



Running jobs in the Vacuum

Andrew McNab
University of Manchester

Mario Ubeda Garcia
& Federico Stagni
CERN

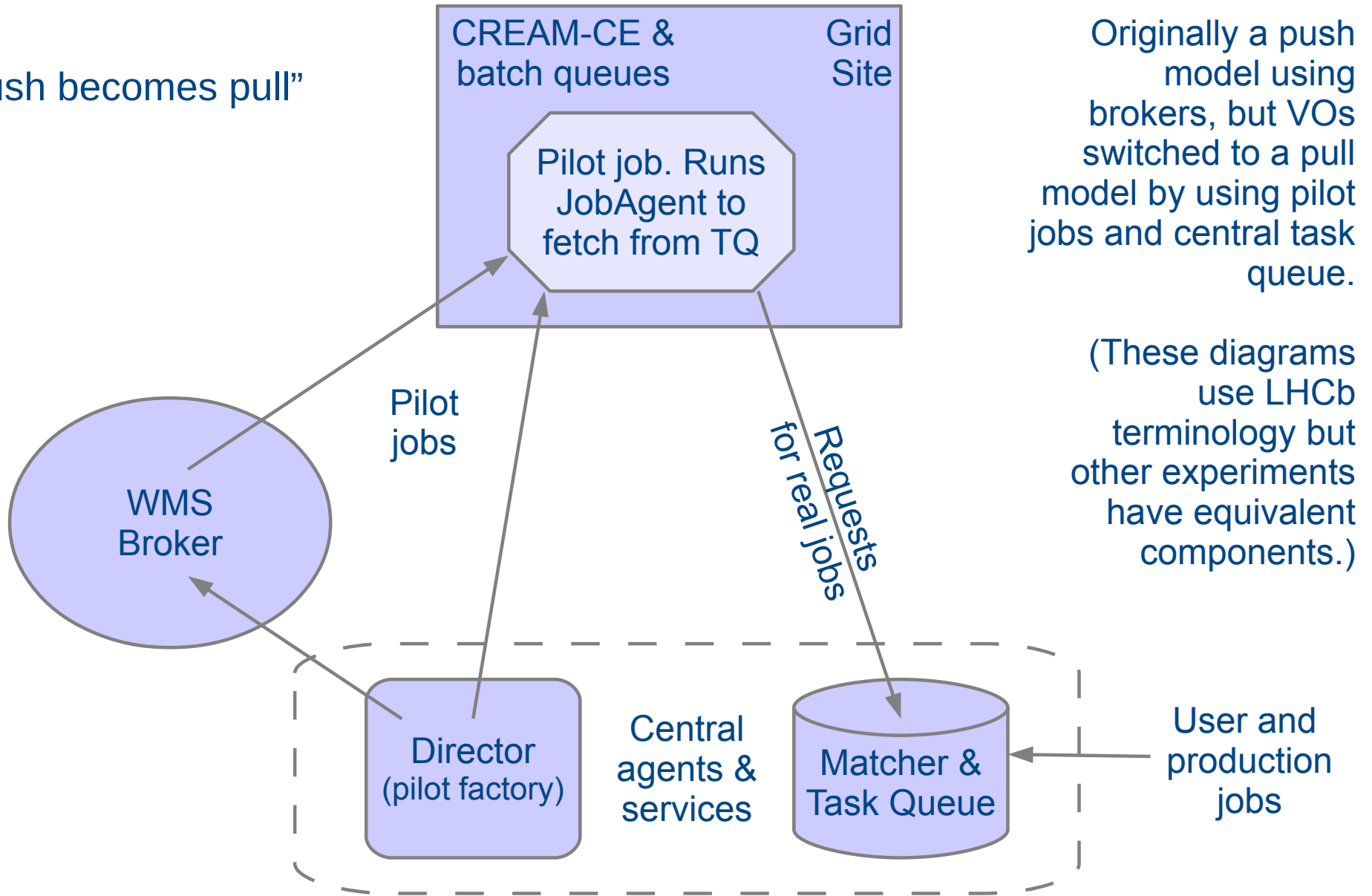


Overview

- The Grid vs The Cloud vs The Vacuum
 - 3 models
 - Vacuum definition
- The Vac implementation
 - UDP protocol
 - Target shares
 - Shutdown messages
- Production use within LHCb
 - Contextualization
 - Production tests
- Summary

