



# Enhanced execution trace abstraction approach using social network analysis methods

Ji Wang, **Naser Ezzati-Jivan**

Oct 11, 2020



Symposium on Software Performance (SSP) 2020

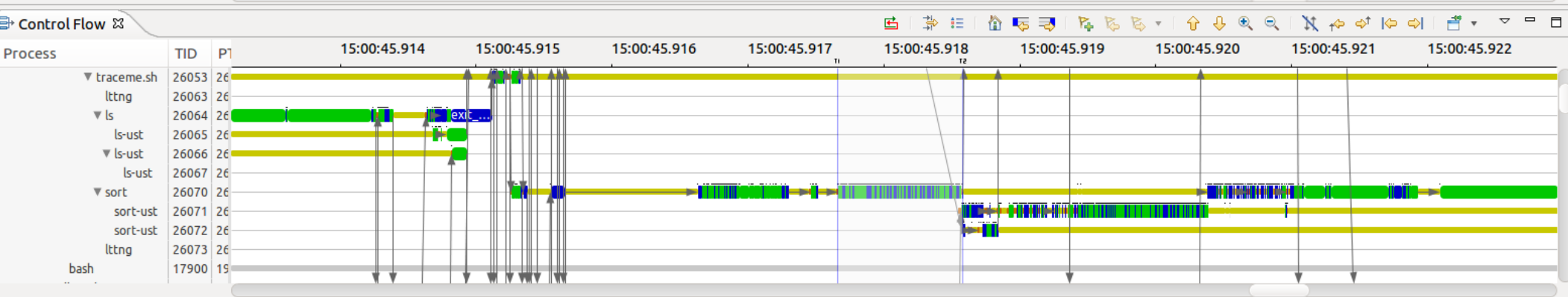
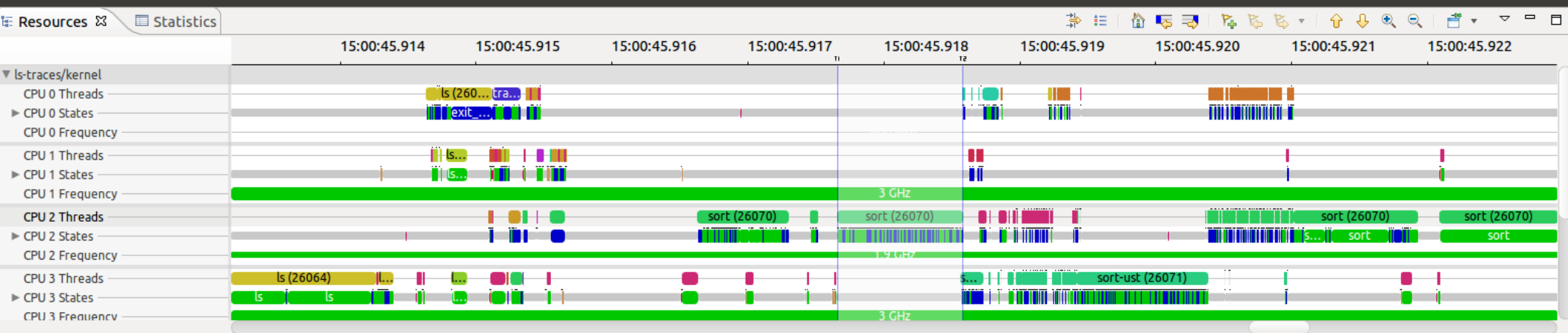
# Introduction

Execution tracing is used for performance analysis and debugging.

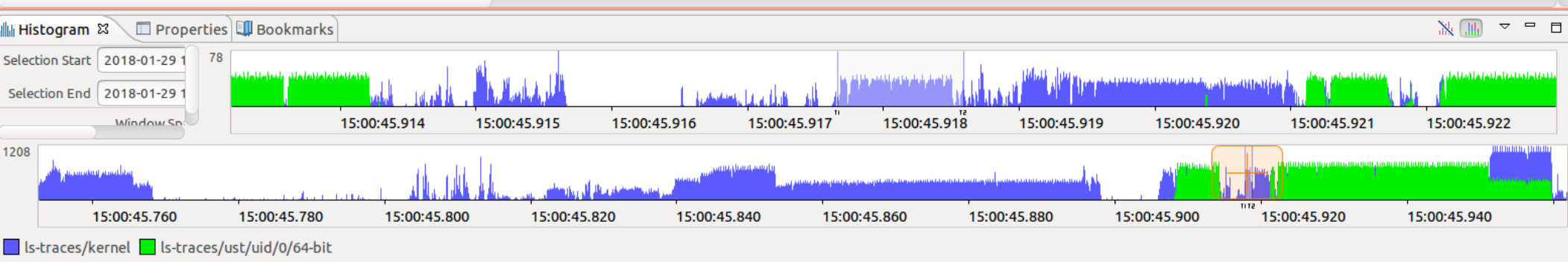
Usually the trace size is huge, and we need some heuristics to cope with it.

Process	TID	PTID	Birth time
gnome-pty-helpe	3609	3596	18:15:38.7695138
▼ bash	3610	3596	18:15:38.7695147
ssh	8972	3610	18:15:38.7695557
▼ bash	5332	3596	18:15:38.7695248
ls	15298	5332	18:15:39.9718121
▼ bash	10234	3596	18:15:38.7695260
ssh	8989	10234	18:15:38.7695567
▼ bash	11275	3596	18:15:38.7695713
▼ bash	15266	11275	18:15:38.7729050
▼ sudo	15291	15266	18:15:38.7729062
lttng	15292	15291	18:15:38.7729070
ab	15294	15266	18:15:38.8254856
▼ sudo	15302	15266	18:15:46.0692929
lttng	15303	15302	18:15:46.0735909
dconf worker	3599	2779	18:15:38.7695104
gmain	3600	2779	18:15:38.7695113
gdbus	3601	2779	18:15:38.7695120

Process	TID	PTID	Birth time
unity-greeter	1082	1070	12:21:27.6053039
llvmpipe-0	1159	1070	12:21:27.6053052
llvmpipe-1	1160	1070	12:21:27.6053064
gmain	1175	1070	12:21:27.6053082
gdbus	1176	1070	12:21:27.6053093
dconf worker	1193	1070	12:21:27.6053111
threaded-ml	1295	1070	12:21:27.6053126
gmain	1037	851	12:21:27.6052899
gdbus	1038	851	12:21:27.6052916
lightdm	1200	851	12:21:27.6053416
gmain	872	1	12:21:27.6052339
gdbus	874	1	12:21:27.6052352
▼ mysqld	861	1	12:21:27.6052366
mysqld	1545	861	12:21:27.6665481
mysqld	941	1	12:21:27.6052379
mysqld	955	1	12:21:27.6052392
mysqld	956	1	12:21:27.6052408



