

# Swimming with the New KernelShark

Yordan Karadzhov

VMware Inc. - OSTC

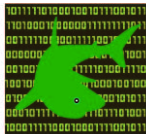
# What is KernelShark?



**KERNELSHARK**

- \* Front end reader of Linux kernel tracing data (Ftrace)

# What is KernelShark?

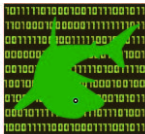


**KERNELSHARK**

- \* Front end reader of Linux kernel tracing data (Ftrace)

- \* The original version - started in 2009.
- \* Written in *Gtk+-2.0*
- \* Main goal: analyse and fully understood the performance of the Real-time scheduler.

# What is KernelShark?



**KERNELSHARK**

- \* **Front end reader of Linux kernel tracing data (Ftrace)**

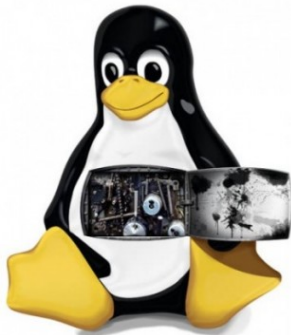
- \* **New KernelShark:** use all lessons learned from the old version.
- \* Completely rewritten to use Qt.
- \* But not only this ...

# The New KernelShark is

- a. Optimized for processing significantly larger amounts of data.
- b. New scalable data model -  $\log(n)$  time complexity.
- c. OpenGL-based visualization.
- d. Preconfigurable - Json config I/O.
- e. User modifiable - plugins.

# Why do we need KernelShark?

- a. Ftrace - the official tracing infrastructure of the Linux kernel.
- b. Extremely powerful instrument.
- c. Allows to see what is happening in the kernel.
- d. But you must know what you are looking for



# Ftrace is extremely powerful but ...

```
Compositor-1124 [001] .... 236364.639085: _cond_resched <-futex wait_queue_me
Compositor-1124 [001] .... 236364.639086: rcu_all_qs <-_cond_resched
chromium-browse-1117 [000] d... 236364.639086: do_syscall_64 <-entry SYSCALL_64 after_hwframe
Compositor-1124 [001] .... 236364.639086: drop_futex_key_refs.isra.14 <-futex wait
chromium-browse-1117 [000] .... 236364.639087: syscall_trace_enter <-do_syscall_64
Compositor-1124 [001] .... 236364.639087: hrtimer_cancel <-futex wait
chromium-browse-1117 [000] .... 236364.639087: _secure_computing <-syscall_trace_enter
Compositor-1124 [001] .... 236364.639087: hrtimer_try_to_cancel <-hrtimer_cancel
chromium-browse-1117 [000] .... 236364.639087: _seccomp_filter <-_secure_computing
Compositor-1124 [001] .... 236364.639088: hrtimer_active <-hrtimer_try_to_cancel
chromium-browse-1117 [000] .... 236364.639088: seccomp_run_filters <-_seccomp_filter
Compositor-1124 [001] .... 236364.639088: lock_hrtimer_base.isra.20 <-hrtimer_try_to_cancel
Compositor-1124 [001] .... 236364.639088: _raw_spin_lock_irqsave <-lock_hrtimer_base.isra.20
chromium-browse-1117 [000] .... 236364.639088: Sys_futex <-do_syscall_64
Compositor-1124 [001] d... 236364.639089: _remove_hrtimer <-hrtimer_try_to_cancel
Compositor-1124 [001] d... 236364.639089: _raw_spin_unlock_irqrestore <-hrtimer_try_to_cancel
chromium-browse-1117 [000] .... 236364.639089: ktime_get <-Sys_futex
chromium-browse-1117 [000] .... 236364.639090: ktime_add_safe <-Sys_futex
chromium-browse-1117 [000] .... 236364.639090: do_futex <-Sys_futex
chromium-browse-1117 [000] .... 236364.639091: futex_wait <-do_futex
chromium-browse-1117 [000] .... 236364.639091: hrtimer_init <-futex wait
chromium-browse-1117 [000] .... 236364.639092: _hrtimer_init <-hrtimer_init
chromium-browse-1117 [000] .... 236364.639092: hrtimer_init_sleep <-futex wait
chromium-browse-1117 [000] .... 236364.639092: ktime_add_safe <-futex wait
chromium-browse-1117 [000] .... 236364.639093: futex_wait_setup <-futex wait
chromium-browse-1117 [000] .... 236364.639093: get_futex_key <-futex wait_setup
chromium-browse-1117 [000] .... 236364.639093: get_futex_key_refs.isra.13 <-get_futex_key
chromium-browse-1117 [000] .... 236364.639094: hash_futex <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: _raw_spin_lock <-futex wait_setup
chromium-browse-1117 [000] .... 236364.639095: get_futex_value_locked <-futex wait_setup
chromium-browse-1117 [000] .... 236364.639095: futex_wait_queue_me <-futex wait
chromium-browse-1117 [000] .... 236364.639096: hrtimer_start_range_ns <-futex wait_queue_me
Compositor-1124 [001] d... 236364.639097: do_syscall_64 <-entry SYSCALL_64 after_hwframe
chromium-browse-1117 [000] .... 236364.639097: lock_hrtimer_base.isra.20 <-hrtimer_start_range_ns
```

