

Advanced Features of the SimPy Language

Norm Matloff

January 27, 2008

©2006-2008, N.S. Matloff

Contents

1 Overview	2
2 Use of SimPy's cancel() Function	2
2.1 Example: Network Timeout	2
2.2 Example: Machine with Breakdown	4
3 Job Interruption	6
3.1 Example: Machine Breakdown Again	6
3.2 Example: Network Timeout Again	9
4 Interthread Synchronization	10
4.1 Example: Yet Another Machine Breakdown Model	10
4.2 Which Comes First?	12
4.3 Waiting for Whichever Action Comes First	13
4.4 The yield queueevent Operation	13
5 Advanced Use of the Resource Class	13
5.1 Example: Network Channel with Two Levels of Service	13
5.2 Example: Call Center	16

1 Overview

In this document we present several advanced features of the SimPy language. These will make your SimPy programming more convenient and enjoyable. In small programs, use of some of these features will produce a modest but worthwhile reduction on programming effort and increase in program clarity. In large programs, the savings add up, and can make a very significant improvement.

2 Use of SimPy's `cancel()` Function

In many simulation programs, a thread is waiting for one of two events; whichever occurs first will trigger a resumption of execution of the thread. The thread will typically want to ignore the other, later-occurring event. We can use SimPy's `cancel()` function to cancel the later event.

2.1 Example: Network Timeout

An example of this is in the program `TimeOut.py`. The model consists of a network node which transmits but also sets a `timeout` period, as follows: After sending the message out onto the network, the node waits for an acknowledgement from the recipient. If an acknowledgement does not arrive within a certain specified period of time, it is assumed that the message was lost, and it will be sent again. We wish to determine the percentage of attempted transmissions which result in timeouts.

The timeout period is assumed to be 0.5, and acknowledgement time is assumed to be exponentially distributed with mean 1.0. Here is the code:

```
1  #!/usr/bin/env python
2
3  # Introductory SimPy example to illustrate the modeling of "competing
4  # events" such as timeouts, especially using SimPy's cancel() method. A
5  # network node sends a message but also sets a timeout period; if the
6  # node times out, it assumes the message it had sent was lost, and it
7  # will send again. The time to get an acknowledgement for a message is
8  # exponentially distributed with mean 1.0, and the timeout period is
9  # 0.5. Immediately after receiving an acknowledgement, the node sends
10 # out a new message.
11
12 # We find the proportion of messages which timeout. The output should
13 # be about 0.61.
14
15 # the main classes are:
16
17 #   Node, simulating the network node, with our instance being Nd
18 #   TimeOut, simulating a timeout timer, with our instance being TO
19 #   Acknowledge, simulating an acknowledgement, with our instance being ACK
20
21 # overview of program design:
22
23 #   Nd acts as the main "driver," with a loop that continually creates
24 #   TimeOuts and Acknowledge objects, passivating itself until one of
25 #   those objects' events occurs; if for example the timeout occurs
26 #   before the acknowledge, the TO object will reactivate Nd and cancel
27 #   the ACK object's event, and vice versa
28
29 from SimPy.Simulation import *
30 from random import Random,expovariate
```

```

31
32 class Node(Process):
33     def __init__(self):
34         Process.__init__(self)
35         self.NMsgs = 0 # number of messages sent
36         self.NTimeOuts = 0 # number of timeouts which have occurred
37         # ReactivatedCode will be 1 if timeout occurred, 2 ACK if received
38         self.ReactivatedCode = None
39     def Run(self):
40         while 1:
41             self.NMsgs += 1
42             # set up the timeout
43             G.TO = Timeout()
44             activate(G.TO,G.TO.Run())
45             # set up message send/ACK
46             G.ACK = Acknowledge()
47             activate(G.ACK,G.ACK.Run())
48             yield passivate,self
49             if self.ReactivatedCode == 1:
50                 self.NTimeOuts += 1
51             self.ReactivatedCode = None
52
53 class Timeout(Process):
54     TOPeriod = 0.5
55     def __init__(self):
56         Process.__init__(self)
57     def Run(self):
58         yield hold,self,Timeout.TOPeriod
59         G.Nd.ReactivatedCode = 1
60         reactivate(G.Nd)
61         self.cancel(G.ACK)
62
63 class Acknowledge(Process):
64     ACKRate = 1/1.0
65     def __init__(self):
66         Process.__init__(self)
67     def Run(self):
68         yield hold,self,G.Rnd.expovariate(Acknowledge.ACKRate)
69         G.Nd.ReactivatedCode = 2
70         reactivate(G.Nd)
71         self.cancel(G.TO)
72
73 class G: # globals
74     Rnd = Random(12345)
75     Nd = Node()
76
77 def main():
78     initialize()
79     activate(G.Nd,G.Nd.Run())
80     simulate(until=10000.0)
81     print 'the percentage of timeouts was', float(G.Nd.NTimeOuts)/G.Nd.NMsgs
82
83 if __name__ == '__main__': main()

```

The main driver here is a class **Node**, whose PEM code includes the lines

```

1 while 1:
2     self.NMsgs += 1
3     G.TO = Timeout()
4     activate(G.TO,G.TO.Run())
5     G.ACK = Acknowledge()
6     activate(G.ACK,G.ACK.Run())
7     yield passivate,self
8     if self.ReactivatedCode == 1:
9         self.NTimeOuts += 1
10    self.ReactivatedCode = None

```

The node creates an object **G.TO** of our **TimeOut** class, which will simulate a timeout period, and creates an object **G.ACK** of our **Acknowledge** class to simulate a transmission and acknowledgement. Then the node passivates itself, allowing **G.TO** and **G.ACK** to do their work. One of them will finish first, and then will call SimPy's **reactivate()** function to "wake up" the suspended node. The node senses whether it was a timeout or acknowledgement which woke it up, via the variable **ReactivatedCode**, and then updates its timeout count accordingly.

