

Running Open vSwitch on a large CPU system

OVS/OVN Conference 2024

Eelco Chaudron Senior Principal Software Engineer



- Objective;
 - Evaluate how high CPU count impacts revalidation and ovs-vswitchd functionality under normal operating conditions.
- Focus Areas;
 - Revalidation process
 - ☐ Interface statistics impact on run time
 - Simulated POD creation
 - ☐ getifaddrs() run time



System Under Test (SUT)

- Hardware:
 - □ Dell PowerEdge R6515
 - AMD EPYC 7702P 64-Core/128-Threads Processor
 - Memory 16GB DDR4
 - ☐ Storage 446.63 GB SATA/SSD
- Software:
 - □ RHEL 9.5
 - Open vSwitch 3.3
- Number of CPUs:
 - □ 128 CPUs; all threads are enabled
 - □ 64 CPUs; the number of cores per CCD has been reduced to FOUR(2+2) in the BIOS.



The Revalidation Process



- Revalidation is managed by the revalidator threads
 - Making sure the installed datapath flows match the configured OpenFlow rules
 - Updating OpenFlow packet and byte statistics
 - ☐ Removing idle datapath flows
- More details in the following blog:

Open vSwitch: The revalidator process explained

- By default, if *n-revalidator-threads* is not configured, for a 128 CPU system, 33 revalidator threads are started.
 - ☐ The number of CPUs available to ovs-vswitchd divided by four, plus one.
- We will compare data with 4, 32 and 128 revalidator threads.



Open vSwitch configuration

- Simple Open vSwitch configuration:
 - ☐ Two Network namespace simulating endpoints
 - veth interface to connect endpoints to OVS
 - □ Scapy for generating the traffic
 - ☐ Maximum datapath flows to exists (200k)
- Using taskset to bind the Scapy script to a specific CPU (set).

```
# ip netns exec ns_h1 bash
```

taskset -c 1 python traffic.py



NETNS configuration details

```
ip link add h1_ovs type veth peer name h1_eth
ip link add h2 ovs type veth peer name h2 eth
ip link set dev h1_ovs up
ip link set dev h2 ovs up
ip netns add ns_h1
ip netns add ns h2
ip link set h1_eth netns ns_h1
ip link set h2_eth netns ns_h2
ip -n ns_h1 link set dev lo up
ip -n ns h1 link set dev h1 eth up
ip -n ns_h1 address add 192.168.9.1/24 dev h1_eth
ip -n ns h2 link set dev lo up
ip -n ns_h2 link set dev h2_eth up
ip -n ns_h2 address add 192.168.9.2/24 dev h2_eth
```



```
ovs-vsctl add-br ovs pvp br0
ovs-vsctl add-port ovs pvp br0 h1 ovs
ovs-vsctl add-port ovs pvp br0 h2 ovs
ovs-appctl upcall/set-flow-limit 200000
ovs-ofctl del-flows ovs pvp br0
ovs-ofctl add-flow ovs pvp br0 in port=h1 ovs,arp,action=h2 ovs
ovs-ofctl add-flow ovs pvp br0 in port=h1 ovs,ip,nw dst=192.168.9.2,action=h2 ovs
ovs-ofctl add-flow ovs pvp br0 in port=h2 ovs,arp,action=h1 ovs
ovs-ofctl add-flow ovs pvp br0 in port=h2 ovs,ip,nw dst=192.168.9.1,action=h1 ovs
ip netns exec ns_h1 ping 192.168.9.2
python -c 'for i in range(0, 200000): \
  print("add in port=h1 ovs,ip,nw dst=10.{}.{}.{}/32,action=h2 ovs".
    format((i >>16) & 0xff, (i >> 8) & 0xff, i & 0xff))' | \
  ovs-ofctl add-flow ovs pvp br0 -
```



```
from scapy.all import *
 def ip from int(decimal):
     ip int = 0x0a000000 + decimal;
     return "{}.{}.{}.".format(ip_int >> 24 & 0xff,
                                 ip_int >> 16 & 0xff,
                                 ip int >> 8 & 0xff,
                                 ip int & 0xff);
  pkt list = []
 eth = Ether(dst="36:c2:88:9e:19:b1", src="00:00:00:00:00:01")
 ip = IP(dst="10.0.0.1", src="10.200.0.1")
 udp = UDP(dport=69)
  print("Building packet list...")
 for i in range(0, 100000):
     ip.dst=ip_from_int(i)
     pkt list.append(eth / ip / udp)
     if (i % 1000) == 0 and i != 0:
          print("Frame: {}".format(i))
  print("Sending packets...")
 sendpfast(pkt_list, iface="h1_eth", loop=0)
```



