# From Network to Application: Understanding Your Distributed System with Trace Compass

**Geneviève Bastien
Research Associate
Dorsal Laboraty
École Polytechnique de Montréal**

# Trace Compass

- Formerly known as TMF, the Linuxtools LTTng Eclipse plugin.

- Trace visualization tool

  - Standalone Rich Client Platform (RCP) application.

  - Also available as Eclipse plugins.

- Extendable framework

  - Add support for new trace types

  - Build trace analysis

  - With data-driven analysis, it's now easier than ever

# Trace Compass

- Now goes beyond Linux-only
  - Trace types:
    - LTTng / CTF
    - BTF
    - Custom text and XML
    - GDB
    - PCAP
    - Windows! (prototype with CTF converter)
  - Analysis:
    - LTTng Kernel: Control Flow View, Resources View, CPU usage
    - LTTng UST: Memory Usage (liblttng-libc-wrapper), CallStack View (-g -finstrument-functions)
    - PCAP: **Network Stream lists**
    - META: **Data-driven analysis, Network trace synchronization, Virtual Machine analysis,** Critical path analysis
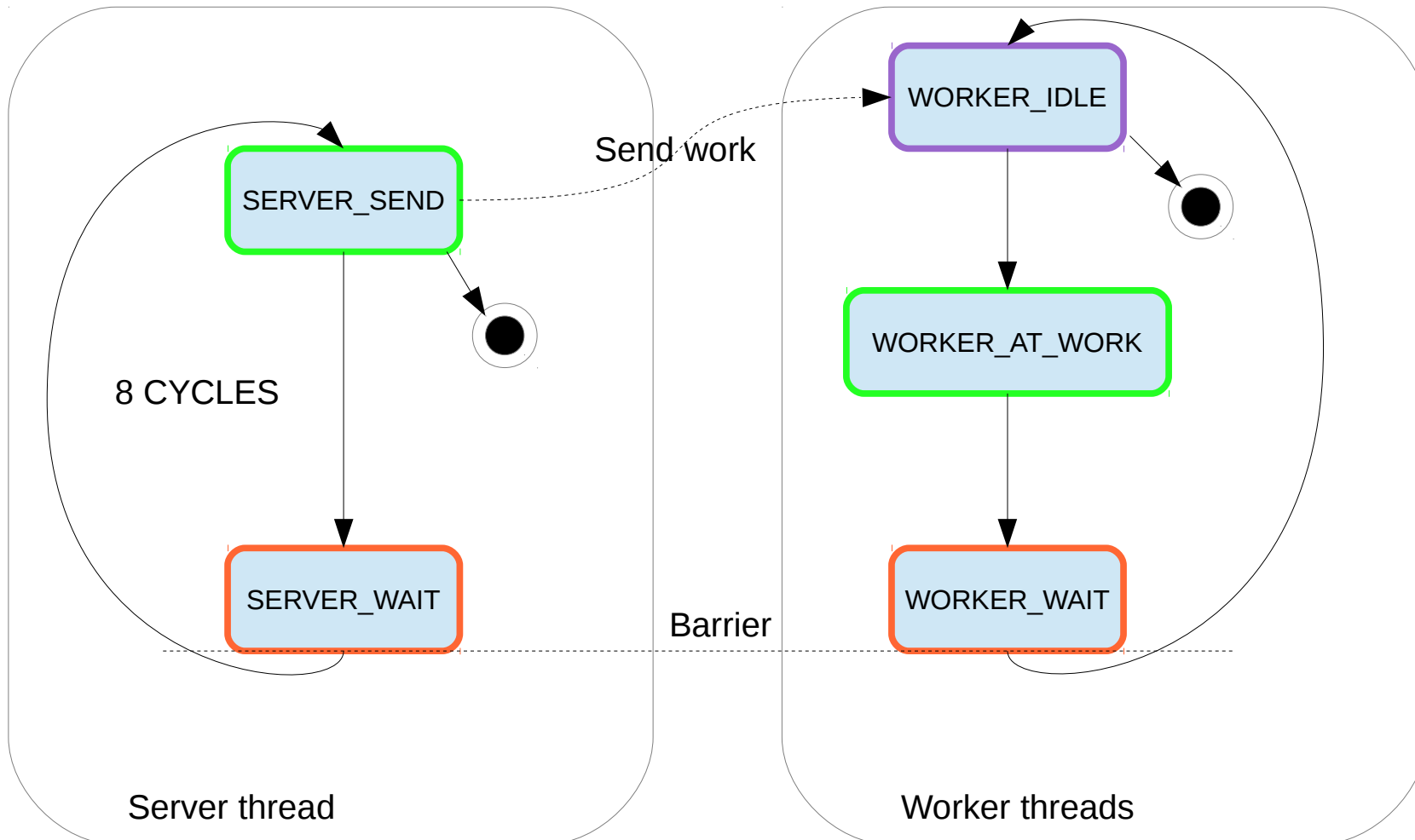
- 1 distributed application : 3 use cases
    - Local only (show data-driven analysis)
    - On 2 machines on the network (show network analysis)
    - On 2 virtual machines on the same host (show virtual machine analysis)
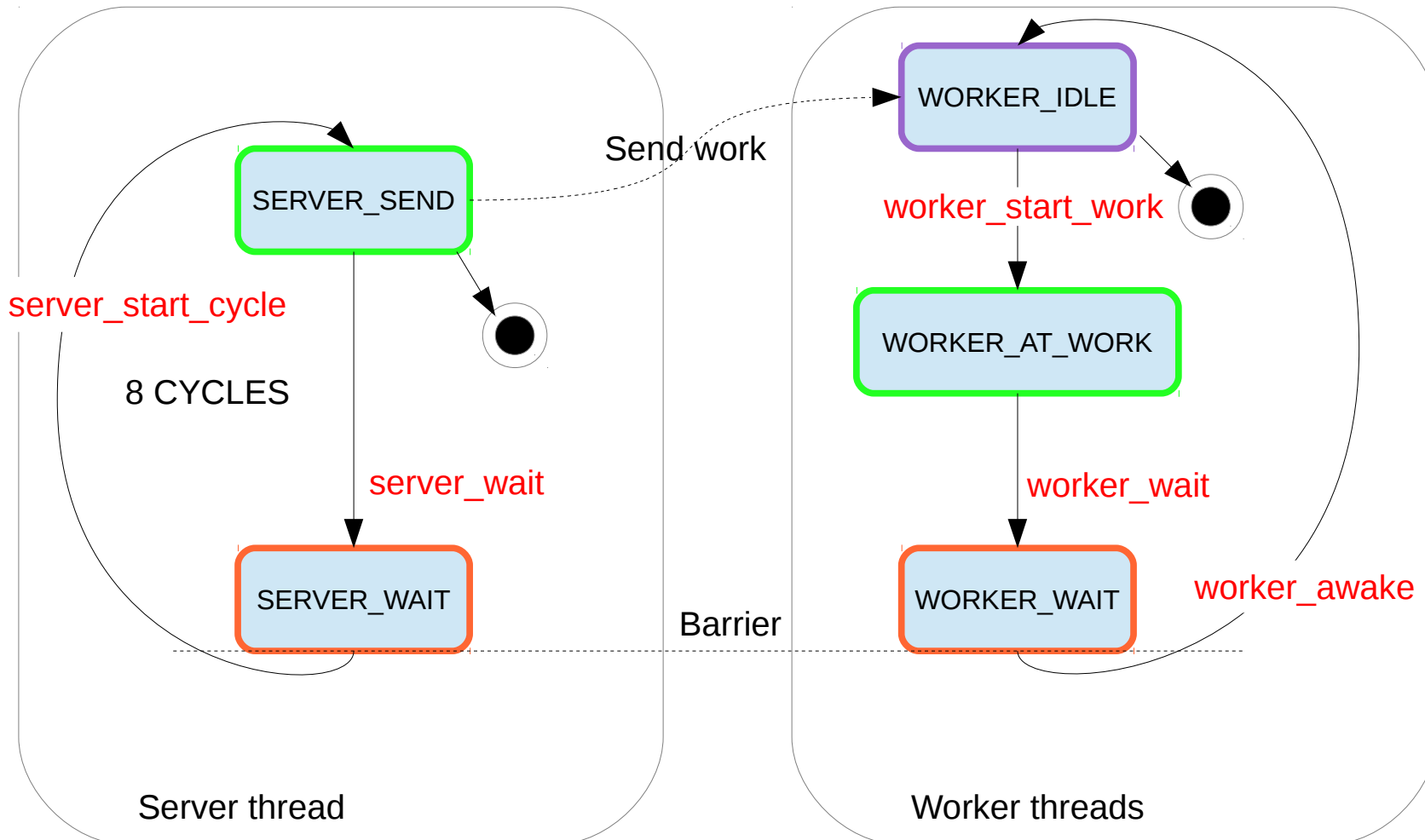
- MPI application: 5 worker threads + 1 server sending imbalanced workload to workers.

