

Leveraging virtio-vsock in the cloud and containers

FOSDEM 2021 Virtual Conference
6th of February 2021

Andra Paraschiv, Amazon Web Services
Stefano Garzarella, Red Hat



Agenda

Introduction

Use Cases

- AWS Nitro Enclaves
- Kata Containers
- libkrun

Developing with AF_VSOCK

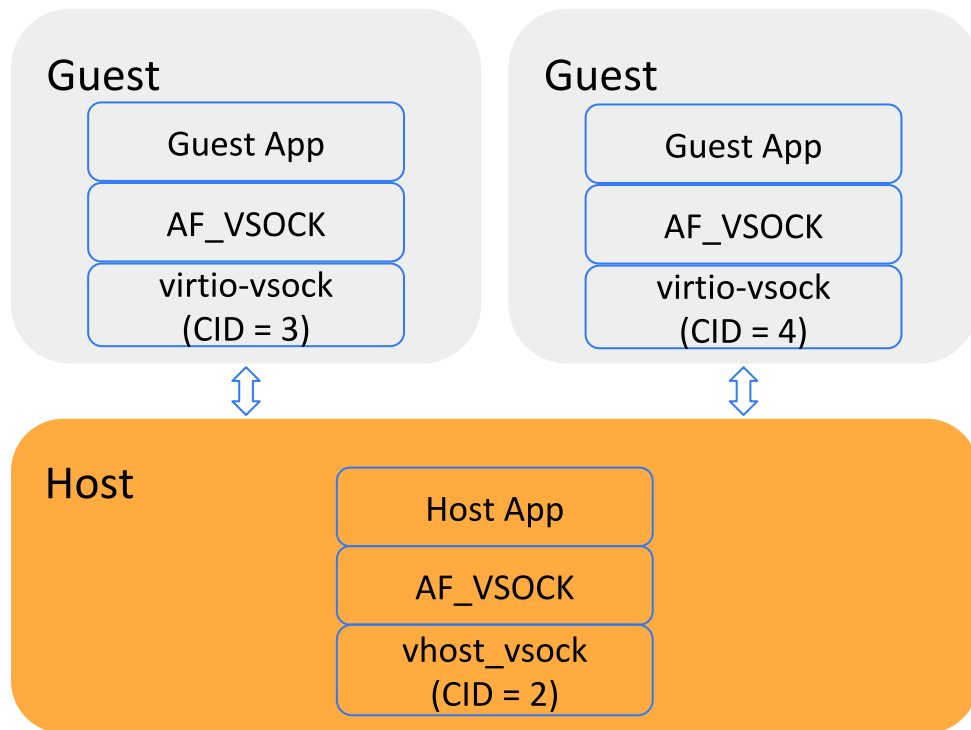
- Languages Bindings
- Communication Debugging
- Performance Evaluation

Next Steps

- SOCK_SEQPACKET
- Namespaces
- Multi-queue
- Shared Memory

Vsock - Overview

- POSIX Socket API (e.g. socket, bind, listen, connect)
- SOCK_STREAM / SOCK_DGRAM
- CID (Context Identifier) and port
 - Well-known CIDs
 - VMADDR_CID_ANY (0xFFFFFFFF)
 - VMADDR_CID_HYPERVISOR (0)
 - VMADDR_CID_LOCAL (1)
 - VMADDR_CID_HOST (2)
 - Privileged ports (port < 1024)
- Multi transports - support for nested VMs

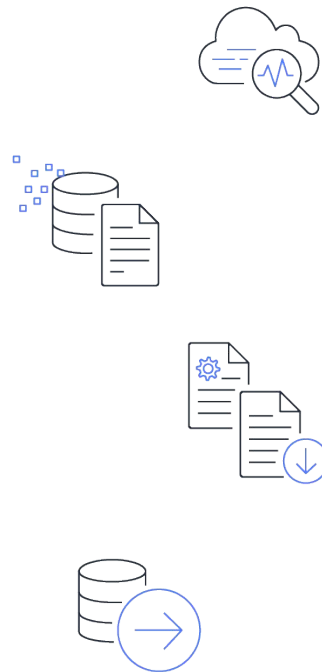


Vsock - Use Cases

Network Applications
(e.g. SOCK_STREAM apps)

Guest Agents
(e.g. QEMU guest agent, Kata containers agent)

Hypervisor Services
(e.g. file sharing)

