# Why DTrace Should Be In Your Toolkit

## Ruud van der Pas

**Director, OpenMP Architecture Review Board**

**Multicore World 2026**

**Christchurch, New Zealand, February 16-20, 2026**

# *Outline*

- *Two Questions*

- *What is DTrace?*

- *Some DTrace Concepts*

- *Scratching the DTrace Surface*

- *Q & (some) A*

*Why DTrace Should Be In Your Toolkit*

# Two Questions

*Why DTrace Should Be In Your Toolkit*

# What's Going On?

- **Serious infrastructure outages**

  ✳ **Large scale, e.g. Crowdstrike**

  ✳ **Very high impact, e.g. an airline check-in system that is down**

- **In all such cases, the outage is very disruptive and expensive**

  ✳ **How much attention is given to prevent these outages**

  ✳ **Are there any tools to help?**

*Why DTrace Should Be In Your Toolkit*

# Conventional Troubleshooting Tools

- **Quite a number of troubleshooting tools exist**

  ✳ **strace - trace system calls**

  ✳ **ftrace - kernel tracing**

  ✳ **netstat - network diagnostics**

  ✳ **iostat - I/O diagnostics**

  ✳ **...**

- **Great tools, but a restricted focus and no customization**

*Why DTrace Should Be In Your Toolkit*

# *What is DTrace?*

# The Origin of DTrace

**"If a complicated system is doing something unexpected, how do you determine what it is doing and why, without taking the system down?"**

## This lead to the development of DTrace

*Why DTrace Should Be In Your Toolkit*

# My One Line Description of DTrace

**"DTrace is your microscope into the Operating System"**

*Why DTrace Should Be In Your Toolkit*

# What is DTrace?

- **A dynamic tracing tool**
  - ✳ **But really more than anything else, a programming language**
- **Traces activities in both userland and the kernel**
  - ✳ **And both at the same time**
  - ✳ **Provides a truly holistic view**
- **Designed to be safe to use on a production system**

*Why DTrace Should Be In Your Toolkit*

*Copyright (©) 2026 by Ruud van der Pas*

# How does it Work?

- **The user writes a program/script in the D language**

  ✳ **Don't worry, "D" is like a shell language with C/awk flavors**

- **Full control over what needs to be traced**

  ✳ **Can trace multiple different events simultaneously**

  ✳ **Also the reporting of events is under user control**

- **Provides a high degree of flexibility and customization**

*Why DTrace Should Be In Your Toolkit*

# Some DTrace Concepts

# DTrace Probes

- **Probes are events that are available to trace**

  ✳ *For example, the entry/return to a specific (kernel) function*

- **Every system with DTrace installed has many probes**

  ✳ *Typically in the order of 100,000 or more*

- **Bind actions to any of the probes being traced**

  ✳ *The D language is used to define the action(s) to be taken*

*Why DTrace Should Be In Your Toolkit*

# *The Structure of a Probe*

```
provider:module:function:name
[/ predicate /]
{
   D statements
}
```

*Why DTrace Should Be In Your Toolkit*

# Examples of Probe Definitions

## An example

```
syscall:vmlinux:bind:entry
    {printf("I am in function %s\n",probefunc);}
```

## Use a blank field, a wildcard, and a predicate:

```
syscall::bind*:entry
/ my_flag == 1 /
{
    printf("I am in function %s\n",probefunc);
}
```

*Why DTrace Should Be In Your Toolkit*

# The pid provider: pid$target

- **The examples later on, use the pid provider: pid$target**

- **Variable $target expands to the process id of the target**

- **This means that all the tracing is restricted to the target**

  ✳ **Without this, all system activities are traced**

  ✳ **That is totally fine, of course, but not what we'll do today**

*Why DTrace Should Be In Your Toolkit*

# Aggregations

- **A very powerful data structure in DTrace**

  ✳ **Stores the result of an aggregation function (e.g. min/max)**

  ✳ **Indexed through an arbitrary key (e.g. integer, string, etc.)**

```
@nicolas-headcount[2026,"christchurch"] = count();
```

*this is an aggregation*          *the key*          *aggregation function*

# *A Starting Point to Learn More*

**https://www.youtube.com/@OracleLearning => Oracle Linux DTrace**

*Why DTrace Should Be In Your Toolkit*

# *Scratching the DTrace Surface*

# Monday Morning 9am - Why Me?

*Hey Rudd,*

*Tomorrow we are launching the world leading BlackBox application.*

*For the all important "go/no go" WhatAreWeDoingHere meeting that starts today at 10am, can you give me the I/O characteristics for this application?*

*We need to know the filenames used, their file descriptors and how many bytes are read and written on a per filename basis.*

*Thanks.*

*A. Hole, VP of AGT (Advanced Generic Technology)*

*Why DTrace Should Be In Your Toolkit*

# For the Impatient - The Results

```
============================
The target functions
============================

Function Module      Count
open     libc.so.6      3
pwrite   libc.so.6      2
read     libc.so.6      5


Function Filename           FD Count
open      /tmp/my-code-Bkp7qS    7    1
open      /tmp/my-code-TddgyV    8    1
open      input-blackbox        3    1
```

```
Function     FD   Bytes to read
read          3        4530176


              Bytes Read
Function     Input     Actual
read        4530176   4529552


Function     FD   Bytes written
pwrite        7          4
pwrite        8          4
Total bytes written by pwrite: 8
```

