解。

支持(0) 反对(0)

#233楼 2020-03-23 21:40 流水晓风

```
/**
 *保证前面的线程都执行完。----文中下面这个1, 要改成2。main线程至少有两个活动线程。
 *在idea中如果是1, 会死循环。
 */
while(Thread.activeCount() > 1 ) {
    Thread.yield();
}
```

支持(0) 反对(0)

#234楼 2020-03-26 22:13 Precious7

对inc例子的理解, 原文中表述可能有不当的地方, 对变量x的inc操作可以分解为以下步骤:

- 从主内存取出x到线程工作内存;
- 临时变量 t = x+1;
- x = t; 因为x更改, 所以立刻将 x 写到主存, 使其他线程的工作内存失效。(这一步是不可拆分的)

x初始值为10, 线程1执行完第二步, 尚未执行第三步时, 切换至线程2, 此时x依然为10, 线程2执行完三步之后, x变更为11, 因为x更改, 线程1的工作内存立刻失效, 此时切换到线程1, 线程1只会继续执行第三步, t之前已经被计算出来是11, 直接赋值给x, 因此x最终还是11。

之所以会这样, 就是因为这不是一步, 而是三步, 也就是非原子的。volatile保证了x的更改立刻被线程知晓, 但无法回退“因为x的改变, 需要重新执行”的步骤。

支持(1) 反对(0)

#235楼 2020-04-06 10:17 海纳分享

牛逼! 写mark, 再研究!

支持(0) 反对(0)

赞
助

注册用户登录后才能发表评论，请 [登录](#) 或 [注册](#)，[访问](#) [网站首页](#)。

Copyright © 2020 Matrix海子
Powered by .NET Core on Kubernetes

