

Real world tips, tricks, and notes of using epoll-based busy polling to reduce latency

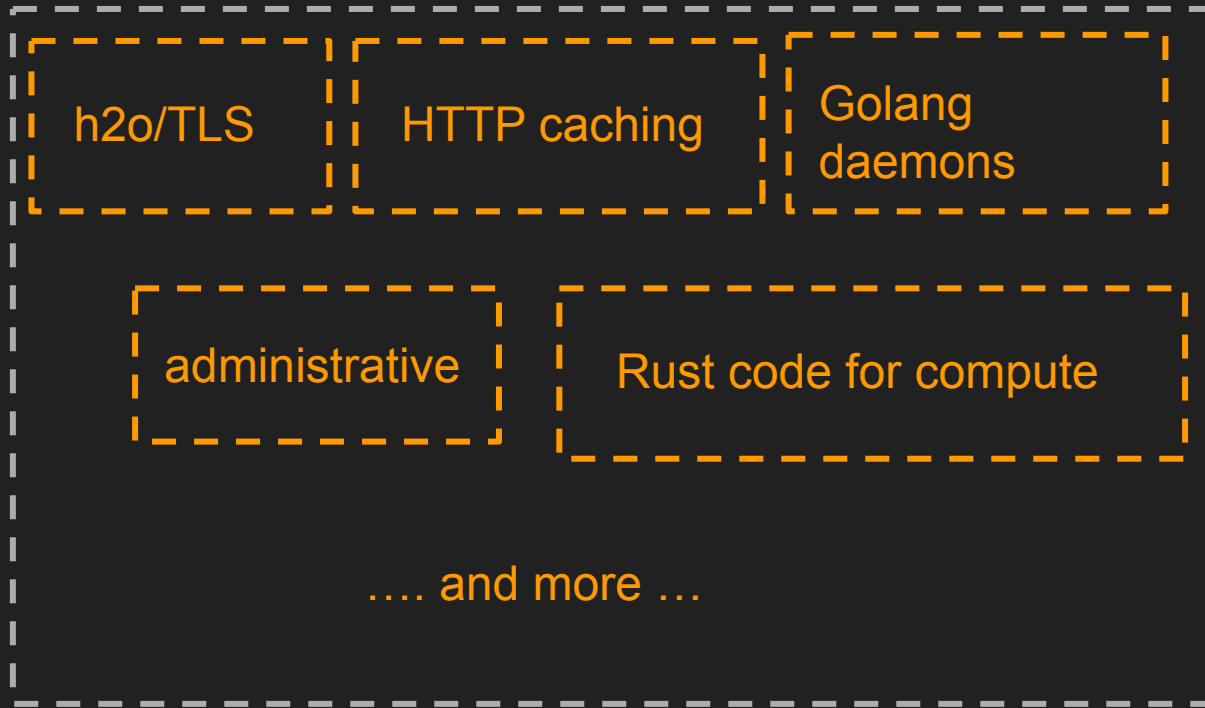
Joe Damato

Hi, my name is Joe.

I work at Fastly.

My opinions are my own.

A Fastly computer looks like this

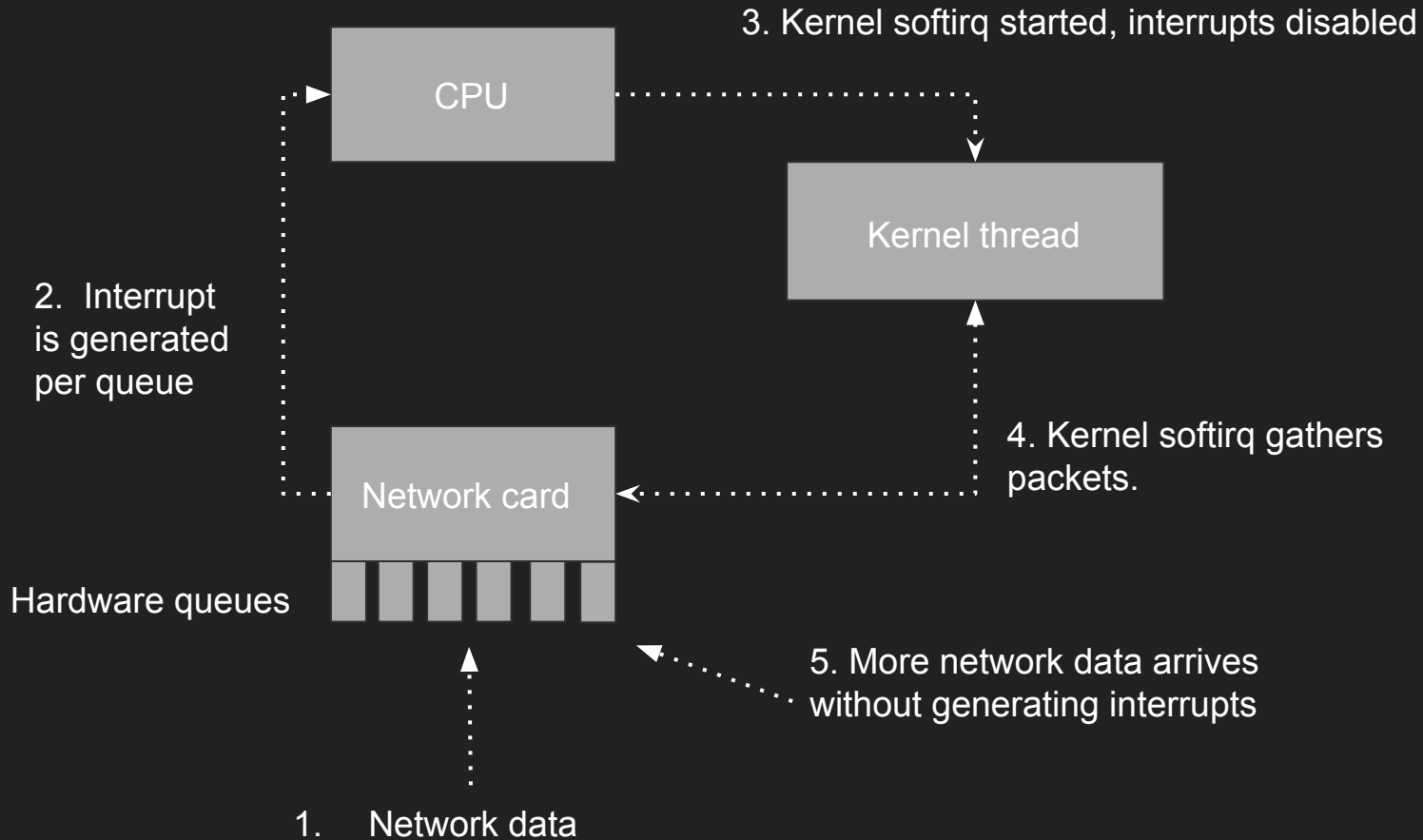


Overall takeaways:

- Busy poll helped us reduce latency in prod
- Needed some tweaks to get us there
- System-wide config is too coarse-grained
- More work to be done
- More docs needed overall

But before we can talk about
busy polling: NAPI.

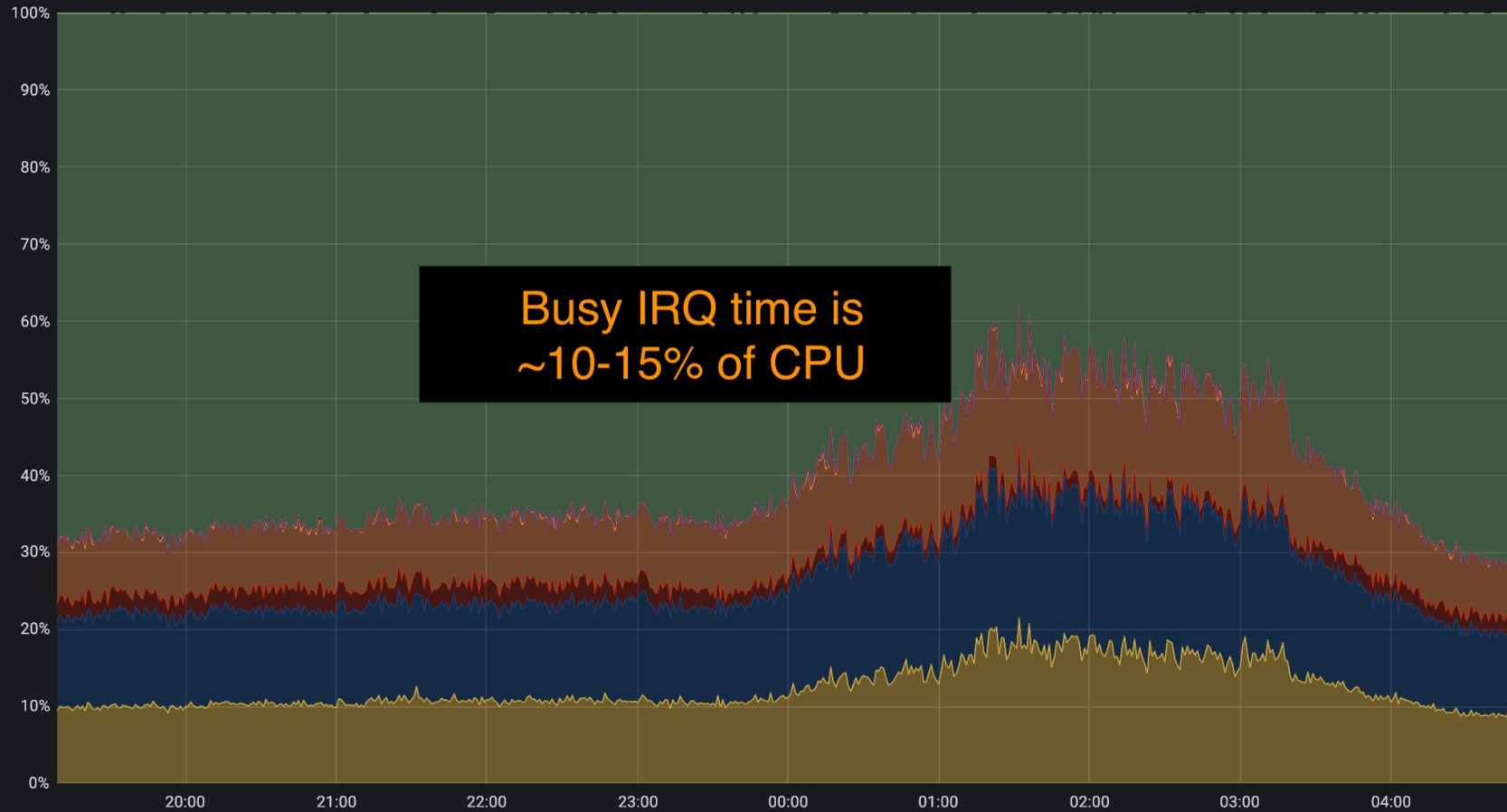
[Very high-level diagram]



NAPI

- Reduce interrupt load
- Improve CPU cache efficiency
- Should benefit high speed NICs (faster speed = faster data arrival = more interrupts)

But ... real world data from
a prod machine with lots of
traffic, suggests....



