



# Seamless Kernel Update

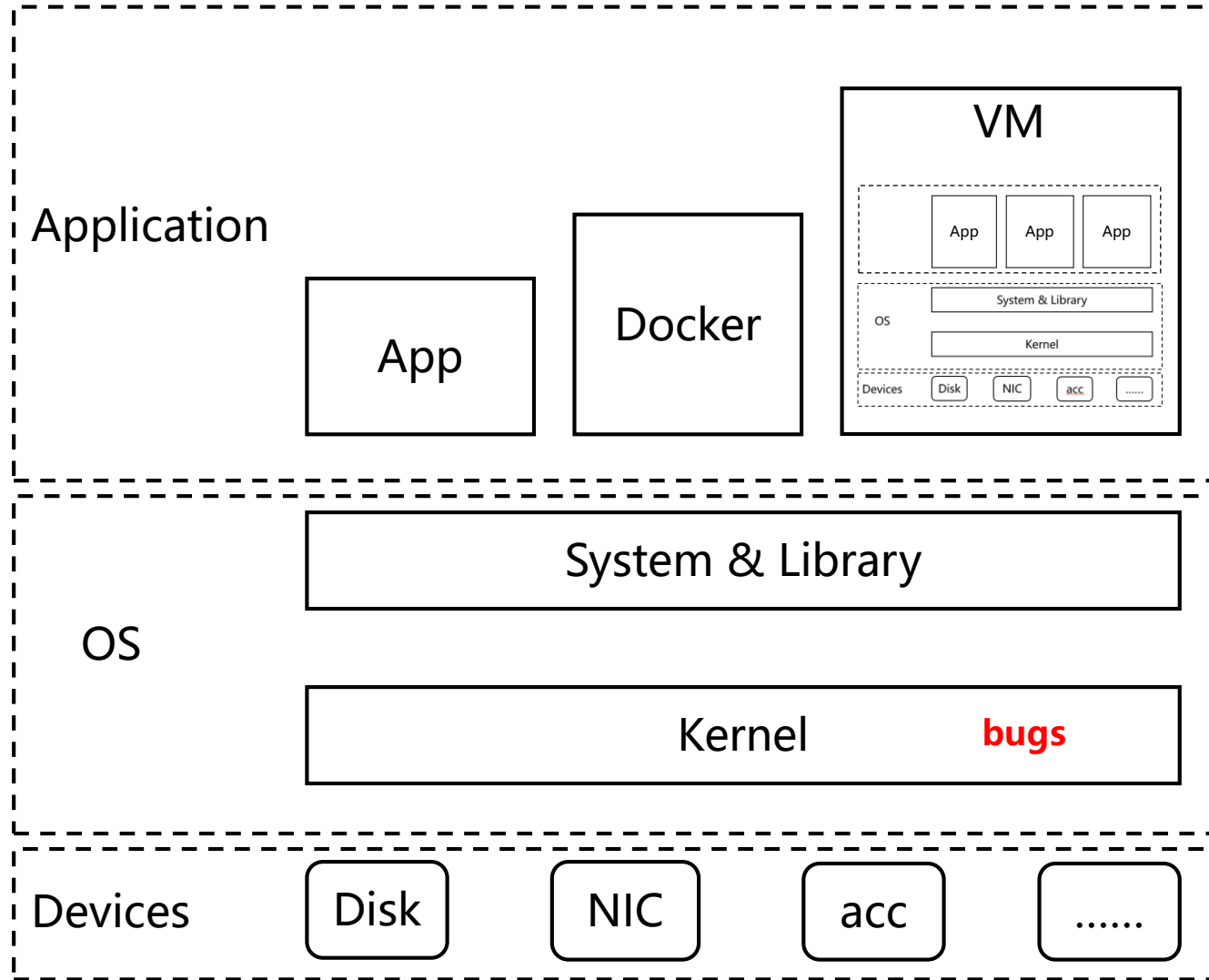
**Author:** 桑琰 [sangyan@huawei.com](mailto:sangyan@huawei.com)  
朱玲 [zhuling8@huawei.com](mailto:zhuling8@huawei.com)





- 01 Background
- 02 Froze/Resume the Application
- 03 Keep Memory
- 04 Kernel Fast boot
- 05 Keep Device State
- 06 Demo - Benchmark
- 07 Todo
- 08 Q&A