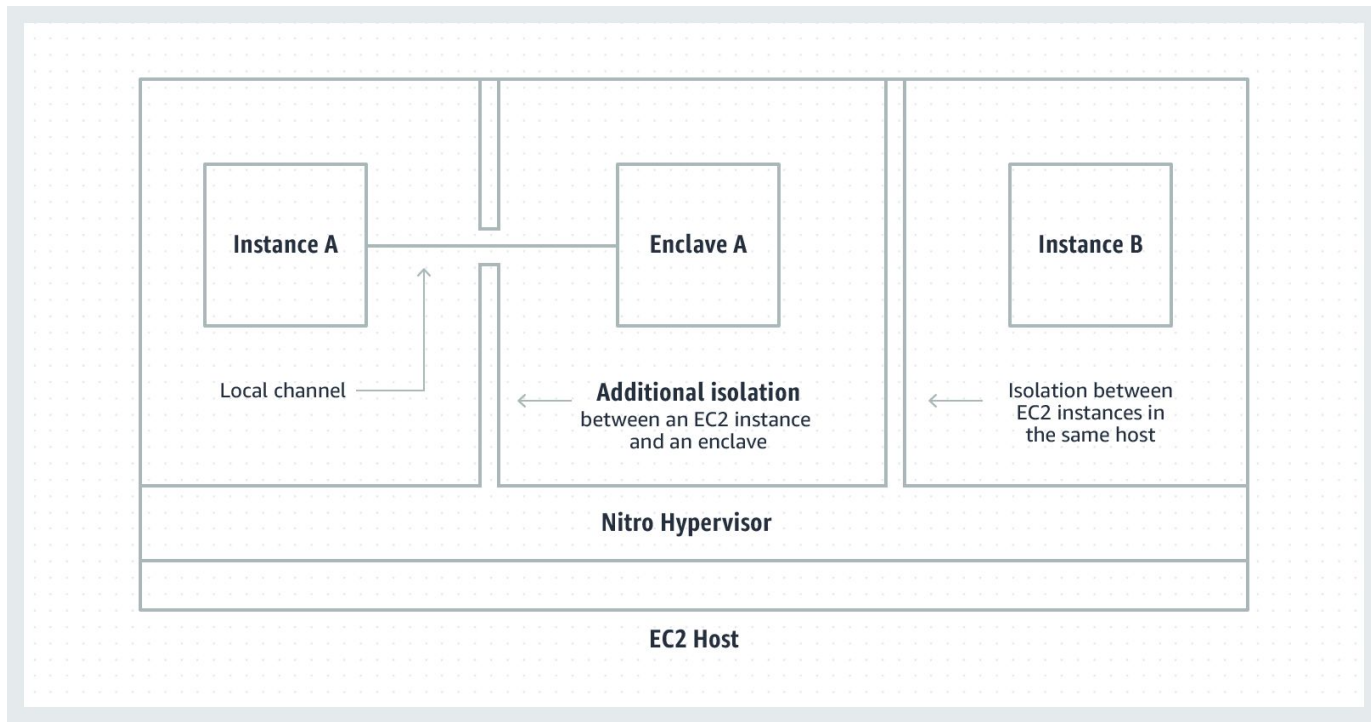


AWS Nitro Enclaves

- Isolated and trusted environment
- Cryptographic attestation
- Open source
- Portable across operating systems and architectures
- Integrated with AWS Key Management Service (KMS) and AWS Certificate Manager (ACM)



AWS Nitro Enclaves (2)



