

# KernelShark 1.0

What's new, and what's coming?

Steven Rostedt  
9/26/2018

vmware®

© 2016 VMware Inc. All rights reserved.

# Correction:

- KernelShark 1.0 was suppose to be released in August
- Still not out yet
- Why?
  - Prototype was written first to get a good idea about the new design
  - The patches were broken up and reviewed in sets
  - Lots of updates were done (we want to get it right!)
    - Rewrite of designs
    - Breaking it up showed better ways to get things done
  - Went on tangents
    - Oh, we could do this too...
    - We are now at: “Just add a TODO comment, and we’ll comeback later”

# What is KernelShark?



**KERNELSHARK**



# What is KernelShark?

- It is a graphical interface to trace-cmd



**KERNELSHARK**



# What is KernelShark?

- It is a graphical interface to trace-cmd
- What is trace-cmd?



**KERNELSHARK**

# What is KernelShark?

- It is a graphical interface to trace-cmd
- What is trace-cmd?
  - A command line interface to ftrace



**KERNELSHARK**

# What is KernelShark?

- It is a graphical interface to trace-cmd
- What is trace-cmd?
  - A command line interface to ftrace
- What is ftrace?



**KERNELSHARK**

# What is KernelShark?

- It is a graphical interface to trace-cmd
- What is trace-cmd?
  - A command line interface to ftrace
- What is ftrace?
  - The official tracing infrastructure of the Linux kernel
  - See any of my other tracing talks



**KERNELSHARK**

# Why update KernelShark?

- Current Version 0.2
  - Written in GTK+ 2.0
  - Needs an update to GTK+ 3.0
  - New version written in Qt (Pronounced “cute”, I am told)
  - Was written as an “idle task”
    - Worked on it between projects
    - A couple of days at a time (may be months between)
  - Now have a full time employee (thanks to VMware)
  - Can make it much more flexible



**KERNELSHARK**

# Why update KernelShark?

- Current Version 0.2
  - Written in GTK+ 2.0
  - Needs an update to GTK+ 3.0
    - Did not want to rewrite again when GTK+ 4.0 comes out
  - New version written in Qt (Pronounced “cute”, I am told)
  - Was written as an “idle task”
    - Worked on it between projects
    - A couple of days at a time (may be months between)
  - Now have a full time employee (thanks to VMware)
  - Can make it much more flexible



**KERNELSHARK**

# Why was KernelShark created?

- trace-cmd can collect a large amount of data
- Hard to visualize from looking at text
- Real example:
  - cyclictst (jitter detecting program) had 250us latency on RT kernel
    - On a totally idle machine (only running cyclictst!)
  - hwlat reported the latency
    - But the HW vendor did not believe it
  - Used trace-cmd to find demonstrate the issue
  - Showed the HW latency with trace-cmd



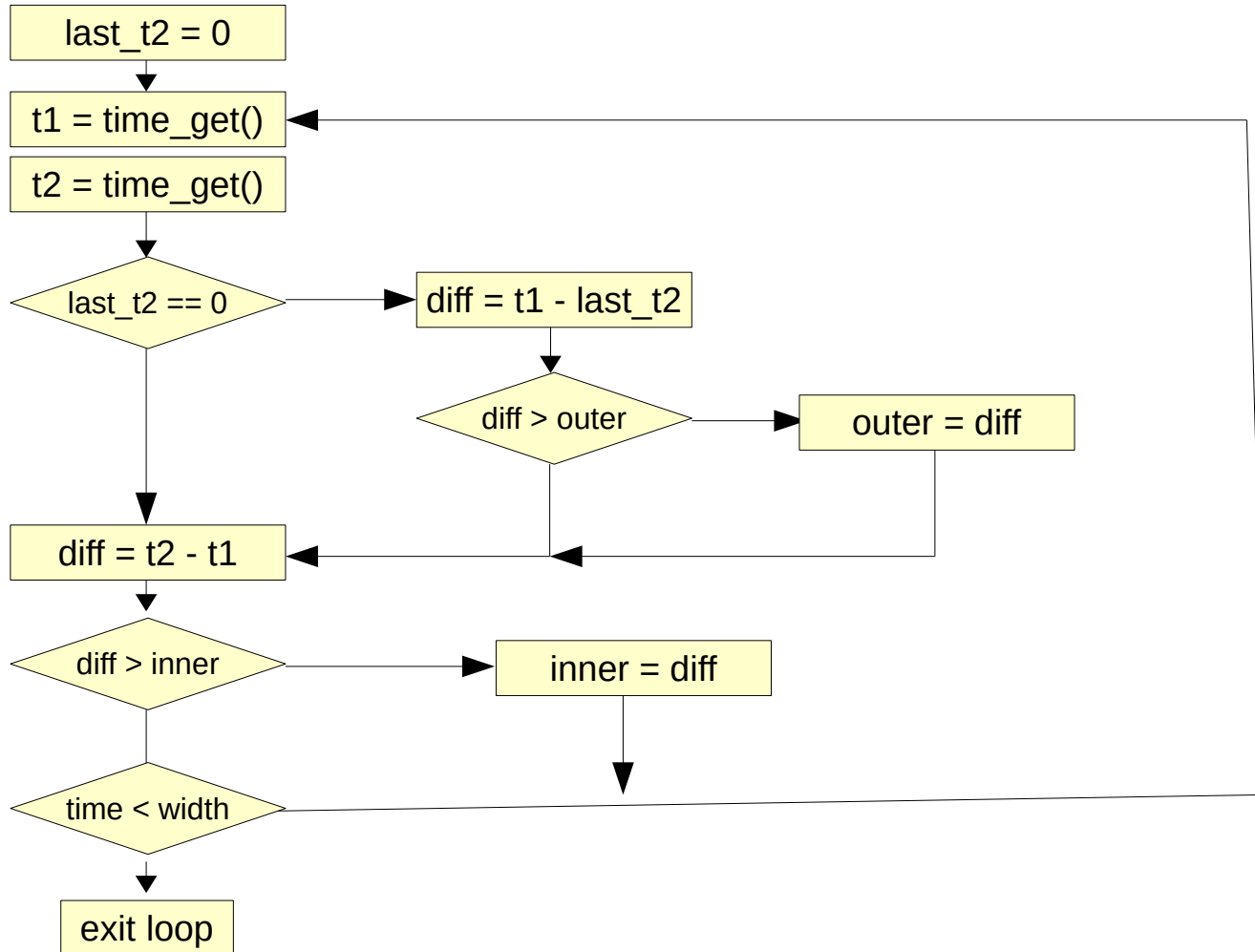
**KERNELSHARK**

# Cyclictest

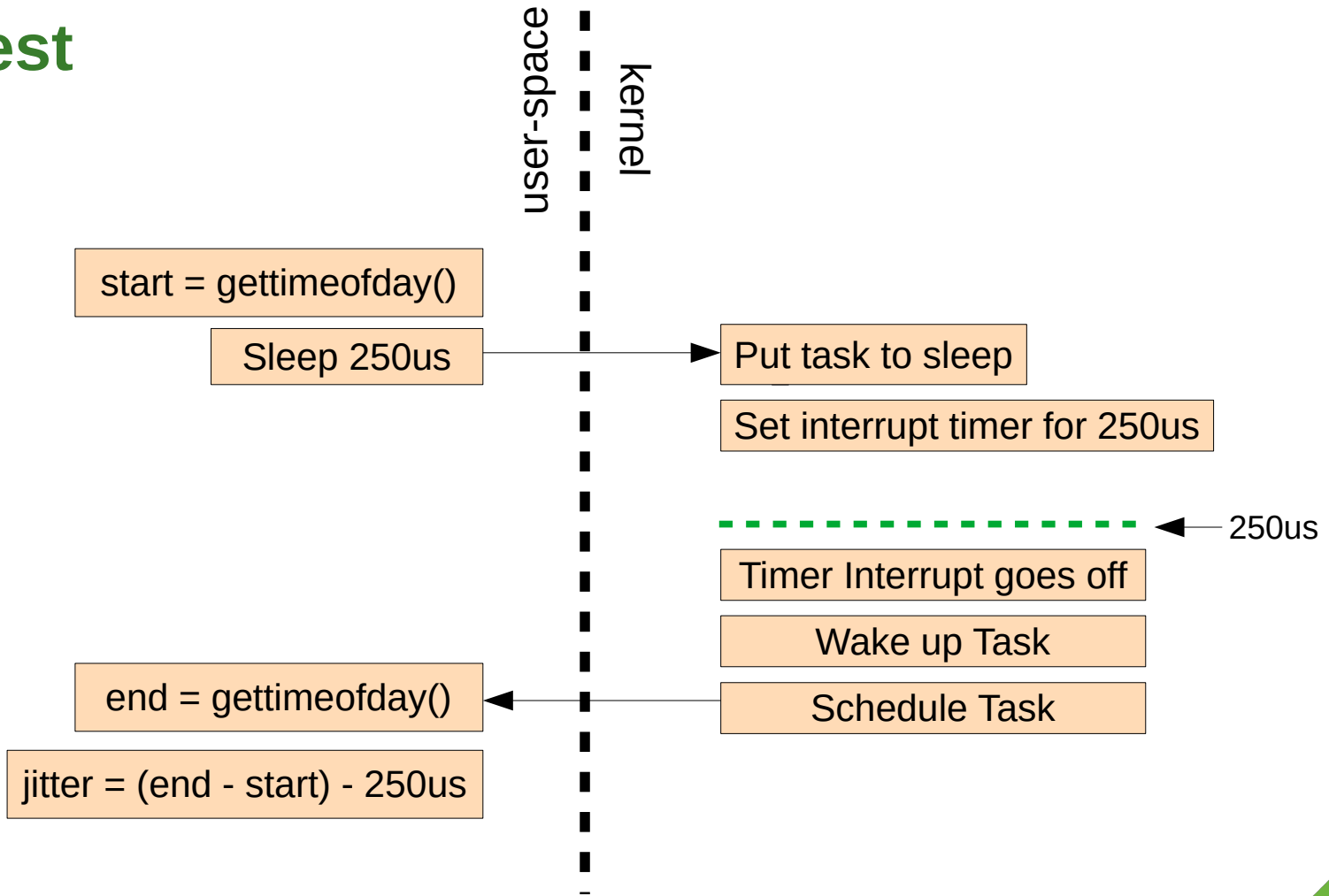
(Recap for those not at Embedded Recipes)

- A tool to measure latency
- Runs a simple loop (in a high priority task)
  - Sleep for a specified time
  - Get timestamp when wakes up
  - Compare the difference
- Favorite application of the Linux RT folks