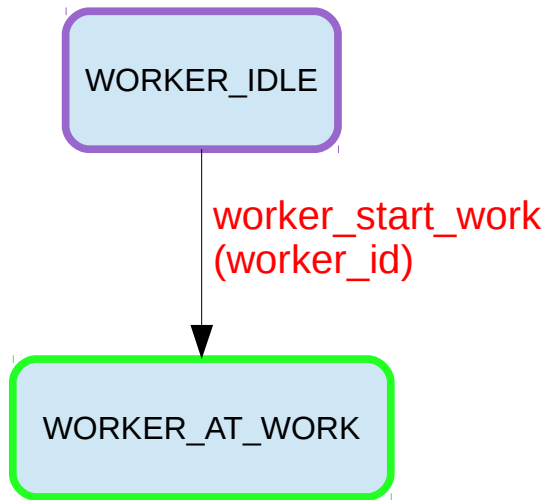- MPI application: 5 worker threads + 1 server sending imbalanced workload to workers.

# Data-Driven Analysis

WORKER_IDLE

worker_start_work
(worker_id)

WORKER_AT_WORK

State change:

Worker/<worker_id> = WORKER_AT_WORK

```
...
<stateProvider id="mpi.imbalance.sp">
  ...
  <definedValue name="WORKER_AT_WORK" value="2" />
  ...
  <eventHandler eventName="mpi_imbalance:worker_start_work">
    <stateChange>
      <stateAttribute type="constant" value="Worker" />
      <stateAttribute type="eventField" value="worker_id" />
      <stateValue type="int" value="$WORKER_AT_WORK" />
    </stateChange>
  </eventHandler>
</stateProvider>
...
```

- Visualization of the thread's states: time graph views or XY views

```
<timeGraphView id="mpi.imbalance.view.timegraph">

    <definedValue name="WORKER_AT_WORK" value="2" color="#66FF33" />
    <definedValue name="WORKER_WAIT" value="3" color="#FF3300" />
    <definedValue name="WORKER_IDLE" value="4" color="#CC66FF" />

    <entry path="Worker/*">
        <display type="self" />
    </entry>
</timeGraphView>
```

# Future work

- Data-driven analysis:
    - Define visually, with state diagrams
    - Smart filters and user-defined actions on those filters
    - And much much more!
- GPU traces and analysis
- Compare traces from different executions, for CPU/Memory usage, etc.
- Live tracing
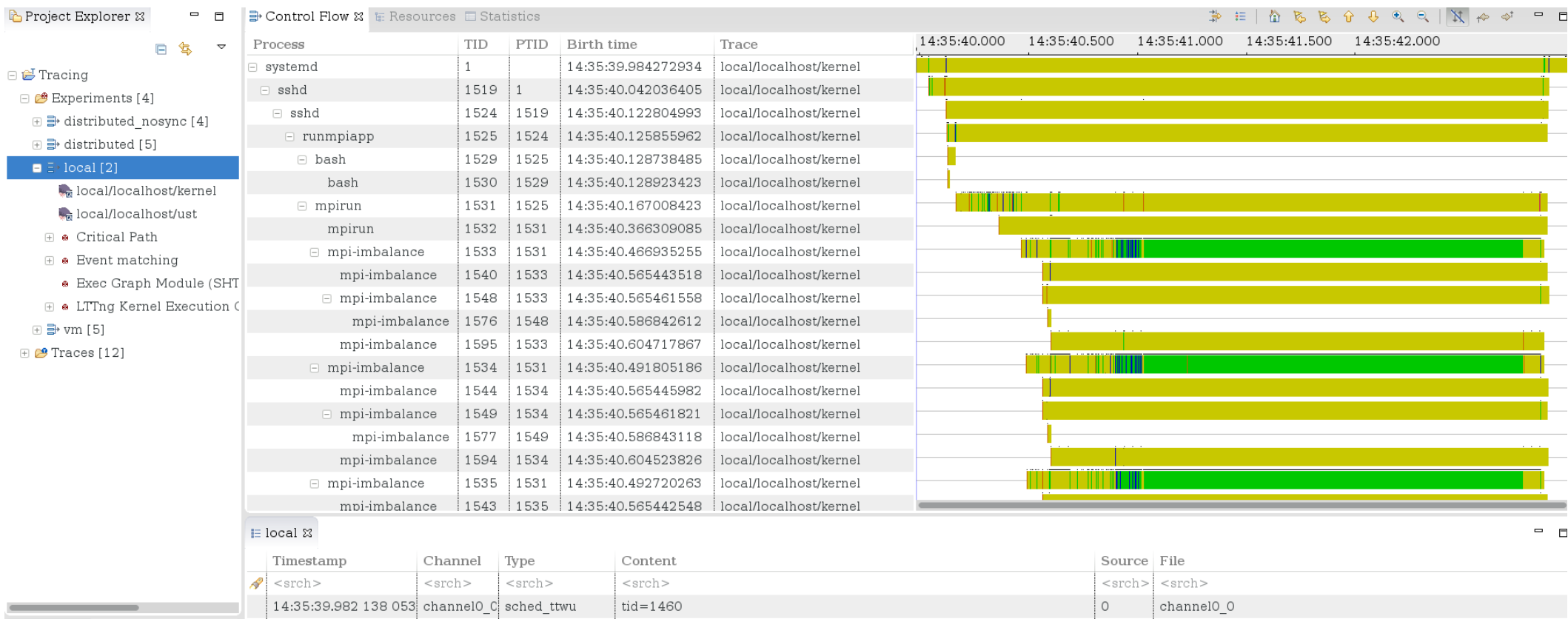- Improve performances with large experiments