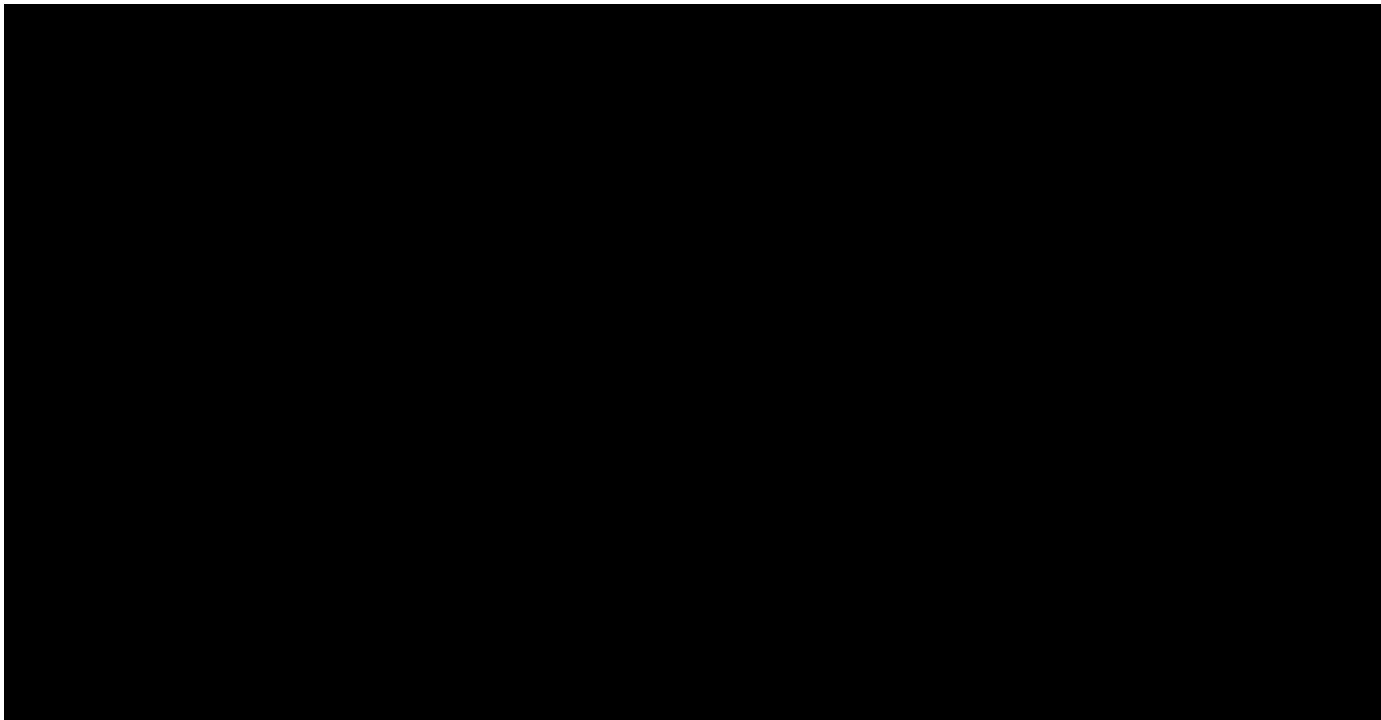
AWS Nitro Enclaves - Serial Console Access (Demo)

```
[[ec2-user@vsock-demo ~]$ lsmod | grep nitro_enclaves
nitro_enclaves 32768 0
[[ec2-user@vsock-demo ~]$ sudo lspci -v -s 00:02.0
00:02.0 Communication controller: Amazon.com, Inc. Device e4c1 (rev 01)
Subsystem: Amazon.com, Inc. Device e4c1
Physical Slot: 2
Flags: fast devsel
Memory at febf0000 (32-bit, non-prefetchable) [size=4K]
Memory at febf1000 (32-bit, non-prefetchable) [size=4K]
Capabilities: [40] MSI-X: Enable+ Count=3 Masked-
Kernel driver in use: nitro_enclaves
Kernel modules: nitro_enclaves

[[ec2-user@vsock-demo ~]$ sudo lspci -v -s 00:02.1
00:02.1 Communication controller: Red Hat, Inc. Device 1053 (rev 01)
Subsystem: Red Hat, Inc. Device 1053
Physical Slot: 2
Flags: fast devsel
[virtual] Memory at febf2000 (32-bit, non-prefetchable) [size=4K]
[virtual] Memory at febe0000 (32-bit, non-prefetchable) [size=32K]
Capabilities: [90] Vendor Specific Information: VirtIO: DeviceCfg
Capabilities: [80] Vendor Specific Information: VirtIO: ISR
Capabilities: [70] Vendor Specific Information: VirtIO: <unknown>
Capabilities: [5c] Vendor Specific Information: VirtIO: Notify
Capabilities: [4c] Vendor Specific Information: VirtIO: CommonCfg
Capabilities: [40] MSI-X: Enable+ Count=4 Masked-
Kernel driver in use: virtio-pci
Kernel modules: virtio_pci

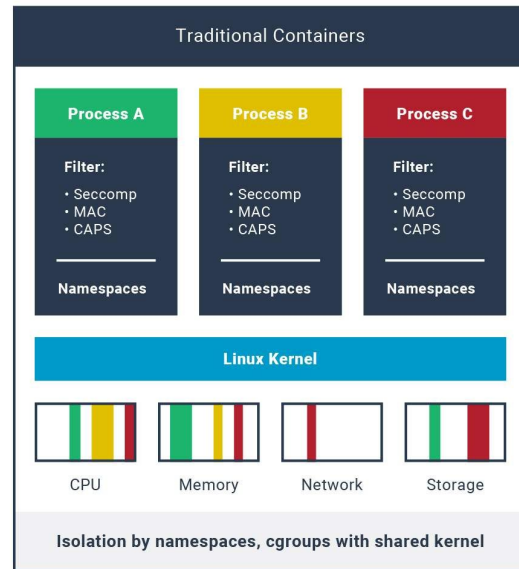
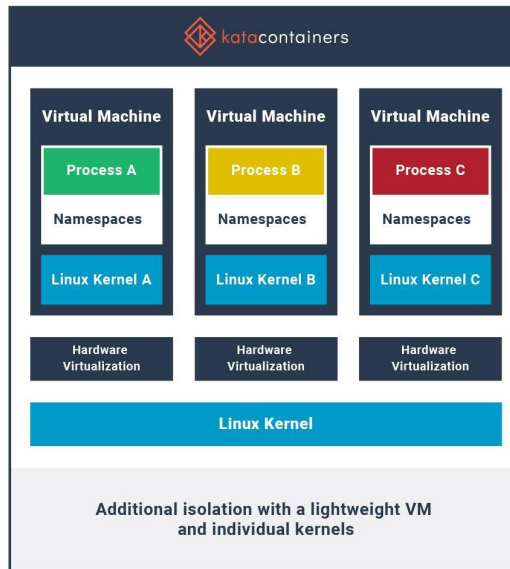
[[ec2-user@vsock-demo ~]$ cat /sys/devices/system/node/node0/hugepages/hugepages-
```

AWS Nitro Enclaves - Vsock Sample (Demo)