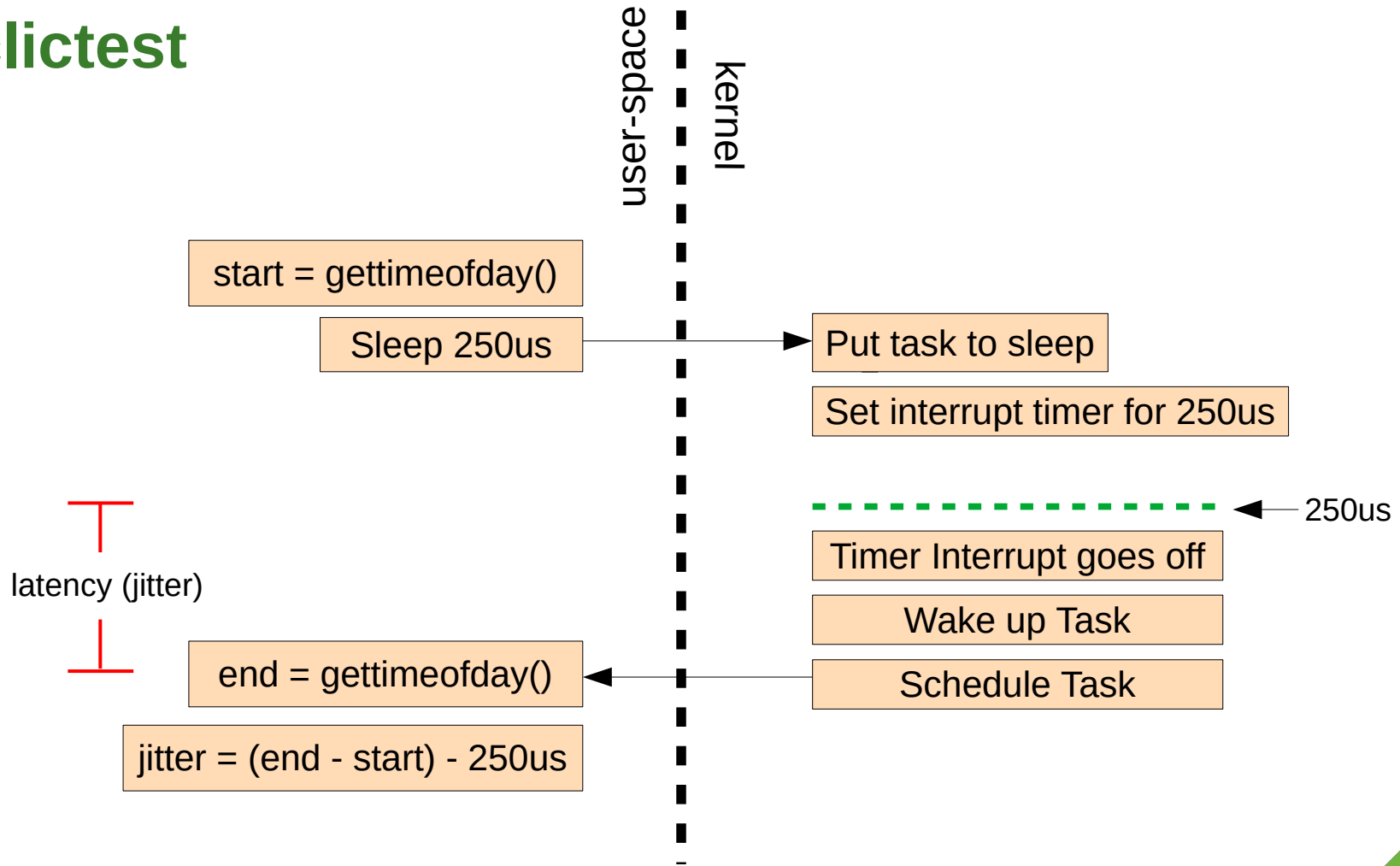




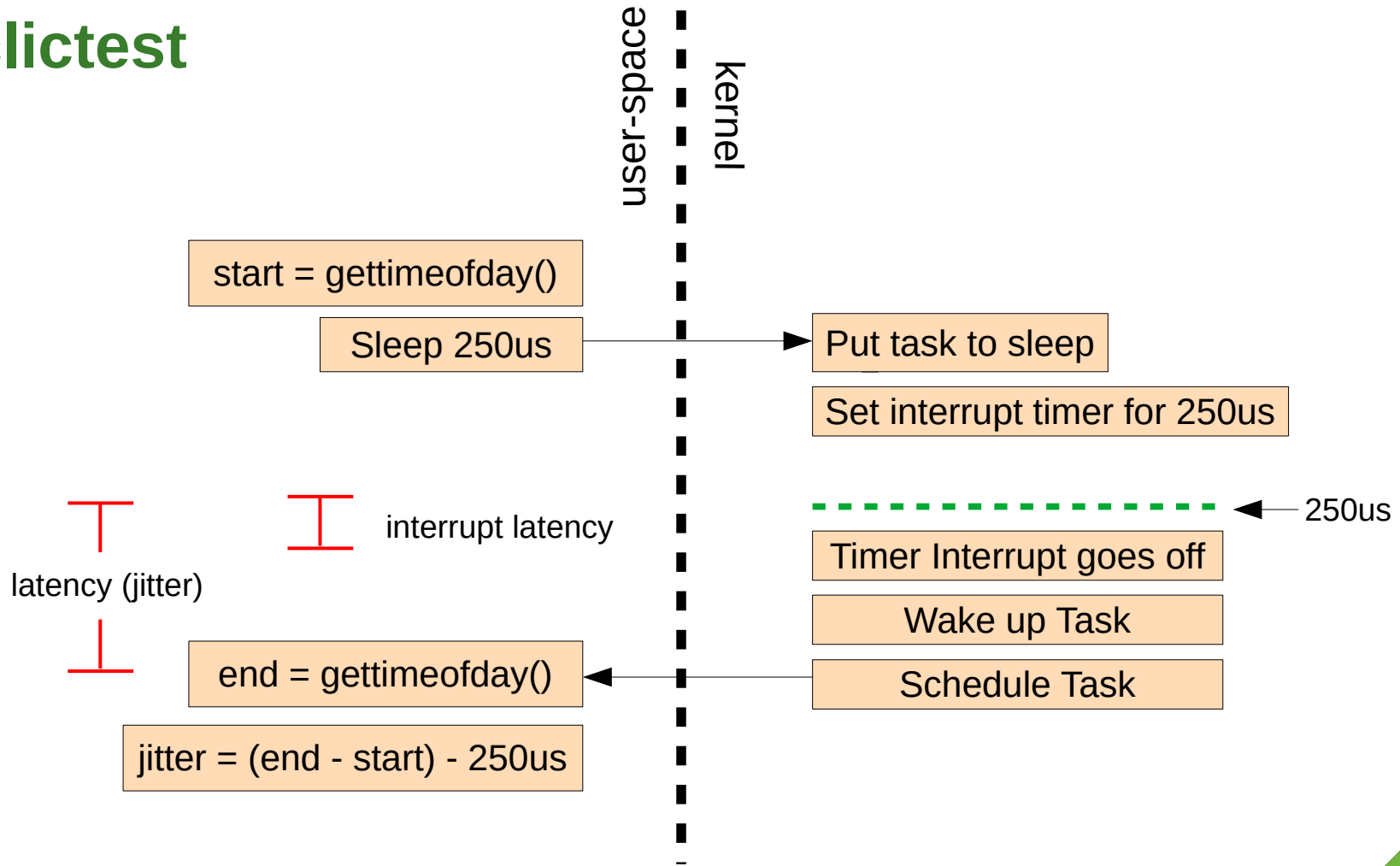
# Cyclictest



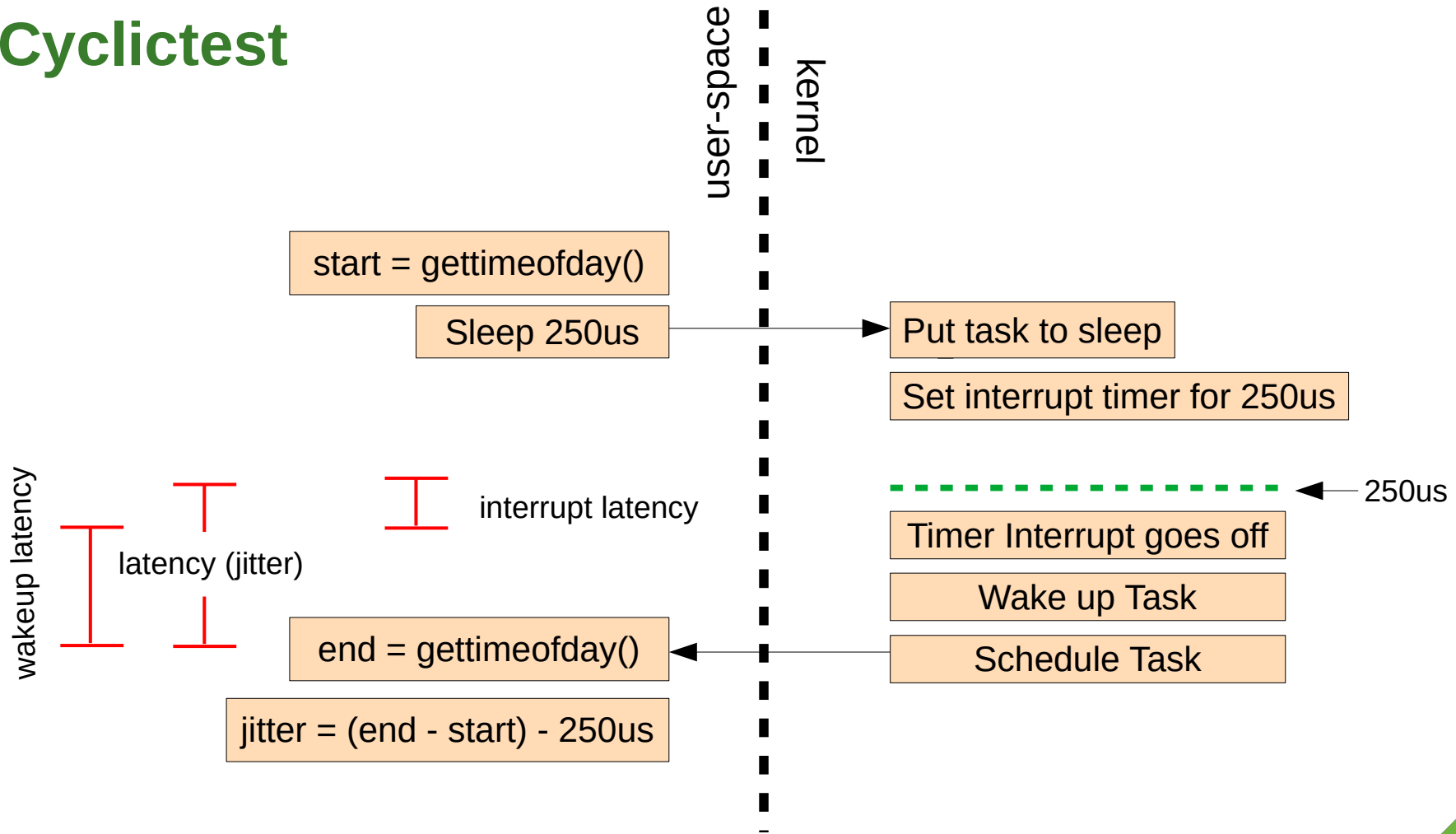
# Cyclictest



# Cyclictest



# Cyclictest



# See the hardware latency?

```
ksoftirqd/33-216 [033] 55597.719935: timer_cancel: timer=0xffff88403f0ce520
ksoftirqd/33-216 [033] 55597.719935: timer_expire_entry: timer=0xffff88403f0ce520 function=delayed_work
ksoftirqd/33-216 [033] 55597.719935: funcgraph_entry: 0.069 us | _raw_spin_lock_irqsave();
ksoftirqd/33-216 [033] 55597.719936: funcgraph_entry: 0.047 us | _raw_spin_lock();
ksoftirqd/33-216 [033] 55597.719936: sched_stat_sleep: comm=kworker/33:1 pid=1222 delay=132870067
ksoftirqd/33-216 [033] 55597.719937: sched_wakeup: kworker/33:1:1222 [120] success=1 CPU:033
ksoftirqd/33-216 [033] 55597.719937: timer_expire_exit: timer=0xffff88403f0ce520
ksoftirqd/33-216 [033] 55597.719942: timer_cancel: timer=0xffff88403f0ce620
ksoftirqd/33-216 [033] 55597.719942: timer_expire_entry: timer=0xffff88403f0ce620 function=delayed_work
ksoftirqd/33-216 [033] 55597.719943: funcgraph_entry: 0.045 us | _raw_spin_lock_irqsave();
ksoftirqd/33-216 [033] 55597.719943: timer_expire_exit: timer=0xffff88403f0ce620
  cyclictst-6110 [007] 55597.719955: funcgraph_entry: 0.194 us | _raw_spin_lock();
  cyclictst-6110 [007] 55597.719956: funcgraph_entry: 0.175 us | _raw_spin_lock_irqsave();
  cyclictst-6110 [007] 55597.719957: funcgraph_entry: 0.175 us | _raw_spin_lock_irqsave();
  cyclictst-6113 [010] 55597.719957: funcgraph_entry: 2.436 us | _raw_spin_lock();
  cyclictst-6110 [007] 55597.719957: funcgraph_entry: 0.203 us | _raw_spin_lock();
  cyclictst-6110 [007] 55597.719958: sched_wakeup: cyclictst:6113 [4] success=1 CPU:010
  cyclictst-6110 [007] 55597.719959: funcgraph_entry: 0.048 us | _raw_spin_lock_irqsave();
  cyclictst-6113 [010] 55597.719960: funcgraph_entry: 0.170 us | _raw_spin_lock_irq();
  cyclictst-6113 [010] 55597.719961: funcgraph_entry: 0.045 us | _raw_spin_lock_irqsave();
  cyclictst-6113 [010] 55597.719962: funcgraph_entry: 0.043 us | _raw_spin_lock_irq();
  cyclictst-6110 [007] 55597.719963: print: ffffffff810e5776 hit latency threshold (247 > 200)
```

