# How do we manage 40K machines in the CERN Computer centre

ZHECHKA TOTEVA (CERN/IT)

CERN – 17/07/2024

# Outline
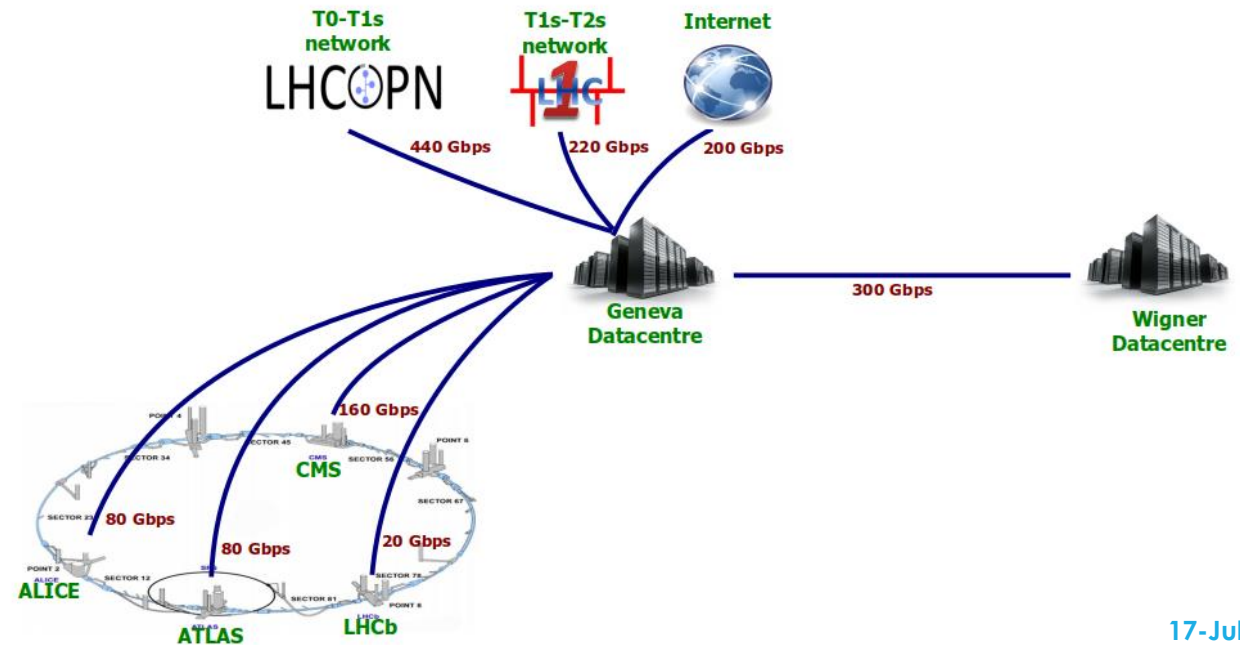
- CERN Computer Centre (CC) in numbers
- Overview of the CERN network and data storage
- Overview of electricity and cooling
- WLCG is couple of numbers
- Configuration management at CERN IT
- CERNMegabus@CERN
  - Architecture
  - Overview of major implemented use cases
  - CERN Computer Centre (CC) power cut management

# CERN Computer Centre (CC) in numbers

# Computing network

- 250 routers, 4100 switches, 1200 Wi-Fi points
- 35 000 km optical fibre (only ~5 000 less than the equator length)

