

MCTrace User Manual (Version 1.0)

Galois (pate-darpa-amp@galois.com)

July 11, 2024

DTrace is a popular tool supporting the definition of small snippets of code which are then inserted into an application at specified locations.¹ For example, these code snippets might be inserted directly before and after every call made to a given function, for example a syscall such as `write()`, in order to track the time elapsed each time the function is called. Such probes are defined in the DTrace scripting language and inserted into the source, which is then compiled.

The MCTrace tool is intended to provide the same functionality as DTrace without requiring access to application source code or to the DTrace run-time. For MCTrace, the probes defined as DTrace scripts are compiled into the appropriate machine language for the target architecture and are then inserted into the binary. Note that the binary must be modified both to include the code for the probes themselves, as well as inserting the calls to those probes at specified locations in the original binary code.

Contents

1	Introduction	2
2	MCTrace	2
2.1	Introduction to MCTrace	2
2.2	Concept of Operations	2
2.3	How MCTrace Works	2
2.4	Supported DTrace Language Features	4
2.5	Current Limitations of MCTrace	5
3	Building MCTrace	5
3.1	Development Build Instructions	5
3.2	Release Build Instructions	6
4	User Release	6
4.1	Docker Image Contents	7
4.2	Using MCTrace in this demonstration	7
5	Acknowledgement	9

¹<https://en.wikipedia.org/wiki/DTrace>

1 Introduction

This is the User Manual for MCTrace. The source for this manual can be found within the repository in the directory `docs/user-manual`.

In addition to providing instructions for building a Docker image suitable for running MCTrace in the development environment provided by the software, `README.md` includes instructions for building a standalone Docker image that is suitable for use by a user who is not interested in development or in the source code. Use of the standalone Docker image is discussed in Section 4.

2 MCTrace

This section describes the MCTrace tool, including its execution environment and requirements. The information in this section can also be found in the repository at the top level in the file `MCTRACE.md`.

For instructions on how to build MCTrace, see Section 3.

2.1 Introduction to MCTrace

The MCTrace tool enables users to insert instrumentation into binaries in order to collect fine-grained tracing information. MCTrace functions similarly to DTrace but does not require any operating system support. The input format of MCTrace is a subset of the DTrace probe script language. Prior knowledge of DTrace concepts and terminology is assumed in this document.

2.2 Concept of Operations

After identifying a program to be instrumented (e.g., for debugging purposes or to collect telemetry), the steps to use MCTrace are as follows:

1. Write a DTrace script that collects the desired data.
2. Run the `mctrace` tool in `instrument` mode to insert probes into the binary as directed by the script.
3. Run the instrumented binary.
4. Collect telemetry data emitted by the instrumented binary and use `mctrace` and associated scripts to decode it.

The rest of this section provides more details on the steps above.

2.3 How MCTrace Works

Using MCTrace requires:

- a PowerPC or x86_64 ELF binary to instrument,
- an object file implementing the MCTrace Platform API (see below), and
- a DTrace probe script containing the probes that will be used to modify the provided ELF binary.

MCTrace works by compiling a DTrace probe script into native code using a compiler backend (e.g., LLVM). It then uses binary rewriting to insert the generated probes into the binary at points indicated by the probe script. Through static analysis, it identifies program locations corresponding to DTrace probe providers; at each provider site, it inserts calls to the compiled probes.

Some supported DTrace language features need access to platform-specific functionality such as memory allocation. Since the DTrace code will run within the context of the modified binary rather than an operating system kernel, MCTrace requires some additional code to provide access to such platform-specific features. The MCTrace tool provides the input program with access to this platform-specific functionality by way of an object file of compiled code called the Platform Implementation. The object code that implements the required functions must conform to a set of C function prototypes called the Platform API. A complete implementation of the Platform API must provide implementations of all of the functions the in header file `mctrace/tests/library/include/platform_api.h` provided in the MCTrace GitHub repository (also in `library/include` in the release Docker image) . Once compiled, the platform API implementation must be provided to `mctrace` as the `--library` argument.

For demonstration purposes, simple platform API implementations for PowerPC and x86_64 Linux user-space are available in the repository in `mctrace/tests/library/` and in the standalone Docker release image in `examples/library/` in both source and object code forms.

The following example demonstrates how an x86_64 binary `foo` would be instrumented and how its produced telemetry would be obtained. In this example, the user wants to count the number of times that the program performs a file read by measuring the number of calls to the `read` system call. The probe script to accomplish this is as follows:

```
int read_calls;

::read:entry {
    read_calls = read_calls + 1;
    send(0);
}
```

This probe increments the `read_calls` DTrace variable just prior to each call to the `read` system call. Furthermore, the complete set of DTrace global variables (containing only `read_calls` in this case) is then transmitted as telemetry according to the platform implementation of `send()`.

Once the probes are written, `mctrace` would be invoked as follows.

- The `foo` binary and the `probes.d` probe script are provided to `mctrace` as inputs.
- The `--output` option tells `mctrace` where the instrumented binary should be written.
- The `--var-mapping` option tells `mctrace` where to record metadata that allows it to later decode telemetry data.

