

Why Userspace Sucks—Or 101 Really Dumb Things Your App Shouldn't Do

Dave Jones

Red Hat

<davej@redhat.com>

Abstract

During the development of Fedora Core 5 I found myself asking the same questions day after day:

- Why does it take so long to boot?
- Why does it take so long to start X?
- Why do I get the opportunity to go fetch a drink after starting various applications and waiting for them to load?
- Why does idling at the desktop draw so much power?
- Why does it take longer to shut down than it does to boot up?

I initially set out to discover if there was something the kernel could do better to speed up booting. A number of suggestions have been made in the past ranging from better read-ahead, to improved VM caching strategies, or better on-disk block layout. I did not get that far however, because what I found in my initial profiling was disturbing.

We have an enormous problem with applications doing unnecessary work, causing wasted time, and more power-drain than necessary.

This talk will cover a number of examples of common applications doing incredibly wasteful things, and will also detail what can be, and what has been done to improve the situation.

I intend to show by example numerous applications doing incredibly dumb things, from silliness such as reloading and re-parsing XML files 50 times each run, to applications that wake up every few seconds to ask the kernel to change the value of something that has not changed since it last woke up.

I created my tests using patches [1] to the Linux kernel, but experimented with other approaches using available tools like strace and systemtap. I will briefly discuss the use of these tools, as they apply to the provided examples, in later sections of this paper.

Our userspace sucks. Only through better education of developers of “really dumb things not to do” can we expect to resolve these issues.

1 Overview

A large number of strange things are happening behind the scenes in a lot of userspace programs. When their authors are quizzed about these discoveries, the responses range from “I had no idea it was doing that,” to “It didn't

do that on my machine.” This paper hopes to address the former by shedding light on several tools (some old, some new) that enable userspace programmers to gain some insight into what is really going on. It addresses the latter by means of showing examples that may shock, scare, and embarrass their authors into writing code with better thought out algorithms.

2 Learning from read-ahead

Improving boot up time has been a targeted goal of many distributions in recent years, with each vendor resorting to a multitude of different tricks in order to shave off a few more seconds between boot and login. One such trick employed by Fedora is the use of a read-ahead tool, which, given a list of files, simply reads them into the page cache, and then exits. During the boot process there are periods of time when the system is blocked on some non-disk I/O event such as waiting for a DHCP lease. Read-ahead uses this time to read in files that are used further along in the boot process. By seeding the page cache, the start-up of subsequent boot services will take less time provided that there is sufficient memory to prevent it from being purged by other programs starting up during the time between the read-ahead application preloaded it, and the real consumer of the data starting up.

The read-ahead approach is a primitive solution, but it works. By amortising the cost of disk IO during otherwise idle periods, we shave off a significant amount of time during boot/login. The bootchart [2] project produced a number of graphs that helped visualise progress during the early development of this tool, and later went on to provide a rewritten version for Fedora Core 5 which improved on the bootup performance even further.

The only remaining questions are, what files do we want to prefetch, and how do we generate a list of them? When the read-ahead service was first added to Fedora, the file list was created using a kernel patch that simply printk'd the filename of every file open()'d during the first five minutes of uptime. (It was necessary to capture the results over a serial console, due to the huge volume of data overflowing the dmesg ring buffer very quickly.)

This patch had an additional use however, which was to get some idea of just what IO patterns userspace was creating.

During Fedora Core 5 development, I decided to investigate these patterns. The hope was that instead of the usual approach of 'how do we make the IO scheduler better', we could make userspace be more intelligent about the sort of IO patterns it creates.

I started by extending the kernel patch to log all file IO, not just open()'s. With this new patch, the kernel reports every stat(), delete(), and path_lookup(), too.

The results were mind-boggling.

- During boot-up, 79576 files were stat()'d. 26769 were open()'d, 1382 commands were exec'd.
- During shutdown, 23246 files were stat()'d, 8724 files were open()'d.

2.1 Results from profiling

Picking through a 155234 line log took some time, but some of the things found were truly spectacular.

Some of the highlights included:

- HAL Daemon.

- Reread and reparsed *dozens* of XML files during startup. (In some cases, it did this 54 times per XML file).
- Read a bunch of files for devices that were not even present.
- Accounted for a total of 1918 `open()`'s, and 7106 `stat()`'s.

