



**Hewlett Packard**  
Enterprise

# **Rev Up Compute Node Reboots** **2x – 5x faster**



Dennis Walker, HPE  
Paul Selwood, MET Office UK

May 4, 2025



## Agenda

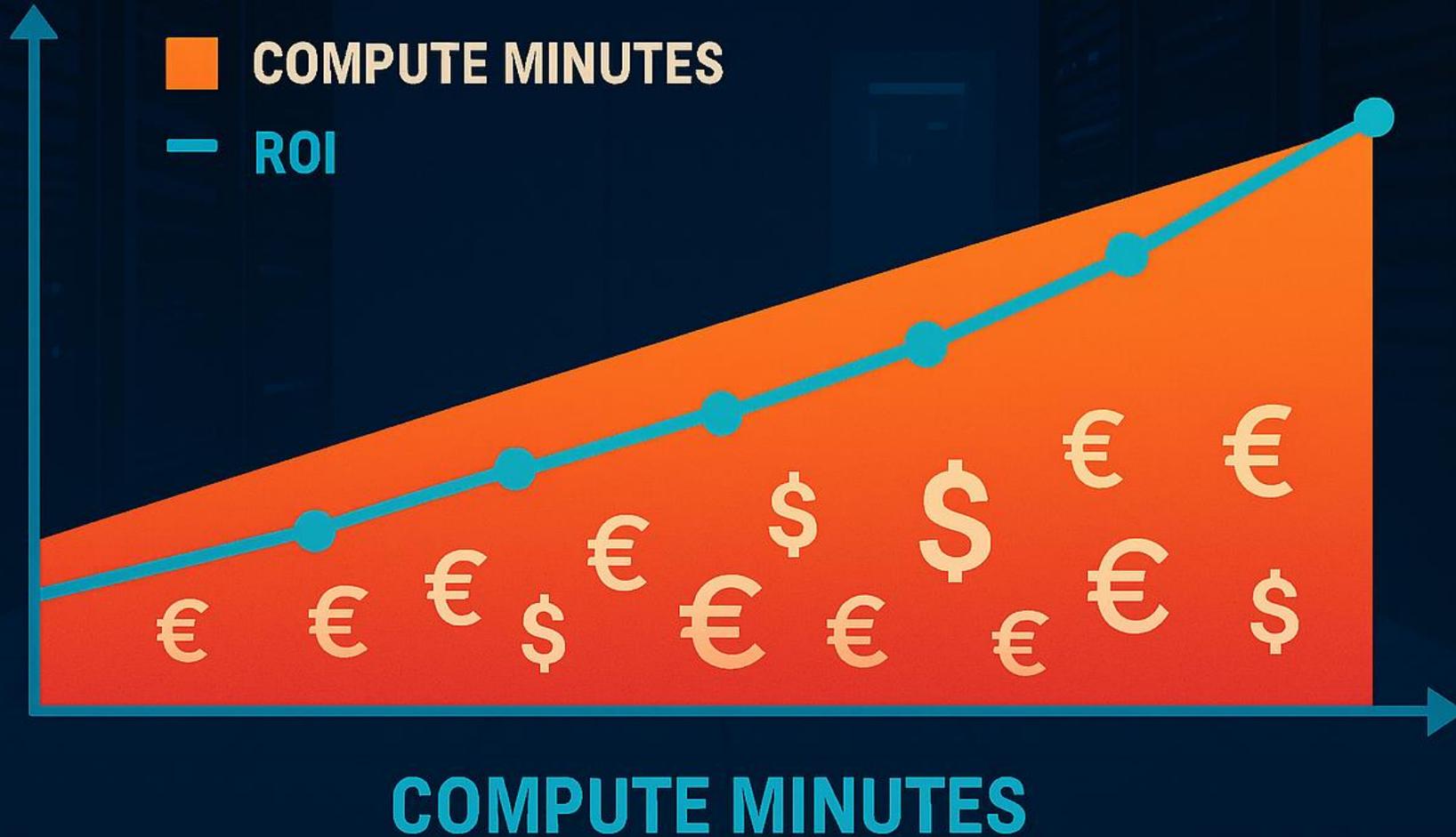
---

- Why optimize reboots?
- Project goals
- Testing & analysis methodology
- Findings
- Remediation
- Results