- Wigner Data centre in Hungary
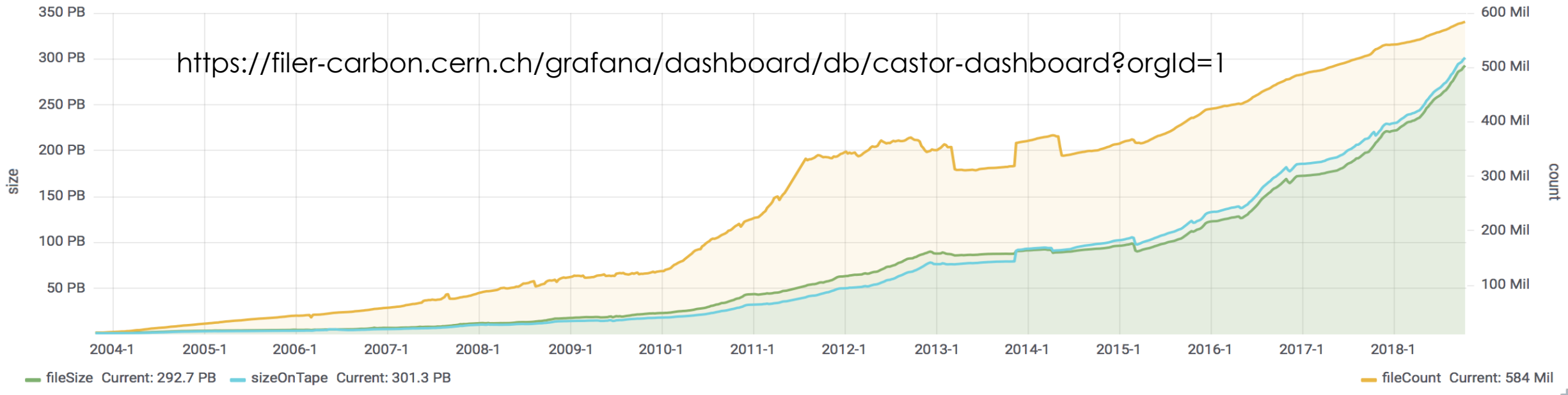  - 1200 km distance
  - with three 100 Gb/s

# Data storage

- 50 Pbytes / year from LHC

  - + 25 Pbytes / year from non-LHC experiments

- **<u>RECORD</u>**: August 2018:  13.8 Pbytes of data written on tape (of which 11.56 is LHC data)

  - More than 2 PB read/write daily

- Tape drives faster than disks; but slower in mounting (latency)

  - 90 K disk drives (of which 10-15% are SSD, providing less than 10% capacity)

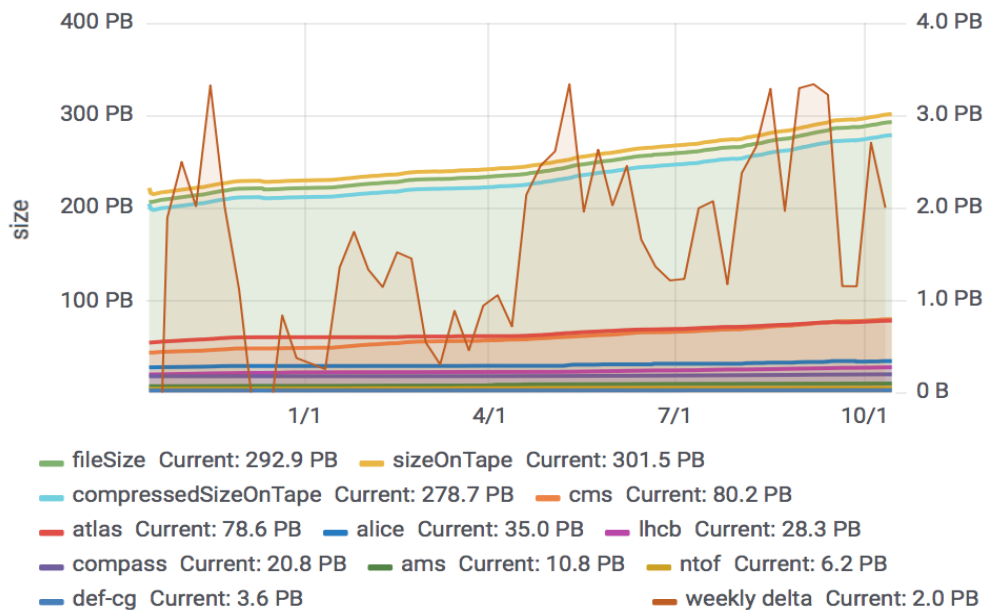  - SSDs are 5-10 times more expensive than spinning disks

www.cern.ch/eos
www.cern.ch/castor

# CASTOR dashboard ▾

https://filer-carbon.cern.ch/grafana/dashboard/db/castor-dashboard?orgId=1

## Physics Data in CASTOR



— fileSize  Current: 292.7 PB   — sizeOnTape  Current: 301.3 PB   — fileCount  Current: 584 Mil

## Total data



— fileSize  Current: 292.9 PB   — sizeOnTape  Current: 301.5 PB
— compressedSizeOnTape  Current: 278.7 PB   — cms  Current: 80.2 PB
— atlas  Current: 78.6 PB   — alice  Current: 35.0 PB   — lhcb  Current: 28.3 PB
— compass  Current: 20.8 PB   — ams  Current: 10.8 PB   — ntof  Current: 6.2 PB
— def-cg  Current: 3.6 PB   — weekly delta  Current: 2.0 PB