# Ftrace is extremely powerful but ...

```
Compositor-1124 [001] .... 236364.639085: cond_resched <-futex wait_queue_me
Compositor-1124 [001] .... 236364.639086: rcu_all_qs <-_cond_resched
chromium-browse-1117 [000] d... 236364.639086: do_syscall_64 <-entry_SYSCALL_64_after_hwframe
Compositor-1124 [001] .... 236364.639086: do_syscall_64 <-entry_SYSCALL_64_after_hwframe
```

... it can be hard to find what you need. You must be  
**True (kernel) detective.**

```
Compositor-1124 [001] .... 236364.639088: nrtimer_active <-nrtimer_try_to_cancel
chromium-browse-1117 [000] .... 236364.639088: seccomp_run_filters <-_seccomp_filter
Compositor-1124 [001] .... 236364.639088: lock_hrtimer_base.isra.20 <-hrtimer_try_to_cancel
Compositor-1124 [001] .... 236364.639088: _raw_spin_lock_irqsave <-lock_hrtimer_base.isra.20
chromium-browse-1117 [000] .... 236364.639088: Sys_futex <-do_syscall_64
Compositor-1124 [001] d... 236364.639089: _remove_hrtimer <-hrtimer_try_to_cancel
Compositor-1124 [001] d... 236364.639089: _raw_spin_unlock_irqrestore <-hrtimer_try_to_cancel
chromium-browse-1117 [000] .... 236364.639089: ktime_get <-Sys_futex
chromium-browse-1117 [000] .... 236364.639090: ktime_add_safe <-Sys_futex
chromium-browse-1117 [000] .... 236364.639090: do_futex <-Sys_futex
chromium-browse-1117 [000] .... 236364.639091: futex_wait <-do_futex
chromium-browse-1117 [000] .... 236364.639091: hrtimer_init <-futex_wait
chromium-browse-1117 [000] .... 236364.639092: _hrtimer_init <-hrtimer_init
chromium-browse-1117 [000] .... 236364.639092: hrtimer_init_sleep <-futex_wait
chromium-browse-1117 [000] .... 236364.639092: ktime_add_safe <-futex_wait
chromium-browse-1117 [000] .... 236364.639093: futex_wait_setup <-futex_wait
chromium-browse-1117 [000] .... 236364.639093: get_futex_key <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639093: get_futex_key_refs.isra.13 <-get_futex_key
chromium-browse-1117 [000] .... 236364.639094: hash_futex <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: _raw_spin_lock <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: get_futex_value_locked <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: futex_wait_queue_me <-futex_wait
chromium-browse-1117 [000] .... 236364.639096: hrtimer_start_range_ns <-futex_wait_queue_me
Compositor-1124 [001] d... 236364.639097: do_syscall_64 <-entry_SYSCALL_64_after_hwframe
chromium-browse-1117 [000] .... 236364.639097: lock_hrtimer_base.isra.20 <-hrtimer_start_range_ns
```



# Ftrace is extremely powerful but ...

```
Compositor-1124 [001] .... 236364.639085: cond_resched <-futex wait_queue_me
Compositor-1124 [001] .... 236364.639086: rcu_all_qs <- cond_resched
chromium-browse-1117 [000] d... 236364.639086: do_syscall_64 <-entry SYSCALL_64 after_hwframe
Compositor-1124 [001] .... 236364.639086: do_syscall_64 <-entry SYSCALL_64 after_hwframe
```

