# Preparing for a new Data Center:

# Automated Management of a 10'000 node bare-metal fleet

**Arne Wiebalck, Luca Atzori, Nikos Papakyprianou, Michał Piszczek, Maryna Savchenko**
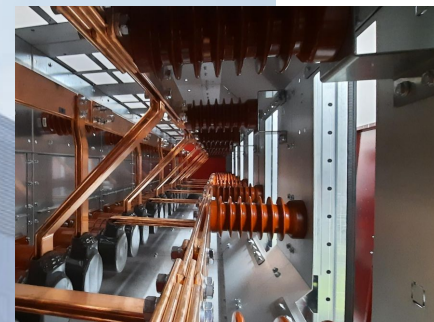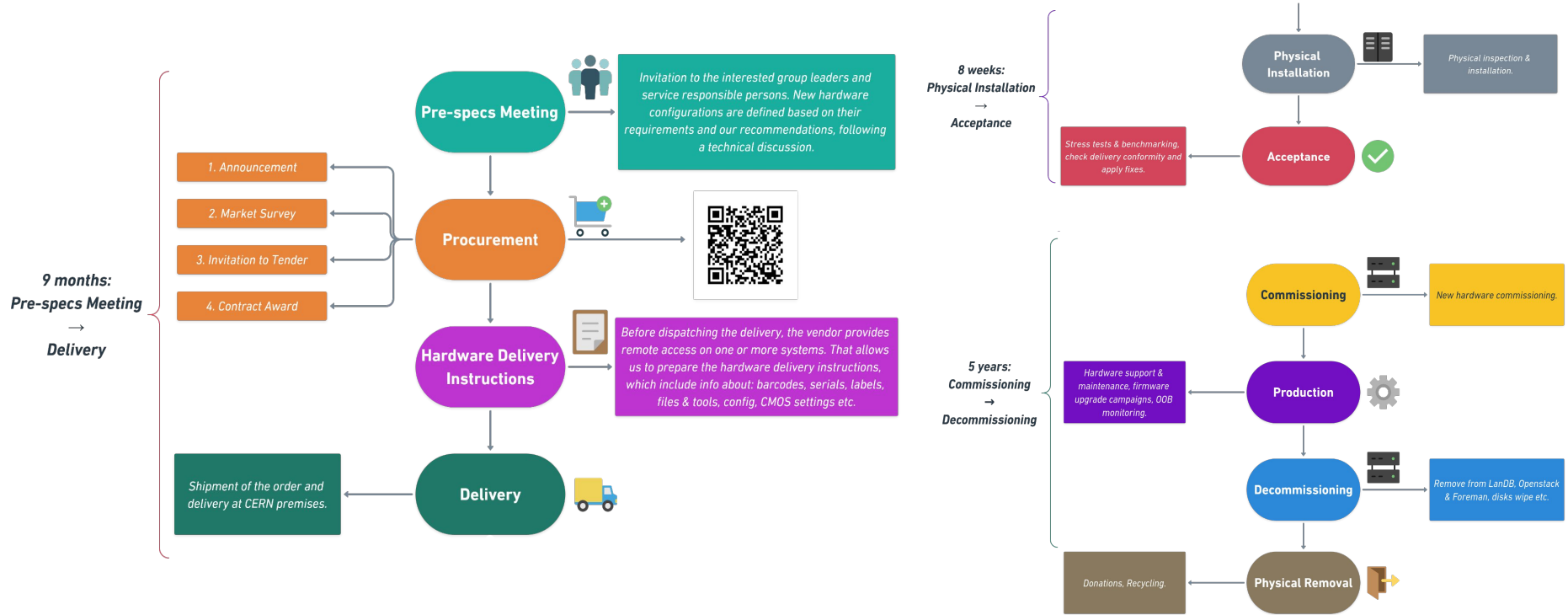
**CERN IT**

