

深入理解Java之线程池

分享到：

原文出处：[海子](#)

在前面的文章中，我们使用线程的时候就去创建一个线程，这样实现起来非常简便，但是就会有一个问题：

如果并发的线程数量很多，并且每个线程都是执行一个时间很短的任务就结束了，这样频繁创建线程就会大大降低系统的效率，因为频繁创建线程和销毁线程需要时间。

那么有没有一种办法使得线程可以复用，就是执行完一个任务，并不被销毁，而是可以继续执行其他的任务？

在Java中可以通过线程池来达到这样的效果。今天我们就来详细讲解一下Java的线程池，首先我们从最核心的ThreadPoolExecutor类中的方法讲起，然后再讲述它的实现原理，接着给出了它的使用示例，最后讨论了一下如何合理配置线程池的大小。

以下是本文的目录大纲：

一.Java中的ThreadPoolExecutor类

二.深入剖析线程池实现原理

三.使用示例

四.如何合理配置线程池的大小

若有不正之处请多多谅解，并欢迎批评指正。

一.Java中的ThreadPoolExecutor类

java.util.concurrent.ThreadPoolExecutor类是线程池中最核心的一个类，因此如果要透彻地了解Java中的线程池，必须先了解这个类。下面我们来看一下ThreadPoolExecutor类的具体实现源码。

在ThreadPoolExecutor类中提供了四个构造方法：

```
1 public class ThreadPoolExecutor extends AbstractExecutorService {
2     ....
3     public ThreadPoolExecutor(int corePoolSize, int maximumPoolSize, long
4         BlockingQueue<Runnable> workQueue);
5
6     public ThreadPoolExecutor(int corePoolSize, int maximumPoolSize, long
7         BlockingQueue<Runnable> workQueue, ThreadFactory threadFa
8
9     public ThreadPoolExecutor(int corePoolSize, int maximumPoolSize, long
10        BlockingQueue<Runnable> workQueue, RejectedExecutionHandl
11
12     public ThreadPoolExecutor(int corePoolSize, int maximumPoolSize, long
```

```

13         BlockingQueue<Runnable> workQueue, ThreadFactory threadFactory,
14         ...
15     }

```

从上面的代码可以得知，ThreadPoolExecutor继承了AbstractExecutorService类，并提供了四个构造器，事实上，通过观察每个构造器的源码具体实现，发现前面三个构造器都是调用的第四个构造器进行的初始化工作。

下面解释一下构造器中各个参数的含义：

- **corePoolSize**：核心池的大小，这个参数跟后面讲述的线程池的实现原理有非常大的关系。在创建了线程池后，默认情况下，线程池中并没有任何线程，而是等待有任务到来才创建线程去执行任务，除非调用了

```
prestartAllCoreThreads()
```

或者

```
prestartCoreThread()
```

方法，从这2个方法的名字就可以看出，是预创建线程的意思，即在没有任务到来之前就创建**corePoolSize**个线程或者一个线程。默认情况下，在创建了线程池后，线程池中的线程数为0，当有任务来之后，就会创建一个线程去执行任务，当线程池中的线程数目达到**corePoolSize**后，就会把到达的任务放到缓存队列当中；
- **maximumPoolSize**：线程池最大线程数，这个参数也是一个非常重要的参数，它表示在线程池中最多能创建多少个线程；
- **keepAliveTime**：表示线程没有任务执行时最多保持多久时间会终止。默认情况下，只有当线程池中的线程数大于**corePoolSize**时，**keepAliveTime**才会起作用，直到线程池中的线程数不大于**corePoolSize**，即当线程池中的线程数大于**corePoolSize**时，如果一个线程空闲的时间达到**keepAliveTime**，则会终止，直到线程池中的线程数不超过**corePoolSize**。但是如果调用了

```
allowCoreThreadTimeOut(boolean)
```

方法，在线程池中的线程数不大于**corePoolSize**时，**keepAliveTime**参数也会起作用，直到线程池中的线程数为0；
- **unit**：参数**keepAliveTime**的时间单位，有7种取值，在TimeUnit类中有7种静态属性：

```

1     TimeUnit.DAYS;           //天
2     TimeUnit.HOURS;         //小时
3     TimeUnit.MINUTES;       //分钟
4     TimeUnit.SECONDS;        //秒
5     TimeUnit.MILLISECONDS;   //毫秒
6     TimeUnit.MICROSECONDS;   //微妙
7     TimeUnit.NANOSECONDS;    //纳秒

```

- **workQueue**：一个阻塞队列，用来存储等待执行的任务，这个参数的选择也很重要，会对线程池的运行过程产生重大影响，一般来说，这里的阻塞队列有以下几种选择：

```

1     ArrayBlockingQueue;
2     LinkedBlockingQueue;
3     SynchronousQueue;