Here's what **TimeOut.Run()** does:

```
1 yield hold,self,TimeOut.TOPeriod
2 G.Nd.ReactivatedCode = 1
3 reactivate(G.Nd)
4 self.cancel(G.ACK)
```

It holds a random timeout time, then sets a flag in **Nd** to let the latter know that it was the timeout which occurred first, rather than the acknowledgement. Then it reactivates **Nd** and cancels **ACK**. **ACK** of course has similar code for handling the case in which the acknowledgement occurs before the timeout.

Note that in our case here, we want the thread to go out of existence when canceled. The **cancel()** function does not make that occur. It simply removes the pending events associated with the given thread. The thread is still there.

However, here the **TO** and **ACK** threads will go out of existence anyway, for a somewhat subtle reason:¹ Think of what happens when we finish one iteration of the **while** loop in **main()**. A new object of type **TimeOut** will be created, and then assigned to **G.TO**. That means that the **G.TO** no longer points to the old **TimeOut** object, and since nothing else points to it either, the Python interpreter will now **garbage collect** that old object.

2.2 Example: Machine with Breakdown

Here is another example of **cancel()**:

```
1 #!/usr/bin/env python
2
3 # JobBreak.py
4
5 # One machine, which sometimes breaks down. Up time and repair time are
6 # exponentially distributed. There is a continuing supply of jobs
7 # waiting to use the machine, i.e. when one job finishes, another
8 # immediately begins. When a job is interrupted by a breakdown, it
9 # resumes "where it left off" upon repair, with whatever time remaining
10 # that it had before.
11
12 from SimPy.Simulation import *
13 from random import Random,expovariate
14
15 import sys
16
17 class G: # globals
18     CurrentJob = None
19     Rnd = Random(12345)
20     M = None # our one machine
21
22 class Machine(Process):
```

¹Thanks to Travis Grathwell for pointing this out.

```

23 def __init__(self):
24     Process.__init__(self)
25 def Run(self):
26     while 1:
27         UpTime = G.Rnd.expovariate(Machine.UpRate)
28         yield hold,self,UpTime
29         CJ = G.CurrentJob
30         self.cancel(CJ)
31         NewNInts = CJ.NInts + 1
32         NewTimeLeft = CJ.TimeLeft - (now()-CJ.LatestStart)
33         RepairTime = G.Rnd.expovariate(Machine.RepairRate)
34         yield hold,self,RepairTime
35         G.CurrentJob = Job(CJ.ID,NewTimeLeft,NewNInts,CJ.OrigStart,now())
36         activate(G.CurrentJob,G.CurrentJob.Run())
37
38 class Job(Process):
39     ServiceRate = None
40     NDone = 0 # jobs done so far
41     TotWait = 0.0 # total wait for those jobs
42     NNoInts = 0 # jobs done so far that had no interruptions
43 def __init__(self,ID,TimeLeft,NInts,OrigStart,LatestStart):
44     Process.__init__(self)
45     self.ID = ID
46     self.TimeLeft = TimeLeft # amount of work left for this job
47     self.NInts = NInts # number of interruptions so far
48     # time this job originally started
49     self.OrigStart = OrigStart
50     # time the latest work period began for this job
51     self.LatestStart = LatestStart
52 def Run(self):
53     yield hold,self,self.TimeLeft
54     # job done
55     Job.NDone += 1
56     Job.TotWait += now() - self.OrigStart
57     if self.NInts == 0: Job.NNoInts += 1
58     # start the next job
59     SrvTm = G.Rnd.expovariate(Job.ServiceRate)
60     G.CurrentJob = Job(G.CurrentJob.ID+1,SrvTm,0,now(),now())
61     activate(G.CurrentJob,G.CurrentJob.Run())
62
63 def main():
64     Job.ServiceRate = float(sys.argv[1])
65     Machine.UpRate = float(sys.argv[2])
66     Machine.RepairRate = float(sys.argv[3])
67     initialize()
68     SrvTm = G.Rnd.expovariate(Job.ServiceRate)
69     G.CurrentJob = Job(0,SrvTm,0,0.0,0.0)
70     activate(G.CurrentJob,G.CurrentJob.Run())
71     G.M = Machine()
72     activate(G.M,G.M.Run())
73     MaxSimtime = float(sys.argv[4])
74     simulate(until=MaxSimtime)
75     print 'mean wait:', Job.TotWait/Job.NDone
76     print '% of jobs with no interruptions:', \
77         float(Job.NNoInts)/Job.NDone
78
79 if __name__ == '__main__': main()

```

Here we have one machine, with occasional breakdown, but we also keep track of the number of jobs done. See the comments in the code for details.

Here we have set up a class **Job**. When a new job starts service, an instance of this class is set up to model that job. If its service then runs to completion without interruption, fine. But if the machine breaks down in the midst of service, this instance of the **Job** class will be discarded, and a new instance will later be created when this job resumes service after the repair. In other words, each object of the class **Job** models one job

to be done, but it can be either a brand new job or the resumption of an interrupted job.