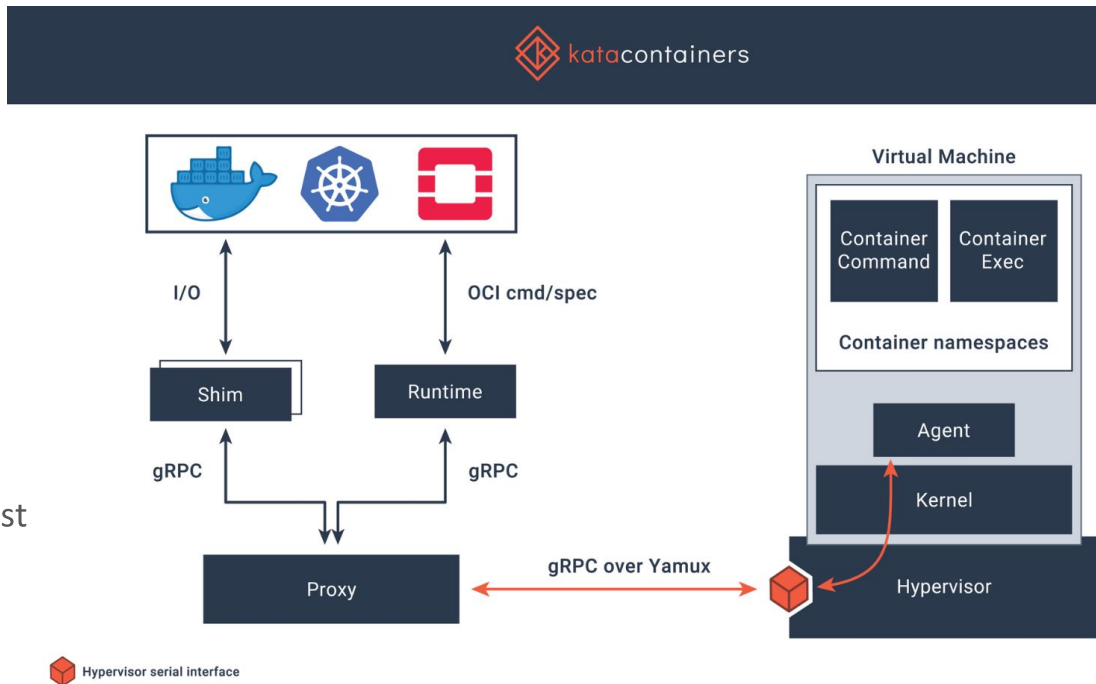
Kata Containers

- **secure** container runtime
- **lightweight** virtual machines
 - feel and perform like containers
- stronger workload **isolation**
 - hardware virtualization technology as a second layer of defense



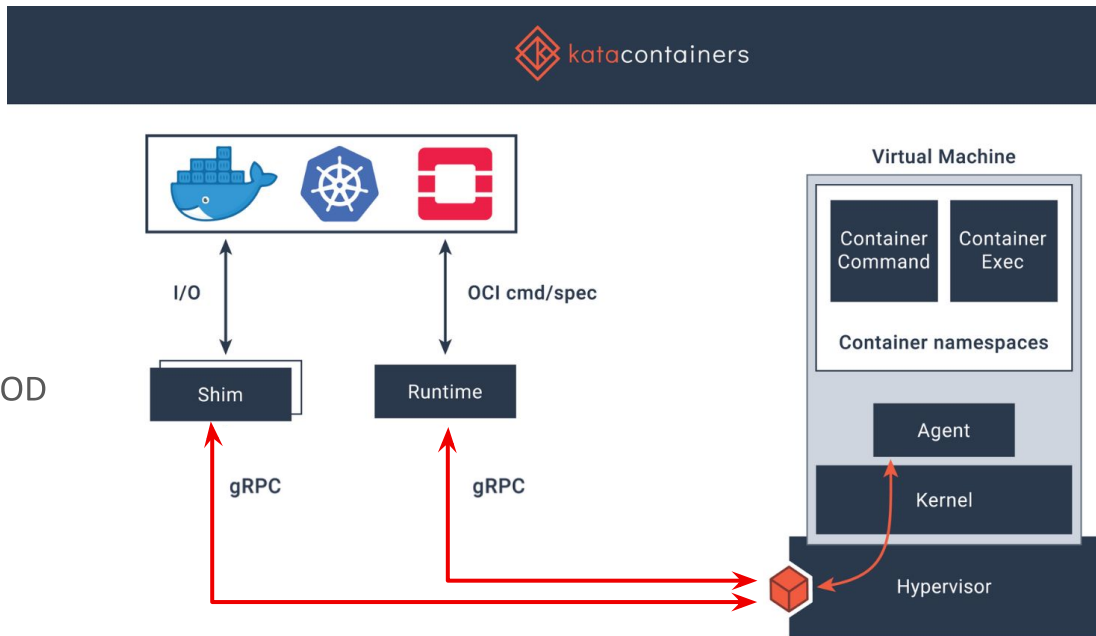
Kata Containers using serial ports

- guest processes
 - serial port device
- host processes
 - Unix socket
- Portable solution
- Cons
 - one process at a time can access serial link
 - serial port and Unix socket must be multiplexed by using kata-proxy and Yamux



Kata Containers using VSOCK

- Listen sockets can accept connections from multiple clients
 - does not require multiplexers (kata-proxy and Yamux)
- High density
 - kata-proxy uses ~4.5MB per POD
- Reliability
 - if kata-proxy dies all POD connections get broken