```

ArrayBlockingQueue和PriorityBlockingQueue使用较少，一般使用LinkedBlockingQueue和Synchronous。线程池的排队策略与BlockingQueue有关。

- **threadFactory**：线程工厂，主要用来创建线程；
- **handler**：表示当拒绝处理任务时的策略，有以下四种取值：

```

1     ThreadPoolExecutor.AbortPolicy: 丢弃任务并抛出RejectedExecutionException
2     ThreadPoolExecutor.DiscardPolicy: 也是丢弃任务，但是不抛出异常。
3     ThreadPoolExecutor.DiscardOldestPolicy: 丢弃队列最前面的任务，然后重新尝试
4     ThreadPoolExecutor CallerRunsPolicy: 由调用线程处理该任务

```

具体参数的配置与线程池的关系将在下一节讲述。

从上面给出的ThreadPoolExecutor类的代码可以知道，ThreadPoolExecutor继承了AbstractExecutorService，我们来看一下AbstractExecutorService的实现：

```
1 public abstract class AbstractExecutorService implements ExecutorService {
2
3     protected <T> RunnableFuture<T> newTaskFor(Runnable runnable, T result) {
4     protected <T> RunnableFuture<T> newTaskFor(Callable<T> callable) {
5     public Future<?> submit(Runnable task) { };
6     public <T> Future<T> submit(Runnable task, T result) { };
7     public <T> Future<T> submit(Callable<T> task) { };
8     private <T> T doInvokeAny(Collection<? extends Callable<T>> tasks,
9                             boolean timed, long nanos)
10        throws InterruptedException, ExecutionException, TimeoutException {
11    };
12    public <T> T invokeAny(Collection<? extends Callable<T>> tasks)
13        throws InterruptedException, ExecutionException {
14    };
15    public <T> T invokeAny(Collection<? extends Callable<T>> tasks,
16                          long timeout, TimeUnit unit)
17        throws InterruptedException, ExecutionException, TimeoutException {
18    };
19    public <T> List<Future<T>> invokeAll(Collection<? extends Callable<T>> tasks)
20        throws InterruptedException {
21    };
22    public <T> List<Future<T>> invokeAll(Collection<? extends Callable<T>> tasks,
23                                       long timeout, TimeUnit unit)
24        throws InterruptedException {
25    };
26 }
```

AbstractExecutorService是一个抽象类，它实现了ExecutorService接口。

我们接着看ExecutorService接口的实现：

```
1 public interface ExecutorService extends Executor {
2
3     void shutdown();
4     boolean isShutdown();
5     boolean isTerminated();
6     boolean awaitTermination(long timeout, TimeUnit unit)
7         throws InterruptedException;
8     <T> Future<T> submit(Callable<T> task);
9     <T> Future<T> submit(Runnable task, T result);
10    Future<?> submit(Runnable task);
11    <T> List<Future<T>> invokeAll(Collection<? extends Callable<T>> tasks)
12        throws InterruptedException;
13    <T> List<Future<T>> invokeAll(Collection<? extends Callable<T>> tasks,
14                              long timeout, TimeUnit unit)
15        throws InterruptedException;
16
17    <T> T invokeAny(Collection<? extends Callable<T>> tasks)
18        throws InterruptedException, ExecutionException;
19    <T> T invokeAny(Collection<? extends Callable<T>> tasks,
20                  long timeout, TimeUnit unit)
21        throws InterruptedException, ExecutionException, TimeoutException;
22 }
```

而ExecutorService又是继承了Executor接口，我们看一下Executor接口的实现：

```
1 public interface Executor {
2     void execute(Runnable command);
3 }
```

到这里，大家应该明白了ThreadPoolExecutor、AbstractExecutorService、ExecutorService和Executor几个之间的关系了。

Executor是一个顶层接口，在它里面只声明了一个方法execute(Runnable)，返回值为void，参数为Runnable类型，从字面意思可以理解，就是用来执行传进去的任务的；

然后ExecutorService接口继承了Executor接口，并声明了一些方法：submit、invokeAll、invokeAny以及shutdown等；

抽象类AbstractExecutorService实现了ExecutorService接口，基本实现了ExecutorService中声明的所有方法；

然后ThreadPoolExecutor继承了类AbstractExecutorService。

在ThreadPoolExecutor类中有几个非常重要的方法：

```
1 execute()  
2 submit()  
3 shutdown()  
4 shutdownNow()
```

execute()方法实际上是Executor中声明的方法，在ThreadPoolExecutor进行了具体的实现，这个方法是ThreadPoolExecutor的核心方法，通过这个方法可以向线程池提交一个任务，交由线程池去执行。

submit()方法是在ExecutorService中声明的方法，在AbstractExecutorService就已经有了具体的实现，在ThreadPoolExecutor中并没有对其进行重写，这个方法也是用来向线程池提交任务的，但是它和execute()方法不同，它能够返回任务执行的结果，去看submit()方法的实现，会发现它实际上还是调用的execute()方法，只不过它利用了Future来获取任务执行结果（Future相关内容将在下一篇讲述）。

shutdown()和shutdownNow()是用来关闭线程池的。

还有很多其他的方法：

比如：getQueue()、getPoolSize()、getActiveCount()、getCompletedTaskCount()等获取与线程池相关属性的方法，有兴趣的朋友可以自行查阅API。

