# Idea of "networkingfy" Linux tracing

## Or "how Linux Tracing can use net/ subsystem"

Alexander Aring

# Why this talk idea?

- -ENOCODE yet – I am sorry!
- It's just an idea – Moonshot
- Any kind of feedback is welcome!
- Helps me: "Can this work? Good idea?"
- Is it too "crazy"?

Red Hat

# Tracing Basics

- Components
- Buffer Handling
- Common Filtering mechanism
- In reality everything is more complex

Red Hat

# Tracefs

- Special Filesystem

- UAPI to control Linux Tracing subsystem

- Low-Level controllable via coreutils

- Usually /sys/kernel/tracing

Red Hat

# TracePoint Declaration

```c
TRACE_EVENT(sk_data_ready,
        TP_PROTO(const struct sock *sk),

        TP_ARGS(sk),

        TP_STRUCT__entry(
                __field(const void *, skaddr)
                __field(__u16, family)
                __field(__u16, protocol)
                __field(unsigned long, ip)
        ),

        TP_fast_assign(
                __entry->skaddr = sk;
                __entry->family = sk->sk_family;
                __entry->protocol = sk->sk_protocol;
                __entry->ip = _RET_IP_;
        ),

        TP_printk("family=%u protocol=%u func=%ps",
                __entry->family, __entry->protocol, (void *)__entry->ip)
);
```