- The following tests where ran with 4, 32, and 128 revalidator threads:
 - 200k datapath flows, with traffic generated from a single CPU.
 - \square 200k datapath flows + 10x a second a flow revalidation sweep.
 - 200k datapath flows, with traffic generated by all CPU cores.
- Existing USDT probes and the <u>reval monitor.py</u> script is used for the initial analysis.



200k datapath flows with traffic arriving on a single CPU



Results with 4 revalidator threads:

```
=> dump duration:
# NumSamples = 641; Min = 158.00; Max = 185.00
# Mean = 169.695788; Variance = 26.277190; SD = 5.126128; Median 170.000000
# each I represents a count of 1
 158.0000 - 160.7000 [
                 163.4000 [
                 160.7000 -
        166.1000 [
                 163.4000 -
 166.1000 -
        168.8000 [
        171.5000 [
 168.8000 -
 171.5000 -
        174.2000 [
        176.9000 [
                 174.2000 -
 176.9000 -
        179.6000 [
                 10]: |||||||||||
        182.3000 [
 179.6000 -
 182.3000 -
                 3]: [[[
        185.0000 [
```



Results with 32 revalidator threads:

```
=> dump duration:
# NumSamples = 626; Min = 171.00; Max = 196.00
# Mean = 180.718850; Variance = 25.441721; SD = 5.043979; Median 180.000000
# each I represents a count of 1
181.0000 - 183.5000 [
        183.5000 - 186.0000 [
        186.0000 - 188.5000 [
188.5000 - 191.0000 [
        9]: |||||||
191.0000 - 193.5000
         3]: [[[
193.5000 - 196.0000 [
```



Results with 128 revalidator threads:

```
=> dump_duration:
# NumSamples = 628; Min = 169.00; Max = 195.00
# Mean = 183.399682; Variance = 16.771783; SD = 4.095337; Median 184.000000
169.0000 - 171.6000 [ 3]: ▮
171.6000 - 174.2000 [ 14]: ||||||
174.2000 - 176.8000 [ 29]: ||||||||||
187.2000 - 189.8000 [
            189.8000 - 192.4000 [
            22]: |||||||||||
            2]: [
192.4000 - 195.0000 [
```



- So more threads, show slight increase in run time.
 - □ 169ms, 181ms, 183ms.
- ☐ Further analysis was done using the <u>kernel_delay.py</u> and <u>analyze_perf_pmd_syscall.py</u> scripts.
- Conclusion; it's due to the a lock contention:

```
#0 ovs_mutex_lock_at() at ./ovs/lib/ovs-thread.c:76
#1 nl_dump_next() at ./ovs/lib/netlink-socket.c:1257
#2 dpif_netlink_flow_dump_next() at ./ovs/lib/dpif-netlink.c:1957
#3 dpif_flow_dump_next() at ./ovs/lib/dpif.c:1145
```



Further research should be done to see if we can optimize nl_dump_next() code:

```
bool
nl dump next(struct nl dump *dump, struct ofpbuf *reply, struct ofpbuf *buffer)
    . . .
    if (!buffer->size) {
        ovs_mutex_lock(&dump->mutex);
       if (!dump->status) {
            /* Take the mutex here to avoid an in-kernel race. If two threads
             * try to read from a Netlink dump socket at once, then the socket
             * error can be set to EINVAL, which will be encountered on the
             * next recv on that socket, which could be anywhere due to the way
             * that we pool Netlink sockets. Serializing the recv calls avoids
             * the issue. */
            dump->status = nl dump refill(dump, buffer);
        retval = dump->status;
        ovs_mutex_unlock(&dump->mutex);
```