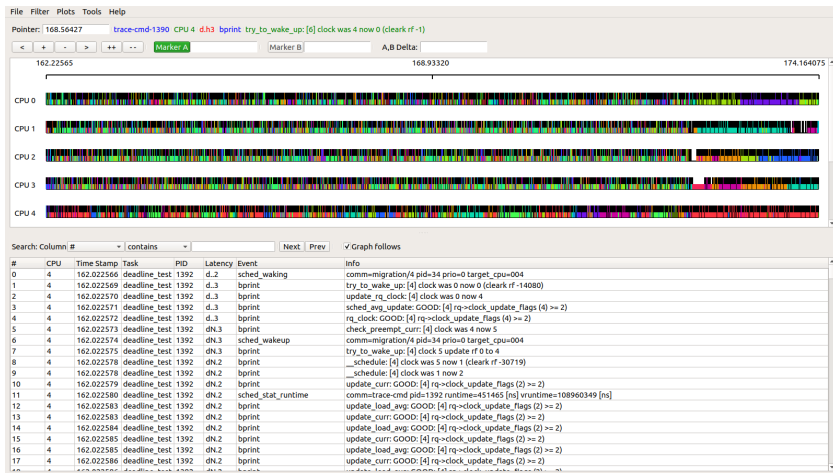
... it can be hard to find what you need. You must be  
**True (kernel) detective.**

```
Compositor-1124 [001] .... 236364.639088: hrtimer_active <-hrtimer_try_to_cancel
chromium-browse-1117 [000] .... 236364.639088: seccomp_run_filters <- seccomp_filter
Compositor-1124 [001] .... 236364.639088: lock_hrtimer_base.isra.20 <-hrtimer_try_to_cancel
Compositor-1124 [001] .... 236364.639088: raw_spin_lock_irqsave <-lock_hrtimer_base.isra.20
chromium-browse-1117 [000] .... 236364.639088: Sys_futex <-do_syscall_64
Compositor-1124 [001] d... 236364.639089: remove_hrtimer <-hrtimer_try_to_cancel
Compositor-1124 [001] d... 236364.639089: remove_hrtimer <-hrtimer_try_to_cancel
```

Interactive visualization can be  
very useful.

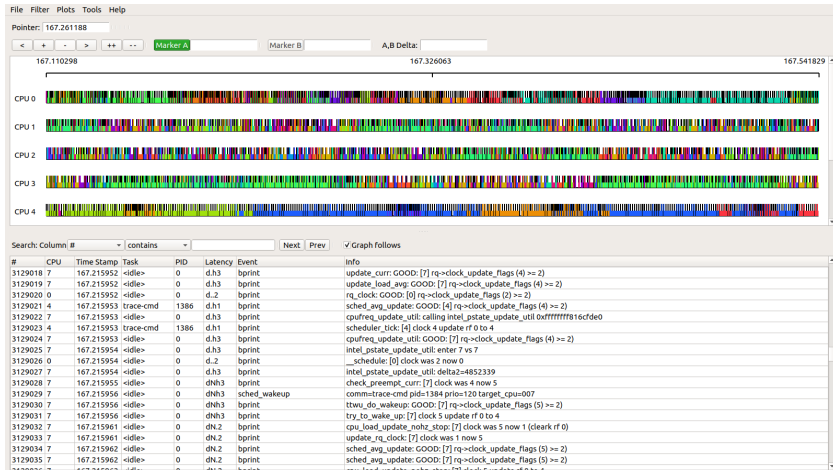
```
chromium-browse-1117 [000] .... 236364.639092: ktime_add_safe <-futex_wait
chromium-browse-1117 [000] .... 236364.639093: futex_wait_setup <-futex_wait
chromium-browse-1117 [000] .... 236364.639093: get_futex_key <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639093: get_futex_key_refs.isra.13 <-get_futex_key
chromium-browse-1117 [000] .... 236364.639094: hash_futex <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: raw_spin_lock <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: get_futex_value_locked <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: futex_wait_queue_me <-futex_wait
chromium-browse-1117 [000] .... 236364.639096: hrtimer_start_range_ns <-futex_wait_queue_me
Compositor-1124 [001] d... 236364.639097: do_syscall_64 <-entry SYSCALL_64 after_hwframe
chromium-browse-1117 [000] .... 236364.639097: lock_hrtimer_base.isra.20 <-hrtimer_start_range_ns
```

# Kernel Shark



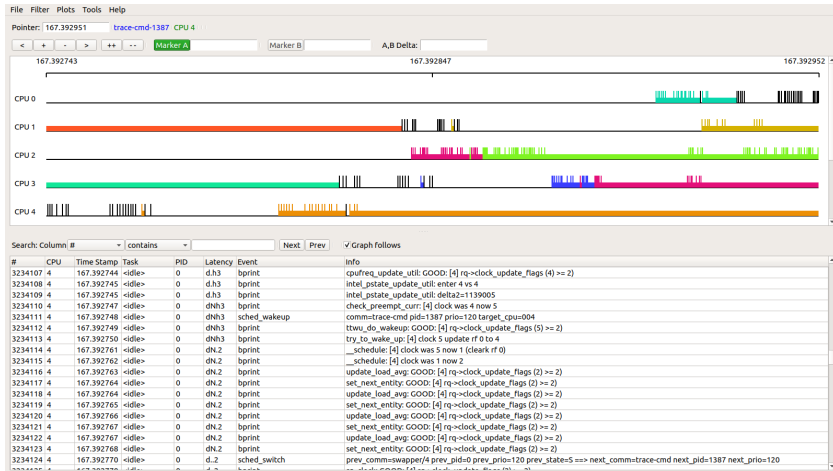
Hmm, why do you think that this complete mess of colors can be of any help?