# How about now?

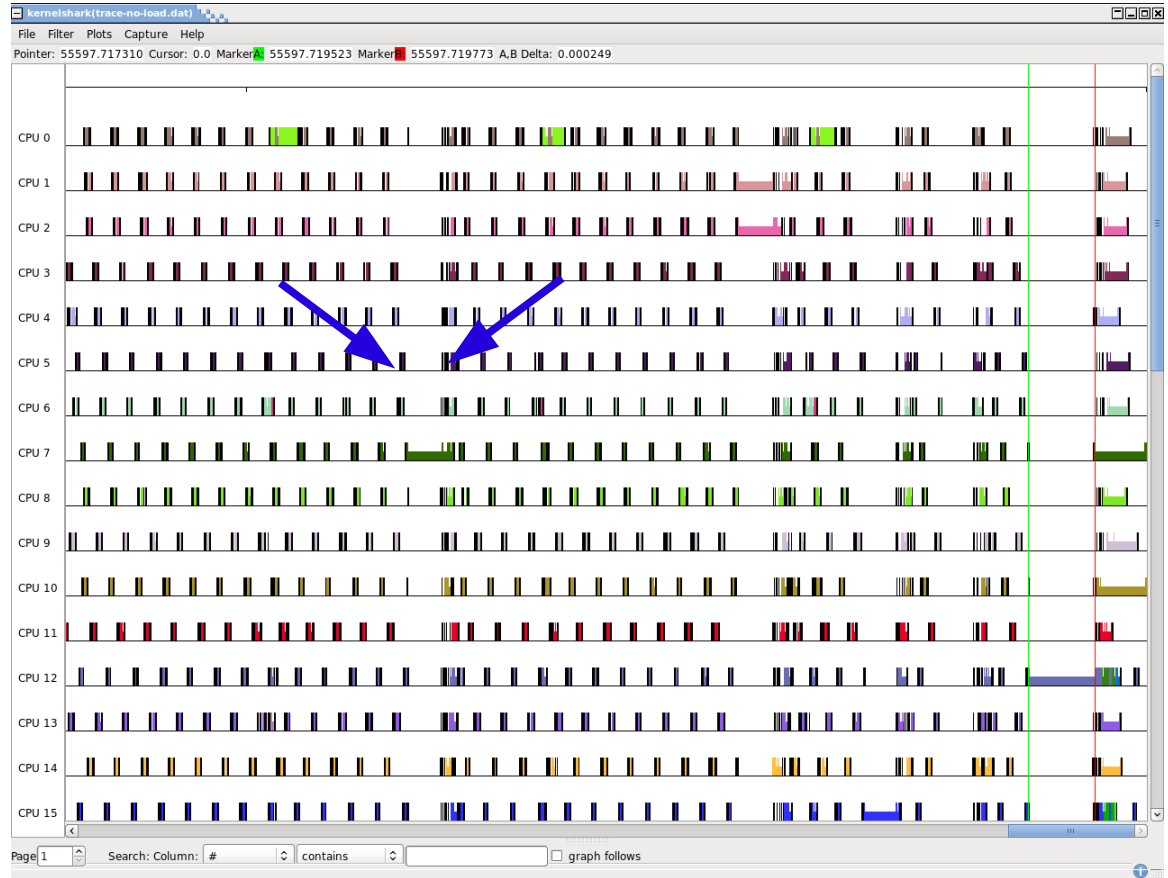


# Better?

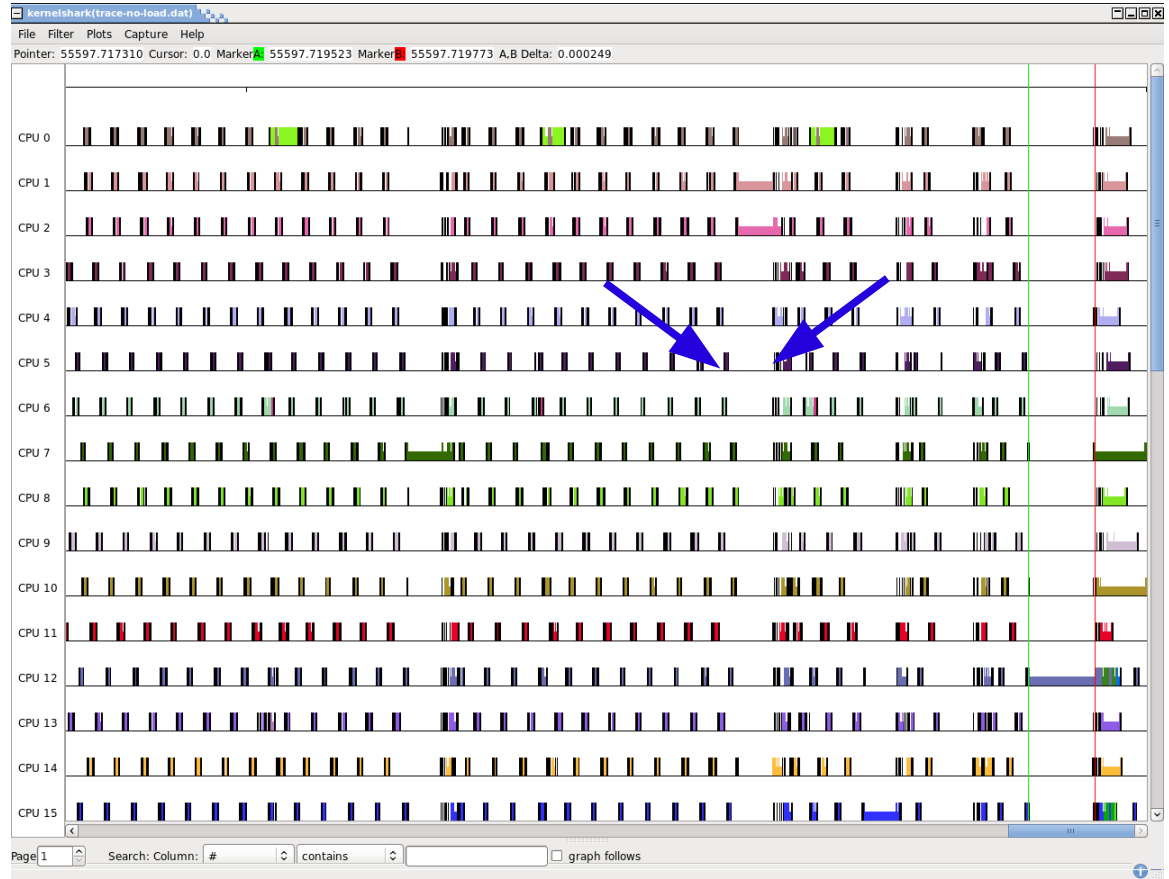




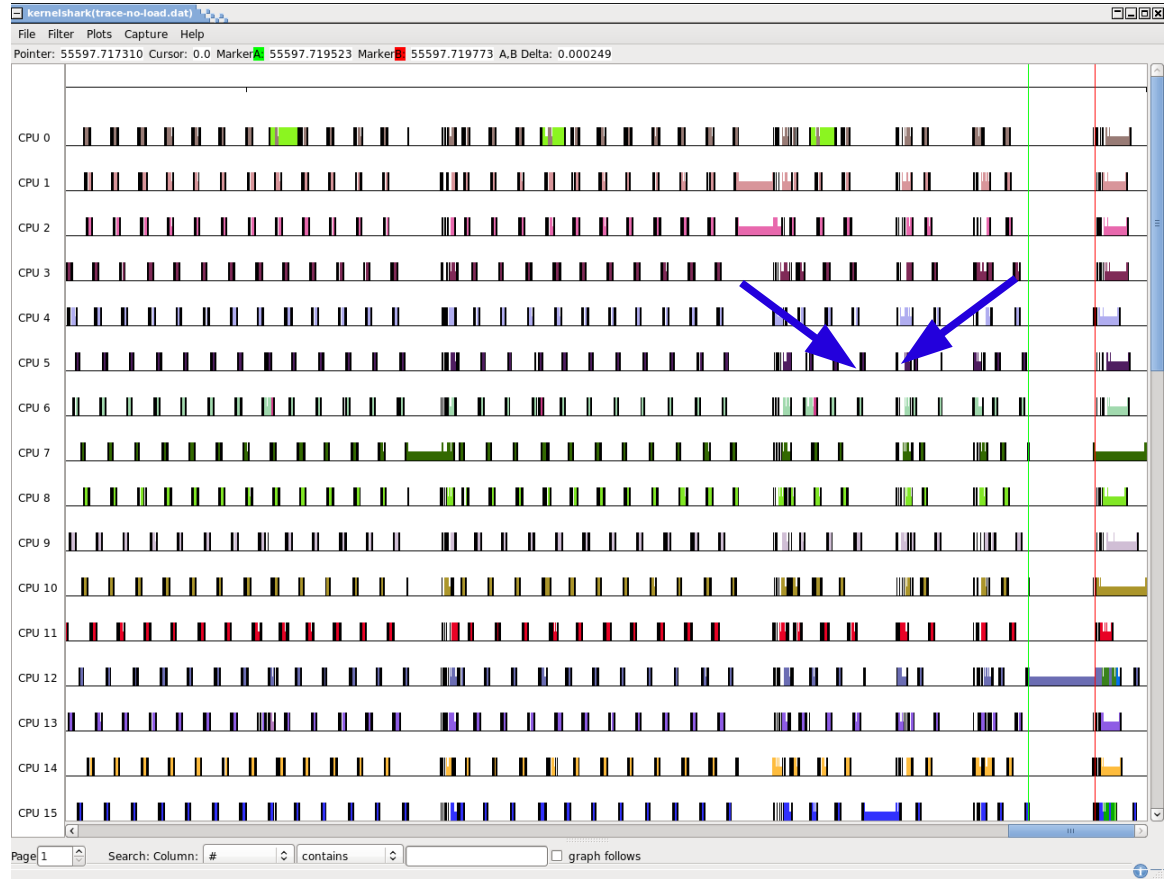
# I do, here



and here



and here



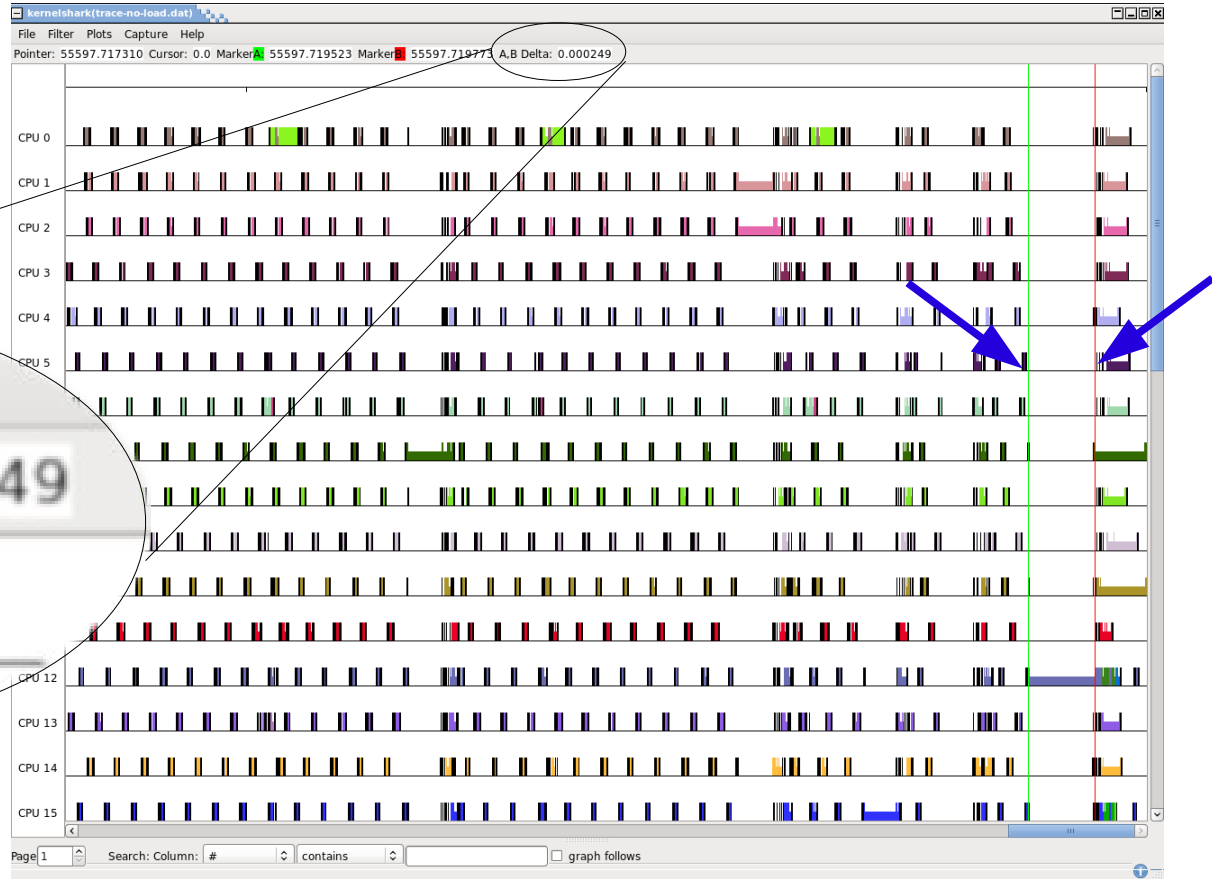
and here



and here



# 249us Latency!



# SCHED\_DEADLINE



# KernelShark 0.2 (GTK version) vs 1.0 (Qt version)

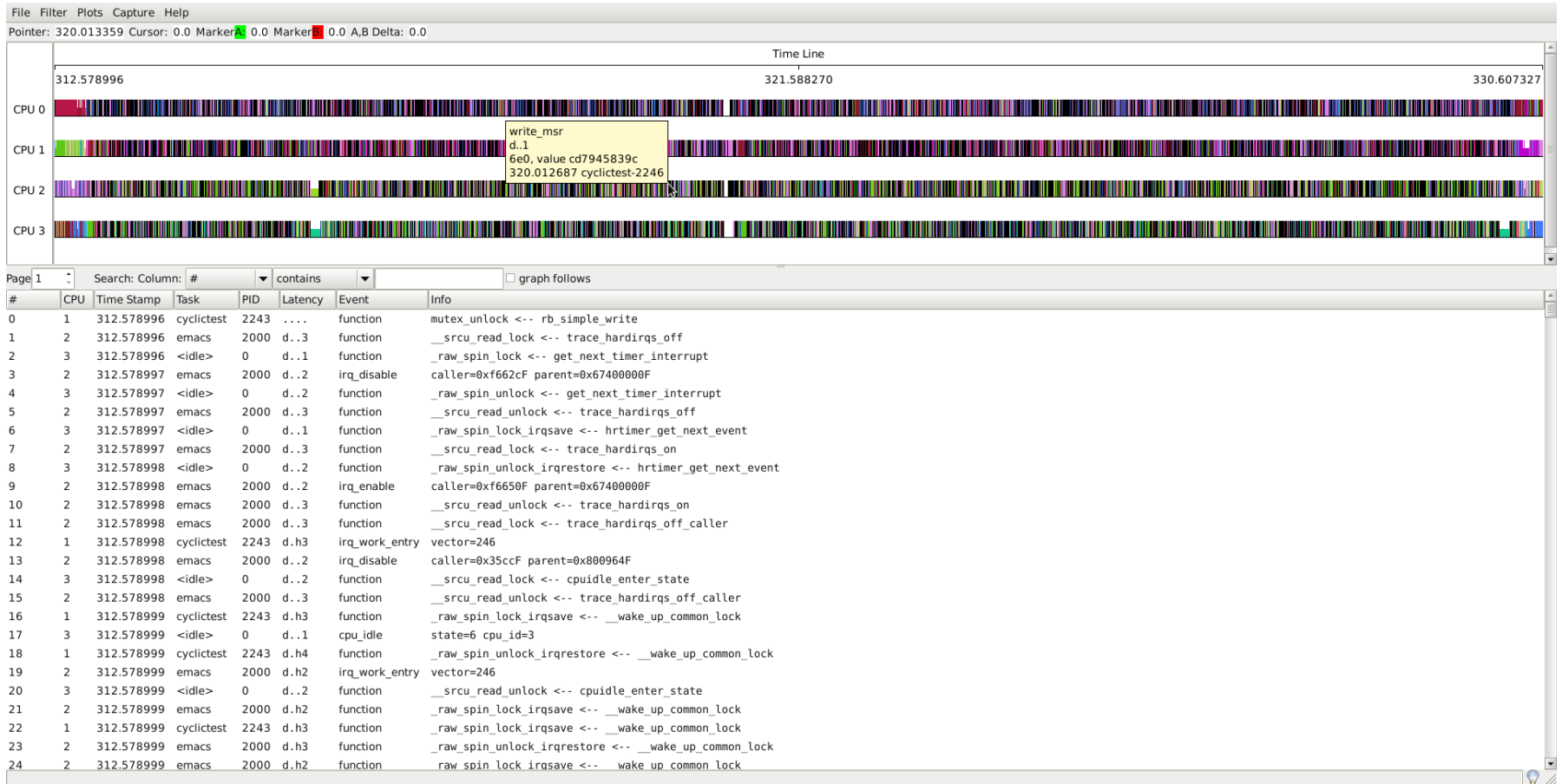
- 1.0 is MUCH faster!
  - New algorithm to fill in pixels
    - With 1,000,000,000( $n$ ) events in 1000( $m$ ) pixels width window
      - 0.2 Would do linear search to find next event to draw the next line
        - 1,000,000 queries per vertical line drawn -  $O((n/m))$
      - 1.0 does a binary search for the event in the next pixel
        - 30 queries per vertical line draw -  $O(\log_2 n)$
        - Switches to linear search when  $\log_2 n \geq n/m$ 
          - $\sim 10,000$  events per 1000 pixels