200k datapath flows + revalidation sweep



- Runtime with 4, 32, and 128 revalidator threads:
 - □ 334ms, 209ms, 224ms
- More revalidation threads will speedup the actual revalidation process.
- However, to many will hit the same nl_dump_next() lock contention.



200k datapath flows with traffic arriving on multiple CPUs



	Wh	Why this test?			
		When experimenting we found that the OOM killer kicked in.			
		Memory was depleted due to massive <code>sw_flow_stat</code> allocations.			
	Wh	at are these <i>sw_flow_stat</i> structures?			
		They hold per-CPU datapath flow statistics			
		32-byte in size			
		Allocated on demand, when traffic is processed on the CPU.			
٥	Ηον	w much memory they consume for 200k datapath flows?			
		On a 128-CPU system: 782 MB			
		On a 256-CPU system: 1.5 GB			
	ls tl	he amount of memory the (main) problem?			
		No			

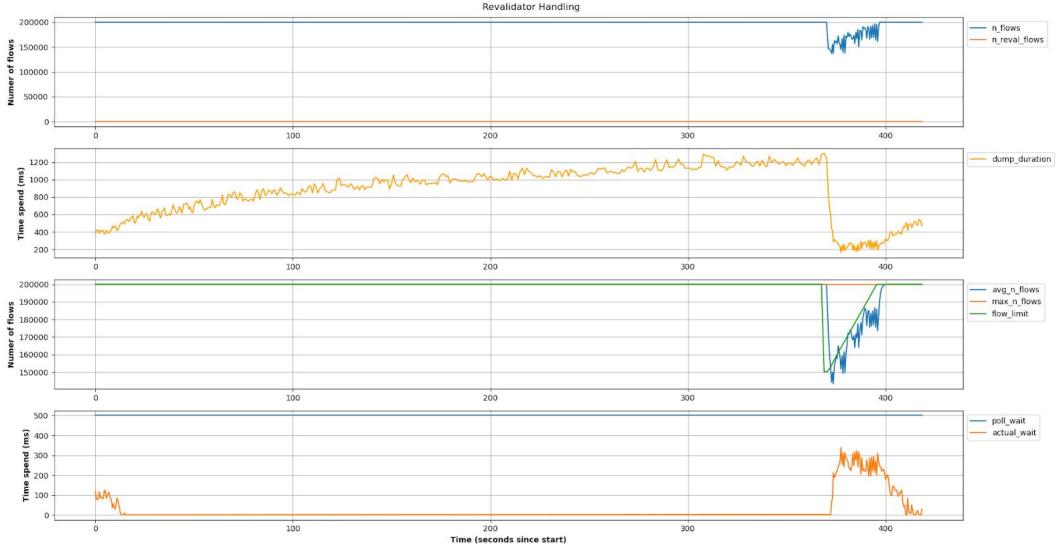


200k datapath flows / Rx on multiple CPUs

	 The main problem; Time spent accumulating data due to non-linear memory storage. Increases latency in processing data, i.e. getting flow statistics from the kernel to userspace. 			
٥	The increased latency when adding more <code>sw_flow_stats</code> structures leads to gradual increase in the <code>dump_duration</code> .			
٥	This continues till we hit a threshold (1.3 second) at which the maximum datapath flows are reduced.			
	Note that this is nothing new, it also exists in a 64-CPU system.			
۵	The next slides will show some graphs while running these tests Interesting observation; it takes longer to reach the threshold with a lower number of revalidator threads. I did not investigate why.			