## File counters



— fileCount  Current: 584 Mil
— segmentsCount  Current: 556 Mil
— maxFileId  Current: 1.727 Bil

## Mean file size



— fileSize / fileCount  Current: 502 MB

# Electricity and cooling

- 2.7 MW consumption (+ ~ 1 MW cooling) from maximum 3.5 MW
  - 480 KW diesel generators
- Protected by UPS
  - Enough to start the diesel generators
  - Enough to shut down non-critical machines*
- Cooling
  - Chilled air via silver ducts enters the false floor and the into the closed server aisles
  - Water-cooled racks in the vault in the basement
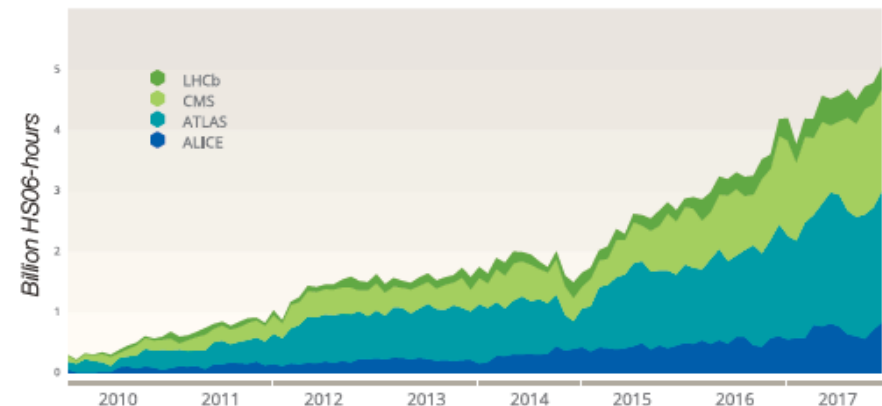
# WLCG – Worldwide LHC computing grid

- More than 170 data centres in 42 countries with about 800,000 cpu cores
  - CERN provides about 20% of the WLCG resources
  - Allows more than 10,000 physicists to access LHC data
  - >250,000 jobs run concurrently on the Grid
  - Storage is about 400 PB disk and 400 PB of tape globally
  - In 2016, global transfer rates have regularly exceeded 35GB/s
- Key facts and numbers (http://information-technology.web.cern.ch)

www.cern.ch/wlcg
www.cern.ch/wlcg-public

Evolution of the global core processor time delivered by the Worldwide LHC Computing Grid (WLCG)

As seen on the graph, the global central processing unit (CPU) time delivered by WLCG (expressed in billions of HS06 hours per month, HS06 being the HEP-wide benchmark for measuring CPU performance) shows a continual increase. In 2017, WLCG combined the computing resources of about 800 000 computer cores.

# Configuration management at CERN IT

FOREMAN

HAPROXY

**Certification Manager**

**CERNMegabus**

*43 000*
*Puppet managed machines*

**TEIGI Tool suite**

puppet

MCOLLECTIVE

python

PuppetDB

**AI-TOOLS**

...

riak KV

django

# The Puppet cycle
## Interactions with the server and the agent

**01** Store manifests into Git
As a first step, manifests (our config) have to be generated and stored in GitLab.

**02** Register a machine
A machine will then be created, in a specific hostgroup (eg. webchat/frontend/atlas). It will be registered in Foreman.

**03** Run Puppet
With the machine ready, the Puppet agent can be executed interactively (or let it run by itself). This will request the catalog (final state) of the machine.

**04** Master asks for hostgroup
The Puppet master handling the request will ask Foreman for the hostgroup of the machine.

**05** Master asks for manifests
Once it has the hostgroup, it will obtain the manifests that we defined in GitLab.

**06** Catalog generation
As a final step, the Puppet master will generate the catalog and return it to the agent, which will apply it to the machine.

