# Linux Instrumentation

Ian Munsie

June 24, 2010

## Outline

1. Overview of available Instrumentation Tools
   - Competition's Instrumentation Tools
   - Overview of Linux Instrumentation Tools
   - Interactions Between Instrumentation Tools
   - Interesting Instrumentation Tools

2. Demonstrations
   - Finding cache misses with perf
   - Locate sources of block I/O with tracepoints
   - Using kprobes to analyse a running kernel

# Outline

What major players are in this area? DTrace

What major players are in this area? DTrace

## DTrace

- All purpose kernel and userspace tracing

- "Better than Linux"

- Originated from Solaris

- Mainstream ports for FreeBSD, NetBSD, Mac OS X

- Licencing Issues with Linux

- No performance impact when probes disabled

- Minimal probe effect

### Definition

Probe Effect: The phenomena where observing a system will change the behaviour of that system

## DTrace

- All purpose kernel and userspace tracing
- "Better than Linux"
- Originated from Solaris
- Mainstream ports for FreeBSD, NetBSD, Mac OS X
- Licencing Issues with Linux
- No performance impact when probes disabled
- Minimal probe effect

### Definition

Probe Effect: The phenomena where observing a system will change the behaviour of that system

## DTrace

- All purpose kernel and userspace tracing
- "Better than Linux"
- Originated from Solaris
- Mainstream ports for FreeBSD, NetBSD, Mac OS X
- Licencing Issues with Linux
- No performance impact when probes disabled
- Minimal probe effect

### Definition

Probe Effect: The phenomena where observing a system will change the behaviour of that system

## DTrace

- All purpose kernel and userspace tracing
- "Better than Linux"
- Originated from Solaris
- Mainstream ports for FreeBSD, NetBSD, Mac OS X
- Licencing Issues with Linux
- No performance impact when probes disabled
- Minimal probe effect

### Definition

Probe Effect: The phenomena where observing a system will change the behaviour of that system

## DTrace Implementation

- Well documented probes

### Definition

Probe: A location or action which can be hooked into in order to perform some arbitrary action

- Probes placed statically in source code
- Probe location well considered
- Probes in Solaris kernel, postgreSQL, x.org, . . .
- Can create new probes dynamically
- "Dynamic tracing framework"
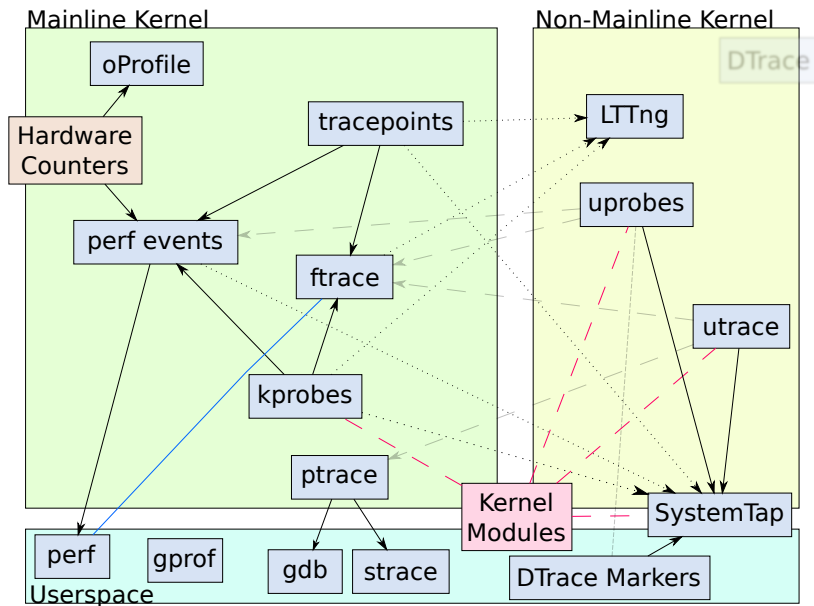- DTrace D programming language

# DTrace Implementation

- Well documented probes

### Definition

Probe: A location or action which can be hooked into in order to perform some arbitrary action

