# Software Debugging and Monitoring for Heterogeneous Many-Core Telecom Systems

## Overview of Research Program

Michel Dagenais, Professor
Dept. of Computer and Software Eng.

POLYTECHNIQUE
MONTRÉAL

AFFILIÉE À
L'UNIVERSITÉ DE MONTRÉAL

# Objectives

- Smart clients everywhere, always online cloud-based services accessed transparently.
- Complex heterogeneous many-core systems.
- Crucial need for better tools!

The objective of this project is to provide better tools that can be used in production, under real load, for the debugging, performance analysis and monitoring of distributed heterogeneous many-core systems. We need to extract and collect information with low overhead, minimize the perturbation on the systems being studied, and efficiently analyse online the gigabytes of data produced.

# Projects Evolution

TDMC Feb. 2009-2012
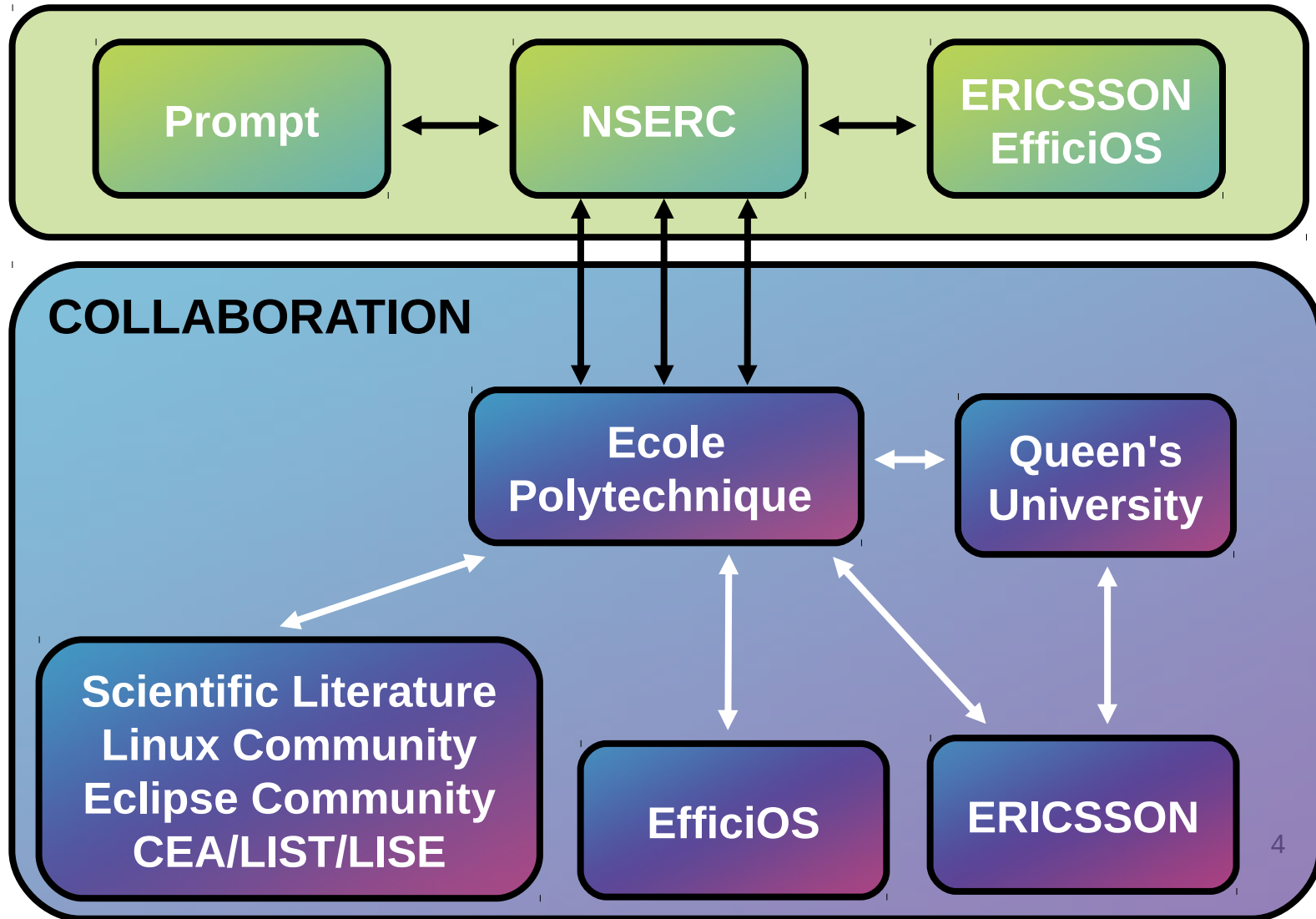
AHLS Feb. 2012-2015

CTPD Sep. 2012-2015

DMHMC Jan. 2015-2018

- TDMC: Ericsson, DRDC, NSERC; static tracepoints for kernel and user-space tracing, multi-session tracing, a posteriori traces synchronization, state system.

- AHLS: Ericsson, DRDC, NSERC; live tracing and synchronization, user-definable state system and views (e.g. Windows), store multiple state subtrees in state system.

- CTPD: Ericsson, EfficiOS, NSERC, Prompt; tracing on Tilera, ARM, Xeon Phi, and GPUs, dynamic tracepoints, VM tracing and view, RPC Debugging, critical path.

- HSDM: Ericsson, EfficiOS, NSERC, Prompt; models (e.g., UML, MARTE) tracing and analysis, Kalray, Adapteva and Intel PT tracing, Cloud and Software Defined Networks tracing, specialized analysis for OpenCL and system resources (I/O, memory, power), tolerating unknown states.

# Structure of Project



Project funding = 4 x industrial contribution

# Modus operandi

- Problems identification and prioritisation with the industrial partners.

- M.Sc., Ph.D. and PostDoc students focus on these difficult applied research problems with input from industry.

- Research associates help the graduate students to integrate their new proposed algorithms in the toolchain for validation and optimisation.

- The best algorithms are added to the LTTng and Trace Compass toolchain with support from industrial partner's R&D engineers.

- LTTng is redistributed by most Linux vendors: Wind River, MontaVista, Linaro, LTSI, Debian, Ubuntu, Suse, Fedora/Red Hat... Code contributions were received from more than 100 individuals and 25 companies including Ericsson, Google, IBM, Red Hat, Ubuntu, Mentor...