# ► 01 Background



## Kernel Bugs

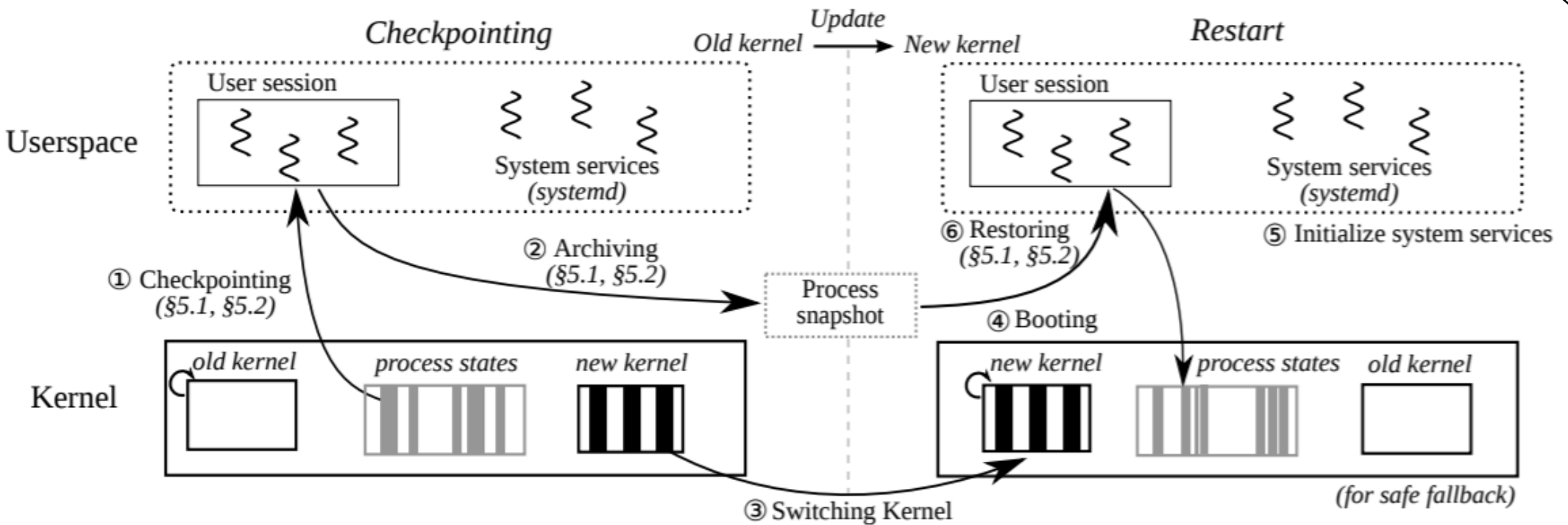
- Kernel Live patch to fix the bug
  - No live patch for some bug
  - Manager Different
- Live migration APP/VM & Reboot
  - No way to the pass-thought device
  - Large memory to transmit

## For Example

Machine: Bare-Metal Server Kunpeng 920  
Memory: 380G  
Application: Mysql (DB)

\* Different to transmit some memory

# ► 1.1 Instant OS Updates via Userspace Checkpoint-and-Restart

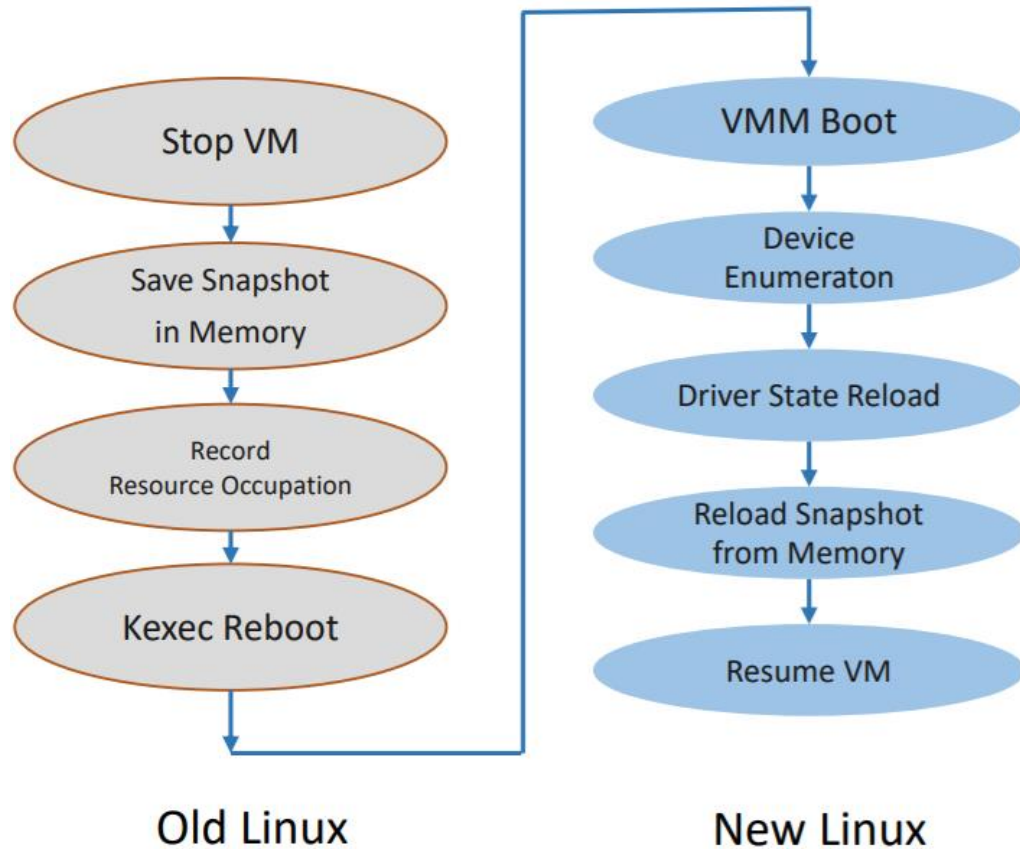
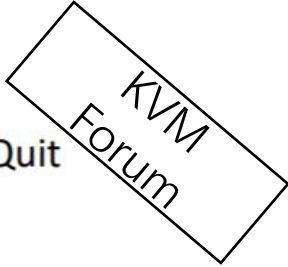


USENIX

**Figure 2:** Overview of KUP's updating procedures. KUP first checkpoints user's processes ①, and archives their snapshots ②. After checkpointing selected processes in a user's current session, KUP replaces the old kernel to the new kernel image ③, and finally switches to the new kernel ④. After the new kernel boots, KUP first initializes its system daemons ⑤, and finally restores snapshots of user applications ⑥.

