# Using kexec to speed-up reboot

Juhyung Park

✉ arter97@dgist.ac.kr
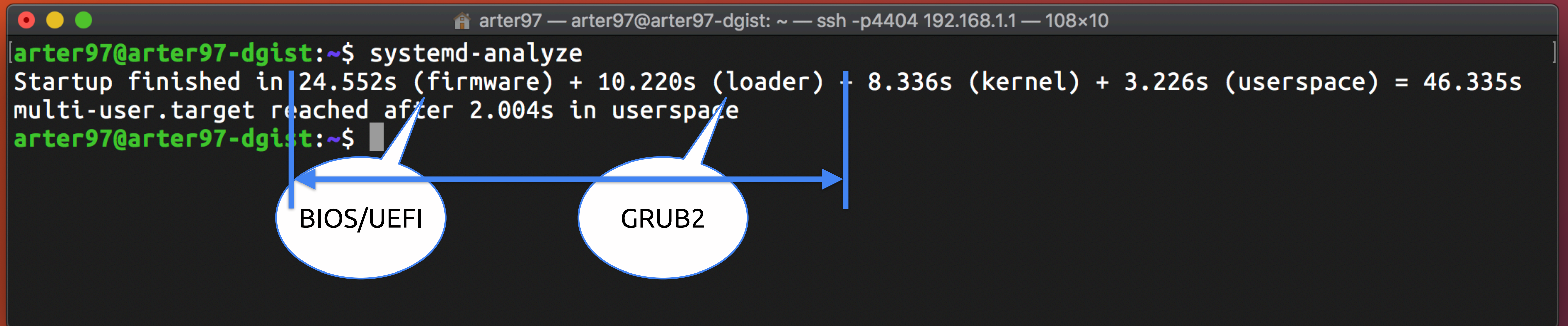
⌂ arter97

🐦 @arter97

# Standard practices to reduce boot time

```
arter97 — arter97@arter97-dgist: ~ — ssh -p4404 192.168.1.1 — 80×13

[arter97@arter97-dgist:~$ systemd-analyze blame
12.421s apt-daily.service
 1.659s nfs-server.service
 1.595s fstrim.service
 1.536s apt-daily-upgrade.service
 1.180s fwupd-refresh.service
  667ms netfilter-persistent.service
  454ms libvirtd.service
  347ms dev-nvme0n1p3.device
  223ms home-arter97-android.mount
  177ms ua-messaging.service
  159ms man-db.service
  144ms proc-fs-nfsd.mount
```

```
$ systemd-analyze blame
$ systemd-analyze plot > plot.svg
```
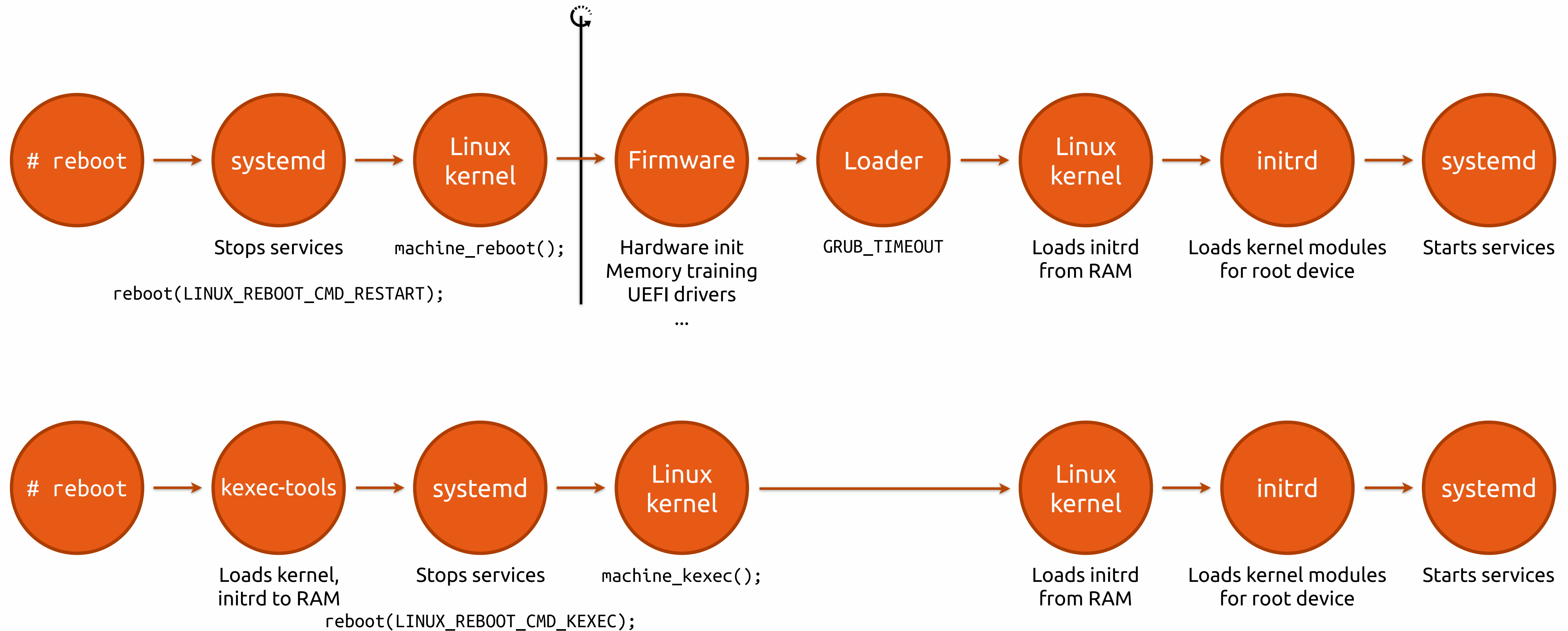
ubuCon ASIA 2021