## The Cost of Downtime (during reboot)

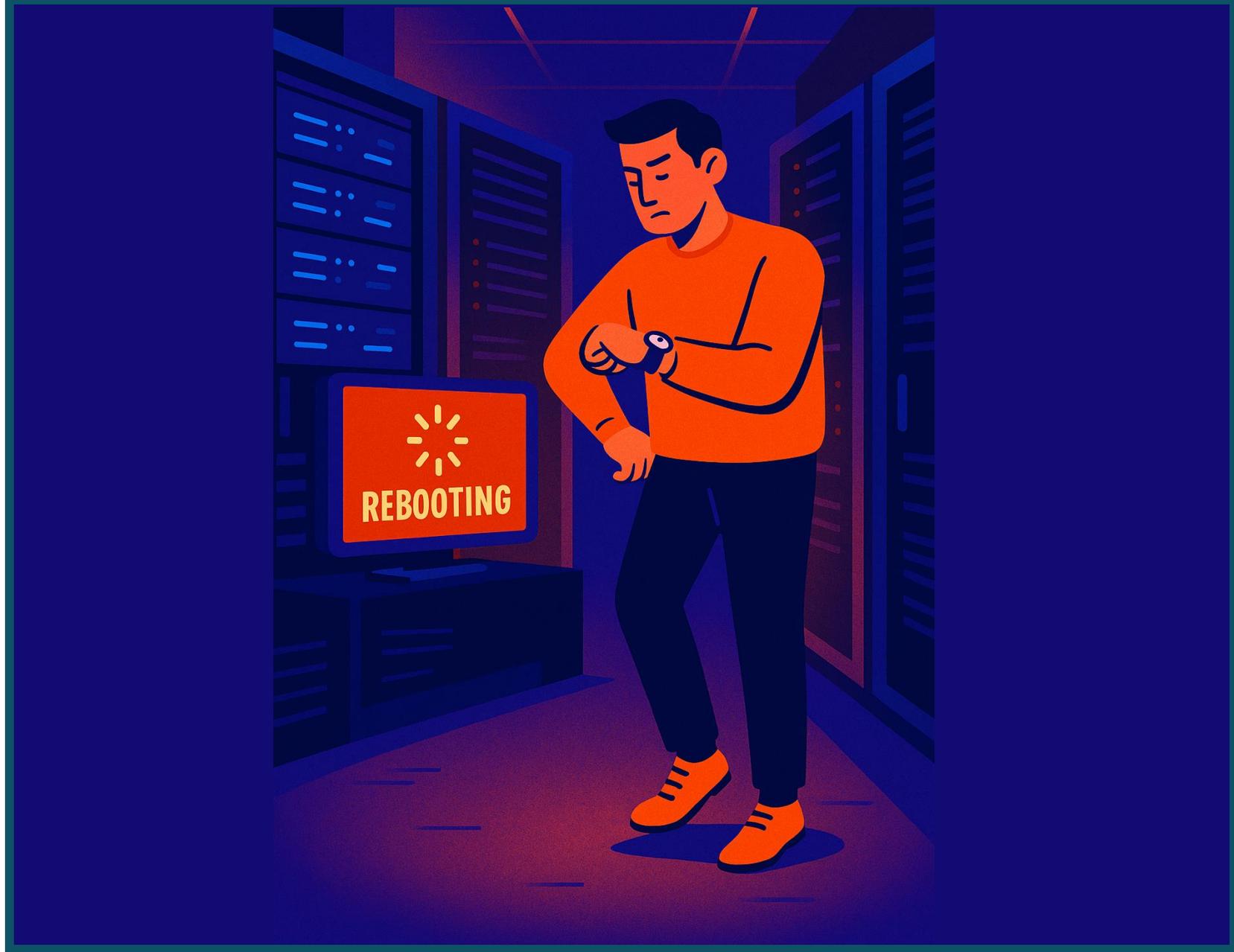
- What is downtime?  
Minutes without job execution availability
- ROI amortization schedule favors early life; systems with new scale and efficiency
- Reboot time multiplies OPEX; rescheduling, iterations, incident response, etc.