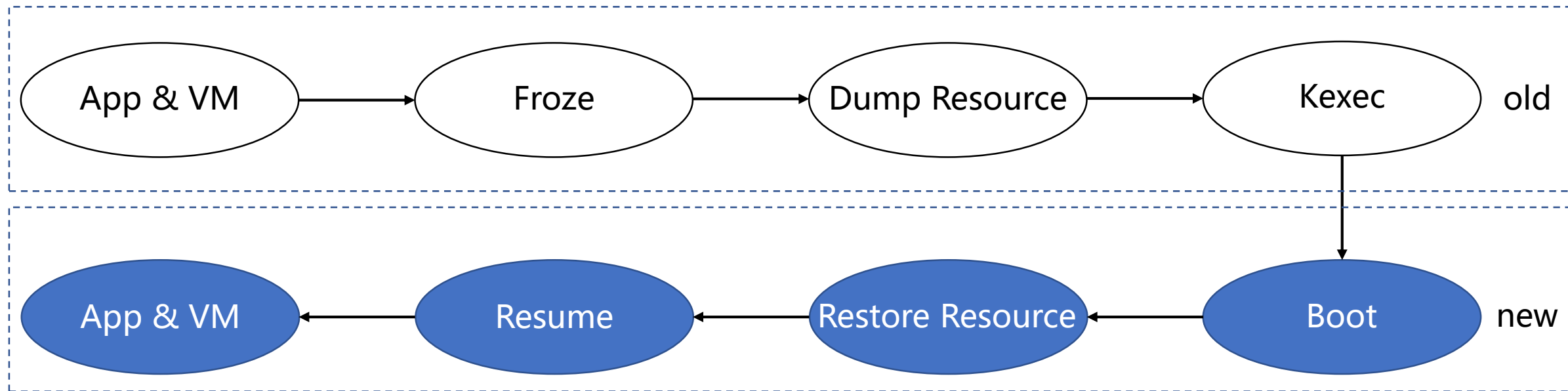
# ▶ 1.2 Seamless Cloud System Upgrade with VMM fast Restart

## High Level Flow



- Save Snapshot in Memory (in old Linux) & VM Quit
  - DAX filesystem in DRAM-as-PMEM
  - Don't free HW resources (IRTE, etc.)
- Record Resource Occupation
  - Device list, memory, etc.
- KEXEC Reboot
  - No hardware clobber in driver shutdown
- VMM Boot (new Linux)
  - Reserve resources
- Device Enumeration
  - No hardware clobber in PCI enumeration
  - No native driver attaching
- Driver State Reload
  - IOMMU driver reload state
- Reload Snapshot from Memory (in new Linux)
  - Re-enable DAX filesystem in DRAM-as-PMEM
- Resume VM
  - Reload DMA mapping
  - Re-enable MSI/MSIX

## ► 02 Froze/Resume the Application



QEMU: qemu save/restore

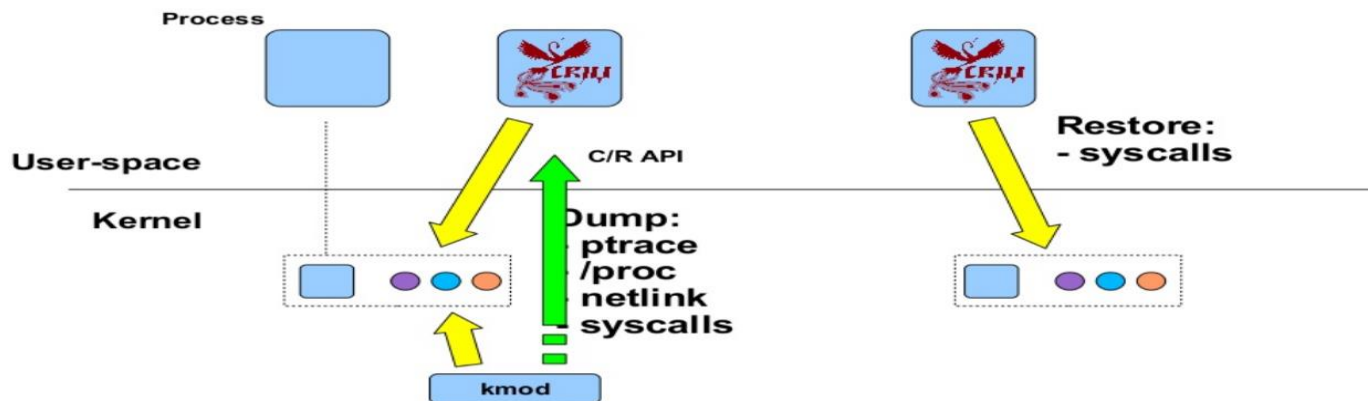
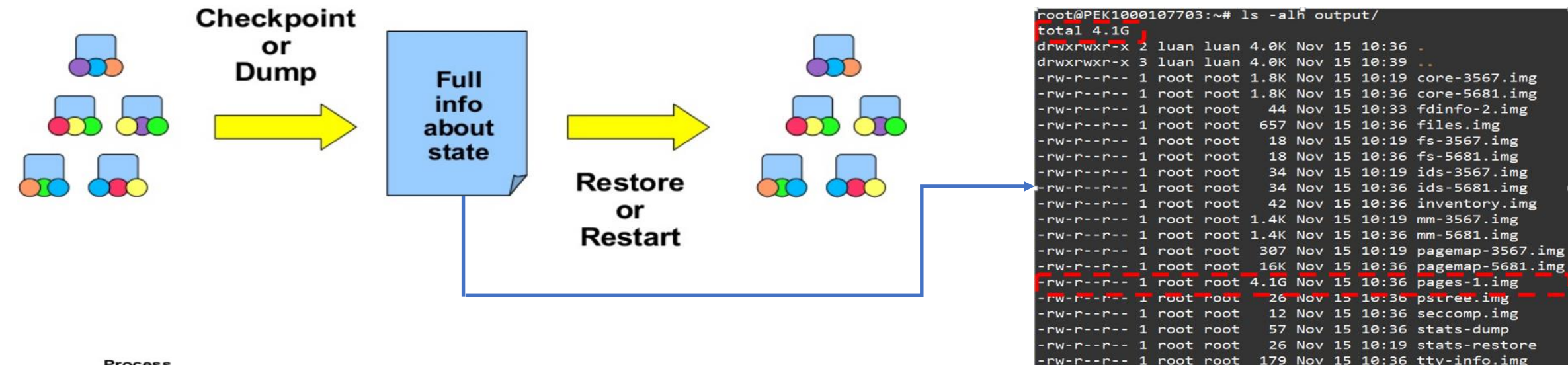
CRIU: criu dump/restore

