

PROOF as a Service on the Cloud

a Virtual Analysis Facility based on the CernVM ecosystem

Dario Berzano

R.Meusel, G.Lestaris, I.Charalampidis, G.Ganis, P.Buncic, J.Blomer

CERN PH-SFT

CHEP2013 - Amsterdam, 15.10.2013

Distributed computing on the Cloud

- Resources are **heterogeneous**
→ *Like the Grid but on a more generic level*
- Virtual machines are **disposable**
→ *Cattle, not puppies: <http://bit.ly/cattlecloud>*



Building a cloud aware application

- **Scale** quickly when resources vary
→ *No prior pinning of data to process to the workers*
- Deal **smoothly** with major failures
→ *Automatic **failover** and clear recovery procedures*

Usual Grid workflow → static job pre-splitting ≠ cloud-aware

A cloud-aware analysis facility



*admins provide
virtual clusters*



*user's workflow
does not change*



geographically distributed independent cloud providers

Virtual Analysis Facility → analysis cluster on the cloud in one click



PROOF: the Parallel ROOT Facility

- Based on unique advanced features of **ROOT**
- Event-based parallelism
- **Automatic merging** and display of results
- Runs on batch systems and Grid with **PROOF on Demand**



PROOF is interactive

- Constant control and feedback of attached resources
- Data is not preassigned to the workers → *pull scheduler*
- Workers **dynamically attached** to a running process

