
Appendix G

What Happened on Mars?

This appendix contains two reprints:

1. "What Happened on Mars," written by Mike Jones of Microsoft;
2. "What Really Happened on Mars," written by Glenn Reeves of the NASA.

Both were widely circulated in the Internet discussion groups.

From: Mike Jones (mbj@microsoft.com)
Sent: Friday, December 05, 1997 4:15 AM

WHAT HAPPENED ON MARS?

The Mars Pathfinder mission was widely proclaimed as "flawless" in the early days after its July 4th, 1997 landing on the Martian surface. Successes included its unconventional "landing"—bouncing onto the Martian surface surrounded by airbags, deploying the Sojourner rover, and gathering and transmitting voluminous data back to Earth, including the panoramic pictures that were such a hit on the Web. But a few days into the mission, not long after Pathfinder started gathering meteorological data, the spacecraft began experiencing total

system resets, each resulting in losses of data. The press reported these failures in terms such as “software glitches” and “the computer was trying to do too many things at once”.

Last night at the IEEE Real-Time Systems Symposium I heard a fascinating keynote address by David Wilner, Chief Technical Officer of Wind River Systems. Wind River makes VxWorks, the real-time embedded systems kernel that was used in the Mars Pathfinder mission. In his talk, he explained in detail the actual software problems that caused the total system resets of the Pathfinder spacecraft, how they were diagnosed, and how they were solved. I wanted to share his story with each of you.

VxWorks provides preemptive priority scheduling of threads. Tasks on the Pathfinder spacecraft were executed as threads with priorities that were assigned in the usual manner reflecting the relative urgency of these tasks.

Pathfinder contained an “information bus”, which you can think of as a shared memory area used for passing information between different components of the spacecraft. A bus management task ran frequently with high priority to move certain kinds of data in and out of the information bus. Access to the bus was synchronized with mutual exclusion locks (mutexes).

The meteorological data gathering task ran as an infrequent, low priority thread, and used the information bus to publish its data. When publishing its data, it would acquire a mutex, do writes to the bus, and release the mutex. If an interrupt caused the information bus thread to be scheduled while this mutex was held, and if the information bus thread then attempted to acquire this same mutex in order to retrieve published data, this would cause it to block on the mutex, waiting until the meteorological thread released the mutex before it could continue. The spacecraft also contained a communications task that ran with medium priority.

Most of the time this combination worked fine. However, very infrequently it was possible for an interrupt to occur that caused the (medium priority) communications task to be scheduled during the short interval while the (high priority) information bus thread was blocked waiting for the (low priority) meteorological data thread. In this case, the long-running communications task, having higher priority than the meteorological task, would prevent it from running, consequently preventing the blocked information bus task from running. After some time had passed, a watchdog timer would go off, notice that the data bus task had not been executed for some time, conclude that something had gone drastically wrong, and initiate a total system reset.

This scenario is a classic case of priority inversion.

How was this debugged?

VxWorks can be run in a mode where it records a total trace of all interesting system events, including context switches, uses of synchronization objects, and interrupts. After the failure, JPL engineers spent hours and hours running the system on the exact spacecraft replica in their lab with tracing turned on, attempting to replicate the precise conditions under which they believed that the reset occurred. Early in the morning, after all but one engineer had gone home, the engineer finally reproduced a system reset on the replica. Analysis of the trace revealed the priority inversion.

How was this problem corrected?

When created, a VxWorks mutex object accepts a boolean parameter that indicates whether priority inheritance should be performed by the mutex. The mutex in question had been initialized with the parameter off; had it been on, the low-priority meteorological thread would

have inherited the priority of the high-priority data bus thread blocked on it while it held the mutex, causing it to be scheduled with higher priority than the medium-priority communications task, thus preventing the priority inversion. Once diagnosed, it was clear to the JPL engineers that using priority inheritance would prevent the resets they were seeing.

VxWorks contains a C language interpreter intended to allow developers to type in C expressions and functions to be executed on the fly during system debugging. The JPL engineers fortuitously decided to launch the spacecraft with this feature still enabled. By coding convention, the initialization parameter for the mutex in question (and those for two others which could have caused the same problem) were stored in global variables, whose addresses were in symbol tables also included in the launch software, and available to the C interpreter. A short C program was uploaded to the spacecraft, which when interpreted, changed the values of these variables from FALSE to TRUE. No more system resets occurred.

Analysis and Lessons

First and foremost, diagnosing this problem as a black box would have been impossible. Only detailed traces of actual system behavior enabled the faulty execution sequence to be captured and identified.

Secondly, leaving the “debugging” facilities in the system saved the day. Without the ability to modify the system in the field, the problem could not have been corrected.