# Kernel Shark



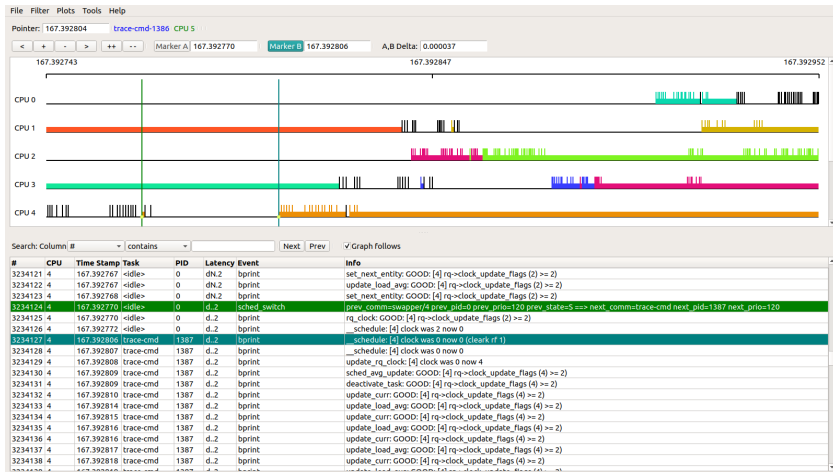
If we **Zoom in** a bit, we will start seeing the time structure of the trace.

# Kernel Shark



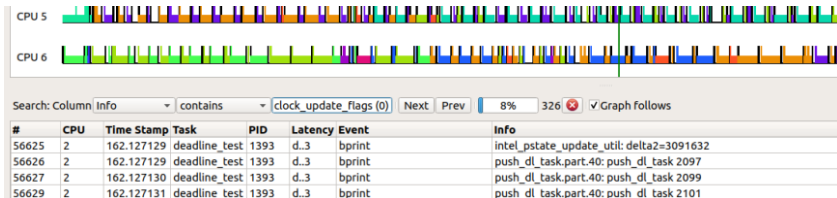
This is a deep zoom. Let's try to understand what is visualized here.

# Kernel Shark



- Double click on the graph selects a trace record.
- The same can be done by clicking on the table row which shows the record as text.

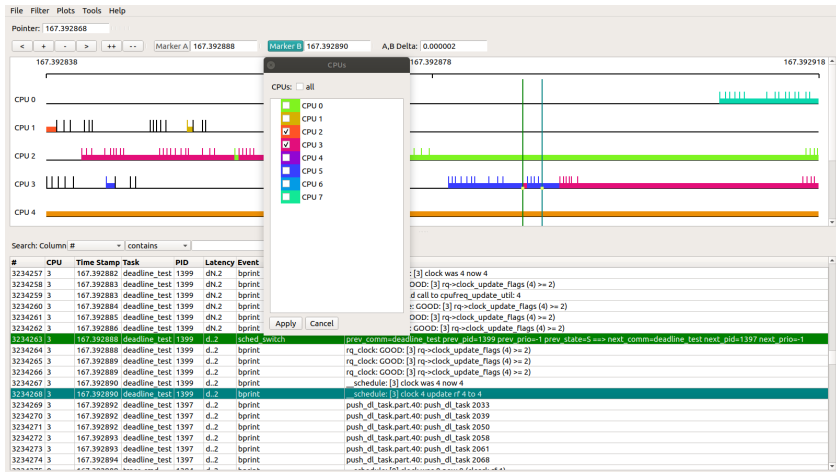
# Kernel Shark



## Search panel

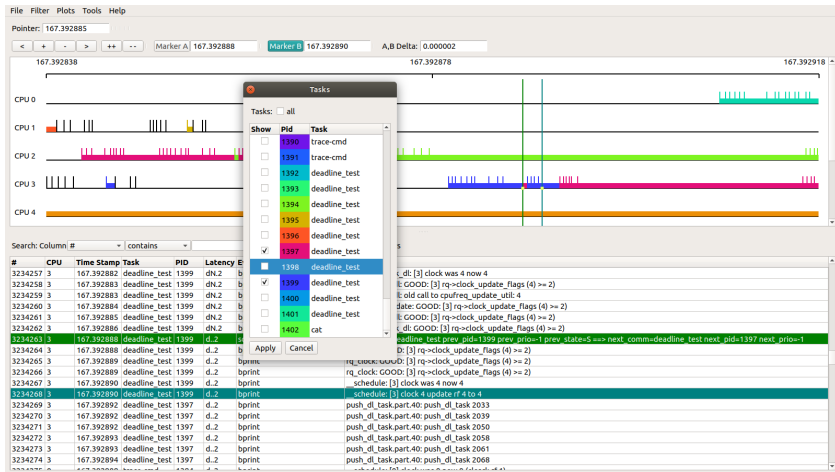
- Provides interactive searching.
- Very useful in combination with the Dual marker.