Interactivity is what makes PROOF cloud-aware

PROOF on Demand: runs PROOF on top of batch systems

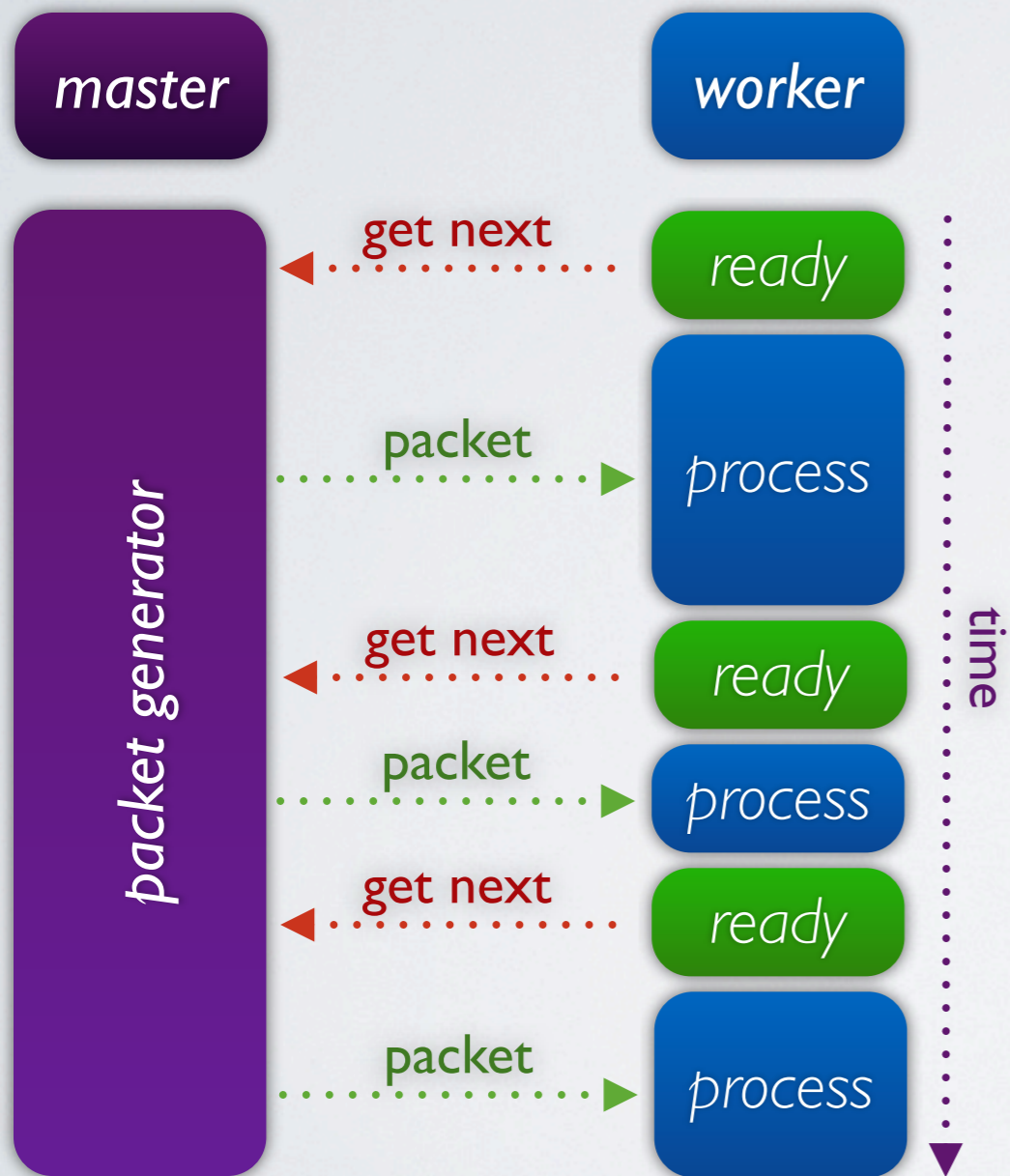
- **Zero configuration**
Per-user, no system-wide installation
- **Stability**
Potential failures affect only one user
- **Self-servicing**
User can restart her PROOF server
- **Advanced scheduling**
Uses scheduling

```
$> pod-server start
Starting PoD server...
updating xproofd configuration file...