DMTCP: dmtcp save/restore

\* Migration the Application/VM to the Same Machine, the Machine exist at the different time

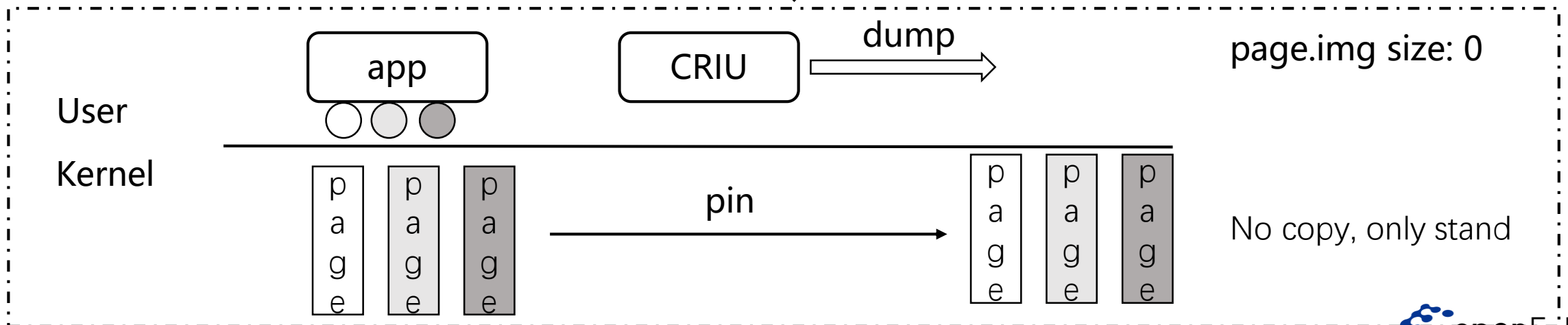
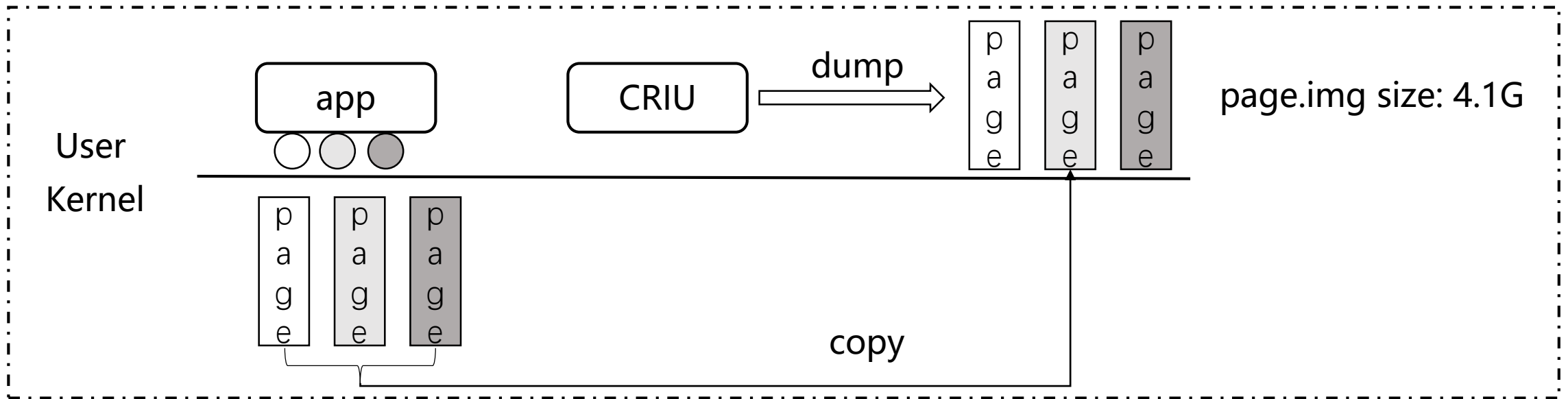