Busy IRQ time is
~10-15% of CPU

This felt like something was wrong.

So much IRQ traffic.

NAPI efficiency - RX packets per IRQ



First thing I tried, reducing the number of queues and configuring IRQ coalescing.

Simplest thing first:
push tx-usecs out

i40e NIC

25gbps NIC

Driver modified to count IRQs

iperf3 test



Zero work readings on i40e during production test

10 second timeout

	test-node-1 (reduced q + TX coalesce)	test-node-2 (tx coalesce)	control
Run 1:	4,013,254	4,659,321	4,558,964
Run 2:	4,059,272	4,774,618	4,833,400
Run 3:	4,043,220	4,858,869	4,841,198
Run 4:	4,124,204	4,781,532	4,947,740
Run 5:	3,730,169	4,199,293	4,184,878

60 second time out

	test-node-1 (reduced q + TX coalesce)	test-node-2 (tx-coalesce)	control
Run 1:	25,634,210	30,087,788	30,231,770
Run 2:	25,340,944	29,948,818	30,862,257
Run 3:	23,767,422	27,684,464	28,243,535
Run 4:	22,551,608	26,401,489	26,392,595
Run 5:	22,625,945	26,109,389	26,695,311

TLDR; reduction in “no-work” IRQs with
i40e, maybe ~15% or so?

	test		control	
Run 1:	7,562,960	vs.	7,849,544	(+286,584 no work IRQs on control)
Run 2:	7,461,928	vs.	7,917,925	(+455,997 no work IRQs on control)
Run 3:	7,477,857	vs.	7,546,702	(+68,845 no work IRQs control)
Run 4:	7,456,894	vs.	7,729,559	(+272,665 no work IRQs on control)

Longer test, timeout 60s:

	test		control	
Run 1:	47,009,259	vs.	49,755,821	(+2,746,562 no work IRQs on control)
Run 2:	46,678,501		vs. 49,547,378	(+2,868,877 no work IRQs on control)
Run 3:	46,313,645		vs. 49,955,651	(+3,642,006 no work IRQs on control)
Run 4:	43,709,993		vs. 47,725,595	(+4,025,592 no work IRQs on control)

TLDR; 3-5% improvement in “no-work”

IRQs on mlx5

Tiny improvement, but still
felt like there was
something I was missing?

Lots of reading of kernel code

Stumbled on `/sys/class/net/*/....`

`napi_defer_hard_irqs`
`gro_flush_timeout`

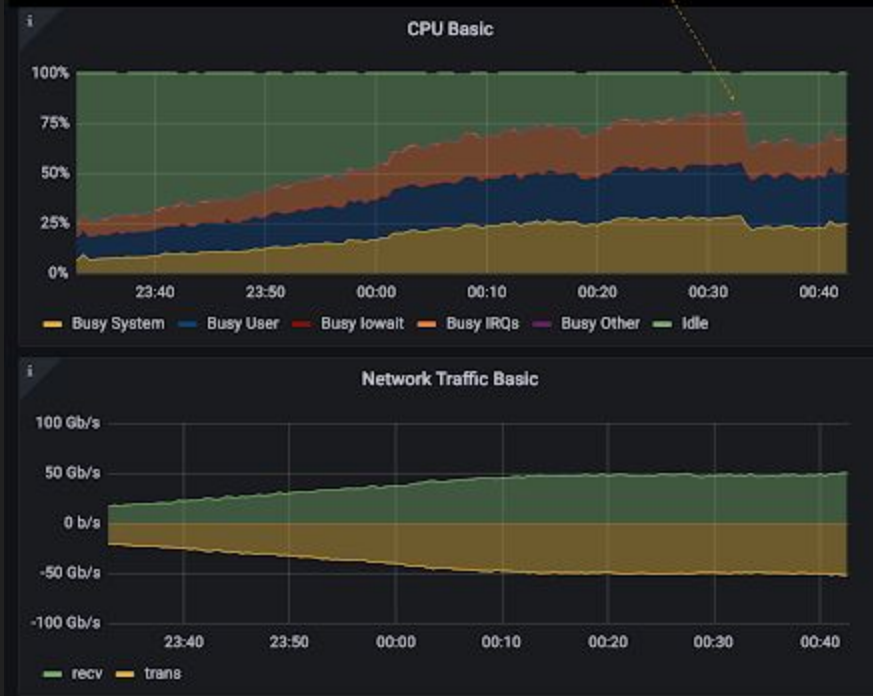
These settings are a
(system-wide) way to keep NIC
IRQs disabled (sw IRQ
coalescing)

And process incoming data via a
timer instead

gro timer + defer hard irq

Take away:

Controlling RX IRQs
reduces busy IRQ
tremendously under
load



Using these settings reduced CPU usage *considerably*.

... which shows that RX IRQs are a significant source of CPU use.

BUT

Latency was introduced.

Timer fires more slowly than the device fires IRQs.

Lots of reading of kernel code

SO_INCOMING_NAPI_ID

No clue what this thing was or what it did.

Stumbled on an email from Cong
Wang from 2019 who also didn't know
what it did

From: Cong Wang <congwang@kernel.org> 2019-02-14 20:15 UTC ([permalink](#) / [raw](#))

To: Alexander Duyck

Cc: Eric Dumazet, sridhar.samudrala, Linux Kernel Network Developers

Hello,

While looking into the busy polling in Linux kernel, three questions come into my mind:

1. In the document[1], it claims `sysctl.net.busy_poll` depends on either `SO_BUSY_POLL` or `sysctl.net.busy_read`. However, from the code in `ep_set_busy_poll_napi_id()`, I don't see such a dependency. It simply checks `sysctl_net_busy_poll` and `sk->sk_napi_id`, but `sk->sk_napi_id` is always set as long as we enable `CONFIG_NET_RX_BUSY_POLL`. So what I am missing here?

2. Why there is no socket option for `sysctl.net.busy_poll`? Clearly `sysctl_net_busy_poll` is global and `SO_BUSY_POLL` only works for `sysctl.net.busy_read`.

3. How is `SO_INCOMING_NAPI_ID` supposed to be used? I can't find any useful documents online. Any example or more detailed doc?

Thanks!

Turns out that Cong Wang writes
kernel code

Actually, a huge amount of
kernel code

If Cong Wang doesn't know
what this thing is.....

then it's *not too surprising* I
had no idea what it did either.

Some documentation in the
man page (man 7 socket)

SO_INCOMING_NAPI_ID (gettable since Linux 4.12)

Returns a system-level unique ID called NAPI ID that is associated with a RX queue on which the last packet associated with that socket is received.

This can be used by an application to split the incoming flows among worker threads based on the RX queue on which the packets associated with the flows are received. It allows each worker thread to be associated with a NIC HW receive queue and service all the connection requests received on that RX queue. This mapping between a app thread and a HW NIC queue streamlines the flow of data from the NIC to the application.

Vaguely seemed like flow
steering?

Supported by [memcached](#)

Introduce NAPI ID based worker thread selection

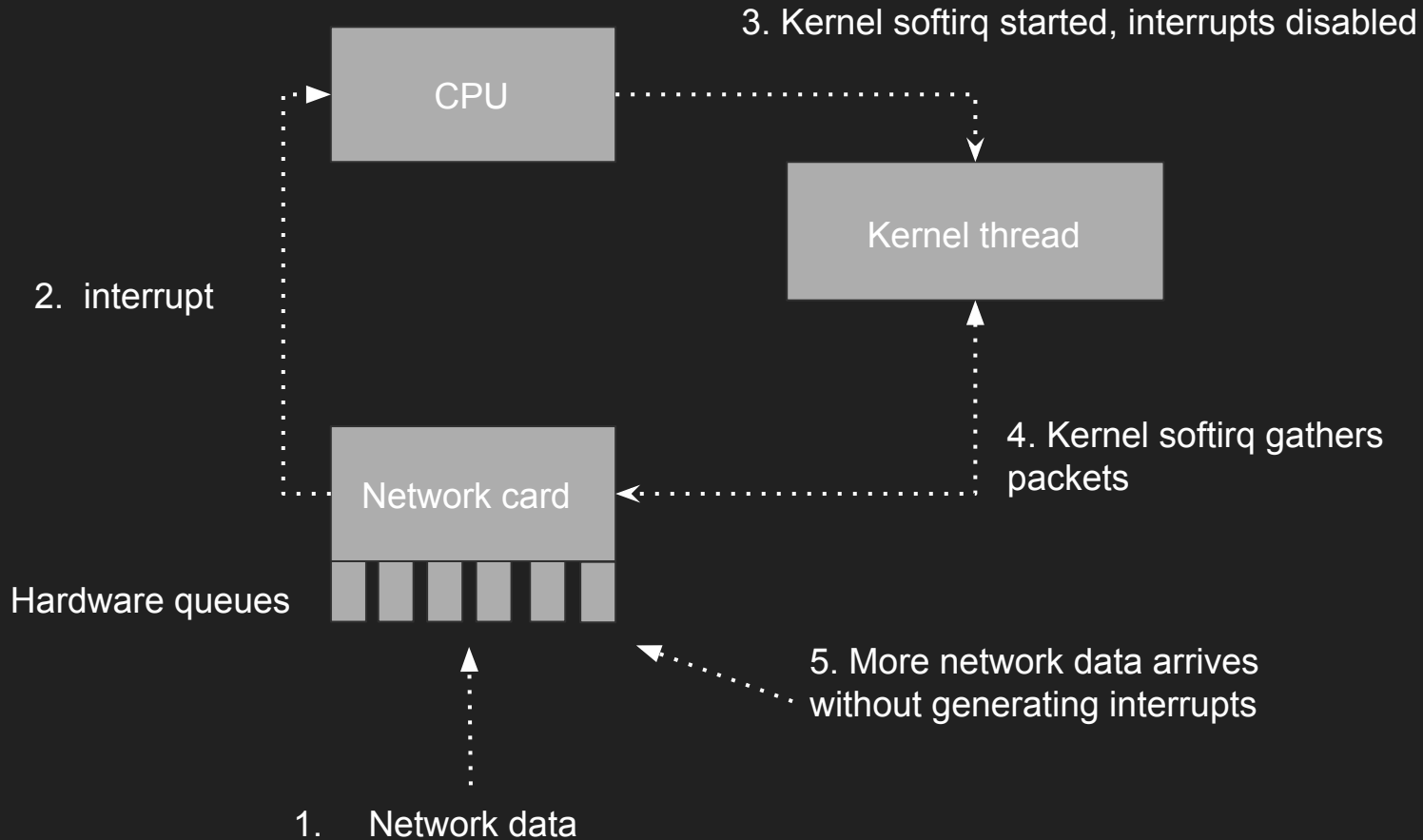
By default memcached assigns connections to worker threads in a round-robin manner. This patch introduces an option to select a worker thread based on the incoming connection's NAPI ID if SO_INCOMING_NAPI_ID socket option is supported by the OS.

