

# 深入理解并实现基本的信号量 (Semaphore) 机制

李睿远

Oct 31, 2025

从并发问题出发，到原理剖析，最后用代码亲手实现一个信号量。

想象一个现实世界中的场景：一个停车场只有十个车位，却有十五辆车需要停放。如果没有有效的管理机制，车辆可能会超停、争抢车位，导致混乱和冲突。在计算机科学中，这种场景对应于多线程或进程争抢有限的共享资源，例如数据库连接池、线程池或打印队列。当多个线程同时访问和修改共享数据时，会出现竞争条件，导致结果不可预测。这种问题的核心在于并发编程的挑战，包括竞争条件、临界区问题以及线程间的协调与同步。临界区是指需要互斥访问的代码段，而线程间不仅需要互斥，有时还需要协作，例如在生产者—消费者模型中。

为了解决这些问题，荷兰计算机科学家 Dijkstra 提出了一种经典的同步工具——信号量。简单来说，信号量是一个计数器，用于管理对多个进程或线程共享的资源池的访问。它通过简单的操作来实现线程间的同步，避免资源冲突和数据不一致。

## 1 信号量核心原理解析

信号量的核心数据结构包括一个整数值和一个等待队列。整数值通常初始化为非负数，代表可用资源的数量；等待队列用于存放因资源不足而阻塞的线程。信号量的基本操作是 P 操作和 V 操作，它们分别对应于等待和发送信号。P 操作源自荷兰语 Proberen，意为尝试；V 操作源自 Verhogen，意为增加。在英语中，它们常被称为 Wait 和 Signal，或 Down 和 Up。

P 操作用于申请资源。其伪代码如下所示：

---

```

1 P(Semaphore s) {
2     s.value--;
3     if (s.value < 0) {
4         将当前线程加入 s 的等待队列 ;
5         阻塞当前线程 ;
6     }
7 }
```

---

这段代码首先减少信号量的值，如果值小于零，表示资源已耗尽，当前线程会被加入等待队列并阻塞。形象地说，这就像进入停车场前查看剩余车位：如果车位充足，直接进入；否则，必须排队等待。

V 操作用于释放资源。其伪代码如下：

---

```

1 V(Semaphore s) {
2     s.value++;
3     if (s.value <= 0) {
```

---

```

5   从 s 的等待队列中移除一个线程 T;
6   唤醒线程 T;
7 }
```

这里，信号量的值增加，如果值小于或等于零，说明有线程在等待，于是从队列中唤醒一个线程。这类似于离开停车场时归还车位，并通知等待的车辆进入。

信号量分为两种类型：计数信号量和二进制信号量。计数信号量的值可以大于一，用于控制对多个实例资源的访问，如连接池。二进制信号量的值只能为零或一，主要用于实现互斥，保护临界区。值得注意的是，互斥锁是二进制信号量的一种特例，它强调所有权概念，即加锁和解锁通常由同一线程执行。

## 2 经典案例：用信号量解决生产者—消费者问题

生产者—消费者问题是一个经典的同步问题，其中生产者线程生产数据并放入缓冲区，消费者线程从缓冲区取出数据。需要确保缓冲区满时生产者等待，缓冲区空时消费者等待，同时保证对缓冲区的操作是互斥的。

解决方案使用三个信号量：empty 代表空槽位数量，初始值为缓冲区大小；full 代表已占用槽位数量，初始值为零；mutex 是二进制信号量，用于互斥访问缓冲区，初始值为一。

生产者线程的伪代码如下：

```

1 while (true) {
2     生产一个数据项 ;
3     P(empty);
4     P(mutex);
5     将数据放入缓冲区 ;
6     V(mutex);
7     V(full);
8 }
```

首先，生产者生产数据后，通过 P(empty) 申请空位；如果无空位，则阻塞。然后，通过 P(mutex) 进入临界区，确保只有一线程操作缓冲区。数据放入后，释放互斥锁并通过 V(full) 通知消费者。

消费者线程的伪代码如下：

```

while (true) {
1     P(full);
2     P(mutex);
3     从缓冲区取出一个数据项 ;
4     V(mutex);
5     V(empty);
6     消费数据项 ;
7 }
8 }
```