Trace	Timestamp	Channel	CPU	Event type	Contents
<srch>	<srch>	<srch>	<srch>	<srch>	<srch>
ls-traces/kernel	2018-01-29 15:00:45.876 797 081	kernel_1	1	kmem_mm_page_alloc	page=0xffffea0009902680, pfn=2506906, order=0, gfp_flags=131290, migratetype=2, context._perf_cpu_mi
ls-traces/kernel	2018-01-29 15:00:45.876 797 698	kernel_1	1	syscall_exit_write	ret=4096, context._perf_cpu_migrations=23
ls-traces/kernel	2018-01-29 15:00:45.876 797 918	kernel_1	1	syscall_entry_write	fd=28, buf=139984728244528, count=4096, context._perf_cpu_migrations=23
ls-traces/kernel	2018-01-29 15:00:45.876 798 223	kernel_1	1	kmem_mm_page_alloc	page=0xffffea000dce7940, pfn=3619301, order=0, gfp_flags=131290, migratetype=2, context._perf_cpu_mi
ls-traces/kernel	2018-01-29 15:00:45.876 799 528	kernel_1	1	syscall_exit_write	ret=4096, context._perf_cpu_migrations=23
ls-traces/kernel	2018-01-29 15:00:45.876 799 754	kernel_1	1	syscall_entry_write	fd=28, buf=139984728244528, count=4096, context._perf_cpu_migrations=23
ls-traces/kernel	2018-01-29 15:00:45.876 800 050	kernel_1	1	kmem_mm_page_alloc	page=0xffffea000e078880, pfn=3677730, order=0, gfp_flags=131290, migratetype=2, context._perf_cpu_mi