This allows a memcached worker thread to be associated with a NIC HW receive queue and service all the connection requests received on a specific RX queue. This mapping between a memcached thread and a HW NIC queue streamlines the flow of data from the NIC to the application. In addition, an optimal path with reduced context switches is possible, if epoll based busy polling (sysctl -w net.core.busy_poll = <non-zero value>) is also enabled.

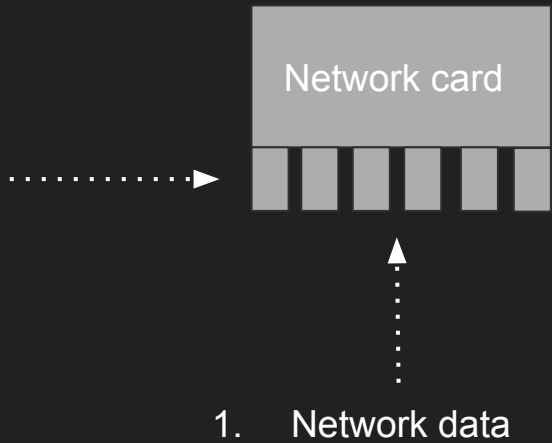
This feature is enabled via a new command line parameter `-N <num>` or `"--napi_ids=<num>"`, where `<num>` is the number of available/assigned NIC hardware RX queues through which the connections can be received. The number of `napi_ids` specified cannot be greater than the number of worker threads specified using `-t/--threads` option. If the option is not specified, or the conditions not met, the code defaults to round robin thread selection.

Signed-off-by: Kiran Patil <kiran.patil@intel.com>

Signed-off-by: Sridhar Samudrala <sridhar.samudrala@intel.com>



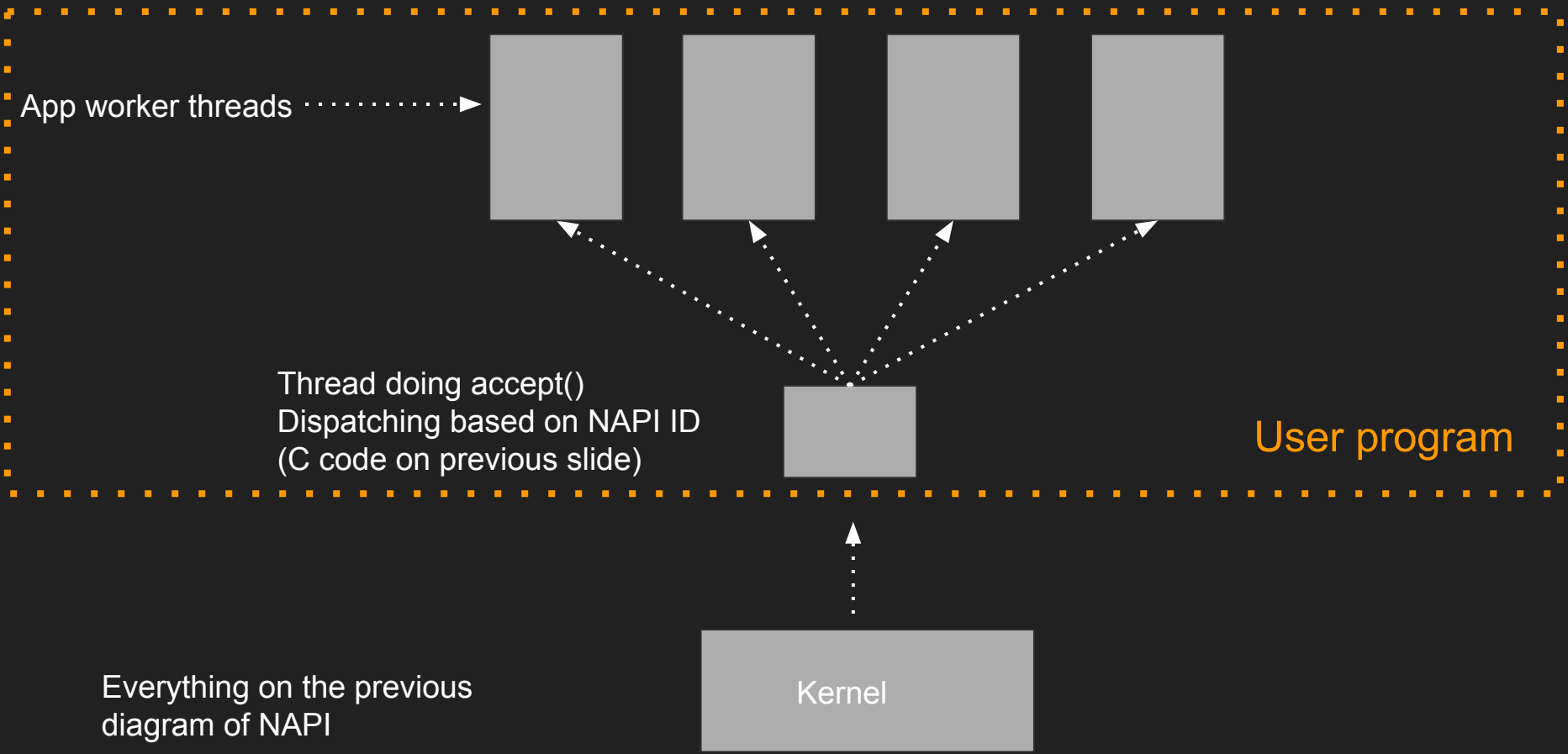
Hardware queues, each with their own system-wide unique ID (aka NAPI ID)

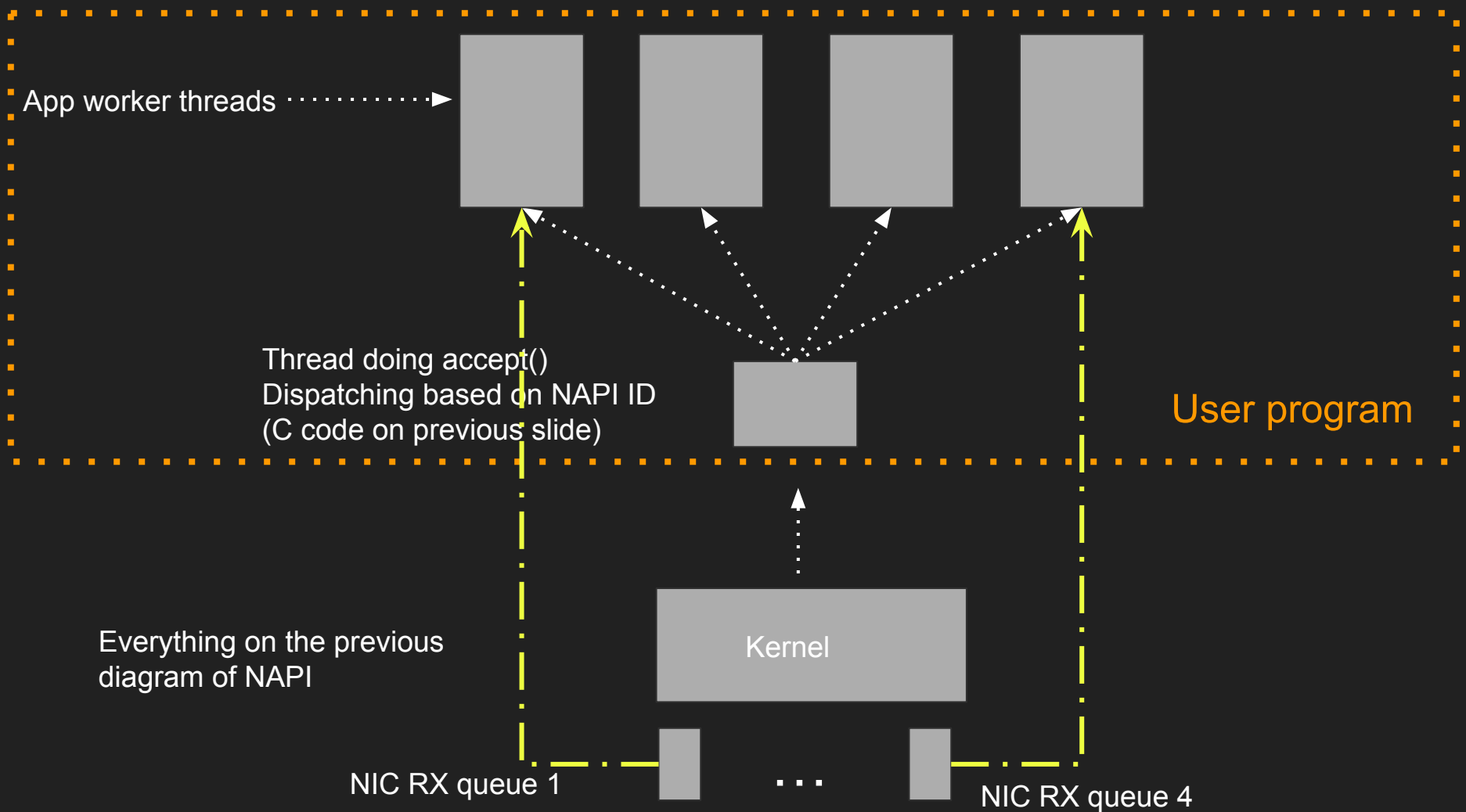


1. Network data

Example (pseudo code) usage

```
fd = accept(server_fd);  
getsockopt(fd, SOL_SOCKET, SO_INCOMING_NAPI_ID, &napi_id, &napi_id_len);  
  
/* now napi_id tells you which hardware queue this connection arrived on */  
dispatch_to_thread_by_napi_id(fd, napi_id);
```