Name

Fields

Pretty Printer

Red Hat

# TraceEvent

```
void my_data_ready(struct sock *sk)
{
        trace_sk_data_ready(sk);
}
```
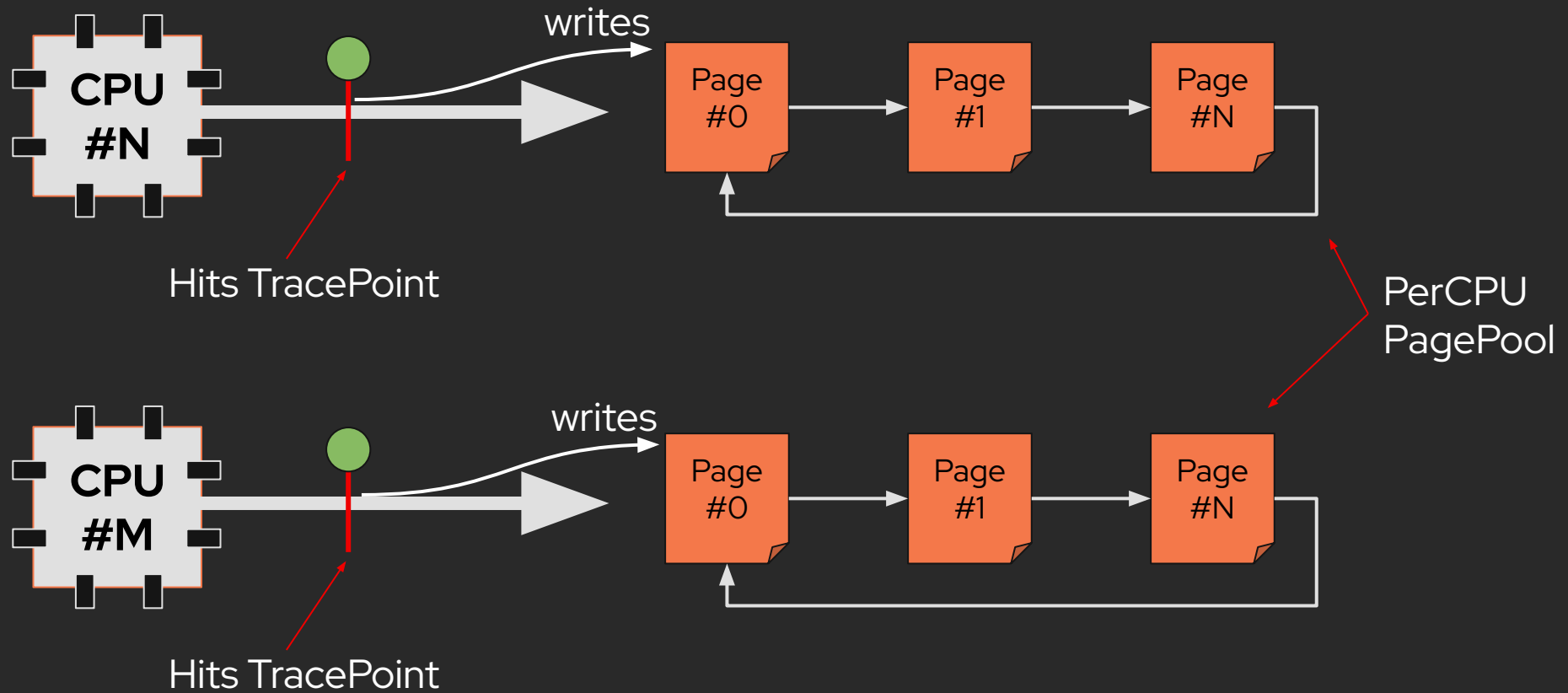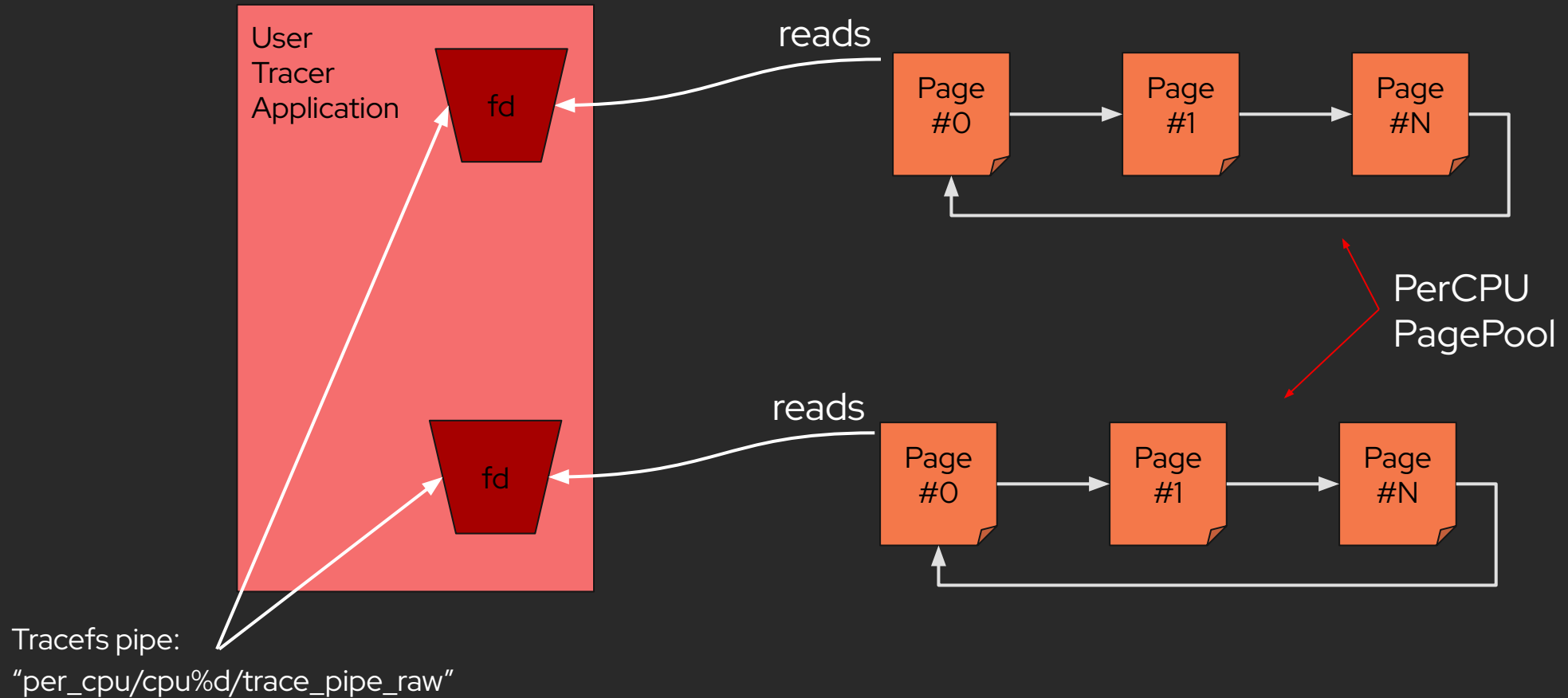
TracePoint generates
multiple Events

| Name | Field TLV #0 | Field TLV #1 | Field TLV #N | Name | Field TLV #0 | Field TLV #1 | Field TLV #N | Name | Field TLV #0 | Field TLV #1 | Field TLV #N |
|------|------|------|------|------|------|------|------|------|------|------|------|
| TraceEvent #1 | | | | TraceEvent #1 | | | | TraceEvent #N | | | |

Simplified: "Stream of Events"

# TraceBuffer – PagePool Producer



writes

CPU #N → Hits TracePoint → writes → Page #0 → Page #1 → Page #N

CPU #M → Hits TracePoint → writes → Page #0 → Page #1 → Page #N

PerCPU PagePool

# TraceBuffer – PagePool Consumer

User
Tracer
Application

fd

fd

reads

reads

Page #0

Page #1

Page #N

Page #0

Page #1

Page #N

PerCPU
PagePool

Tracefs pipe:
"per_cpu/cpu%d/trace_pipe_raw"

Red Hat

# In-Kernel Event Filtering

Common solution

In-Kernel Interpreter (yes there is one!)