# Resources

– Home Page: http://www.eclipse.org/tracecompass

– Mailing List: https://dev.eclipse.org/mailman/listinfo/tracecompass-dev

– Trace Compass standalone application used in this presentation:
http://secretaire.dorsal.polymtl.ca/~gbastien/TracingRCP/DorsalExperimental/

– Sources:

  • Master (coming soon): http://git.eclipse.org/c/tracecompass/org.eclipse.tracecompass.git

  • TMF in Linuxtools: (under the lttng folder)
git://git.eclipse.org/gitroot/linuxtools/org.eclipse.linuxtools.git

  • Experimental: branch dorsal_experimental
http://git.dorsal.polymtl.ca/~gbastien?p=linuxtools-tmf.git;a=summary

– Used in this demo:

  • Sample MPI traces and XML analysis:
http://secretaire.dorsal.polymtl.ca/~gbastien/tracingSummit2014/

  • MPI-imbalance source code: branch cluster (folder cluster/mpi-imbalance)
http://git.dorsal.polymtl.ca/~gbastien?p=workload-kit.git;a=summary
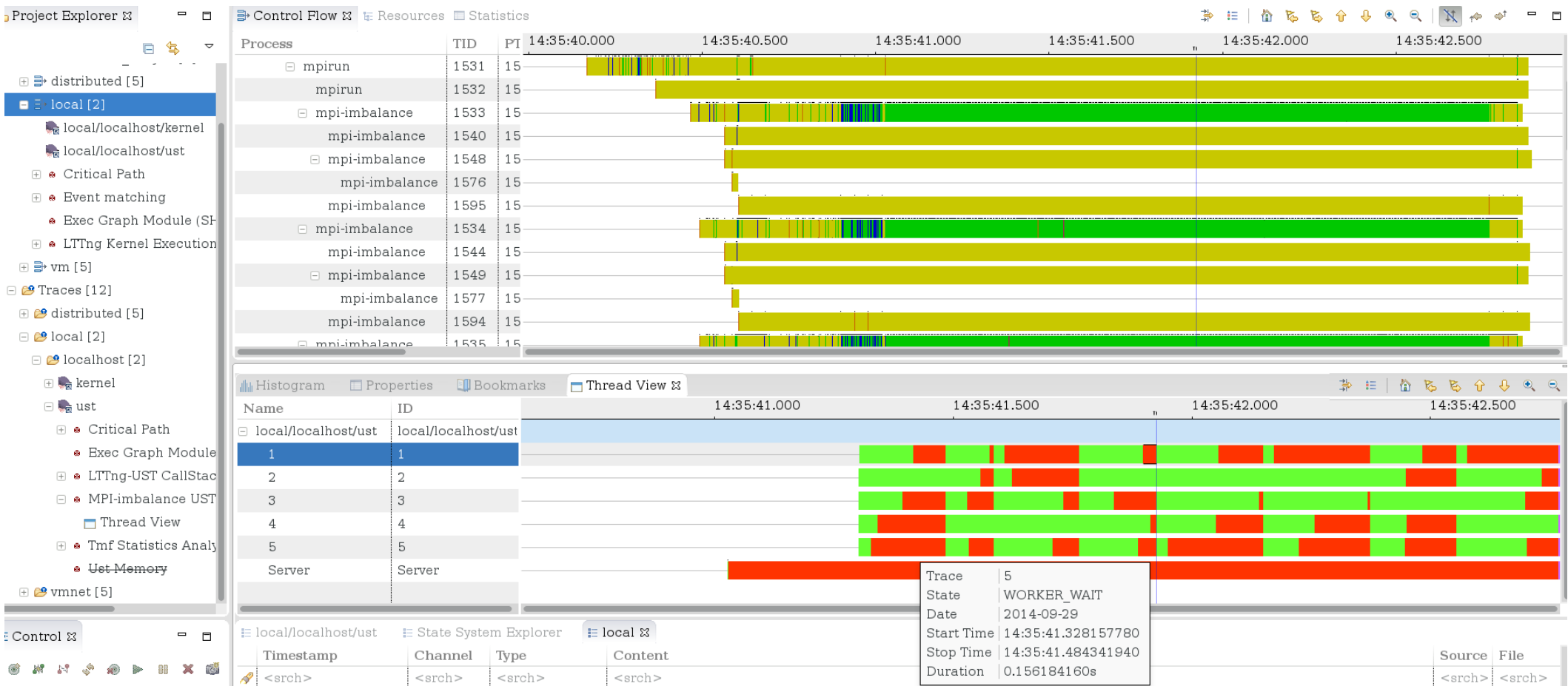
– IRC: #lttng on oftc