## Reboot Use Cases

---

- Update hardware/software while mitigating downtime
- Job specifies another image
- Incident response



## Project Goals

---

Reduce Node Reboot Time

Original contract: 10 minutes

Amended contract: 15 minutes

-----

Actual time as of 2023: 35.5 minutes

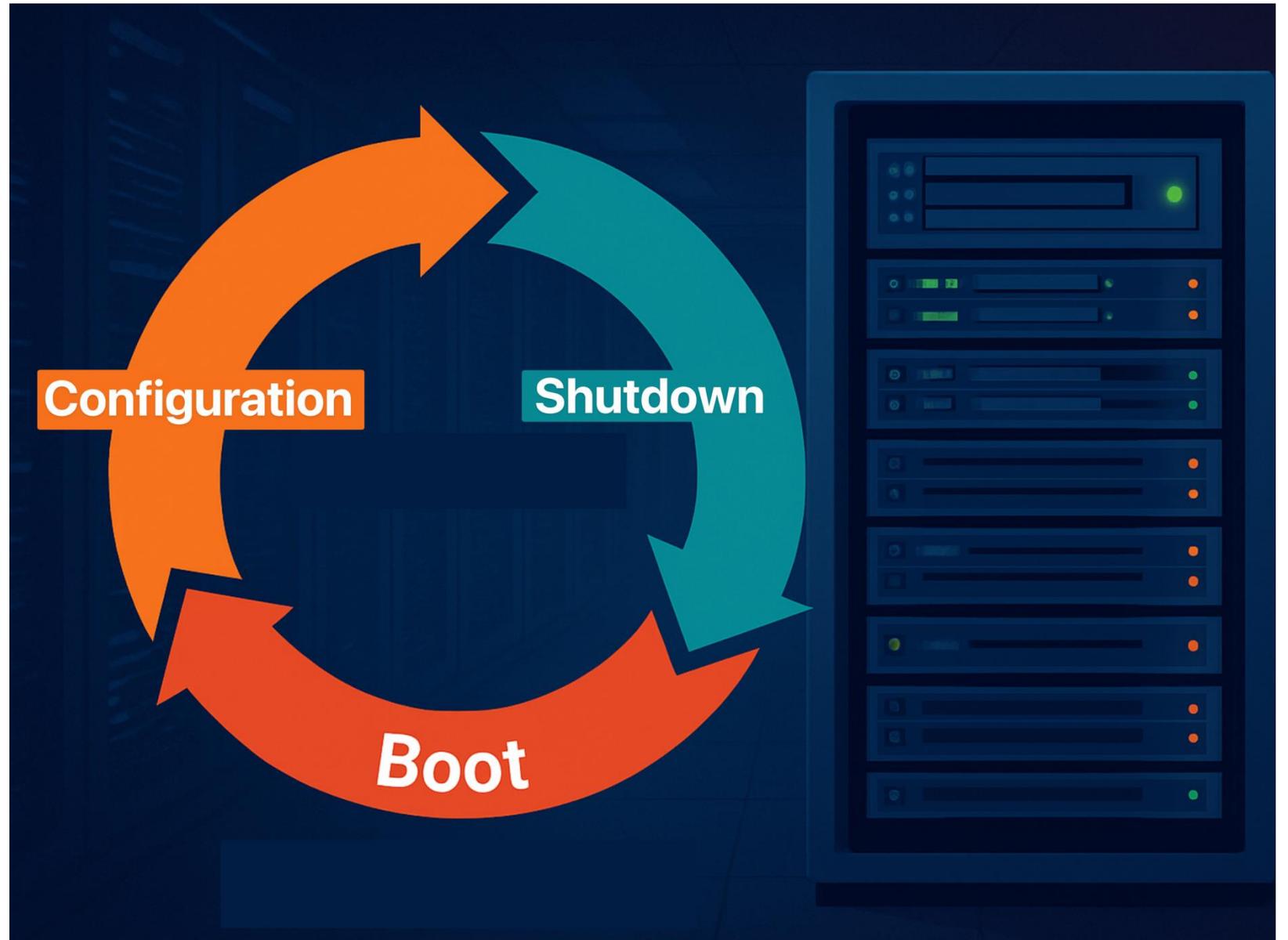
Actual time as of 2024: (end of deck)