Finally, the engineer’s initial analysis that “the data bus task executes very frequently and is time-critical—we shouldn’t spend the extra time in it to perform priority inheritance” was exactly wrong. It is precisely in such time critical and important situations where correctness is essential, even at some additional performance cost.

Human nature, Deadline Pressures

David told us that the JPL engineers later confessed that one or two system resets had occurred in their months of pre-flight testing. They had never been reproducible or explainable, and so the engineers, in a very human-nature response of denial, decided that they probably weren’t important, using the rationale “it was probably caused by a hardware glitch”.

Part of it too was the engineers’ focus. They were extremely focused on ensuring the quality and flawless operation of the landing software. Should it have failed, the mission would have been lost. It is entirely understandable for the engineers to discount occasional glitches in the less-critical land-mission software, particularly given that a spacecraft reset was a viable recovery strategy at that phase of the mission.

The Importance of Good Theory/Algorithms

David also said that some of the real heroes of the situation were some people from CMU who had published a paper he’d heard presented many years ago who first identified the priority inversion problem and proposed the solution. He apologized for not remembering the precise details of the paper or who wrote it. Bringing things full circle, it turns out that the three authors of this result were all in the room, and at the end of the talk were encouraged by the program chair to stand and be acknowledged. They were Lui Sha, John Lehoczky, and Raj Rajkumar. When was the last time you saw a room of people cheer a group of computer science theorists for their significant practical contribution to advancing human knowledge? :-)

It was quite a moment.

Postlude

For the record, the paper was:

L. Sha, R. Rajkumar, and J. P. Lehoczky, Priority Inheritance Protocols: An Approach to Real-Time Synchronization. In *IEEE Transactions on Computers*, vol. 39, pp. 1175-1185, Sep. 1990.

From: Glenn Reeves (glenn.e.reeves@jpl.nasa.gov)

Sent: Monday, December 15, 1997 10:28 AM

WHAT REALLY HAPPENED ON MARS?¹

By now most of you have read Mike's (mbj@microsoft.com) summary of Dave Wilner's comments given at the IEEE Real-Time Systems Symposium. I don't know Mike and I didn't attend the symposium (though I really wish I had now) and I have not talked to Dave Wilner since before the talk. However, I did lead the software team for the Mars Pathfinder spacecraft. So, instead of trying to find out what was said I will just tell you what happened. You can make your own judgments.

I sent this message out to everyone who was a recipient of Mike's original that I had an email address for. Please pass it on to anyone you sent the first one to. Mike, I hope you will post this wherever you posted the original.

Since I want to make sure the problem is clearly understood I need to step through each of the areas which contributed to the problem.

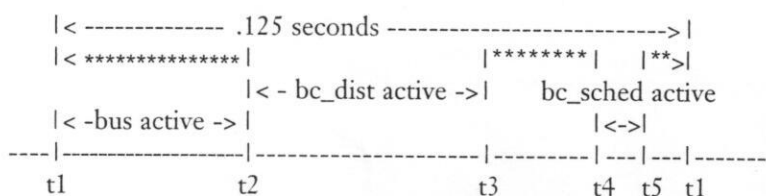
The Hardware

The simplified view of the Mars Pathfinder hardware architecture looks like this. A single CPU controls the spacecraft. It resides on a VME bus which also contains interface cards for the radio, the camera, and an interface to a 1553 bus. The 1553 bus connects to two places: The "cruise stage" part of the spacecraft and the "lander" part of the spacecraft. The hardware on the cruise part of the spacecraft controls thrusters, valves, a sun sensor, and a star scanner. The hardware on the lander part provides an interface to accelerometers, a radar altimeter, and an instrument for meteorological science known as the ASI/MET. The hardware which we used to interface to the 1553 bus (at both ends) was inherited from the Cassini spacecraft. This hardware came with a specific paradigm for its usage: the software will schedule activity at an 8 Hz rate. This **feature** dictated the architecture of the software which controls both the 1553 bus and the devices attached to it.

The Software Architecture

The software to control the 1553 bus and the attached instruments was implemented as two tasks. The first task controlled the setup of transactions on the 1553 bus (called the bus scheduler or bc_sched task) and the second task handled the collection of the transaction results i.e. the data. The second task is referred to as the bc_dist (for distribution) task. A typical timeline for the bus activity for a single cycle is shown below. It is not to scale. This cycle was constantly repeated.

¹Provided through the courtesy of the Jet Propulsion Laboratory, California Institute of Technology, Pasadena, California.



The *** are periods when tasks other than the ones listed are executing. Yes, there is some idle time.