- CUPS

- Read in `ppd` files describing every printer known to man. (Even though there was not even a printer connected.)
- Responsible for around 2500 `stat()`'s, and around 500 `open()`'s.

- Xorg

A great example of how not to do PCI bus scanning.

- Scans through `/proc/bus/pci/` in order.
- **Guesses** at random bus numbers, and tries to open those devices in `/proc/bus/pci/`.
- Sequentially probes for devices on busses `0xf6` through `0xfb` (even though they may not exist).
- Retries entries that it has already attempted to scan regardless of whether they succeeded or not.

Aside from this, when it is not busy scanning non-existent PCI busses, X really likes to `stat` and `reopen` lot of files it has already opened, like `libGLcore.so`. A weakness of its dynamic loader perhaps?

- XFS

- Was rebuilding the font cache every time it booted, even if no changes had occurred in the fonts directories.

- `gdm / gnome-session`.

- Tried to open a bunch of non-existent files with odd-looking names like `/usr/share/pixmaps/Bluecurve/cursors/00000000000000000000000000000000`
- Suffers from font madness (See below).

2.2 Desktop profiling

Going further, removing the “first 5 minutes” check of the patch allowed me to profile what was going on at an otherwise idle desktop.

- `irqbalance`.

- Wakes up every 10 seconds to re-balance interrupts in a round-robin manner. Made a silly mistake where it was re-balancing interrupts where no IRQs had ever occurred. A three line change saved a few dozen syscalls.
- Was also re-balancing interrupts where an IRQ had not occurred in some time every 10 seconds.
- Did an `open/write/close` of each `/proc/irq/n/smp_affinity` file each time it rebalanced, instead of keeping the `fd`'s open, and doing 1/3rd of syscalls.

Whilst working with `/proc` files does not incur any I/O, it does trigger a transition to-and-from kernel space for each system call, adding up to a lot of unneeded work on an otherwise ‘idle’ system.

- `gamin`

- Was `stat()`'ing a bunch of `gnome` menu files every few seconds for no apparent reason.

% time	seconds	usecs/call	calls	errors	syscall
32.98	0.003376	844	4		clone
27.87	0.002853	4	699	1	read
23.50	0.002405	32	76		getdents
10.88	0.001114	0	7288	10	stat
1.38	0.000141	0	292		munmap
1.31	0.000134	0	785	382	open

Figure 1: `strace -c` output of `gnome-terminal` with lots of fonts.

- nautilus
 - Was `stat()`'ing `$HOME/Templates`, `/usr/share/applications`, and `$HOME/.local/share/applications` every few seconds even though they had not changed.
- More from the unexplained department...
 - `mixer_applet2` did a `real_lookup` on `libgstffmjpegcolorspace.so` for some bizarre reason.
 - Does `trashapplet` really need to `stat` the `svg` for every size icon when it is rarely resized?

2.3 Madness with fonts

I had noticed through reviewing the log, that a lot of applications were `stat()`'ing (and occasionally `open()`'ing) a bunch of fonts, and then never actually using them. To try to make problems stand out a little more, I copied 6000 TTF's to `$HOME/.fonts`, and reran the tests. The log file almost doubled in size.

Lots of bizarre things stood out.

- `gnome-session` `stat()`'d 2473 and `open()`'d 2434 ttf's.
- `metacity` `open()`'d another 238.

- Just to be on the safe side, `wnck-applet` `open()`'d another 349 too.
- Nautilus decided it does not want to be left out of the fun, and `open()`'d another 301.
- `mixer_applet` rounded things off by `open()`'ing 860 ttf's.

`gnome-terminal` was another oddball. It `open()`'ed 764 fonts and `stat()`'d another 770 including `re-stat()`'ing many of them multiple times. The vast majority of those fonts were not in the system-wide fonts preferences, nor in `gnome-terminals` private preferences. `strace -c` shows that `gnome-terminal` spends a not-insignificant amount of its startup time, `stat()`'ing a bunch of fonts that it never uses. (See Figure 1.)

Another really useful tool for parsing huge `strace` logs is Morten Wellinders `strace-account` [3] which takes away a lot of the tedious parsing, and points out some obvious problem areas in a nice easy-to-read summary.

Whilst having thousands of fonts is a somewhat pathological case, it is not uncommon for users to install a few dozen (or in the case of arty types, a few hundred). The impact of this defect will be less for most users, but it is still doing a lot more work than it needs to.