# T1: Tracing, Debugging and Profiling Mechanisms and Architecture on Many-Core Systems

**Objective:** develop common efficient mechanisms for tracing and debugging heterogeneous many-core systems.
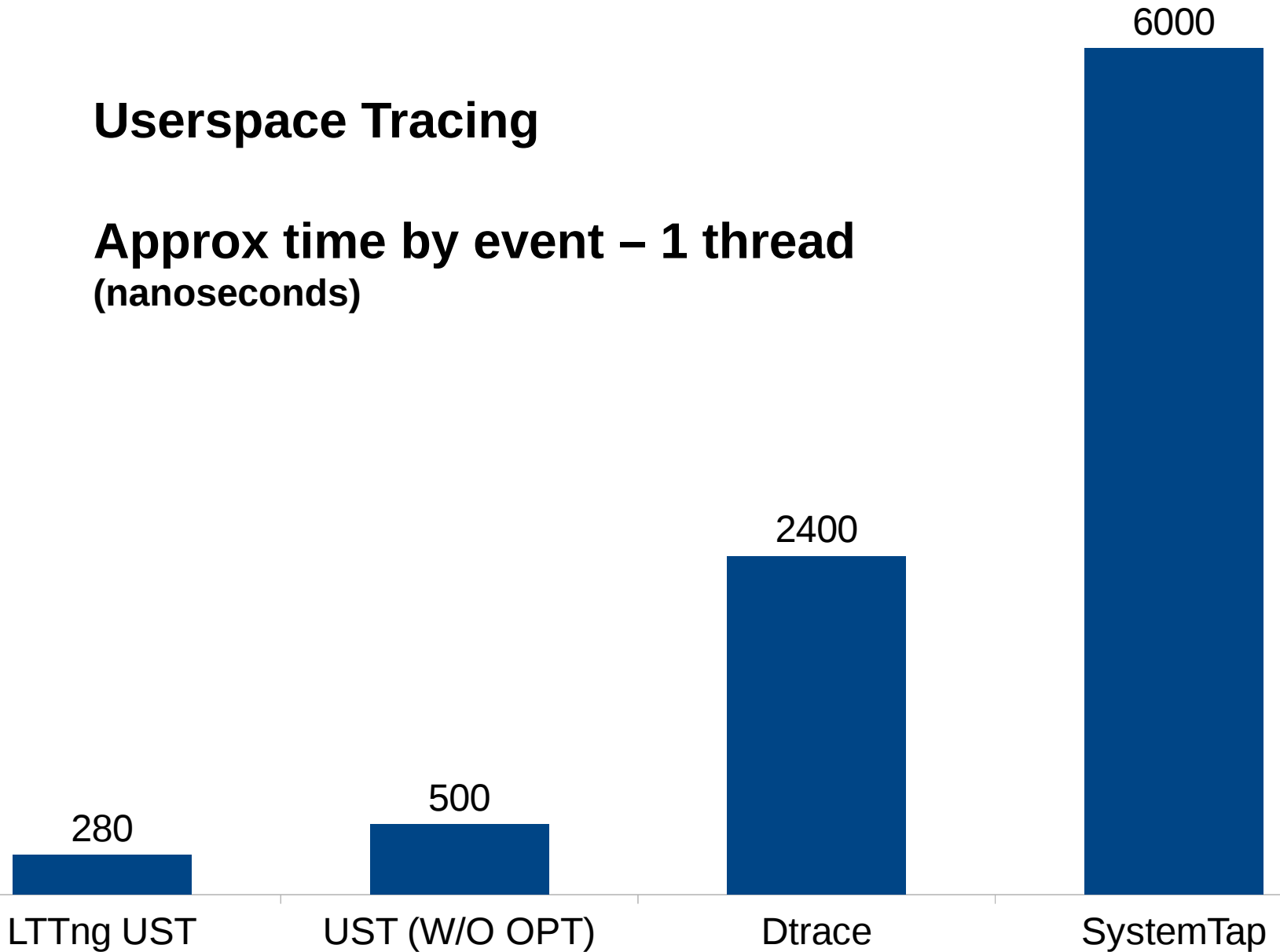
- T1M1: tracing new architectures (Adapteva Epiphany, Kalray MPPA).

- T1M2: tracing new architectures (ARM CoreSight, Intel PT, bare-metal processors).

- T1D1: efficient algorithms to control debugging, profiling and tracing on many-core processors and to retrieve data.

- T1D2: scalable tracing / debugging mechanisms for dynamic instrumentation, checking conditions, aggregating values...

**Key:** low disturbance, low overhead, no impact on system, dynamically adjust the overhead / level of details compromise
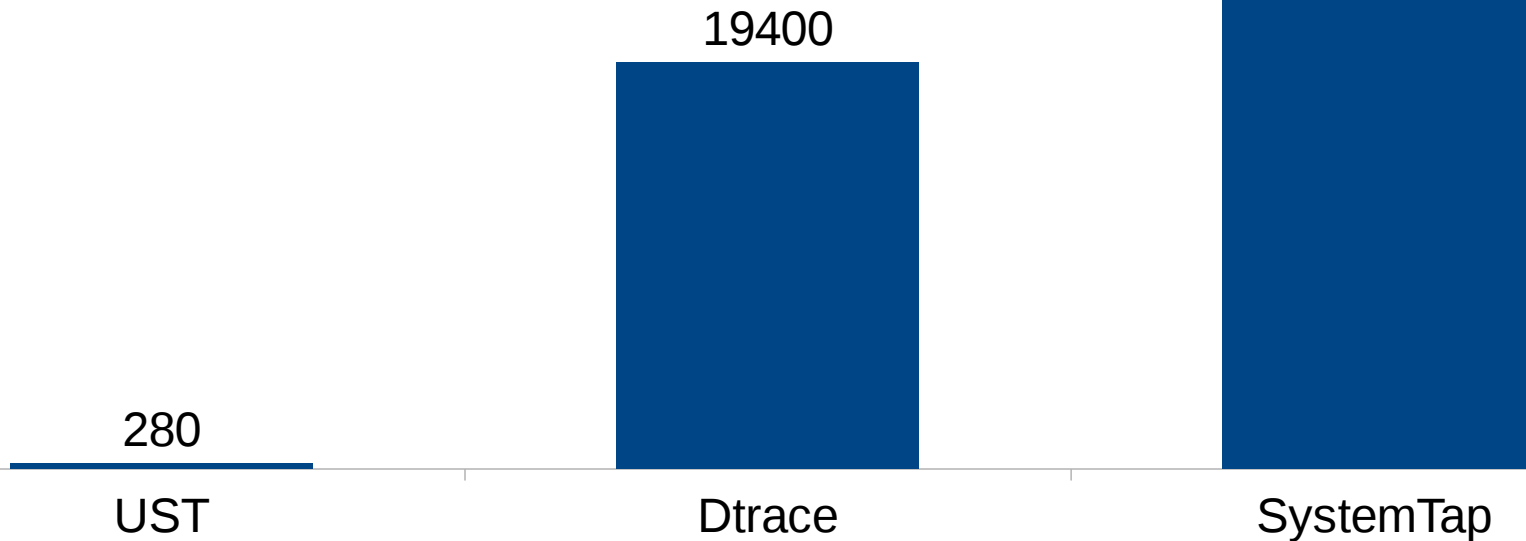
**Userspace Tracing**

**Approx time by event – 1 thread**
**(nanoseconds)**

| | | | |
|---|---|---|---|
| 280 | 500 | 2400 | 6000 |
| LTTng UST | UST (W/O OPT) | Dtrace | SystemTap |

# Trace collection performance

- User-space tracing in user space, no system call.

- "Unlikely if" for tracepoint activation.

- Efficient binary format being optimised and standardized as the Multi-Core Association Common Trace Format. No formatting!

- Per CPU buffers with local lockless atomic operations.

- Read Copy Update (RCU) synchronisation for configuration information (tracepoint activation, multiple sessions…).

- Zero-copy trace recording on disk.

- Tracepoints may be inserted even in interrupt and NMI contexts.

- Efficient timestamping (rdtsc) even from virtual machines.

- Most efficient and flexible tracer available!

# Latency-tracker

- Kernel module to track down latency problems at run-time.

- Simple API that can be called from anywhere in the kernel (tracepoints, kprobes, netfilter hooks, hardcoded in other module or the kernel tree source code).

- Keep track of entry/exit events and calls a callback (e.g., to dump stack or take a trace snapshot) if the delay between the two events is higher than a threshold, or no exit before timeout.

```
latency_tracker_event_in(tracker, key,
threshold, timeout, callback);
....
latency_tracker_event_out(tracker, key);
```

# Sample usage

- Block layer latency

  - Delay between block request issue and complete

- Wake-up latency

  - Delay between sched_wakeup and sched_switch

- Network latency

- IRQ latency

- System call latency

  - Delay between the entry and exit of a system call

- Offcpu latency

  - How long a process has been scheduled out

```
on sched_switch:
    event = latency_tracker_get_event(next_pid);
    if event && ((now – event->start) > threshold):
        dump_stack(next_pid);
```
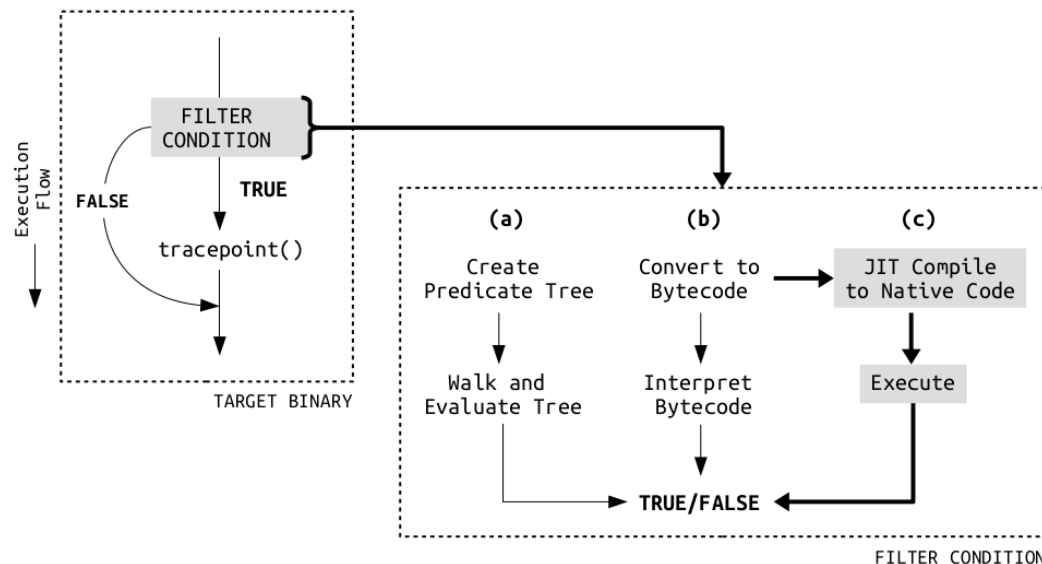
# Tracker performance