Norfolk, 9 May 2023

# The new CERN Prévessin Data Center

- **Three phases on three floors**
  - 4MW – 2nd floor
  - 8MW – 1st + 2nd floor
  - 12MW – ground floor + 1st + 2nd floor

- **Air-cooled racks with hot-aisle containment**

- **Two redundant power feeds: red and blue**
  - Red feed: 20% UPS coverage

# From Specification to Removal

# The Ironic Bare Metal API

➜ **Idea: Extend the cloud approach to bare metal servers!**

  ➢ Bare Metal offering to complement VMs and containers

  ➢ Provided via the same interface

  ➢ Simplification: workflows, approval, accounting, …

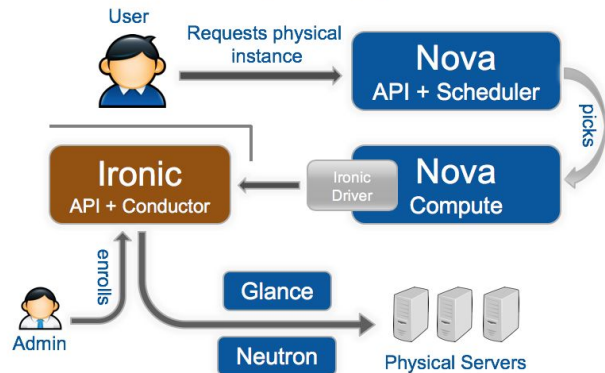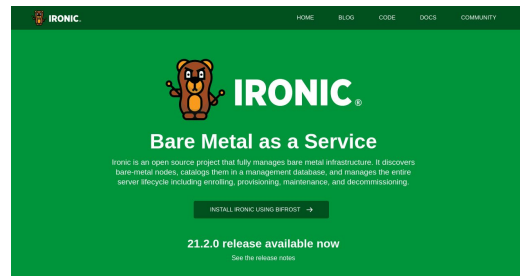➜ **API service to manage/interact with physical servers**