Let's take a look at **Job.Run()**:

```
1 yield hold, self, self.TimeLeft
2 Job.NDone += 1
3 Job.TotWait += now() - self.OrigStart
4 if self.NInts == 0: Job.NNoInts += 1
5 SrvTm = G.Rnd.expovariate(Job.ServiceRate)
6 G.CurrentJob = Job(G.CurrentJob.ID+1, SrvTm, 0, now(), now())
7 activate(G.CurrentJob, G.CurrentJob.Run())
```

This looks innocuous enough. We hold for the time it takes to finish the job, then update our totals, and launch the next job. What is not apparent, though, is that we may actually never reach that second line,

```
Job.NDone += 1
```

The reason for this is that the machine may break down before the job finishes. In that case, what we have set up is that **Machine.Run()** will cancel the pending job completion event,

```
self.cancel(CJ)
```

simulate the repair of the machine,

```
RepairTime = G.Rnd.expovariate(Machine.RepairRate)
yield hold, self, RepairTime
```

and then create a new instance of **Job** which will simulate the processing of the remainder of the interrupted job (which may get interrupted too):

```
NewNInts = CJ.NInts + 1
NewTimeLeft = CJ.TimeLeft - (now() - CJ.LatestStart)
...
G.CurrentJob = Job(CJ.ID, NewTimeLeft, NewNInts, CJ.OrigStart, now())
activate(G.CurrentJob, G.CurrentJob.Run())
```

There are other ways of doing this, in particular by using SimPy's **interrupt()** and **interrupted()** functions, but we defer this to Section 3.

3 Job Interruption

SimPy allows one thread to interrupt another, which can be very useful.

3.1 Example: Machine Breakdown Again

In Section 2.2 we had a program **JobBreak.py**, which modeled a machine with breakdown on which we collected job time data. We presented that program as an example of **cancel()**. However, it is much more easily handled via the function **interrupt()**. Here is a new version of the program using that function:

```

1  #!/usr/bin/env python
2
3  # JobBreakInt.py: illustration of interrupt() and interrupted()
4
5  # One machine, which sometimes breaks down. Up time and repair time are
6  # exponentially distributed. There is a continuing supply of jobs
7  # waiting to use the machine, i.e. when one job finishes, the next
8  # begins. When a job is interrupted by a breakdown, it resumes "where
9  # it left off" upon repair, with whatever time remaining that it had
10 # before.
11
12 from SimPy.Simulation import *
13 from random import Random,expovariate
14
15 import sys
16
17 class G: # globals
18     CurrentJob = None
19     Rnd = Random(12345)
20     M = None # our one machine
21
22 class Machine(Process):
23     def __init__(self):
24         Process.__init__(self)
25     def Run(self):
26         from SimPy.Simulation import _e
27         while 1:
28             UpTime = G.Rnd.expovariate(Machine.UpRate)
29             yield hold,self,UpTime
30             self.interrupt(G.CurrentJob)
31             RepairTime = G.Rnd.expovariate(Machine.RepairRate)
32             yield hold,self,RepairTime
33             reactivate(G.CurrentJob)
34
35 class Job(Process):
36     ServiceRate = None
37     NDone = 0 # jobs done so far
38     TotWait = 0.0 # total wait for those jobs
39     NNoInts = 0 # jobs done so far that had no interruptions
40     NextID = 0
41     def __init__(self):
42         Process.__init__(self)
43         self.ID = Job.NextID
44         Job.NextID += 1
45         # amount of work left for this job
46         self.TimeLeft = G.Rnd.expovariate(Job.ServiceRate)
47         self.NInts = 0 # number of interruptions so far
48         # time this job originally started
49         self.OrigStart = now()
50         # time the latest work period began for this job
51         self.LatestStart = now()
52     def Run(self):
53         from SimPy.Simulation import _e
54         while True:
55             yield hold,self,self.TimeLeft
56             # did the job run to completion?
57             if not self.interrupted(): break
58             self.NInts += 1
59             self.TimeLeft -= now() - self.LatestStart
60             yield passivate,self # wait for repair
61             self.LatestStart = now()
62         Job.NDone += 1
63         Job.TotWait += now() - self.OrigStart
64         if self.NInts == 0: Job.NNoInts += 1
65         # start the next job
66         G.CurrentJob = Job()
67         activate(G.CurrentJob,G.CurrentJob.Run())
68

```

```

69 def main():
70     Job.ServiceRate = float(sys.argv[1])
71     Machine.UpRate = float(sys.argv[2])
72     Machine.RepairRate = float(sys.argv[3])
73     initialize()
74     G.CurrentJob = Job()
75     activate(G.CurrentJob,G.CurrentJob.Run())
76     G.M = Machine()
77     activate(G.M,G.M.Run())
78     MaxSimtime = float(sys.argv[4])
79     simulate(until=MaxSimtime)
80     print 'mean wait:', Job.TotWait/Job.NDone
81     print '% of jobs with no interruptions:', \
82           float(Job.NNoInts)/Job.NDone
83
84 if __name__ == '__main__': main()

```

The first key part of **Machine.Run()** is

```

yield hold,self,UpTime
self.interrupt(G.CurrentJob)

```

A call to **interrupt()** cancels the pending **yield hold** operation of its “victim,” i.e. the thread designated in the argument.² A new artificial event will be created for the victim, with event time being the current simulated time, **now()**. The caller does not lose control of the CPU, and continues to execute, but when it hits its next **yield** statement (or **passivate()** etc.) and thus loses control of the CPU, the victim will probably be next to run, as its (new, artificial) event time will be the current time.