– More doc and links: http://lttng.org/eclipse

# Annexes

(Screenshots in case Eclipse refuses to cooperate)
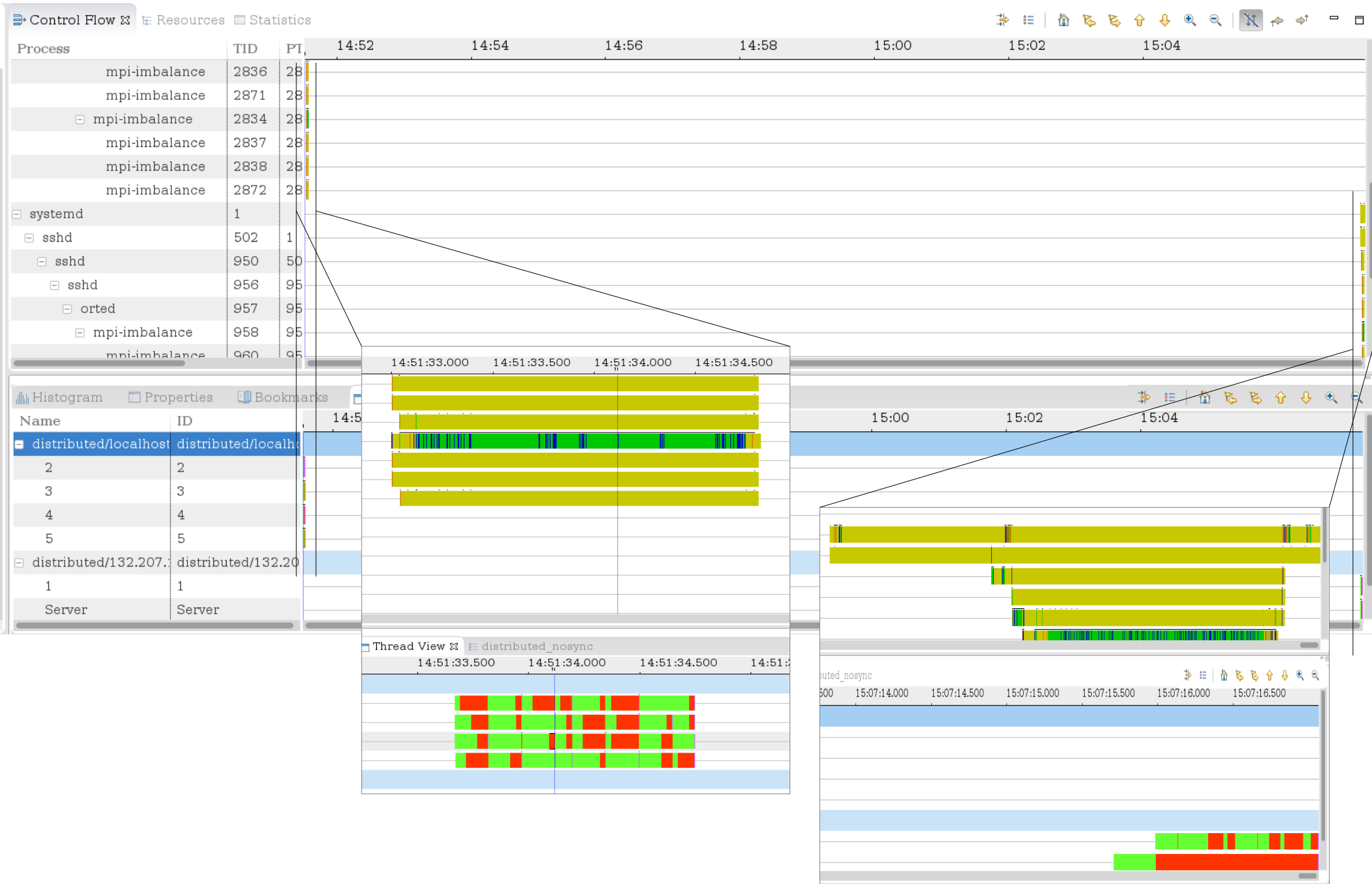