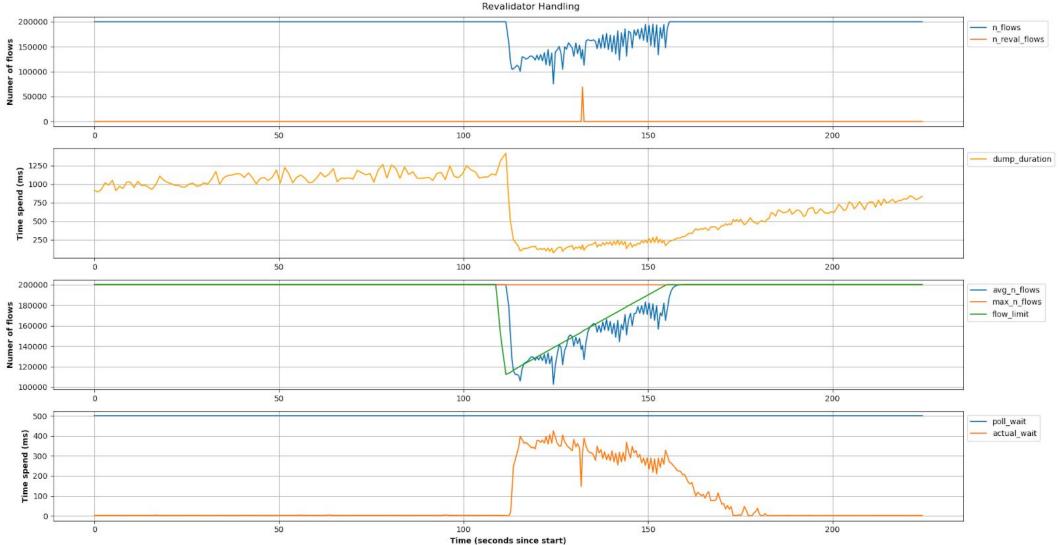


4 revalidator threads:



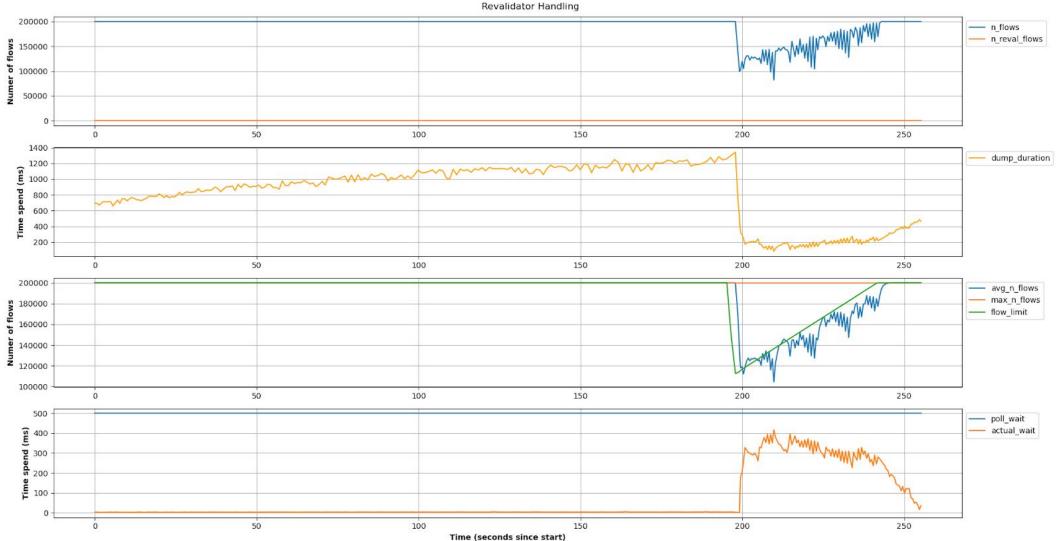


□ 32 revalidator threads:





□ 128 revalidator threads:





Kernel Interface Statistics Impact



- ovs-vswitchd periodically gathers all interface statistics (± every 5 seconds).
- Different kernel net drivers utilize different per-CPU statistic counters
 - We use veth drivers, which have specific per-CPU counters.
- Per-CPU netdev statistics are linear in memory.
- The Linux kernel also maintains certain exception counters in the net device core stats.