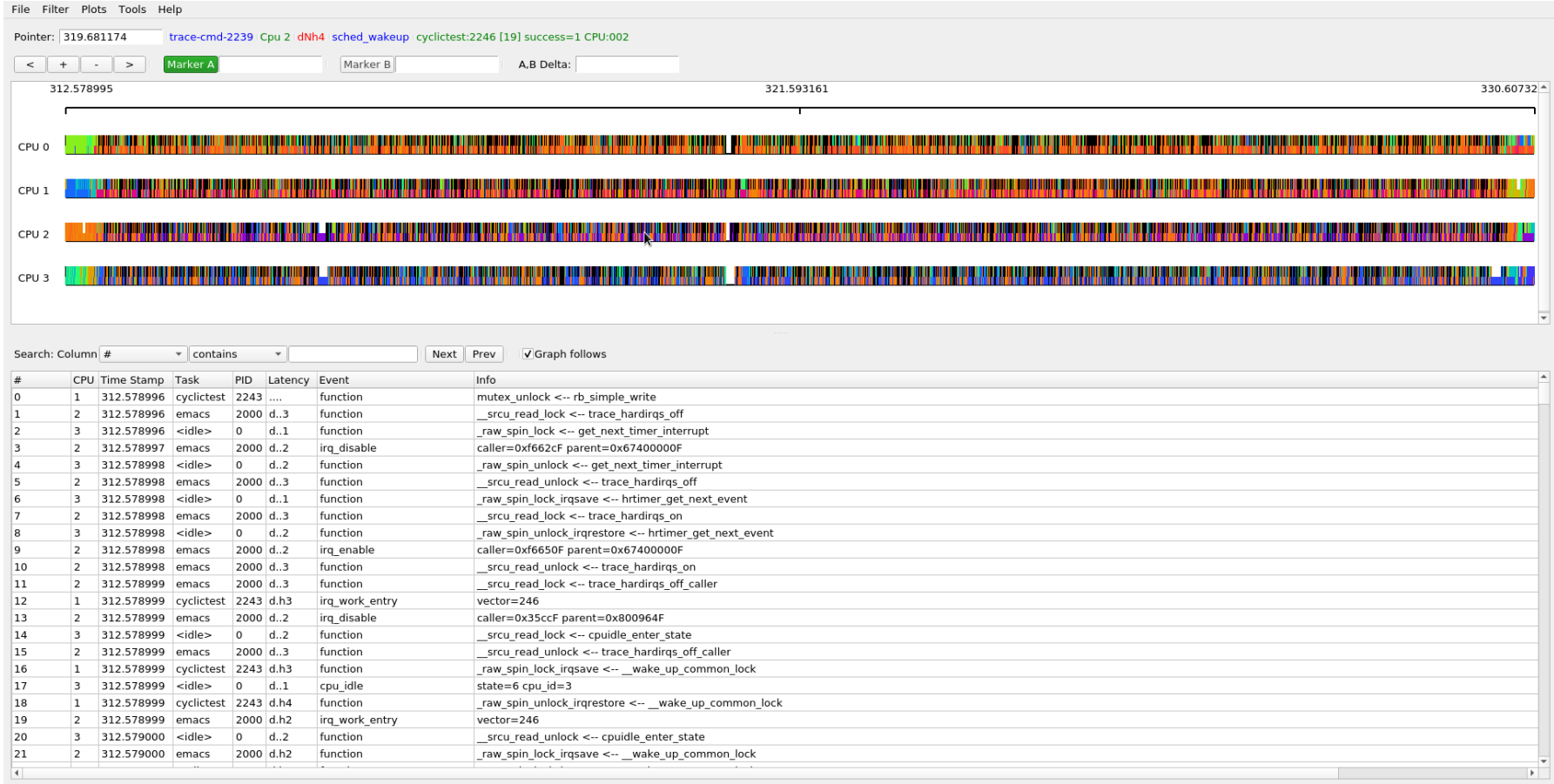
# KernelShark 1.0 (Qt version)

- Written by Yordan Karadzov
  - Although I'm still the maintainer (reviewer of the code)
- Rewritten from scratch
  - I let him play
    - The requirement was to have all features of KernelShark 0.2
    - Then I revisited to help him
    - Break it up into a patch series
      - Why it's late :-p
  - New Look!
  - New design!
    - Becoming a library (more on that later)
- Much faster than KernelShark 0.2 (must keep stressing this!)

# KernelShark 0.2



# KernelShark 1.0

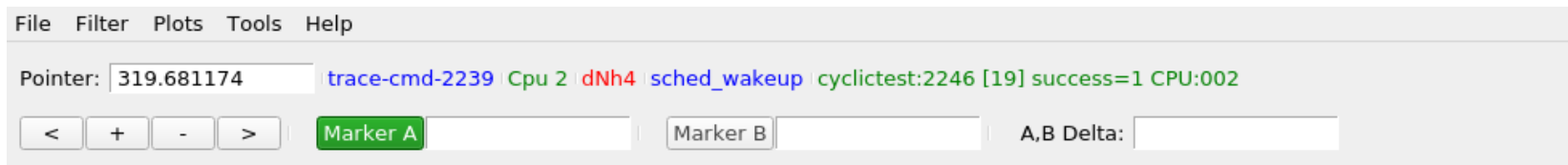


# Headers

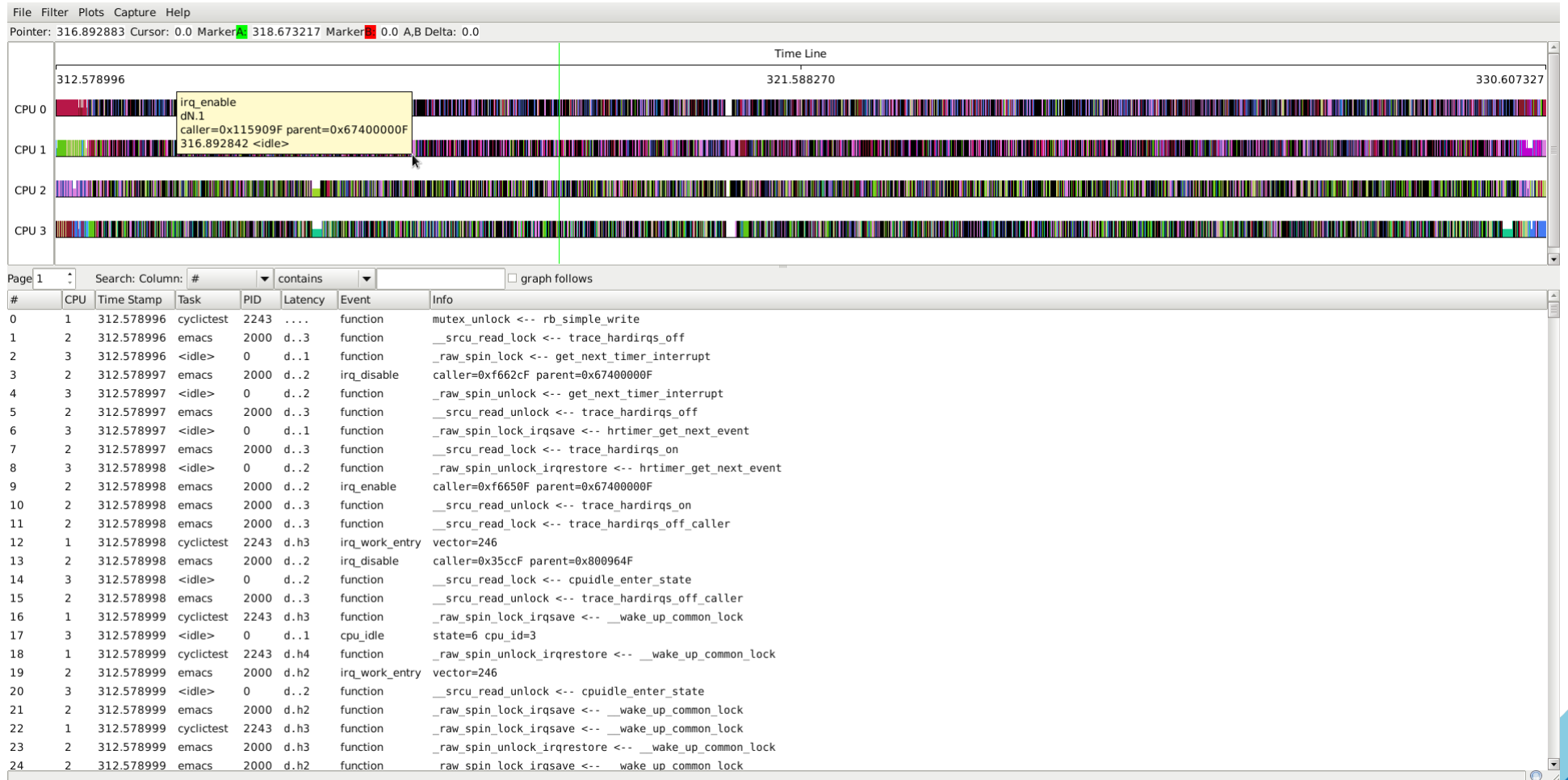
## KernelShark 0.2



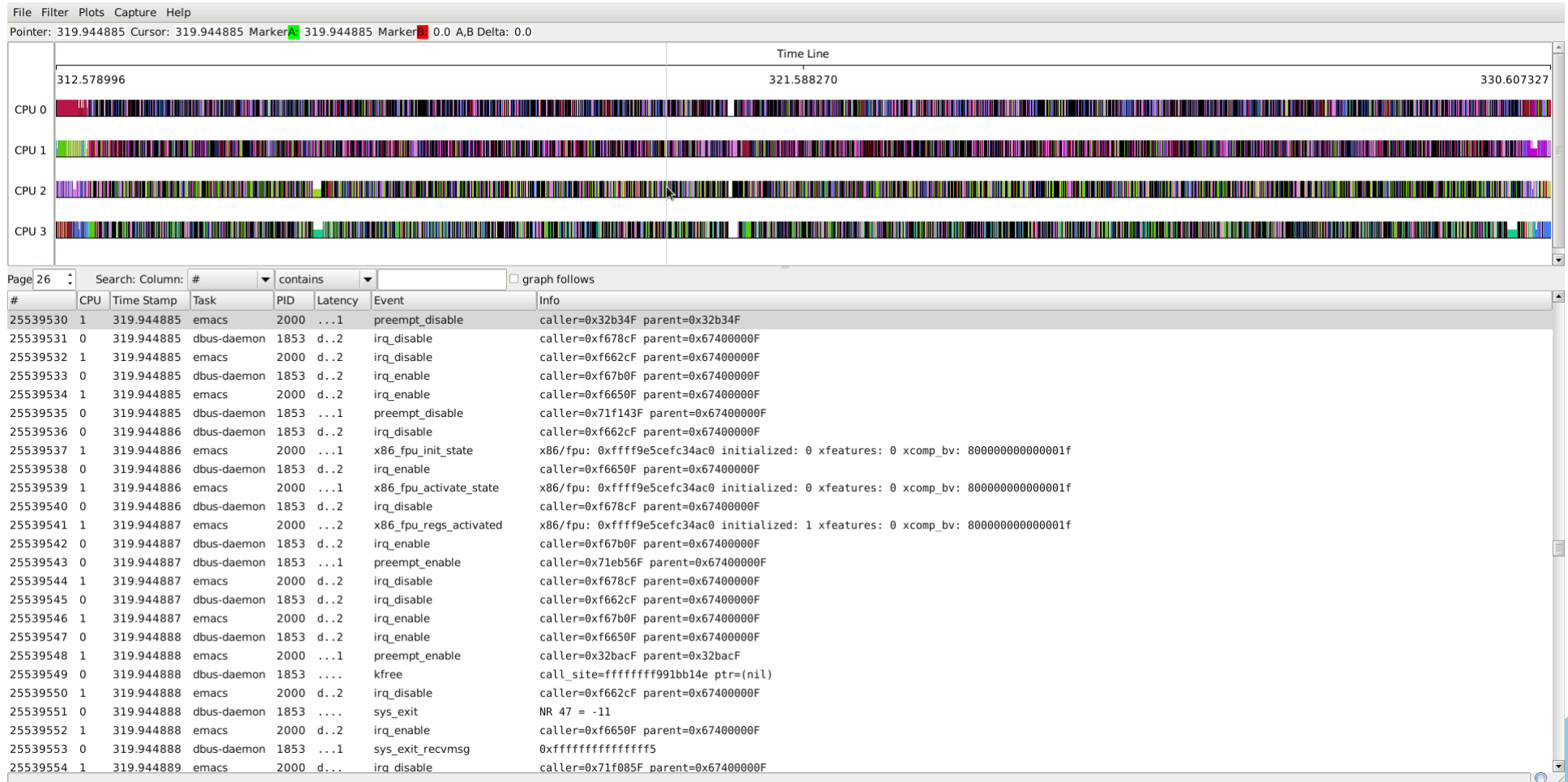
## KernelShark 1.0



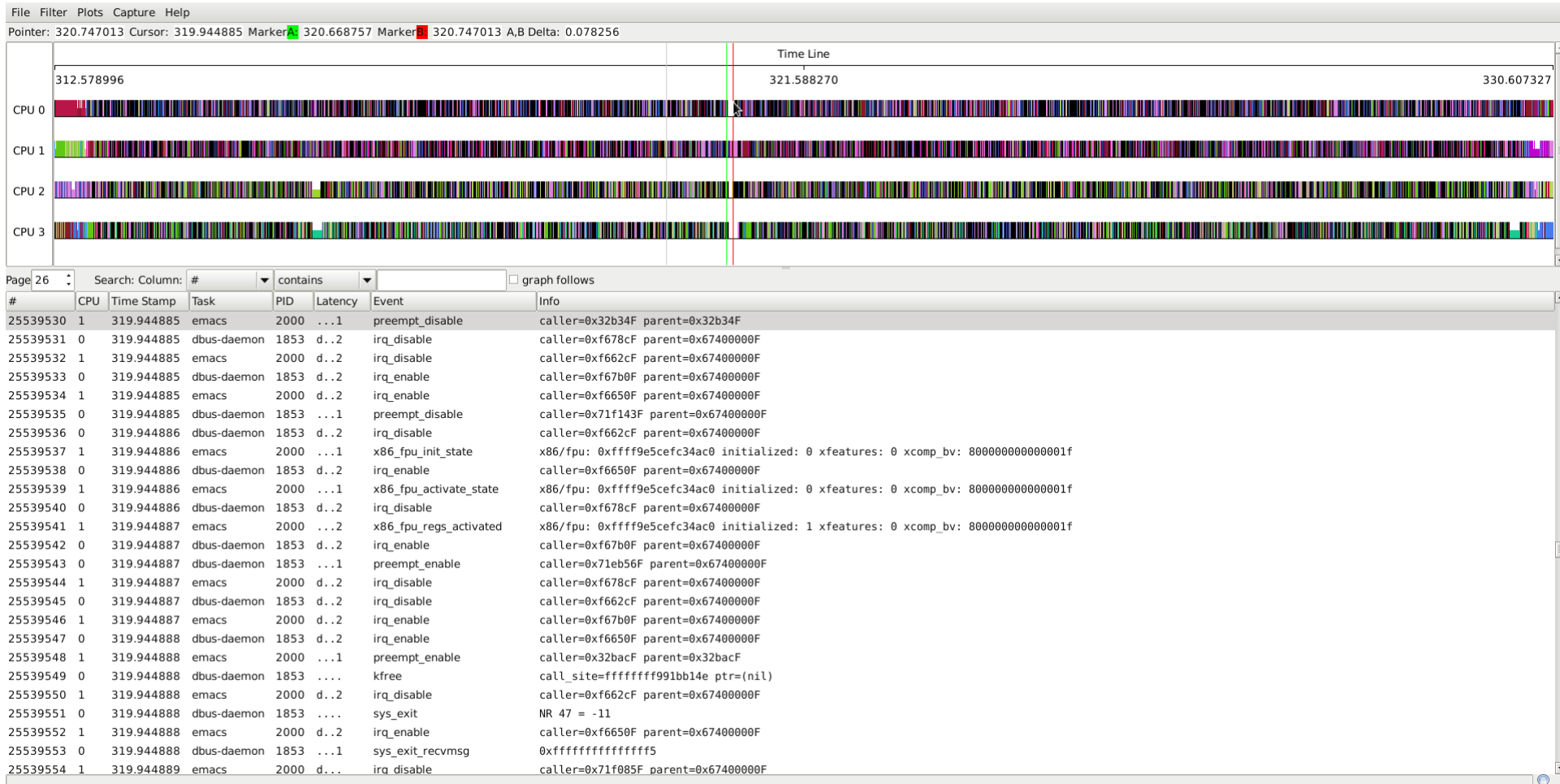
# KernelShark 0.2 Markers ("Marker A")