- The `--library` option tells `mctrace` where to find the platform implementation object code.
- The `--var-mapping` option tells `mctrace` where to write its metadata describing the encoding of the telemetry data stream.

```
$ mctrace instrument --binary=foo \
    --output=foo.instrumented \
    --var-mapping=foo.mapping.json \
    --library=mctrace/tests/library/platform/_impl.o \
    --script=probes.d
```

The resulting binary can then be run:

```
$ ./foo.instrumented 2>telemetry.bin
```

Once the instrumented binary has finished running, the file `telemetry.bin` will contain the binary telemetry data corresponding to the value of `read.calls` after each invocation of the probe. The telemetry data must then be decoded:

```
$ extractor.py foo.mapping.json --columns < telemetry.bin
```

The `extractor.py` script uses the telemetry mapping file to decode the written telemetry data.

NOTE: In this example, the instrumented binary was run with a `stderr` redirect. This is because the `x86_64` platform implementation for `send()` used in this example writes the DTrace global variable data to `stderr`. In practice, `send()` would trigger some other platform-specific mode of data transmission such as sending packets on a network connection, writing to a local bus, etc.

For full details of the MCTrace tool's command line usage, run `mctrace instrument --help`.

2.4 Supported DTrace Language Features

This information is also available in the repository, in `DTRACE.md`.

MCTrace uses the DTrace language as the means for expressing how it should modify its input binary. While MCTrace does not implement all of the DTrace language, the following core DTrace language features are supported:

- Probe name pattern-matching. Supported metacharacters are `*`, `.`, `[...]`, and `\`
- Global variables
- Constants with suffixes `'u'`, `'U'`, `'l'`, `'L'`, `'ul'`, `'UL'`, `'ll'`, `'LL'`, `'ull'`, and `'ULL'`
- Data types:
 - `'int'`
 - `'char'`
 - `'short'`
 - `'long'`
 - `'long long'`

- Integer arithmetic operators ‘+’, ‘-’, and ‘*’
- Builtins:
 - ‘int arg0’
 - ‘timestamp’
 - ‘long ucallee’

Example DTrace probe scripts demonstrating these features can be found in `mctrace/tests/eval/`.

In addition to the core DTrace language features listed above, MCTrace also supports some additional DTrace constructs specific to MCTrace:

- A `send` action of the form: `send(<numeric channel ID>)`. Invoking `send` will result in an invocation of the `platform_send` function in the Platform API with that channel ID and the global data store.
- A `copyint32(<address>)` subroutine that returns the 32-bit value from the specified location.

2.5 Current Limitations of MCTrace

MCTrace has the following limitations:

- MCTrace supports only statically-linked input binaries.
- `arg0` always returns a 32-bit value (on both PowerPC and x86_64 platforms). Similarly, `copyint32` takes a 32-bit address as its input on both platforms (and returns a 32-bit value). As a result, these work best on the PowerPC 32-bit platform since the argument and return value width match the architecture.
- Platform API implementations are subject to the following restrictions:
- Functions in the Platform API implementation must be self-contained and cannot call other functions even in the same object file.
- Functions in the Platform API (as opposed to DTrace scripts) cannot make use of global variables.

3 Building MCTrace

MCTrace can be built for one of two purposes: either for local development in a Haskell build environment or for release as a Docker image. Instructions for each method are detailed below.

3.1 Development Build Instructions

The development environment setup and build processes are automated. The build process requires Ubuntu 20.04. To perform a one-time setup of the development environment including the installation of LLVM, cross compilers, and other required tools, run the development setup script:

```
./dev_setup.sh
```

Once the development environment is set up and the required tools are installed, MCTrace can be built with the build script:

```
./build.sh
```

After the build has completed, various cross compilers and other tools can be added to the `PATH` environment variable for easier access with:

```
. env.sh
```

To build the example test programs for x86_64 and instrument them using various testing probes, run:

```
make -C mctrace/tests/full
```

To do the same for PowerPC, run:

```
make -C mctrace/tests/full ARCH=PPC
```

MCTRACE can then be run manually:

```
cabal run mctrace <args>
```

3.2 Release Build Instructions

To build the stand-alone release Docker image, execute the following from the root of the repository:

```
cd release
./build.sh
```

This will build two docker images:

- A self-contained image that contains MCTrace, its dependencies, associated tools, and examples. For information on using this image, see Section 4.
- A minimal image containing just MCTrace and its dependencies. A helper script, `release/mctrace` has been provided to run the command in a container. Note that paths passed to this script should be relative to the root of the repository and paths outside of the repository will not be accessible.

4 User Release

This section describes how to run the standalone Docker image for MCTrace described in Section 1. This information is also available in the repository in `release/README.md` or in the standalone Docker image itself in the file `README.md` at the top level.

To load and run the image, run the following commands from a directory containing the image, using the filename `mctrace.tar.gz`.

```
docker image load -i mctrace.tar.gz
docker run -it -w /mctrace-test mctrace
```

This will give a `bash` shell within the Docker container in the directory `mctrace-test` where `mctrace` can be used to instrument binaries. All relative paths mentioned in this section are relative to `mctrace-test`.