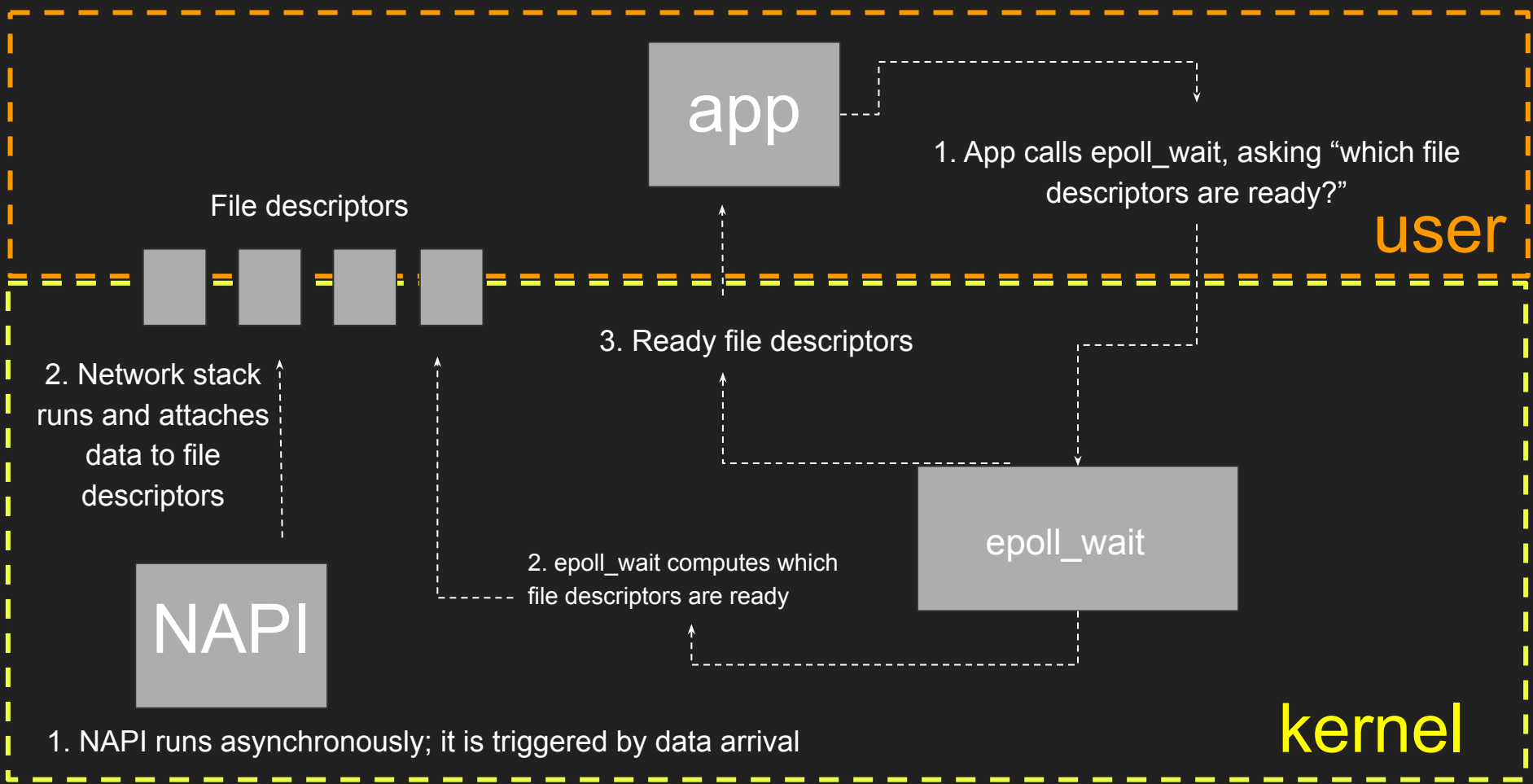


SO_INCOMING_NAPI_ID

- Provides hardware queue ID associated with connection
- Application can select worker thread based on queue ID
- Allows apps to map hardware NIC queues to threads
- Should improve cache hit rates and memory locality

But... there's more.

Before we can proceed, let's look at how
`epoll_wait` + NAPI work normally at a high level



app

File descriptors

1. App calls `epoll_wait`, asking "which file descriptors are ready?"

user

3. Ready file descriptors

2. Network stack runs and attaches data to file descriptors

NAPI

2. `epoll_wait` computes which file descriptors are ready

`epoll_wait`

1. NAPI runs asynchronously; it is triggered by data arrival

kernel

We'll assume this all happens on the same CPU

That's how “normal” `epoll_wait` works,
but...

An epoll specific modification was
added to the kernel for
SO_INCOMING_NAPI_ID

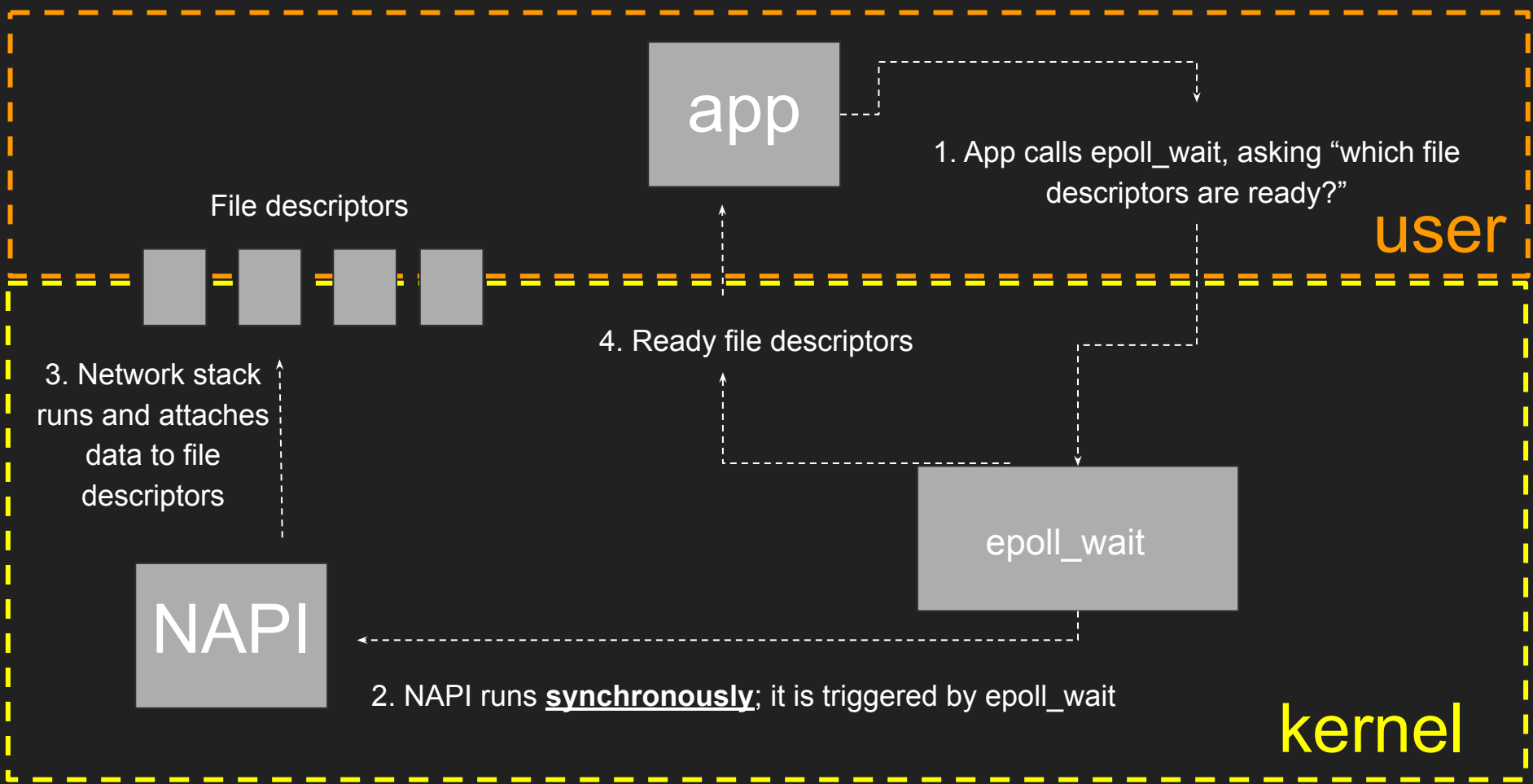
The `SO_INCOMING_NAPI_ID` of the last fd added to the epoll set can be busy polled.

To do so, you'd want to ensure all fds added are from the same RX queue

Set a *system-wide* sysctl or proc
value

`/proc/sys/net/core/busy_poll`

Kind of looks like this



This is all happening on the same CPU

epoll + SO_INCOMING_NAPI_ID

1. All FDs must have same NAPI_ID
2. epoll_wait drives packet processing
3. Probably better cache efficiency
4. The app can decide to do network processing vs its own work
5. *But this is limited by device IRQs which will still fire*

But it's system wide... so this makes
everything suddenly busy poll

So.... I wrote a patch to make it per
epoll-context specific via an ioctl

@ 2024-02-13 0:10 Joe Damato
2024-02-13 6:16 [PATCH net-next v8 1/4] eventpoll: support busy poll per epoll instance Joe Damato
(4 more replies)
0 siblings, 5 replies; 7+ messages in thread
From: Joe Damato @ 2024-02-13 6:16 UTC (permalink / raw)
To: linux-kernel, netdev
Cc: chuck.lever, jlayton, linux-api, brauner, edumazet, davem,
alexander.duyck, sridhar.samudrala, kuba, willemdebruijn.kernel,
weiwan, David.Laight, arnd, sdf, amritha.nambiar, Joe Damato,
Alexander Viro, Greg Kroah-Hartman, Helge Deller, Jan Kara,
Jiri Slaby, Jonathan Corbet, Julien Panis,
open list:DOCUMENTATION,
open list:FILESYSTEMS (VFS and infrastructure),
Michael Ellerman, Nathan Lynch, Palmer Dabbelt, Steve French,
Thomas Huth, Thomas Zimmermann

Greetings:

Welcome to v8.

TL;DR This builds on commit bf3b9f6372c4 ("epoll: Add busy poll support to epoll with socket fds.") by allowing user applications to enable epoll-based busy polling, set a busy poll packet budget, and enable or disable prefer busy poll on a per epoll context basis.

This makes epoll-based busy polling much more usable for user applications than the current system-wide sysctl and hardcoded budget.