**Thanks Config team fro the slide**

# CERNMegabus at CERN IT

**A service** that provides for **instant communication between services**

**CERNMegabus**

- The CERNMegabus architecture is based on the **publisher**-**consumer** **model and** utilises the **CERN IT messaging infrastructure**
- The publisher and the consumer services comprises of building blocks
    - **configured with Puppet**
    - to use the **CERNMegabus python libraries**
- **Installed on all** Puppet managed **machines** in the **CERN CC**

# CERNMegabus architecture



Change service "A"

2. Publish message

Central CERN IT ActiveMQ brokers

1. Subscribe for message

3. New message

Affected service "B"

4. Execute reactive actions

External program

The CERNMegabus architecture is based on the **publisher**-**consumer** **model and** utilises the **CERN IT** **messaging infrastructure**

# Already our clients

**CASTOR**

**HAPROXY**

**BATCH**

**EOS**

**CERNMegabus**

**43 000**
*Puppet managed machines*

**CLOUD**

**TEIGI Tool suite (roger)**

**DNS Load Balancing**

**IT Monitoring**

**CERN Computer Centre Power Cut Management**

# From roger to EOS/CASTOR/Puppet HAProxy

1. Subscribe for topic /topic/roger.**group.**castor** and hostgroup selector in message

**Central CERN IT ActiveMQ brokers**

2. Publish message to /topic/roger.**group.**castor**

Set roger state of *castor-lhcb-disknode-X* to disabled

3. Message de-queued

*castor-lhcb-headnode-Y*

**In practice - I**

hostgroup header: `castor-lhcb-diskservers`

{"**new**": {"update_time": "1538633786", "updated_by": "blueuser", "hostname": "*castor-lhcb-disknode-X*", …, "appstate": "**disabled**"}, "**old**": {"update_time": "1538633774", "updated_by": "somebody", "hostname": "*castor-lhcb-disknode-X* ", …, "appstate": "**production**"}