- Controlled memory allocation.

- Lockless free-list.

- Out-of-context reallocation of memory if needed/enabled.

- Using userspace-rcu hashtable (kernel version) for lockless insert and lookup.

- Custom call_rcu thread to avoid the variable side-effects of the builtin one.

# Dynamic instrumentation

- When we cannot recompile or restart program.

- When dynamic decisions can help optimise tracing.
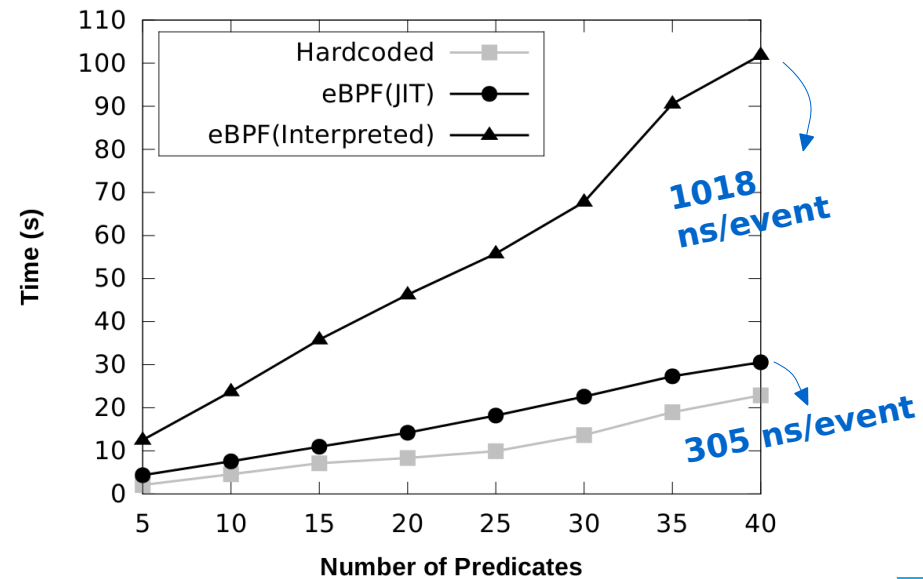
- Comparison of different kernel and user-space solutions, using bytecode and JIT compilation.

# Using eBPF bytecode or JIT for tracing



**Pure eBPF Filter Performance with 50 Predicates**

**Pure eBPF Filter Performance with Increasing Number of Predicates**

# T2: Cloud Debugging and Monitoring

**Objective:** provide efficient tools to monitor and understand the performance of clouds and cloud-based applications.

- T2M1: data collection, analysis and monitoring for clouds (OpenStack).

- T2M2: data collection, analysis and monitoring for Software Defined Networks (OpenDaylight).

- T2D1: monitoring of migrating virtual machines in the cloud, with ability to navigate through information at multiple levels.
  **Key:** combine information from multiple levels, maintain scalability through sampling and varying levels of details.

# Real Performance of Virtual Machines

# T3: Advanced Analysis

**Objective:** propose a number of new advanced analysis and views, in particular for specialized parallel processors.

- T3M1: specialized programming models such as OpenCL, OpenMP and MPI.

- T3M2: specialized analysis and views for tasks such as virtual / physical memory usage, and power consumption.

- T3D1: sophisticated pattern matching on the state history and tracing events.

- T3D2: trace correlation for regression testing, highlighting the differences between similar runs of different versions, or in different contexts.

**Key:** tackle huge traces efficiently, help the user model the system and automate problem identification, provide contextualized analysis.

17

## Trace Analysis

- Present events from different traces on a common reference time.

- Build a model of the traced system (process table, opened file descriptors…).

- Allow easy and efficient navigation and zooming in the trace.

- Present statistics and histograms for different intervals in the trace.

- Help identify the dependencies among events and trace the critical path.

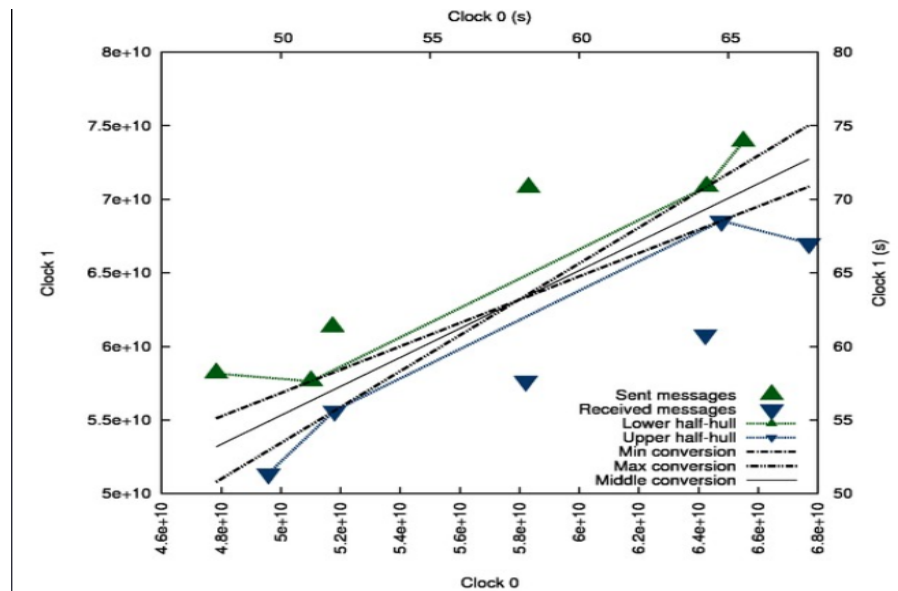- Traces may contain hundreds of GBs of data!

# Eclipse Linux Tools

# Distributed Traces Synchronization

- Each node, VM or heterogeneous processor may have its own independent clock, yet correlating events from different traces on a common time reference is needed.

- First order effect: different initial time and clock frequency.

- Second order effects: clock jitter, latency to read clock value, clock frequency drift with time.

- Identify corresponding send and receive events for packets to estimate clock differences.

# Synchronisation

- New linear incremental algorithm to compute the clock offset and drift between two traces.

- Send and received events are matched by unique id (e.g., TCP sequence number for a given connection) and are incrementally added to build the upper and lower convex hull bounds.



Convex Hull

# Modeled State

- Current State Tree expressed as a generic attribute tree.

- For the Linux kernel, the Current State Tree contains: process table, file descriptors, sockets...

- Some events cause a state change (based on event type, event fields, current state).

- State builder expressed as custom java code or declarative XML expression.

- Efficient State History Tree database for quickly navigating through traces.

- The same declarative XML expressions can be used to define views, filters, abstract events and metrics.