In our case here, at the time

```

self.interrupt(G.CurrentJob)

```

is executed by the **Machine** thread, the current job is in the midst of being serviced. The call interrupts that service, to reflect the fact that the machine has broken down. At this point, the current job’s event is canceled, with the artificial event being created as above. The current job’s thread won’t run yet, and the **Machine** thread will continue. But when the latter reaches the line

```

yield hold,self,RepairTime

```

the **Machine** thread loses control of the CPU and the current job’s thread runs. The latter executes

```

if not self.interrupted(): break
self.NInts += 1
self.TimeLeft -= now() - self.LatestStart
yield passivate,self # wait for repair

```

The interruption will be sensed by **self.interrupted()** returning True. The job thread will then do the proper bookkeeping, and then passivate itself, waiting for the machine to come back up. When the latter event occurs, the machine’s thread executes

```

reactivate(G.CurrentJob)

```

²The function **interrupt()** should not be called unless the thread to be interrupted is in the midst of **yield hold**.

to get the interrupted job started again.

Note that a job may go through multiple cycles of run, interruption, run, interruption, etc., depending on how many breakdowns the machine has during the lifetime of this job. This is the reason for the **while** loop in **Job.Run()**:

```
1 while True:
2     yield hold,self,self.TimeLeft
3     # did the job run to completion?
4     if not self.interrupted(): break
5     self.NInts += 1
6     self.TimeLeft -= now() - self.LatestStart
7     yield passivate,self # wait for repair
8     self.LatestStart = now()
9     Job.NDone += 1
10    Job.TotWait += now() - self.OrigStart
11    ...
```

In the job's final cycle (which could be its first), the **yield hold** will not be interrupted. In this case the call to **interrupted()** will inform the thread that it had *not* been interrupted. The loop will be exited, the final bookkeeping for this job will be done, and the next job will be started.

By the way, we did not have to have our instance variable **TimeLeft** in **Job**. SimPy's **Process** class has its own built-in instance variable **interruptLeft** which records how much time in the **yield hold** had been remaining at the time of the interruption.

3.2 Example: Network Timeout Again

Use of interrupts makes our old network node acknowledgement/timeout program **TimeOut.py** in Section 2.1 considerably simpler:

```
1 #!/usr/bin/env python
2
3 # TimeOutInt.py
4
5 # Same as TimeOut.py but using interrupts. A network node sends a message
6 # but also sets a timeout period; if the node times out, it assumes the
7 # message it had sent was lost, and it will send again. The time to get
8 # an acknowledgement for a message is exponentially distributed with
9 # mean 1.0, and the timeout period is 0.5. Immediately after receiving
10 # an acknowledgement, the node sends out a new message.
11
12 # We find the proportion of messages which timeout. The output should
13 # be about 0.61.
14
15 from SimPy.Simulation import *
16 from random import Random,expovariate
17
18 class Node(Process):
19     def __init__(self):
20         Process.__init__(self)
21         self.NMsgs = 0 # number of messages sent
22         self.NTimeOuts = 0 # number of timeouts which have occurred
23     def Run(self):
24         from SimPy.Simulation import _e
25         while 1:
26             self.NMsgs += 1
27             # set up the timeout
28             G.TO = TimeOut()
```

```

29         activate(G.TO,G.TO.Run())
30         # wait for ACK, but could be timeout
31         yield hold,self,G.Rnd.expovariate(1.0)
32         if self.interrupted():
33             self.NTimeOuts += 1
34         else: self.cancel(G.TO)
35
36 class Timeout(Process):
37     TOPeriod = 0.5
38     def __init__(self):
39         Process.__init__(self)
40     def Run(self):
41         from SimPy.Simulation import _e
42         yield hold,self,Timeout.TOPeriod
43         self.interrupt(G.Nd)
44
45 class G: # globals
46     Rnd = Random(12345)
47     Nd = Node()
48
49 def main():
50     initialize()
51     activate(G.Nd,G.Nd.Run())
52     simulate(until=10000.0)
53     print 'the percentage of timeouts was', float(G.Nd.NTimeOuts)/G.Nd.NMsgs
54
55 if __name__ == '__main__': main()

```

Use of interrupts allowed us to entirely eliminate our old **ACK** class. Moreover, the code looks more natural now, as a timeout could be thought of as “interrupting” the node.

4 Interthread Synchronization

In our introductory SimPy document, in cases in which one thread needed to wait for some other thread to take some action,³ we made use of **passivate()** and **reactivate()**. Those can be used in general, but more advanced constructs would make our lives easier.

For example, suppose many threads are waiting for the same action to occur. The thread which triggered that action would then have to call **reactivate()** on all of them. Among other things, this would mean we would have to have code which kept track of which threads were waiting. We could do that, but it would be nicer if we didn’t have to.

In fact, actions like **yield waitevent** alleviate us of that burden. This makes our code easier to write and maintain, and easier to read.

4.1 Example: Yet Another Machine Breakdown Model

Below is an example, again modeling a machine repair situation. It is similar to **MachRep3.py** from our introductory document, but with R machines instead of two, and a policy that the repairperson is called if the number of operating machines falls below K.