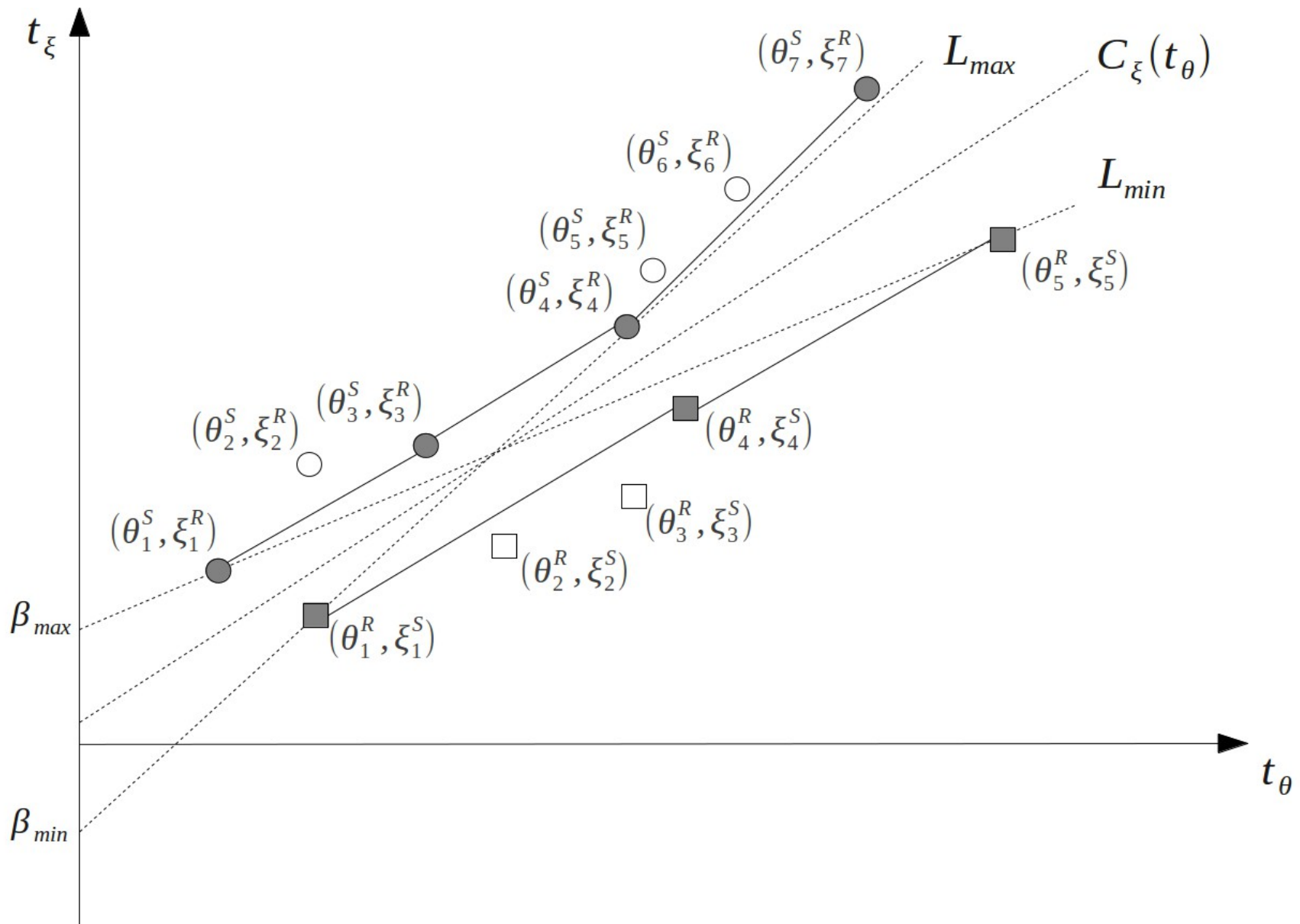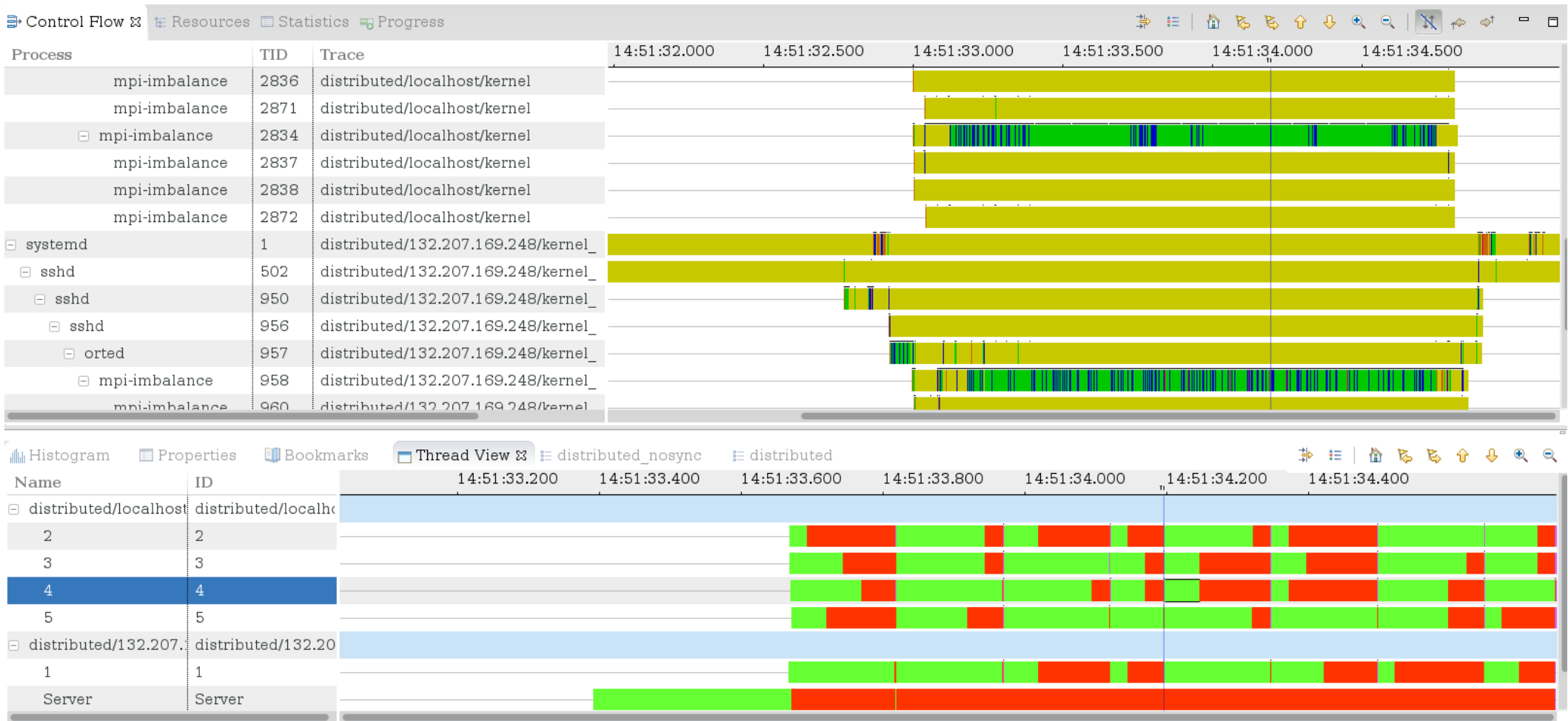
# Experiment 1: Local: Thread View