# Kernel Shark



Selecting which CPU to visualize.

# Kernel Shark



Selecting which Task to visualize.



# Kernel Shark

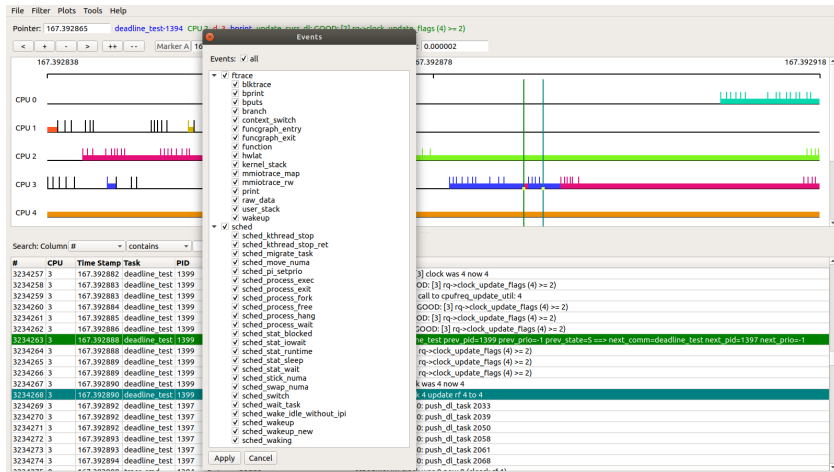


## Filter menu

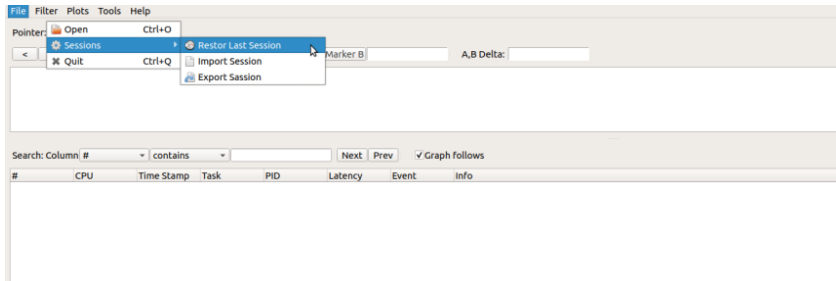
- Filter out tasks
- Filter in tasks
- Filter events
- Advanced (content-based) filtering