starting xproofd...
starting PoD agent...
preparing PoD worker package...
select user defined environment script...
selecting pre-compiled bins to be added...
PoD worker package: /Users/volpe/.PoD/worker
-----
XPROOFD [1675] port: 21001
PoD agent [1680] port: 22001
PROOF connection string: proof@pcphsft1
-----

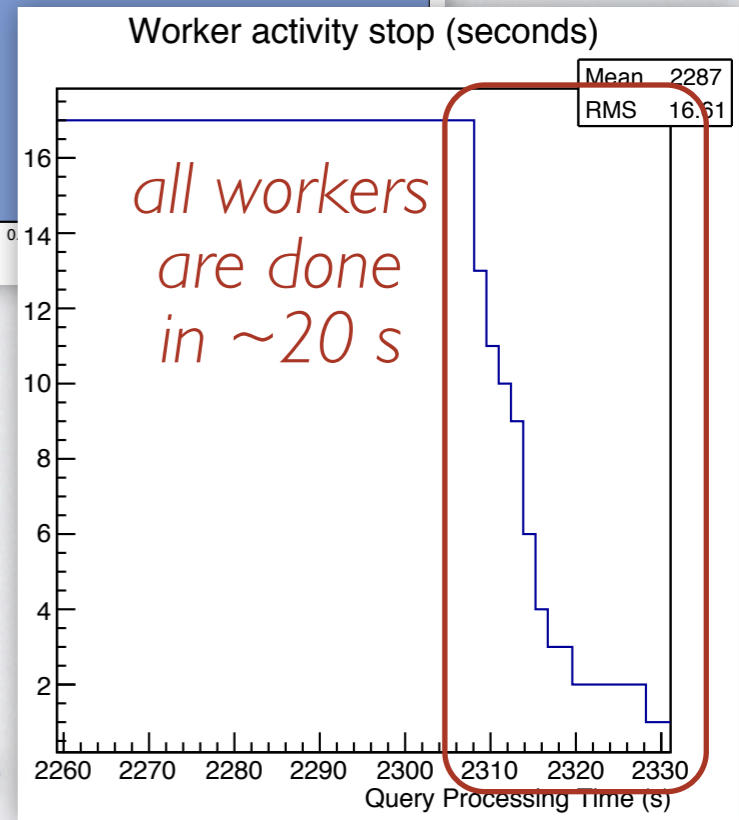
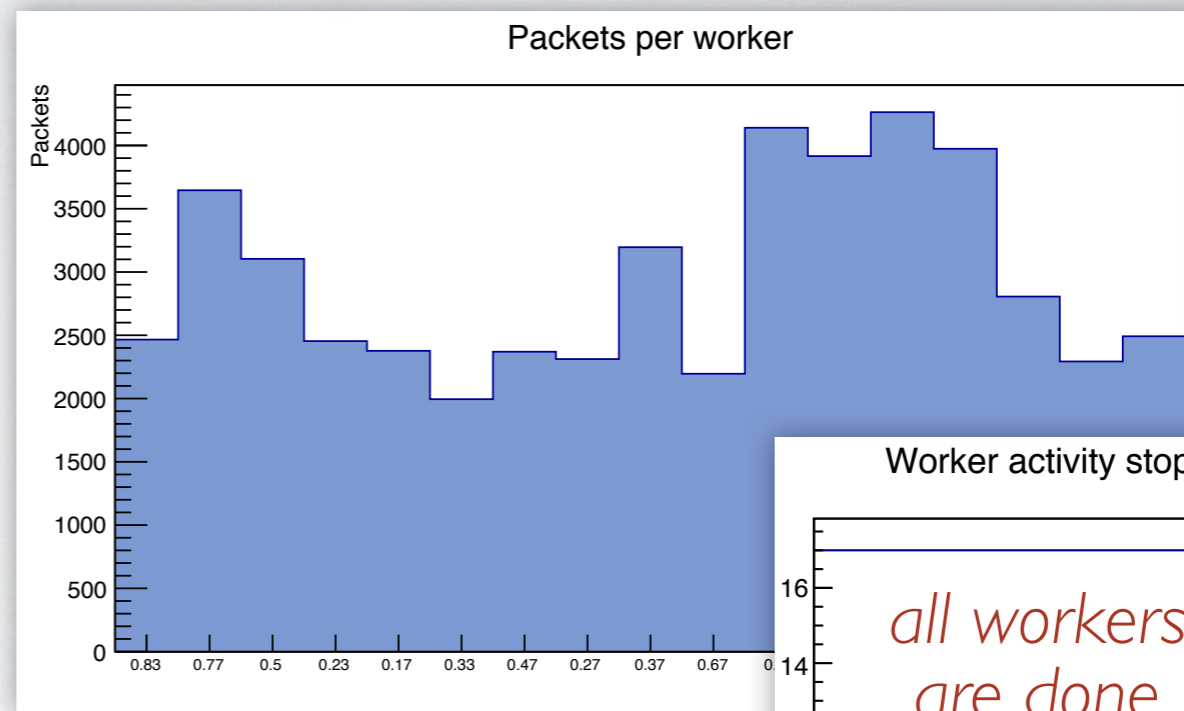
$> pod-submit -r condor -q proof -n 100
Job ID: 12345
```