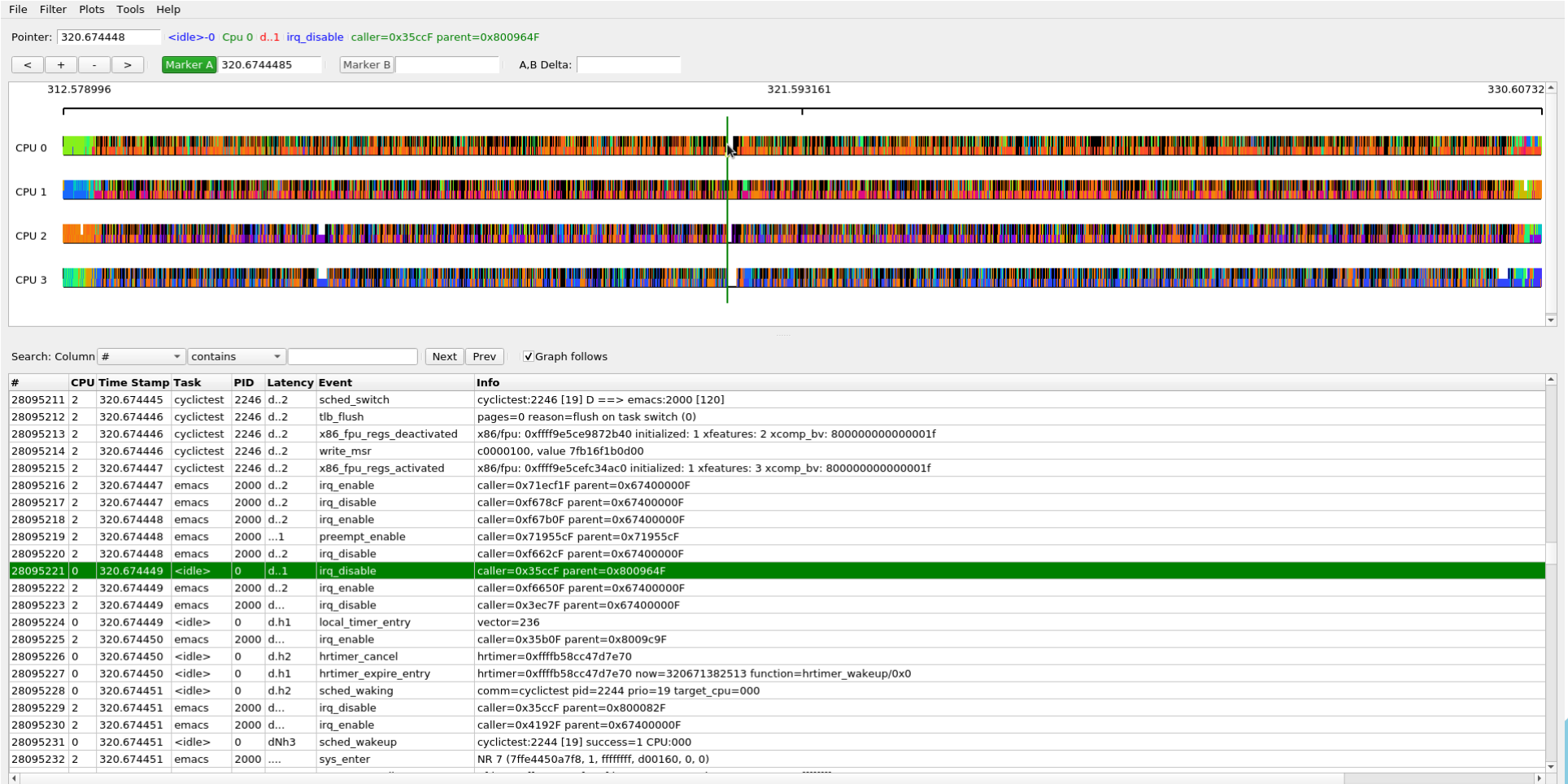
# KernelShark 0.2 Markers (“Cursor”)



# KernelShark 0.2 Markers ("Marker B")

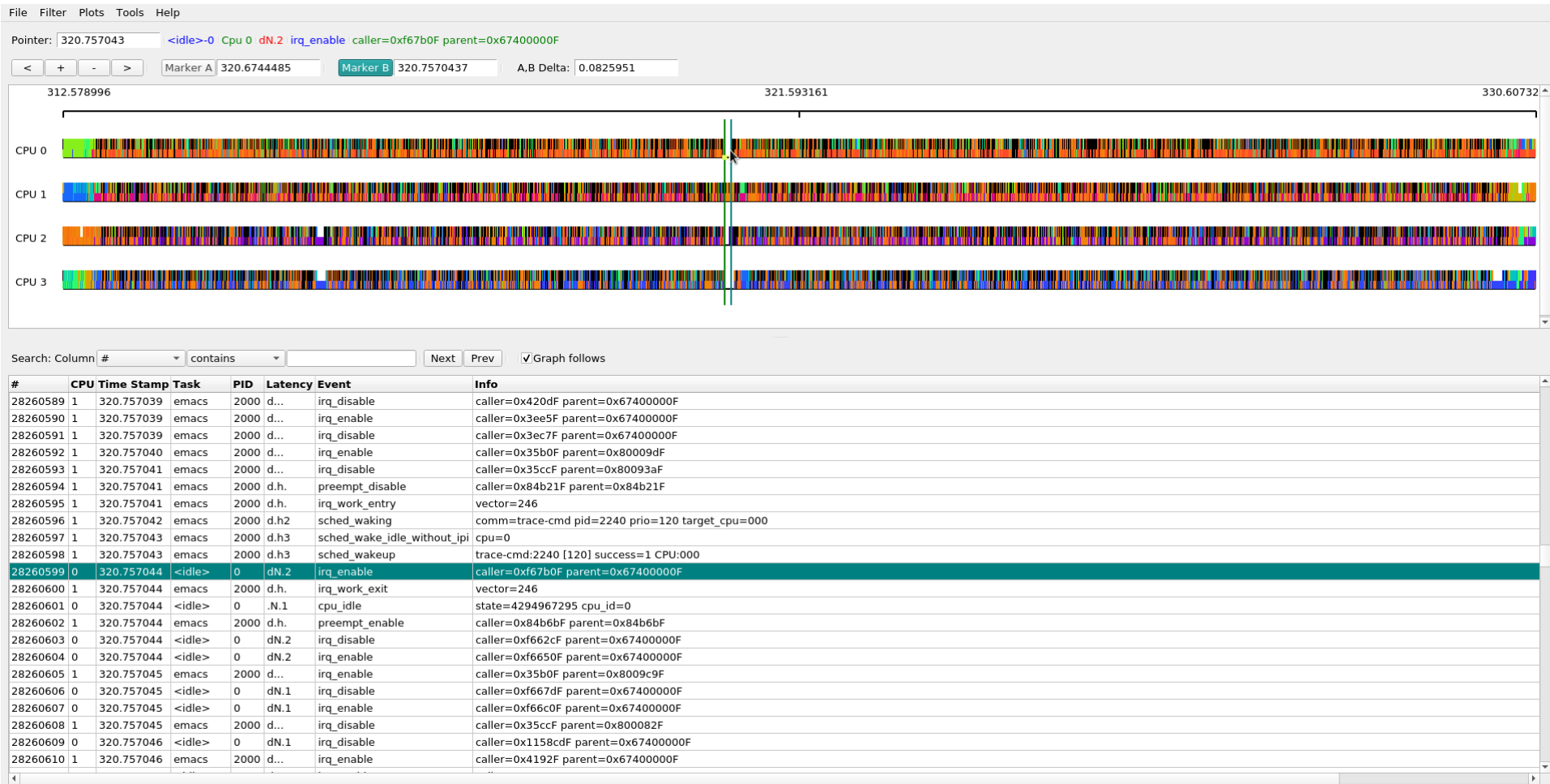


# KernelShark 1.0 Markers (“Marker A”)

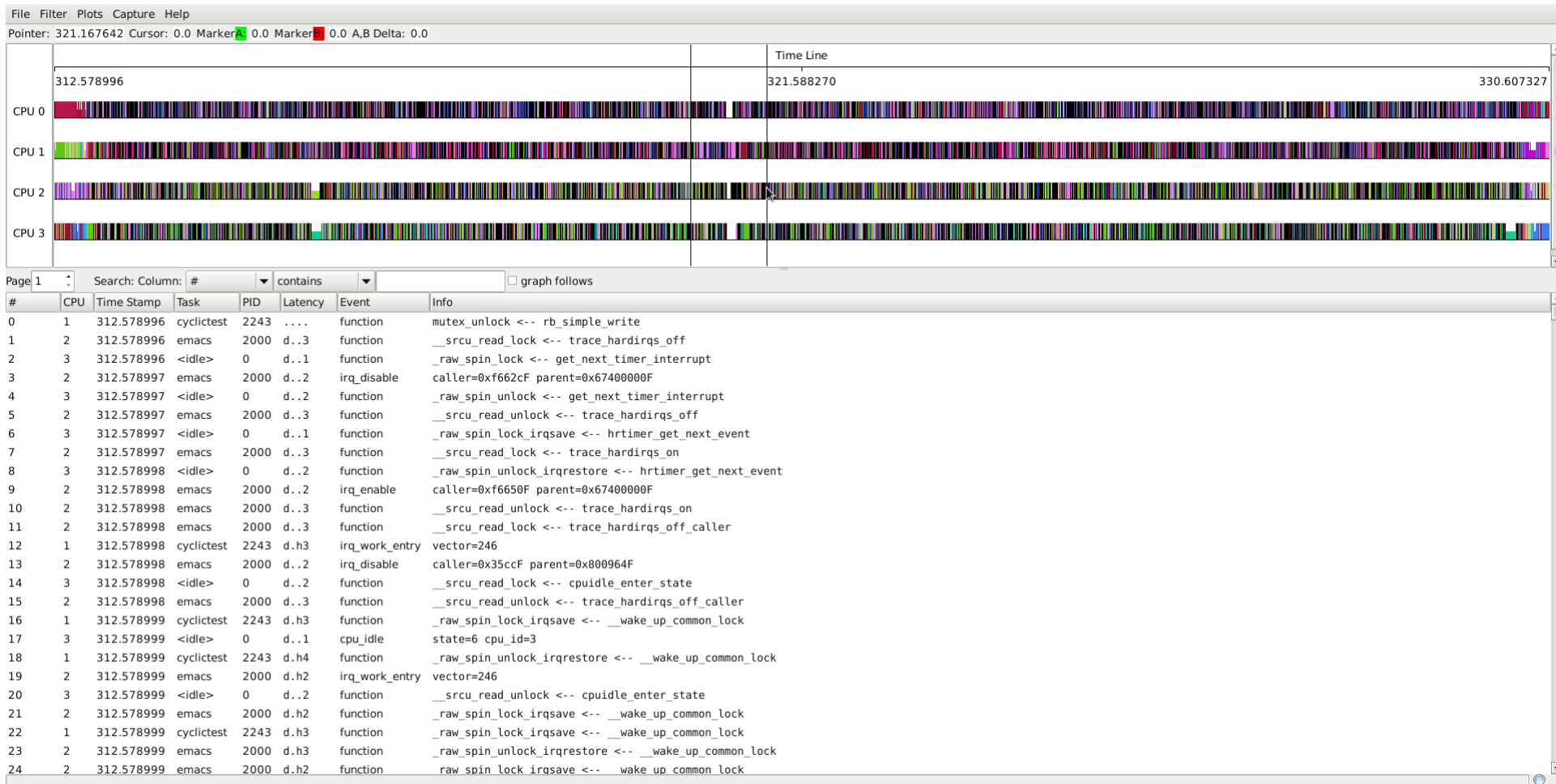




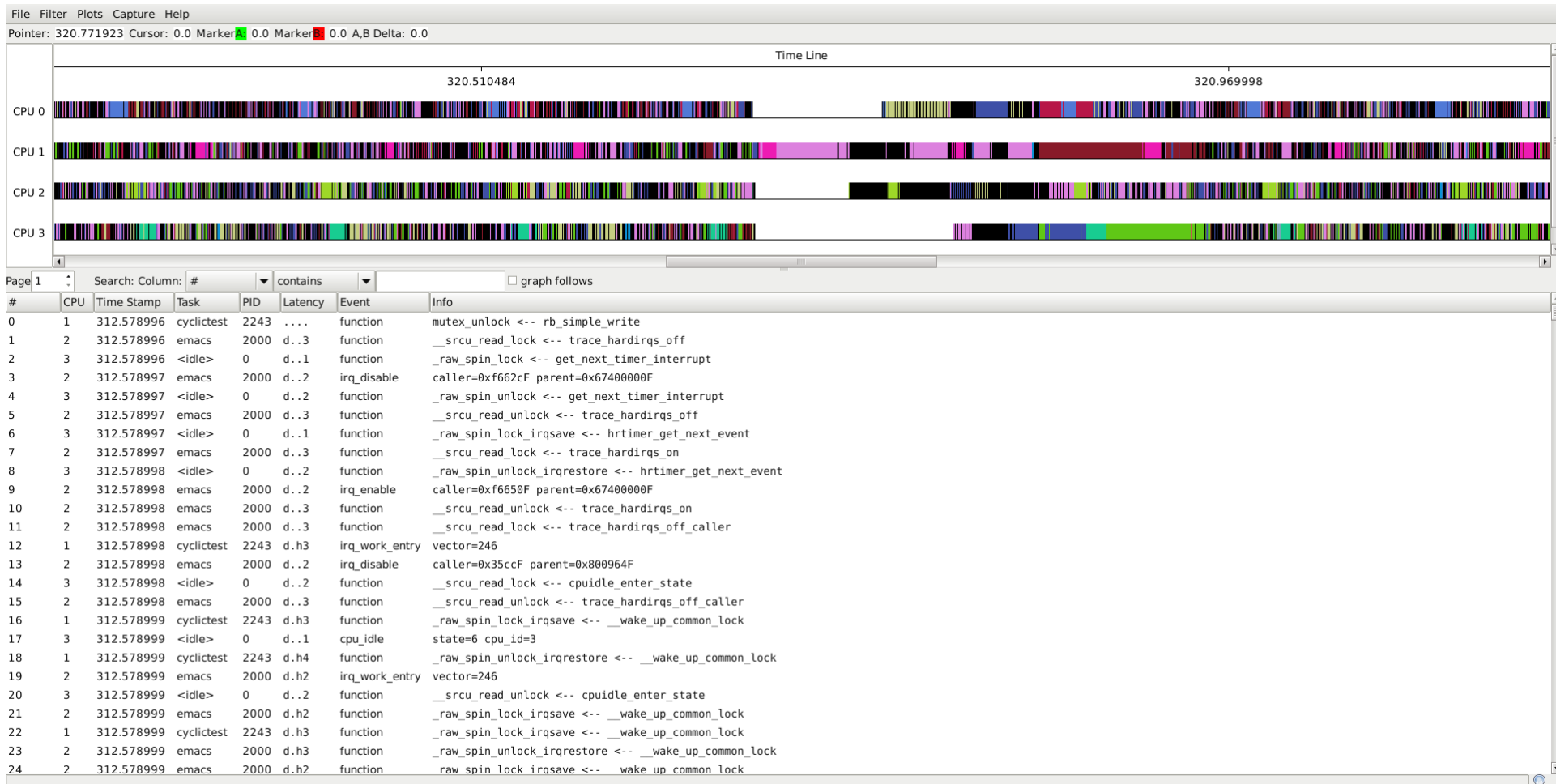
# KernelShark 1.0 Markers (“Marker B”)