# Kernel Shark

## Event filters

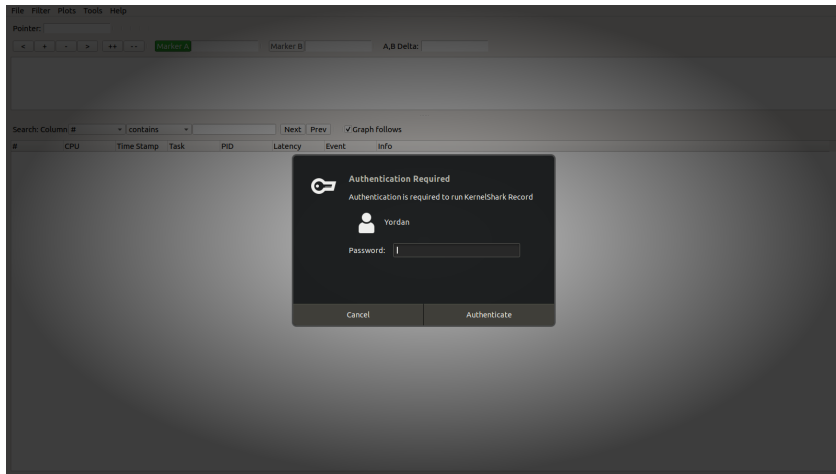


# Kernel Shark - Sessions



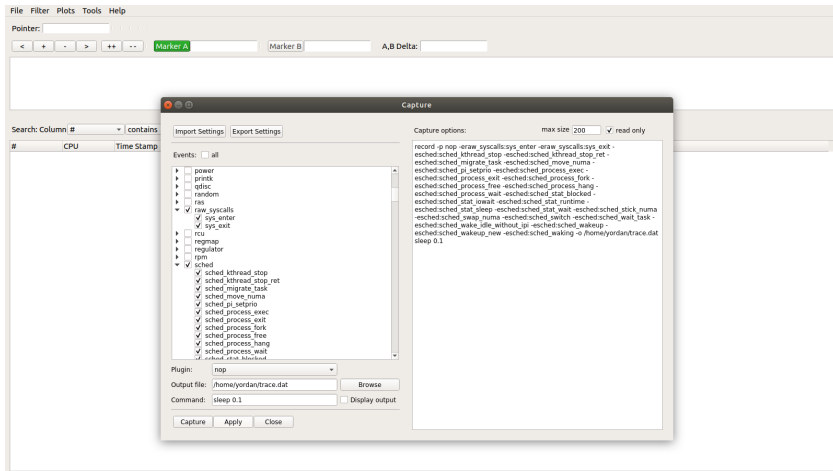
- a. Import/Export session
- b. Restor last session.

# Kernel Shark - Recording



- a. Trace data can be recorded directly from Kernel Shark.
- b. Root password is required.

# Kernel Shark - Recording



# Ftrace is extremely powerful but ...

```
Compositor-1124 [001] .... 236364.639085: _cond_resched <-futex wait_queue_me
Compositor-1124 [001] .... 236364.639086: rcu_all_qs <-_cond_resched
chromium-browse-1117 [000] d... 236364.639086: do_syscall_64 <-entry SYSCALL_64 after_hwframe
Compositor-1124 [001] .... 236364.639086: drop_futex_key_refs.isra.14 <-futex wait
chromium-browse-1117 [000] .... 236364.639087: syscall_trace_enter <-do_syscall_64
Compositor-1124 [001] .... 236364.639087: hrtimer_cancel <-futex wait
chromium-browse-1117 [000] .... 236364.639087: _secure_computing <-syscall_trace_enter
Compositor-1124 [001] .... 236364.639087: hrtimer_try_to_cancel <-hrtimer_cancel
chromium-browse-1117 [000] .... 236364.639087: _seccomp_filter <-_secure_computing
Compositor-1124 [001] .... 236364.639088: hrtimer_active <-hrtimer_try_to_cancel
chromium-browse-1117 [000] .... 236364.639088: seccomp_run_filters <-_seccomp_filter
Compositor-1124 [001] .... 236364.639088: lock_hrtimer_base.isra.20 <-hrtimer_try_to_cancel
Compositor-1124 [001] .... 236364.639088: _raw_spin_lock_irqsave <-lock_hrtimer_base.isra.20
chromium-browse-1117 [000] .... 236364.639088: Sys_futex <-do_syscall_64
Compositor-1124 [001] d... 236364.639089: _remove_hrtimer <-hrtimer_try_to_cancel
Compositor-1124 [001] d... 236364.639089: _raw_spin_unlock_irqrestore <-hrtimer_try_to_cancel
chromium-browse-1117 [000] .... 236364.639089: ktime_get <-Sys_futex
chromium-browse-1117 [000] .... 236364.639090: ktime_add_safe <-Sys_futex
chromium-browse-1117 [000] .... 236364.639090: do_futex <-Sys_futex
chromium-browse-1117 [000] .... 236364.639091: futex_wait <-do_futex
chromium-browse-1117 [000] .... 236364.639091: hrtimer_init <-futex wait
chromium-browse-1117 [000] .... 236364.639092: _hrtimer_init <-hrtimer_init
chromium-browse-1117 [000] .... 236364.639092: hrtimer_init_sleep <-futex wait
chromium-browse-1117 [000] .... 236364.639092: ktime_add_safe <-futex wait
chromium-browse-1117 [000] .... 236364.639093: futex_wait_setup <-futex wait
chromium-browse-1117 [000] .... 236364.639093: get_futex_key <-futex wait_setup
chromium-browse-1117 [000] .... 236364.639093: get_futex_key_refs.isra.13 <-get_futex_key
chromium-browse-1117 [000] .... 236364.639094: hash_futex <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: _raw_spin_lock <-futex wait_setup
chromium-browse-1117 [000] .... 236364.639095: get_futex_value_locked <-futex wait_setup
chromium-browse-1117 [000] .... 236364.639095: futex_wait_queue_me <-futex wait
chromium-browse-1117 [000] .... 236364.639096: hrtimer_start_range_ns <-futex wait_queue_me
Compositor-1124 [001] d... 236364.639097: do_syscall_64 <-entry SYSCALL_64 after_hwframe
chromium-browse-1117 [000] .... 236364.639097: lock_hrtimer_base.isra.20 <-hrtimer_start_range_ns
```