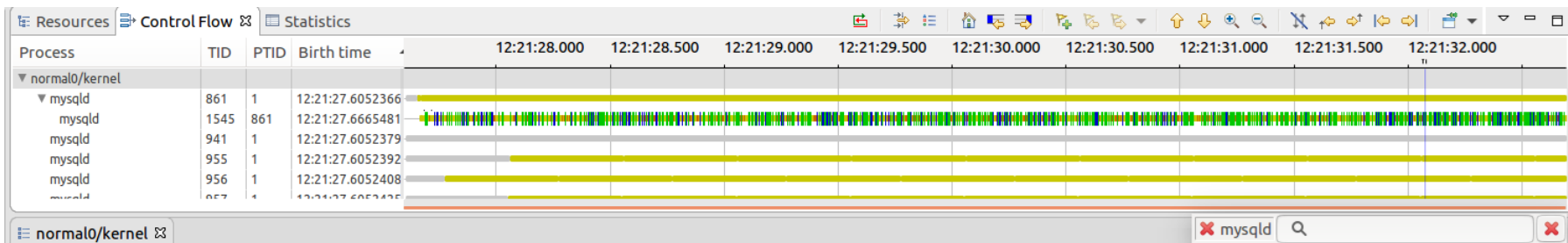


# Current methods



- ✓ Sampling or event filtering.
- ✓ Time view filtering.
- ✓ Global filtering

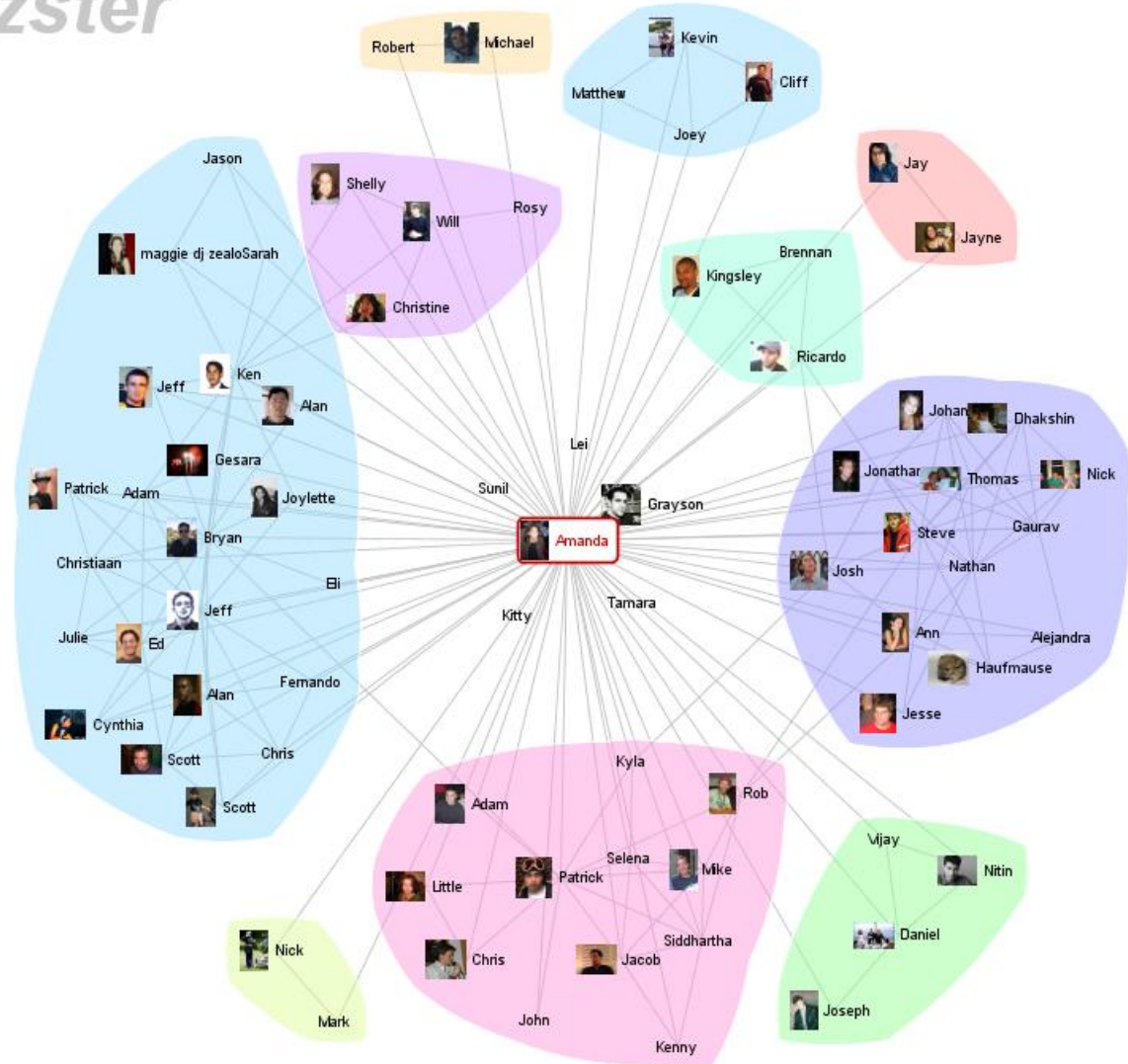
Timestamp	Channel	CPU	Event type	Contents
<srch>	<srch>	<srch>	<srch>	<b>tid=861</b>
12:21:27.650 850 349	more-subbuf_0	0	syscall_entry_connect	fd=3, servaddr=0x7ffe5564c030, addrlen=0x6e, family=1, dport=0, v4addr_length=0, v4addr=[], v6addr_length=0, v6addr=[], contex
12:21:27.650 863 113	more-subbuf_0	0	sched_waking	comm=mysql, tid=861, prio=20, target_cpu=0, context_tid=1544, context_pid=1544, context_procname=sysbench
12:21:27.650 868 663	more-subbuf_0	0	sched_wakeup	comm=mysql, tid=861, prio=20, target_cpu=0, context_tid=1544, context_pid=1544, context_procname=sysbench
12:21:27.650 870 198	more-subbuf_0	0	syscall_exit_connect	ret=0, context_tid=1544, context_pid=1544, context_procname=sysbench
12:21:27.650 873 554	more-subbuf_0	0	sched_switch	prev_comm=sysbench, prev_tid=1544, prev_prio=20, prev_state=0, next_comm=mysql, next_tid=861, next_prio=20, context_tid=1544
12:21:27.650 886 455	more-subbuf_0	0	syscall_exit_poll	ret=1, nfds=2, fds_length=1, fds=[[fd=23, raw_events=0x1, events_POLLIN=1, events_POLLPRI=0, events_POLLOUT=0, events_POLL
12:21:27.650 895 764	more-subbuf_0	0	syscall_entry_accept	fd=23, upeer_sockaddr=0x7fff9470d50, upeer_addrlen=128, context_tid=861, context_pid=861, context_procname=mysql
12:21:27.650 904 076	more-subbuf_0	0	syscall_exit_accept	ret=39, upeer_addrlen=2, family=1, sport=0, v4addr_length=0, v4addr=[], v6addr_length=0, v6addr=[], context_tid=861, context_pid=861
12:21:27.650 918 648	more-subbuf_0	0	syscall_entry_futex	uaddr=31402148, op=133, val=1, utime=1, uaddr2=31402144, val3=67108865, context_tid=861, context_pid=861, context_procname=mysql
12:21:27.650 923 327	more-subbuf_0	0	sched_waking	comm=mysql, tid=1291, prio=20, target_cpu=0, context_tid=861, context_pid=861, context_procname=mysql
12:21:27.650 926 771	more-subbuf_0	0	sched_wakeup	comm=mysql, tid=1291, prio=20, target_cpu=0, context_tid=861, context_pid=861, context_procname=mysql
12:21:27.650 928 588	more-subbuf_0	0	syscall_exit_futex	ret=1, uaddr=31402148, uaddr2=31402144, context_tid=861, context_pid=861, context_procname=mysql



# Idea

## Graph theory

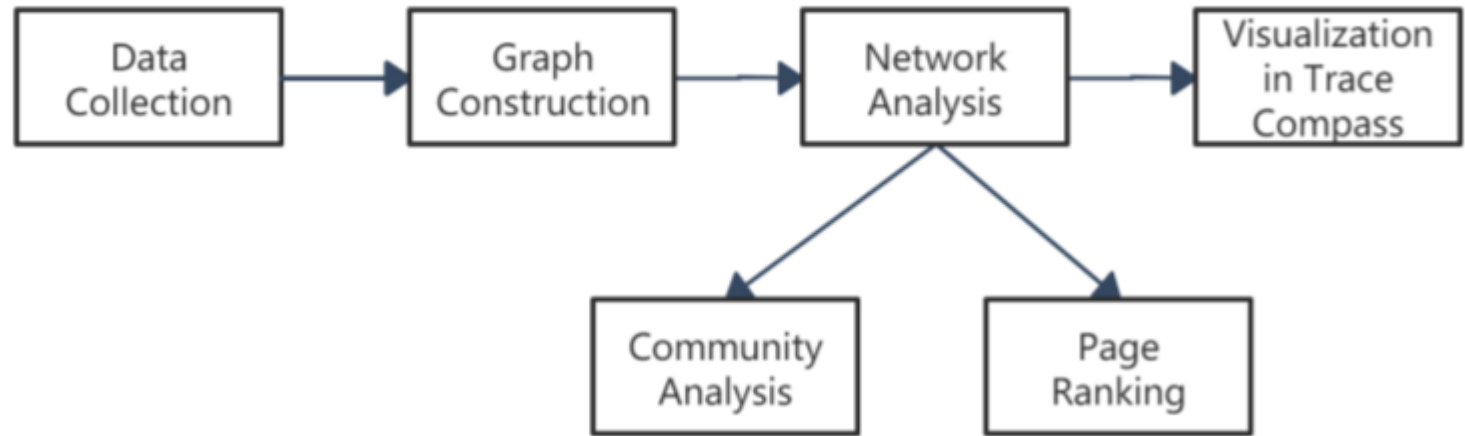
*vizster*



Picture from: <http://vis.stanford.edu/jheer/projects/vizster/>



# Method



I.

✓ Extract the processes interactions based on the execution graph.

✓ Building a graph  $G = (V, E)$

✓  $V$  includes all the processes running on the system

✓  $E$  is a set of edges defined based on the interactions between the processes.

or kernel level lock methods.