*Why DTrace Should Be In Your Toolkit*

*Copyright (©) 2026 by Ruud van der Pas*

# The Target Application

```
$ ./blackbox -f input-blackbox
$
```

**That's correct. We have no idea what is going on in this code.**

**Let's find out!**

*Why DTrace Should Be In Your Toolkit*

# The Strategy

- **Cast a wide net first: find all potentially relevant functions**
  - ✳ *Leverage the support for empty fields and wildcards*
- **Write and test probes for each piece of information needed**
  - ✳ *It is easy to do this and focus on one thing at a time*
- **Merge the scripts into one overall script**

*Why DTrace Should Be In Your Toolkit*

```
1   pid$target::*read*:entry,
2   pid$target::*write*:entry,
3   pid$target::*open*:entry
4   {
5      @target_calls[probefunc,probemod] = count();
6   }
```

*The probes
(use a blank field for the module
and leverage wildcards)*

*function name*

*module name*

**Trace all functions with read, write, open, or close in the name**

*Why DTrace Should Be In Your Toolkit*

# DTrace in Action

```
$ sudo dtrace -s scan-wide-net.d -c "./blackbox -f input-blackbox"
```

```
dtrace: script './scan-wide-net.d' matched 391 probes
    _IO_file_close                      libc.so.6                   1
    _IO_file_open                       libc.so.6                   1
    _IO_new_fclose                      libc.so.6                   1
    _IO_new_file_close_it               libc.so.6                   1
    _IO_new_file_fopen                  libc.so.6                   1
    __GI___pthread_mutex_lock           libc.so.6                   1
    __GI___pthread_mutex_unlock         libc.so.6                   1
    __GI___pthread_mutex_unlock_usercnt libc.so.6                   1
    __close_nocancel                    libc.so.6                   1
    __fopen_internal                    libc.so.6                   1
    fopen64                             libc.so.6                   1
    omp_set_num_threads                 libgomp.so.1                1
    read_graph_data_from_file           blackbox                    1
    __GI___pthread_key_delete           libc.so.6                   2
    __GI___pthread_setspecific          libc.so.6                   2
    close                               libc.so.6                   2
    pwrite                              libc.so.6                   2
    open                                libc.so.6                   3
    _IO_file_read                       libc.so.6                   5
    read                                libc.so.6                   5
    fread                               libc.so.6                   6
    _dl_audit_objclose                  ld-linux-aarch64.so.1       7
    omp_get_num_threads                 libgomp.so.1              196
    omp_get_thread_num                  libgomp.so.1              196
```

**A long list with functions
Let's select those of interest**

**pwrite, open, read**

**24**

*Why DTrace Should Be In Your Toolkit*

## *From the man pages of open, read, and pwrite:*

```
fd = open(const char *pathname, int flags,/* mode_t mode */ );

bytes_read = read(int fd, void buf[.count], size_t count);

bytes_written = pwrite(int fd, const void buf[.count],

                       size_t count, off_t offset);
```

*Why DTrace Should Be In Your Toolkit*

# *The First Six Lines*

**Every script shown next, starts with these lines:**

```
1   #!/usr/sbin/dtrace –s

2

3   #pragma D option quiet
4   #pragma D option aggsortkey=1
5   #pragma D option aggsortkeypos=0

6
```

*execute script under control of dtrace*

*suppress default output*

*Print aggregations sorted by the first colum*

*Why DTrace Should Be In Your Toolkit*

*Copyright (©) 2026 by Ruud van der Pas*

```
 7 pid$target:libc.so:open:entry,
 8 pid$target:libc.so:read:entry,
 9 pid$target:libc.so:pwrite:entry
10 {
11   @target_calls[probefunc,probemod] = count();
12 }
13 pid$target:libc.so:open:entry
14 {
15   @files_opened[probefunc,stringof(arg0)] = count();
16 }
```

The 3 probes
(no more blank fields and wildcards)

probe for the open() call

arg0 contains the filename

**There are 3 functions that we are interested in**

*Why DTrace Should Be In Your Toolkit*

```
17  END
18  {
19    printf("%s\n","=============================");
20    printf("The target functions\n");
21    printf("%s\n","=============================");
22    printf("%-8s %-10s %5s\n","Function","Module","Count");
23    printa("%-8s %-10s %@5d\n",@target_calls);
24
25    printf("\n%-8s %20s %6s\n","Function","Filename","Count");
26    printa("%-8s %20s %@6d\n",@files_opened);
27  }
```