4. Execute `**modifydiskserver Disab castor-lhcb-disknode-X** `

Set castor-lhcb-disknode-X in read-only mode

**~1 sec**

# CERN CC Power Cut event



1. Preserve data

**Power back**

2. Shutdown
(all machines which we can)

t

# CERN CC Power cut event detection

ccpcoX programmatically detects power cut/power back event

# UPS and PLC

# CERN CC Power cut event detection algorithm



Data collection

Decision making algorithm

# CERN CC Power cut tests



- During mid-annual power cut test on the 2nd of July, 2018

  - Detected power cut

  - Notified the subscribed machines

  - Shutdown the machines, which had been predefined to be shutdown

  - Detected the power back

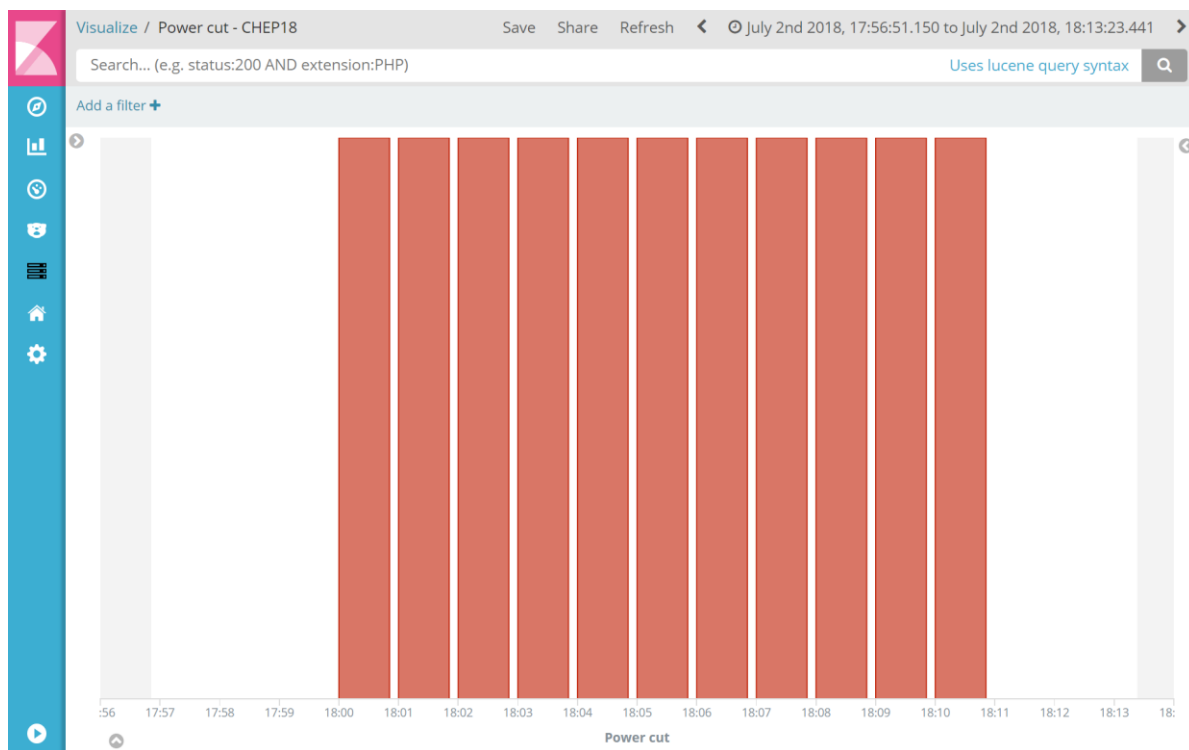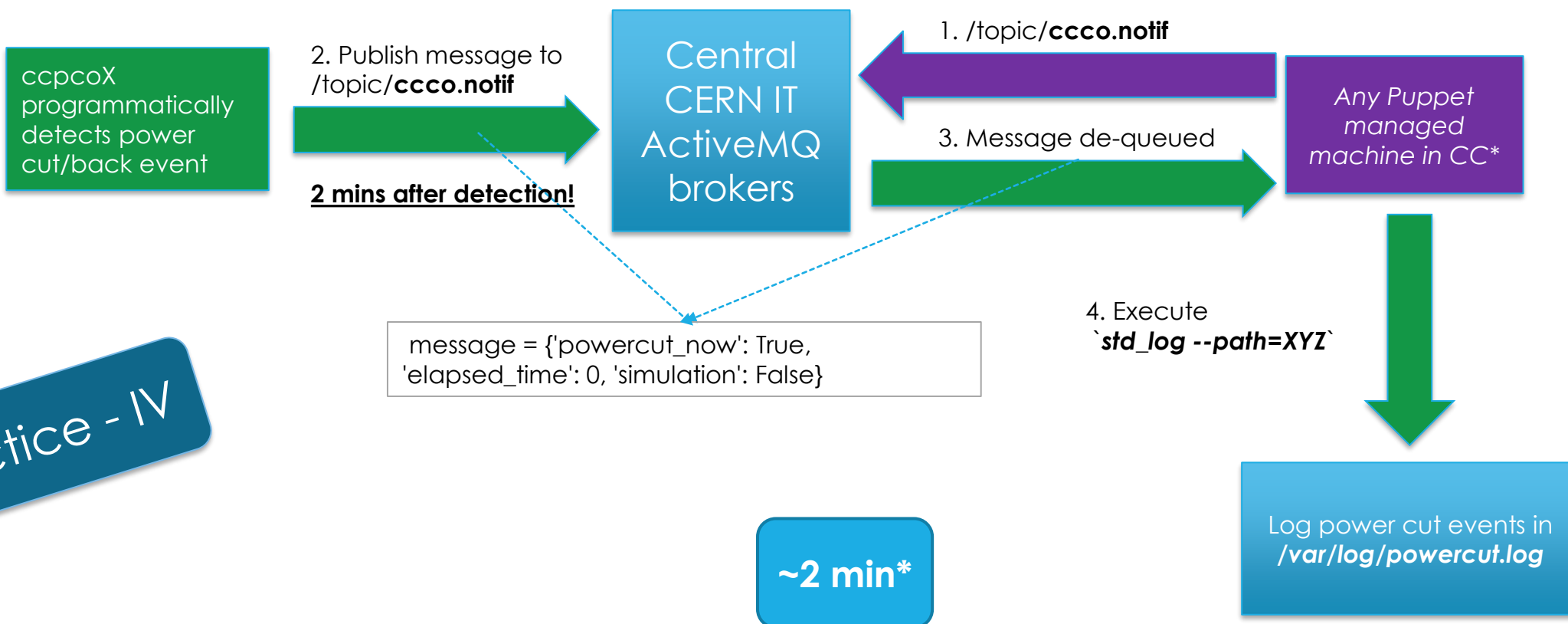  - Notified the machines, which predefined to wait

*Presented at CHEP'18*

# From CERN CC UPS PLC to CERN CC shutdown

In practice - IV

ccpcoX programmatically detects power cut/back event

2. Publish message to /topic/**ccco.notif**

**2 mins after detection!**

1. /topic/**ccco.notif**

Central CERN IT ActiveMQ brokers

3. Message de-queued

*Any Puppet managed machine in CC\**

message = {'powercut_now': True, 'elapsed_time': 0, 'simulation': False}

4. Execute `**std_log --path=XYZ**`

~2 min*

Log power cut events in **/var/log/powercut.log**

# Thanks

**THANKS for listening and ENJOY your visit at CERN**

**Zhechka**