# ► 03 Keep Memory



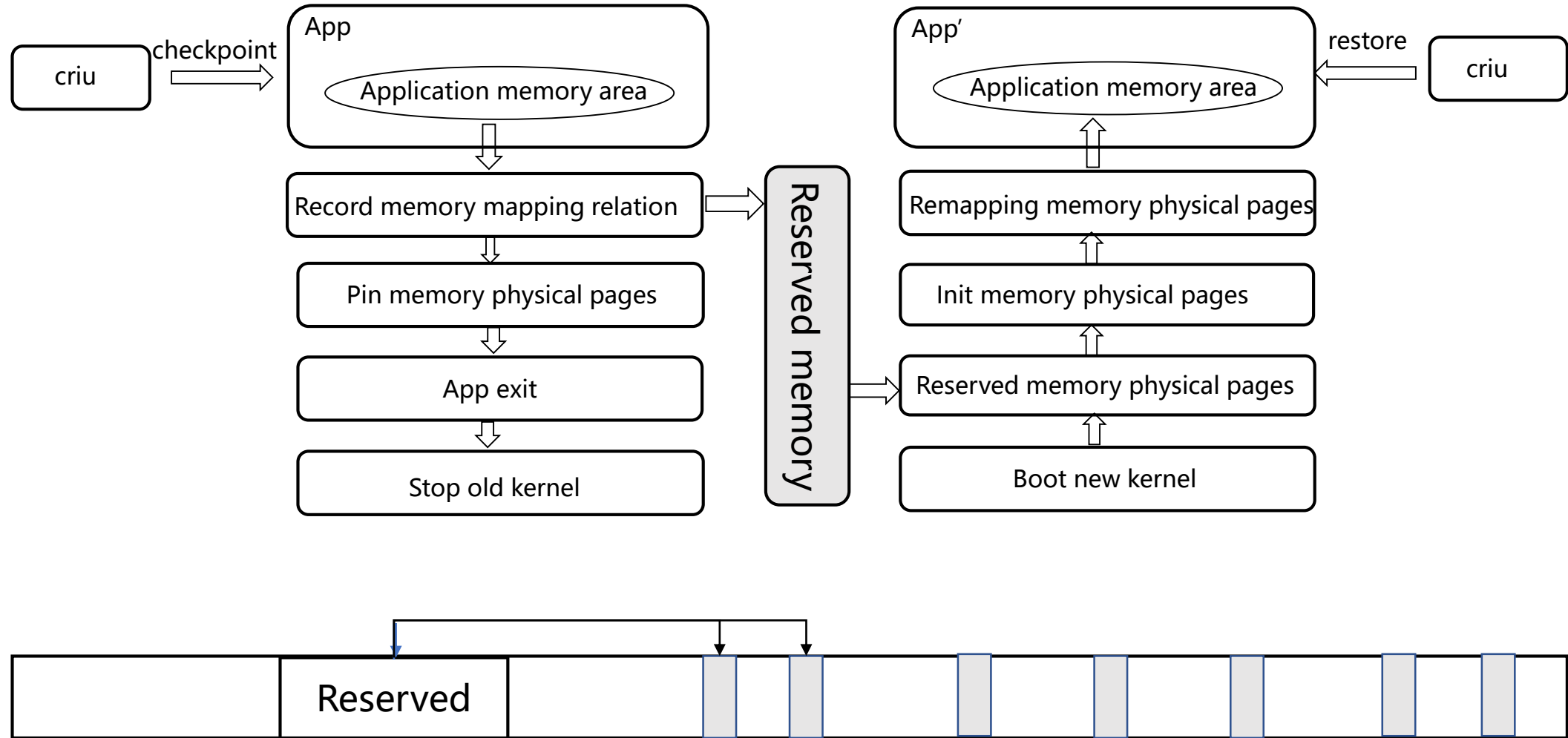
Write the copy of application memory into disk file, data will be large and get poor performance. For example will write 4.1G memory copy into disk. No copy and no write, only keep the app into memory (Pin memory).