POD pod.gsi.de
PROOF on Demand

Adaptive workload: workers ask for data to process when ready

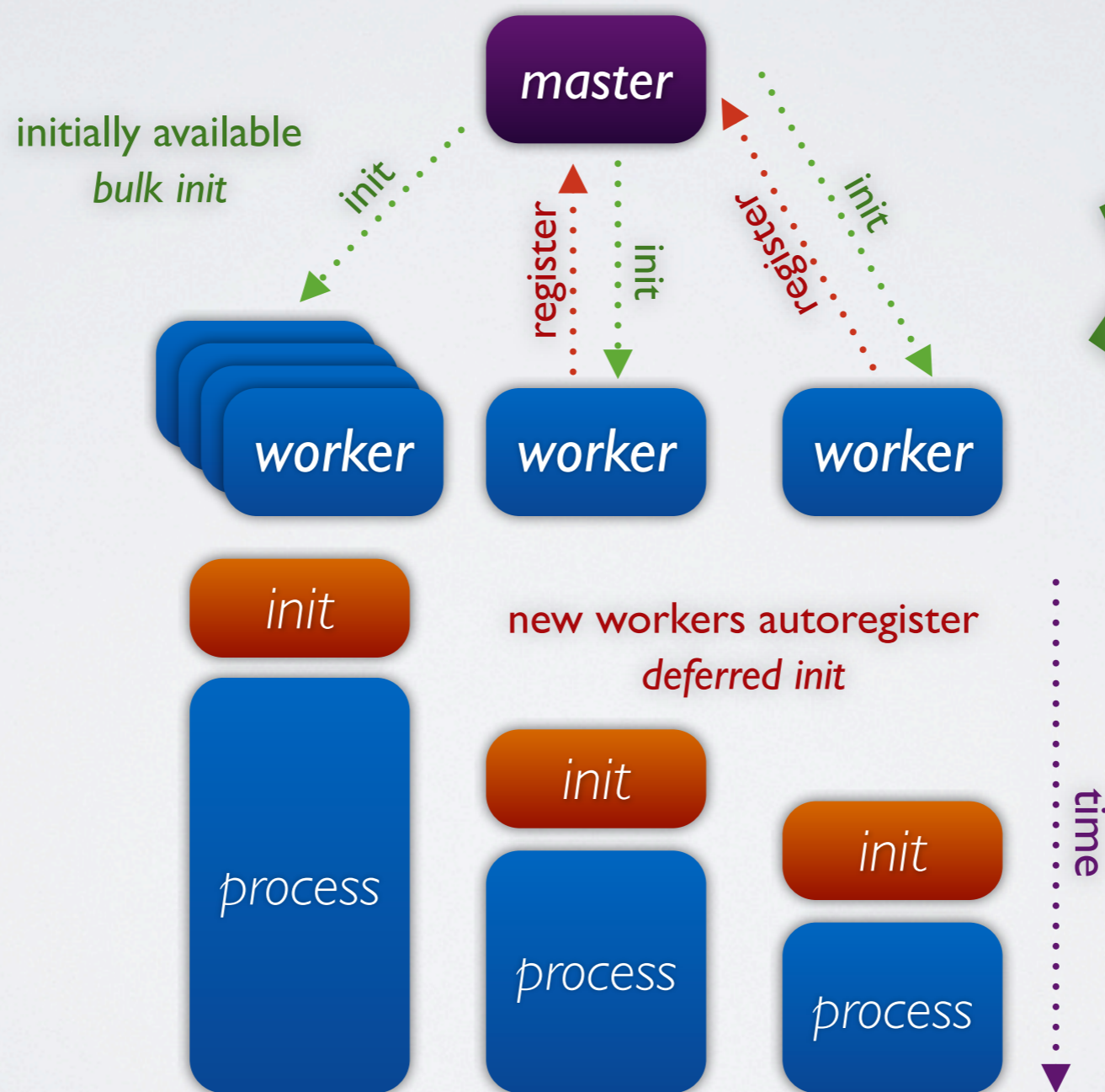


nonuniform workload distribution



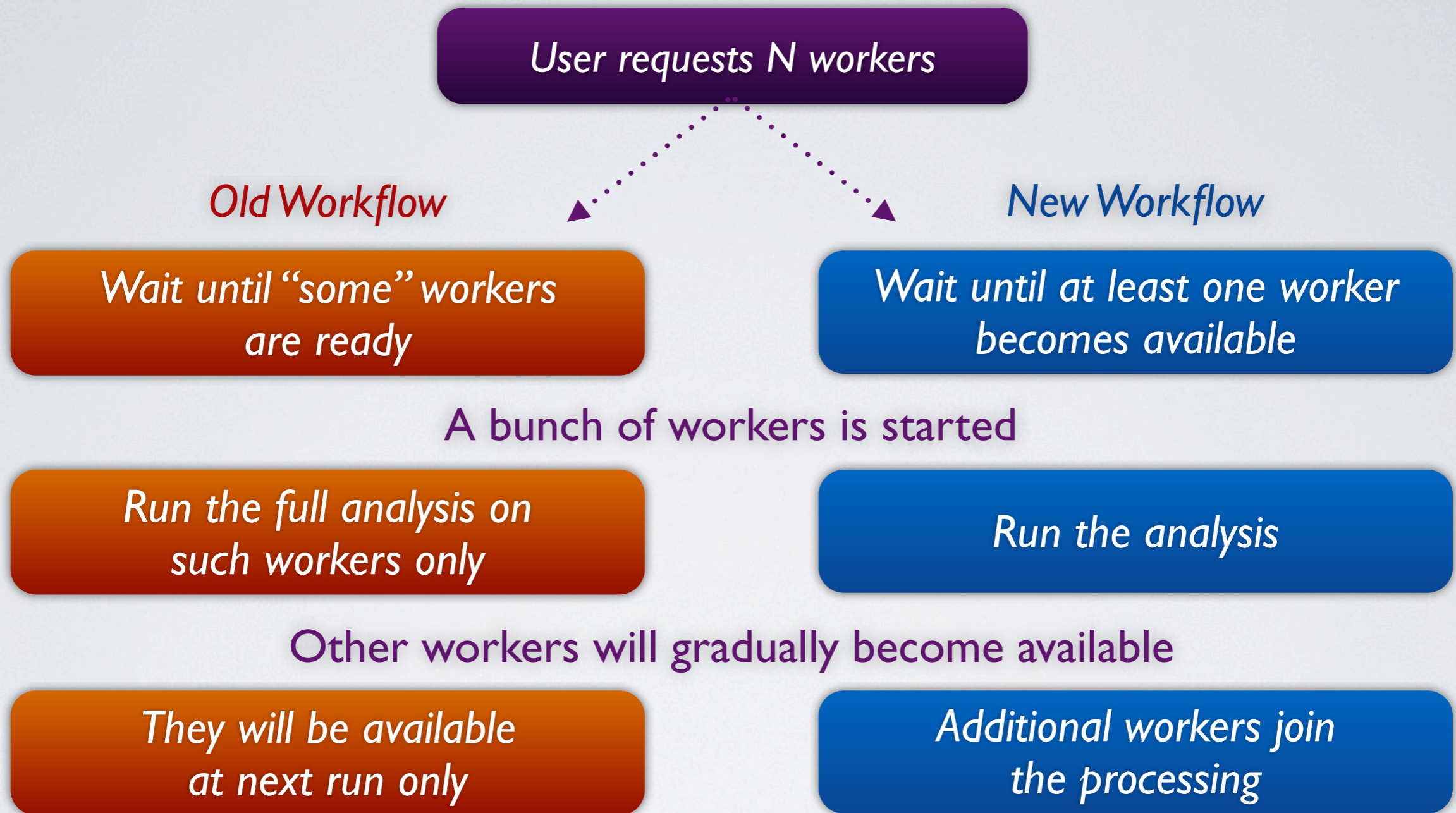
uniform completion time

Dynamic addition of workers *new workers can join and offload a running process*

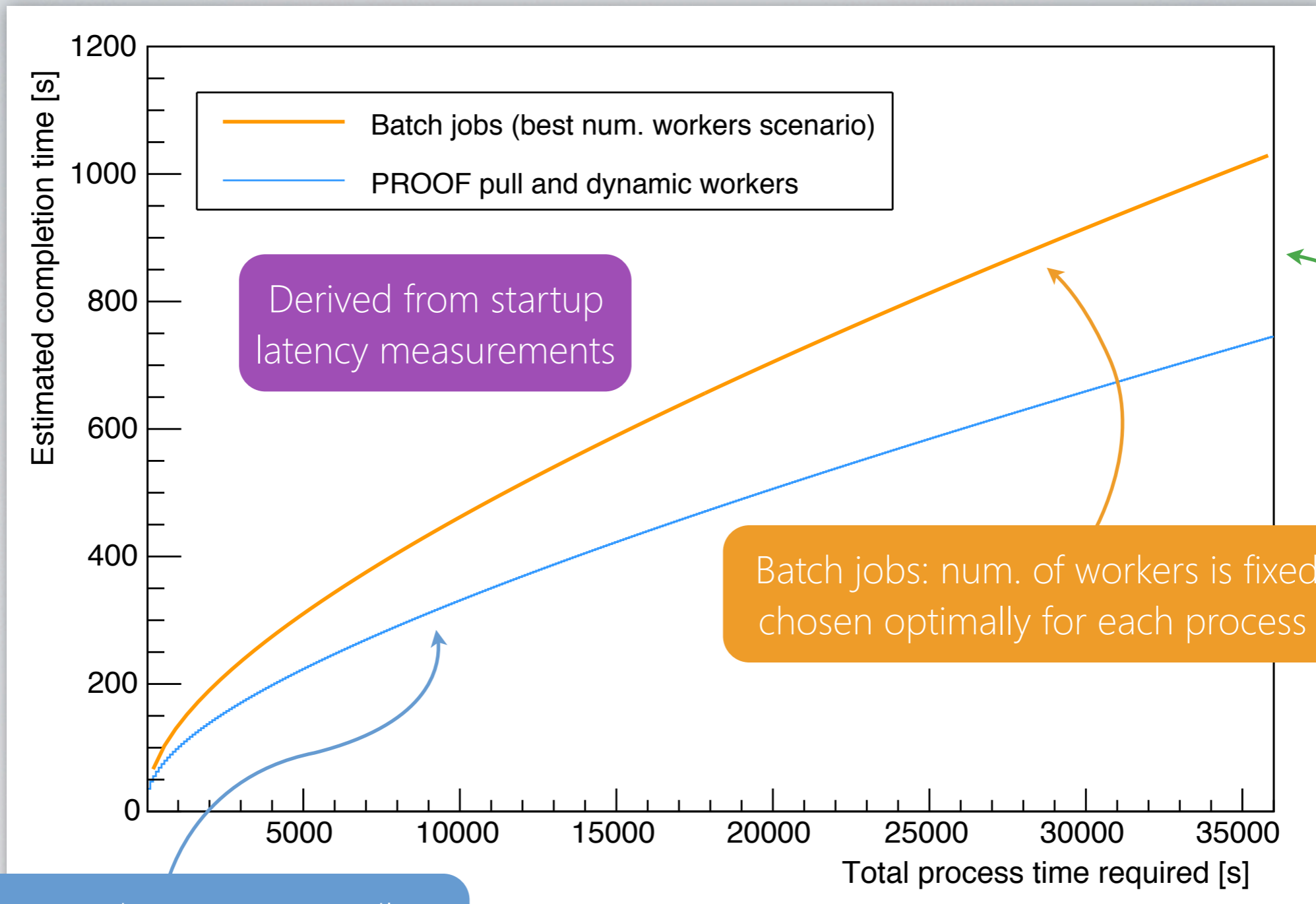


**NEW IN
ROOT
v5.34.10**

PROOF dynamic workers



Minimal latency and optimal resources usage
See ATLAS use case here: <http://chep2013.org/contrib/256>



Dynamic workers are never idle: uptime always spent in computing

Batch jobs: num. of workers is fixed and chosen optimally for each process time

Dynamic workers reduce completion time of a constant ~30%



PROOF

PoD

μ CernVM

Elastiq

HTCondor

CVM online

CernVM-FS

authn/authz

- **What:** a cluster of μ CernVMs with HTCondor
→ *One head node plus a scalable number of workers*
- **How:** contextualization configured on the Web
→ *Simple web interface: <http://cernvm-online.cern.ch>*
- **Who:** so easy that can be spawned by end users
→ *You can have your **personal analysis facility***
- **When:** scales up/down automatically whenever it's needed
→ *Frees unused resources*