# Rev Up Compute Node Reboots: 2x to 5x Faster



## Anatomy of a Reboot

- Shutdown
  - **BOS session (start)**, job completes, data sync, services stop, power down
- Boot
  - POST, PXE, Dracut
- Configuration
  - CFS Ansible, health check, **node ready for jobs (stop)**



## Fully Automated Approach

30-60 Minute Iterations

- Branch VCS repositories and push
- Create/update CFS configs
- Build images
- Boot 2-4 nodes
- Capture and graph results



## Measurement Framework

- Pulls logs and data of a reboot together
- Parses statistics, exports to csv
- Renders dashboards (Jupyter notebook)



## Console Log Boot Stat Parsing

- Fetches console logs
- Parses start and stop strings
- Calculates duration between

```
echo "BOS_Session_Id,Test_Type,Stage,Minutes" >"$this_csv_file"
{
  get_shutdown_time
  get_duration_between "BIOS" \
    "PeiCore.Entry" \
    "NBP filename is ipxe.efi"
  get_duration_between "DVS_generate_node_map" \
    "DVS: loaded dvsproc module." \
    "DVS: node map generated."
  get_duration_between "Dracut" \
    "wicked: net0: Request to acquire DHCPv4 lease with UUID" \
    "mount is: mount -t squashfs -o loop /tmp/cps/rootfs /new_root/lower"
  get_duration_between "iPXE" \
    "NBP filename is ipxe.efi" \
    "wicked: net0: Request to acquire DHCPv4 lease with UUID"
  get_duration_between "Login-Prompt" \
    "mount is: mount -t squashfs -o loop /tmp/cps/rootfs /new_root/lower" \
    " login:"
} >>"$this_csv_file"
```

Figure 3.1.2 - Resulting reboot stage metrics in minutes

1	BOS_Session_Id,Test_Type,Stage,Minutes
2	b13.11-new-window,New_Window,Shutdown,4.18
3	b13.11-new-window,New_Window,BIOS,1.81
4	b13.11-new-window,New_Window,DVS_generate_node_map,.16
5	b13.11-new-window,New_Window,Dracut,1.00
6	b13.11-new-window,New_Window,iPXE,1.68
7	b13.11-new-window,New_Window,Login-Prompt,1.00
8	b13.11-new-window,New_Window,Total_Preboot,5.50

## Ansible Log & Rollup Stats

- Fetches console logs, BOS session status, and PBS pro node status
- Fetches Ansible status (top right)
- Creates rollup csv (bottom right)

Figure 3.2.2 - Ansible playbook execution times in seconds

```
1 | BOS_Session_Id,Test_Type,Playbook,Repo,Seconds
2 | b13.11-new-antero,New_Antero,shs_cassini_install.yml,slingshot-host-software-
  | config-management,1.34
3 | b13.11-new-antero,New_Antero,cos-compute.yml,cos-config-management,78.43
4 | b13.11-new-antero,New_Antero,sma-ldms-compute.yml,sma-config-management,37.31
5 | b13.11-new-antero,New_Antero,cos-compute-last.yml,cos-config-management,6.85
6 | b13.11-new-antero,New_Antero,site.yml,custo-post-boot,20.86
7 | b13.11-new-antero,New_Antero,site.yml,mo-config-management,74.25
```