- Probes placed statically in source code
- Probe location well considered
- Probes in Solaris kernel, postgreSQL, x.org, ...
- Can create new probes dynamically
- "Dynamic tracing framework"
- DTrace D programming language

## DTrace Implementation

- Well documented probes

### Definition

Probe: A location or action which can be hooked into in order to perform some arbitrary action

- Probes placed statically in source code
- Probe location well considered
- Probes in Solaris kernel, postgreSQL, x.org, . . .
- Can create new probes dynamically
- "Dynamic tracing framework"
- DTrace D programming language

## DTrace on Linux?

- DTrace approach generally unsuitable for Linux kernel
- Linux kernel evolves rapidly
- Resistance to placement of static probes
- Unofficial Linux port does exist
- SystemTap can use DTrace markers

## DTrace on Linux?

- DTrace approach generally unsuitable for Linux kernel
- Linux kernel evolves rapidly
- Resistance to placement of static probes
- Unofficial Linux port does exist
- SystemTap can use DTrace markers

Enough DTrace envy, what about Linux?

Static probes? tracepoints

Static probes? tracepoints

# Tracepoints

- Static probe points in the kernel source
- Replace old Kernel Markers
- Infrastructure mainline
- Low performance impact
- TRACE_EVENT() macro

Kernel dynamic probes? kprobes

Kernel dynamic probes? kprobes

## Kprobes

- Insert a new probe into kernel at runtime
- Does not do userspace probes
- Heavily architecture specific
- Can be used by loadable kernel modules
- Used by other instrumentation tools
- Provides three types of probes
  - kprobes—Probe instruction
  - jprobes—Probe function call
  - kretprobes—Probe function return
- Mainline

Userspace dynamic probes? uprobes
Not mainline yet.

Userspace dynamic probes? uprobes

Not mainline yet.

Userspace dynamic probes? uprobes
Not mainline yet.

## Uprobes

- Kprobes for userspace
- Not mainline yet—predict that may soon change
- Used by SystemTap
- Handlers run in task context
- Background page replacement mechanism
- Probed instruction single stepped in XOL area
- No longer depends on utrace
- Handlers implemented as kernel module
- perf interface

Userspace profiling? gprof

Userspace profiling? gprof

## gprof

- Legacy, but deserves a mention
- Counts and profiles calls to functions
- Function granuality
- Compile executable with -pg flag
- Produce function call graph
- Does NOT profile libraries or kernel
- Kprof provides GUI to visualise call graph

# gprof Flat Profile

```
uxterm
Flat profile:

Each sample counts as 0.01 seconds.
  %   cumulative   self              self     total
 time   seconds   seconds    calls   s/call   s/call  name
36.07      4.79      4.79   247575     0.00     0.00  fastscale
15.51      6.85      2.06     2034     0.00     0.00  disp3d
10.69      8.27      1.42     1418     0.00     0.00  draw_plasma
 8.81      9.44      1.17      315     0.00     0.00  do_julia
 6.85     10.35      0.91  1183665     0.00     0.00  mand_calc
 3.09     10.76      0.41     6186     0.00     0.00  move_starfield
 2.94     11.15      0.39     5072     0.00     0.00  mkrealloc_table
 1.96     11.41      0.26      359     0.00     0.00  mydraw1
 1.51     11.61      0.20      199     0.00     0.00  toblack1
 1.43     11.80      0.19  8712256     0.00     0.00  bbupdate
 1.36     11.98      0.18      378     0.00     0.00  mydraw
 1.36     12.16      0.18      176     0.00     0.00  drawfire
 1.28     12.33      0.17    51536     0.00     0.00  drawline
 0.98     12.46      0.13 17414911     0.00     0.00  tl_process_group
 0.98     12.59      0.13     2187     0.00     0.00  bbwait
 0.90     12.71      0.12     2536     0.00     0.00  do_fractal
 0.75     12.81      0.10 26164967     0.00     0.00  tl_update_time
 0.75     12.91      0.10  8650042     0.00     0.00  tl_sleep
 0.68     13.00      0.09     6178     0.00     0.00  draw_starfield
 0.45     13.06      0.06      615     0.00     0.00  fastcscale
 0.45     13.12      0.06       39     0.00     0.00  morphdraw
 0.23     13.15      0.03  8752246     0.00     0.00  tl_lookup_timer
:
```

# gprof Call Graph

# kprof Profile Visualisation Tool

Userspace tracing? ptrace ...
A bit ugly.