  ➢ Originally a provisioning driver in Nova

➜ **Can be used with OpenStack or stand-alone**

  ➢ ironicbaremetal.org

➜ **Leverages OSS standards and tooling**

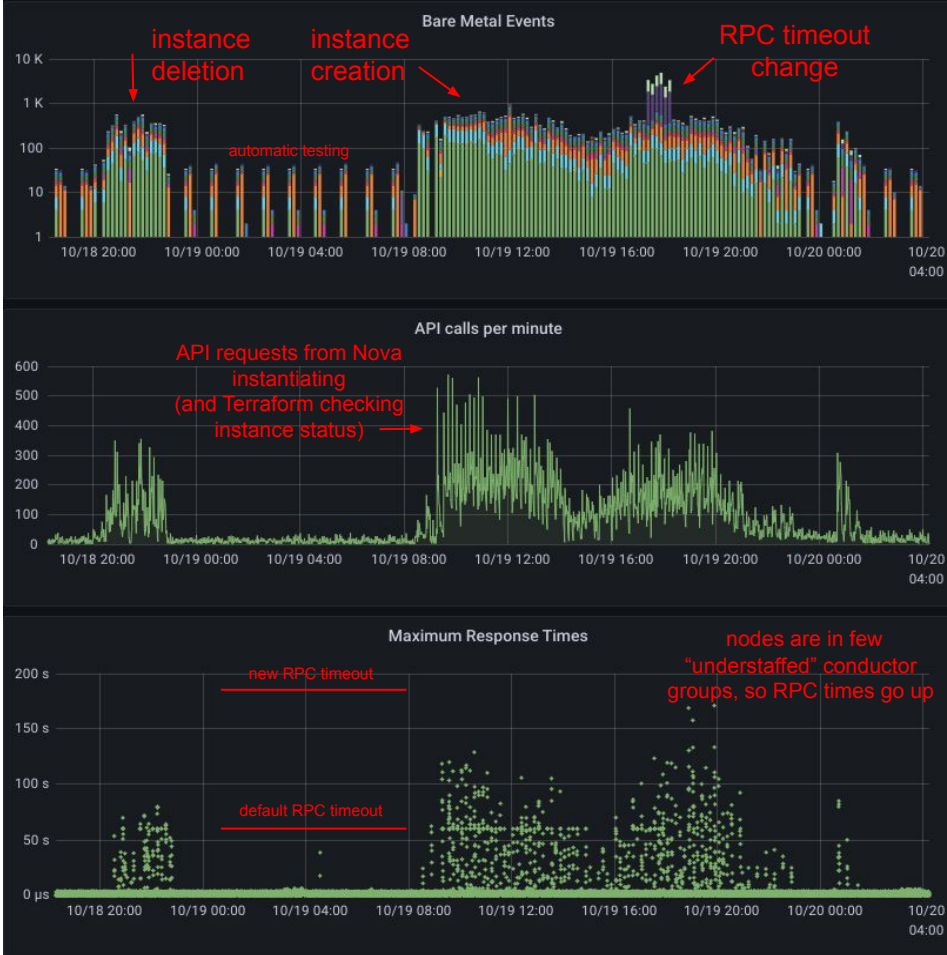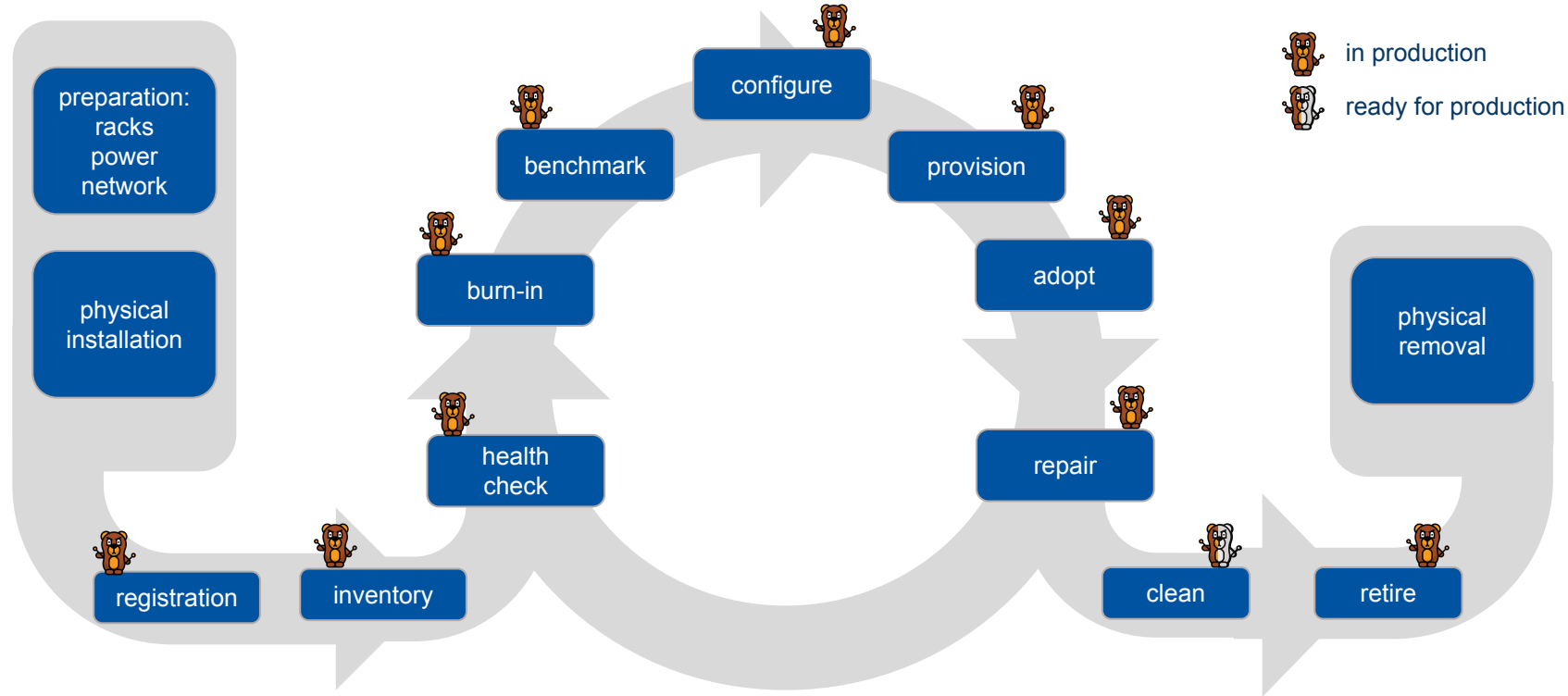  ➢ IPMI/Redfish, PXE, DHCP, … but allows for vendor plugins