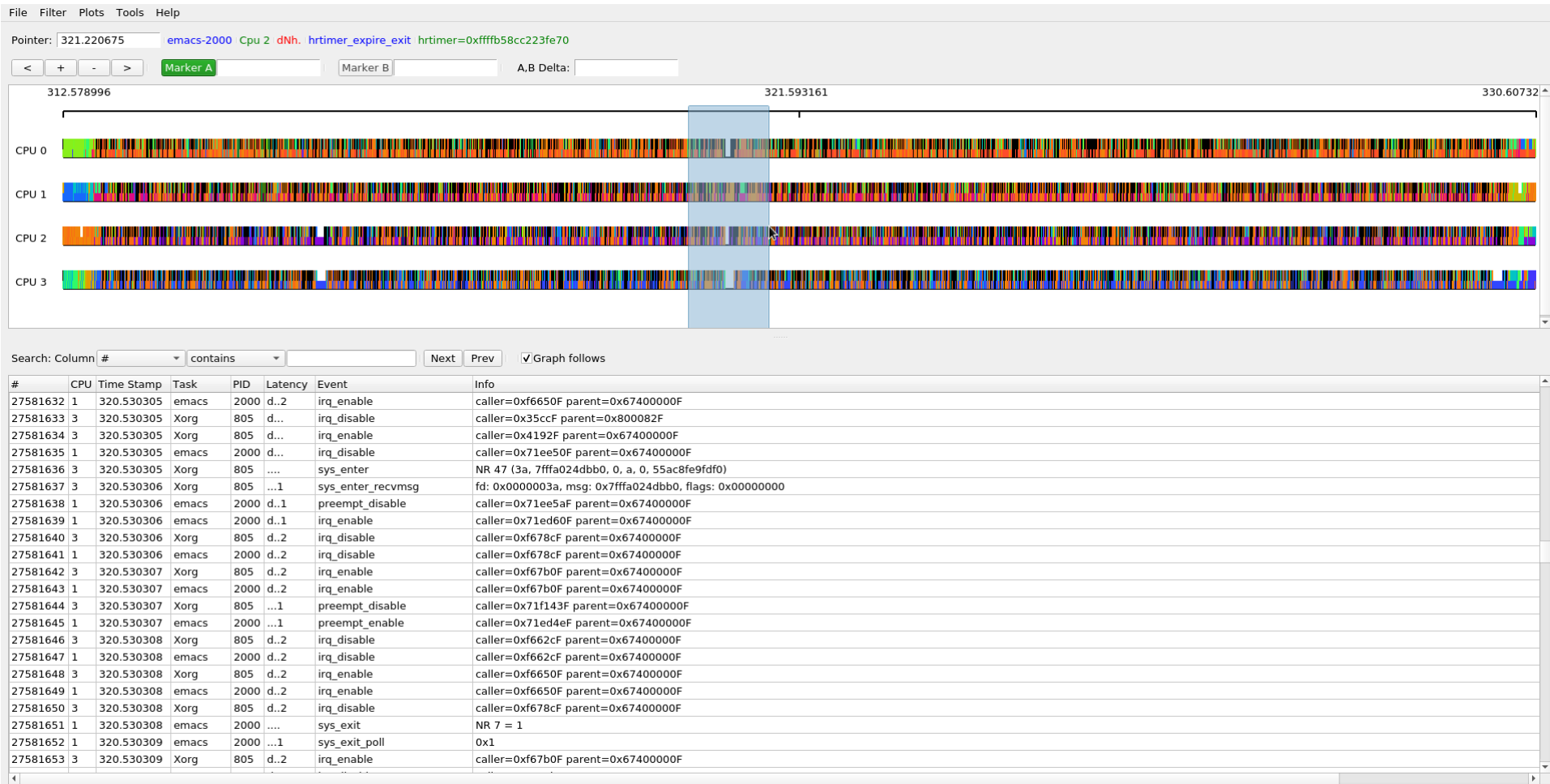
# KernelShark 0.2 Zooming



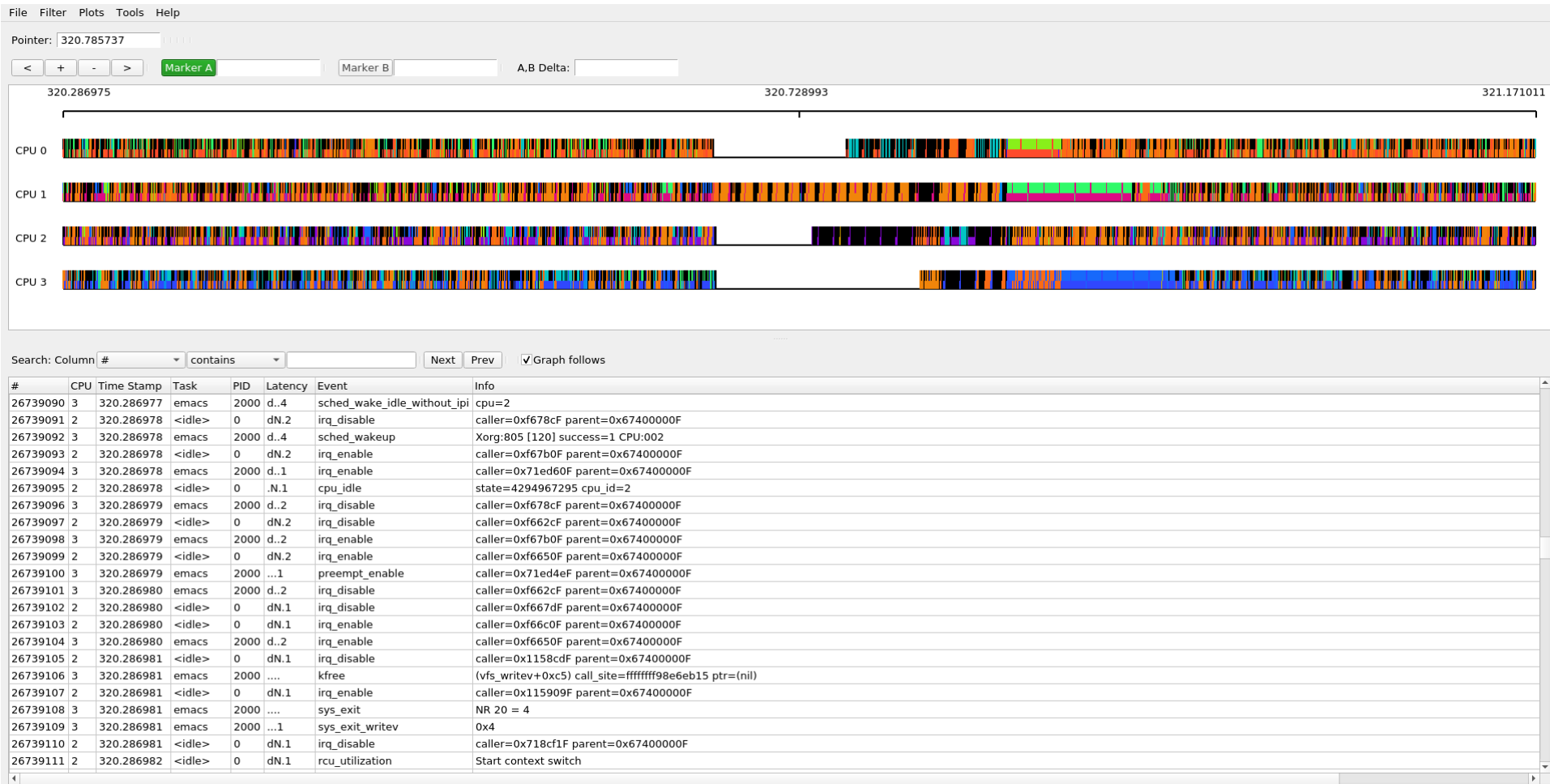
# KernelShark 0.2 Zooming



# KernelShark 1.0 Zooming



# KernelShark 1.0 Zooming



# Selecting Tasks

Plot	PID	Task
<input type="checkbox"/>	1938	gmain
<input type="checkbox"/>	1943	panel-22-batter
<input type="checkbox"/>	1945	panel-16-pulsea
<input type="checkbox"/>	1949	gmain
<input type="checkbox"/>	1952	panel-20-cpugra
<input type="checkbox"/>	1953	panel-17-netloa
<input type="checkbox"/>	1954	panel-18-system
<input type="checkbox"/>	1970	gdbus
<input type="checkbox"/>	1977	xfce4-terminal
<input type="checkbox"/>	2000	emacs
<input type="checkbox"/>	2037	dconf
<input type="checkbox"/>	2051	nm-applet
<input type="checkbox"/>	2118	gdbus
<input type="checkbox"/>	2239	trace-cmd
<input type="checkbox"/>	2240	trace-cmd
<input type="checkbox"/>	2241	trace-cmd
<input type="checkbox"/>	2242	trace-cmd
<input checked="" type="checkbox"/>	2243	cyclictest
<input checked="" type="checkbox"/>	2244	cyclictest
<input checked="" type="checkbox"/>	2245	cyclictest
<input checked="" type="checkbox"/>	2246	cyclictest
<input checked="" type="checkbox"/>	2247	cyclictest
<input type="checkbox"/>	2248	docker-containe