# Ftrace is extremely powerful but ...

```
Compositor-1124 [001] .... 236364.639085: cond_resched <-futex_wait_queue_me
Compositor-1124 [001] .... 236364.639086: rcu_all_qs <- cond_resched
chromium-browse-1117 [000] d... 236364.639086: do_syscall_64 <-entry SYSCALL_64 after_hwframe
Compositor-1124 [001] .... 236364.639086: drop_futex_key_refs.isra.14 <-futex_wait
chromium-browse-1117 [000] .... 236364.639087: syscall_trace_enter <-do_syscall_64
Compositor-1124 [001] .... 236364.639087: hrtimer_cancel <-futex_wait
```

## What do we do if:

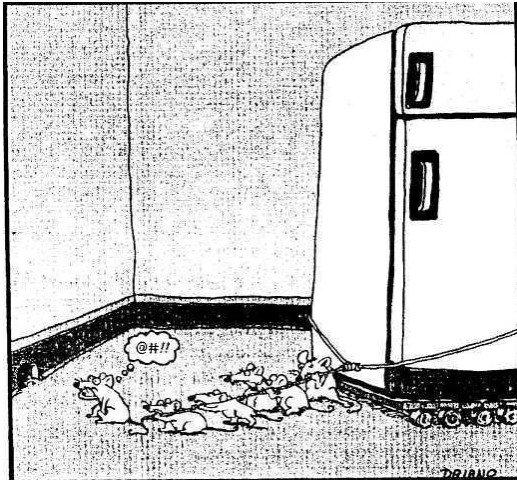
- We have a very large data-set of trace records.
- We are limited by the number of screen pixels available.

```
Compositor-1124 [001] d... 236364.639089: _raw_spin_unlock_irqrestore <-hrtimer_try_to_cancel
chromium-browse-1117 [000] .... 236364.639089: ktime_get <-SyS_futex
chromium-browse-1117 [000] .... 236364.639090: ktime_add_safe <-SyS_futex
chromium-browse-1117 [000] .... 236364.639090: do_futex <-SyS_futex
chromium-browse-1117 [000] .... 236364.639091: futex_wait <-do_futex
chromium-browse-1117 [000] .... 236364.639091: hrtimer_init <-futex_wait
chromium-browse-1117 [000] .... 236364.639092: hrtimer_init <-hrtimer_init
chromium-browse-1117 [000] .... 236364.639092: hrtimer_init_clopper <-futex_wait
```

## And we have to process all this in a reasonable amount of time.

```
chromium-browse-1117 [000] .... 236364.639093: get_futex_key_refs.isra.13 <-get_futex_key
chromium-browse-1117 [000] .... 236364.639094: hash_futex <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: _raw_spin_lock <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: get_futex_value_locked <-futex_wait_setup
chromium-browse-1117 [000] .... 236364.639095: futex_wait_queue_me <-futex_wait
chromium-browse-1117 [000] .... 236364.639096: hrtimer_start_range_ns <-futex_wait_queue_me
Compositor-1124 [001] d... 236364.639097: do_syscall_64 <-entry SYSCALL_64 after_hwframe
chromium-browse-1117 [000] .... 236364.639097: lock_hrtimer_base.isra.20 <-hrtimer_start_range_ns
```

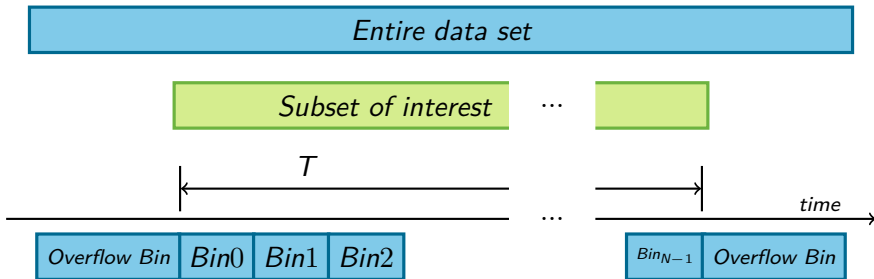
Fitting something large inside something small ...



... and we have to do this quickly!