# : Physical Batch

➤ **Conversion of virtual to physical batch**
  ➤ with the availability of a bare metal API, we revisited the virtualisation tax

➤ **3'775 hypervisors recreated as physical batch instances**
  ➤ done in multiple chunks over several months

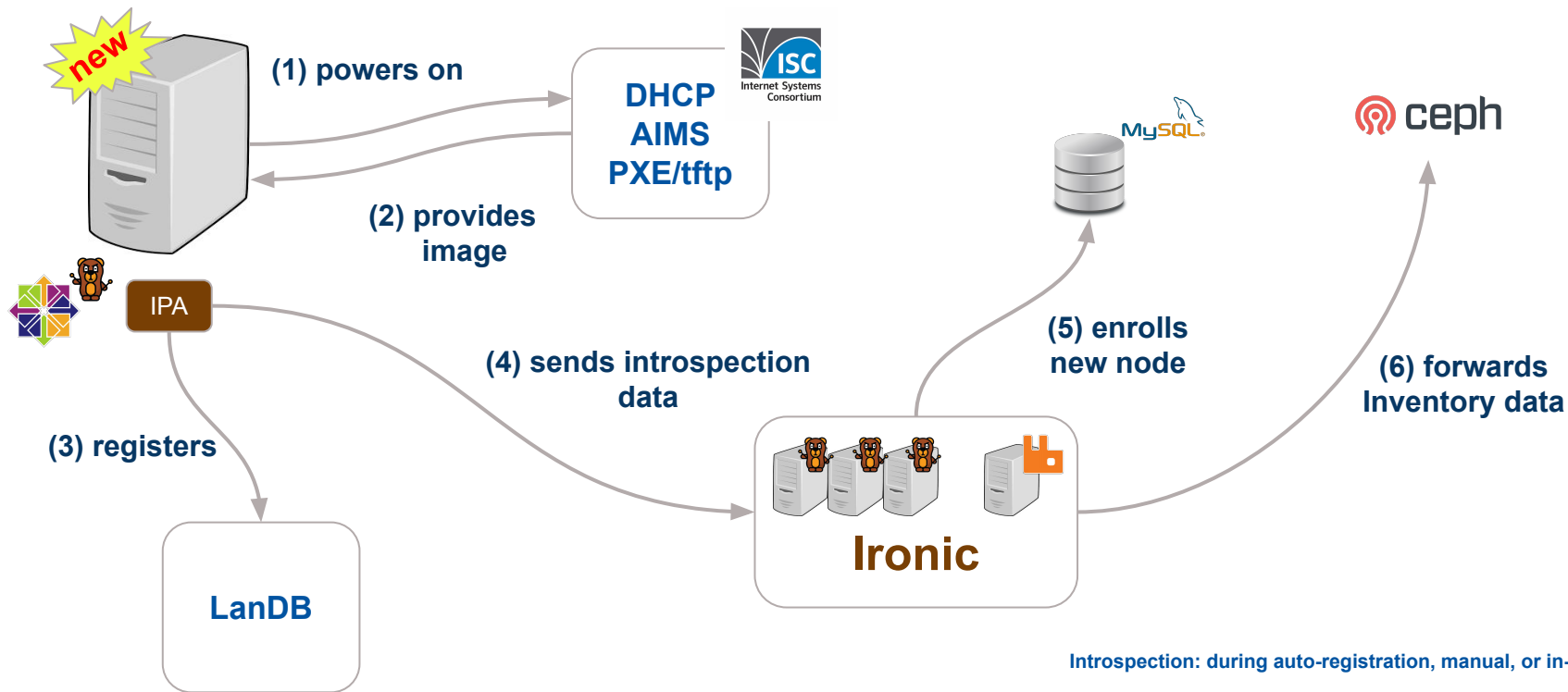➤ **Terraform as the 'Infrastructure-as-Code' tool to interface with OpenStack/Ironic**