# Declarative XML expression example

**XML structure**

```xml
<filterHandler filterName="sched_switch">
    <initialFsm id="sched_switch" />
    <transitionInput id="sched_switch">
        <event eventName="sched_switch"/>
    </transitionInput>
    <action id="update Current_thread">
        <stateChange>
            <stateAttribute type="location" value="CurrentCPU" />
            <stateAttribute type="constant" value="Current_thread" />
            <stateValue type="eventField" value="next_tid" />
        </stateChange>
    </action>
    <fsm id="sched_switch" multiple="false">
        <precondition input="sched_switch"/>
        <stateTable>
            <stateDefinition name="sched_switch">
                <transition input="sched_switch" next="sched_switch" action="update Current_thread" />
                <transition input="#other" next="sched_switch" />
            </stateDefinition>
        </stateTable>
        <initialState id="sched_switch"/>
    </fsm>
</filterHandler>
```
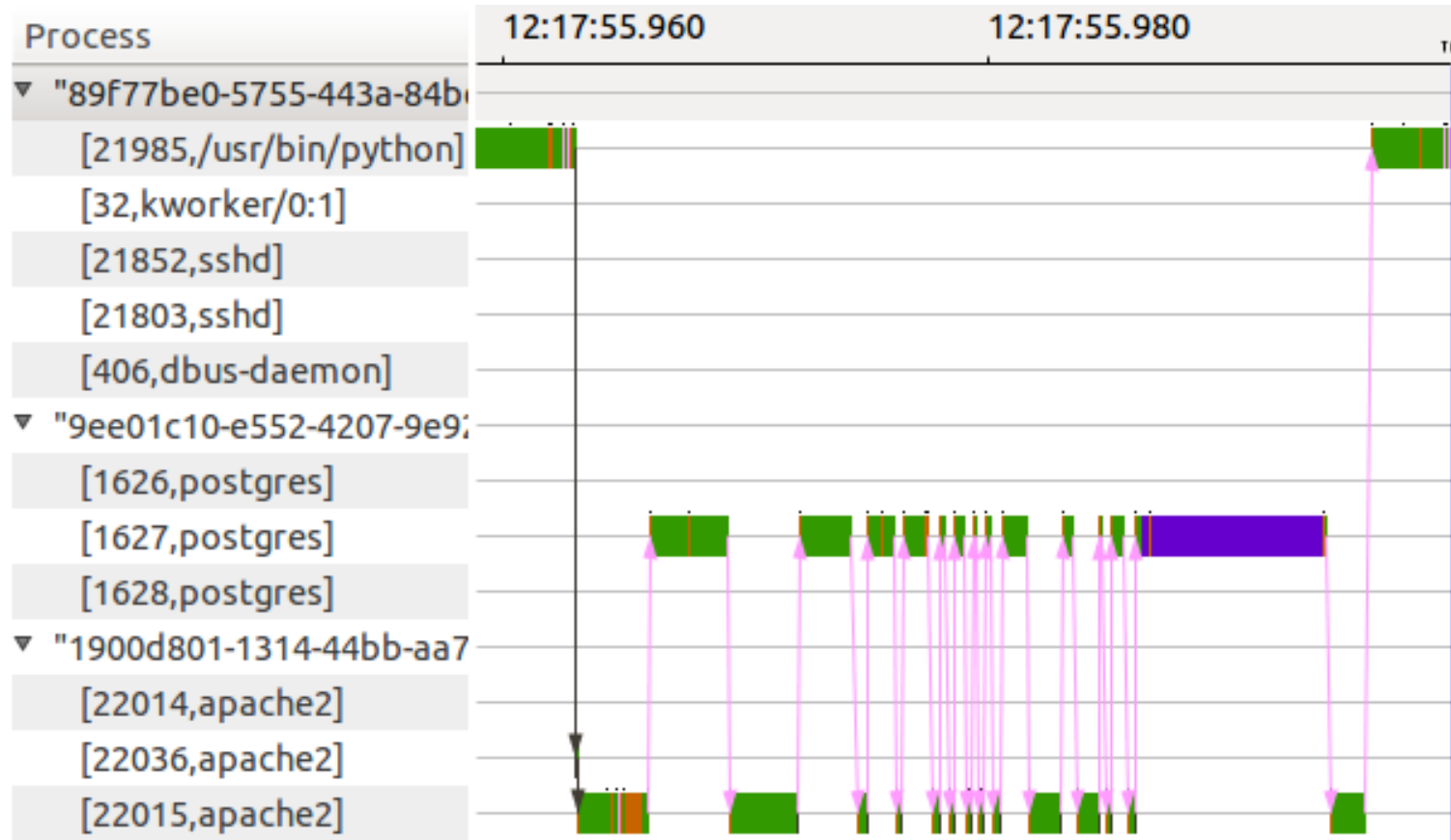
Transition

Action

FSM

# Critical Path Analysis

- For each process, the elapsed time is broken down by state: running, waiting for resource (file, pipe, socket, CPU) or sleeping on a timer.

- The critical path for reaching a point in a process is the chain of unblocking events (wakeup) and their origin (device ready, lock freed or packet received from another process).

- Find and display the critical path with a breakdown of elapsed time per process and by state (e.g. web request took 3s of Firefox, 2s of X11, 5s of Apache, 3s of MySQL).

# Critical Path Analysis view

# Profile a path segment by call stack

# Call stack computation



Unwind time according to call stack depth

CPython traceback time according to call stack depth

# Real-time systems analysis

- Develop a model to define real-time task executions in a trace.

- Identify common problems in real-time systems and useful information for their analysis.

- Develop a method to analyze the trace segment corresponding to an execution to identify if the execution presents a problem.

# Overview

1) Control Flow View

2) Define executions

3) Stackbars View

4) Critical Flow View with
   CP Complement view

5) Other views

# Stackbar view: comparing task executions



*Figure 10 : Stackbars view*
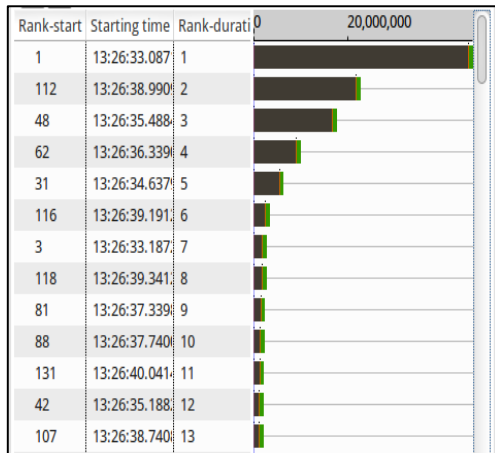
# Distribution: duration and along the time axis
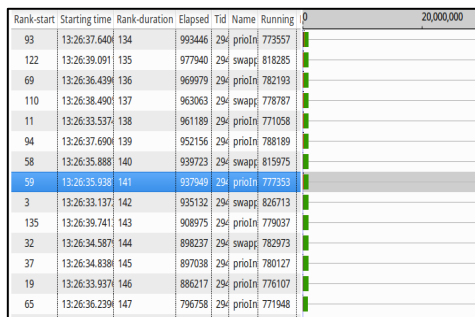