## ▶ 3.1 Pin application memory

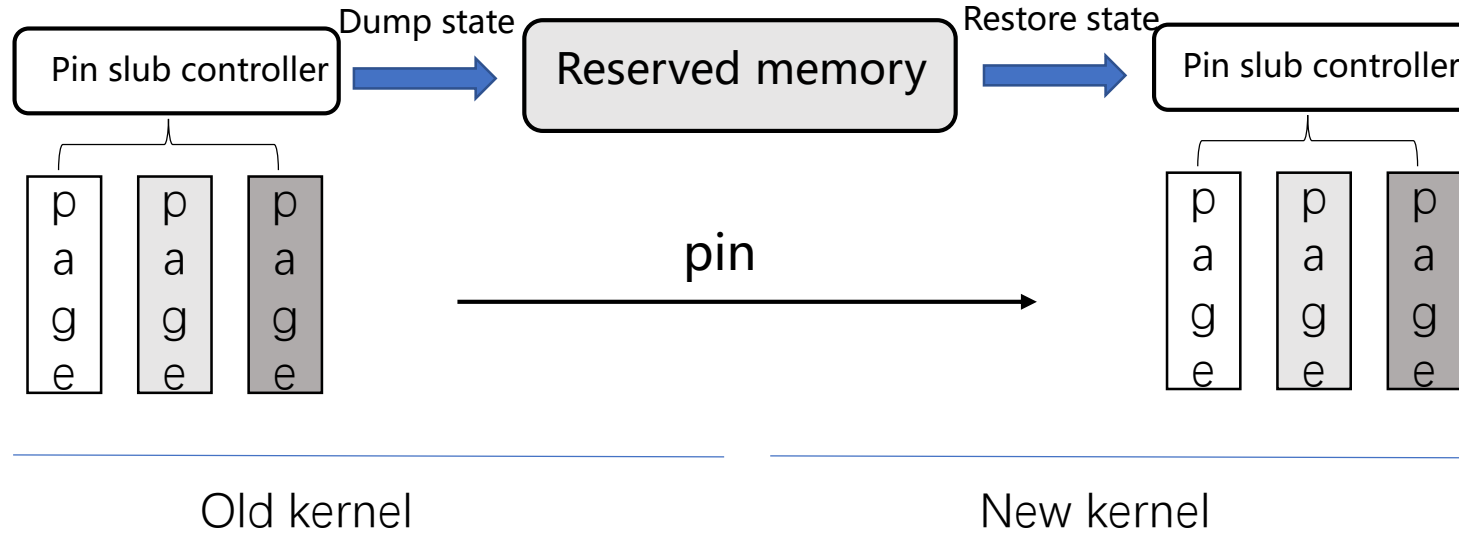




## ▶ 3.2 keep user memory unchanged in new kernel



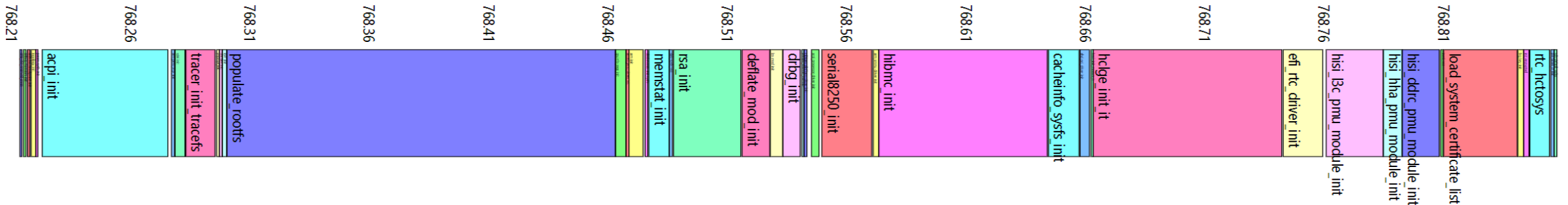
## ▶ 3.3 keep kernel memory unchanged in new kernel



\*Create pin slub controller to manage the old kernel pages which need to keep constant While booting the new kernel.

## ► 04 Kernel Fast Boot

- CPU currently booting
- Pre decompress kernel & initrd
- Memory Defer initialization
- Deferring device driver probe
- Deferring device initialization
- Avoid unless device scanning
- NIC PHY issue



## ▶ 4.1 CPU Parallel booting

- Current CPU sequence boot, each ARM core boot will elapse 0.036s.

```
[ 2.796042] Detected VIPT I-cache
[ 2.796078] GICv3: CPU81: found re
[ 2.796097] GICv3: CPU81: using al
[ 2.796155] CPU81: Booted secondar
[ 2.836491] Detected VIPT I-cache
[ 2.836527] GICv3: CPU82: found re
[ 2.836546] GICv3: CPU82: using al
[ 2.836603] CPU82: Booted secondar
[ 2.876940] Detected VIPT I-cache
[ 2.876977] GICv3: CPU83: found re
[ 2.876996] GICv3: CPU83: using al
[ 2.877053] CPU83: Booted secondar
[ 2.917566] Detected VIPT I-cache
```

- Modify CPU booting. Parallel CPU initiation will shorter the duration (0.0004s).

```
[ 0.049353] Detected VIPT I-cach
[ 0.049358] GICv3: CPU6: found
[ 0.049364] GICv3: CPU6: using
[ 0.049396] CPU6: Booted seconda
[ 0.049502] park_text 0xffff0000
[ 0.049510] Write cpu 7 entry 0:
[ 0.049681] Detected VIPT I-cach
[ 0.049687] GICv3: CPU7: found
[ 0.049692] GICv3: CPU7: using
[ 0.049723] CPU7: Booted seconda
[ 0.049827] park_text 0xffff0000
[ 0.049835] Write cpu 8 entry 0:
[ 0.050011] Detected VIPT I-cach
[ 0.050018] GICv3: CPU8: found
```

## ▶ 4.2 Pre decompress kernel & initrd

- Decompress bzImage to Image(vmlinux) before kexec
- Enhance kexec-tools to support for original kernel image
- Decompress initramfs.gz before kexec
- Unpack initramfs to ramfs in RAM before kexec

## ► 4.3 Memory Defer initialization

mm: parallelize deferred\_init\_memmap()

```
mm: parallelize deferred_init_memmap()
```

Deferred struct page init is a significant bottleneck in kernel boot. Optimizing it maximizes availability for large-memory systems and allows spinning up short-lived VMs as needed without having to leave them running. It also benefits bare metal machines hosting VMs that are sensitive to downtime. In projects such as VMM Fast Restart[1], where guest state is preserved across kexec reboot, it helps prevent application and network timeouts in the guests.