After my initial experiments were over, Dan Berrange wrote a set of `systemtap` scripts [4]

to provide similar functionality to my tracing kernel patch, without the need to actually patch and rebuild the kernel.

3 Learn to use tools at your disposal

Some other profiling techniques are not as intrusive as to require kernel modifications, yet remarkably, they remain under-utilised.

3.1 valgrind

For some unexplained reason, there are developers that still have not tried (or in many cases, have not heard of) valgrind [5]. This is evident from the number of applications that still output lots of scary warnings during runtime.

Valgrind can find several different types of problems, ranging from memory leaks, to the use of uninitialised memory. Figure 2 shows an example of mutt running under valgrind.

The use of uninitialised memory can be detected without valgrind, by setting the environment variable `_MALLOC_PERTURB_` [6] to a value that will cause glibc to poison memory allocated with `malloc()` to the value the variable is set to, without any need to recompile the program.

Since Fedora Core 5 development, I run with this flag set to `%RANDOM` in my `.bashrc`. It adds some overhead to some programs which call `malloc()` a lot, but it has also found a number of bugs in an assortment of packages. A `gdb` backtrace is usually sufficient to spot the area of code that the author intended to use a `calloc()` instead of a `malloc()`, or in some cases, had an incorrect `memset` call after the `malloc` returns.

3.2 oprofile

Perceived by many as complicated, oprofile is actually remarkably trivial to use. In a majority of cases, simply running

```
opcontrol --start
(do application to
 be profiled)
opcontrol --shutdown
opreport -l
```

is sufficient to discover the functions where time is being spent. Should you be using a distribution which strips symbols out to separate packages (for example, Fedora/RHEL's `-debuginfos`), you will need to install the relevant `-debuginfo` packages for the applications and libraries being profiled in order to get symbols attributed to the data collected.

3.3 Heed the warnings

A lot of developers ignore, or even suppress warnings emitted by the compiler, proclaiming "They are just warnings." On an average day, the Red Hat package-build system emits around 40–50,000 warnings as part of its daily use giving some idea of the scale of this problem.

Whilst many warnings are benign, there are several classes of warnings that can have undesirable effects. For example, an implicit declaration warning may still compile and run just fine on your 32-bit machine, but if the compiler assumes the undeclared function has `int` arguments when it actually has long arguments, unusual results may occur when the code is run on a 64-bit machine. Leaving warnings unfixed makes it easier for real problems to hide amongst the noise of the less important warnings.

```

==20900== Conditional jump or move depends on uninitialised value(s)
==20900==    at 0x3CDE59E76D: re_compile_fastmap_iter (in /lib64/libc-2.4.so)
==20900==    by 0x3CDE59EBFA: re_compile_fastmap (in /lib64/libc-2.4.so)
==20900==    by 0x3CDE5B1D23: regcomp (in /lib64/libc-2.4.so)
==20900==    by 0x40D978: ??? (color.c:511)
==20900==    by 0x40DF79: ??? (color.c:724)
==20900==    by 0x420C75: ??? (init.c:1335)
==20900==    by 0x420D8F: ??? (init.c:1253)
==20900==    by 0x422769: ??? (init.c:1941)
==20900==    by 0x42D631: ??? (main.c:608)
==20900==    by 0x3CDE51D083: __libc_start_main (in /lib64/libc-2.4.so)

```

Figure 2: Mutt under valgrind

Bonus warnings can be enabled with compiler options `-Wall` and `-Wextra` (this option used to be `-W` in older gcc releases)

For the truly anal, static analysis tools such as splint [7], and sparse [8] may turn up additional problems.

In March 2006, a security hole was found in Xorg [9] by the Coverity Prevent scanner [10]. The code looked like this.

```
if (getuid() == 0 || geteuid != 0)
```

No gcc warnings are emitted during this compilation, as it is valid code, yet it does completely the wrong thing. Splint on the other hand indicates that something is amiss here with the warning:

```

Operands of != have incompatible types
([function (void) returns
__uid_t], int): geteuid != 0
Types are incompatible.
(Use -type to inhibit warning)

```

Recent versions of gcc also allow programs to be compiled with the `-D_FORTIFY_SOURCE=2` which enables various security checks in various C library functions. If the size of memory passed to functions such as `memcpy` is known at compile time, warnings will be emitted if the `len` argument overruns the buffer being passed.