二.深入剖析线程池实现原理

在上一节我们从宏观上介绍了ThreadPoolExecutor，下面我们来深入解析一下线程池的具体实现原理，将从下面几个方面讲解：

1.线程池状态

2.任务的执行

3.线程池中的线程初始化

4.任务缓存队列及排队策略

5.任务拒绝策略

6.线程池的关闭

7.线程池容量的动态调整

1.线程池状态

在ThreadPoolExecutor中定义了一个volatile变量，另外定义了几个static final变量表示线程池的各个状态：

```
1 volatile int runState;
2 static final int RUNNING = 0;
3 static final int SHUTDOWN = 1;
4 static final int STOP = 2;
5 static final int TERMINATED = 3;
```

runState表示当前线程池的状态，它是一个volatile变量用来保证线程之间的可见性；

下面的几个static final变量表示runState可能的几个取值。

当创建线程池后，初始时，线程池处于RUNNING状态；

如果调用了shutdown()方法，则线程池处于SHUTDOWN状态，此时线程池不能够接受新的任务，它会等待所有任务执行完毕；

如果调用了shutdownNow()方法，则线程池处于STOP状态，此时线程池不能接受新的任务，并且会去尝试终止正在执行的任务；

当线程池处于SHUTDOWN或STOP状态，并且所有工作线程已经销毁，任务缓存队列已经清空或执行结束后，线程池被设置为TERMINATED状态。

2.任务的执行

在了解将任务提交给线程池到任务执行完毕整个过程之前，我们先来看一下ThreadPoolExecutor类中其他的一些比较重要成员变量：

```
1 private final BlockingQueue<Runnable> workQueue; //任务队列
2 private final ReentrantLock mainLock = new ReentrantLock(); //线程池锁
3 //、runState
4 private final HashSet<Worker> workers = new HashSet<Worker>(); //用
5
6 private volatile long keepAliveTime; //线程存活时间
7 private volatile boolean allowCoreThreadTimeOut; //是否允许为核心线程
8 private volatile int corePoolSize; //核心池的大小（即线程池中的线
9 private volatile int maximumPoolSize; //线程池最大能容忍的线程数
10
11 private volatile int poolSize; //线程池中当前的线程数
12
13 private volatile RejectedExecutionHandler handler; //任务拒绝策略
14
15 private volatile ThreadFactory threadFactory; //线程工厂，用来创建线
16
17 private int largestPoolSize; //用来记录线程池中曾经出现过的最大线程数
18
19 private long completedTaskCount; //用来记录已经执行完毕的任务个数
```

每个变量的作用都已经标明出来了，这里要重点解释一下corePoolSize、maximumPoolSize、largestPoolSize三个变量。

corePoolSize在很多地方被翻译成核心池大小，其实我的理解这个就是线程池的大小。举个简单的例子：

假如有一个工厂，工厂里面有10个工人，每个工人同时只能做一件任务。

因此只要当10个工人中有工人是空闲的，来了任务就分配给空闲的工人做；

当10个工人都有任务在做时，如果还来了任务，就把任务进行排队等待；

如果说新任务数目增长的速度远远大于工人做任务的速度，那么此时工厂主管可能会想补救措施，比如重新招4个临时工人进来；

然后将任务也分配给这4个临时工人做；

如果说着14个工人做任务的速度还是不够，此时工厂主管可能就要考虑不再接收新的任务或者抛弃前面的一些任务了。

当这14个工人当中有人空闲时，而新任务增长的速度又比较缓慢，工厂主管可能就考虑辞掉4个临时工了，只保持原来的10个工人，毕竟请额外的工人是要花钱的。

这个例子中的corePoolSize就是10，而maximumPoolSize就是14（10+4）。

也就是说corePoolSize就是线程池大小，maximumPoolSize在我看来是线程池的一种补救措施，即任务量突然过大时的一种补救措施。

不过为了方便理解，在本文后面还是将corePoolSize翻译成核心池大小。

largestPoolSize只是一个用来起记录作用的变量，用来记录线程池中曾经有过的最大线程数目，跟线程池的容量没有任何关系。

下面我们进入正题，看一下任务从提交到最终执行完毕经历了哪些过程。

在ThreadPoolExecutor类中，最核心的任务提交方法是execute()方法，虽然通过submit也可以提交任务，但是实际上submit方法里面最终调用的还是execute()方法，所以我们只需要研究execute()方法的实现原理即可：

```
1 public void execute(Runnable command) {
2     if (command == null)
3         throw new NullPointerException();
4     if (poolSize >= corePoolSize || !addIfUnderCorePoolSize(command))
5         if (runState == RUNNING && workQueue.offer(command)) {
6             if (runState != RUNNING || poolSize == 0)
7                 ensureQueuedTaskHandled(command);
8         }
9     else if (!addIfUnderMaximumPoolSize(command))
10        reject(command); // is shutdown or saturated
11 }
12 }
```