Userspace tracing? ptrace . . .
A bit ugly.

Userspace tracing? ptrace . . .
A bit ugly.

## ptrace

- Means of one process observing/controlling another
- Prominently used by gdb, strace, ltrace
- Ugly, limited interface
- Processes cannot be traced by multiple processes
- Signal oriented architecture
- Large overhead

# ptrace—ltrace

# ptrace—strace



```
\0"..., 512) = 512
fstat64(3, {st_mode=S_IFREG|0755, st_size=1331684, ...}) = 0
mmap2(NULL, 1337704, PROT_READ|PROT_EXEC, MAP_PRIVATE|MAP_DENYWRITE, 3,
0) = 0xf7605000
mmap2(0xf7746000, 12288, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_FIXED|MAP
_DENYWRITE, 3, 0x141) = 0xf7746000
mmap2(0xf7749000, 10600, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_FIXED|MAP
_ANONYMOUS, -1, 0) = 0xf7749000
close(3)                                = 0
mmap2(NULL, 4096, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_ANONYMOUS, -1, 0
) = 0xf7604000
set_thread_area({entry_number:-1 -> 12, base_addr:0xf76046c0, limit:1048
575, seg_32bit:1, contents:0, read_exec_only:0, limit_in_pages:1, seg_no
t_present:0, useable:1}) = 0
mprotect(0xf7746000, 8192, PROT_READ)   = 0
mprotect(0xf777d000, 4096, PROT_READ)   = 0
munmap(0xf774c000, 69737)               = 0
fstat64(1, {st_mode=S_IFIFO|0600, st_size=0, ...}) = 0
mmap2(NULL, 4096, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_ANONYMOUS, -1, 0
) = 0xf775d000
write(1, "Initialising array...\n", 22Initialising array...
) = 22
write(1, "Trying approach 1: ", 19Trying approach 1: )      = 19
write(1, "0x18fff7eb26fd058\nTrying approac"..., 370x18fff7eb26fd058
Trying approach 2: ) = 37
write(1, "0x18fff7eb26fd058\n", 180x18fff7eb26fd058
)      = 18
exit_group(0)                           = ?
(END)
```

Improved userspace tracing? utrace
Not mainline yet.

Improved userspace tracing? utrace

Not mainline yet.

Improved userspace tracing? utrace
Not mainline yet.

## utrace

- Improved userspace tracing
- Not mainline
- Prominently used by SystemTap
- Utrace clients run in kernel, typically in module
- Ptrace reimplemented as a utrace client
- Infrastructure to monitor threads
- Establish "engine" for each monitored thread
    - Event reporting
    - Thread control
    - Thread machine state access

Kernel tracing? ftrace

Kernel tracing? ftrace

## ftrace

- Mainline
- Kernel tracing infrastructure
- Originated from the realtime efforts
- Hijacks mcount (from gprof) for own purposes
- Low performance impact
- Tracing Plugins
    - Function
    - Funtion graph
    - Context switches
    - Time interrupts disabled
    - Time preemtion disabled
    - Delay for high priority tasks
    - Power state transitions
    - Branch prediction
    - . . .
- Exposed through debugfs—Manipulate with echo & cat
- Gained event support—perf events favoured for this

# ftrace—Function Tracer



```
delenn# mount -t debugfs debugfs /sys/kernel/debug                              ~
delenn# echo 'ext3*' > /sys/kernel/debug/tracing/set_ftrace_filter            ~
delenn# echo function > /sys/kernel/debug/tracing/current_tracer               ~
delenn# echo 1 > /sys/kernel/debug/tracing/tracing_on                          ~
delenn# head /sys/kernel/debug/tracing/trace                                   ~
# tracer: function
#
#           TASK-PID    CPU#    TIMESTAMP  FUNCTION
#              | |        |          |         |
       kjournald-1309   [000]  1362.725446: ext3_bmap <-bmap
       kjournald-1309   [000]  1362.725449: ext3_get_block <-generic_block_bmap
       kjournald-1309   [000]  1362.725450: ext3_get_blocks_handle <-ext3_get_block
       kjournald-1309   [000]  1362.725451: ext3_block_to_path <-ext3_get_blocks_hand
le
       kjournald-1309   [000]  1362.725452: ext3_get_branch <-ext3_get_blocks_handle
       kjournald-1309   [000]  1362.725472: ext3_bmap <-bmap
delenn#                                                                         ~
```