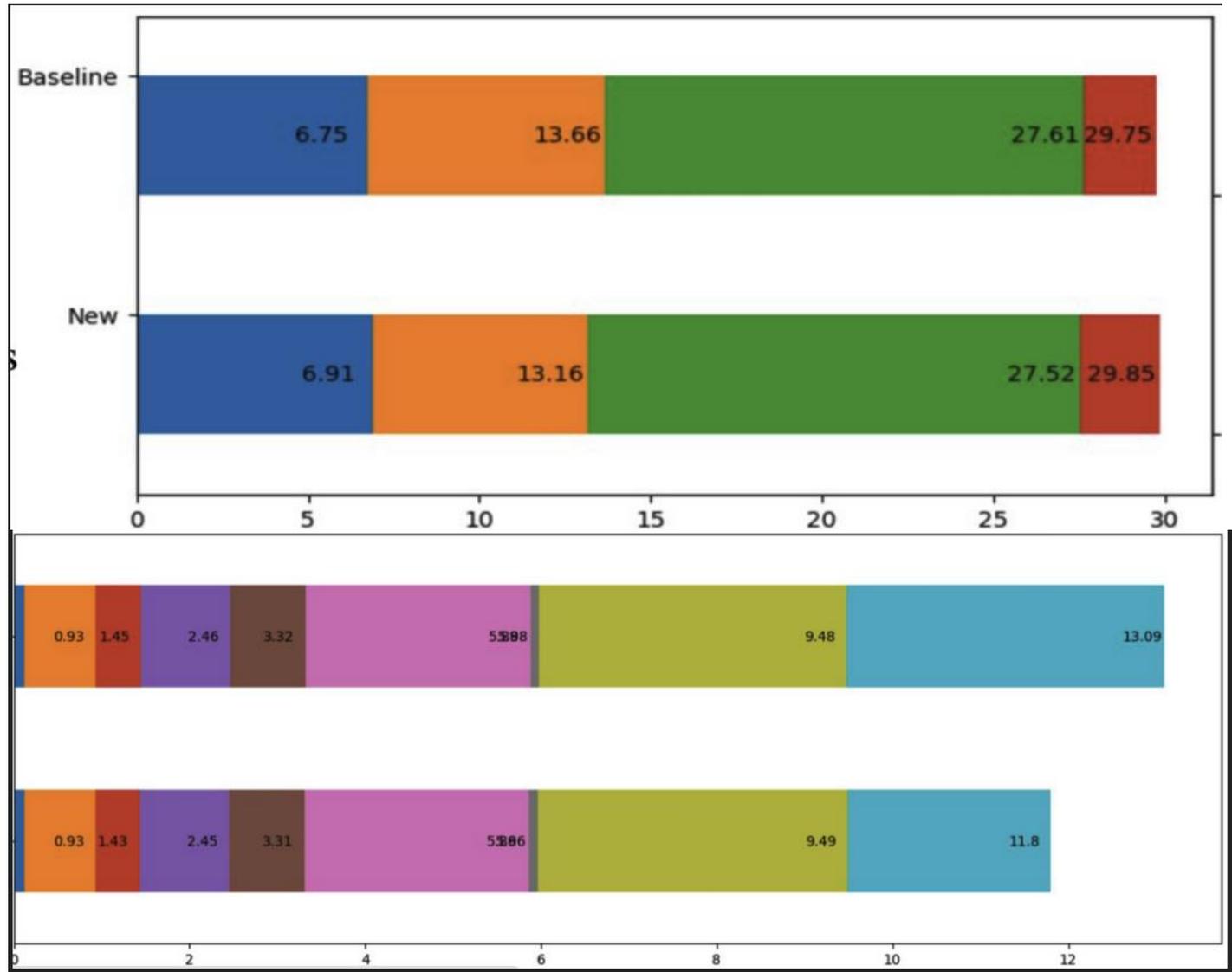
```
BOS_Session_Id      b13.11-new-windom
Test_Type           New_Windom
Start               2024-02-27T00:27:34
Stop                2024-02-27T00:43:03
Minutes             15.48
Preboot_Minutes     5.66
Num_CFS_Sessions    1
CFS_Sessions_Minutes ;5.21
CFS_Sessions_Start_Times ;2024-02-27T00:37:50
CFS_Sessions_Stop_Times ;2024-02-27T00:43:03
```

# Jupyter Notebook

- Visually displays breakdown of timing for analysis
- Compares before and after change for current iteration

Top right, a total reboot time

Bottom right, CFS Ansible breakdown



## Remediation

- Control for non-HPE activity
  - job termination, data sync, customer personalization
- Reduce boot time
  - System management tuning
- Reduce configuration time
  - Implement Systemd units for each CFS layer