HashiCorp
**Terraform**

**Bonus: 16'000 VMs less than one year ago … 10k+ IPv4 addresses free'd up.**

# Server Life-cycle Management

# Auto-Registration with Ironic



(1) powers on

DHCP
AIMS
PXE/tftp

ISC Internet Systems Consortium

MySQL

ceph

(2) provides image

IPA

(4) sends introspection data

(5) enrolls new node

(6) forwards Inventory data

(3) registers

LanDB

Ironic

Introspection: during auto-registration, manual, or in-band!
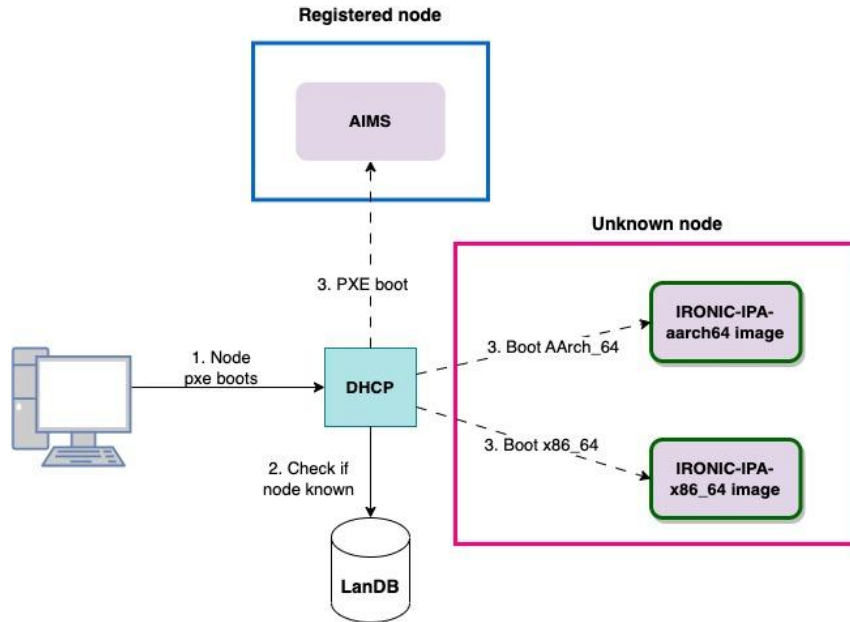
# Hardware Inventory



→ **Pool data from Ironic + Infor EAM**

→ **Combine with OpenDCIM**

→ **Provide info via CLI or webUI**

# Detour: Auto-registration (ARM)



- ➜ **IPA image is architecture dependent!**

- ➜ **Automatically build for x86 & ARM**

- ➜ **DHCP decides based on PXE data**

# Hardware Burn-in

→ **Provoke early failure via stress-test**
   ➢ CPU (stress-ng)
   ➢ Memory (stress-ng)
   ➢ Disk (fio)
   ➢ Network (fio)