4.1 Docker Image Contents

The docker image contains PowerPC and x86_64 test programs and example probes that can be used to exercise MCTrace. Important directories are as follows:

- `examples/eval` contains a collection of probes primarily derived from those provided by WebSensing.
- `examples/full` contains source code and binaries for bundled test programs.
- `examples/binaries` contains binaries from a statically compiled version of GNU `coreutils` for use with `mctrace`.

4.2 Using MCTrace in this demonstration

The `mctrace` tool is in the shell `PATH` when running within Docker so no special steps are needed to be able to invoke the executable within Docker. For instructions on how to run the `mctrace` tool with the appropriate command-line arguments, see `MCTRACE.md` included in the Docker image.

As an example, the `read-write-syscall-PPC.4.inst` binary in this distribution is the instrumented version of the PowerPC binary `read-write-syscall-PPC` and was instrumented with the following `mctrace` command:

```
mctrace instrument --binary=/mctrace-test/examples/full/read-write-syscall-PPC \
  --output=/mctrace-test/examples/full/read-write-syscall-PPC.4.inst \
  --library=/mctrace-test/examples/library/PPC/platform_impl.o \
  --var-mapping=/mctrace-test/examples/full/read-write-syscall-PPC.4.json \
  --script=/mctrace-test/examples/eval/write-timing-probe.d
```

- The `--binary` and the `--script` options tell `mctrace` to instrument the specified binary with the given probe script.
- The `--output` option specifies the name for the instrumented binary.
- The `--library` option specifies the path to the Platform API implementation.
- The `--var-mapping` option tells `mctrace` where to record metadata that allows it to later interpret the collected telemetry.

The above command instruments the binary with probes that triggers at the start and end of the `write` function and computes timing information for the call. Note that the instrumentation command produces a significant amount of debugging log output that can be ignored.

When probes call the DTrace `send` action, the test implementation of `send` pushes the set of telemetry variables, in a compact binary format, to the standard error file descriptor. A script `extractor.py` has been included with the image to interpret this data.

To invoke the instrumented binary and use the `extractor.py` script to decode any emitted telemetry:

```
/mctrace-test/examples/full/read-write-syscall-PPC.4.inst 2>&1
>/dev/null | \
    extractor.py /mctrace-test/examples/full/read-write-syscall-PPC.4.json \
    --extract --big-endian
```

This produces output similar to the following:

```
{"write_count":1,"write_elapsed":162240,"write_ts":1681222607714740774}
{"write_count":2,"write_elapsed":1740,"write_ts":1681222607714800756}
{"write_count":3,"write_elapsed":1309,"write_ts":1681222607714803400}
{"write_count":4,"write_elapsed":1344,"write_ts":1681222607714805742}
{"write_count":5,"write_elapsed":1228,"write_ts":1681222607714807992}
{"write_count":6,"write_elapsed":1222,"write_ts":1681222607714810214}
{"write_count":7,"write_elapsed":1257,"write_ts":1681222607714812555}
```

Note that `2>&1 >/dev/null` has the effect of piping the standard error to the next command while suppressing the standard output of the command. This is done because the provided platform API implementations write `send()` data to `stderr` and that data needs to be piped to the extractor script.

When extracting telemetry data from instrumented PowerPC binaries, the flag `--big-endian` must be passed to the extractor script as in the command above. The flag should be elided when working with x86_64 binaries.

The `extractor.py` script offers a few other conveniences when extracting data from instrumented programs; for example, it can produce columnar outputs and filter columns. See `extractor.py --help` for details on these options.

See the file `README.md` for a list of other binaries for PowerPC and x86_64 included in the Docker image as well some example probes that can be used to instrument each binary. Note that many other combinations of example programs and probes can work together; the full list of combinations can be found in `examples/full/Makefile`.

5 Acknowledgement

This material is based upon work supported by the United States Air Force AFRL/SBRK under Contract No. FA8649-21-P-0293, and by the Defense Advanced Research Projects Agency (DARPA) and Naval Information Warfare Center Pacific (NIWC Pacific) under Contract Number N66001-20-C-4027 and 140D0423C0063. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the DARPA, NIWC Pacific, or its Contracting Agent, the U.S. Department of the Interior, Interior Business Center, Acquisition Services Directorate, Division III.

SBIR DATA RIGHTS

Contract No 140D0423C0063

Contractor Name Galois, Inc.

Contractor Address 421 SW Sixth Ave., Suite 300, Portland, OR 97204

Expiration of SBIR Data Protection Period 06/07/2042

The Government's rights to use, modify, reproduce, release, perform, display, or disclose technical data or computer software marked with this legend are restricted during the period shown as provided in paragraph (b)(5) of the Rights in Noncommercial Technical Data and Computer Software-Small Business Innovation Research (SBIR) Program clause contained in the above identified contract. After the expiration date shown above, the Government has perpetual government purpose rights as provided in paragraph (b)(5) of that clause. Any reproduction of technical data, computer software, or portions thereof marked with this legend must also reproduce the markings.