Stage	Before – 12.2023
Shutdown	8.9
Boot	6.2
Configuration	20.4
Total	35.5

## System Management Tuning

- Update PXE script to select correct NIC for node type (left)
- Tune BOS for shorter shutdown tolerance (a control)
- Tune CFS to start more quickly (saves 30-60 seconds)

```
Content
---
bss.ipxe: "#!ipxe\n\nnecho Chaining to BSS; trying interfaces net2, net0, net1, net3,
net4, net5...\n\nset attempt:int32 0\nset maxattempts:int32 1024\nset sleepytime:int32
0\nset ceiling:int32 64\nset sp:hex 20 && set sp ${sp:string}\n\n\n:start\ninc
63 ||\niseq ${attempt} ${maxattempts} && goto debug_retry ||\nnecho ${manufacturer:hexhyp}
${product} \n\niseq ${product} HPE_CRAY_EX425 && goto net1 || \niseq ${product}
ProLiant${sp}DL385${sp}Gen10${sp}Plus${sp}v2 && goto net0 || \n\n:start2\niseq
```



## Systemd Local Ansible

During Image Build (1):

- Install Ansible
- Create Ansible inventory
- Rsync CFS layers into image
- Wrap in Systemd units
  - Specifying order where necessary (4)

Later, CFS personalization play checks  
Systemd status to maintain visibility (2)

Ad-hoc play for pushing updates (3)

```
.
├── files
│   ├── ansible_hosts
│   └── ansible_hosts.collab
├── tasks
│   ├── check_systemd.yml 2
│   ├── main.yml 1
│   └── update_and_restart_systemd_units.yml 3
├── templates
│   ├── personalization_layer.service.j2
│   └── personalization_layer_init.sh.j2
└── vars
    └── main.yml 4
```

# Systemd Local Ansible Results

- Parallelized Playbook Execution
- Earlier Execution
- From 20 minutes to 2 minutes

File: main.yml

```
---
- name: Install Ansible
  ansible.builtin.package:
    name: ansible
    state: present
- name: Make local playbooks directory
  ansible.builtin.file:
    path: '{{ playbooks_dir }}'
    state: directory
    mode: '0640'
- name: Make systemd init script directory
  ansible.builtin.file:
    path: '{{ personalization_init_dir }}'
    state: directory
    mode: '0640'
- name: Rsync all playbooks to IMS host
  shell: |
    set -eou pipefail
    IMS_HOST=$(grep "ansible_host" /inventory/hosts/01-cfs-generated.yaml | sed -rn "s|^s*.*:(.*)$|\1|p" | xargs)
    rsync -avz /inventory/ -e "ssh -i /etc/ansible/ssh/id_image" "$IMS_HOST": "{{ playbooks_dir }}"
  args:
    creates: '{{ playbooks_dir }}/layer0'
    delegate_to: localhost
- name: Remove build credentials and info
  ansible.builtin.file:
    state: absent
    path: '{{ item }}'
  with_items:
    - '{{ playbooks_dir }}/ssh'
    - '{{ playbooks_dir }}/hosts'
    - '{{ playbooks_dir }}/image_to_job.yaml'
    - '{{ playbooks_dir }}/ansible.cfg'
    - '{{ playbooks_dir }}/complete'
```

## Final Results

- Tested on 4 random nodes in 4 production environments
- Observed live by 4 organizations

• **10.4 Minutes**

Table 5.0.1 - Individual Node Reboot Times

System	Node 1	Node 2	Node 3	Node 4
Quad A	10.40	10.23	10.26	10.68
Quad B	12.53	12.30	11.11	10.93
Quad C	10.46	10.45	11.48	11.46
Quad D	11.65	12.28	10.00	10.40

Table 5.0.2 - Final Comparative Results

Stage	Before – 12.2023	After – 07.2024
Shutdown	8.9	2.3
Boot	6.2	5.5
Configuration	20.4	2.6
Total	35.5	10.4

# Thank you

---

[dennis.walker@hpe.com](mailto:dennis.walker@hpe.com), [paul.selwood@metoffice.gov.uk](mailto:paul.selwood@metoffice.gov.uk)

Images on slides 2-8 generated with prompts to ChatGPT CUG 2025 | © HPE