**Algorithm 1** Execution Graph Construction algorithm, for events collected from lock library .

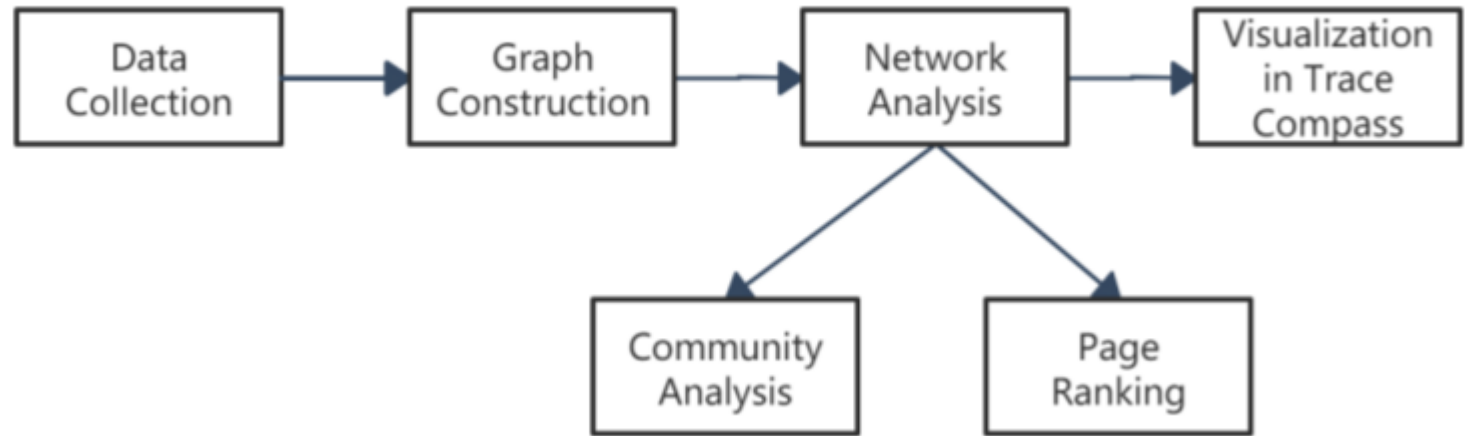
**Input Input:** Trace  $T$ , Threads =  $\{t_1, t_2, \dots, t_n\}$

Set1 =  $\{*_lock\_req\}$ ,  
Set2 =  $\{*_lock\_acq\}$ ,  
Set3 =  $\{*_unlock\}$

**Output Output:** execution graph  $G$

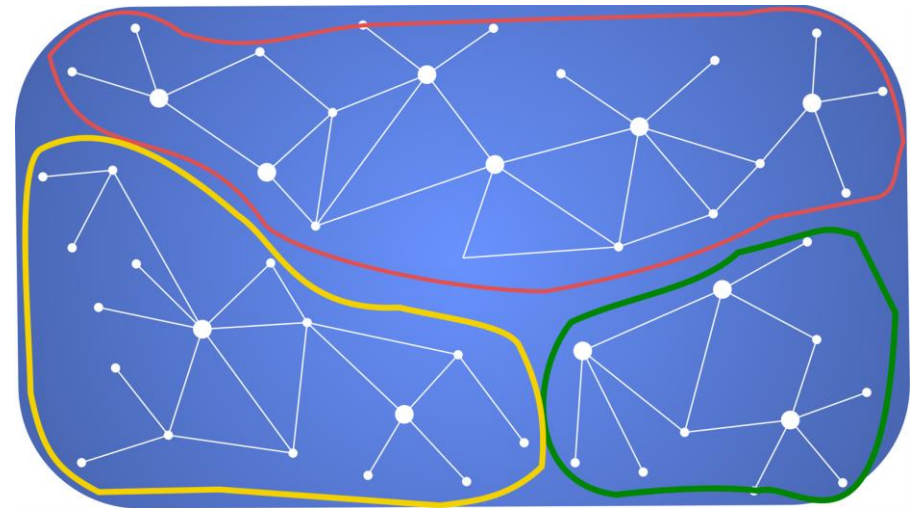
```
1:  $THREADS \leftarrow$  initial set of threads active in the multi-  
thread application Declarations  
2:  $LOCKHOLDER \leftarrow \emptyset$   
3: for all thread  $t \in THREADS$  do // Initialization  
4:   Create initial vertex of thread  $t$  with timestamp  $t.begin$   
5: end for  
6: for all event  $e \in T$  do // Mainprocedure  
7:   if  $e \in Set1$  then  
8:     new_h_edge( $e.tid, e.ts, blocked$ )  
9:     new_v_edge( $e.tid, LOCKHOLDER$ )  
10:   end if  
11:   if  $e \in Set2$  then  
12:     new_v_edge( $LOCKHOLDER, e.tid$ )  
13:     new_h_edge( $e.tid, e.ts, running$ )  
14:      $LOCKHOLDER \leftarrow e.tid$   
15:   end if  
16:   if  $e \in Set3$  then  
17:      $LOCKHOLDER \leftarrow e.tid$   
18:   end if
```

# Method



## II.

Compute connected components on the graph to extract all the subgraphs. (e.g., BFS)