Figure 18 : Problematic executions
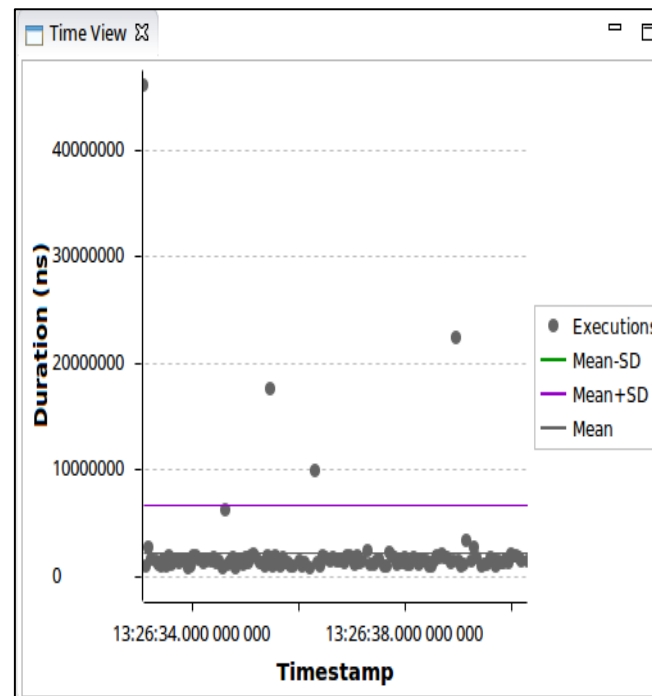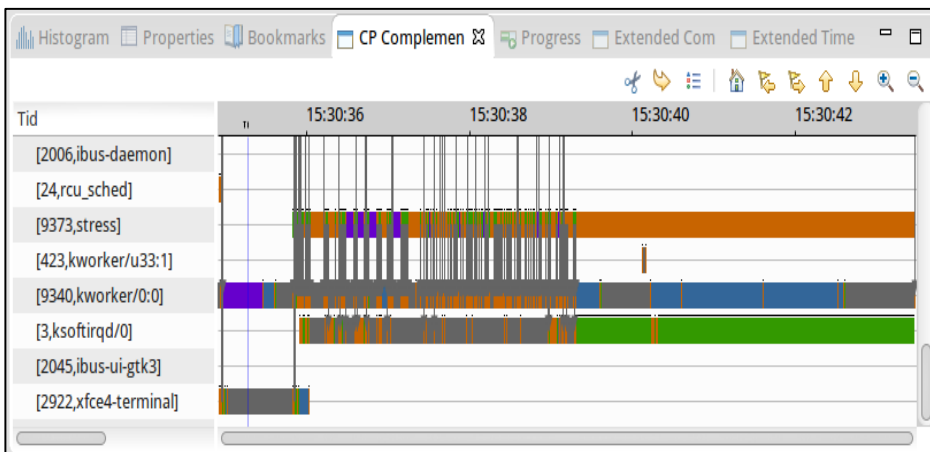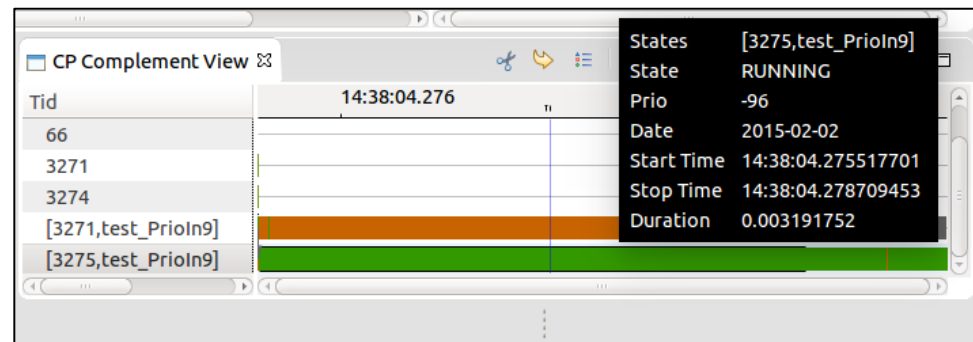


Figure 19 : Normal executions
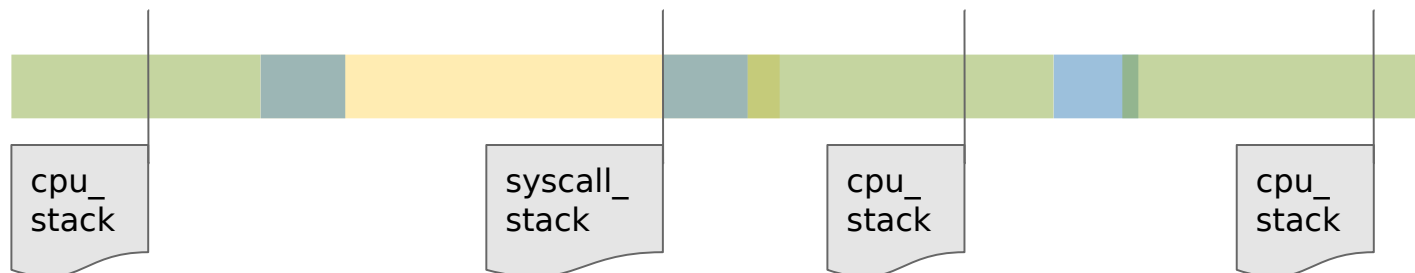


Figure 20 : Time View

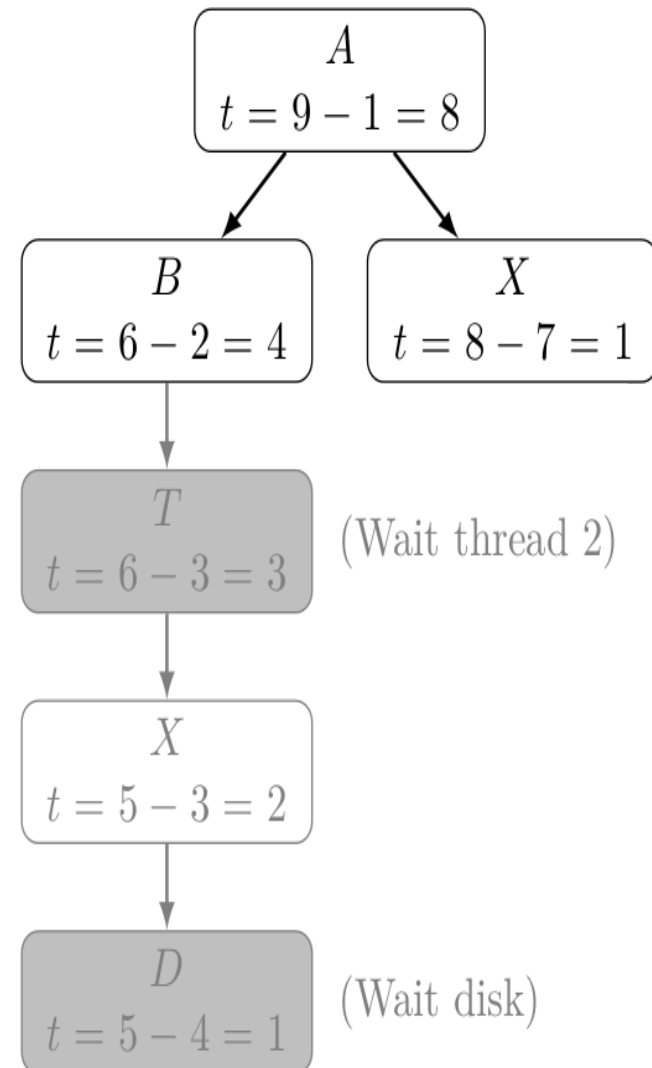# Critical path view for a task execution

# TraceCompare: diagnose performance variations

- Good versus bad run: program, library or OS version, hardware, network load?

- Collect execution traces and compute different metrics for different tasks (e.g., time for complete process, search request, or frame redraw).