# ftrace—Task scheduling

# ftrace—Measuring task wakeup latency

Kernel and Userspace tracing? LTTng
Not mainline.

Kernel and Userspace tracing? LTTng
Not mainline.

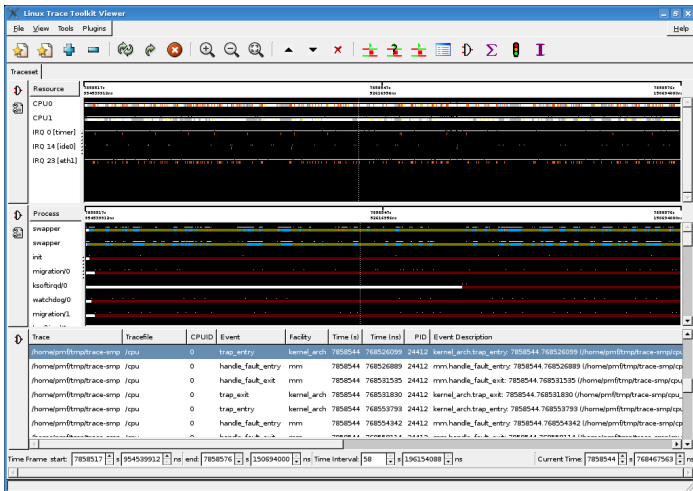Kernel and Userspace tracing? LTTng
Not mainline.

# LTTng

- Kernel & Userspace tracing package
- Not mainline
- Adds instrumentation points into kernel
- Uses tracepoints, ftrace, kprobes
- Can plot traces
- Completely failed to work for me in every way

# LTTng

- Kernel & Userspace tracing package
- Not mainline
- Adds instrumentation points into kernel
- Uses tracepoints, ftrace, kprobes
- Can plot traces
- Completely failed to work for me in every way

# LTTng Screenshot

LTTng did not work for me, but here's what it's supposed to do:

Hardware performance counters? oProfile

Hardware performance counters? oProfile

## oProfile

- System wide sampling profiler
- Mainline
- Interface to Performance Measurement Unit
- GUI eases selection of performance counters
- Can generate gprof style call graph

# oProfile GUI

# oProfile Report

# oProfile Annotate

# oProfile Call Graph



```
uxterm
delenn% opreport -c                                          ~/tests
CPU: Core 2, speed 1200 MHz (estimated)
Counted LLC_MISSES events (L2 cache demand requests from this core that missed t
he L2) with a unit mask of 0x41 (No unit mask) count 400000
samples   %          app name                    symbol name
-------------------------------------------------------------------------------
55        96.4912  cachetest                    sumArrayNaive
  55      100.000  cachetest                     sumArrayNaive [self]
-------------------------------------------------------------------------------
  1       100.000  vmlinux                      run_local_timers
1         1.7544   vmlinux                      raise_softirq
  1       100.000  vmlinux                       raise_softirq [self]
-------------------------------------------------------------------------------
  1       100.000  vmlinux                      mangle
1         1.7544   vmlinux                      seq_escape
  1       100.000  vmlinux                       seq_escape [self]
-------------------------------------------------------------------------------
  1       100.000  vmlinux                      tick_sched_timer
0            0     vmlinux                      __run_hrtimer
  1       100.000  vmlinux                      update_process_times
  0            0   vmlinux                       __run_hrtimer [self]
-------------------------------------------------------------------------------
  1       33.3333  vmlinux                      cpuidle_idle_call
  2       66.6667  vmlinux                      acpi_idle_enter_bm
0            0     vmlinux                      acpi_idle_enter_bm
  2       66.6667  vmlinux                      acpi_idle_enter_bm
  1       33.3333  vmlinux                      cpuidle_idle_call
  0            0   vmlinux                       acpi_idle_enter_bm [self]
-------------------------------------------------------------------------------
  1       100.000  vmlinux                      cpuidle_idle_call
0            0     vmlinux                      apic_timer_interrupt
  1       100.000  vmlinux                       smp_apic_timer_interrupt
```

Something to pull everything together? SystemTap
Not mainline.