³I’ve used the word *action* here rather than *event*, as the latter term refers to items in SimPy’s internal event list, generated by **yield hold** operations. But this won’t completely remove the confusion, as the SimPy keyword **waitevent** will be introduced below. But again, that term will refer to what I’m describing as *actions* here. The official SimPy term is a *SimEvent*.

```

1  #!/usr/bin/env python
2
3  # MachRep4.py
4
5  # SimPy example: R machines, which sometimes break down. Up time is
6  # exponentially distributed with rate UpRate, and repair time is
7  # exponentially distributed with rate RepairRate. The repairperson is
8  # summoned when fewer than K of the machines are up, and reaches the
9  # site after a negligible amount of time. He keeps repairing machines
10 # until there are none that need it, then leaves.
11
12 # usage:  python MachRep4.py R UpRate RepairRate K MaxSimTime
13
14 from SimPy.Simulation import *
15 from random import Random,expovariate
16
17 class G: # globals
18     Rnd = Random(12345)
19     RepairPerson = Resource(1)
20     RepairPersonOnSite = False
21     RPArrive = SimEvent()
22
23 class MachineClass(Process):
24     MachineList = [] # list of all objects of this class
25     UpRate = None # reciprocal of mean up time
26     RepairRate = None # reciprocal of mean repair time
27     R = None # number of machines
28     K = None # threshold for summoning the repairperson
29     TotalUpTime = 0.0 # total up time for all machines
30     NextID = 0 # next available ID number for MachineClass objects
31     NUp = 0 # number of machines currently up
32     # create an event to signal arrival of repairperson
33     def __init__(self):
34         Process.__init__(self)
35         self.StartupTime = None # time the current up period started
36         self.ID = MachineClass.NextID # ID for this MachineClass object
37         MachineClass.NextID += 1
38         MachineClass.MachineList.append(self)
39         MachineClass.NUp += 1 # start in up mode
40     def Run(self):
41         from SimPy.Simulation import _e
42         while 1:
43             self.StartupTime = now()
44             yield hold,self,G.Rnd.expovariate(MachineClass.UpRate)
45             MachineClass.TotalUpTime += now() - self.StartupTime
46             MachineClass.NUp -= 1
47             # if the repairperson is already onsite, just request him;
48             # otherwise, check whether fewer than K machines are up
49             if not G.RepairPersonOnSite:
50                 if MachineClass.NUp < MachineClass.K:
51                     G.RPArrive.signal()
52                     G.RepairPersonOnSite = True
53                 else: yield waitevent,self,G.RPArrive
54             yield request,self,G.RepairPerson
55             yield hold,self,G.Rnd.expovariate(MachineClass.RepairRate)
56             MachineClass.NUp += 1
57             # if no more machines waiting for repair, dismiss repairperson
58             if G.RepairPerson.waitQ == []:
59                 G.RepairPersonOnSite = False
60             yield release,self,G.RepairPerson
61
62 def main():
63     initialize()
64     MachineClass.R = int(sys.argv[1])
65     MachineClass.UpRate = float(sys.argv[2])
66     MachineClass.RepairRate = float(sys.argv[3])
67     MachineClass.K = int(sys.argv[4])
68     for I in range(MachineClass.R):

```

```

69     M = MachineClass()
70     activate(M,M.Run())
71     MaxSimtime = float(sys.argv[5])
72     simulate(until=MaxSimtime)
73     print 'proportion of up time was', \
74           MachineClass.TotalUpTime/(MachineClass.R*MaxSimtime)
75
76 if __name__ == '__main__': main()

```

Here we make use of a new SimPy class, **SimEvent**:

```

RepairPersonOnSite = False
RPArrive = SimEvent()

```

We also set up a variable **RepairPersonOnSite** to keep track of whether the repairperson is currently available; more on this point below.

Here is the core code, executed when a machine goes down:

```

MachineClass.NUp -= 1
if not G.RepairPersonOnSite:
    if MachineClass.NUp < MachineClass.K:
        G.RPArrive.signal()
        G.RepairPersonOnSite = True
    else: yield waitevent,self,G.RPArrive
yield request,self,G.RepairPerson

```

If the repairperson is on site already, then we go straight to the **yield request** to queue up for repair. If the repairperson is not on site, and the number of working machines has not yet dropped below **K**, our machine executes **yield waitevent** on our action **G.RPArrive**, which basically passivates this thread. If on the other hand our machine's failure does make the number of working machines drop below **K**, we execute the **signal()** function, which reactivates all the machines which had been waiting.

Again, all of that could have been done via explicit **passivate()** and **reactivate()** calls, but it's much more convenient to let SimPy do that work for us, behind the scenes.

One of the member variables of **SimEvent** is **occurred**, which of course is a boolean variable stating whether the action has occurred yet. Note that as soon as a wait for an event finishes, this variable reverts to **False**. This is why we needed a separate variable above, **G.RepairPersonOnSite**.

4.2 Which Comes First?

In general thread terminology, we say that we **post** a signal when we call **signal()**. One of the issues to resolve when you learn any thread system concerns what happens when a signal is posted before any waits for it are executed. In many thread systems, that posting will be completely ignored, and subsequent waits will thus last forever, or at least until another signal is posted. This obviously can cause bugs and makes programming more difficult.

In SimPy it's the opposite: If a signal is posted first, before any waits are started, subsequent waits will return immediately. That was not an issue in this program, but it's important to keep in mind in general.

4.3 Waiting for Whichever Action Comes First

You can also use **yield waiteventj** to wait for several actions, producing a “whichever comes first” operation. To do this, instead of using the form

```
yield waitevent, self action_name
```

use

```
yield waitevent, self tuple_or_list_of_action_names
```

Then whenever a signal is invoked on any one of the specified actions occurs, all actions queued will be reactivated.

4.4 The yield queueevent Operation

This works just like **yield waitevent**, but when the signal is invoked, only the action at the head of the queue will be reactivated.

5 Advanced Use of the Resource Class

The default queuing **discipline**, i.e. priority policy, for the **Resource** class is First Come, First Served (FCFS). The alternative is to assign different priorities to threads waiting for the resource, which you do by the named argument **qType**. For example,

```
R = Resource(8, qType=PriorityQ)
```

creates a resource **R** with eight service units, the queue for which has priorities assigned. The priorities are specified in the **yield request** statement. For instance,

```
yield request, self, R, 88
```

requests to use the resource **R**, with priority 88. The priorities are user-defined.