To allow for this, two ioctls have been added for epoll contexts for getting and setting a new struct, struct epoll_params.

ioctl was chosen vs a new syscall after reviewing a suggestion by Willem de Bruijn [1]. I am open to using a new syscall instead of an ioctl, but it seemed that:

- Busy poll affects all existing epoll_wait and epoll_pwait variants in the same way, so new versions of many syscalls might be needed. It seems much simpler for users to use the correct epoll_wait/epoll_pwait for their app and add a call to ioctl to enable or disable busy_poll as needed. This also probably means less work to get an existing epoll app using busy poll.
- previously added epoll_pwait2 helped to bring epoll closer to existing syscalls (like pselect and ppoll) and this busy poll change

Threaded it through to:

- [glibc](#)
- [uclibc-ng](#)
- musl ([patch sent](#), waiting)

[git://sourceware.org](https://sourceware.org) / [glibc.git](https://sourceware.org/glibc.git) / commit

[summary](#) | [shortlog](#) | [log](#) | [commit](#) | [commitdiff](#) | [tree](#)
(parent: [400bdb5](#)) | [patch](#)

commit

Linux: Add epoll ioctls

```
author      Joe Damato <jdamato@fastly.com>
            Tue, 28 May 2024 17:37:06 +0000 (17:37 +0000)
committer   Noah Goldstein <goldstein.w.n@gmail.com>
            Tue, 4 Jun 2024 17:09:15 +0000 (12:09 -0500)
commit      92c270d32caf3f8d5a02b8e46c7ec5d9d0315158
tree        cd19d4bbfe7031daa99539773a5d18daa934c1db      tree
parent      400bdb5c85af5a52b3f5653357c9fca87f036bd3      commit | diff
```

Linux: Add epoll ioctls

As of Linux kernel 6.9, some ioctls and a parameters structure have been introduced which allow user programs to control whether a particular epoll context will busy poll.

Update the headers to include these for the convenience of user apps.

The ioctls were added in Linux kernel 6.9 commit [18e2bf0edf4dd](#) ("eventpoll: Add epoll ioctl for epoll_params") [1] to include/uapi/linux/eventpoll.h.

[1]: <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/diff/?h=v6.9&id=18e2bf0edf4dd>

Signed-off-by: Joe Damato <jdamato@fastly.com>

Reviewed-by: Adhemerval Zanella <adhemerval.zanella@linaro.org>

And, added a man page!

NAME

ioctl_eventpoll, EPIOCSPARAMS, EPIOCGPARAMS - ioctl() operations for epoll file descriptors

LIBRARY

Standard C library (libc, -lc)

SYNOPSIS

```
#include <sys/epoll.h> /* Definition of EPIOC* constants */
#include <sys/ioctl.h>
```

```
int ioctl(int fd, EPIOCSPARAMS, const struct epoll_params *argp);
int ioctl(int fd, EPIOCGPARAMS, struct epoll_params *argp);
```

```
#include <sys/epoll.h>
```

```
struct epoll_params {
    uint32_t busy_poll_usecs; /* Number of usecs to busy poll */
    uint16_t busy_poll_budget; /* Max packets per poll */
    uint8_t prefer_busy_poll; /* Boolean preference */

    /* pad the struct to a multiple of 64bits */
    uint8_t __pad; /* Must be zero */
};
```

Around this same time, an [academic paper](#)
[was published](#):

Kernel vs. User-Level Networking: Don't Throw Out the Stack with the Interrupts

PETER CAI, University of Waterloo, Canada

MARTIN KARSTEN, University of Waterloo, Canada

This paper reviews the performance characteristics of network stack processing for communication-heavy server applications. Recent literature often describes kernel-bypass and user-level networking as a silver bullet to attain substantial performance improvements, but without providing a comprehensive understanding of how exactly these improvements come about. We identify and quantify the direct and indirect costs of asynchronous hardware interrupt requests (IRQ) as a major source of overhead. While IRQs and their handling have a substantial impact on the effectiveness of the processor pipeline and thereby the overall processing efficiency, their overhead is difficult to measure directly when serving demanding workloads. This paper

5.2.3 *Kernel Polling*. Kernel polling does not suffer from the issues that IRQ packing and IRQ suppression face, because the decision whether to poll or enable interrupts is made automatically based on the application's workload. Its performance is strong in both maximum throughput and tail latency, as evident by Table 5 and Figure 5. While it is difficult to compare throughput numbers for specific tail latencies with this methodology, it is clear that kernel polling outperforms the vanilla configuration by at least 30%.

The paper was showing promising results, but would we see those in our production workloads?

Adding

`SO_INCOMING_NAPI_ID`

support was much trickier than I expected

At Fastly:

- We use an [open source web server](#) called h2o, which uses **SO_REUSEPORT**
- Our machines have 1 or 2 dual port NICs
- We run 1 process which listens on all NICs
- We listen on lots of ports (80, 443, ...)
- We use Mellanox ConnectX-5 Ex NICs

Rough diagram

h2o threads



Listening ports

All threads

Listen all ports

via

SO_REUSEPORT

80

80

80

80

Reuseport group 1

443

443

443

443

Reuseport group 2

NIC1

NIC2

NIC3

NIC4

Connections come in on any NIC and kernel balances amongst sockets in matching reuseport group

I know what to do!

1. Custom RSS contexts to steer flows to CPUs where webserver worker threads run
2. Push IRQs out (we are going to poll)
3. Reuseport bpf filter inserted into the webserver code
4. Enable busy poll just for the web server via the ioctl !

The NIC settings via ethtool look like this

(sorry if this is small, download the slides later to read it)

```
NICS='eth1 eth2 eth3 eth4'
```

```
for DEV in ${NICS}; do
  # 16 queues per device
  sudo ethtool -L $DEV combined 16

  # for queues in the mask 0xff, turn off adaptive IRQ algorithms; we plan to busy poll so push IRQs out as far
  # as hardware allows
  sudo ethtool --per-queue $DEV queue_mask 0xff --coalesce adaptive-rx off adaptive-tx off
  sudo ethtool --per-queue $DEV queue_mask 0xff --coalesce rx-usecs 4095 rx-frames 65535 tx-usecs 4095 tx-frames 65535

  # create a custom RSS context which sends all flows to queue 0-7, this is context 1
  sudo ethtool -X $DEV weight 1 1 1 1 1 1 1 1 context new

  # set the default RSS context to send all flows to queues 8-15
  sudo ethtool -X $DEV weight 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1

  # attach a filter rule to send all flows with a dst-port of $PORT to the queue in RSS context 1 (e.g. queues 0-7)
  sudo ethtool -U $DEV flow-type tcp4 dst-port 80 context 1
  sudo ethtool -U $DEV flow-type tcp4 dst-port 443 context 1

  # attach a filter rule to send all flows with a dst-port of $PORT to the queue in RSS context 1 (e.g. queues 0-7)
  sudo ethtool -U $DEV flow-type tcp6 dst-port 80 context 1
  sudo ethtool -U $DEV flow-type tcp6 dst-port 443 context 1

  # pin IRQs to CPUs which are NUMA local to the devices (via lstopo --whole-io)
  if [ "${DEV}" = "vlan101" ]; then
    sudo ./set_irq_affinity 9,73,4,68,24,88,28,92,3,10,19,26,34,40,48,56 $DEV
  elif [ "${DEV}" = "vlan201" ]; then
    sudo ./set_irq_affinity 12,76,32,96,36,100,42,106,67,74,83,90,98,104,112,120 $DEV
  elif [ "${DEV}" = "vlan301" ]; then
    sudo ./set_irq_affinity 16,80,44,108,50,114,54,118,7,15,20,31,39,45,55,63 $DEV
  else
    sudo ./set_irq_affinity 22,86,58,122,62,126,64,82,71,79,84,95,103,109,119,127 $DEV
  fi

  # setup transmit steering based on receive-queues map
  # /sys/class/net/<dev>/queues/tx-<n>/xps_rxqs
  sudo ./set_xps_rxqs $DEV
done
```

I know what to do!

- ~~1. Custom RSS contexts to steer flows to CPUs where webserver worker threads run~~
- ~~2. Push IRQs out (we are going to poll)~~
3. Reuseport bpf filter inserted into the webserver code
4. Enable busy poll just for the web server via the ioctl !

```
struct sock_filter code[] = {
    /* A = skb->queue_mapping */
    { BPF_LD | BPF_W | BPF_ABS, 0, 0, SKF_AD_OFF + SKF_AD_QUEUE },
    /* A = A % n */
    { BPF_ALU | BPF_MOD, 0, 0, n },
    /* return A */
    { BPF_RET | BPF_A, 0, 0, 0 },
};

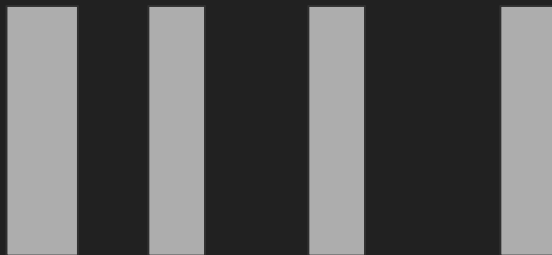
struct sock_fprog p = {
    .len = ARRAY_SIZE(code),
    .filter = code,
};

if (setsockopt(fd, SOL_SOCKET, SO_ATTACH_REUSEPORT_CBPF, &p, sizeof(p))) {
```

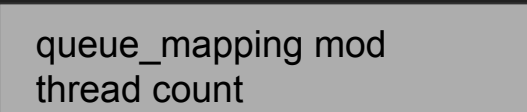

This actually doesn't work when you have multiple NICs.