Something to pull everything together? SystemTap
Not mainline.

Something to pull everything together? SystemTap
Not mainline.

## SystemTap

- Powerful kernel & userspace analysis suite
- Scripting language
- Generates, Compiles and loads kernel modules
- Primarily uses kprobes
- Utrace for userspace tracing
- Uprobes support
- Many userspace dependencies
- Hard to use
- More of interest to enterprise distros
- Version 1.2 recently released
    - Prototypical support for perf events
    - Hardware breakpoint probe support
- Not mainline and . . .

## SystemTap

*Let's face it, system tap isn't going to be merged, so why even bring it up? Every kernel developer I have _ever_ seen agrees that all the new tracing is a million times superior. I'm sure there are system tap people who disagree, but quite frankly, I don't see it being merged considering how little the system tap people ever did for the kernel.*

*So if things like system tap and "security models that go behind the kernel by tying into utrace" are the reasons for utrace, color me utterly uninterested. In fact, color me actively hostile. I think that's the worst possible situation that we'd ever be in as kernel people (namely exactly the "do things in kernel space by hiding behind utrace without having kernel people involved")*

*Linus*

Something else to pull everything together? perf events

The future of Linux Instrumentation.

Something else to pull everything together? perf events

The future of Linux Instrumentation.

Something else to pull everything together? perf events
The future of Linux Instrumentation.

## perf events

- All in one Performance tool
- Newcomer, but already mainline
- Interface to Performance Measurement Unit
- Virtual counters
- Integrates with tracepoints
- Kprobe support
- Scripting support
- SLAB subsystem profiling
- Lock profiling
- Scheduler analyser
- Timechart

# Sample perf report

# Sample perf call graph

# perf—Available events



```
root@Ego2: ~
List of pre-defined events (to be used in -e):

  cpu-cycles OR cycles                          [Hardware event]
  instructions                                  [Hardware event]
  cache-references                              [Hardware event]
  cache-misses                                  [Hardware event]
  branch-instructions OR branches               [Hardware event]
  branch-misses                                 [Hardware event]
  bus-cycles                                    [Hardware event]

  cpu-clock                                     [Software event]
  task-clock                                    [Software event]
  page-faults OR faults                         [Software event]
  minor-faults                                  [Software event]
  major-faults                                  [Software event]
  context-switches OR cs                        [Software event]
  cpu-migrations OR migrations                  [Software event]
  alignment-faults                              [Software event]
  emulation-faults                              [Software event]

  L1-dcache-loads                               [Hardware cache event]
  L1-dcache-load-misses                         [Hardware cache event]
  L1-dcache-stores                              [Hardware cache event]
  L1-dcache-store-misses                        [Hardware cache event]
  L1-dcache-prefetches                          [Hardware cache event]
  L1-dcache-prefetch-misses                     [Hardware cache event]
  L1-icache-loads                               [Hardware cache event]
  L1-icache-load-misses                         [Hardware cache event]
:
```