- t1—bus hardware starts via hardware control on the 8 Hz boundary. The transactions for the this cycle had been set up by the previous execution of the bc_sched task.
- t2—1553 traffic is complete and the bc_dist task is awakened.
- t3—bc_dist task has completed all of the data distribution.
- t4—bc_sched task is awakened to setup transactions for the next cycle.
- t5—bc_sched activity is complete.

The bc_sched and bc_dist tasks check each cycle to be sure that the other had completed its execution. The bc_sched task is the highest priority task in the system (except for the vxWorks “tExec” task). The bc_dist is third highest (a task controlling the entry and landing is second). All of the tasks which perform other spacecraft functions are lower. Science functions, such as imaging, image compression, and the ASI/MET task are still lower.

Data is collected from devices connected to the 1553 bus only when they are powered. Most of the tasks in the system that access the information collected over the 1553 do so via a double buffered shared memory mechanism into which the bc_dist task places the latest data. The exception to this is the ASI/MET task which is delivered its information via an inter-process communication mechanism (IPC). The IPC mechanism uses the vxWorks pipe() facility. Tasks wait on one or more IPC “queues” for messages to arrive. Tasks use the select() mechanism to wait for message arrival. Multiple queues are used when both high and lower priority messages are required. Most of the IPC traffic in the system is not for the delivery of real-time data. However, again, the exception to this is the use of the IPC mechanism with the ASI/MET task. The cause of the reset on Mars was in the use and configuration of the IPC mechanism.

The Failure

The failure was identified by the spacecraft as a failure of the bc_dist task to complete its execution before the bc_sched task started. The reaction to this by the spacecraft was to reset the computer. This reset reinitializes all of the hardware and software. It also terminates the execution of the current ground commanded activities. No science or engineering data is lost that has already been collected (the data in RAM is recovered so long as power is not lost). However, the remainder of the activities for that day were not accomplished until the next day.

The failure turned out to be a case of priority inversion (how we discovered this and how we fixed it are covered later). The higher priority bc_dist task was blocked by the much lower priority ASI/MET task that was holding a shared resource. The ASI/MET task had acquired this resource and then been preempted by several of the medium priority tasks. When the bc_sched task was activated, to setup the transactions for the next 1553 bus cycle, it detected that the bc_dist task had not completed its execution. The resource that caused this problem was a mutual exclusion semaphore used within the select() mechanism to control access to the list of file descriptors that the select() mechanism was to wait on.

The select mechanism creates a mutual exclusion semaphore to protect the “wait list” of file descriptors for those devices which support select. The vxWorks pipe() mechanism is such a device and the IPC mechanism we used is based on using pipes. The ASI/MET task had called select, which had called pipeIoctl(), which had called selNodeAdd(), which was in the process of giving the mutex semaphore. The ASI/MET task was preempted and semGive() was not completed. Several medium priority tasks ran until the bc_dist task was activated. The bc_dist task attempted to send the newest ASI/MET data via the IPC mechanism which called pipeWrite(). pipeWrite() blocked, taking the mutex semaphore. More of the medium priority tasks ran, still not allowing the ASI/MET task to run, until the bc_sched task was awakened. At that point, the bc_sched task determined that the bc_dist task had not completed its cycle (a hard deadline in the system) and declared the error that initiated the reset.

How We Found it

The software that flies on Mars Pathfinder has several debug features within it that are used in the lab but are not used on the flight spacecraft (not used because some of them produce more information than we can send back to Earth). These features were not “fortuitously” left enabled but remain in the software by design. We strongly believe in the “test what you fly and fly what you test” philosophy.

One of these tools is a trace/log facility which was originally developed to find a bug in an early version of the vxWorks port (Wind River ported vxWorks to the RS6000 processor for use for this mission). This trace/log facility was built by David Cummings who was one of the software engineers on the task. Lisa Stanley, of Wind River, took this facility and instrumented the pipe services, msgQ services, interrupt handling, select services, and the tExec task. The facility initializes at startup and continues to collect data (in ring buffers) until told to stop. The facility produces a voluminous dump of information when asked.

After the problem occurred on Mars we did run the same set of activities over and over again in the lab. The bc_sched was already coded so as to stop the trace/log collection and dump the data (even though we knew we could not get the dump in flight) for this error. So, when we went into the lab to test it we did not have to change the software.

In less than 18 hours we were able to cause the problem to occur. Once we were able to reproduce the failure the priority inversion problem was obvious.

How Was the Problem Corrected

Once we understood the problem the fix appeared obvious: change the creation flags for the semaphore so as to enable the priority inheritance. The Wind River folks, for many of their services, supply global configuration variables for parameters such as the “options” parameter for the semMCreate used by the select service (although this is not documented and those who do not have vxWorks source code or have not studied the source code might be unaware of this feature). However, the fix is not so obvious for several reasons:

1. The code for this is in the selectLib() and is common for all device creations. When you change this global variable all of the select semaphores created after that point will be created with the new options. There was no easy way in our initialization logic to only modify the semaphore associated with the pipe used for bc_dist task to ASI/MET task communications.
2. If we make this change, and it is applied on a global basis, how will this change the behavior of the rest of the system?