上面的代码可能看起来不是那么容易理解，下面我们一句一句解释：

首先，判断提交的任务command是否为null，若是null，则抛出空指针异常；

接着是这句，这句要好好理解一下：

```
1 if (poolSize >= corePoolSize || !addIfUnderCorePoolSize(command))
```

由于是或条件运算符，所以先计算前半部分的值，如果线程池中当前线程数不小于核心池大小，那么就会直接进入下面的if语句块了。

如果线程池中当前线程数小于核心池大小，则接着执行后半部分，也就是执行


```
1 | addIfUnderCorePoolSize(command)
```

如果执行完addIfUnderCorePoolSize这个方法返回false，则继续执行下面的if语句块，否则整个方法就直接执行完毕了。

如果执行完addIfUnderCorePoolSize这个方法返回false，然后接着判断：

```
1 | if (runState == RUNNING && workQueue.offer(command))
```

如果当前线程池处于RUNNING状态，则将任务放入任务缓存队列；如果当前线程池不处于RUNNING状态或者任务放入缓存队列失败，则执行：

```
1 | addIfUnderMaximumPoolSize(command)
```

如果执行addIfUnderMaximumPoolSize方法失败，则执行reject()方法进行任务拒绝处理。

回到前面：

```
1 | if (runState == RUNNING && workQueue.offer(command))
```

这句的执行，如果说当前线程池处于RUNNING状态且将任务放入任务缓存队列成功，则继续进行判断：

```
1 | if (runState != RUNNING || poolSize == 0)
```

这句判断是为了防止在将此任务添加进任务缓存队列的同时其他线程突然调用shutdown或者shutdownNow方法关闭了线程池的一种应急措施。如果是这样就执行：

```
1 | ensureQueuedTaskHandled(command)
```

进行应急处理，从名字可以看出是保证 添加到任务缓存队列中的任务得到处理。

我们接着看2个关键方法的实现：addIfUnderCorePoolSize和addIfUnderMaximumPoolSize：

```
1 | private boolean addIfUnderCorePoolSize(Runnable firstTask) {
2 |     Thread t = null;
3 |     final ReentrantLock mainLock = this.mainLock;
4 |     mainLock.lock();
5 |     try {
6 |         if (poolSize < corePoolSize && runState == RUNNING)
7 |             t = addThread(firstTask); //创建线程去执行firstTask
8 |         } finally {
9 |             mainLock.unlock();
10 |        }
11 |        if (t == null)
12 |            return false;
13 |        t.start();
14 |        return true;
15 |    }
```

这个是addIfUnderCorePoolSize方法的具体实现，从名字可以看出它的意图就是当低于核心池大小时执行的方法。下面看其具体实现，首先获取到锁，因为这地方涉及到线程池状态的变化，先通过if语句判断当前线程池中的线程数目是否小于核心池大小，有朋友也许会有疑问：前面在execute()方法中不是已经判断过了吗，只有线程池当前线程数目小于核心池大小才会执行addIfUnderCorePoolSize方法的，为何这地方还要继续判断？原因很简单，前面的判断过程中并没有加锁，因此可能在execute方法判断的时候poolSize小于corePoolSize，而判断完之后，在其他线程中又向线程池提交了任务，就可能导致poolSize不小于corePoolSize了，所以需要在这个地方继续判断。然后接着判

断线程池的状态是否为RUNNING，原因也很简单，因为有可能在其他线程中调用了shutdown或者shutdownNow方法。然后就是执行

```
1 | t = addThread(firstTask);
```

这个方法也非常关键，传进去的参数为提交的任务，返回值为Thread类型。然后接着在下面判断t是否为空，为空则表明创建线程失败（即poolSize>=corePoolSize或者runState不等于RUNNING），否则调用t.start()方法启动线程。

我们来看一下addThread方法的实现：

```
1 | private Thread addThread(Runnable firstTask) {
2 |     Worker w = new Worker(firstTask);
3 |     Thread t = threadFactory.newThread(w); //创建一个线程，执行任务
4 |     if (t != null) {
5 |         w.thread = t; //将创建的线程的引用赋值为w的成员变量
6 |         workers.add(w);
7 |         int nt = ++poolSize; //当前线程数加1
8 |         if (nt > largestPoolSize)
9 |             largestPoolSize = nt;
10 |    }
11 |    return t;
12 | }
```