5.1 Example: Network Channel with Two Levels of Service

Below is an example of a model in which we use the non-FCFS version of **Resource**. Here we have a shared network channel on which both video and data are being transmitted. The two types of traffic act in complementary manners:

- We can tolerate a certain percentage of lost video packets, as small loss just causes a bit of jitter on the screen. But we can't have any noticeable delay.
- We can tolerate a certain increase in delay for data packets. We won't care about or even notice a small increase in delay. But we can't lose packets.

Accordingly,

- We discard video packets that are too “old,” with threshold being controlled by the design parameter L explained in the comments in the program below.
- We don’t discard data packets.

For a fixed level of data traffic, we can for example use simulation to study the tradeoff arising from our choice of the value of L. Larger L means more lost video packets but smaller delay for data, and vice versa.

Here is the program:

```
1  #!/usr/bin/env python
2
3  # QoS.py:  illustration of non-FCFS priorities in Resource class
4
5  # Communications channel, shared by video and data.  Video packets
6  # arrive every 2.0 amount of time, and have transmission time 1.0.  Data
7  # packet interarrivals are exponentially distributed with rate DArrRate,
8  # and their transmission time is uniformly distributed on {1,2,3,4,5}.
9  # Video packets have priority over data packets but the latter are not
10 # pre-emptable.  A video packet is discarded upon arrival if it would be
11 # sent L or more amount of time late.
12
13 # usage:  python QoS.py DArrRate L MaxSimTime
14
15 from SimPy.Simulation import *
16 from random import Random,expovariate
17
18 class G: # globals
19     Rnd = Random(12345)
20     Chnl = None # our one channel
21     VA = None # our one video arrivals process
22     DA = None # our one video arrivals process
23
24 class ChannelClass(Resource):
25     def __init__(self):
26         # note arguments to parent constructor:
27         Resource.__init__(self,capacity=1,qType=PriorityQ)
28         # if a packet is currently being sent, here is when transmit will end
29         self.TimeEndXMit = None
30         self.NWaitingVid = 0 # number of video packets in queue
31
32 class VidJob(Process):
33     def __init__(self):
34         Process.__init__(self)
35     def Run(self):
36         Lost = False
37         # if G.Chnl.TimeEndXMit is None, then no jobs in the system
38         # now, so this job will start right away (handled below);
39         # otherwise:
40         if G.Chnl.TimeEndXMit != None:
41             # first check for loss
42             TimeThisPktStartXMit = G.Chnl.TimeEndXMit + G.Chnl.NWaitingVid
43             if TimeThisPktStartXMit - now() > VidArrivals.L:
44                 Lost = True
45                 VidArrivals.NLost += 1
46             return
47         G.Chnl.NWaitingVid += 1
48         yield request,self,G.Chnl,1 # higher priority
49         G.Chnl.NWaitingVid -= 1
50         G.Chnl.TimeEndXMit = now() + 0.999999999999
51         yield hold,self,0.999999999999 # to avoid coding "ties"
```

```

52     G.Chnl.TimeEndXMit = None
53     yield release,self,G.Chnl
54
55 class VidArrivals(Process):
56     L = None # threshold for discarding packet
57     NArrived = 0 # number of video packets arrived
58     NLost = 0 # number of video packets lost
59     def __init__(self):
60         Process.__init__(self)
61     def Run(self):
62         while 1:
63             yield hold,self,2.0
64             VidArrivals.NArrived += 1
65             V = VidJob()
66             activate(V,V.Run())
67
68 class DataJob(Process):
69     def __init__(self):
70         Process.__init__(self)
71         self.ArrivalTime = now()
72     def Run(self):
73         yield request,self,G.Chnl,0 # lower priority
74         XMitTime = G.Rnd.randint(1,6) - 0.000000000001
75         G.Chnl.TimeEndXMit = now() + XMitTime
76         yield hold,self,XMitTime
77         G.Chnl.TimeEndXMit = None
78         DataArrivals.NSent += 1
79         DataArrivals.TotWait += now() - self.ArrivalTime
80         yield release,self,G.Chnl
81
82 class DataArrivals(Process):
83     DArrRate = None # data arrival rate
84     NSent = 0 # number of video packets arrived
85     TotWait = 0.0 # number of video packets lost
86     def __init__(self):
87         Process.__init__(self)
88     def Run(self):
89         while 1:
90             yield hold,self,G.Rnd.expovariate(DataArrivals.DArrRate)
91             D = DataJob()
92             activate(D,D.Run())
93
94 # def ShowStatus():
95 #     print 'time', now()
96 #     print 'current xmit ends at', G.Chnl.TimeEndXMit
97 #     print 'there are now',len(G.Chnl.waitQ), 'in the wait queue'
98 #     print G.Chnl.NWaitingVid, 'of those are video packets'
99
100 def main():
101     initialize()
102     VidArrivals.L = float(sys.argv[1])
103     DataArrivals.DArrRate = float(sys.argv[2])
104     G.Chnl = ChannelClass()
105     G.VA = VidArrivals()
106     activate(G.VA,G.VA.Run())
107     G.DA = DataArrivals()
108     activate(G.DA,G.DA.Run())
109     MaxSimtime = float(sys.argv[3])
110     simulate(until=MaxSimtime)
111     print 'proportion of video packets lost:', \
112           float(VidArrivals.NLost)/VidArrivals.NArrived
113     print 'mean delay for data packets:', \
114           DataArrivals.TotWait/DataArrivals.NSent
115
116 if __name__ == '__main__': main()

```

We have chosen to make a subclass of **Resource** for channels. In doing so, we do have to be careful when our subclass' constructor calls **Resource**'s constructor:

```
Resource.__init__(self, capacity=1, qType=PriorityQ)
```

The named argument **capacity** is the number of resource units, which is 1 in our case. I normally don't name it in my **Resource** calls, as it is the first argument and thus doesn't need to be named, but in this case I've used the name for clarity. And of course I've put in the **qType** argument.