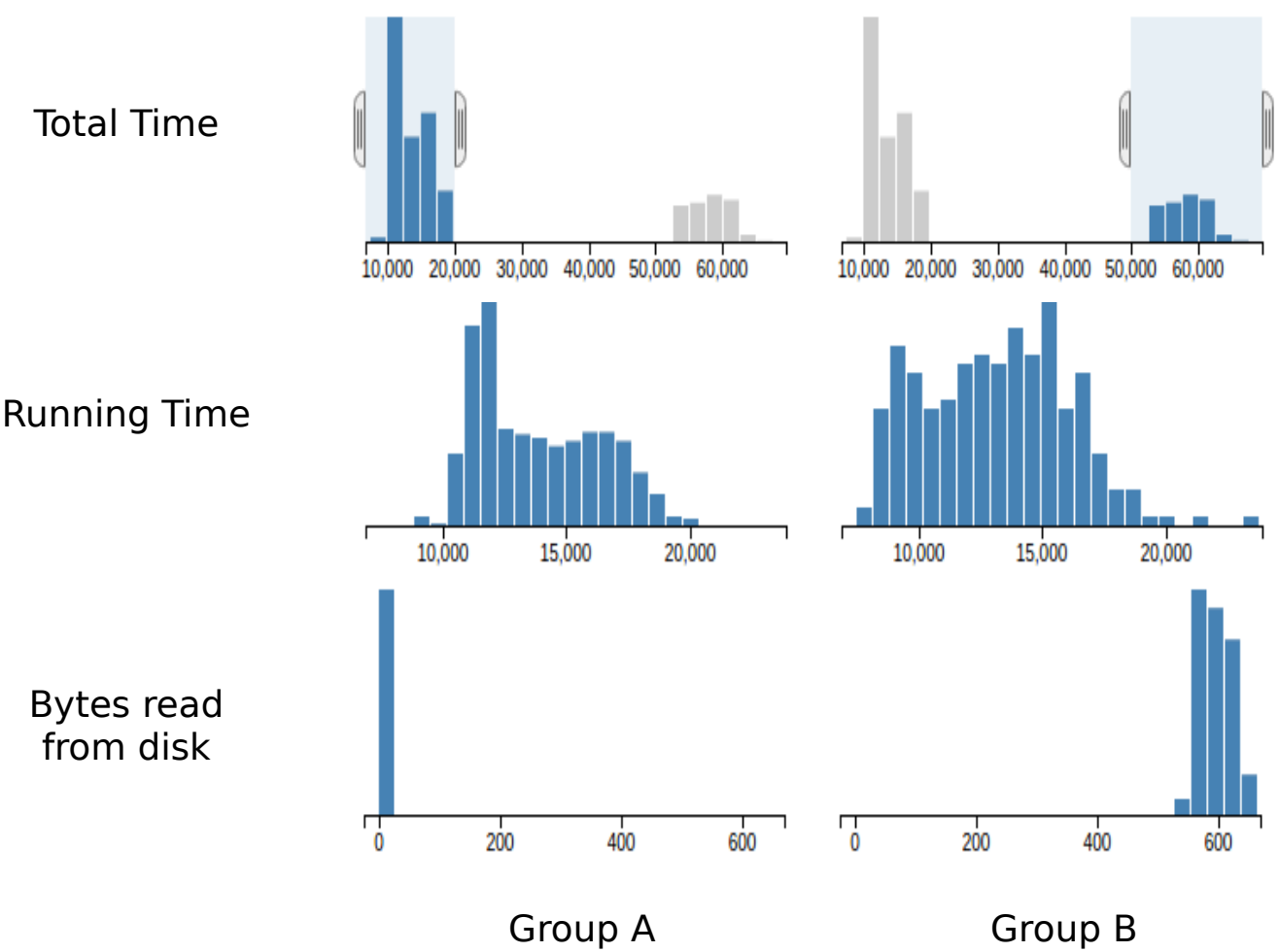- Add call stack samples collected periodically, or upon a very long system call.

cpu_
stack

syscall_
stack

cpu_
stack

cpu_
stack

# Call context tree with nested calls and waits

| Time | Thread 1 |
|------|----------|
| 1    | Call A   |
| 2    | Call B   |
| 3    |          |
| 4    |          |
| 5    |          |
| 6    | Return B |
| 7    | Call X   |
| 8    | Return X |
| 9    | Return A |

$A$
$t = 9 - 1 = 8$

$B$
$t = 6 - 2 = 4$

$X$
$t = 8 - 7 = 1$

$T$
$t = 6 - 3 = 3$ (Wait thread 2)

$X$
$t = 5 - 3 = 2$

$D$
$t = 5 - 4 = 1$ (Wait disk)

# Select groups of executions based on metrics



Total Time

Running Time

Bytes read from disk
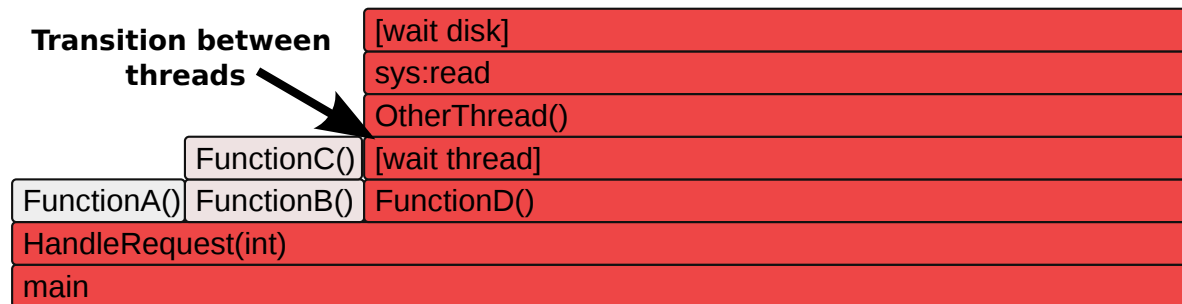
Group A        Group B

# Differential flame graph of the two groups

- View of call stack over time for the tasks (assuming relatively stable call tree).

- Average duration of each call in group A versus group B, highlight where the differences are significant.
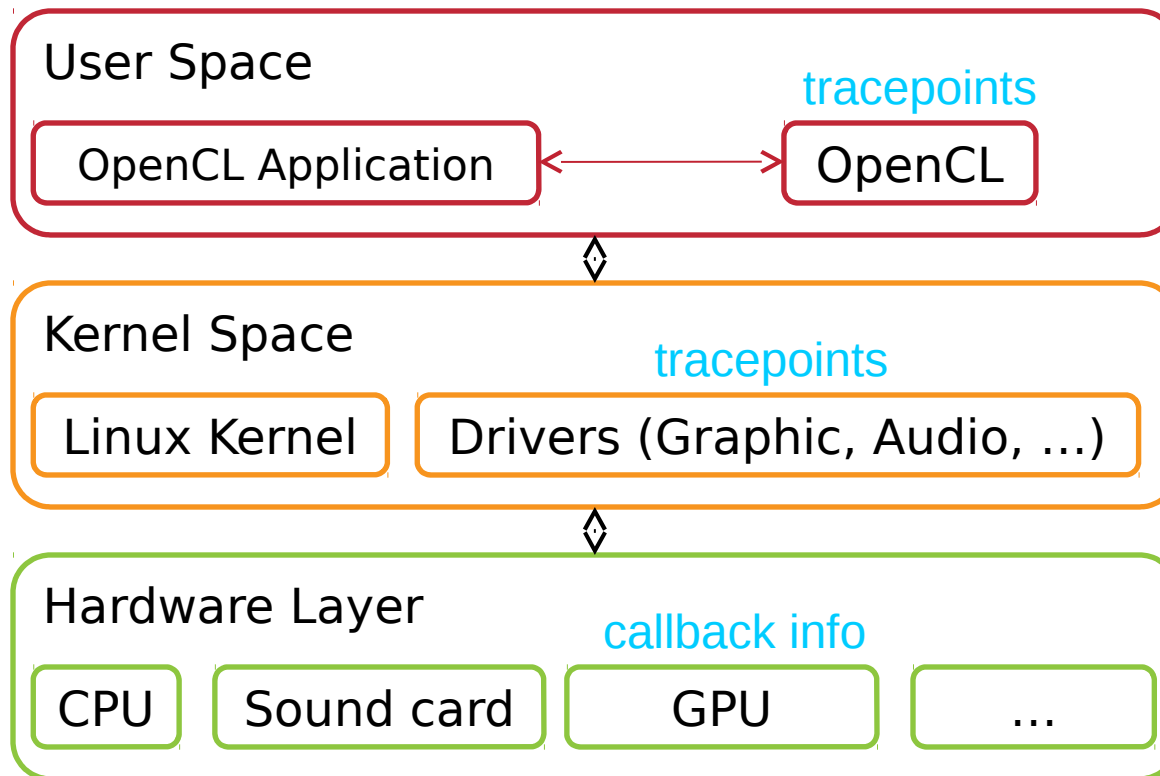
**Transition between threads**

| | | |
|---|---|---|
| | | [wait disk] |
| | | sys:read |
| | | OtherThread() |
| | FunctionC() | [wait thread] |
| FunctionA() | FunctionB() | FunctionD() |
| HandleRequest(int) | | |
| main | | |