### Multithread to take full advantage of system memory bandwidth.

Intel(R) Xeon(R) Platinum 8167M CPU @ 2.00GHz (Skylake, bare metal)  
2 nodes \* 26 cores \* 2 threads = 104 CPUs  
384G/node = 768G memory

Intel(R) Xeon(R) CPU E5-2699 v3 @ 2.30GHz (Haswell, bare metal)  
2 nodes \* 18 cores \* 2 threads = 72 CPUs  
128G/node = 256G memory

kernel boot				deferred init			kernel boot				deferred init		
node%	(thr)	speedup	time_ms (stdev)	speedup	time_ms (stdev)	node%	(thr)	speedup	time_ms (stdev)	speedup	time_ms (stdev)		
	( 0)	---	4089.7 ( 8.1)	---	1785.7 ( 7.6)		( 0)	---	1680.0 ( 4.6)	---	627.0 ( 4.0)		
2%	( 1)	1.7%	4019.3 ( 1.5)	3.8%	1717.7 ( 11.8)	3%	( 1)	0.3%	1675.7 ( 4.5)	-0.2%	628.0 ( 3.6)		
12%	( 6)	34.9%	2662.7 ( 2.9)	79.9%	359.3 ( 0.6)	11%	( 4)	25.6%	1250.7 ( 2.1)	67.9%	201.0 ( 0.0)		
25%	( 13)	39.9%	2459.0 ( 3.6)	91.2%	157.0 ( 0.0)	25%	( 9)	30.7%	1164.0 ( 17.3)	81.8%	114.3 ( 17.7)		
37%	( 19)	39.2%	2485.0 ( 29.7)	90.4%	172.0 ( 28.6)	36%	( 13)	31.4%	1152.7 ( 10.8)	84.0%	100.3 ( 17.9)		
50%	( 26)	39.3%	2482.7 ( 25.7)	90.3%	173.7 ( 30.0)	50%	( 18)	31.5%	1150.7 ( 9.3)	83.9%	101.0 ( 14.1)		
75%	( 39)	39.0%	2495.7 ( 5.5)	89.4%	190.0 ( 1.0)	75%	( 27)	31.7%	1148.0 ( 5.6)	84.5%	97.3 ( 6.4)		
100%	( 52)	40.2%	2443.7 ( 3.8)	92.3%	138.0 ( 1.0)	100%	( 36)	32.0%	1142.3 ( 4.0)	85.6%	90.0 ( 1.0)		



## ▶ 4.4 Deferring device driver probe

drivercore: add driver probe deferral mechanism

```
drivercore: Add driver probe deferral mechanism
```

```
Allow drivers to report at probe time that they cannot get all the resources required by the device, and should be retried at a later time.
```

```
module.async_probe [KNL]
```

```
Enable asynchronous probe on this module.
```

```
driver_async_probe= [KNL]
```

```
List of driver names to be probed asynchronously.
```

```
Format: <driver_name1>, <driver_name2>...
```

```
deferred_probe_timeout=
```

```
[KNL] Debugging option to set a timeout in seconds for deferred probe to give up waiting on dependencies to probe. Only specific dependencies (subsystems or drivers) that have opted in will be ignored. A timeout of 0 will timeout at the end of initcalls. This option will also dump out devices still on the deferred probe list after retrying.
```



## ▶ 4.5 Deferring device initialization

- Deferring device driver, only defer the driver for the same, but sometime we only want to defer some device probe, for example:

```
0.860672] hns3 0000:7d:00.0: The firmware version is 1.8.12.3
0.940779] hns3 0000:7d:00.0: hclge driver initialization finished.
0.941897] hns3 0000:7d:00.1: The firmware version is 1.8.12.3
1.020563] hns3 0000:7d:00.1: hclge driver initialization finished.
1.021641] hns3 0000:7d:00.2: The firmware version is 1.8.12.3
1.100787] hns3 0000:7d:00.2: hclge driver initialization finished.
1.101882] hns3 0000:7d:00.3: The firmware version is 1.8.12.3
1.180567] hns3 0000:7d:00.3: hclge driver initialization finished.
```

- Deferring device initialization when the device add system, we only active one hns3 NIC, other will be defer.

## ▶ 4.6 Avoid unless device scanning

Some device no change, no rescan the device. We store some initialize information into reserve memory, use the information when reinitialize the device.

- PCI Scan information
- SATA disk scan information
- .....

## ▶ 4.7 NIC PHY issue

- When NIC reset PHY and reload firmware, the NIC transports no package input system until 3-5s later. No network, the downtime of the application will be poor.