Additionally, use of certain functions without checking their return code will also result in a warning. Some 30-40 or so C runtime functions have had such checks added to them.

It also traps a far-too-common¹ bug: `memset` with size and value arguments transposed. Code that does this:

```
memset(ptr, sizeof(foo), 0);
```

now gets a compile time warning which looks like this:

```

warning: memset used with
constant zero length parameter;
this could be due to transposed
parameters

```

Even the simpler (and deprecated) `bzero` function is not immune from screwups of the size parameter it seems, as this example shows:

```
bzero(pages + npagesmax, npagesmax
- npagesmax);
```

Another useful gcc feature that was deployed in Fedora Core 5 was the addition of a stack overflow detector. With all applications compiled with the flags

¹Across 1282 packages in the Fedora tree, 50 of them had a variant of this bug.

```
-fstack-protector --param=
ssp-buffer-size=4
```

any attempt at overwriting an on-stack buffer results in the program being killed with the following message:

```
*** stack smashing detected ***:
./a.out terminated
Aborted (core dumped)
```

This turned up a number of problems during development, which were usually trivial to fix up.

4 Power measurement

The focus of power management has traditionally been aimed at mobile devices; lower power consumption leads to longer battery life. Over the past few years, we have seen increased interest in power management from data centers, too. There, lowering power consumption has a direct affect on the cost of power and cooling. The utility of power savings is not restricted to costs though, as it will positively affect up-time during power outages, too.

We did some research into power usage during Fedora Core 5 development, to find out exactly how good/bad a job we were doing at being idle. To this end, I bought a ‘kill-a-watt’ [11] device (and later borrowed a ‘Watts-up’ [12] which allowed serial logging). The results showed that a completely idle EM64T box (Dell Precision 470) sucked a whopping 153 Watts of power. At its peak, doing a kernel compile, it pulled 258W, over five times as much power as its LCD display. By comparison, a VIA C3 Nehemiah system pulled 48 Watts whilst idle. The lowest power usage I measured on modern hardware was 21W idle on a mobile-Athlon-based Compaq laptop.

Whilst vastly lower than the more heavyweight systems, it was still higher than I had anticipated, so I investigated further as to where the power was being used. For some time, people have been proclaiming the usefulness of the ‘dynamic tick’ patch for the kernel, which stops the kernel waking up at a regular interval to check if any timers have expired, instead idling until the next timer in the system expires.

Without the patch, the Athlon XP laptop idled at around 21W. With dynticks, after settling down for about a minute, the idle routine auto-calibrates itself and starts putting off delays. Suddenly, the power meter started registering. . . 20,21,19,20,19,20,18,21,19,20,22 changing about once a second. Given the overall average power use went down below its regular idle power use, the patch does seem like a win. (For reference, Windows XP does not do any tricks similar to the results of the dynticks patch, and idles at 20W on the same hardware). Clearly the goal is to spend longer in the lower states, by not waking up so often.

Another useful side-effect of the dyntick patch was that it provides a /proc file that allows you to monitor which timers are firing, and their frequency. Watching this revealed a number of surprises. Figure 3 shows the output of this file (The actual output is slightly different, I munged it to include the symbol name of the timer function being called.)

- Kernel problems Whilst this paper focuses on userspace issues, for completeness, I will also enumerate the kernel issues that this profiling highlighted.
 - USB. Every 256ms, a timer was firing in the USB code. Apparently the USB 2.0 spec mandates this timer, but if there are no USB devices connected (as was the case when I measured), it does call into the question

peer_check_expire	181	crond
dst_run_gc	194	syslogd
rt_check_expire	251	auditd
process_timeout	334	hald
it_real_fn	410	automount
process_timeout	437	kjournald
process_timeout	1260	
it_real_fn	1564	rpc.idmapd
commit_timeout	1574	
wb_timer_fn	1615	init
process_timeout	1652	sendmail
process_timeout	1653	
process_timeout	1833	
neigh_periodic_timer	1931	
process_timeout	2218	hald-addon-stor
process_timeout	3492	cpuspeed
delayed_work_timer_fn	4447	
process_timeout	7620	watchdog/0
it_real_fn	7965	Xorg
process_timeout	13269	gdmgreeter
process_timeout	15607	python
cursor_timer_handler	34096	
i8042_timer_func	35437	
rh_timer_func	52912	

Figure 3: `/proc/timertop`

what exactly it is doing. For the purposes of testing, I worked around this with a big hammer, and `rmmod`'d the USB drivers.