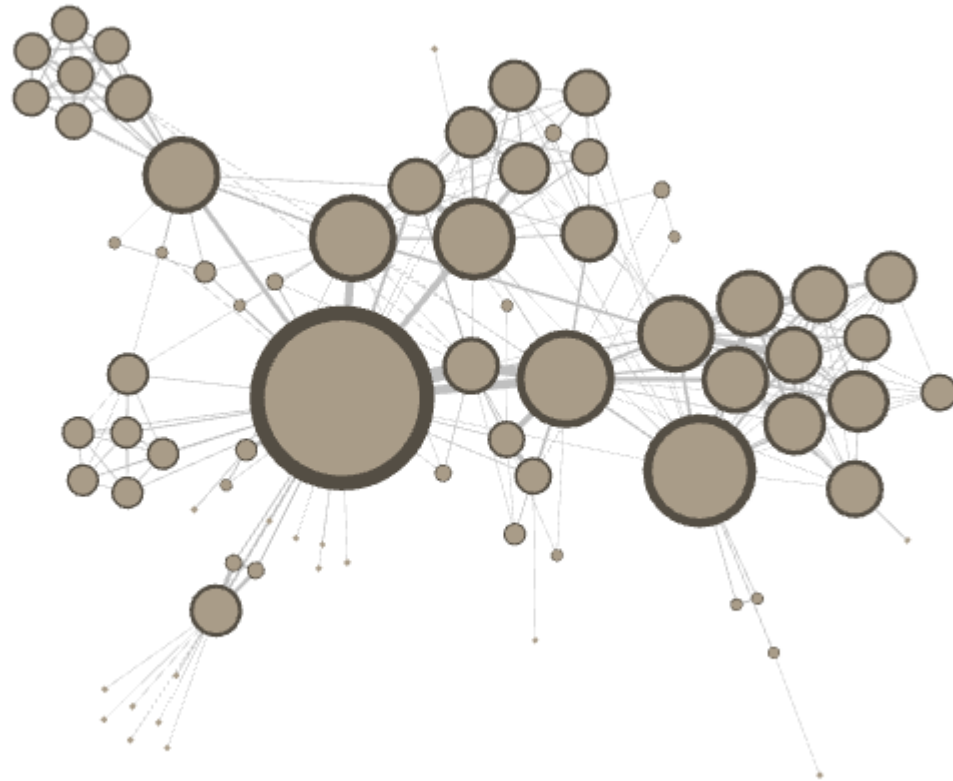


# Method

---

NB. Computing connected components on the graph

Other approaches:  
e.g., Graph Modularity



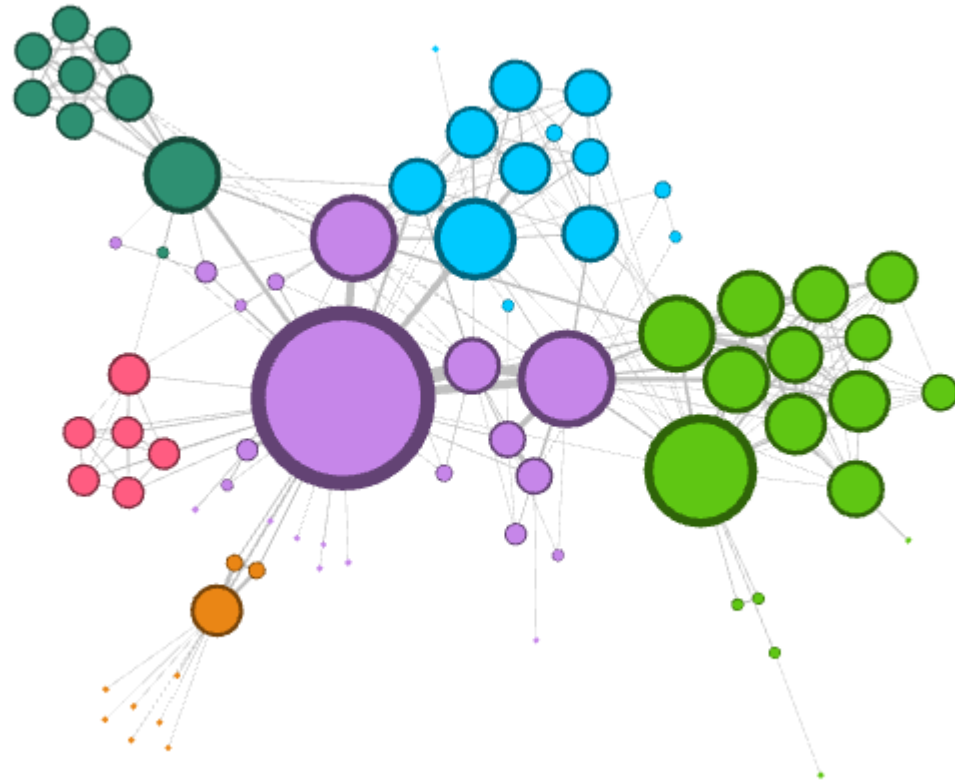


# Method

---

NB. Computing connected components on the graph

Other approaches:  
e.g., Graph Modularity

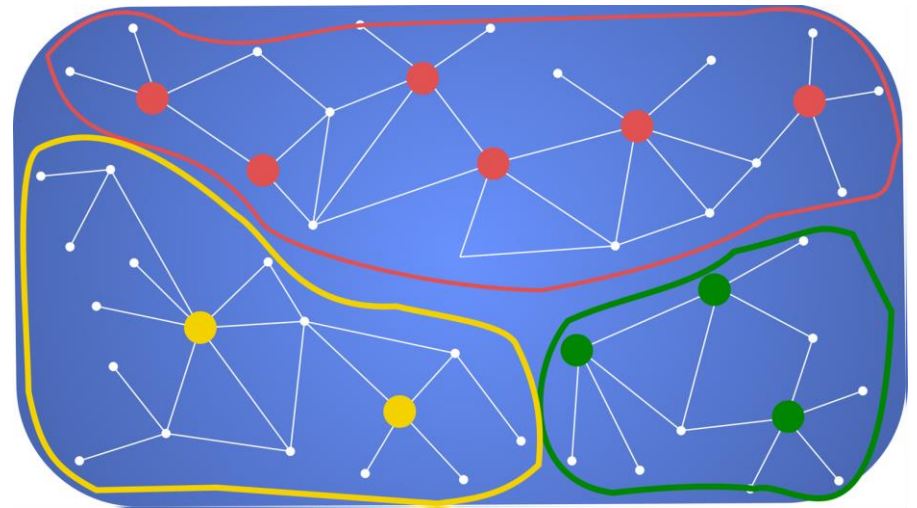


# Method

---

## III.

- ✓Apply a ranking algorithm on each connected component to assign a rank to each process
- ✓Choose the top-K processes..



# Method

---

**Modified PageRank** algorithm: Includes only **one** iteration

Node rank initialization: equal weights, or even CPU/IO usage

$$R(u) = c \sum_{v \in B_u} \frac{R(v)}{N_v}$$

$u$ : a process

$B_u$ : the set of  $u$ 's in-edges

$N_v$ : the number of out-edges of process  $v$

$c$ : the normalization factor

---

**Algorithm 2:** Process Ranking Algorithm (PRA)

---

```
1 if event == sched_ttwu then
2   | j = getVMvCPU(wakee_tid);
3   | k = getVMvCPU(waker_tid);
4   | wakerCR3 = getLastCR3(vCPUk);
5   | updateWaker(vCPUj, wakerCR3);
6 else if event == vm_inj_virq then
7   | j = getVMvCPU(tid);
8   | pCR3 = getVMvCPU(vCPUj);
9   | if vec == IPI then
10  |   | wakerCR3 = queryWaker(vCPUj);
11  |   | LINK_HORIZONTAL(wakerCR3, pCR3);
12 for all process cr3i ∈ CR3 do
13   | R(cr3i);
14 connected_subgraph = breadth-first-search(CR3);
15 customized_graph = denoise(connected_subgraph, R(CR3));
```

---

# Use case I

Trace Compass:

Extract process interactions using EASE scripting.

Providing a view showing top-ranked processes in each subgraph.

Filtering top-ranked processes in **event** view.