消费者通过 P(full) 申请数据项；如果无数据，则阻塞。然后获取互斥锁，取出数据后释放锁，并通过 V(empty)

通知生产者。关键点在于，`P(empty)` 和 `P(mutex)` 的顺序不能颠倒，否则可能造成死锁。例如，如果生产者先获取互斥锁再申请空位，而缓冲区已满，它可能阻塞并持有锁，导致消费者无法释放资源。

### 3 动手实现：用代码构建我们自己的信号量

为了深入理解信号量，我们可以用代码实现一个简单的 `Semaphore` 类。这里以 Python 风格伪代码为例，设计一个包含计数器和等待队列的类。

首先，定义 `Semaphore` 类的成员变量：`count` 表示信号量的计数值，`waitingQueue` 是一个队列用于存放等待线程，`lock` 是一个锁用于保护对 `count` 和 `waitingQueue` 的修改，确保 P 和 V 操作的原子性。

实现 P 操作（`wait` 方法）的代码如下：

```
def wait(self):
    1   with self.lock:
    2       self.count -= 1
    3       if self.count < 0:
    4           current_thread = get_current_thread()
    5           self.waitingQueue.enqueue(current_thread)
    6           sleep(current_thread)
```

这段代码使用 `with` 语句获取锁，确保操作原子性。首先减少 `count` 值，如果值小于零，将当前线程加入等待队列并使其睡眠。睡眠操作会释放锁，允许其他线程执行。

实现 V 操作（`signal` 方法）的代码如下：

```
def signal(self):
    1   with self.lock:
    2       self.count += 1
    3       if self.count <= 0:
    4           thread_to_wakeup = self.waitingQueue.dequeue()
    5           wakeup(thread_to_wakeup)
```

这里，同样先获取锁，增加 `count` 值。如果值小于或等于零，说明有线程在等待，于是从队列中取出一个线程并唤醒。被唤醒的线程会重新尝试获取锁，并从睡眠点继续执行。

这种实现采用了阻塞等待，而非忙等待，从而避免 CPU 资源浪费。原子性保证是关键，因为 P 和 V 操作本身必须是不可分割的，否则可能导致竞争条件。

### 4 现代编程语言中的信号量

在实际开发中，我们通常使用现代编程语言提供的成熟信号量实现，而不是自行构建。例如，在 Java 中，可以使用 `java.util.concurrent.Semaphore` 类。初始化时指定许可数量，如 `Semaphore sem = new Semaphore(5);`；然后通过 `sem.acquire()` 执行 P 操作，`sem.release()` 执行 V 操作。这些方法经过优化和测试，能有效处理高并发场景。

在 Python 中，`threading` 模块提供了 `Semaphore` 类。创建实例如 `sem = threading.Semaphore(5)`，然

后使用 `sem.acquire()` 和 `sem.release()` 进行操作。C++ 从 C++20 标准开始，在 `<semaphore>` 头文件中提供了信号量支持。使用这些库可以避免常见错误，提高代码可靠性和性能。

建议开发者在项目中优先使用语言内置的信号量实现，因为它们集成了底层系统优化，并能处理边缘情况，如超时和中断。

信号量是一种强大的同步原语，通过计数器和等待队列管理共享资源访问。其核心操作 P 和 V 实现了线程间的协调，计数信号量适用于资源池管理，而二进制信号量常用于互斥。与互斥锁相比，信号量更灵活，P 和 V 操作可由不同线程执行，但也更容易出错，例如如果 V 操作多于 P 操作，可能导致信号量值异常。

进阶挑战包括读者—写者问题和哲学家就餐问题，这些经典问题需要更复杂的信号量应用来避免死锁和确保公平性。例如，在读者—写者问题中，需要平衡读写线程的访问权限，防止写者饥饿。信号量使用不当可能引入风险，因此在实际应用中需仔细设计测试。

## 5 互动与参考资料

读者可以尝试用自实现的 `Semaphore` 类重写生产者—消费者问题，以加深理解。另外，考虑如何扩展实现，添加超时等待功能，例如 `wait(timeout)` 方法，这能提高程序的响应性。参考资料包括 Dijkstra 的原始论文、经典教材如《现代操作系统》和《操作系统概念》，以及 Java 和 Python 的官方文档，这些资源提供了更深入的理论和实践指导。