"kernel/trace/trace_events_filter.c"

Mostly work on TLVs (numbers, strings)

Red Hat

# Kernel Tracer vs User Tracer

- **Kernel Tracer**
  - Lives in the Kernel
  - E.g. ftrace "kernel/trace/trace_functions.c"
- **User Tracer**
  - Lives in the Userspace "trace-pipe-raw"
  - Kernelshark 1), trace-cmd 2)
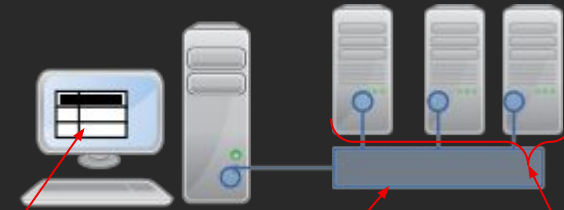
# Local vs Remote Tracing

## Local



Local Tracer

Generates
TraceEvents
for Local
Tracer

## Remote



Remote Tracer
receives
generated
TraceEvents

Switch

Generating
TraceEvents

Time
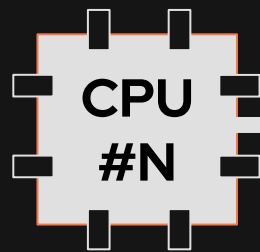Synchronized

Red Hat

# How we can adapt this to net/ ?

- Producer and Consumer

  - Operate directly on NIC DMA Rings
  - Abstract a Tracing Interface/Sockets?
  - Use net/core/page_pool, AF_XDP?
- Offload Filtering

  - Use existing Filtering Infrastructure
  - TC, eBPF Action, even P4?

Red Hat

# How we can adapt this to net/ ?

- Encapsulate LinkLayer e.g. Ethernet around raw data?

- Even TCP/IP based? If necessary?

- Local Tracing

  ○ Internal Loopback to RX Buffer?
  ○ Loopback cable?

- Remote Tracing

  ○ Send the data directly to remote Machine
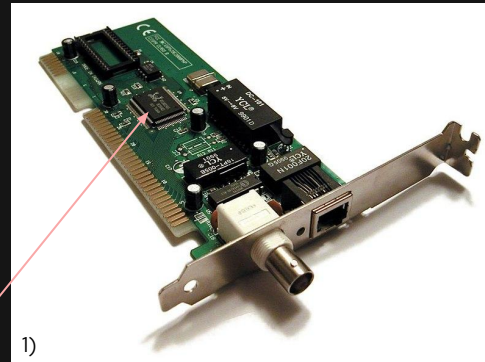  ○ Time Synchronized Tracing
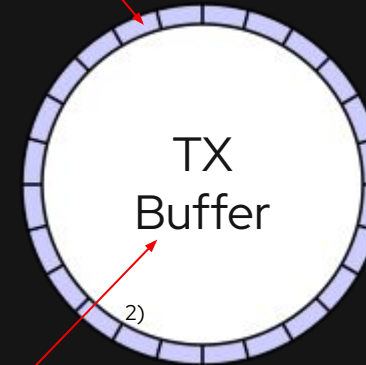
# Tracing – Producer



NUMA Node

Writes TraceEvent
Raw Data into

**CPU #N**

Hits TracePoint

ASIC

Pinned to a CPU?
Possible?