The screenshot displays the Eclipse IDE interface for Trace Compass. The main window is titled "runtime-EclipseApplication - Tracing/extractGraphInfo.js - Eclipse Platform". The interface is divided into several panes:

- Project Explorer:** Shows a tree view of the project structure, including "Active Thread", "Context CallStacks", "Context switch", "Counters", "CPU usage", "File Access", "Futex Contention Analysis", "Input/Output", "IRQ Analysis", "Kernel memory usage", "Linux Kernel", "OS Execution Graph", "Critical Flow View", "Scripted analyses", "State Machine Automatic Analy", "State Machine Backend State Sy", "Statistics", "System Call Latency", "System call stats", "Test a builtin XML module file", "External Analyses", "Reports", "extractGraphInfo.js", "sample.js", and "test.js".
- Control:** A control panel with a "type filter text" input field.
- Process Table:** A table listing processes with columns for Process, TID, PTID, and Birth time. The data is as follows:

Process	TID	PTID	Birth time
mysqld	1211	1	18:15:38.765
mysqld	1212	1	18:15:38.765
mysqld	1213	1	18:15:38.765
mysqld	1214	1	18:15:38.765
mysqld	14241	1	18:15:38.765
mysqld	14316	1	18:15:38.765
mysqld	14339	1	18:15:38.765
mysqld	14368	1	18:15:38.765
mysqld	14369	1	18:15:38.765
mysqld	14371	1	18:15:38.765
mysqld	14372	1	18:15:38.765
mysqld	14373	1	18:15:38.765
sshd	1096	1	18:15:38.765
dnsmasq	2265	1	18:15:38.765
dnsmasq	2266	2265	18:15:38.765
- Scripted Analyses:** A code editor showing the EASE script for "extractGraphInfo.js":

```
// normal0/kernel  ranking-apache-sql/kernel  extractGraphInfo.js
// load additional UI commands
loadModule('/TraceCompass/Analysis');
loadModule('/TraceCompass/View');

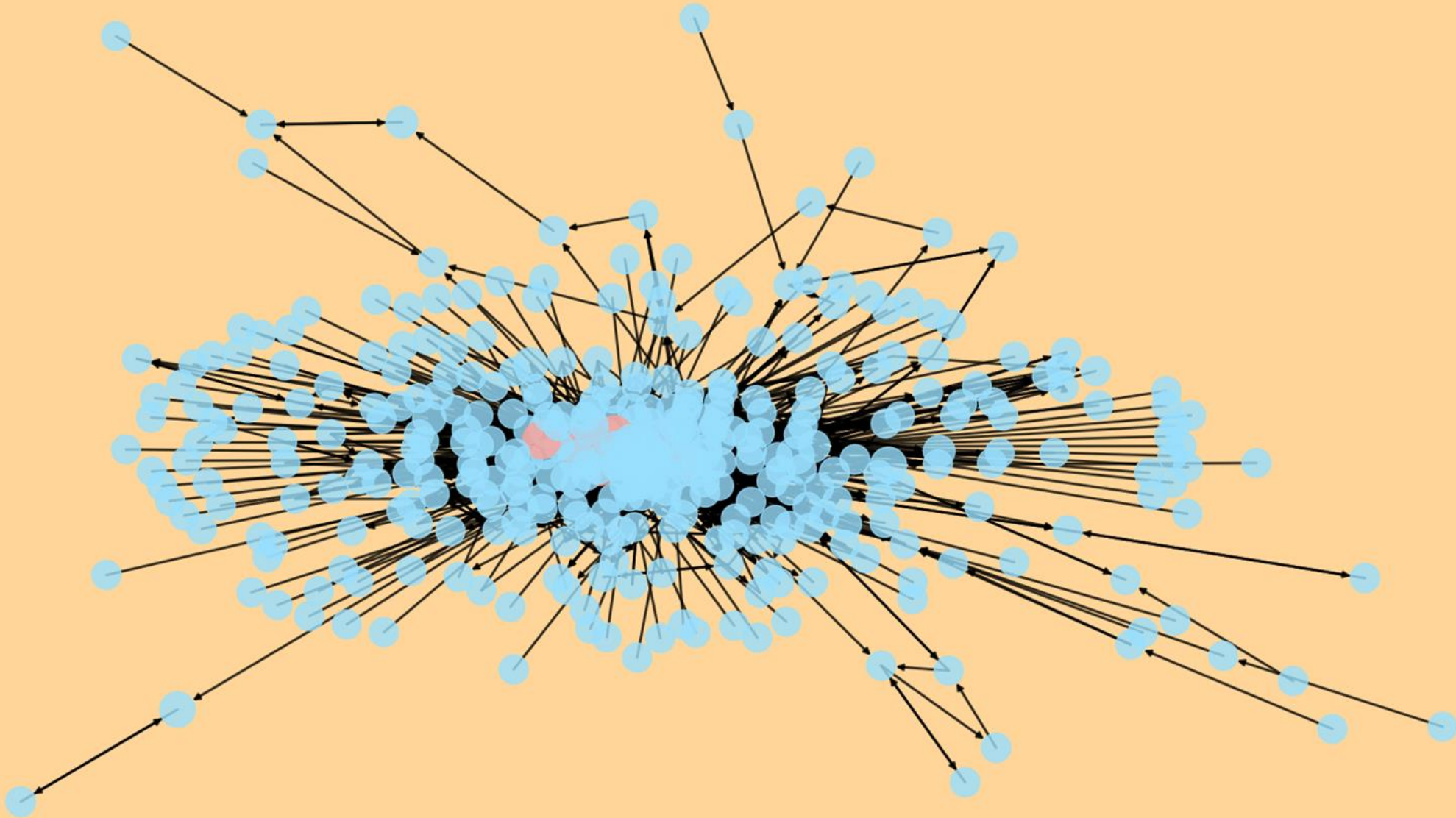
// open message box
var analysis = getAnalysis('activetid.js');
if (analysis == null) {
    print("Trace is null");
    exit();
}

var ss = analysis.getStateSystem(false);
var stackString = "";
var session = 0;
var dict = {};
var sessions = {};
var processes = {};
var pids = {};
var names = {};
```
- Process Interaction View:** A timeline view showing process interactions over time, with a table below it:

Process	Elapsed	Pi
"/usr/bin/termin.3596"	10.586743570	14
[Xorg.2771]	17.696122680	24
ProcessThread.77761	10.586743570	14
	0.791744436	14