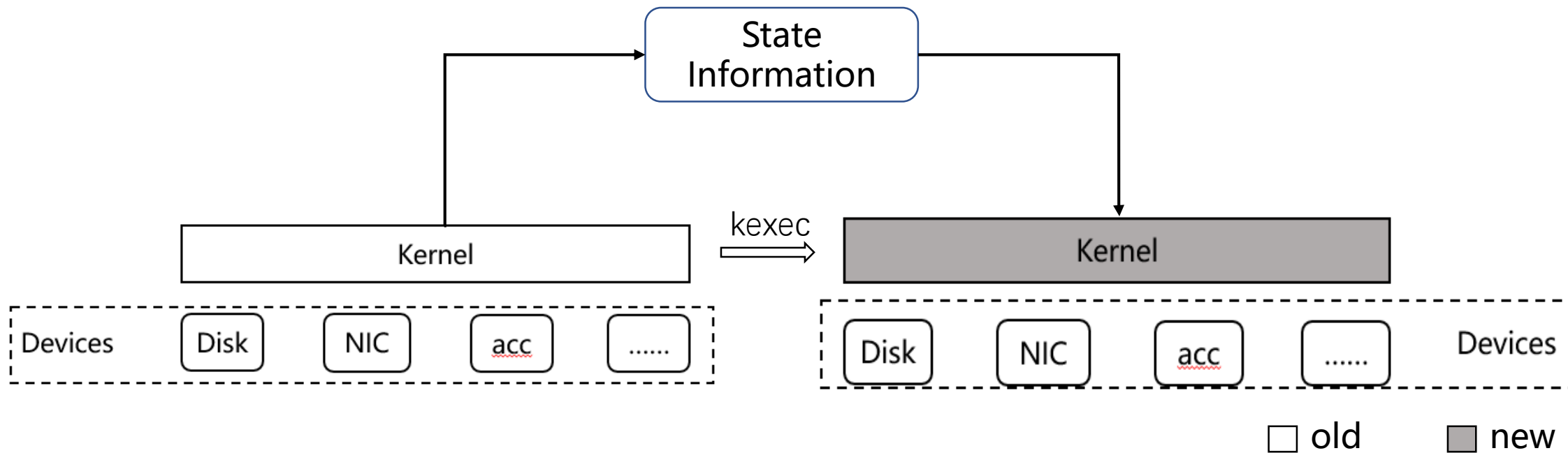
```
Link detected: yes 16:18:42
Link detected: no 16:18:43
Link detected: no 16:18:44
Link detected: no 16:18:45
Link detected: no 16:18:46
Link detected: no 16:18:47
Link detected: no 16:18:48
Link detected: no 16:18:49
Link detected: no 16:18:50
Link detected: no 16:18:51
Link detected: no 16:18:52
Link detected: no 16:18:53
Link detected: yes 16:18:54
```

- Modify the driver of the NIC, no reset PHY and reload firmware when quick reboot.

```
[ 4263.380716] icmp: icmp_echo: id=52763, seq=19903
[ 4263.484684] icmp: icmp_echo: id=52763, seq=19904
[ 4263.588715] icmp: icmp_echo: id=52763, seq=19905
[ 4263.692706] icmp: icmp_echo: id=52763, seq=19906
[ 4263.800706] icmp: icmp_echo: id=52763, seq=19907
[ 4263.904708] icmp: icmp_echo: id=52763, seq=19908
[ 4264.008702] icmp: icmp_echo: id=52763, seq=19909
[ 4264.116695] icmp: icmp_echo: id=52763, seq=19910
[ 4264.224690] icmp: icmp_echo: id=52763, seq=19911
[ 4264.328679] icmp: icmp_echo: id=52763, seq=19912
```

\* Should sure the state correct, as no reset PHY

## ► 05 Keep Device State



State information: reserve memory to save device state information

\* Caution: must concurrent device state and kernel driver information

# ▶ 06 Demo -- Benchmark

The image shows a screenshot of a terminal window and a Prometheus Redis Dashboard. The terminal window, titled "10.44.142.93 - Xshell 6", displays two sessions. The left session is connected to "Euler:./ #". The right session is connected to "[root@localhost redis\_demo]#". The Prometheus Redis Dashboard, titled "Redis Dashboard for Prometheus Redis Exporter 1.x", shows the instance "10.44.142.93:9121". The dashboard displays a line graph titled "Commands / Ops per Sec" with a y-axis from 0 to 1.0 and an x-axis from 15:37:20 to 15:38:18. The graph shows a very low, stable value near 0.0 throughout the time period.

## ▶ 07 Todo

- Open Source to openEuler
- Replace CRIU with kernel module
- Standard Device State Process
- Faster Kernel boot
- Multi Kernel parallel initialization

# Reference

- 1 [https://criu.org/Seamless\\_kernel\\_upgrade](https://criu.org/Seamless_kernel_upgrade)
- 2 [https://www.usenix.org/system/files/conference/atc16/atc16\\_paper-kashyap.pdf](https://www.usenix.org/system/files/conference/atc16/atc16_paper-kashyap.pdf)
- 3 <https://kvmforum2019.sched.com/event/TmvJ/seamless-cloud-system-upgrade-with-vmm-fast-restart-jason-zeng-intel>
- 4 [https://www.linuxplumbersconf.org/event/4/contributions/281/attachments/216/617/LPC\\_2019\\_kernel\\_fastboot\\_on\\_the\\_way.pdf](https://www.linuxplumbersconf.org/event/4/contributions/281/attachments/216/617/LPC_2019_kernel_fastboot_on_the_way.pdf)
- 5 <http://dmtcp.sourceforge.net/papers/dmtcp.pdf>



# Q&A

## 欢迎关注

新浪微博: openEuler社区

Twitter: openEuler

B站: openEuler

微信公众号: openEuler

### 社区网站



### 代码托管平台



### 微信交流群



添加小助手微信号“openeuler123”拉你进群



THANKS