TX Buffer
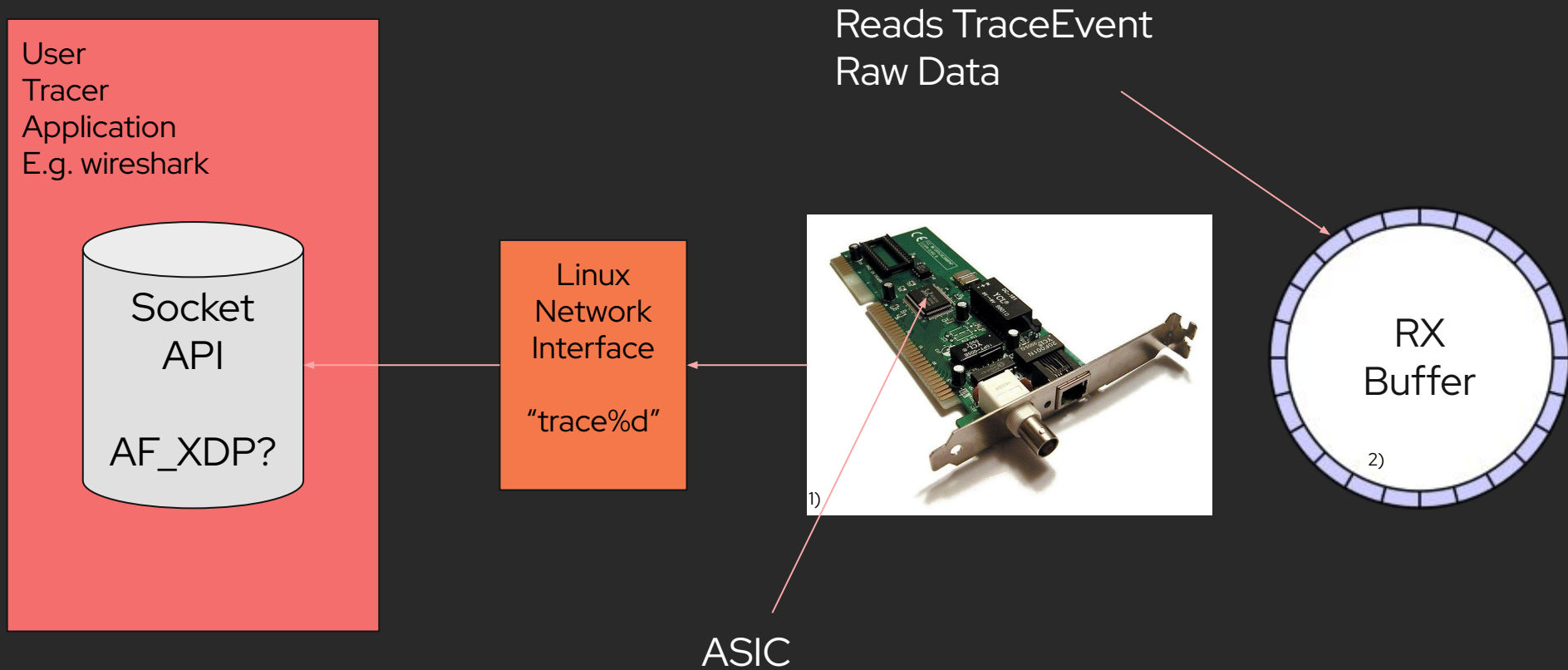
1)

2)

# DMA Ring as Tracing Buffer

- Pin DMA Ring to CPU? (Our Page Pool)

- NUMA Node requirement

  ○ CPU and NIC on same NUMA Node
  ○ Avoid Traversing over internal Bus
  ○ See other Netdev Talk 1)

1) https://netdevconf.info/0x18/sessions/talk/multi-pf-single-netdev.html

# Tracing – Consumer

User
Tracer
Application
E.g. wireshark

Socket
API

AF_XDP?

Linux
Network
Interface

"trace%d"

Reads TraceEvent
Raw Data

RX
Buffer

2)

1)

ASIC

# Tracing Filtering and net/

- Kernel knows the Metadata!

- Can be "discovered" by Tracefs

- Userspace Application only (Control Plane)
  - Easy to use Key-Value (Events) pairs
  - u32 (offload?), eBPF -> Action DROP
  - In Software - May faster than Interpreter?

# Traceevent Dropping and net/

- TraceEvent Dropping

  - It is a topic in Linux Tracing (avoid buffer bloating)
  - Unreliable e.g. currently PagePool is Full
  - Reliable Protocols for Tracing Data?
  - Which event to drop?
- Qdisc for Linux Tracing?

Red Hat

# Next steps? Proof of Concept?

- Ignore the DMA Ring-Buffer for now

- Focus on the virtual Tracing networking Interface

- Local Tracing only (put Events in a skb)

- Avoid recursion tracing cases

- Wireshark as Linux Tracer (AF_PACKET)

    - Shows Traceevent TLVs
    - Dissector configurable during runtime?

Red Hat

# Future steps? Try to Filtering!

- Create user space app to configure Tracing Filter

  - Operates on Tracefs
  - Reads Metadata configure existing Networking Filtering techniques to apply filtering
  - Observe Wireshark Tracer
- Simulate "Remote Tracing" over veth?

**Red Hat**

# Future steps? Look for Performance?

- Try to use real hardware

- Use DMA TX/RX Rings
  (AF_XDP/"net/core/page_pool.c")?

- Try to offload Filtering on NICs ASIC

Red Hat

# Future steps? Time synchronization?

- SO_TIMESTAMPING 1) ? Willem de Bruijn

- Additional metadata TLV required?

- Causality requirement (Events in Order as they appeared in the Network)

Red Hat

# Thank you

Red Hat