# Use case I

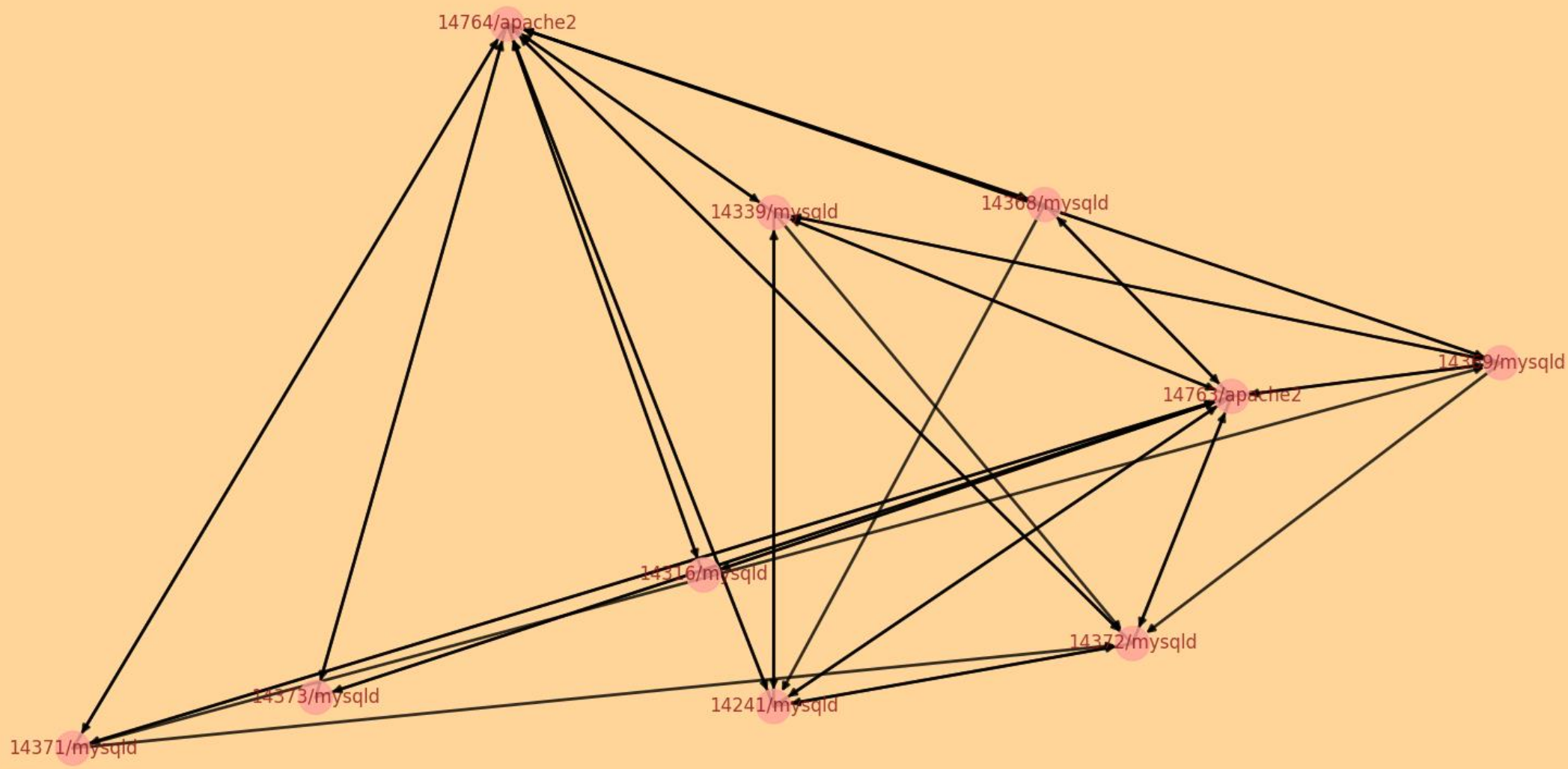
---



# Use case I

---

Top-ranked processes in subgraph I

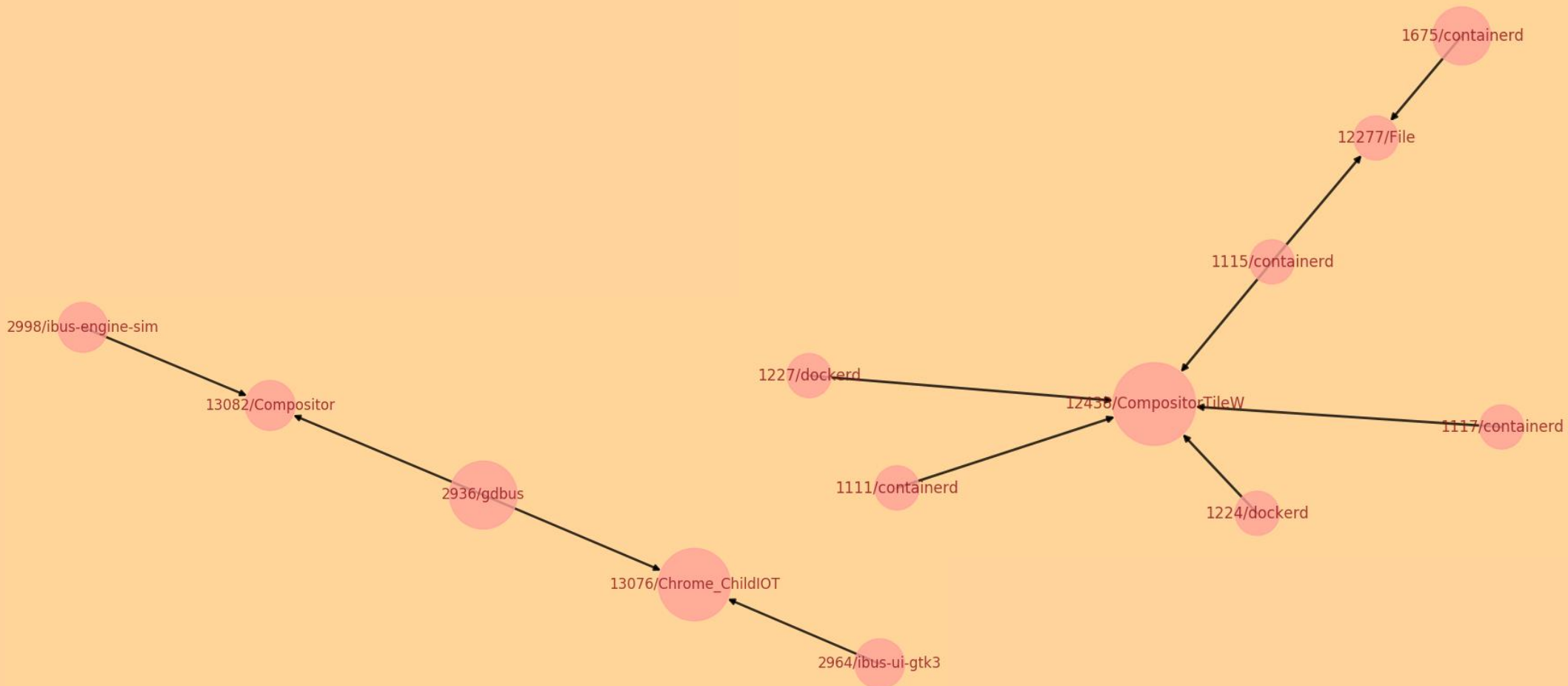




# Use case I

---

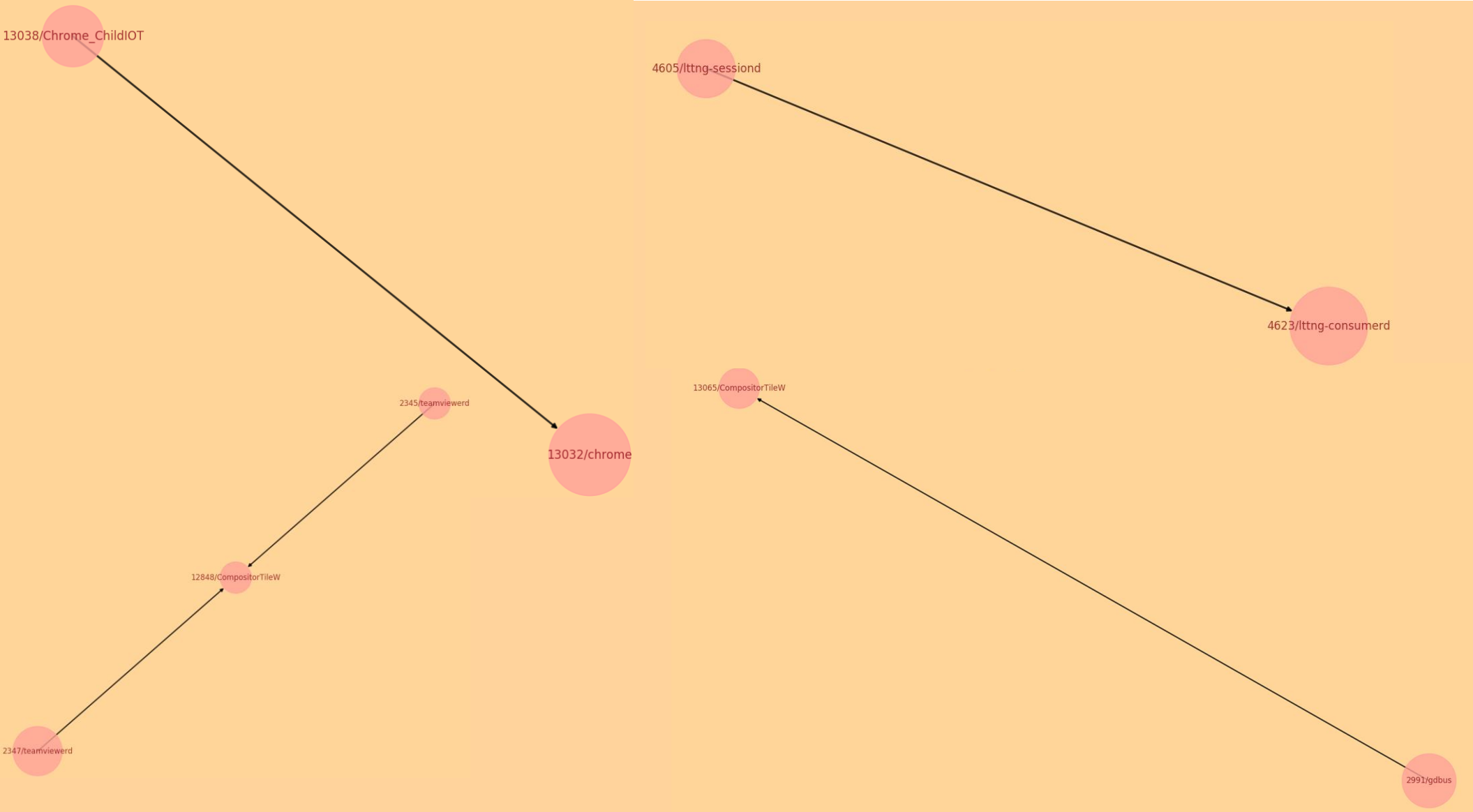
Top-ranked processes in subgraph II and III

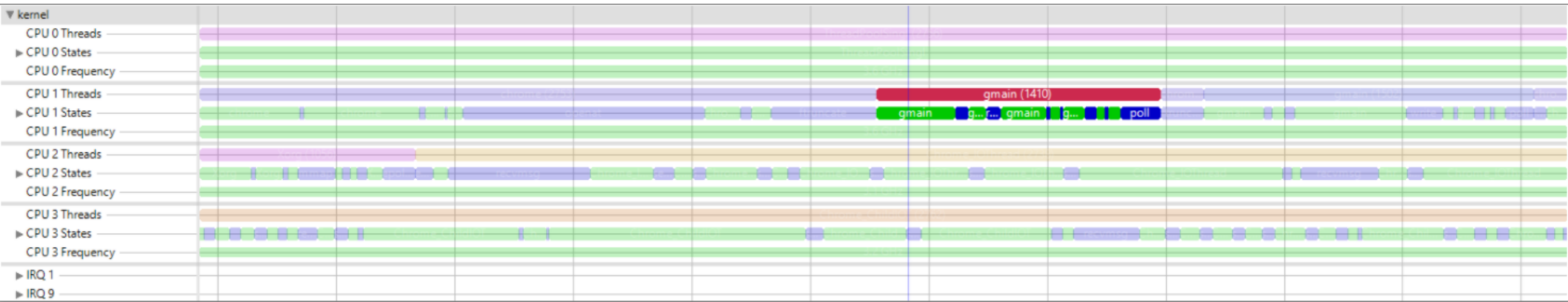
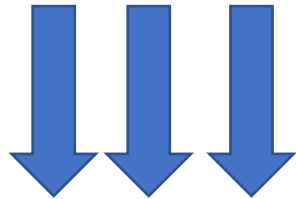


# Use case I

---

Top-ranked processes in other subgraphs:





# Conclusions & Future work

---

- Filtering based on process ranking
- Top interacting processes in one system
- Interaction between VM's processes in a cloud environment
- Time-window feature extraction for machine learning applications
- Enhance graph structure analysis, e.g., graph modularity

# Questions

---

*Thanks for your attention!*  
*nezzati@brocku.ca*