- keyboard controller. At HZ/20, the `i8042` code polls the keyboard controller to see if someone has hot-plugged a keyboard/mouse or not.
- Cursor blinking. hilariously, at HZ/5 we wake up to blink the cursor. (Even if we are running X, and not sat at a VT)
- gdm
 - For some reason, `gdm` keeps getting scheduled to do work, even when it is not the active tty.
- Xorg

- X is hitting `it_real_fn` a *lot* even if it is not the currently active VT. Ironically, this is due to X using its 'smart scheduler', which hits `SIGALRM` regularly, to punish X clients that are hogging the server. Running X with `-dumsched` made this completely disappear. At the time it was implemented, `itimer` was considered the fastest way of getting a timer out of the kernel. With advances from recent years speeding up `gettimeofday()` through the use of `vsyscalls`, this may no longer be the most optimal way for it to go about things.

- python
 - The python process that kept waking up belonged to `hpssd.py`, a part of `hplip`. As I do not have a printer,

this was completely unnecessary.

By removing the unneeded services and kernel modules, power usage dropped another watt. Not a huge amount, but significant enough to be measured.

Work is continuing in this area for Fedora Core 6 development, including providing better tools to understand the huge amount of data available. Current gnome-power-manager CVS even has features to monitor `/proc/acpi` files over time to produce easy-to-parse graphs.

5 Conclusions.

The performance issues discussed in this paper are not typically reported by users. Or, if they are, the reports lack sufficient information to root-cause the problem. That is why it is important to continue to develop tools such as those outlined in this paper, and to run these tools against the code base moving forward. The work is far from over; rather, it is a continual effort that should be engaged in by all involved parties.

With increased interest in power management, not only for mobile devices, but for desktops and servers too, a lot more attention needs to be paid to applications to ensure they are “optimised for power.” Further development of monitoring tools such as the current gnome-power-manager work is key to understanding where the problem areas are.

Whilst this paper pointed out a number of specific problems, the key message to be conveyed is that the underlying problem does not lie with any specific package. The problem is that developers need to be aware of the tools that are available, and be informed of the new tools being developed. It is through the use of these tools that we can make Linux not suck.

References

- [1] <http://people.redhat.com/davej/filemon>
- [2] <http://www.bootchart.org>
- [3] <http://www.gnome.org/~mortenw/files/strace-account>
- [4] <http://people.redhat.com/berrange/systemtap/bootprobe>
- [5] <http://valgrind.org>
- [6] <http://people.redhat.com/drepper/defprogramming.pdf>
- [7] <http://www.splint.org>
- [8] <http://www.codemonkey.org.uk/projects/git-snapshots/sparse>
- [9] <http://lists.freedesktop.org/archives/xorg/2006-March/013992.html>
http://blogs.sun.com/roller/page/alanc?entry=security_hole_in_xorg_6
- [10] <http://www.coverity.com>
- [11] **kill-a-watt:**
<http://www.thinkgeek.com/gadgets/electronic/7657>
- [12] **watts-up:**
<http://www.powermeterstore.com/plugin/wattsup.php>

Proceedings of the Linux Symposium

Volume One

July 19th–22nd, 2006
Ottawa, Ontario
Canada

Conference Organizers

Andrew J. Hutton, *Steamballoon, Inc.*
C. Craig Ross, *Linux Symposium*

Review Committee

Jeff Garzik, *Red Hat Software*
Gerrit Huizenga, *IBM*
Dave Jones, *Red Hat Software*
Ben LaHaise, *Intel Corporation*
Matt Mackall, *Selenic Consulting*
Patrick Mochel, *Intel Corporation*
C. Craig Ross, *Linux Symposium*
Andrew Hutton, *Steamballoon, Inc.*

Proceedings Formatting Team

John W. Lockhart, *Red Hat, Inc.*
David M. Fellows, *Fellows and Carr, Inc.*
Kyle McMartin

Authors retain copyright to all submitted papers, but have granted unlimited redistribution rights to all as a condition of submission.