For measuring the impact:

- Configure 1024 veth pairs; one side in its own netns, the other side in OVS.
- Add two USDT probes to ovs-vswitchd (see next slide) measuring the time it takes to update the interface statistics.
- Use a modified version of the <u>bride loop.bt</u> script.
- Run for five minutes on a 64-core and 128-core configuration.
- One run without the exception stats being present, one with them being present.



Additional USDT probes:

```
--- a/vswitchd/bridge.c
+++ b/vswitchd/bridge.c
@@ -48,6 +48,7 @@
#include "openvswitch/meta-flow.h"
#include "openvswitch/ofp-print.h"
#include "openvswitch/ofpbuf.h"
+#include "openvswitch/usdt-probes.h"
#include "openvswitch/vconn.h"
#include "openvswitch/vlog.h"
#include "ovs-lldp.h"
@@ -3131,6 +3132,7 @@ run stats update(void)
             struct bridge *br;
             stats txn = ovsdb idl txn create(idl);
            OVS_USDT_PROBE(run_stat_update, br_loop_start);
             HMAP_FOR_EACH (br, node, &all_bridges) {
                 struct port *port;
                 struct mirror *m;
@@ -3148,6 +3150,7 @@ run_stats_update(void)
                    mirror refresh stats(m);
             OVS USDT PROBE(run stat update, br loop done);
             refresh_controller_status();
```



□ 64-CPU result:

128-CPU result:

- Doubling the number of CPU result in slightly less than doubling the run time (40%).
- This is outside the control of OVS as these per-CPU counters are a driver specific implementation.



Comparing with and without exception counters on a 128-CPU system

```
Average time: 0.104204051 seconds (without)
@statistics_loop_time in micro seconds:
[90000, 100000)
                                                                                                                                       oxed{1} , oxed{1}
[100000, 110000)
[110000, 120000)
                                                                                                                                       \mid @@@@@@@@@@@@@@@@@@@@@@@@@@@@@
Average time: 0.104665879 seconds (with exception counters)
@poll_block_wait_time:
[90000, 100000)
                                                                                                                                       oxedge කෙනෙක් කෙන්න කෙනෙක් කෙනෙන කෙනෙක් කෙනෙන කෙනෙක් කෙනෙන කෙනෙන කෙනෙන කෙනෙන කෙනෙන කෙනෙනි කෙනෙනි කෙනෙන කෙනෙන කෙනෙන කෙනෙන කෙනෙන කෙනෙන කෙනෙන කෙනෙන කොන කොනෙන කෙනෙන කොන කොනෙන කෙනෙන කෙනෙන කොනෙන කොනෙන කොනෙන කොනෙන කොන කොනෙන කොනෙන කොන කොනෙන කොනෙන කෙනෙන කොනෙන කොනෙන කොන කොනෙන කොනෙන කොනෙන කොනෙන කොනෙන කොනෙන කොනෙන කෙනෙන කොනෙන කෙනෙන කොනෙන කොන කොනෙන කොනෙන කොනෙන කොනෙන කොන කොනෙන කොනෙන කොනෙන කොනෙන කොන කොනෙන කොනෙන කොනෙන කෙනෙන කොන කොනෙන කොනෙන කොන කොනෙන කොනෙන කොනෙන කොනෙන කොනෙන කොනෙන ක
 [100000, 110000)
                                                                                                                                       [110000, 120000)
 [120000, 130000)
[130000, 140000)
[140000, 150000)
[150000, 160000)
```

- There's a difference of 461,828ns on average ≈ 451ns per interface.
- On a 64-Core system this was even slower, however I did not further research why. 5,089,948ns (5ms!) on average ≈ 4970ns per interface.