Buggy example

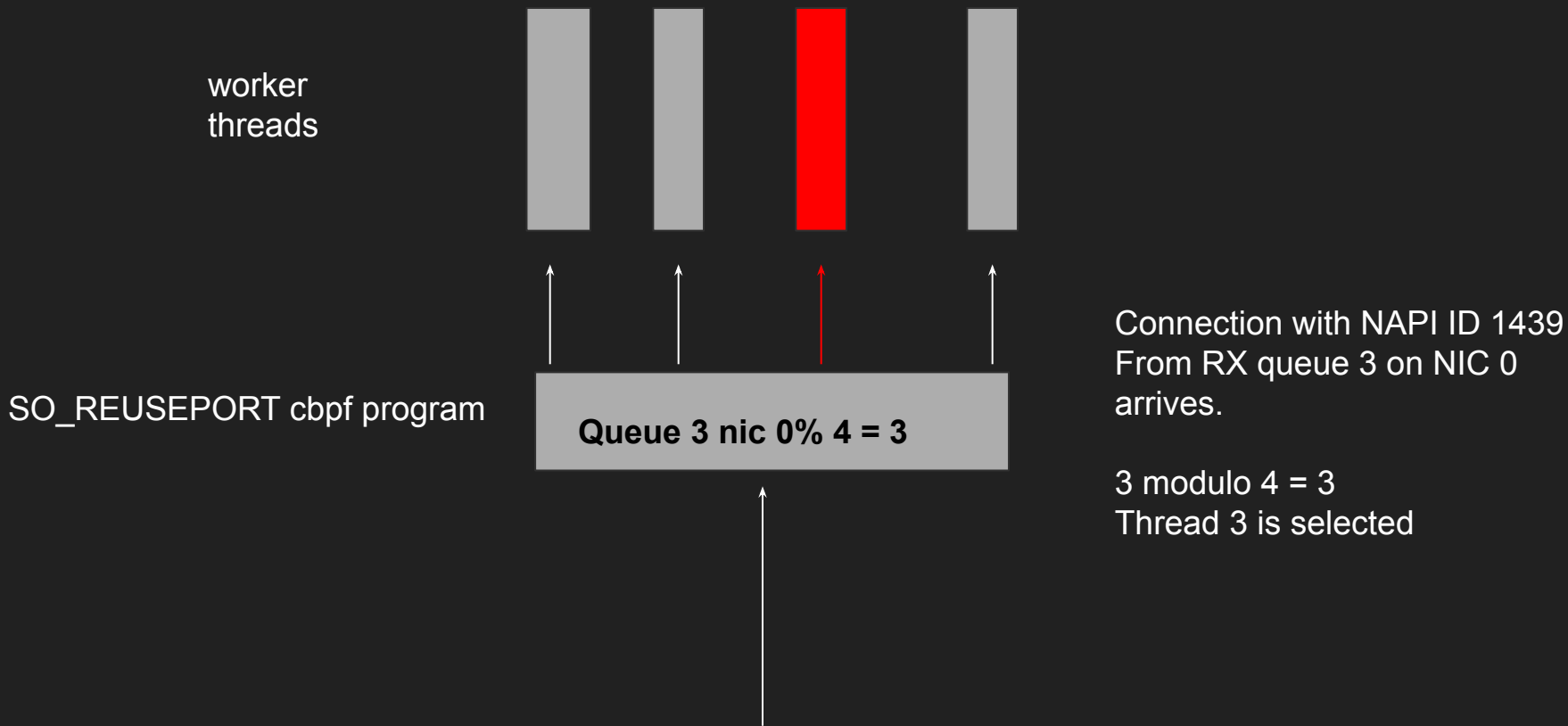
worker threads



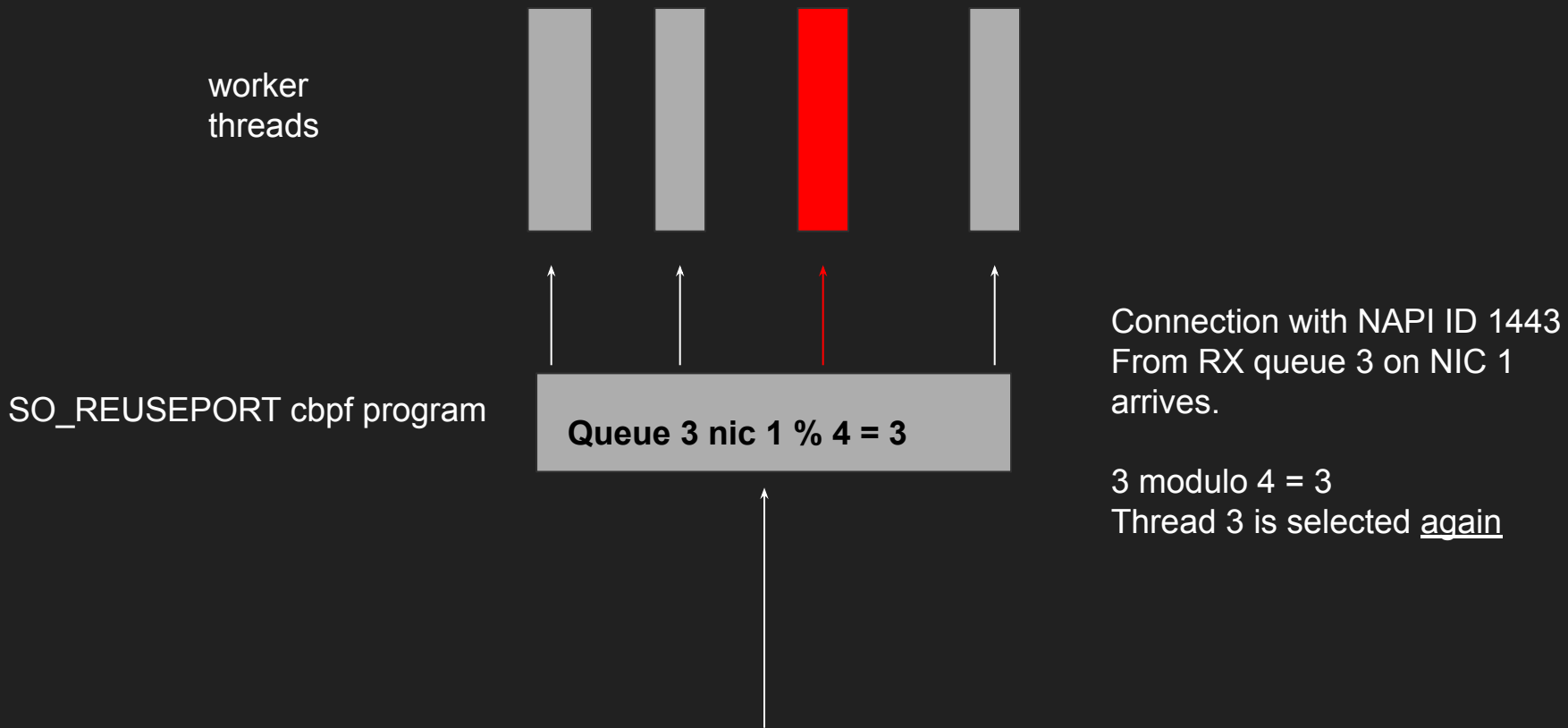
SO_REUSEPORT cbpf program
Installed



Buggy example part 1



Buggy example part 2



Now thread 3 has been given two connections with different NAPI IDs:

This breaks busy poll and is not allowed.

So... how to
fix this?

A few ways come
to mind, but this is
what I did

Read a lot of kernel code


```
static int inet_reuseport_add_sock(struct sock *sk,
                                   struct inet_listen_hashbucket *ilb)
{
    struct inet_bind_bucket *tb = inet_csk(sk)->icsk_bind_hash;
    const struct hlist_nulls_node *node;
    struct sock *sk2;
    kuid_t uid = sock_i_uid(sk);

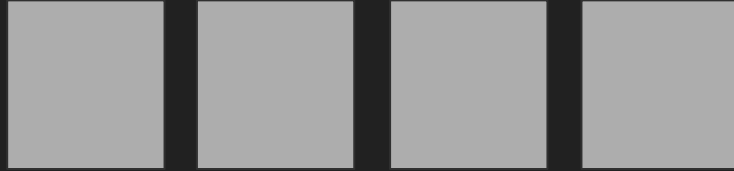
    sk_nulls_for_each_rcu(sk2, node, &ilb->>nulls_head) {
        if (sk2 != sk &&
            sk2->sk_family == sk->sk_family &&
            ipv6_only_sock(sk2) == ipv6_only_sock(sk) &&
            sk2->sk_bound_dev_if == sk->sk_bound_dev_if &&
            inet_csk(sk2)->icsk_bind_hash == tb &&
            sk2->sk_reuseport && uid_eq(uid, sock_i_uid(sk2)) &&
            inet_rcv_saddr_equal(sk, sk2, false))
            return reuseport_add_sock(sk, sk2,
                                       inet_rcv_saddr_any(sk));
    }

    return reuseport_alloc(sk, inet_rcv_saddr_any(sk));
}
```

What if... listen sockets
were grouped per NIC with
`SO_BINDTODEVICE`?

Rough diagram

h2o threads



Listening ports

All threads

Listen all ports

via

SO_REUSEPORT

80

80

80

80

Reuseport group 1

443

443

443

443

Reuseport group 2

NIC1

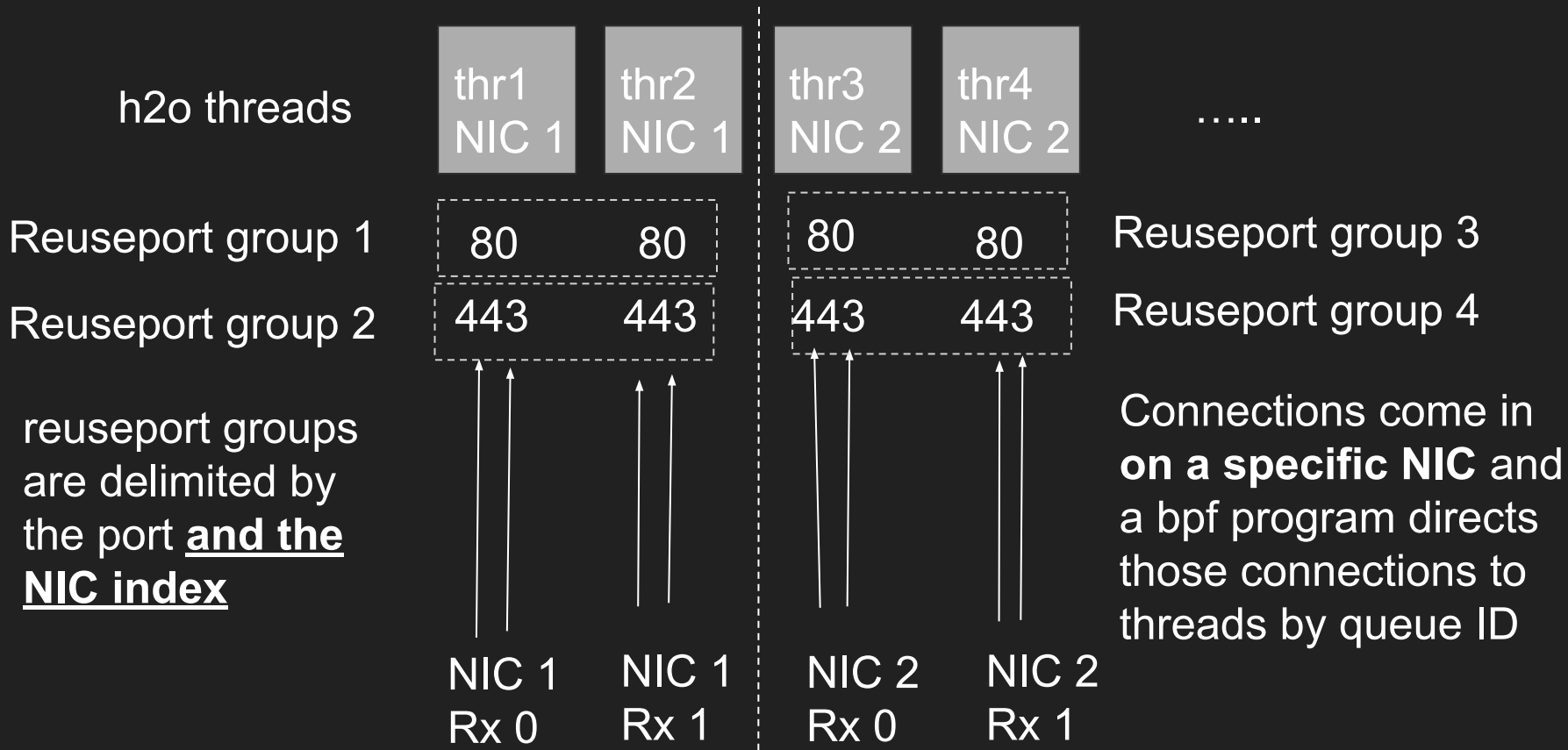
NIC2

NIC3

NIC4

Connections come in on any NIC and kernel balances amongst sockets in matching reuseport group