# Experiment 2: Distributed Network: Control Flow View and Worker View

# Convex-Hull Synchronization Algorithm

# Experiment 2: Distributed Network: Synchronized View

# Experiment 2: Distributed Network: PCap traces

Histogram    Properties    Bookmarks    Thread View    distributed_nosync    distributed    **distributed/tcpdump.out** ⊠

| Timestamp | Source | Destination | File | Protocol | Content |
|---|---|---|---|---|---|
| <srch> | <srch> | <srch> | <srch> | <srch> | <srch> |
| 14:51:32.769 934 000 | 00:22:4d:86:a8:09/132.207.72.9/56547 | d8:24:bd:90:00:40/132.207.169.248/22 | tcpdump.out | TCP | 56547 > 22 [SYN] Seq=4048612383 Len=40 |
| 14:51:32.770 525 000 | d8:24:bd:90:00:40/132.207.169.248/22 | 00:22:4d:86:a8:09/132.207.72.9/56547 | tcpdump.out | TCP | 22 > 56547 [SYN, ACK] Seq=383585901 Ack=4048612384 Len=40 |
| 14:51:32.770 558 000 | 00:22:4d:86:a8:09/132.207.72.9/56547 | d8:24:bd:90:00:40/132.207.169.248/22 | tcpdump.out | TCP | 56547 > 22 [ACK] Seq=4048612384 Ack=383585902 Len=32 |
| 14:51:32.770 816 000 | 00:22:4d:86:a8:09/132.207.72.9/56547 | d8:24:bd:90:00:40/132.207.169.248/22 | tcpdump.out | TCP | 56547 > 22 [ACK, PSH] Seq=4048612384 Ack=383585902 Len=32 |
| 14:51:32.771 089 000 | d8:24:bd:90:00:40/132.207.169.248/22 | 00:22:4d:86:a8:09/132.207.72.9/56547 | tcpdump.out | TCP | 22 > 56547 [ACK] Seq=383585902 Ack=4048612407 Len=32 |
| 14:51:32.784 895 000 | d8:24:bd:90:00:40/132.207.169.248/22 | 00:22:4d:86:a8:09/132.207.72.9/56547 | tcpdump.out | TCP | 22 > 56547 [ACK, PSH] Seq=383585902 Ack=4048612407 Len=32 |
| 14:51:32.784 980 000 | 00:22:4d:86:a8:09/132.207.72.9/56547 | d8:24:bd:90:00:40/132.207.169.248/22 | tcpdump.out | TCP | 56547 > 22 [ACK] Seq=4048612407 Ack=383585925 Len=32 |
| 14:51:32.785 382 000 | 00:22:4d:86:a8:09/132.207.72.9/56547 | d8:24:bd:90:00:40/132.207.169.248/22 | tcpdump.out | TCP | 56547 > 22 [ACK] Seq=4048612407 Ack=383585925 Len=32 |
| 14:51:32.785 386 000 | 00:22:4d:86:a8:09/132.207.72.9/56547 | d8:24:bd:90:00:40/132.207.169.248/22 | tcpdump.out | TCP | 56547 > 22 [ACK, PSH] Seq=4048613855 Ack=383585925 Len=32 |
| 14:51:32.786 736 000 | d8:24:bd:90:00:40/132.207.169.248/22 | 00:22:4d:86:a8:09/132.207.72.9/56547 | tcpdump.out | TCP | 22 > 56547 [ACK, PSH] Seq=383585925 Ack=4048612407 Len=32 |
| 14:51:32.786 756 000 | 00:22:4d:86:a8:09/132.207.72.9/56547 | d8:24:bd:90:00:40/132.207.169.248/22 | tcpdump.out | TCP | 56547 > 22 [ACK] Seq=4048614375 Ack=383587573 Len=32 |
| 14:51:32.786 759 000 | d8:24:bd:90:00:40/132.207.169.248/22 | 00:22:4d:86:a8:09/132.207.72.9/56547 | tcpdump.out | TCP | 22 > 56547 [ACK] Seq=383587573 Ack=4048614375 Len=32 |

local/localhost/ust    State System Explorer    local    **Stream List** ⊠

Ethernet II | **Internet Protocol Version 4** | Transmission Control Protocol | User Datagram Protocol

| ID | Endpoint A | Endpoint B | Packets | Bytes | Packets A -: | Bytes A -> I | Packets B -: | Bytes B -> / | Start Time |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 00:22:4d:86:a8:09/132.207.72.9 | d8:24:bd:90:00:40/132.207.169.248 | 361 | 42848 | 190 | 20689 | 171 | 22159 | 14:51:32.769 934 000 |
| 1 | 00:22:4d:86:a8:09/132.207.72.9 | d8:24:bd:90:00:40/74.125.226.134 | 2 | 132 | 1 | 66 | 1 | 66 | 14:51:33.898 411 000 |
| 2 | 00:22:4d:86:a8:09/132.207.72.9 | d8:24:bd:90:00:40/74.125.226.159 | 2 | 132 | 1 | 66 | 1 | 66 | 14:51:34.351 755 000 |
| 3 | d8:24:bd:90:00:40/132.207.180.14 | 00:22:4d:86:a8:09/132.207.72.9 | 6 | 540 | 3 | 228 | 3 | 312 | 14:51:34.483 763 000 |

# Experiment 2: Distributed Network: PCap stream filter



| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Histogram | Properties | Bookmarks | Thread View | distributed_nosync | **distributed** ✕ | distributed/tcpdump.out | | |

| Timestamp | Source | Destination | File | Protocol | Content | Channel | Type | Source |
|---|---|---|---|---|---|---|---|---|
| <srch> | <srch> | <srch> | <srch> | <srch> | <srch> | <srch> | <srch> | <srch> |
| 14:51:32.770 768 048 | | | channel0_0 | | sk=0xffff8800785b8000, seq=0xf150ec20, ack_seq=0x16dd0e6e, chec | channel0_0 | inet_sock_local_in | 0 |
| 14:51:32.770 813 569 | | | channel0_5 | | sk=0xffff8803486faa00, seq=0xf150ec20, ack_seq=0x16dd0e6e, chec | channel0_5 | inet_sock_local_out | 5 |
| 14:51:32.770 816 000 | 00:22:4 | d8:24:bd:90: | tcpdump.out | TCP | 56547 > 22 [ACK, PSH] Seq=4048612384 Ack=383585902 Len=32 | tcpdump.out | packet:tcp | linktype:ethernet |
| 14:51:32.770 951 238 | | | channel0_0 | | sk=0xffff8800785bdb00, seq=0xf150ec20, ack_seq=0x16dd0e6e, chec | channel0_0 | inet_sock_local_in | 0 |
| 14:51:32.770 967 022 | | | channel0_0 | | sk=0xffff8800785bdb00, seq=0x16dd0e6e, ack_seq=0xf150ec37, chec | channel0_0 | inet_sock_local_out | 0 |
| 14:51:32.771 089 000 | d8:24:b | 00:22:4d:86: | tcpdump.out | TCP | 22 > 56547 [ACK] Seq=383585902 Ack=4048612407 Len=32 | tcpdump.out | packet:tcp | linktype:ethernet |
| 14:51:32.771 099 672 | | | channel0_0 | | sk=0xffff8803486faa00, seq=0x16dd0e6e, ack_seq=0xf150ec37, chec | channel0_0 | inet_sock_local_in | 0 |
| 14:51:32.784 570 292 | | | channel0_0 | | sk=0xffff8800785bdb00, seq=0x16dd0e6e, ack_seq=0xf150ec37, chec | channel0_0 | inet_sock_local_out | 0 |
| 14:51:32.784 895 000 | d8:24:b | 00:22:4d:86: | tcpdump.out | TCP | 22 > 56547 [ACK, PSH] Seq=383585902 Ack=4048612407 Len=32 | tcpdump.out | packet:tcp | linktype:ethernet |
| 14:51:32.784 913 380 | | | channel0_0 | | sk=0xffff8803486faa00, seq=0x16dd0e6e, ack_seq=0xf150ec37, chec | channel0_0 | inet_sock_local_in | 0 |
| 14:51:32.784 973 742 | | | channel0_5 | | sk=0xffff8803486faa00, seq=0xf150ec37, ack_seq=0x16dd0e85, chec | channel0_5 | inet_sock_local_out | 5 |
| 14:51:32.784 980 000 | 00:22:4 | d8:24:bd:90: | tcpdump.out | TCP | 56547 > 22 [ACK] Seq=4048612407 Ack=383585925 Len=32 | tcpdump.out | packet:tcp | linktype:ethernet |

| | | | | | |
|---|---|---|---|---|---|
| State System Explorer | Stream List | **Filters** ✕ | distributed/localhost/kernel | | |