 virtio-vsock device

<https://github.com/kata-containers/documentation/blob/master/design/Vsocks.md>

Kata Containers demo

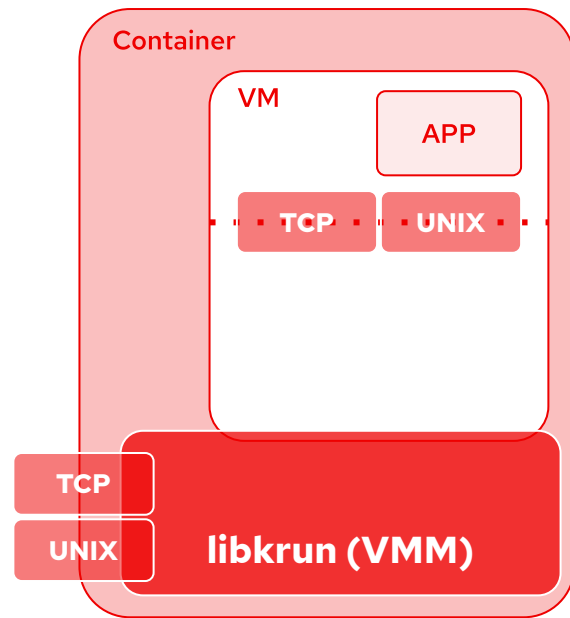
```
1: @ada0d224b72c/ 63x15 2: ./demo-host2.sh 63x15
00:1f.6 Ethernet controller: Intel Corporation Ethernet Connect
ion (4) I219-LM (rev 21)
04:00.0 PCI bridge: Intel Corporation JHL6240 Thunderbolt 3 Bri
dge (Low Power) [Alpine Ridge LP 2016] (rev 01)
05:00.0 PCI bridge: Intel Corporation JHL6240 Thunderbolt 3 Bri
dge (Low Power) [Alpine Ridge LP 2016] (rev 01)
05:01.0 PCI bridge: Intel Corporation JHL6240 Thunderbolt 3 Bri
dge (Low Power) [Alpine Ridge LP 2016] (rev 01)
05:02.0 PCI bridge: Intel Corporation JHL6240 Thunderbolt 3 Bri
dge (Low Power) [Alpine Ridge LP 2016] (rev 01)
06:00.0 System peripheral: Intel Corporation JHL6240 Thunderbol
t 3 NHI (Low Power) [Alpine Ridge LP 2016] (rev 01)
07:00.0 PCI bridge: Intel Corporation JHL6540 Thunderbolt 3 Bri
dge (C step) [Alpine Ridge 4C 2016] (rev 02)
08:00.0 PCI bridge: Intel Corporation JHL6540 Thunderbolt 3 Bri
dge (C step) [Alpine Ridge 4C 2016] (rev 02)
08:01.0 PCI bridge: Intel Corporation JHL6540 Thunderbolt 3 Bri
dge (C step) [Alpine Ridge 4C 2016] (rev 02)
08:02.0 PCI bridge: Intel Corporation JHL6540 Thunderbolt 3 Bri
dge (C step) [Alpine Ridge 4C 2016] (rev 02)

3: @781ca166b285/ 63x15 4: ./demo-host4.sh 63x15
ls[root@781ca166b285 /]# lspci
00:00.0 Host bridge: Intel Corporation 82G33/G31/P35/P31 Expres
s DRAM Controller
00:01.0 Communication controller: Red Hat, Inc. Virtio console
00:02.0 PCI bridge: Red Hat, Inc. QEMU PCI-PCI bridge
00:03.0 SCSI storage controller: Red Hat, Inc. Virtio SCSI
00:04.0 Unclassified device [00ff]: Red Hat, Inc. Virtio RNG
00:05.0 Communication controller: Red Hat, Inc. Virtio socket
00:06.0 Mass storage controller: Red Hat, Inc. Device 105a (rev
01)
00:07.0 Ethernet controller: Red Hat, Inc. Virtio network devic
e
00:1f.0 ISA bridge: Intel Corporation 82801IB (ICH9) LPC Interf
ace Controller (rev 02)
00:1f.2 SATA controller: Intel Corporation 82801IR/IO/IH (ICH9R
/DO/DH) 6 port SATA Controller [AHCI model] (rev 02)

har-e84208feddfb6791,tag=kataShared,romfile=-netdev tap,id=net
work-0,vhost=on,vhostfds=4,fds=5 -device driver=virtio-net-pci,
netdev=network-0,mac=5e:92:72:f5:45:3a,disable-modern=false,mq=
on,vectors=4,romfile=-global kvm-pit.lost_tick_policy=discard
-vga none -no-user-config -nodefaults -nographic -daemonize -ob
ject memory-backend-file,id=dimm1,size=2048M,mem-path=/dev/shm,
share=on -numa node,memdev=dimm1 -kernel /usr/lib/modules/5.9.1
6-100.fc32.x86_64/vmlinuz -initrd /var/cache/kata-containers/os
builder-images/5.9.16-100.fc32.x86_64/fedora-kata-5.9.16-100.fc
32.x86_64.initrd -append tsc=reliable no_timer_check rcupdate.r
cu expedited=1 i8042.direct=1 i8042.dumbkbd=1 i8042.nopnp=1 i80
42.noaux=1 noreplace-smp reboot=k console=hvc0 console=hvc1 iom
mu=off cryptomgr.notevents net.ifnames=0 pci=lastbus=0 quiet pani
c=1 nr_cpus=8 agent.use_vsock=true scsi_mod.scan=none -pidfile
/run/vc/vm/781ca166b285f08190347aaa87be1a2159c077f5ce45834922d
c8732482214e/nid -smp 1 cores=1 threads=1 sockets=8 maxcpus=8
```

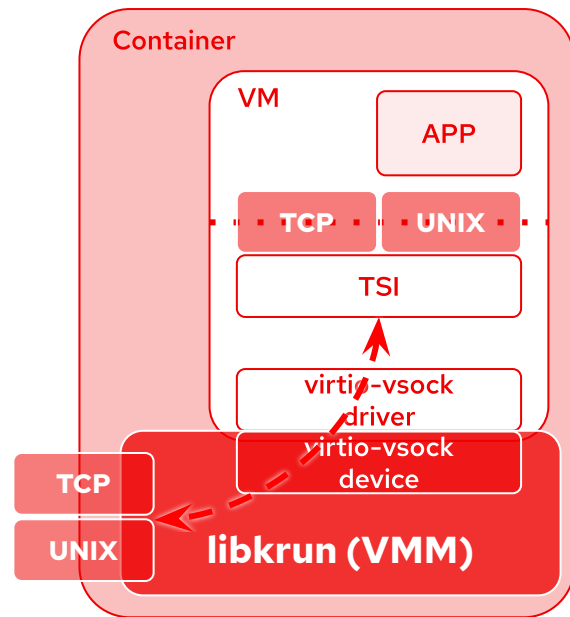
libkrun

- **Dynamic library** to run processes in a partially isolated environment
 - HW Virtualization
 - Linux KVM
 - macOS Hypervisor.Framework
 - VMM with a minimum amount of emulated devices
 - virtio-console
 - virtio-fs
 - virtio-sock
 - virtio-balloon
 - Transparent Socket Impersonation (TSI)



libkrun: Transparent Socket Impersonation (TSI)