在addThread方法中，首先用提交的任务创建了一个Worker对象，然后调用线程工厂threadFactory创建了一个新的线程t，然后将线程t的引用赋值给了Worker对象的成员变量thread，接着通过workers.add(w)将Worker对象添加到工作集中。

下面我们看一下Worker类的实现：

```
1 | private final class Worker implements Runnable {
2 |     private final ReentrantLock runLock = new ReentrantLock();
3 |     private Runnable firstTask;
4 |     volatile long completedTasks;
5 |     Thread thread;
6 |     Worker(Runnable firstTask) {
7 |         this.firstTask = firstTask;
8 |     }
9 |     boolean isActive() {
10 |         return runLock.isLocked();
11 |     }
12 |     void interruptIfIdle() {
13 |         final ReentrantLock runLock = this.runLock;
14 |         if (runLock.tryLock()) {
15 |             try {
16 |                 if (thread != Thread.currentThread())
17 |                     thread.interrupt();
18 |             } finally {
19 |                 runLock.unlock();
20 |             }
21 |         }
22 |     }
23 |     void interruptNow() {
24 |         thread.interrupt();
25 |     }
26 |
27 |     private void runTask(Runnable task) {
28 |         final ReentrantLock runLock = this.runLock;
29 |         runLock.lock();
30 |         try {
31 |             if (runState < STOP &&
32 |                 Thread.interrupted() &&
33 |                 runState >= STOP)
34 |                 boolean ran = false;
35 |             beforeExecute(thread, task); //beforeExecute方法是Thread
36 |             //自己需要重载这个方法和后面的afterExecute方法来进行一些统计
37 |             try {
38 |                 task.run();
```



```

39         ran = true;
40         afterExecute(task, null);
41         ++completedTasks;
42     } catch (RuntimeException ex) {
43         if (!ran)
44             afterExecute(task, ex);
45         throw ex;
46     }
47 } finally {
48     runLock.unlock();
49 }
50 }
51
52 public void run() {
53     try {
54         Runnable task = firstTask;
55         firstTask = null;
56         while (task != null || (task = getTask()) != null) {
57             runTask(task);
58             task = null;
59         }
60     } finally {
61         workerDone(this); //当任务队列中没有任务时，进行清理工作
62     }
63 }
64 }

```

它实际上实现了Runnable接口，因此上面的Thread t = threadFactory.newThread(w);效果跟下面这句的效果基本一样：

```

1 Thread t = new Thread(w);

```

相当于传进去了一个Runnable任务，在线程t中执行这个Runnable。

既然Worker实现了Runnable接口，那么自然最核心的方法便是run()方法了：

```

1 public void run() {
2     try {
3         Runnable task = firstTask;
4         firstTask = null;
5         while (task != null || (task = getTask()) != null) {
6             runTask(task);
7             task = null;
8         }
9     } finally {
10         workerDone(this);
11     }
12 }

```

从run方法的实现可以看出，它首先执行的是通过构造器传进来的任务firstTask，在调用runTask()执行完firstTask之后，在while循环里面不断通过getTask()去取新的任务来执行，那么去哪里取呢？自然是从任务缓存队列里面去取，getTask是ThreadPoolExecutor类中的方法，并不是Worker类中的方法，下面是getTask方法的实现：

```

1 Runnable getTask() {
2     for (;;) {
3         try {
4             int state = runState;
5             if (state > SHUTDOWN)
6                 return null;
7             Runnable r;
8             if (state == SHUTDOWN) // Help drain queue
9                 r = workQueue.poll();
10            else if (poolSize > corePoolSize || allowCoreThreadTimeOut
11                //则通过poll取任务，若等待一定的时间取不到任务，则返回null
12                r = workQueue.poll(keepAliveTime, TimeUnit.NANOSECON
13            else
14                r = workQueue.take();
15            if (r != null)
16                return r;
17            if (workerCanExit()) { //如果没取到任务，即r为null，则判

```

```

18         if (runState >= SHUTDOWN) // Wake up others
19             interruptIdleWorkers(); //中断处于空闲状态的work
20         return null;
21     }
22     // Else retry
23 } catch (InterruptedException ie) {
24     // On interruption, re-check runState
25 }
26 }
27 }

```

在getTask中，先判断当前线程池状态，如果runState大于SHUTDOWN（即为STOP或者TERMINATED），则直接返回null。

如果runState为SHUTDOWN或者RUNNING，则从任务缓存队列取任务。

如果当前线程池的线程数大于核心池大小corePoolSize或者允许为核心池中的线程设置空闲存活时间，则调用poll(time,timeUnit)来取任务，这个方法会等待一定的时间，如果取不到任务就返回null。