# perf—Available events



```
root@Ego2: ~
 dTLB-prefetch-misses                            [Hardware cache event]
 iTLB-loads                                      [Hardware cache event]
 iTLB-load-misses                                [Hardware cache event]
 branch-loads                                    [Hardware cache event]
 branch-load-misses                              [Hardware cache event]

 rNNN (see 'perf list --help' on how to encode it) [Raw hardware event

 mem:<addr>[:access]                             [Hardware breakpoint]

 skb:skb_copy_datagram_iovec                     [Tracepoint event]
 skb:kfree_skb                                   [Tracepoint event]
 bkl:unlock_kernel                               [Tracepoint event]
 bkl:lock_kernel                                 [Tracepoint event]
 block:block_rq_remap                            [Tracepoint event]
 block:block_remap                               [Tracepoint event]
 block:block_split                               [Tracepoint event]
 block:block_unplug_io                           [Tracepoint event]
 block:block_unplug_timer                        [Tracepoint event]
 block:block_plug                                [Tracepoint event]
 block:block_sleeprq                             [Tracepoint event]
 block:block_getrq                               [Tracepoint event]
 block:block_bio_queue                           [Tracepoint event]
 block:block_bio_frontmerge                      [Tracepoint event]
 block:block_bio_backmerge                       [Tracepoint event]
 block:block_bio_complete                        [Tracepoint event]
 block:block_bio_bounce                          [Tracepoint event]
 block:block_rq_issue                            [Tracepoint event]
:
```

# perf top

# perf timechart

# perf timechart

How does it all fit together?

# Interesting Technologies for the Future

- DTrace
- SystemTap
- LTTng
- gprof
- ptrace
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap
- LTTng
- gprof
- ptrace
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap
- LTTng
- gprof
- ptrace
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng
- gprof
- ptrace
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng
- gprof
- ptrace
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof
- ptrace
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof
- ptrace
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events
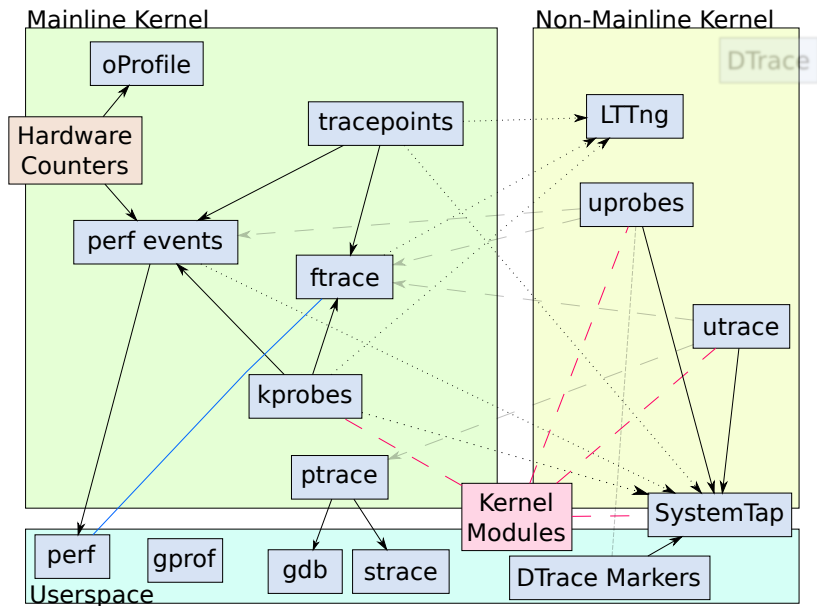
# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace—Trouble going mainline
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace—Trouble going mainline
- oProfile
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace—Trouble going mainline
- oProfile—Not going anywhere
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace—Trouble going mainline
- oProfile—Not going anywhere
- uprobes
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace—Trouble going mainline
- oProfile—Not going anywhere
- uprobes—Not mainline yet, likely soon...
- kprobes
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace—Trouble going mainline
- oProfile—Not going anywhere
- uprobes—Not mainline yet, likely soon. . .
- kprobes—Mainline, Used by others
- tracepoints
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace—Trouble going mainline
- oProfile—Not going anywhere
- uprobes—Not mainline yet, likely soon. . .
- kprobes—Mainline, Used by others
- tracepoints—Mainline, Used by others
- ftrace
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace—Trouble going mainline
- oProfile—Not going anywhere
- uprobes—Not mainline yet, likely soon. . .
- kprobes—Mainline, Used by others
- tracepoints—Mainline, Used by others
- ftrace—Mainline, Active
- perf events