→ **Integrated upstream**
   ➢ Released with Xena
   ➢ Implemented via cleaning steps

→ **Real-time log handling**



→ **Network requires pairing**
   ➢ Initial version: static pairs
   ➢ Works, but clumsy

→ **Dynamic pairing: distributed arbiter**
   ➢ OpenStack/tooz with ZooKeeper
   ➢ delivery and interface separation
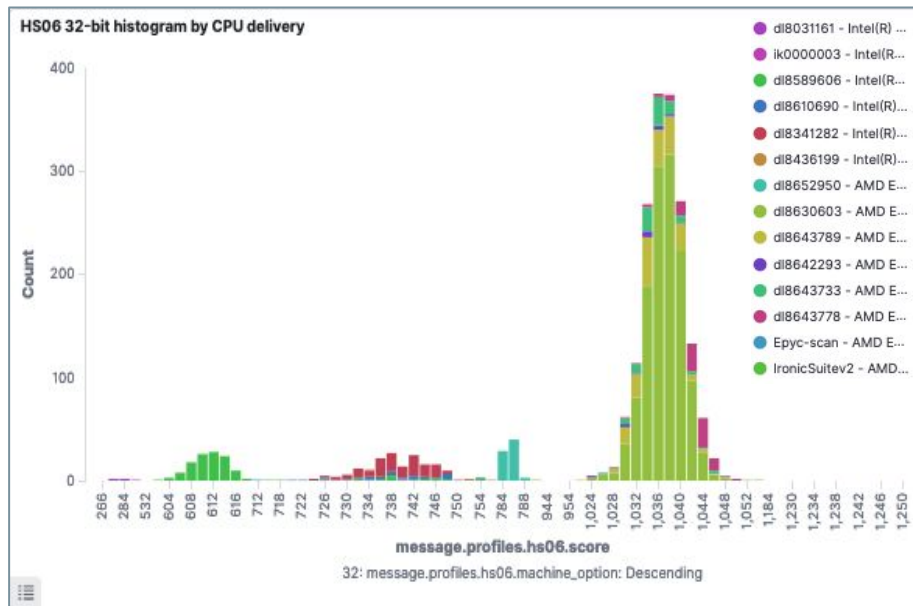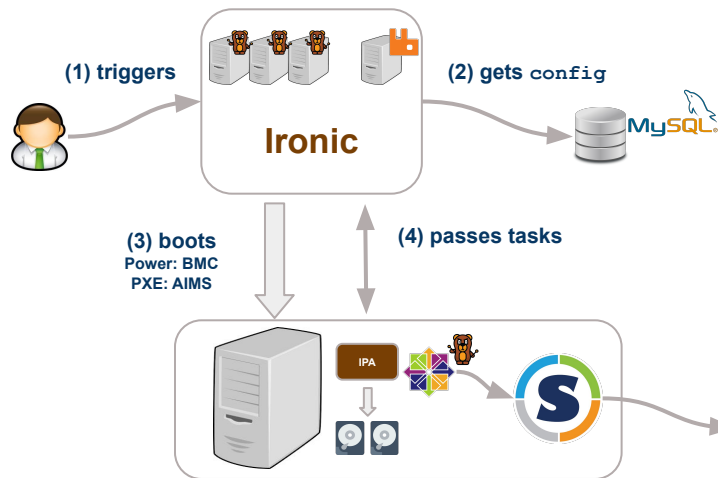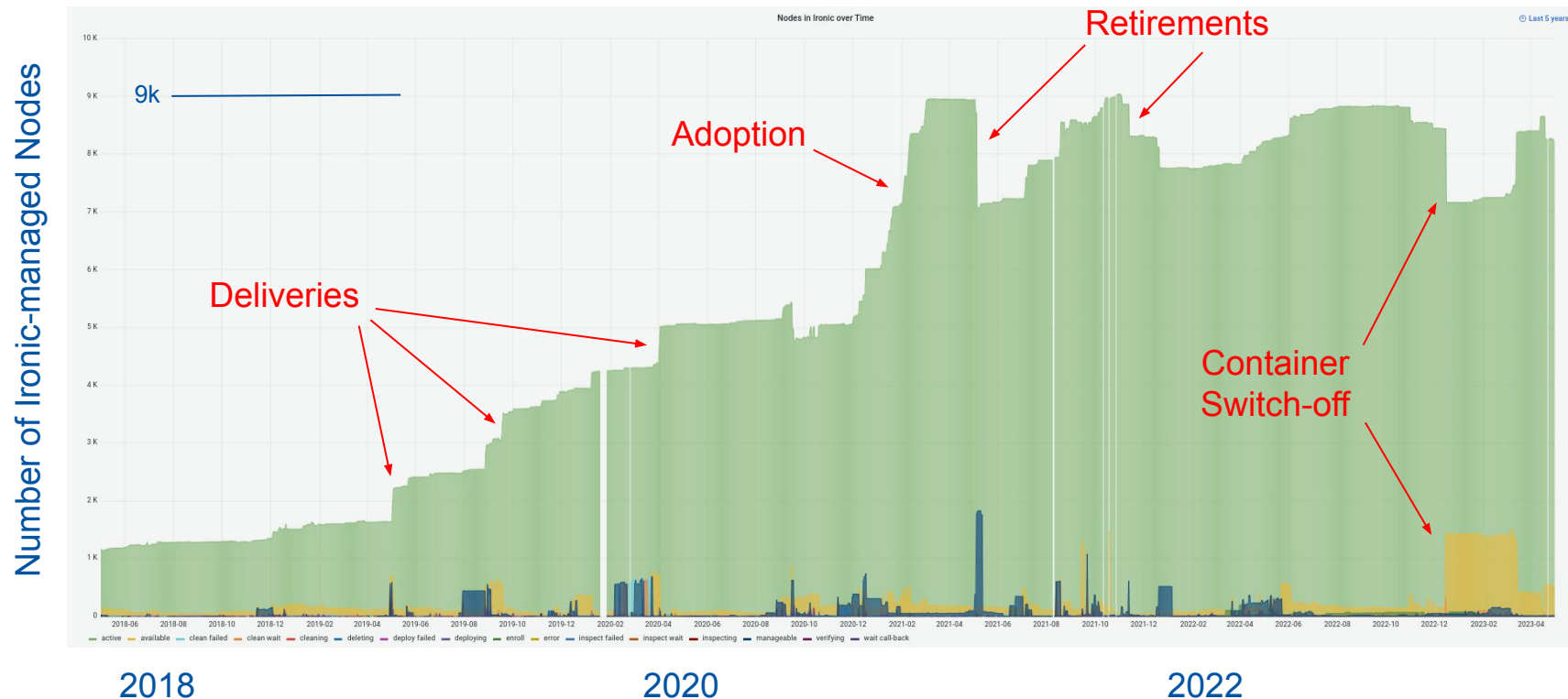   ➢ merged upstream, in-prod downstream

# Benchmarking

➜ **Based on a cleaning step**

- ➤ Downloads singularity image
- ➤ Executes configured benchmarks
- ➤ Sends results into OpenSearch



**HEP Spec06 32bit by delivery**

# Ironic Evolution over Time

# The GRUBsetta Stone



**loading initrd ...error: out of memory**

| Where it appears | Suspected reason | Resolution |
|---|---|---|
| Loading image with GRUB2 | Image too big. Hitting 4G memory limits | Reduce image size |

**loading initrd ...error: cannot allocate memory**

| Where it appears | Suspected reason | Resolution |
|---|---|---|
| Loading image with GRUB2 | Broken BIOS | Reinstall BIOS firmware and reload the BIOS config |

# Thanks!

Arne.Wiebalck@cern.ch