FILTER stream ipv4 00:22:4d:86:a8:09/132.207.72.9 <--> d8:24:bd:90:00:40/132.207.169.248

name: TCP between the 2 hosts

**FILTER TCP between the 2 hosts**
- OR
  - AND
    - Internet Protocol Version 4 CONTAINS
      - OR
        - AND
          - :packetsource: CONTAINS "00:22:4d:86:a8:09/132.207.72.9"
          - :packetdestination: CONTAINS "d8:24:bd:90:00:40/132.207.169.248"
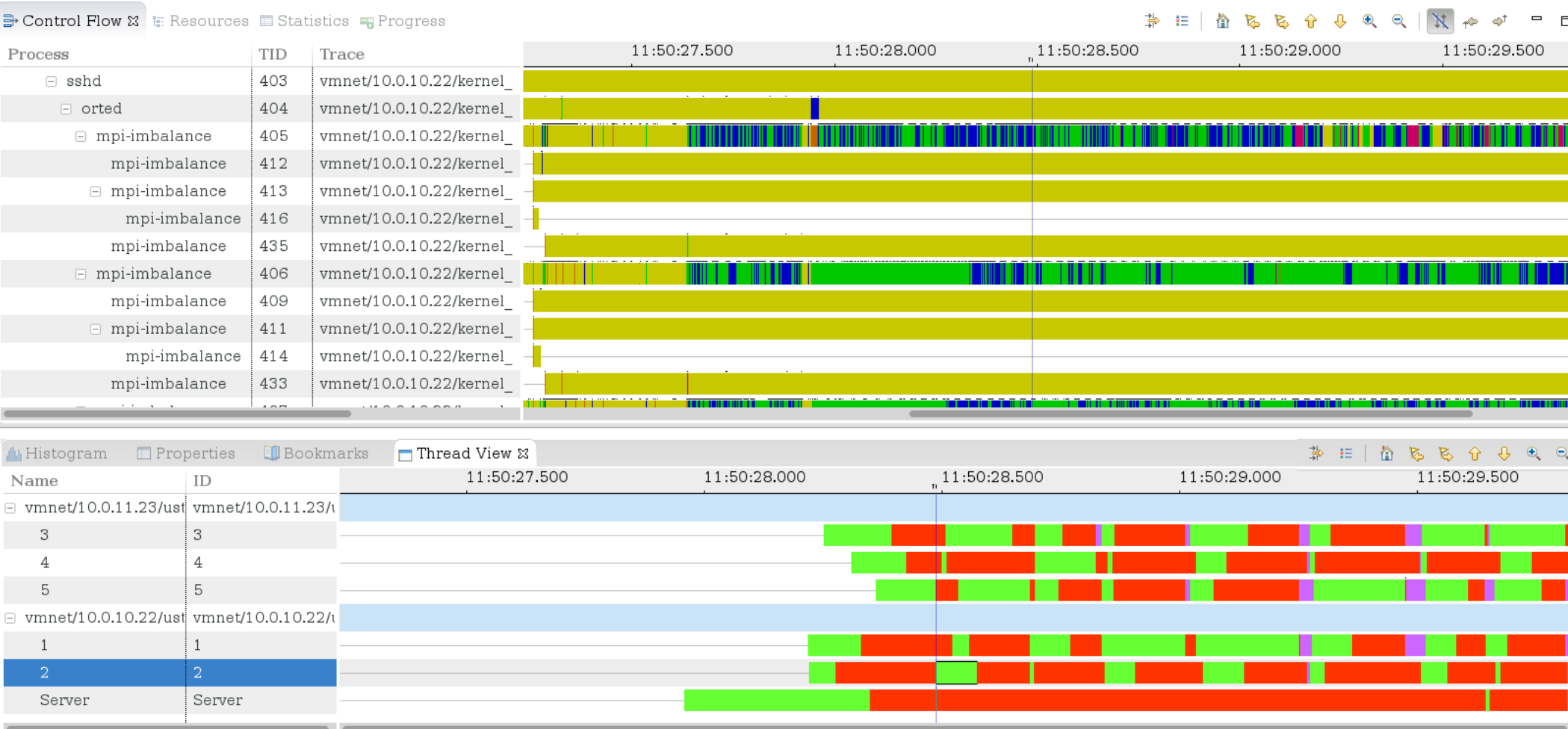        - AND
  - WITH EVENTTYPE Common Trace Format : LTTng Kernel Trace