# Tracing the GPU as an asynchronous resource
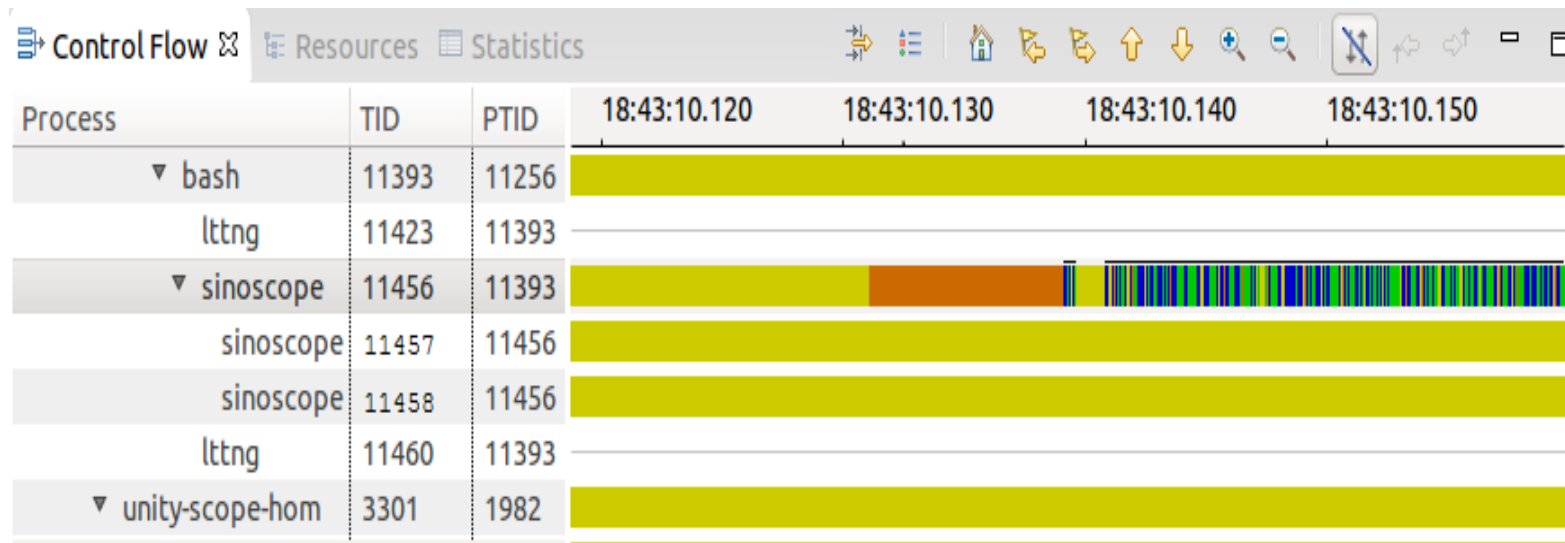


Source: nvidia.com

# Combining info from OpenCL, driver and HW
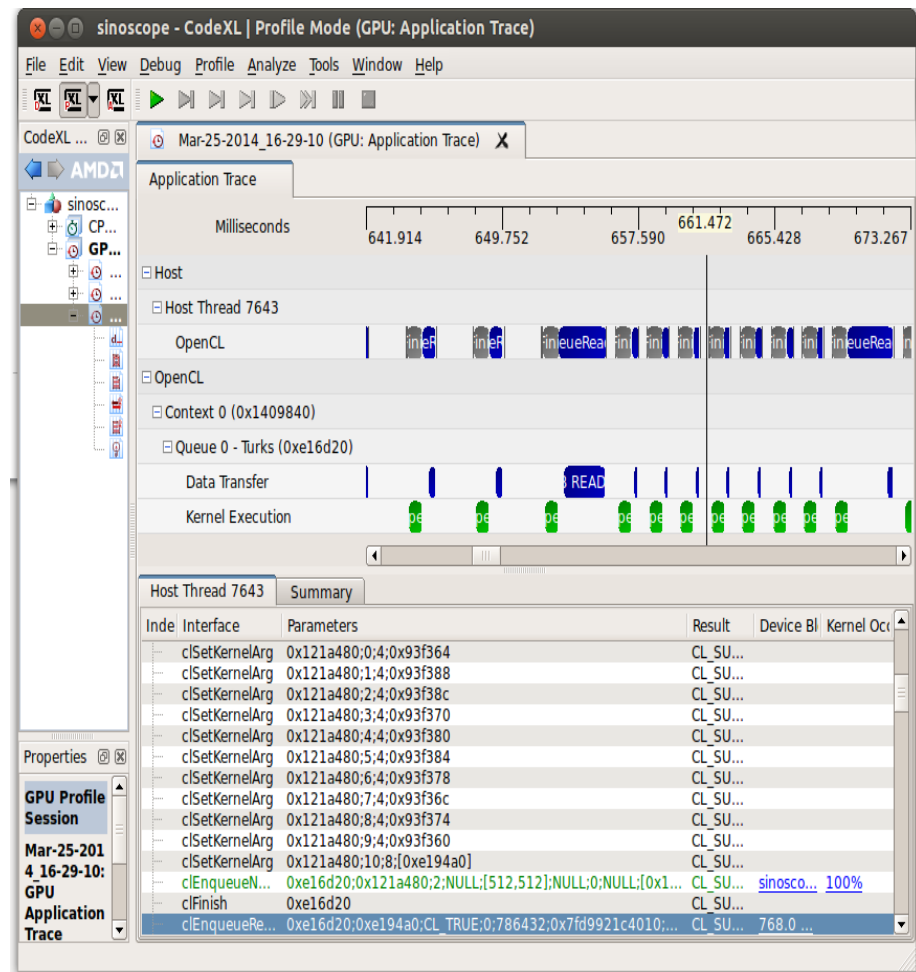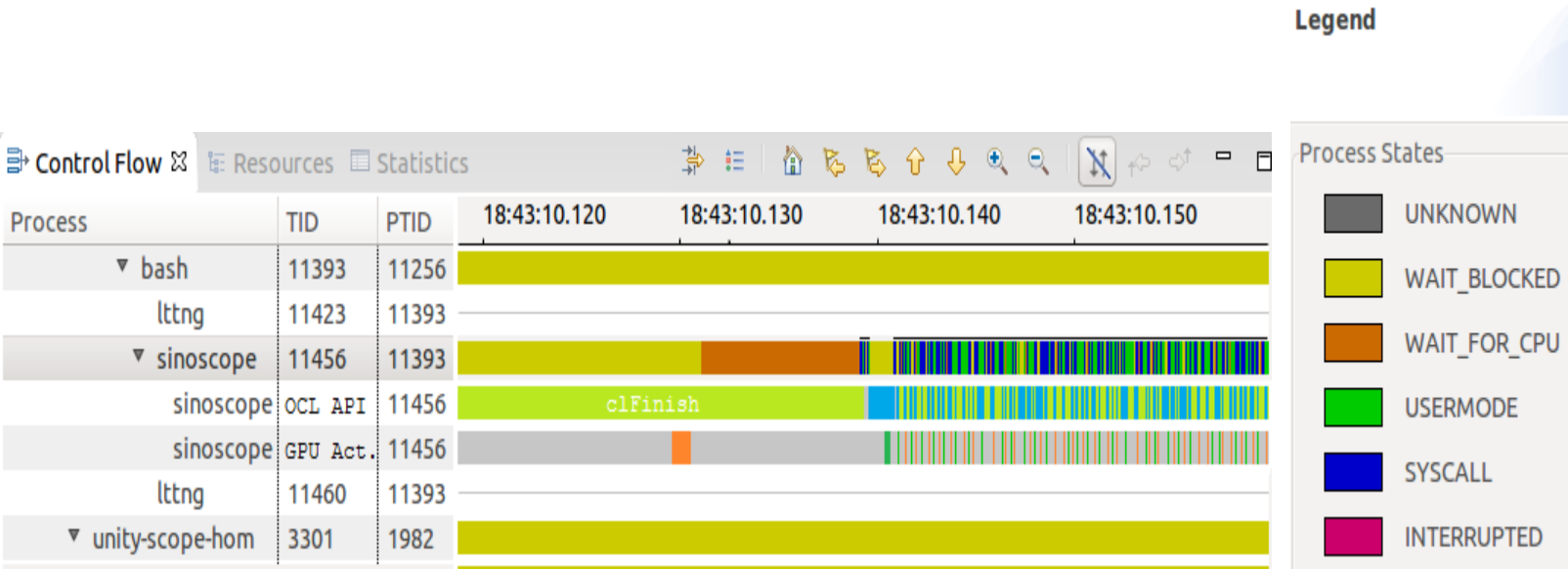
# CPU trace

