# Beyond kernel tracing

**Pierre-Marc Fournier**

ÉCOLE
**POLYTECHNIQUE**
M O N T R É A L

# LTTng Tracers

- Kernel tracer
  - Record events in the kernel.
  - Project started in 2005.

- Userspace tracer
  - Port to userspace of kernel tracer.
  - Project started in 2009.
  - Preceded by several generations of experimental userspace tracers.

# Why userspace tracing?

- The kernel knows about: scheduling, blocking, IRQs, softIRQs, traps, system calls, signals.

- But it's not enough.


- Essential for application that:

  - Talk directly to hardware

  - Handle lots of sessions

  - Are complex in general

# Objectives

- Lowest performance impact possible.

- Optimized for huge event rates and traces.

- Low intrusiveness.

- Maximize traceability.

Combined tracing of the kernel and applications with LTTng

# Existing userspace tracers

- DTrace, SystemTap
- Historical LTTng userspace tracing methods
- Ftrace
- "printf"

- **Need for better performance.**

Combined tracing of the kernel and applications with LTTng

# Performance

- No system call, no trap.

- Minimize impact on instruction cache.

- Minimize impact on data cache.

- Never copy events.

- Markers / tracepoints.

Combined tracing of the kernel and applications with LTTng

# Scalability

- Use of userspace RCU.

- Non-blocking atomic operations.

- Per-CPU buffers.

Combined tracing of the kernel and applications with LTTng

# Trace format

- Binary.

- Saved in system byte order.

- Compact.

- Same format as LTTng kernel tracer.

Combined tracing of the kernel and applications with LTTng

# Combined Tracing

- Tracing both the kernel and userspace at the same time.

    - Using a kernel tracer.
    - And using a userspace tracer.

- This provides a full tracing solution that produces events covering the whole system.

- Traces are merged at analysis time.

Combined tracing of the kernel and applications with LTTng

# Maximize traceability

- Early tracing.

    - Early activation of instrumentation.

- Crash recovery.

- Support multi-threaded applications.

- Signal safe.

- Follow forks.

- Do not require initialization. (URCU BP)

- Trace executable, dynamically linked libs, even

    dynamically loaded

# Time

- Timestamp must be accurate and fast to read.

- Combined tracing: must be synchronized with kernel time source.

- x86: using cycle counter (TSC)

- Fallback: gettimeofday()

- Need userspace trace clock API

Combined tracing of the kernel and applications with LTTng
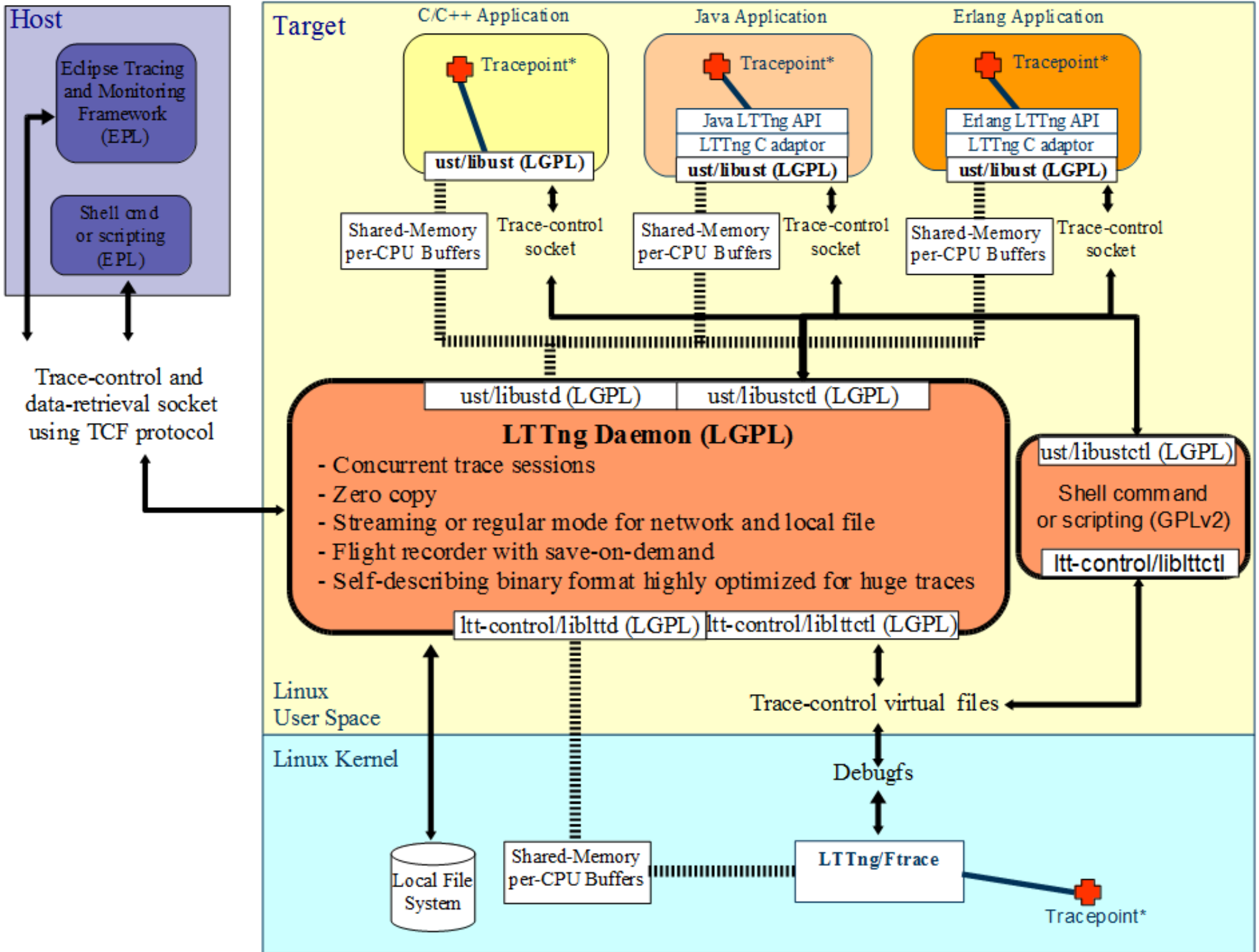
# Language support

- Native API in C (markers, tracepoints)
- Trivial port to any language that can call C code.

- Special support needed for
  - Passing fundamental types in payload.
  - Efficient individual activation/deactivation.

Combined tracing of the kernel and applications with LTTng

# System awareness

- Quotas

- Permissions

- Multi-process sessions

Combined tracing of the kernel and applications with LTTng

# Performance

- Early performance measurements.

  - Dtrace: 5 us / event

  - UST: 698 ns / event

- Not yet in optimization phase.

Combined tracing of the kernel and applications with LTTng

# Status

- UST 0.1 to be released soon.

- Working with Suse Linux for integration.

- GDB static tracepoints integration in progress.

- Git for userspace tracer:
  **http://git.dorsal.polymtl.ca/?p=ust.git**

Combined tracing of the kernel and applications with LTTng