然后判断取到的任务r是否为null，为null则通过调用workerCanExit()方法来判断当前worker是否可以退出，我们看一下workerCanExit()的实现：

```

1 private boolean workerCanExit() {
2     final ReentrantLock mainLock = this.mainLock;
3     mainLock.lock();
4     boolean canExit;
5     //如果runState大于等于STOP，或者任务缓存队列为空了
6     //或者 允许为核心池线程设置空闲存活时间并且线程池中的线程数目大于1
7     try {
8         canExit = runState >= STOP ||
9                 workQueue.isEmpty() ||
10                (allowCoreThreadTimeOut &&
11                 poolSize > Math.max(1, corePoolSize));
12     } finally {
13         mainLock.unlock();
14     }
15     return canExit;
16 }

```

也就是说如果线程池处于STOP状态、或者任务队列已为空或者允许为核心池线程设置空闲存活时间并且线程数大于1时，允许worker退出。如果允许worker退出，则调用interruptIdleWorkers()中断处于空闲状态的worker，我们看一下interruptIdleWorkers()的实现：

```

1 void interruptIdleWorkers() {
2     final ReentrantLock mainLock = this.mainLock;
3     mainLock.lock();
4     try {
5         for (Worker w : workers) //实际上调用的是worker的interruptIfIdle
6             w.interruptIfIdle();
7     } finally {
8         mainLock.unlock();
9     }
10 }

```

从实现可以看出，它实际上调用的是worker的interruptIfIdle()方法，在worker的interruptIfIdle()方法中：

```

1 void interruptIfIdle() {
2     final ReentrantLock runLock = this.runLock;
3     if (runLock.tryLock()) { //注意这里，是调用tryLock()来获取锁的，
4                             //如果成功获取了锁，说明当前worker处于空闲状态
5         try {
6             if (thread != Thread.currentThread())
7                 thread.interrupt();
8         } finally {

```

```

9         runLock.unlock();
10    }
11 }
12 }

```

这里有一个非常巧妙的设计方式，假如我们来设计线程池，可能会有一个任务分派线程，当发现有线程空闲时，就从任务缓存队列中取一个任务交给空闲线程执行。但是在这里，并没有采用这样的方式，因为这样会要额外地对任务分派线程进行管理，无形地会增加难度和复杂度，这里直接让执行完任务的线程去任务缓存队列里面取任务来执行。

我们再看addIfUnderMaximumPoolSize方法的实现，这个方法的实现思想和addIfUnderCorePoolSize方法的实现思想非常相似，唯一的区别在于addIfUnderMaximumPoolSize方法是在线程池中的线程数达到了核心池大小并且往任务队列中添加任务失败的情况下执行的：

```

1 private boolean addIfUnderMaximumPoolSize(Runnable firstTask) {
2     Thread t = null;
3     final ReentrantLock mainLock = this.mainLock;
4     mainLock.lock();
5     try {
6         if (poolSize < maximumPoolSize && runState == RUNNING)
7             t = addThread(firstTask);
8     } finally {
9         mainLock.unlock();
10    }
11    if (t == null)
12        return false;
13    t.start();
14    return true;
15 }

```

看到没有，其实它和addIfUnderCorePoolSize方法的实现基本一模一样，只是if语句判断条件中的poolSize < maximumPoolSize不同而已。

到这里，大部分朋友应该对任务提交给线程池之后到被执行的整个过程有了一个基本的了解，下面总结一下：

- 1) 首先，要清楚corePoolSize和maximumPoolSize的含义；
- 2) 其次，要知道Worker是用来起到什么作用的；
- 3) 要知道任务提交给线程池之后的处理策略，这里总结一下主要有4点：
 - 如果当前线程池中的线程数目小于corePoolSize，则每来一个任务，就会创建一个线程去执行这个任务；
 - 如果当前线程池中的线程数目 ≥ corePoolSize，则每来一个任务，会尝试将其添加到任务缓存队列当中，若添加成功，则该任务会等待空闲线程将其取出去执行；若添加失败（一般来说是任务缓存队列已满），则会尝试创建新的线程去执行这个任务；
 - 如果当前线程池中的线程数目达到maximumPoolSize，则会采取任务拒绝策略进行处理；
 - 如果线程池中的线程数量大于 corePoolSize时，如果某线程空闲时间超过keepAliveTime，线程将被终止，直至线程池中的线程数目不大于corePoolSize；如果允许为核心池中的线程设置存活时间，那么核心池中的线程空闲时间超过keepAliveTime，线程也会被终止。