Here is where I set the priorities:

```
yield request, self, G.Chnl, 1 # video
...
yield request, self, G.Chnl, 0 # data
```

I chose the values 1 and 0 arbitrarily. Any values would have worked, as long as the one for video was higher, to give it a higher priority.

Note that I have taken transmission times to be 0.000000000001 lower than an integer, so as to avoid "ties," in which a transmission would end exactly when messages might arrive. This is a common issue when **yield hold** times are integers.

5.2 Example: Call Center

This program simulates the operation of a call-in advice nurse system, such as the one in Kaiser Permanente. The key issue here is that the number of servers (nurses) varies through time, as the policy here is to take nurses off the shift when the number of callers is light, and to add more nurses during periods of heavy usage.

```
1 #!/usr/bin/env python
2
3 # CallCtr.py: simulation of call-in advice nurse system
4
5 # patients call in, with exponential interarrivals with rate Lambda1;
6 # they queue up for a number of advice nurses which varies through time
7 # (initially MOL); service time is exponential with rate Lambda2; if the
8 # system has been empty (i.e. no patients in the system, either being
9 # served or in the queue) for TO amount of time, the number of nurses
10 # is reduced by 1 (but it can never go below 1); a new TO period is then
11 # begun; when a new patient call comes in, if the new queue length is
12 # at least R the number of nurses is increased by 1, but it cannot go
13 # above MOL; here the newly-arrived patient is counted in the queue
14 # length
15
16 # usage:
17
18 # python CallCtr.py MOL, R, TO, Lambda1, Lambda2, MaxSimtime, Debug
19
20 from SimPy.Simulation import *
21 from random import Random, expovariate
22 import sys
23 import PeriodicSampler
24
25 # globals
26 class G:
```



```

27     Rnd = Random(12345)
28     NrsPl = None # nurse pool
29
30 class NursePool(Process):
31     def __init__(self,MOL,R,TO):
32         Process.__init__(self)
33         self.Rsrc = Resource(capacity=MOL,qType=PriorityQ) # the nurses
34         self.MOL = MOL # maximum number of nurses online
35         self.R = R
36         self.TO = TO
37         self.NrsCurrOnline = 0 # current number of nurses online
38         self.TB = None # current timebomb thread, if any
39         self.Mon = Monitor() # monitors numbers of nurses online
40         self.PrSm = PeriodicSampler.PerSmp(1.0,self.Mon,self.MonFun)
41         activate(self.PrSm,self.PrSm.Run())
42     def MonFun(self):
43         return self.NrsCurrOnline
44     def Wakeup(NrsPl,Evt): # wake nurse pool manager
45         reactivate(NrsPl)
46         # state the cause
47         NrsPl.WakingEvent = Evt
48         if G.Debug: ShowStatus(Evt)
49     def StartTimeBomb(self):
50         self.TB = TimeBomb(self.TO,self)
51         activate(self.TB,self.TB.Run())
52     def Run(self):
53         self.NrsCurrOnline = self.MOL
54         # system starts empty, so start timebomb
55         self.StartTimeBomb()
56         # this thread is a server, usually sleeping but occasionally being
57         # wakened to handle an event:
58         while True:
59             yield passivate,self # sleep until an event occurs:
60             if self.WakingEvent == 'arrival':
61                 # if system had been empty, cancel timebomb
62                 if PtClass.NPtsInSystem == 1:
63                     self.cancel(self.TB)
64                     self.TB = None
65                 else: # check for need to expand pool
66                     # how many in queue, including this new patient?
67                     NewQL = len(self.Rsrc.waitQ) + 1
68                     if NewQL >= self.R and self.NrsCurrOnline < self.MOL:
69                         # bring a new nurse online
70                         yield release,self,self.Rsrc
71                         self.NrsCurrOnline += 1
72                     continue # go back to sleep
73             if self.WakingEvent == 'departure':
74                 if PtClass.NPtsInSystem == 0:
75                     self.StartTimeBomb()
76                     continue # go back to sleep
77             if self.WakingEvent == 'timebomb exploded':
78                 if self.NrsCurrOnline > 1:
79                     # must take 1 nurse offline
80                     yield request,self,self.Rsrc,100
81                     self.NrsCurrOnline -= 1
82                 self.StartTimeBomb()
83                 continue # go back to sleep
84
85 class TimeBomb(Process):
86     def __init__(self,TO,NrsPl):
87         Process.__init__(self)
88         self.TO = TO # timeout period
89         self.NrsPl = NrsPl # nurse pool
90         self.TimeStarted = now() # for debugging
91     def Run(self):
92         yield hold,self,self.TO
93         NursePool.Wakeup(G.NrsPl,'timebomb exploded')
94