# Interesting Technologies for the Future

- DTrace—Licence issues, not suited to Linux
- SystemTap—Not mainline, enterprise focus
- LTTng—Not mainline, enterprise focus
- gprof—Limitations, not going anywhere
- ptrace—Ugly, set in stone
- utrace—Trouble going mainline
- oProfile—Not going anywhere
- uprobes—Not mainline yet, likely soon. . .
- kprobes—Mainline, Used by others
- tracepoints—Mainline, Used by others
- ftrace—Mainline, Active
- perf events—Mainline, Very active

# Outline

2. Demonstrations
   - Finding cache misses with perf
   - Locate sources of block I/O with tracepoints
   - Using kprobes to analyse a running kernel

Finding cache misses with perf

# Finding cache misses with perf

- Fairly common scenario
- High cache misses is bad for performance
- Use a very simple C program as an example
- Task: sum every element in a large array
- Complication: array & program too big for L2 cache
- Two approaches with different array access order
- Use perf to observe cache misses

## Finding cache misses with perf

- Fairly common scenario
- High cache misses is bad for performance
- Use a very simple C program as an example
- Task: sum every element in a large array
- Complication: array & program too big for L2 cache
- Two approaches with different array access order
- Use perf to observe cache misses

# Finding cache misses with perf

# Finding cache misses with perf

```
cachetest.c (~/eeeconfig/tests) - VIM
#include <stdio.h>
#include <stdlib.h>

//4 MiB of L2 cache with a 4MiB array...
#define WIDTH 1024
#define HEIGHT 1024
#define ITERATIONS 100
int bigArray[WIDTH*HEIGHT];

void initArray(int *array, int w, int h) {
        for (int n = 0; n < w*h; n++)
                array[n] = rand();
}
unsigned long long sumArrayNaive(int *array, int w, int h, int n) {
        unsigned long long ret = 0;
        for (int i = 0; i < n; i++)
                for (int x = 0; x < w; x++)
                        for (int y = 0; y < h; y++)
                                ret += array[y*w + x];
        return ret;
}
unsigned long long sumArrayOptimal(int *array, int w, int h, int n) {
        unsigned long long ret = 0;
        for (int i = 0; i < n; i++)
                for (int y = 0; y < h; y++)
                        for (int x = 0; x < w; x++)
                                ret += array[y*w + x];
        return ret;
}
```

# Finding cache misses with perf

$ perf record -e cache-misses ./cachetest
$ perf report

# Finding cache misses with perf

$ perf record -e cache-misses ./cachetest
$ perf report

Locate sources of block I/O with perf and tracepoints

## perf—Locate sources of block I/O

- System is under heavy disk I/O causing slowdowns
- Traditional top does not reveal offending processes
- Even atop only shows I/O at process granularity
- block:block_rq_issue tracepoint available
- Kernel must be compiled with block I/O tracing
- perf can produce call graphs to track exact origin

# perf—Locate sources of block I/O

# perf record -g -a -e block:block_rq_issue sleep 10
# perf report

# perf—Locate sources of block I/O

# perf record -g -a -e block:block_rq_issue sleep 10
# perf report

Analyse a running kernel with perf and kprobes

# Analyse a running kernel with perf and kprobes

- Want to pull information out of a running kernel
- Information is not already revealed to userspace
- No tracepoints provide the information
- Live system—cannot afford downtime
- Kernel debugging information and source is available
- Let's probe the scheduler as an example. . .

# Analysing a running kernel with perf and kprobes

# perf probe –line schedule

# Analysing a running kernel with perf and kprobes

# Analysing a running kernel with perf and kprobes

# perf report

Thankyou all for listening

- perf can be found in the Linux kernel sources under /tools/perf
- perf is being actively developed in the "tip" kernel tree
- Linux Weekly News has some good further reading on ftrace—google *site:lwn.net ftrace*
- perf articles not as easy to find yet, here are some:
  - http://lwn.net/Articles/339361—"Perfcounters added to the mainline"
  - http://lwn.net/Articles/373842—"Scripting support for perf"
  - http://lwn.net/Articles/382554—"A "live mode" for perf"
  - http://lwn.net/Articles/346470—"Fun with tracepoints"
  - http://lkml.org/lkml/2010/3/17/91—Ingo Molnar on database I/O latency