3. The priority inversion option was deliberately left out by Wind River in the default `selectLib()` service for optimum performance. How will performance degrade if we turn the priority inversion on?
4. Was there some intrinsic behavior of the select mechanism itself that would change if the priority inversion was enabled?

We did end up modifying the global variable to include the priority inversion. This corrected the problem. We asked Wind River to analyze the potential impacts for (3) and (4). They concluded that the performance impact would be minimal and that the behavior of `select()` would not change so long as there was always only one task waiting for any particular file descriptor. This is true in our system. I believe that the debate at Wind River still continues on whether the priority inversion option should be on as the default. For (1) and (2) the change did alter the characteristics of all of the select semaphores. We concluded, both by analysis and test, that there was no adverse behavior. We tested the system extensively before we changed the software on the spacecraft.

How We Changed the Software on the Spacecraft

No, we did not use the vxWorks shell to change the software (although the shell is usable on the spacecraft). The process of “patching” the software on the spacecraft is a specialized process. It involves sending the differences between what you have onboard and what you want (and have on Earth) to the spacecraft. Custom software on the spacecraft (with a whole bunch of validation) modifies the onboard copy. If you want more info you can send me email.

Why Didn't We Catch it Before Launch?

The problem would only manifest itself when ASI/MET data was being collected and intermediate tasks were heavily loaded. Our before launch testing was limited to the “best case”; high data rates and science activities. The fact that data rates from the surface were higher than anticipated and the amount of science activities proportionally greater served to aggravate the problem. We did not expect nor test the “better than we could have ever imagined” case.

Human Nature, Deadline Pressures

We did see the problem before landing but could not get it to repeat when we tried to track it down. It was not forgotten nor was it deemed unimportant. Yes, we were concentrating heavily on the entry and landing software. Yes, we considered this problem lower priority. Yes, we would have liked to have everything perfect before landing. However, I don't see any problem here other than we ran out of time to get the lower priority issues completed.

We did have one other thing on our side; we knew how robust our system was because that is the way we designed it.

We knew that if this problem occurred we would reset. We built in mechanisms to recover the current activity so that there would be no interruptions in the science data (although this wasn't used until later in the landed mission). We built in the ability (and tested it) to go through multiple resets while we were going through the Martian atmosphere. We designed the software to recover from radiation induced errors in the memory or the processor. The spacecraft would have even done a 60 day mission on its own, including deploying the rover, if the radio receiver had broken when we landed. There are a large number of safeguards in

the system to ensure robust, continued operation in the event of a failure of this type. These safeguards allowed us to designate problems of this nature as lower priority.

We had our priorities right.

Analysis and Lessons

Did we (the JPL team) make an error in assuming how the select/pipe mechanism would work? Yes, probably. But there was no conscious decision to not have the priority inversion enabled. We just missed it. There are several other places in the flight software where similar protection is required for critical data structures and the semaphores do have priority inversion protection. A good lesson when you fly COTS stuff—make sure you know how it works.

Mike is quite correct in saying that we could not have figured this out ****ever**** if we did not have the tools to give us the insight. We built many of the tools into the software for exactly this type of problem. We always planned to leave them in. In fact, the shell (and the stdout stream) were very useful the entire mission. If you want more detail send me a note.

Setting the Record Straight

First, I want to make sure that everyone understands how I feel in regard to Wind River. These folks did a fantastic job for us. They were enthusiastic and supported us when we came to them and asked them to do an affordable port of vxWorks. They delivered the alpha version in 3 months. When we had a problem they put some of the brightest engineers I have ever worked with on the problem. Our communication with them was fantastic. If they had not done such a professional job the Mars Pathfinder mission would not have been the success that it is.

Second, Dave Wilner did talk to me about this problem before he gave his talk. I could not find my notes where I had detailed the description of the problem. So, I winged it and I sure did get it wrong. Sorry Dave.

Acknowledgments

First, thanks to Mike for writing a very nice description of the talk. I think I have had probably 400 people send me copies. You gave me the push to write the part of the Mars Pathfinder End-of-Mission report that I had been procrastinating doing.

A special thanks to Steve Stolper for helping me do this.

The biggest thanks should go to the software team that I had the privilege of leading and whose expertise allowed us to succeed:

Pam Yoshioka
 Dave Cummings
 Don Meyer
 Karl Schneider
 Greg Welz

Rick Achatz
 Kim Gostelow
 Dave Smyth
 Steve Stolper