# Reboots are still slow



*Firmware and loader accounts 75%*

## Things get worse for workstations/servers

# Enter kexec

kexec can bypass the firmware and loader for reboots

# Regular reboots vs. kexec

# How to use kexec on Ubuntu (1/2)

- `apt install kexec-tools`

- `apt remove finalrd`

  - With finalrd, systemd fails to find kexec binary and performs a full reboot instead

  - ubuntu-server installs mdadm, which installs finalrd - mdadm still works without finalrd

  - This was fixed (workaround) in systemd v246 - Ubuntu 20.04 ships with systemd v245 :(

  - … if you really need finalrd, add a new rule to /usr/share/finalrd/kexec.finalrd

# How to use kexec on Ubuntu (2/2)

- Custom reboot helper

```bash
#!/bin/bash

# Store at /usr/local/sbin/reboot with 755 permission

echo ""
echo "Using kexec for faster reboot."
echo ""
echo "If you want to perform a full reboot,"
echo "use 'systemctl reboot' instead."
echo ""

sudo kexec \
    -l $(ls /boot/vmlinuz-* | sort -V | tail -n1) \
    --initrd=$(ls /boot/vmlinuz-* | sort -V | tail -n1 | sed s/vmlinuz/initrd.img/g) \
    --reuse-cmdline && sudo systemctl kexec
```

- This script loads the latest kernel and initrd to the RAM and calls systemd for kexec reboot, reusing cmdline

- https://gist.github.com/arter97

ubuCon ASIA 2021

# How effective is kexec? - Testing how long it takes for the network to get back online



▲ kexec reboots in 8s

▲ Regular reboot takes 36s

# Limitations

- efifb may not work
  - Display output may not work until the new kernel loads the graphics drivers

- BIOS implementations may break kexec
  - Especially on AMD platforms, updating your BIOS may help

- Some hardware changes may still require traditional reboots
  - e.g., Newly connected HDDs/SSDs may not appear on kexec reboots

- So, test kexec when you have physical access (or IPMI) before deployment!

# Extra

- qboot is also extremely helpful for virtualized kernel development

  - New login shell appears within 3s under QEMU

  - https://github.com/bonzini/qboot

  - https://twitter.com/arter97/status/1295671273784147969

- kexec in action

  - https://youtu.be/60Gh1NnK0MA