Busy poll h2o



The code to get there is reasonable

Rough pseudo code of
how I approached this:

```
char *ifaces[2] = ["eth0", "eth1"];
struct iface_listen_sockets listen_sockets[2];
char *iface;
int thread_count = 4;
int reuse = 1;
int *fd;

/* create 1 listen socket per worker thread per iface */
for (int iff = 0; iff < 2; iff++) {
    iface = ifaces[iff];
    listen_sockets = listen_sockets[iff];

    for (int i = 0; i < thread_count; i++) {
        /* create the socket */
        listen_sockets->fd[i] = socket(domain, type, protocol);

        /* bind it to the device */
        setsockopt(listen_sockets->fd[i], SOL_SOCKET, SO_BINDTODEVICE, iface, strlen(iface));

        /* set reuseport */
        setsockopt(listen_sockets->fd[i], SOL_SOCKET, SO_REUSEPORT, &reuse, sizeof(reuse));

        /* apply the BPF filter to the first socket in the group */
        if (i == 0) {
            apply_bpf_filter(fd, thread_count);
        }
    }
}
```

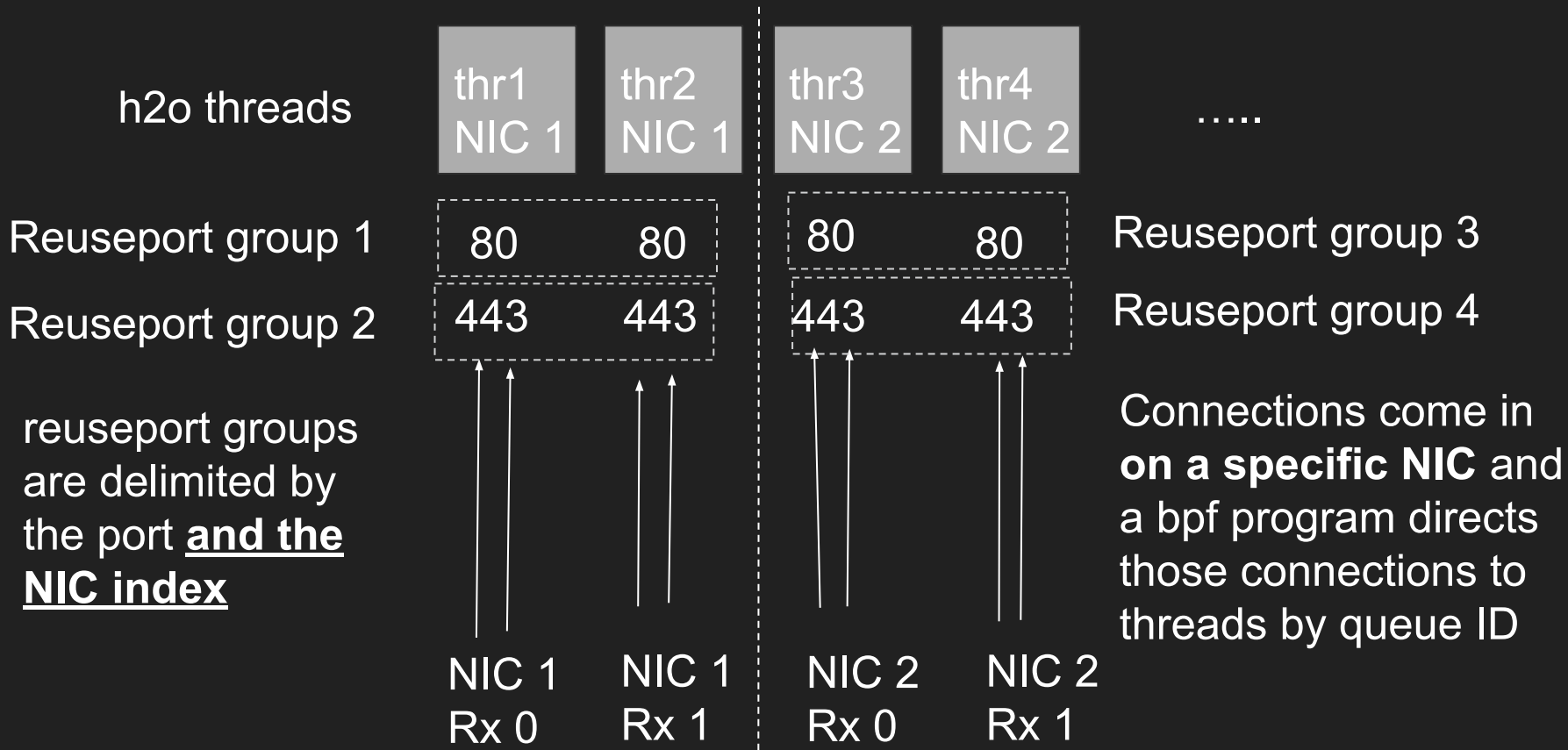
```
/* now that the reuseport groups are created, call listen on all
 * the sockets
 */
for (int iff = 0; iff < 2; iff++) {
    listen_sockets = listen_sockets[iff];

    for (int i = 0; i < thread_count; i++) {
        listen(listen_sockets->fd[i], 1024);
    }
}
```


Turned out this worked!

So now, incoming connections are going to worker threads while maintaining the constraint that each thread maps to one NAPI ID.

Busy poll h2o



I think it's important to mention
what the accept path looks like.

This is what I did, but there many
other ways, too!

App config read at startup maps NICs to NUMA local CPUs

busy-poll-cpu-map:

eth1:

- 9

- 10

- 11

-

eth2:

- 18

-

So, accept path looks like this
(psuedo code):

```
static __thread int cpu_claimed = 0;
static __thread int napi_id_claimed = 0;

static void handle_busy_poll_accept(int sockfd)
{
    unsigned napi_id = get_napi_id(sockfd);

    /* has this thread claimed a CPU? */
    if (cpu_claimed) {
        /* it has... does the claimed NAPI ID match the new
         * connection?
         */
        if (napi_id_claimed == napi_id)
            return;

        /* here we have a mismatch. this is what happens if you
         * have multiple NICs and you dont setup the reuseport
         * groups in the kernel correctly
         */
        if (napi_id_claimed != napi_id)
            fire_an_alert();
    } else {
        struct ifreq ifr = {};

        /* use a hack to get the NIC name from NAPI ID */
        get_nic_name_by_napi(sockfd, napi_id, &ifr);

        /* pick an un-used CPU from the list provided in config */
        int cpu_id = assign_cpu_from_map(ifr.ifr_name);

        /* thread claims that CPU and NAPI ID */
        cpu_claimed = cpu_id;
        napi_id_claimed = napi_id;

        /* pin the thread to that CPU*/
        pin_thread_to_cpu(cpu_id);
    }
}
```