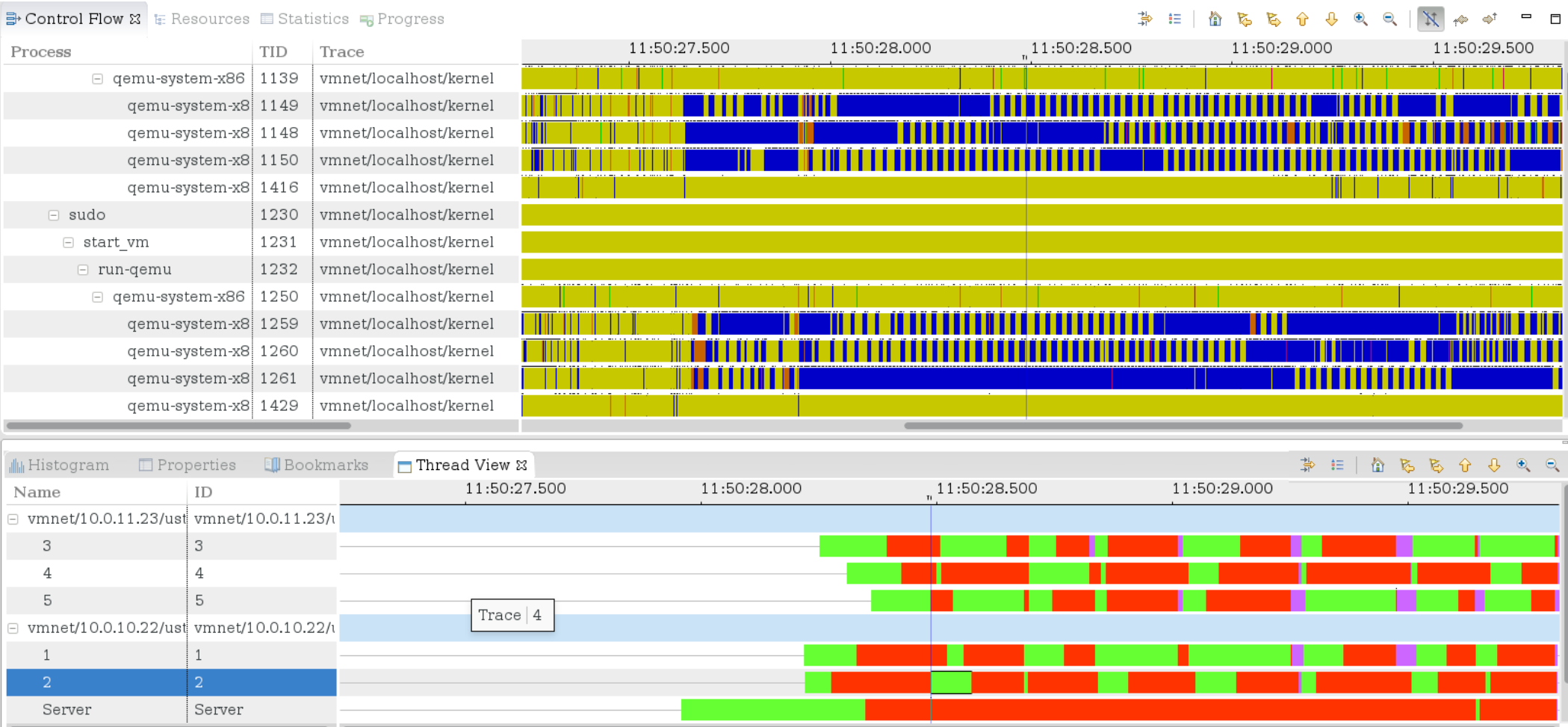    - OR
      - :type: CONTAINS "inet_sock_local_in"
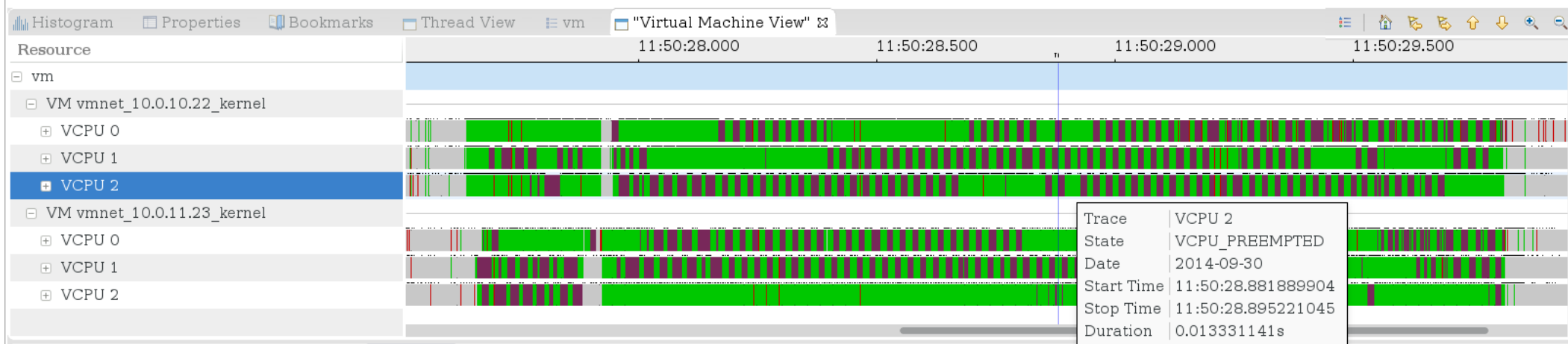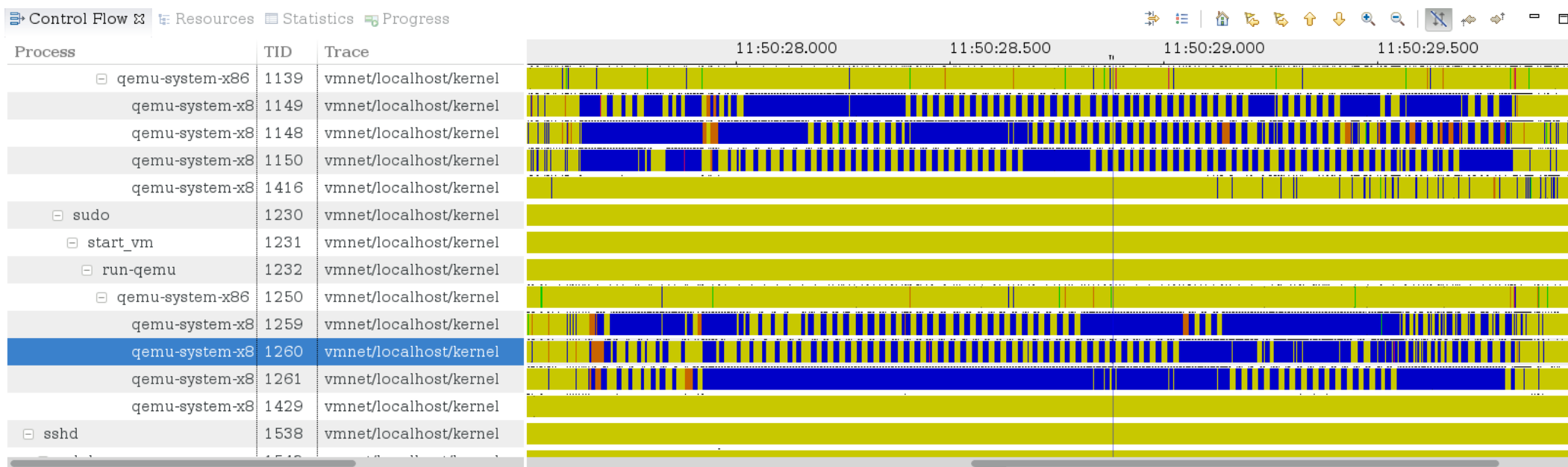      - :type: CONTAINS "inet_sock_local_out"

# Experiment 3: Virtual Machines: Control Flow View and Thread View

# Experiment 3: Virtual Machines: qemu processes view

# Experiment 3: Virtual Machines: VCPUs view

# Experiment 3: Virtual Machines: 1 VCPU view

# Experiment 3: Virtual Machines: VM Preempt View