- VMs have network connectivity without a virtual NIC
- Network connections encapsulated using VSOCK sockets
 - Guest kernel
 - Intercepts sockets syscall
 - Forward requests to the host using virtio-vsock
 - libkrun (VMM in the host)
 - Receives requests through virtio-vsock
 - behaves like the application in the guest
 - making the same sockets syscall in the host
- Fully transparent for application running in the VMs
- Proof of Concept
 - limited to TCP/IPv4 and UNIX domain sockets



libkrun goals

- Simple API to easily allow other applications to run processes in a VM-based isolated environment
 - Small footprint
 - RAM
 - CPU
 - Boot time
 - Compatible with a reasonable amount of use cases
 - Stronger isolation for containers
 - Fully encrypted containers (using SEV/TDX)
 - Self-isolating microservices
 - Running Lightweight VMs based on OCI images on Linux and macOS
- <https://twitter.com/slpnix/status/1348406770990002178?s=20>



libkrun demo

```
1: @70870d2e0079:/ 63x15 ▾
00001000-0009fbff : System RAM
000f0000-000fffff : System ROM
00100000-1fffffff : System RAM
01000000-01a00bd0 : Kernel code
01c00000-01d71fff : Kernel rodata
01e00000-01e9d37f : Kernel data
0203d000-025fffff : Kernel bss
d0000000-d0000fff : virtio-mmio.0
d0000000-d0000fff : virtio-mmio.0
d0001000-d0001fff : virtio-mmio.1
d0001000-d0001fff : virtio-mmio.1
d0002000-d0002fff : virtio-mmio.2
d0002000-d0002fff : virtio-mmio.2
fee00000-fee00fff : Local APIC
[root@70870d2e0079 /]# ip li

2: ./demo-host2.sh 63x15 ▾
host $

3: @781ca166b285:/ 63x15 ▾
host $

4: ./demo-host4.sh 63x15 ▾
15:02:47.323610 VIRTIO 3319058279.1024 > 2.52193 PAYLOAD, length 115
15:02:47.324211 VIRTIO 2.52193 > 3319058279.1024 PAYLOAD, length 93
15:02:47.348178 VIRTIO 3319058279.1024 > 2.52193 PAYLOAD, length 146
15:02:47.348550 VIRTIO 2.52193 > 3319058279.1024 PAYLOAD, length 106
15:02:47.348842 VIRTIO 2.52193 > 3319058279.1024 DISCONNECT, length 76
15:02:47.348979 VIRTIO 3319058279.1024 > 2.52193 DISCONNECT, length 76
15:02:47.349724 VIRTIO 3319058279.1024 > 2.52193 DISCONNECT, length 76
```


Developing with AF_VSOCK

Languages providing AF_VSOCK bindings

- C
 - glibc >= 2.18 [2013-08-10]
- Python
 - python >= 3.7 alpha 1 [2017-09-19]
- Golang
 - <https://github.com/mdlayher/vsock>
- Rust
 - libc crate >= 0.2.59 [2019-07-08]
 - struct sockaddr_vm
 - VMADDR_* macros
 - nix crate >= 0.15.0 [2019-08-10]
 - VSOCK supported in the socket API (nix::sys::socket)



Python example

GUEST

```
# Client running in the guest
import socket

s = socket.socket(socket.AF_VSOCK, socket.SOCK_STREAM)
s.connect((socket.VMADDR_CID_HOST, 1234))

s.send(b'Hello, world')
```



VSOCK

HOST

```
# Server running in the host
import socket

s = socket.socket(socket.AF_VSOCK, socket.SOCK_STREAM)
s.bind((socket.VMADDR_CID_ANY, 1234))
s.listen()
client, addr = s.accept()

data = client.recv(1024)
print("CID: {} port: {} data: {}".format(addr[0], addr[1], data))
```

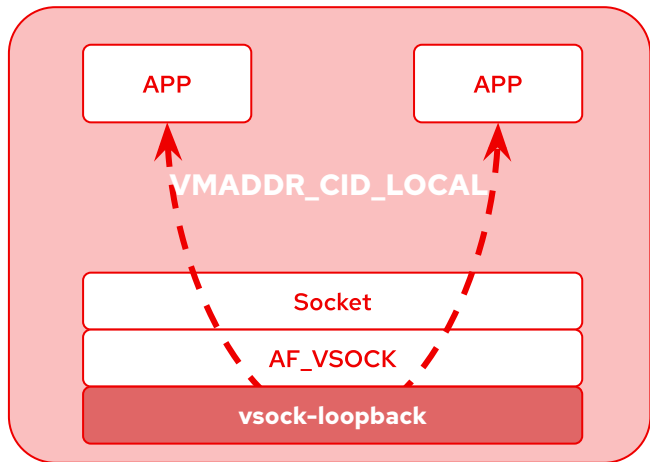
```
1. @70870d2e0079/63x15 ~
host $ python
Python 3.8.6 (default, Sep 25 2020, 00:00:00)
[GCC 10.2.1 20200723 (Red Hat 10.2.1-1)] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> import socket
>>> s = socket.socket(socket.AF_VSOCK, socket.SOCK_STREAM)
>>> s.bind((socket.VMADDR_CID_ANY, 1234))
>>> s.listen()
>>> client, addr = s.accept()
>>>

3. @781c166b285/63x15 ~
host $

2. /demo-host2.sh 63x15 ~
guest $ python
Python 3.8.6 (default, Sep 25 2020, 00:00:00)
[GCC 10.2.1 20200723 (Red Hat 10.2.1-1)] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> import socket
>>> s = socket.socket(socket.AF_VSOCK, socket.SOCK_STREAM)
>>> s.connect((socket.VMADDR_CID_HOST, 1234))
>>> s.send(b'Hello, world')
12
>>>

4. /demo-host4.sh 63x15 ~
h 93
15:02:47.348178 VIRTIO 3319058279.1024 > 2.52193 PAYLOAD, length
h 146
15:02:47.348550 VIRTIO 2.52193 > 3319058279.1024 PAYLOAD, length
h 106
15:02:47.348842 VIRTIO 2.52193 > 3319058279.1024 DISCONNECT, length
76
15:02:47.348979 VIRTIO 3319058279.1024 > 2.52193 DISCONNECT, length
76
15:02:47.349724 VIRTIO 3319058279.1024 > 2.52193 DISCONNECT, length
76
15:07:37.803519 VIRTIO 42.50667 > 2.1234 CONNECT, length 76
15:07:37.803617 VIRTIO 2.1234 > 42.50667 CONNECT, length 76
15:07:41.302608 VIRTIO 42.50667 > 2.1234 PAYLOAD, length 88
```