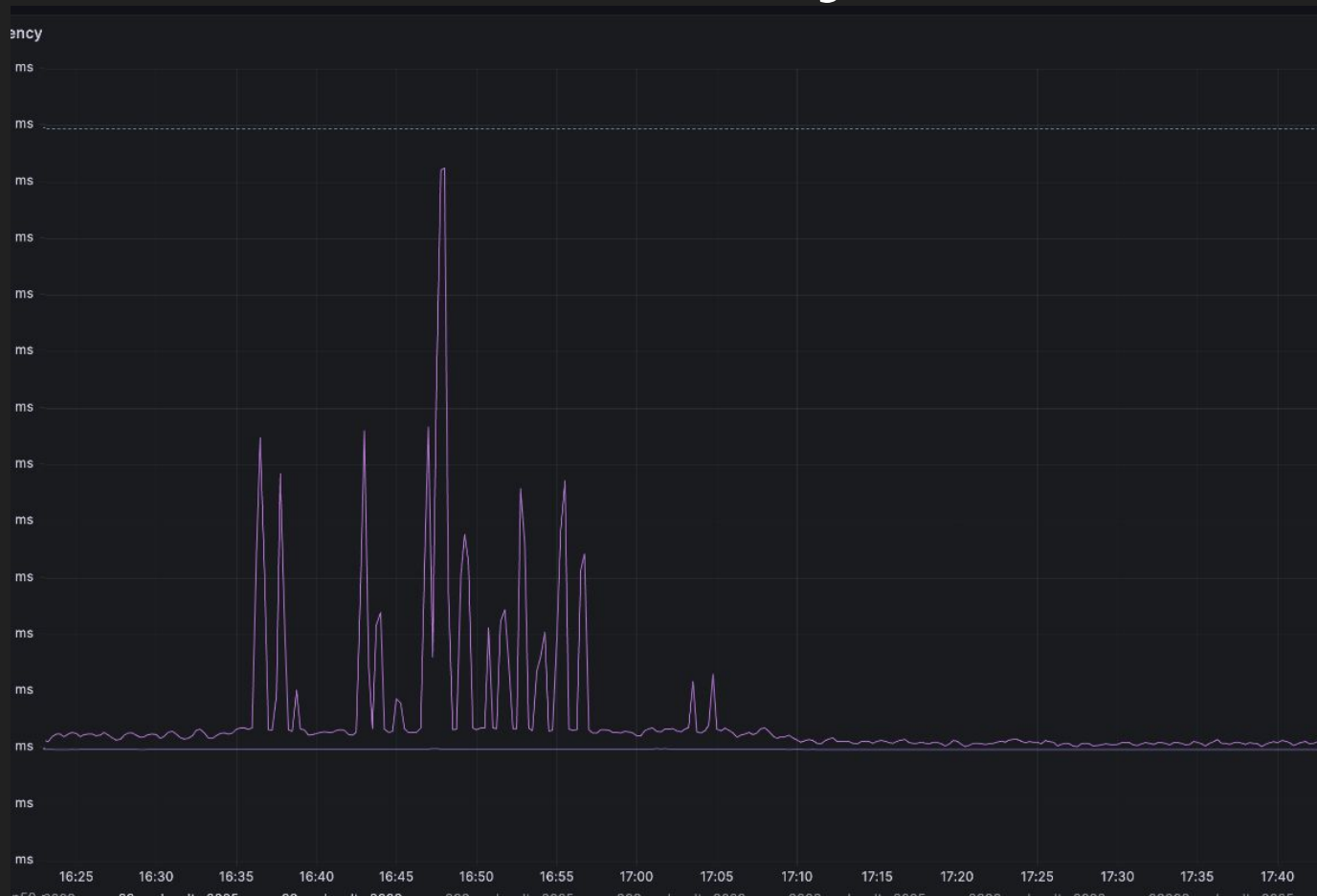
Internal hack that we use:

```
/* this is an internal kernel hack; works like SIOCGIFNAME but given a NAPI
 * ID (instead of an ifindex) fills in the interface name.
 *
 * instead, drivers should implement netdev-genl to map NAPI IDs to
 * ifindexes
 *
 * then user apps can use that real interface instead of this hack
 */
static int get_nic_name_by_napi(int fd, unsigned int napi_id, struct ifreq *ifr)
{
    ifr->ifr_ifru.ifru_ival = napi_id;

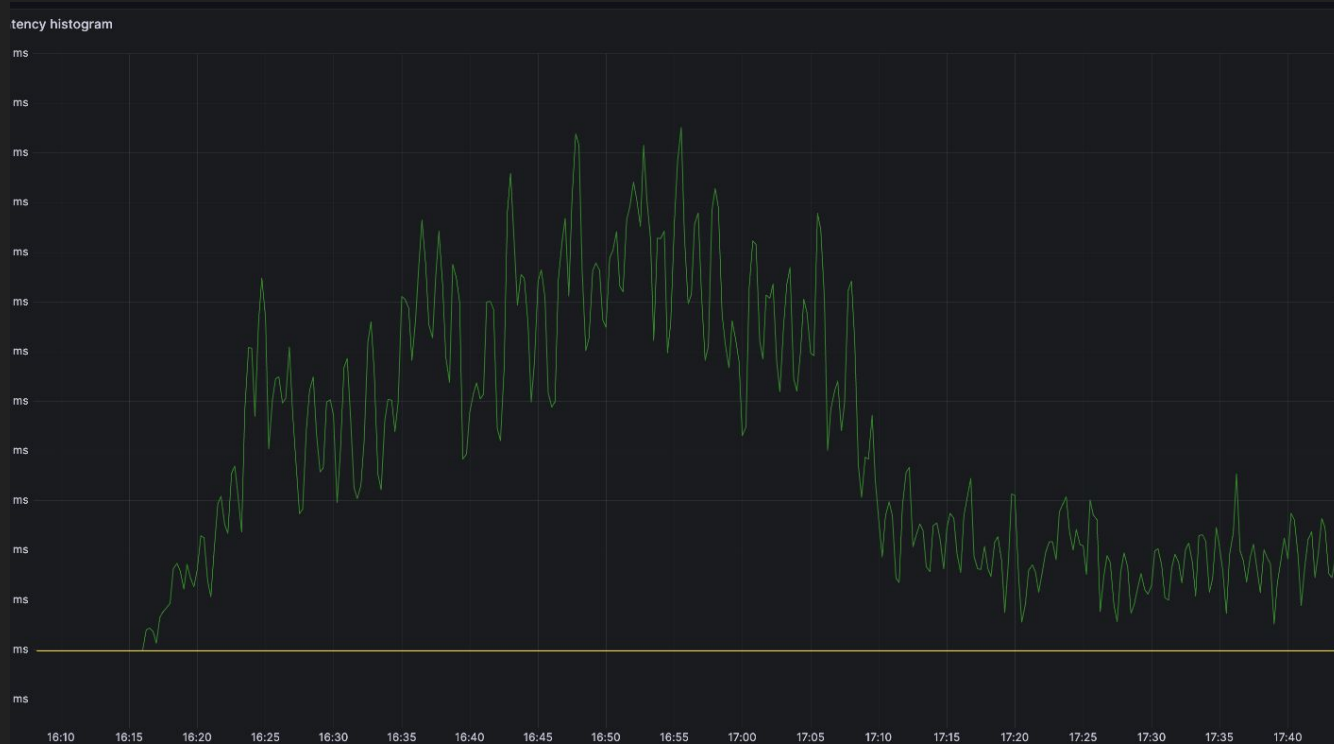
    return ioctl(fd, SIOCGIFNAME_BY_NAPI_ID, ifr);
}
```


OK So we did all that,
what does it look like in
prod?

P99 latency better



P999 latency better



P9999 latency better

tency histogram

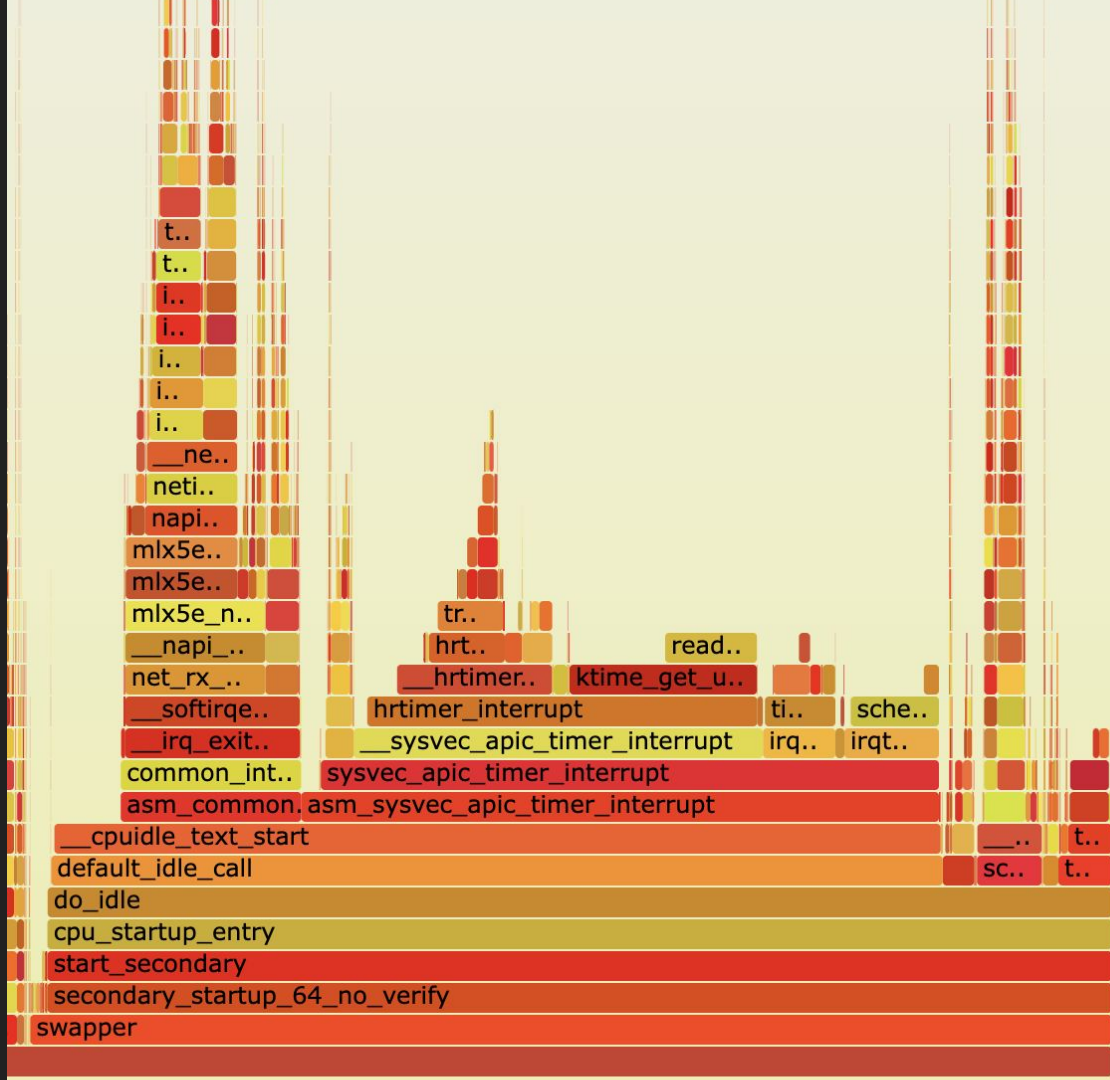


Cool, so latency looks much better for the higher percentiles.

BUT

1. we use `epoll_pwait2`
2. we set the busy poll usecs to match the `pwait2` time
3. Occasionally, `pwait2` runs idle

perf measurement suggests
IRQs from the device still arrive
which kick off NAPI and induce
latency



Can we maybe set maxevents
lower and defer IRQs with:

```
napi_defer_hard_irqs  
gro_flush_timeout
```

We can't.

These settings are system-wide. Other latency sensitive apps that are not busy-poll compatible will suffer.

Turns out choosing these two values is really, really hard anyway.

That's where the future comes in:

Some things you may want to work on

Some things I and Fastly am/are working on

You (driver maintainers) can:

- Add netdev-genl support to your driver
- Already supported by:
 - mlx4
 - mlx5
 - ice
 - bnxt

I am currently working on:

A collaboration with Martin Karsten of U Waterloo:

- A kernel patch which:
 - disables interrupts during busy poll **but**
 - re-enables IRQs if busy poll has no data
 - Avoids needing to pick “the right numbers” for `napi_defer_hard_irqs` and `gro_flush_timeout` which might change if the system is under load or idle
 - `gro_flush_timeout` serves as a “backstop” if userland is slow
- Initial data measured against memcached looks *very promising*.
- Enabled via a busy poll ioctl per context

Hopefully
submitting an RFC
with data soon.

Fastly is working on

Upstreaming our busy poll implementation to [open source h2o](#).

Not ready yet, but we are working on it and planning to open source it.

I am *hoping* to work on:

napi_defer_hard_irqs

gro_flush_timeout

per NAPI

(via netdev-genl hopefully?)