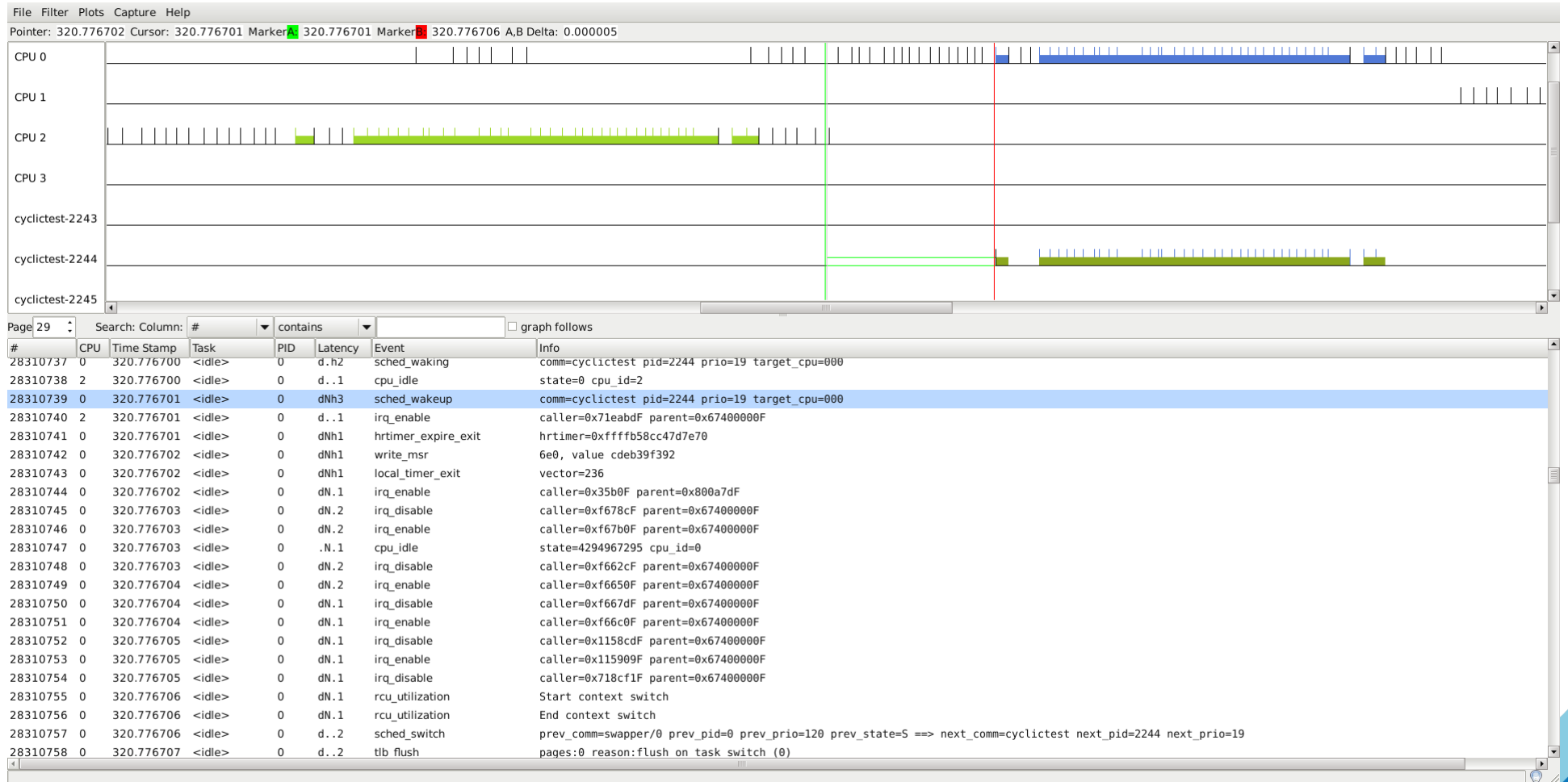
Apply

Tasks:  all

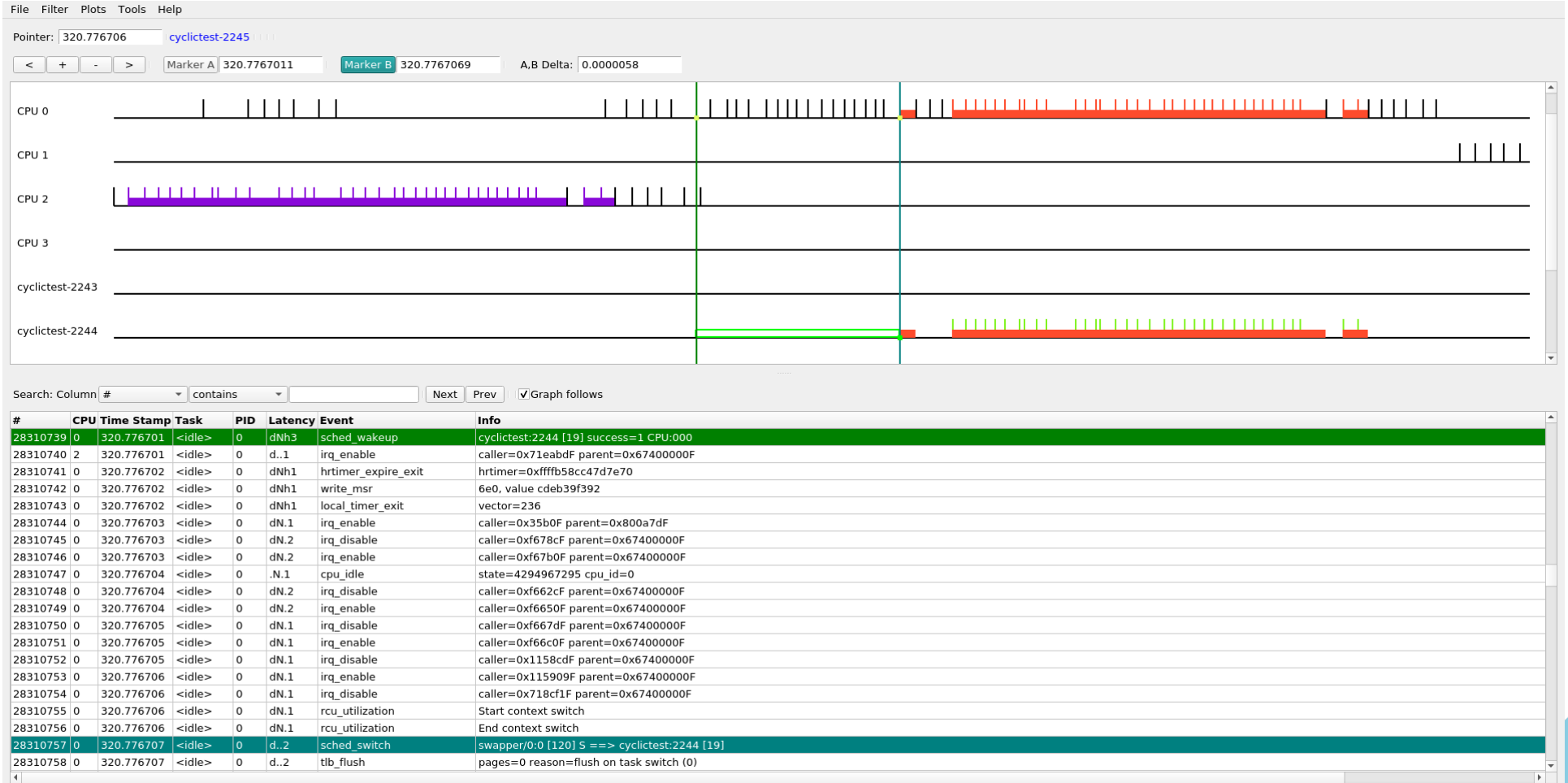
Hide	Pid	Task
<input type="checkbox"/>	1970	gdbus
<input type="checkbox"/>	1977	xfce4-terminal
<input type="checkbox"/>	2000	emacs
<input type="checkbox"/>	2037	dconf
<input type="checkbox"/>	2051	nm-applet
<input type="checkbox"/>	2118	gdbus
<input type="checkbox"/>	2239	trace-cmd
<input type="checkbox"/>	2240	trace-cmd
<input type="checkbox"/>	2241	trace-cmd
<input type="checkbox"/>	2242	trace-cmd
<input checked="" type="checkbox"/>	2243	cyclictest
<input checked="" type="checkbox"/>	2244	cyclictest
<input checked="" type="checkbox"/>	2245	cyclictest
<input checked="" type="checkbox"/>	2246	cyclictest
<input checked="" type="checkbox"/>	2247	cyclictest
<input type="checkbox"/>	2248	docker-containe