3.线程池中的线程初始化

默认情况下，创建线程池之后，线程池中是没有线程的，需要提交任务之后才会创建线程。

在实际中如果需要线程池创建之后立即创建线程，可以通过以下两个方法办到：

- `prestartCoreThread()`：初始化一个核心线程；
- `prestartAllCoreThreads()`：初始化所有核心线程

下面是这2个方法的实现：

```
1 public boolean prestartCoreThread() {
2     return addIfUnderCorePoolSize(null); //注意传进去的参数是null
3 }
4
5 public int prestartAllCoreThreads() {
6     int n = 0;
7     while (addIfUnderCorePoolSize(null)) //注意传进去的参数是null
8         ++n;
9     return n;
10 }
```

注意上面传进去的参数是null，根据第2小节的分析可知如果传进去的参数为null，则最后执行线程会阻塞在`getTask`方法中的

```
1 r = workQueue.take();
```

即等待任务队列中有任务。

4.任务缓存队列及排队策略

在前面我们多次提到了任务缓存队列，即`workQueue`，它用来存放等待执行的任务。

`workQueue`的类型为`BlockingQueue<Runnable>`，通常可以取下面三种类型：

- 1) `ArrayBlockingQueue`：基于数组的先进先出队列，此队列创建时必须指定大小；
- 2) `LinkedBlockingQueue`：基于链表的先进先出队列，如果创建时没有指定此队列大小，则默认为`Integer.MAX_VALUE`；
- 3) `synchronousQueue`：这个队列比较特殊，它不会保存提交的任务，而是将直接新建一个线程来执行新来的任务。

5.任务拒绝策略

当线程池的任务缓存队列已满并且线程池中的线程数目达到`maximumPoolSize`，如果还有任务到来就会采取任务拒绝策略，通常有以下四种策略：

```
1 ThreadPoolExecutor.AbortPolicy: 丢弃任务并抛出RejectedExecutionException
2 ThreadPoolExecutor.DiscardPolicy: 也是丢弃任务，但是不抛出异常。
3 ThreadPoolExecutor.DiscardOldestPolicy: 丢弃队列最前面的任务，然后重新尝试
4 ThreadPoolExecutor.CallersRunsPolicy: 由调用线程处理该任务
```

6.线程池的关闭

ThreadPoolExecutor提供了两个方法，用于线程池的关闭，分别是shutdown()和shutdownNow()，其中：

- shutdown()：不会立即终止线程池，而是要等所有任务缓存队列中的任务都执行完后才终止，但再也不会接受新的任务
- shutdownNow()：立即终止线程池，并尝试打断正在执行的任务，并且清空任务缓存队列，返回尚未执行的任务

7.线程池容量的动态调整

ThreadPoolExecutor提供了动态调整线程池容量大小的方法：
setCorePoolSize()和setMaximumPoolSize()，

- setCorePoolSize：设置核心池大小
- setMaximumPoolSize：设置线程池最大能创建的线程数目大小

当上述参数从小变大时，ThreadPoolExecutor进行线程赋值，还可能立即创建新的线程来执行任务。

三.使用示例

前面我们讨论了关于线程池的实现原理，这一节我们来看一下它的具体使用：

```
1 public class Test {
2     public static void main(String[] args) {
3         ThreadPoolExecutor executor = new ThreadPoolExecutor(5, 10,
4             new ArrayBlockingQueue<Runnable>(5));
5
6         for(int i=0;i<15;i++){
7             MyTask myTask = new MyTask(i);
8             executor.execute(myTask);
9             System.out.println("线程池中线程数目: "+executor.getPoolS
10                 executor.getQueue().size()+", 已执行玩别的任务数目: "+exe
11         }
12         executor.shutdown();
13     }
14 }
15
16 class MyTask implements Runnable {
17     private int taskNum;
18
19     public MyTask(int num) {
20         this.taskNum = num;
21     }
22
23     @Override
24     public void run() {
25         System.out.println("正在执行task "+taskNum);
26         try {
27             Thread.currentThread().sleep(4000);
28         } catch (InterruptedException e) {
29             e.printStackTrace();
30         }
31         System.out.println("task "+taskNum+"执行完毕");
32     }
33 }
```