Local communication with vsock-loopback



- Local communication without VMs
 - Tests
 - Debug
- vsock-loopback
 - New transport available from Linux v5.6
 - `modprobe vsock-loopback`
- CIDs
 - VMADDR_CID_LOCAL (1)
 - Well-known CID for loopback

Useful tools with AF_VSOCK support

- **wireshark** >= 2.40 [2017-07-19]
- **iproute2** >= 4.15 [2018-01-28]
 - ss
- **tcpdump** >= 4.99 [2020-12-30]
 - libpcap >= 1.9 [2018-06-24]
 - Fedora backported vsock patches on tcpdump 4.9.3
- **nmap** >= 7.80 [2019-08-10]
 - ncat
- **nbd**
 - nbdkit >= 1.15.5 [2019-10-19]
 - libnbd >= 1.1.6 [2019-10-19]
- **socat** >= 1.7.4 [2021-01-04]
- **iperf-vsock**
 - iperf3 fork
 - <https://github.com/stefano-garzarella/iperf-vsock>

Concatenate and redirect sockets

- Ncat - <https://nmap.org/ncat/>
 - new parameter
 - --vsock
 - Examples
 - ncat --vsock -l 4321
 - ncat --vsock 1 4321
- socat - <http://www.dest-unreach.org/socat/>
 - New address types
 - VSOCK-LISTEN
 - VSOCK-CONNECT
 - Examples
 - socat - VSOCK-LISTEN:4321
 - socat - VSOCK-CONNECT:1:4321



```
1: @70870d2e0079f/63x15
root $ ncat --vsock -l 1234
Hello in local
root $ socat - VSOCK-CONNECT:42:1234
Hello from the host!
root $ socat TCP4:LISTEN:4321,reuseaddr,fork VSOCK-CONNECT:42:22
[ ]

2: /demo-host2.sh 63x15
guest $ socat - VSOCK-LISTEN:1234
Hello from the host!
guest $ socat VSOCK-LISTEN:22,reuseaddr,fork TCP:localhost:22
[ ]

3: root@fosdemguest-63x15
Process
v_str LISTEN 0 0 *:1234 *:
users:({"ncat",pid=31842,fd=3})
host $ ss --vsock -ap
NetId State Recv-Q Send-Q Local Address:Port
Peer Address:Port Process
v_str ESTAB 0 0 1:3233074146
1:1234 users:({"ncat",pid=31891,fd=3})
v_str ESTAB 0 0 1:1234
1:3233074146 users:({"ncat",pid=31842,fd=4})
host $ ssh -p 4321 root@localhost
Warning: Permanently added '[localhost]:4321' (ECDSA) to the list of known hosts.
Last login: Thu Jan 14 09:07:04 2021
[root@fosdemguest ~]# uname -a
Linux localhost 3.10.0-112.el7.x86_64 #1 SMP Mon Aug 14 22:03:11 UTC 2017; root:x86_64 Linux

4: /demo-host4.sh 63x15
15:10:49.580220 VIRTIO 2.52196 > 42.22 PAYLOAD, length 112
15:10:49.582090 VIRTIO 2.52196 > 42.22 PAYLOAD, length 112
15:10:49.571055 VIRTIO 2.52196 > 42.22 PAYLOAD, length 112
15:10:49.579185 VIRTIO 42.22 > 2.52196 PAYLOAD, length 112
15:10:49.708094 VIRTIO 2.52196 > 42.22 PAYLOAD, length 112
15:10:49.780445 VIRTIO 42.22 > 2.52196 PAYLOAD, length 112
15:10:49.771137 VIRTIO 2.52196 > 42.22 PAYLOAD, length 112
15:10:49.773294 VIRTIO 42.22 > 2.52196 PAYLOAD, length 112
15:10:49.986275 VIRTIO 2.52196 > 42.22 PAYLOAD, length 112
15:10:49.907814 VIRTIO 42.22 > 2.52196 PAYLOAD, length 112
15:10:50.062690 VIRTIO 2.52196 > 42.22 PAYLOAD, length 112
15:10:50.063685 VIRTIO 42.22 > 2.52196 PAYLOAD, length 112
15:10:50.313136 VIRTIO 2.52196 > 42.22 PAYLOAD, length 112
15:10:50.316752 VIRTIO 42.22 > 2.52196 PAYLOAD, length 112
```