Apply

# KernelShark 0.2 Wake-up Latency

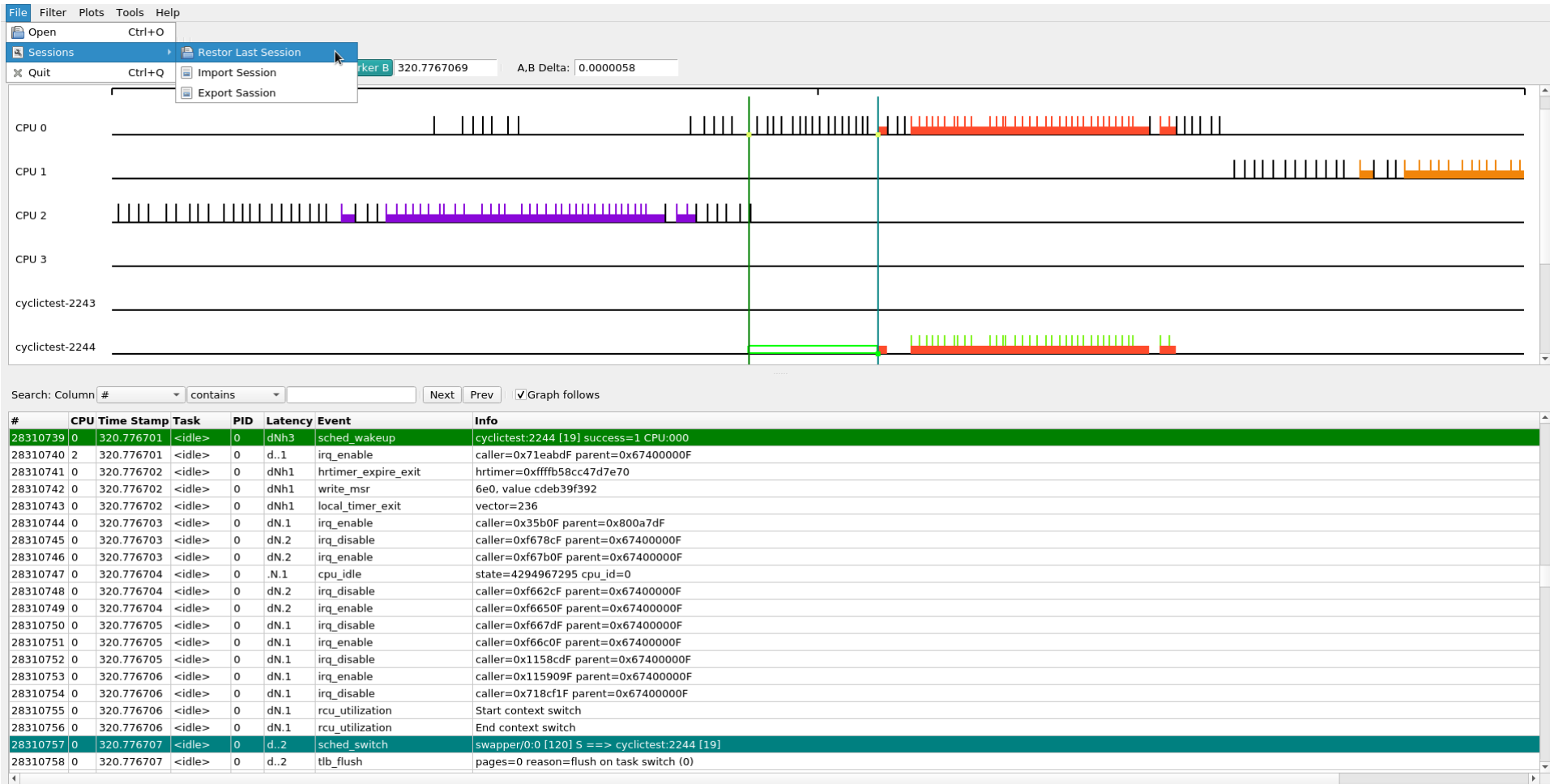


# KernelShark 1.0 Wake-up Latency

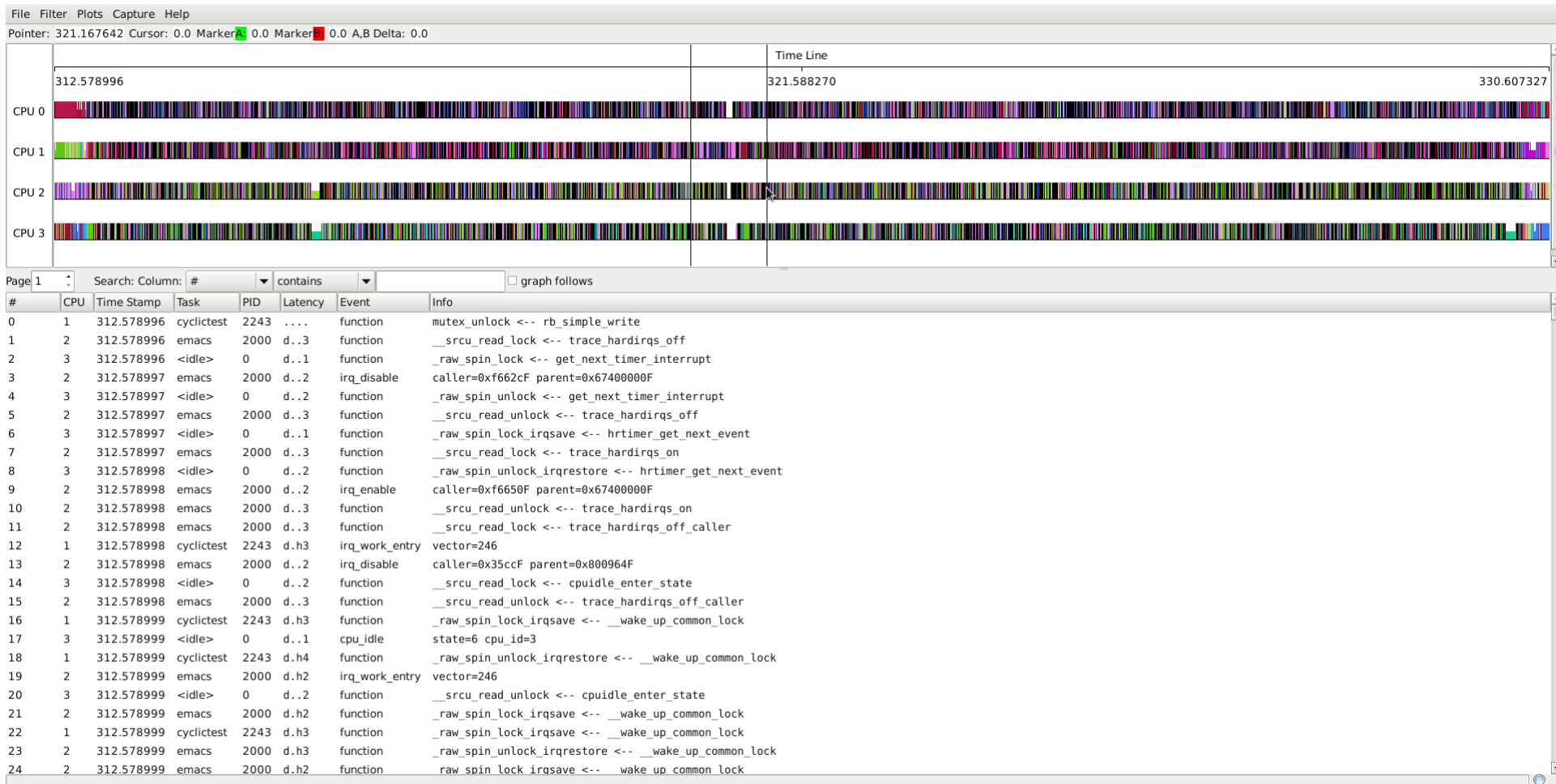




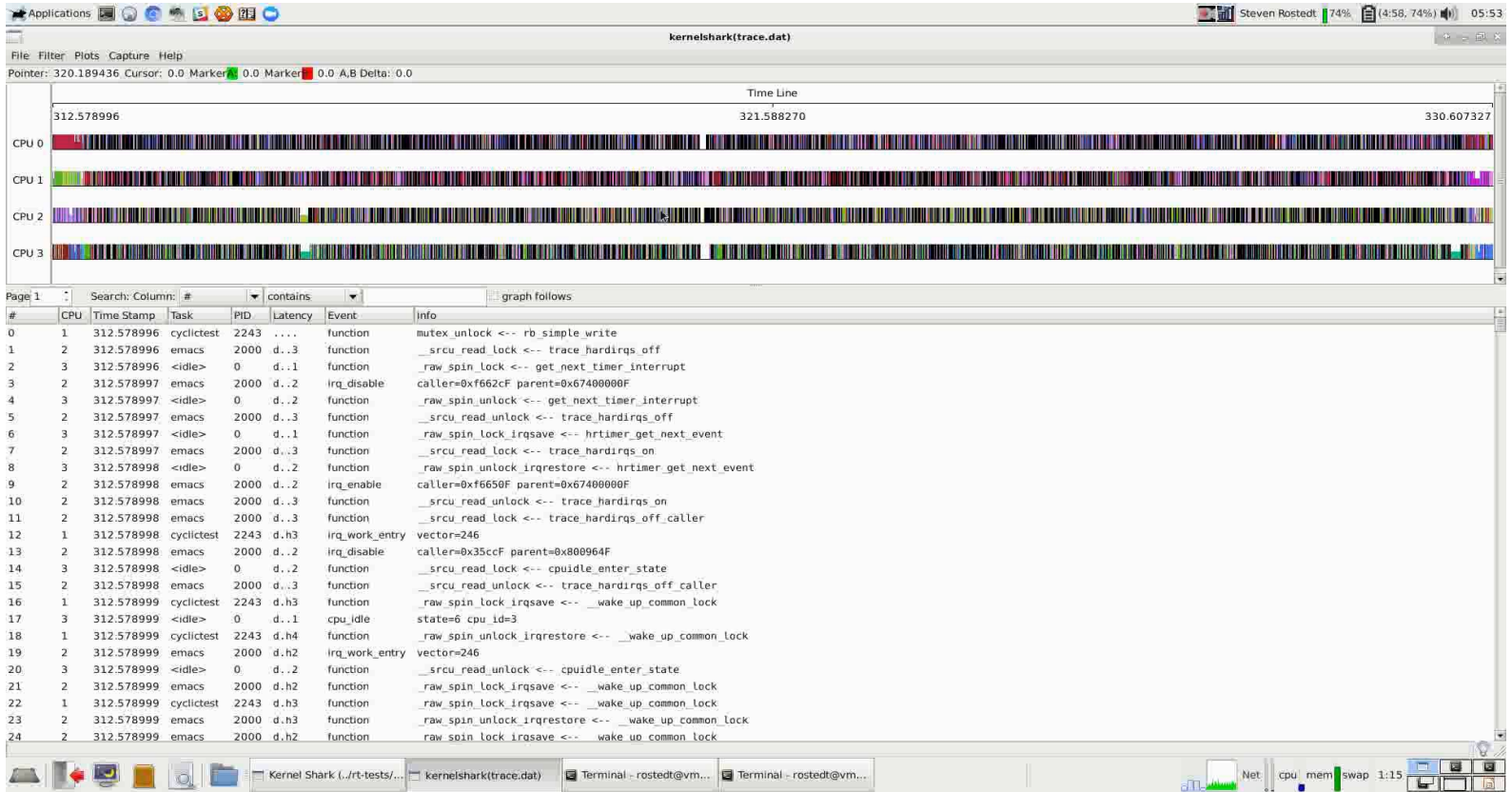
# KernelShark 1.0 Sessions



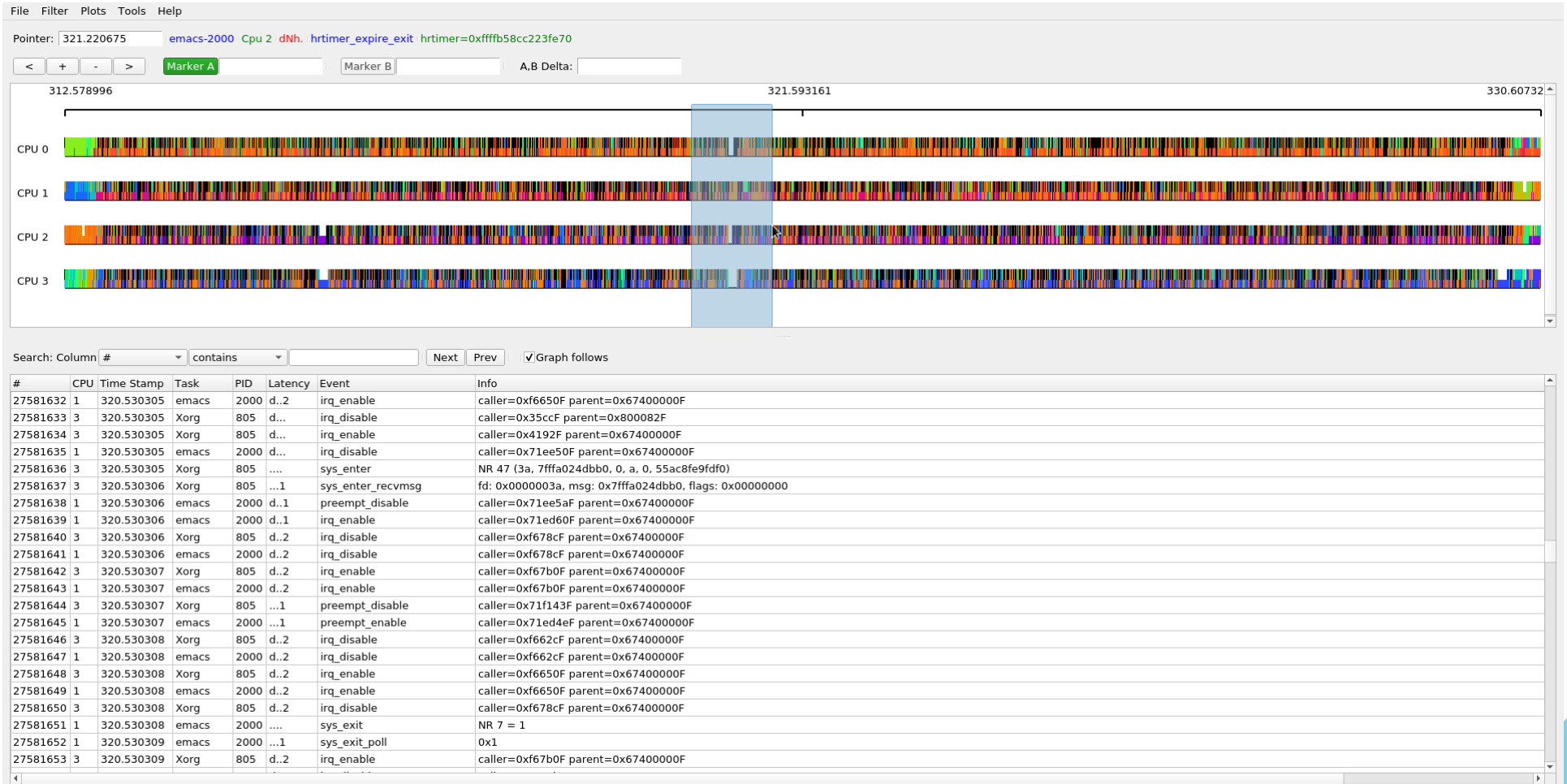
# KernelShark 0.2 Zooming



# KernelShark 0.2 Zooming



# KernelShark 1.0 Zooming



# KernelShark 1.0 Zooming

Kernel Shark (../rt-tests/trace.dat)

Pointer: 320.197580 <idle>-0 Cpu:2 d..2 irq\_enable caller=0x71ed4eF parent=0x67400000F

312.578996 321.593161 330.60732

CPU 0  
CPU 1  
CPU 2  
CPU 3

Search: Column # contains Next Prev  Graph follows

#	CPU	Time Stamp	Task	PID	Latency	Event	Info
28146170	1	320.700362	emacs	2000	d..2	irq_disable	caller=0xf678cF parent=0x67400000F
28146171	1	320.700362	emacs	2000	d..2	irq_enable	caller=0xf67b0F parent=0x67400000F
28146172	1	320.700363	emacs	2000	...1	preempt_enable	caller=0x71ed4eF parent=0x67400000F
28146173	1	320.700363	emacs	2000	d..2	irq_disable	caller=0xf662cF parent=0x67400000F
28146174	1	320.700363	emacs	2000	d..2	irq_enable	caller=0xf6650F parent=0x67400000F
28146175	1	320.700364	emacs	2000	d..2	irq_disable	caller=0xf678cF parent=0x67400000F
28146176	1	320.700364	emacs	2000	d..2	irq_enable	caller=0xf67b0F parent=0x67400000F
28146177	1	320.700364	emacs	2000	...1	preempt_disable	caller=0x2619e9F parent=0x2619e9F
28146178	1	320.700364	emacs	2000	d..2	irq_disable	caller=0xf662cF parent=0x67400000F
28146179	1	320.700365	emacs	2000	d..2	irq_enable	caller=0xf6650F parent=0x67400000F
28146180	1	320.700365	emacs	2000	d..2	irq_disable	caller=0xf678cF parent=0x67400000F
28146181	1	320.700365	emacs	2000	d..2	irq_enable	caller=0xf67b0F parent=0x67400000F
28146182	1	320.700366	emacs	2000	...1	preempt_enable	caller=0x2619ffF parent=0x2619ffF
28146183	1	320.700366	emacs	2000	d..2	irq_disable	caller=0xf662cF parent=0x67400000F
28146184	1	320.700366	emacs	2000	d..2	irq_enable	caller=0xf6650F parent=0x67400000F
28146185	1	320.700366	emacs	2000	....	kmem_cache_alloc_node	( _alloc_skb+0x57) call_site=fffff991c5527 ptr=0xffff9e5d3a8e8b00 bytes_req=232 bytes_alloc=256 gfp_flags=GFP_KERNEL_ACCOUNT _GFP_RETRY_MAYFAIL node=-1
28146186	1	320.700367	emacs	2000	d..2	irq_disable	caller=0xf678cF parent=0x67400000F
28146187	1	320.700367	emacs	2000	d..2	irq_enable	caller=0xf67b0F parent=0x67400000F
28146188	1	320.700367	emacs	2000	...1	preempt_disable	caller=0x25c06eF parent=0x25c06eF
28146189	1	320.700368	emacs	2000	d..2	irq_disable	caller=0xf662cF parent=0x67400000F
28146190	1	320.700368	emacs	2000	d..2	irq_enable	caller=0xf6650F parent=0x67400000F
28146191	1	320.700368	emacs	2000	d..2	irq_disable	caller=0xf678cF parent=0x67400000F

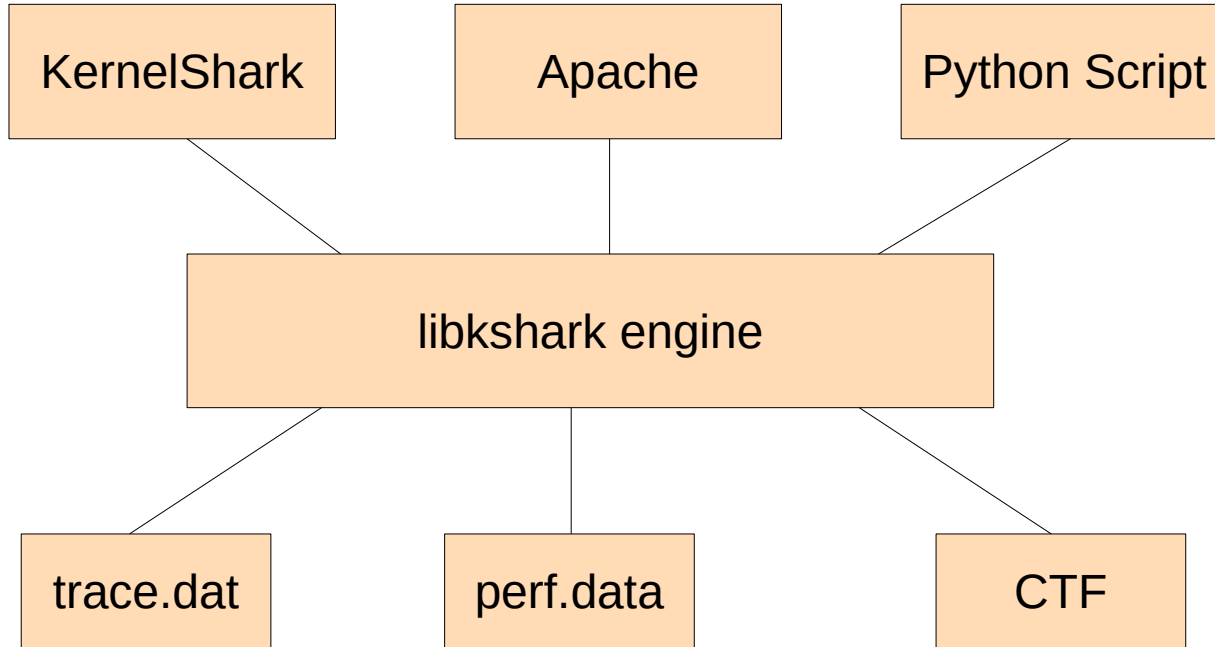
Kernel Shark (../rt-tests/...) kernelshark(trace.dat) Terminal - rostedt@vm... Terminal - rostedt@vm...

Net | cpu mem | swap 1:16

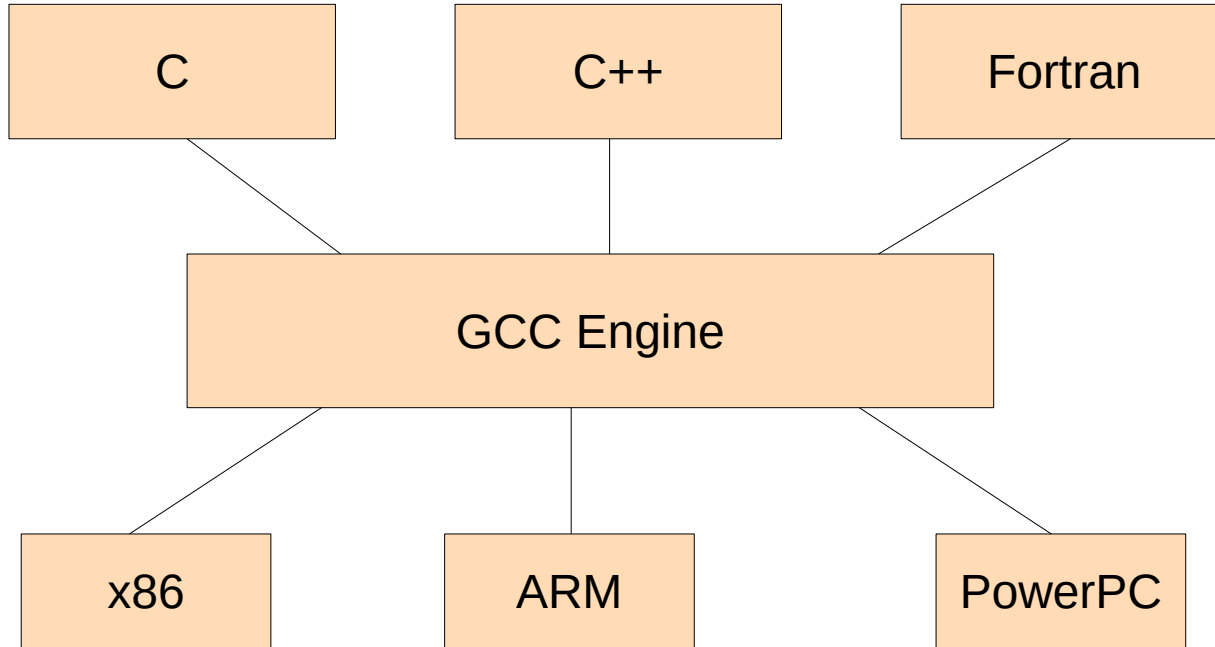
# What's next after 1.0?

- libkshark.so
  - The guts of KernelShark is being “modularized”
  - The KernelShark application is just a “shell”
  - Any tool will be able to use it
  - Also will be used by Python applications (libkshark.py)
- Will read multiple formats
  - Not just trace.dat
  - perf.data
  - Common Trace Format (CTF)
- Plugins
  - C, C++, Python

# What's next after 1.0?



# What's next after 1.0?





# What's next after 1.0?

- More than one view
  - Plot view (what it currently is)
  - Graphs
  - Flame graphs (see Brendan Gregg's presentations)
- Will show multiple machines
  - Show interactions between hosts and multiple guests

**Questions?**

**We have time!**

**DEMO!**