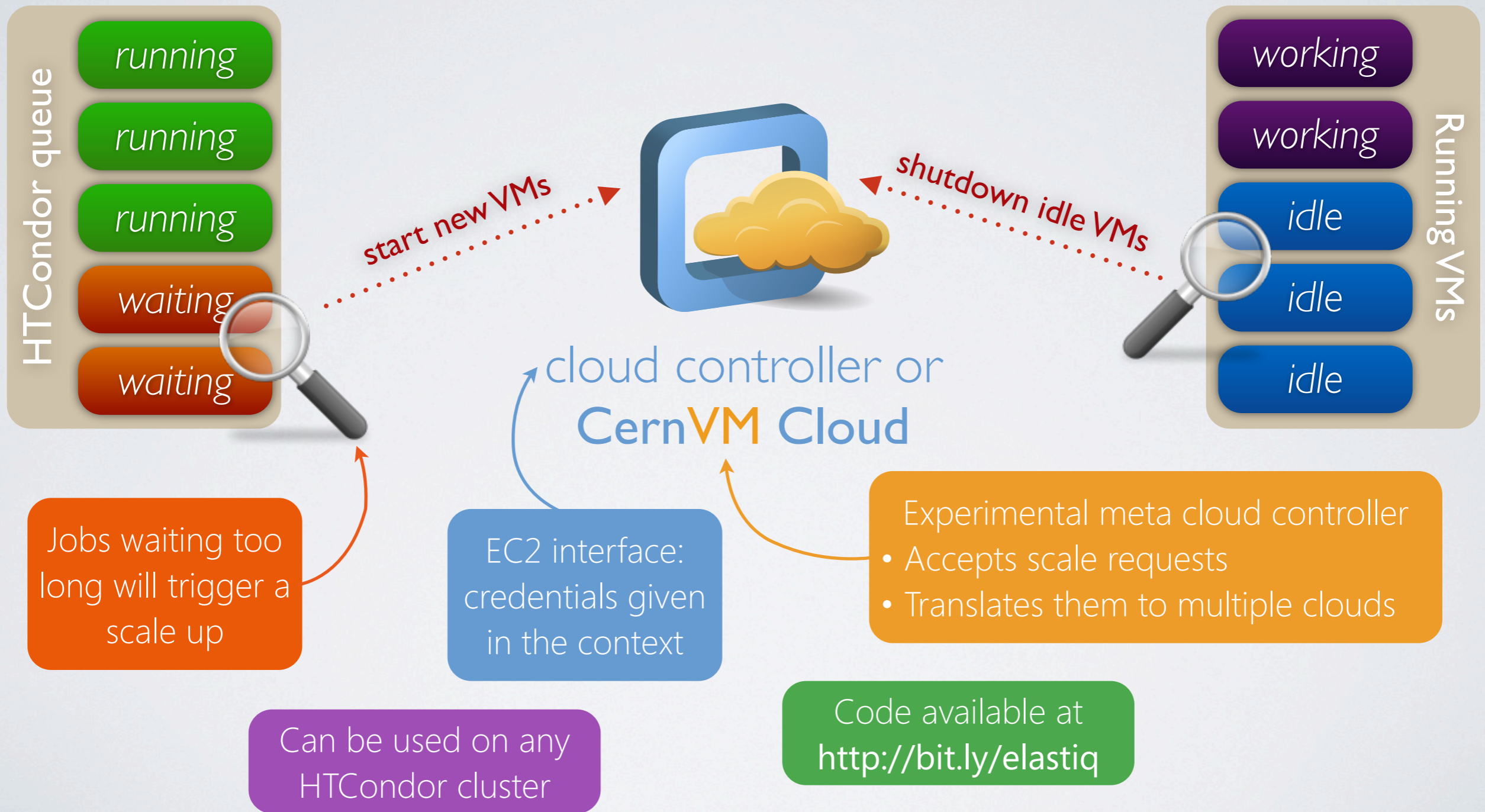
Integration to the CernVM ecosystem and HTCondor

- **μCernVM:** one-flavor SLC6 operating system on demand
→ See previous talk: <http://chep2013.org/contrib/213>
- **CernVM-FS:** HTTP-based cached FUSE filesystem
→ Both OS and experiments software downloaded on demand
- **CernVM Online:** safe context GUI and repository
→ See previous talk: <http://chep2013.org/contrib/185>
- **HTCondor:** light and stable workload management system
→ Workers auto-register to head node: no static resources configuration



The full stack of components is cloud-aware

Python app to monitor HTCondor and scale up or down



Context creation with CernVM Online: <http://cernvm-online.cern.ch>

Dashboard

Your context definitions

Name
VAF Ibex Master v1
VAF Ibex Node v1

Create new context...