Simulated POD Creation



Simulated POD creation

- ☐ The objective is to assess whether the creation and teardown of PODs have any effect on the main thread.
- Our SUT does not have enough memory for a OCP cluster so we use standalone OVN to simulate this.
- The test is defined as follows:
 - We create 128 OVN switch ports simulating 128 PODs.
 - ☐ We continuously add and delete 10 additional PODs.
 - ☐ The cycle takes about 1.3 seconds (typical large OCP cluster creates around 6 PODs per seconds).
 - ☐ This test does not include any actual traffic.
- For the initial analysis we use <u>kernel_delay.py</u> on ovs-vswitchd's main thread.



64-Cores (condensed output):

4472 ovs-vswitchd	[SYSCAL	L STATISTICS		
NAME	NUMBER	COUNT	TOTAL ns	MAX ns
sendmsg	46	99109	4,958,977,988	89,797,200
recvmsg	47	115014	1,666,039,446	47,172,400
sendto	44	915	863,941,291	42,862,807
TOTAL(- poll):		412283	8,295,257,784	

□ 128-Cores (condensed output):

3765 ovs-vswitchd	[SYSCALL	STATISTICS]		
NAME	NUMBER	COUNT	TOTAL ns	MAX ns
sendmsg	46	91052	4,461,157,621	58,967,791
recvmsg	47	107575	3,307,665,973	57,098,125
sendto	44	932	485,015,325	8,727,121
• • •				
TOTAL(- poll):		377964	8,721,298,017	

- So the 128-Cores instance spends more time on recvmesg()
 - **□** ≈ 50%



☐ We used <u>analyze perf pmd syscall.py</u> script to the call trace:

```
#0 netdev_get_addrs() at ./ovs/lib/netdev.c:2314
#1 netdev_linux_get_addr_list() at ./ovs/lib/netdev-linux.c:3540
#2 netdev_get_addr_list() @ netdev_get_addr_list at ./ovs/lib/netdev.c:1518
#3 verify_prefsrc() at ./ovs/lib/ovs-router.c:182
#4 ovs router insert__() at ./ovs/lib/ovs-router.c:292
```

64-Cores (condensed output):

□ 128-Cores (condensed output):



- ☐ I did not further investigate this more 2x increase.
- Could this be related to the additional statistics cost in getifaddrs()?
 - □ It's not likely, see next section (\approx 6%)
- This needs further investigation!



getifaddrs() Statistics Gathering



□ Glibc's getifaddrs() function also retrieves all interface statistics
 □ Leads to unsesacry overhead for systems with a large number of CPU due to per-CPU IPv6 statistics.
 □ Can be avoided with the RTEXT_FILTER_SKIP_STATS flag
 □ We aim to test the impact on getifaddrs() with 1024 interfaces.
 □ We create 1024 veth interfaces and add them to OVS (see Kernel Interfaces Statistics Impact)
 □ Add USDT probes around getifaddrs() in netdev_get_addrs().
 □ Measure for 5 minutes doing 10 revalidations per second.



Results:

System	getifaddrs() runtime (ns)
64-CPUs	3985.7
128-CPUs	4227.6
Delta(%)	241.9(6.1%)

Next step would be to patch **getifaddrs()** in glibc to not gather the statistics.



Conclusion / Follow up



TO DO:

- \square Explore removing mutex in $nl_dump_next()$.
- Optimize the kmod's sw_flow_stats for faster combined stats gathering.
- Analyze netdev_get_addrs() delays in POD simulation.
- Check if unnecessary getifaddrs() stats cause delays; submit a glibc patch if needed.
- Asses of upcall-related impact measurements, particularly if ct() and nat() are involved.



Questions?



Thank you

Red Hat is the world's leading provider of enterprise open source software solutions. Award-winning support, training, and consulting services make Red Hat a trusted adviser to the Fortune 500.

- in linkedin.com/company/red-hat
- youtube.com/user/RedHatVideos
- facebook.com/redhatinc
- twitter.com/RedHat