**Print the results in the END probe (this is almost half of the code)**

*Why DTrace Should Be In Your Toolkit*

# Get the Results!

```
$ sudo scan.d -c "./blackbox -f input-blackbox"
```

```
===========================

The target functions
===========================

Function Module      Count

open     libc.so.6      3

pwrite   libc.so.6      2

read     libc.so.6      5


Function              Filename   Count

open       /tmp/my-code-Famb8X      1

open       /tmp/my-code-fgUczk      1

open         input-blackbox         1
```

**Now we know the filenames in use**

**How about the file descriptors?**

**29**

MulticoreWorldXIII 2026

```
 7 pid$target:libc.so:open:entry
 8 {
 9    self->fname_open = copyinstr(arg0);
10 }
11 pid$target:libc.so:open:return
12 / self->fname_open != 0 /
13 {
14   @files_and_fd_opened[probefunc,
15                        self->fname_open,
16                        arg1] = count();
17   self->fname_open = 0;
18 }
```

**capture the filename**
*(copy from user space to a DTrace buffer)*

**reference the filename**

**arg1 contains the file descriptor**

**In the return probe we link the filename and descriptor**

```
19 END
20 {
21   printf("\nThe mapping between filenames and descriptorss\n");
22   printf("%-8s %-20s %2s %5s\n","Function","Filename","FD","Count");
23   printa("%-8s %-20s %2d %@5d\n",@files_and_fd_opened);
24 }
```

**Print the header and the aggregation**

# Get the Results!

```
$ sudo file-mapping.d -c "./blackbox -f input-blackbox"
```

```
The mapping between filenames and descriptorss
Function Filename                FD Count
open      /tmp/my-code-6mMFMC     8      1
open      /tmp/my-code-y89cUY     7      1
open      input-blackbox          3      1
```

**Now we know the filenames in use and their file descriptors**

*Why DTrace Should Be In Your Toolkit*

*Copyright (©) 2026 by Ruud van der Pas*

```
 7 pid$target:libc.so:read:entry
 8 {
 9    @total_bytes_read[probefunc,arg0] = sum(arg2);
10    @total_bytes_input[probefunc]     = sum(arg2);
11 }
12 pid$target:libc.so:read:return
13 {
14    @total_bytes_actually_read[probefunc] = sum(arg1);
15 }
```

*arg0 contains the file descriptor*

*arg2 contains the number of bytes to be read*

*arg1: the number of bytes actually read*

**This tells us how many bytes have been read**

# *Print the Results in the END Probe*

```
16 END
17 {
18   printf("\n%-8s %5s  %13s\n","Function","FD","Bytes to read");
19   printa("%-8s %5d  %@13d\n",@total_bytes_read);
20
21   printf("\n%23s\n","Bytes Read");
22   printf("%-8s %8s %8s\n","Function","Input","Actual");
23   printa("%-8s %@8d %@8d\n",@total_bytes_input,
24                            @total_bytes_actually_read);
25 }
```

**Both the number of bytes to be read, and delivered, are shown**

*Why DTrace Should Be In Your Toolkit*

# Get the Results!

```
$ sudo read.d -c "./blackbox -f input-blackbox"
```

```
Function       FD   Bytes to read
read            3         4530176


              Bytes Read
Function     Input     Actual
read         4530176   4529552
```

We know how many bytes have been read

But why a difference between the 2 numbers?

*Why DTrace Should Be In Your Toolkit*

# Number of Bytes Written

```
 7 pid$target:libc.so:pwrite:entry
 8 {
 9    @pwrite_bytes[probefunc,arg0]  = sum(arg2);
10    @pwrite_total_bytes[probefunc] = sum(arg2);
11 }
12 END
13 {
14   printf("\n%-8s %5s  %s\n","Function","FD","Bytes written");
15   printa("%-8s %5d  %@8d\n",@pwrite_bytes);
16   printa("Total bytes written by %s: %@d\n",@pwrite_total_bytes);
17 }
```

**arg0 contains the file descriptor**

**arg2 contains the number of bytes to be written**

**This probe is very similar to the one for the read() function**

*Why DTrace Should Be In Your Toolkit*

# Get the Results!

```
$ sudo write.d "./blackbox -f input-blackbox"
```

```
Function        FD   Bytes written

pwrite           7           4

pwrite           8           4

Total bytes written by pwrite: 8
```

**Only 4 bytes written per file**

*Why DTrace Should Be In Your Toolkit*

# *Putting It All Together (plus s... ...leanup)*

```
========================
The target functions
========================

Function Module        Count
open      libc.so.6      3
pwrite    libc.so.6      2
read      libc.so.6      5


Function Filename              FD Count
open     /tmp/my-code-Bkp7qS   7    1
open     /tmp/my-code-TddgyV   8    1
open     input-blackbox        3    1
```

```
Fun...
read

                   Bytes Read
Function     Input      Actual
read        4530176    4529552


Function      FD    Bytes written
pwrite        7             4
pwrite        8             4
Total bytes written by pwrite: 8
```

**Phew! It's done and only 9:55am. Time for a quick coffee :-)**

*Why DTrace Should Be In Your Toolkit*

# Q & (some) A

*Why DTrace Should Be In Your Toolkit*

# Thank You And ... Stay Tuned!