New context

Virtual Analysis Facility node

Create new special context

Customize a few options

Context template

Please fill the following parameters and click create in order to create a new virtual machine context definition

User configuration

Context name:

Role:

Auth method:

Num. pool accounts:

Proxy for CVMFS:

HTCondor shared secret:

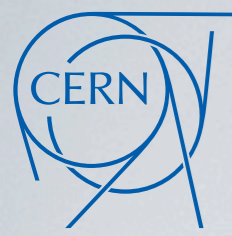
Context password:

Get generated user-data

Dashboard

Your context definitions

Name	ID	Operations
My PROOF Master	5e81170380b0432aaed63e40e1d90bbd	<input type="button" value="Clone"/> <input type="button" value="Publish"/> <input type="button" value="Clone"/> <input type="button" value="Clone"/>
VAF Ibex Master v1	dd3d44092b094f898eca4464e3d65124	<input type="button" value="Clone"/> <input type="button" value="Clone with full options"/> <input type="button" value="Get rendered context"/> <input type="button" value="Get raw user data"/>
VAF Ibex Node v1	3a7336b3485b4ba2918202b640a1c6c5	<input type="button" value="Clone"/>



Elastic cloud computing in action



The screenshot shows the OpenNebula Sunstone web interface. The main content area displays a table of Virtual Machines. A modal dialog box is open in the foreground, showing the execution progress of a task. The dialog includes a progress bar, a speedometer, and various status metrics.

ID	Owner	Group	Name	Status	IPs	VNC Access
3737	proof	ec2	ec2-m1-medium	RUNNING	172.16.212.121	VNC
3738	proof	ec2	ec2-m1-large	RUNNING	172.16.212.122	VNC
3739	proof	ec2	ec2-m1-large	RUNNING	172.16.212.123	VNC
3740	proof	ec2	ec2-m1-large	RUNNING	172.16.212.124	VNC
3741	proof	ec2	ec2-m1-large	RUNNING	172.16.212.125	VNC
3742	proof	ec2	ec2-m1-large	RUNNING	172.16.212.126	VNC
3743	proof	ec2	ec2-m1-large	RUNNING	172.16.212.127	VNC
3744	proof	ec2	ec2-m1-large	RUNNING	172.16.212.128	VNC
3745	proof	ec2	ec2-m1-large	RUNNING	172.16.212.129	VNC
3746	proof	ec2	ec2-m1-large	RUNNING	172.16.212.130	VNC

Dialog Box Content:

```

PROOF Query Progress: dberzano@localhost
Executing on PROOF cluster "localhost" with 78 parallel workers:
Selector: ProofEvent.C
0 files, number of events 100000000, starting event 0
[Progress Bar: 17%]
Initialization time: 0.4 secs
Estimated time left: 3 min 7 sec
Processing status: 47684014 / 100000000 events - 0.00 MB
Processing rate: 827183.6 evts/sec
                  avg: 279466.8 evts/sec (0.0 MB/sec)
 Close dialog when processing is complete
 Smooth speedometer update
  
```

Speedometer: Ev/s x10⁻², Proc Time (ms) 0000

Buttons: Show Logs, Performance plot, Memory Plot, Enable speedometer, Run in background, Stop, Cancel, Close

Screencast:
<http://youtu.be/t1s2jFaXf64>

Measured the delay before requested resources become available

Target clouds:

- Small: **OpenNebula @ INFN Torino**
- Large: **OpenStack @ CERN (Agile)**

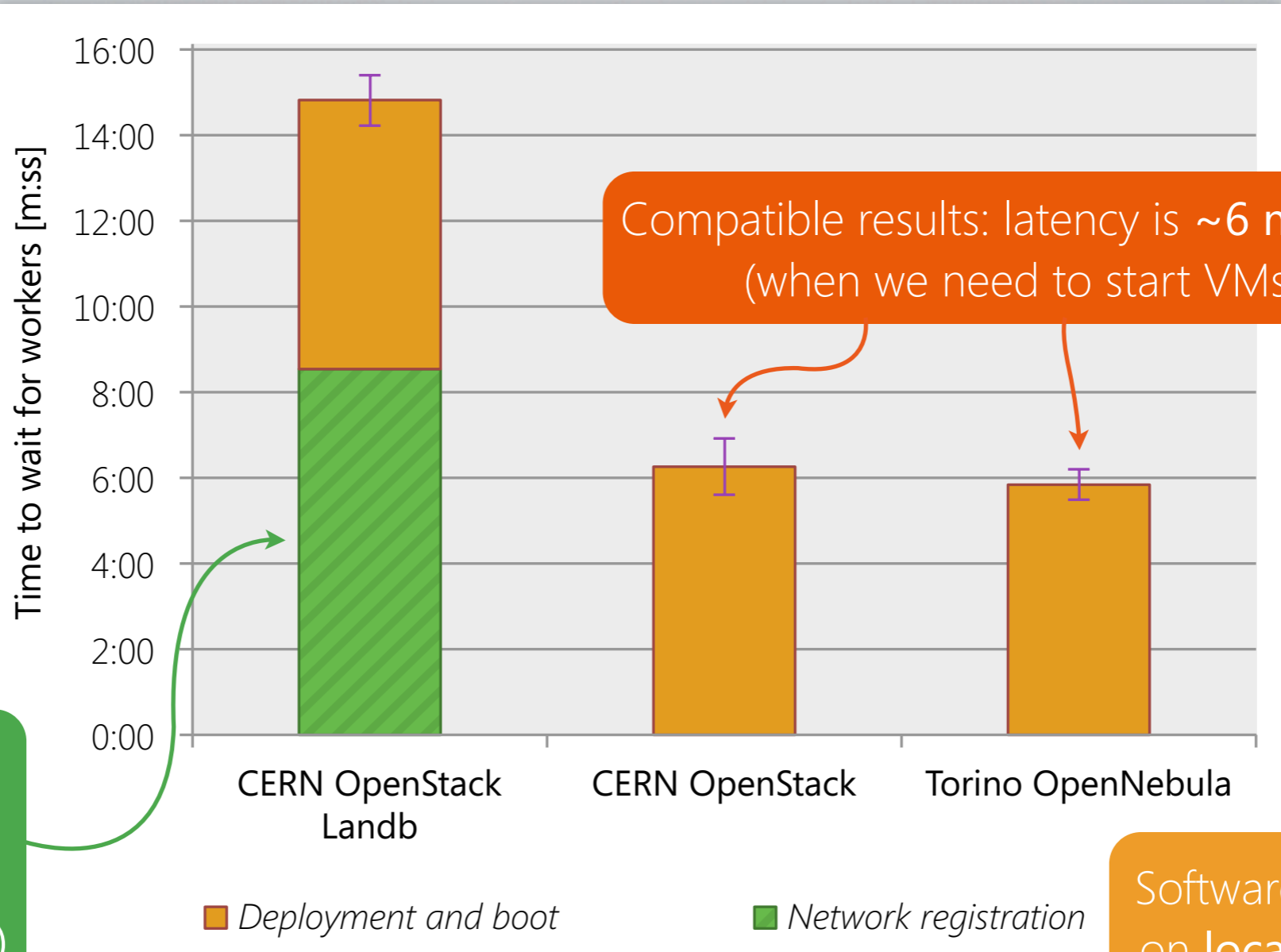
Test conditions:

- μ CernVM use a HTTP caching proxy
→ **Precaching** via a dummy boot
- μ CernVM image is 12 MB big
→ image **transfer time negligible**
- VMs deployed when resources are **available**
→ rule out failures and delay due to lack of resources

Note: not comparing cloud infrastructures. Only measuring μ CernVM+PROOF latencies.

Measuring latency due to:

- μ CernVM boot time
- HTCondor automatic registration of new nodes
- PoD and PROOF reaction time



Latency due to legacy hostname registration issues (this test ignored)

Compatible results: latency is ~6 minutes (when we need to start VMs)

Software and OS precached on local Squid proxies by a previous test run

Measured time elapsed between PoD worker's request and availability:
`pod-info -l`

10 VMs started in the test

Every VAF layer is cloud-aware and elastic

- **PROOF+HTCondor** deal with dynamic addition/removal of workers
- **μCernVM** is very **small** and fast to deploy
- **CernVM-FS** downloads only what is needed

Consistent configuration of solid and **independent** components

- **No login** to configure: all done via **CernVM Online** context
- PROOF+PoD also work dynamically **on the Grid**
- **Elastiq** can scale **any HTCondor cluster**, not PROOF-specific
- **Reused** existing components wherever possible

Thank you for your attention!



- **PROOF (the Parallel ROOT Framework)**
<http://root.cern.ch/drupal/content/proof>
- **Virtual Analysis Facility client and Elastiq**
<https://github.com/dberzano/virtual-analysis-facility>
- **The CernVM Ecosystem**
<http://cernvm.cern.ch/portal/publications>
- **Cloud @ INFN Torino**
<http://chep2013.org/contrib/474>
- **CERN Agile Infrastructure**
<http://chep2013.org/contrib/86>