Trace Compass:

# GPU trace

CodeXL
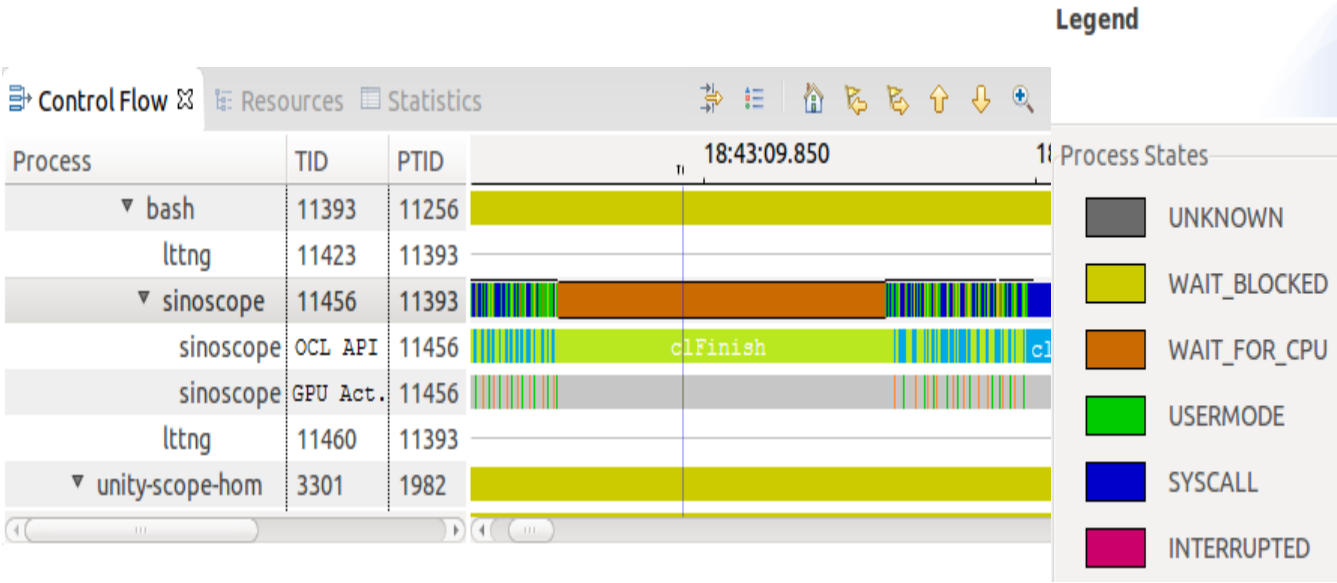(AMD's GPU tracing tool)

Limitations:

- OpenCL application has to be launched by CodeXL

- Recording trace performance for large traces

- AMD GPU only

# Combined trace: GPU contention

# Combined trace: CPU preemption
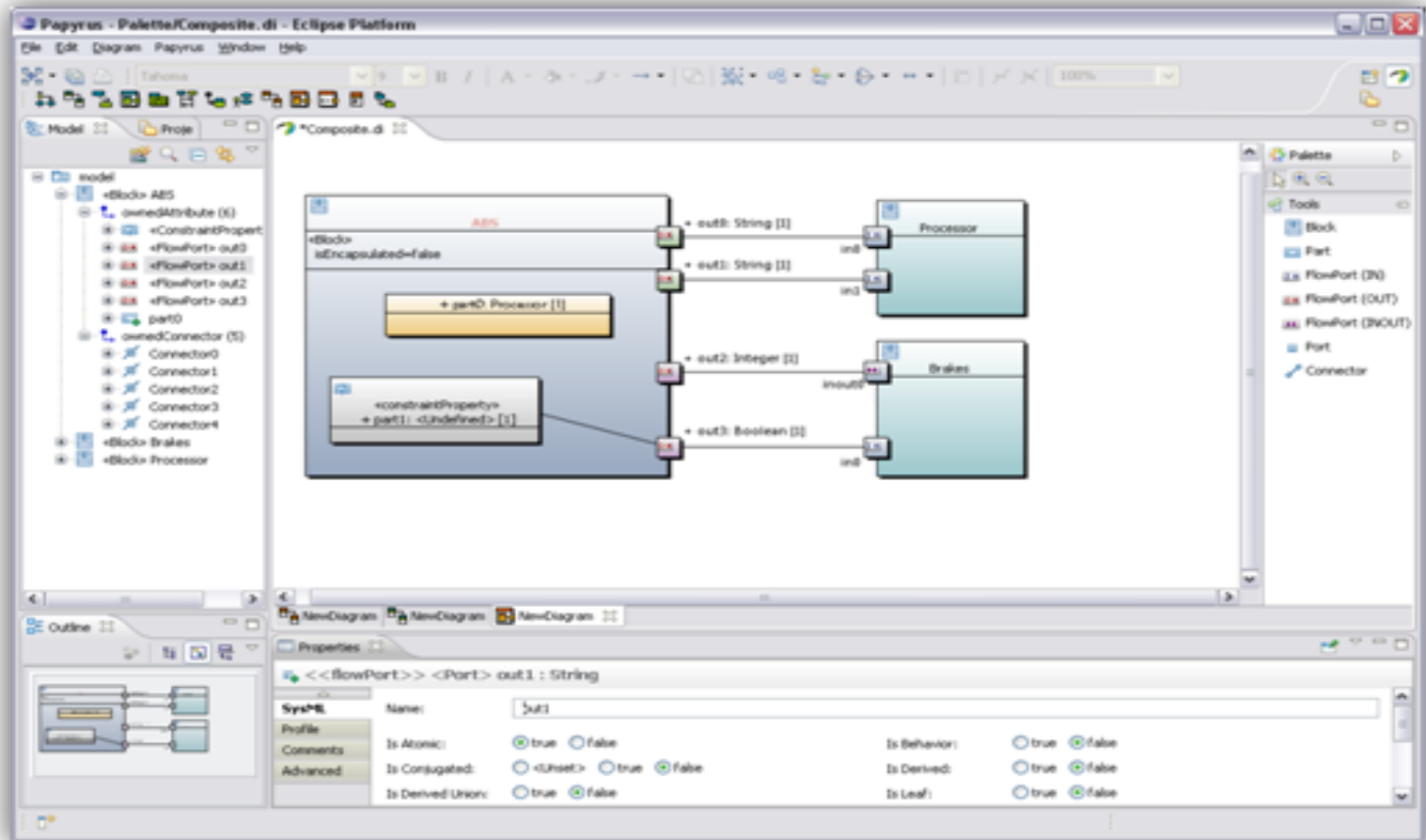
# T4: Model-Driven Engineering support

**Objective:** fully support model-driven engineering with tracing, debugging and testing tools.

**Personnel:** Prof. Juergen Dingel, Francis Bordeleau Ph.D., Sébastien Gérard Ph.D., Bernd Hufmann B.Eng., RA1.

- T5M1, T5M2: extend the model-level to store as attribute values, current states, transitions taken and messages sent.

- T5D1: new algorithms and techniques for the automatic correction and refinement of environment assumptions based on runtime information.

- T5M3, T5M4: extend the model-level to specify test cases, execute the tests and annotate the results.

- T5D2: new automatic test generation algorithms for correcting, refining and completing the model-level information.

# Model-Driven Engineering support

# Conclusion

- Efficient tracing is an essential tool for high performance heterogeneous distributed systems.

- Free software as infrastructure: natural collaboration, rich user community, easy customization, no surprise (e.g. product or supplier end of life, high cost per seat).

- Viewpoints of multiple clients: benefit from the experience, use cases and best practices of several technology leaders.

- Technology transfer: research and prototyping in Polytechnique, collaborative development of mainline implementations.

- LTTng delivers low disturbance, scalable tracing.