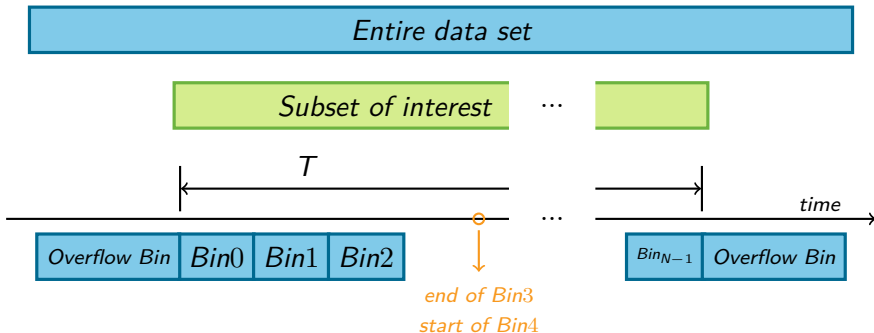


# Visualization model - How does it work?



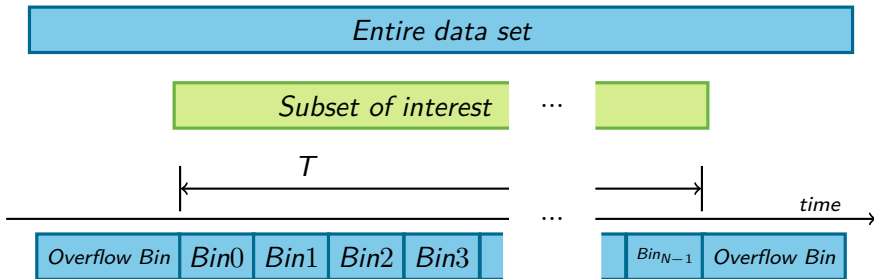
- Break the data-set into **time-bins**  $\mapsto$  like a histogram.
- Check only the records at the beginning and at the end of each bin.  
 $\mapsto$  constant time.

# Visualization model - How does it work?



- c. Have the trace records, sorted in time.
- d. Knowing the index of the first record in each *Bin* determines the state of the model.
- e. But the first element can be found with a binary search  $\mapsto \log(n)$  complexity.

# Visualization model - How does it work?

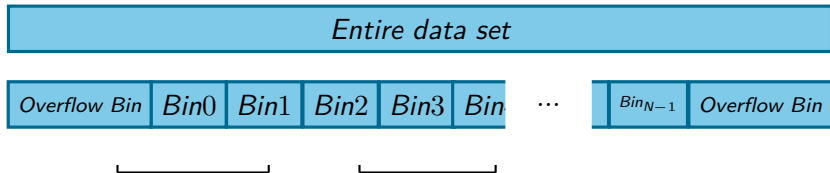


**Data Binning** provides  $O(\log_2(n))$  average time complexity of all operations of the model.

# Visualization model & tracing data formats

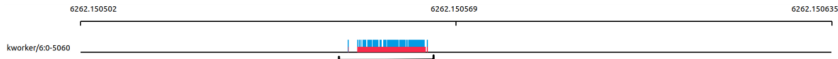
- The KernelShark Visualization model is not coupled to a particular data format.
- Uses KernelShark-specific data structure.
- Contains only the absolute minimum of information need by the model.
- The rest of the information - available on demand (can be slow)

# Visualization model & tracing data



- Only one model (data structure) for all graphs.
- Worst-case complexity becomes linear.
- Solution - Data collections.

# Visualization model & Data collections.



- Only one model (data structure) for all graphs.
- Worst-case complexity becomes linear.
- Solution - Data collections.

# DEMO

## KernelShark: current version 0.9

<https://git.kernel.org/pub/scm/utils/trace-cmd/trace-cmd.git/>

To build the code follow the instructions in

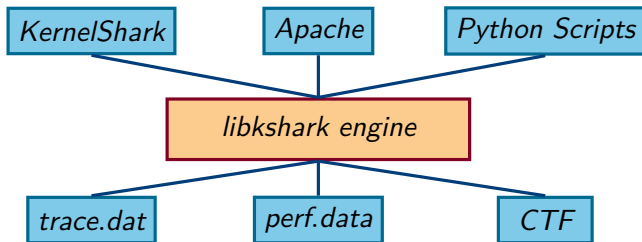
[/trace-cmd/kernel-shark-qt/README](#)

and

[/trace-cmd/README](#)



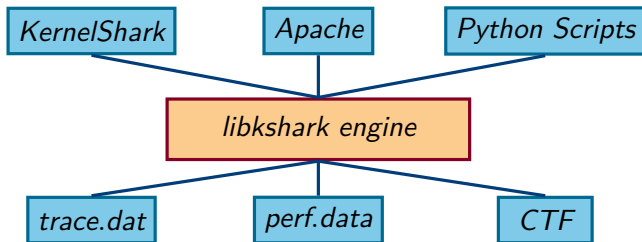
KernelShark is not a GUI. KernelShark is a toolkit.



## What's next after KernelShark 1.0?

- a. KernelShark engine (libkshark.so)
- b. Available under GNU LGPL v2.1
- c. Highly customizable (via plugins)
- d. Will read multiple data formats

KernelShark is not a GUI. KernelShark is a toolkit.



## What's next after KernelShark 1.0?

- a. Any tool will be able to use the library
- b. Available for Python applications (libkshark.py)
- c. The KernelShark application is just a “shell”.