```

```

95 class PtClass(Process): # simulates one patient
96     SrvRate = None # service rate
97     NPtsInSystem = 0
98     Mon = Monitor()
99     def __init__(self):
100         Process.__init__(self)
101         self.ArrivalTime = now()
102     def Run(self):
103         # changes which trigger expansion or contraction of the nurse pool
104         # occur at arrival points and departure points
105         PtClass.NPtsInSystem += 1
106         NursePool.Wakeup(G.NrsPl,'arrival')
107         # dummy to give nurse pool thread a chance to wake up, possibly
108         # change the number of nurses, and reset the timebomb:
109         yield hold,self,0.000000000000001
110         yield request,self,G.NrsPl.Rsrc,1
111         if G.Debug: ShowStatus('service starts')
112         yield hold,self,G.Rnd.expovariate(PtClass.SrvRate)
113         yield release,self,G.NrsPl.Rsrc
114         PtClass.NPtsInSystem -= 1
115         Wait = now() - self.ArrivalTime
116         PtClass.Mon.observe(Wait)
117         NursePool.Wakeup(G.NrsPl,'departure')
118
119 class ArrivalClass(Process): # simulates patients arrivals
120     ArvRate = None
121     def __init__(self):
122         Process.__init__(self)
123     def Run(self):
124         while 1:
125             yield hold,self,G.Rnd.expovariate(ArrivalClass.ArRate)
126             Pt = PtClass()
127             activate(Pt,Pt.Run())
128
129 def ShowStatus(Evt): # for debugging and code verification
130     print
131     print Evt, 'at time', now()
132     print G.NrsPl.NrsCurrOnline, 'nurse(s) online'
133     print PtClass.NPtsInSystem, 'patient(s) in system'
134     if G.NrsPl.TB:
135         print 'timebomb started at time', G.NrsPl.TB.TimeStarted
136     else: print 'no timebomb ticking'
137
138 def main():
139     MOL = int(sys.argv[1])
140     R = int(sys.argv[2])
141     TO = float(sys.argv[3])
142     initialize()
143     G.NrsPl = NursePool(MOL,R,TO)
144     activate(G.NrsPl,G.NrsPl.Run())
145     ArrivalClass.ArRate = float(sys.argv[4])
146     PtClass.SrvRate = float(sys.argv[5])
147     A = ArrivalClass()
148     activate(A,A.Run())
149     MaxSimTime = float(sys.argv[6])
150     G.Debug = int(sys.argv[7])
151     simulate(until=MaxSimTime)
152     print 'mean wait =',PtClass.Mon.mean()
153     print 'mean number of nurses online =',G.NrsPl.Mon.mean()
154
155 if __name__ == '__main__': main()

```

Since the number of servers varies through time, we cannot use the SimPy **Resource** class in a straightforward manner, as that class assumes a fixed number of servers. However, by making use of that class' priorities capability, we can achieve the effect of a varying number of servers. Here we make use of an idea from a page on the SimPy Web site, <http://simpy.sourceforge.net/changingcapacity.htm>.

The way this works is that we remove a server from availability by performing a **yield request** with a very high priority level, a level higher than is used for any real request. In our case here, a real request is done via the line

```
yield request,self,G.NrsPl.Rsrc,1
```

with priority 1. By contrast, in order to take one nurse off the shift, we perform

```
yield request,self,self.Rsrc,100
self.NrsCurrOnline -= 1
```

The high priority ensures that this bogus “request” will prevail over any real one, with the effect that the nurse is taken offline. Note, though, that existing services are not pre-empted, i.e. a nurse is not removed from the shift in the midst of serving someone.

Note the necessity of the line

```
self.NrsCurrOnline -= 1
```

The **n** member variable of SimPy’s **Resource** class, which records the number of available resource units, would not tell us here how many nurses are available, because some of the resource units are held by the bogus “requests” in our scheme here. Thus we need a variable of our own, **NrsCurrOnline**.

As you can see from the call to **passivate()** in **NursePool.Run()**, the thread **NursePool.Run()** is mostly dormant, awakening only when it needs to add or delete a nurse from the pool. It is awakened for this purpose by the patient and “timebomb” classes, **PtClass** and **TimeBomb**, which call this function in **NursePool**:

```
def Wakeup(NrsPl,Evt):
    reactivate(NrsPl)
    NrsPl.WakingEvent = Evt
```

It wakes up the **NursePool** thread, which will then decide whether it should take action to change the size of the nurse pool, based on the argument **Evt**.

For example, when a new patient call arrives, generated by the **ArrivalClass** thread, the latter creates a **PtClass** thread, which simulates that one patient’s progress through the system. The first thing this thread does is

```
NursePool.Wakeup(G.NrsPl,'arrival')
```

so as to give the **NursePool** thread a chance to check whether the pool should be expanded.

We also have a **TimeBomb** class, which deals with the fact that if the system is devoid of patients for a long time, the size of the nurse pool will be reduced. After the given timeout period, this thread awakens the **NursePool** thread with the event ‘timebomb exploded’.

By the way, since **activate()** requires that its first argument be a class instance rather than a class, we are forced to create an instance of **NursePool**, **G.NrsPl**, even though we only have one nurse pool. That leads to

the situation we have with the function **NursePool.Wakeup()** being neither a class method nor an instance method.

Note the use of monitors, including in our **PeriodicSampler** class.

I have included a function **ShowStatus()** to help with debugging, and especially with verification of the program. Here is some sample output:

```
timebomb exploded at time 0.5
6 nurse(s) online
0 patient(s) in system
timebomb started at time 0

arrival at time 0.875581049552
5 nurse(s) online
1 patient(s) in system
timebomb started at time 0.5

service starts at time 0.875581049552
5 nurse(s) online
1 patient(s) in system
no timebomb ticking

departure at time 1.19578373243
5 nurse(s) online
0 patient(s) in system
no timebomb ticking

timebomb exploded at time 1.69578373243
5 nurse(s) online
0 patient(s) in system
timebomb started at time 1.19578373243

timebomb exploded at time 2.19578373243
4 nurse(s) online
0 patient(s) in system
timebomb started at time 1.69578373243

timebomb exploded at time 2.69578373243
3 nurse(s) online
0 patient(s) in system
timebomb started at time 2.19578373243

timebomb exploded at time 3.19578373243
2 nurse(s) online
0 patient(s) in system
timebomb started at time 2.69578373243

timebomb exploded at time 3.69578373243
1 nurse(s) online
0 patient(s) in system
timebomb started at time 3.19578373243
```