Dump and analyze AF_VSOCK traffic

- Prerequisites

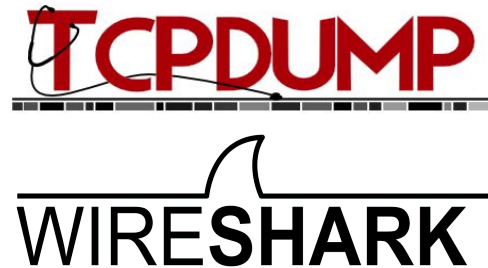
- vsockmon kernel module available (CONFIG_VSOCKMON=m)
- Create vsockmon virtual device to monitor AF_VSOCK sockets
 - `ip link add type vsockmon`
 - `ip link set vsockmon0 up`

- tcpdump

- `tcpdump -i vsockmon0`

- Wireshark

- `wireshark -k -i vsockmon0`



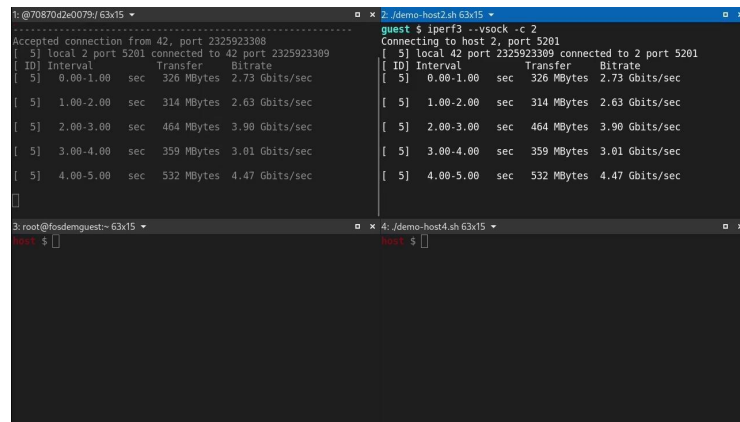
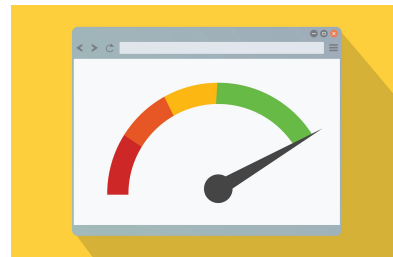
```
1: @70870d2e0079f/63x15
host $ ncat --vsock 42 1234
Ncat: Connection reset by peer.
host $

2: /demo-host2.sh 63x15
guest $ ncat --vsock -l 1234
[]

3: root@f0admguest-63x15
host $
host $ sudo ip link add type vsockmon
[sudo] password for stefano:
host $ sudo ip link set vsockmon0 up
host $ tcpdump -i vsockmon0
tcpdump: verbose output suppressed, use -v or -vv for full protocol decode
listening on vsockmon0, link-type VSOCK (Linux vsock), capture size 262144 bytes
15:11:48.457868 VIRTIO 2.52197 > 42.1234 CONNECT, length 76
15:11:48.458658 VIRTIO 42.1234 > 2.52197 DISCONNECT, length 76
^C
2 packets captured
2 packets received by filter
0 packets dropped by kernel
host $ wireshark -k -i vsockmon0
```

Performance evaluation

- iperf-vsock: iperf3 fork with AF_VSOCK support
 - <https://github.com/stefano-garzarella/iperf-vsock>
 - new parameter
 - --vsock
 - Examples
 - host\$ iperf3 --vsock -s
 - guest\$ iperf3 --vsock -c 2
 - Firecracker's hybrid VSOCK over AF_UNIX
 - Host runs iperf server
 - iperf3 --vsock -s -B /tmp/vm.vsock
 - Host runs iperf client
 - iperf3 --vsock -c /tmp/vm.vsock



Next Steps

- SOCK_SEQPACKET
 - Sequenced, reliable, two-way connection-based data transmission path for datagrams
 - [PATCH 0/5] virtio/vsock: introduce SOCK_SEQPACKET support
 - <https://lore.kernel.org/netdev/20210103195454.1954169-1-arseny.krasnov@kaspersky.com>
- Network namespaces
 - Useful for partitioning VMs or in a nested environment
 - [PATCH net-next 0/3] vsock: support network namespace
 - <https://lore.kernel.org/lkml/20200116172428.311437-1-sgarzare@redhat.com/>

Next Steps (2)

- Multi-queue
 - Currently one TX / RX queue per vsock device
 - I/O intensive workloads
 - Multiple endpoints communication
- Shared memory
 - Vsock communication using shared memory regions
 - Less buffer copies
 - Fast communication

Q&A

Andra Paraschiv <andraps@amazon.com>

Stefano Garzarella <sgarzare@redhat.com>

Blog: <https://stefano-garzarella.github.io/>

IRC: sgarzare on #qemu irc.oftc.net

Thank you!

References

- [1] <https://man7.org/linux/man-pages/man7/vsock.7.html>
- [2] <https://aws.amazon.com/ec2/nitro/nitro-enclaves/>
- [3] <https://docs.aws.amazon.com/enclaves/latest/user/nitro-enclave.html>
- [4] https://www.kernel.org/doc/html/latest/virt/ne_overview.html
- [5] <https://github.com/aws/aws-nitro-enclaves-cli>
- [6] <https://github.com/aws/aws-nitro-enclaves-samples>

References (2)

[7] <https://katacontainers.io>

[8] <https://github.com/kata-containers/documentation/blob/master/design/Vsocks.md>

[9] <https://katacontainers.io/learn>

[10] <https://github.com/kata-containers/documentation/blob/master/design/architecture.md>

[11] <https://github.com/containers/libkrun>

[12] <https://github.com/containers/libkrunfw>