执行结果：

```
1 正在执行task 0
2 线程池中线程数目: 1, 队列中等待执行的任务数目: 0, 已执行玩别的任务数目: 0
3 线程池中线程数目: 2, 队列中等待执行的任务数目: 0, 已执行玩别的任务数目: 0
4 正在执行task 1
5 线程池中线程数目: 3, 队列中等待执行的任务数目: 0, 已执行玩别的任务数目: 0
6 正在执行task 2
7 线程池中线程数目: 4, 队列中等待执行的任务数目: 0, 已执行玩别的任务数目: 0
```

```

8  正在执行task 3
9  线程池中线程数目: 5, 队列中等待执行的任务数目: 0, 已执行玩别的任务数目: 0
10 正在执行task 4
11 线程池中线程数目: 5, 队列中等待执行的任务数目: 1, 已执行玩别的任务数目: 0
12 线程池中线程数目: 5, 队列中等待执行的任务数目: 2, 已执行玩别的任务数目: 0
13 线程池中线程数目: 5, 队列中等待执行的任务数目: 3, 已执行玩别的任务数目: 0
14 线程池中线程数目: 5, 队列中等待执行的任务数目: 4, 已执行玩别的任务数目: 0
15 线程池中线程数目: 5, 队列中等待执行的任务数目: 5, 已执行玩别的任务数目: 0
16 线程池中线程数目: 6, 队列中等待执行的任务数目: 5, 已执行玩别的任务数目: 0
17 正在执行task 10
18 线程池中线程数目: 7, 队列中等待执行的任务数目: 5, 已执行玩别的任务数目: 0
19 正在执行task 11
20 线程池中线程数目: 8, 队列中等待执行的任务数目: 5, 已执行玩别的任务数目: 0
21 正在执行task 12
22 线程池中线程数目: 9, 队列中等待执行的任务数目: 5, 已执行玩别的任务数目: 0
23 正在执行task 13
24 线程池中线程数目: 10, 队列中等待执行的任务数目: 5, 已执行玩别的任务数目: 0
25 正在执行task 14
26 task 3执行完毕
27 task 0执行完毕
28 task 2执行完毕
29 task 1执行完毕
30 正在执行task 8
31 正在执行task 7
32 正在执行task 6
33 正在执行task 5
34 task 4执行完毕
35 task 10执行完毕
36 task 11执行完毕
37 task 13执行完毕
38 task 12执行完毕
39 正在执行task 9
40 task 14执行完毕
41 task 8执行完毕
42 task 5执行完毕
43 task 7执行完毕
44 task 6执行完毕
45 task 9执行完毕

```

从执行结果可以看出，当线程池中线程的数目大于5时，便将任务放入任务缓存队列里面，当任务缓存队列满了之后，便创建新的线程。如果上面程序中，将for循环中改成执行20个任务，就会抛出任务拒绝异常了。

不过在java doc中，并不提倡我们直接使用ThreadPoolExecutor，而是使用Executors类中提供的几个静态方法来创建线程池：

```

1  Executors.newCachedThreadPool();           //创建一个缓冲池，缓冲池容量大小
2  Executors.newSingleThreadExecutor();       //创建容量为1的缓冲池
3  Executors.newFixedThreadPool(int);         //创建固定容量大小的缓冲池

```

下面是这三个静态方法的具体实现:

```

1  public static ExecutorService newFixedThreadPool(int nThreads) {
2      return new ThreadPoolExecutor(nThreads, nThreads,
3                                     0L, TimeUnit.MILLISECONDS,
4                                     new LinkedBlockingQueue<Runnable>())
5  }
6  public static ExecutorService newSingleThreadExecutor() {
7      return new FinalizableDelegatedExecutorService
8          (new ThreadPoolExecutor(1, 1,
9                                   0L, TimeUnit.MILLISECONDS,
10                                  new LinkedBlockingQueue<Runnable>()))
11  }
12 public static ExecutorService newCachedThreadPool() {
13     return new ThreadPoolExecutor(0, Integer.MAX_VALUE,
14                                    60L, TimeUnit.SECONDS,
15                                    new SynchronousQueue<Runnable>());
16 }

```

从它们的具体实现来看，它们实际上也是调用了ThreadPoolExecutor，只不过参数都已配置好了。

newFixedThreadPool创建的线程池corePoolSize和maximumPoolSize值是相等的，它使用的LinkedBlockingQueue；

newSingleThreadExecutor将corePoolSize和maximumPoolSize都设置为1，也使用的LinkedBlockingQueue；

newCachedThreadPool将corePoolSize设置为0，将maximumPoolSize设置为Integer.MAX_VALUE，使用的SynchronousQueue，也就是说来了任务就创建线程运行，当线程空闲超过60秒，就销毁线程。

实际中，如果Executors提供的三个静态方法能满足要求，就尽量使用它提供的三个方法，因为自己去手动配置ThreadPoolExecutor的参数有点麻烦，要根据实际任务的类型和数量来进行配置。

另外，如果ThreadPoolExecutor达不到要求，可以自己继承ThreadPoolExecutor类进行重写。

四.如何合理配置线程池的大小

本节来讨论一个比较重要的话题：如何合理配置线程池大小，仅供参考。

一般需要根据任务的类型来配置线程池大小：

如果是CPU密集型任务，就需要尽量压榨CPU，参考值可以设为 $N_{CPU}+1$

如果是IO密集型任务，参考值可以设置为 $2*N_{CPU}$

当然，这只是一个参考值，具体的设置还需要根据实际情况进行调整，比如可以先将线程池大小设置为参考值，再观察任务运行情况和系统负载、资源利用率来进行适当调整。