The Grid

“Push becomes pull”



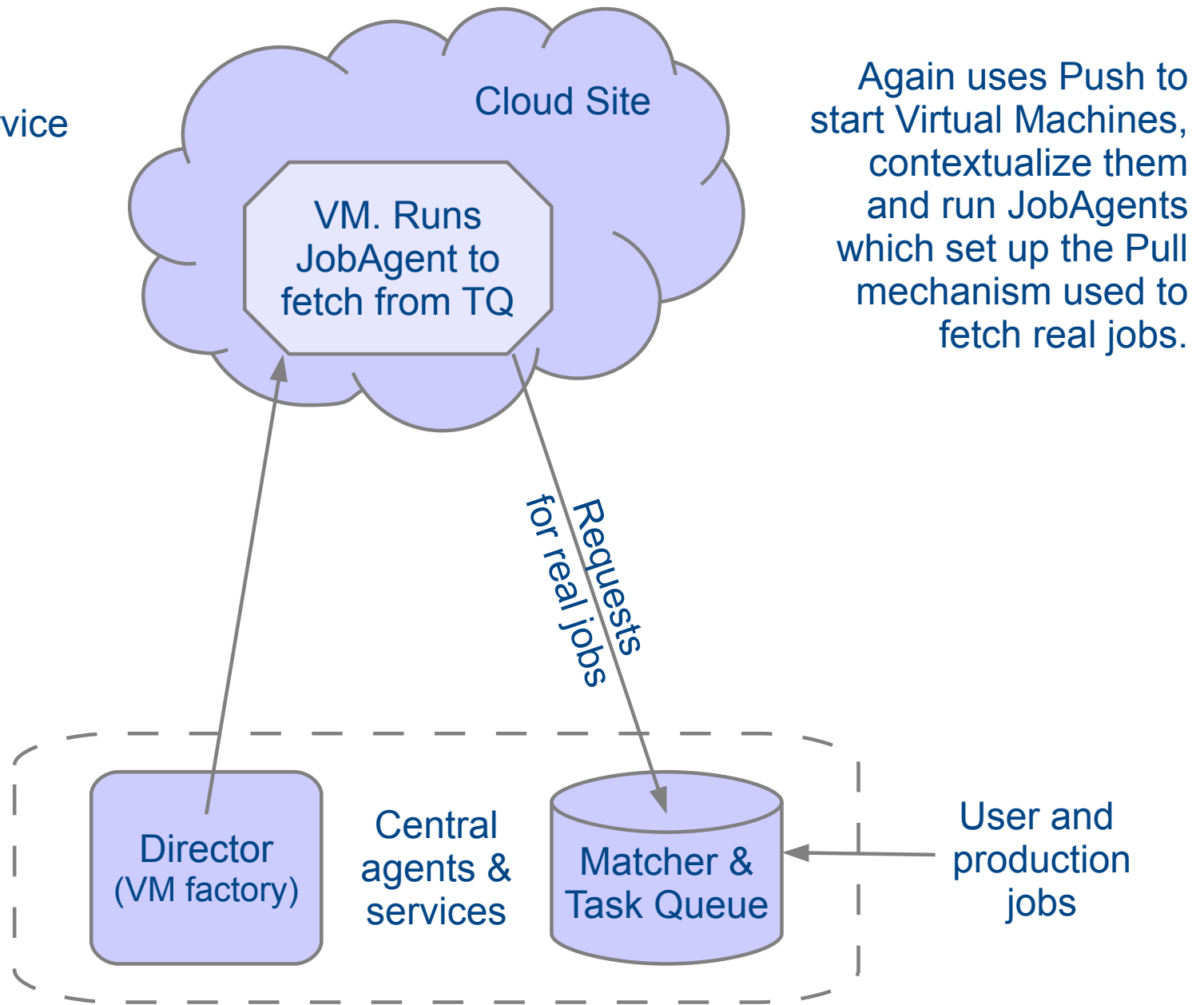
Originally a push model using brokers, but VOs switched to a pull model by using pilot jobs and central task queue.

(These diagrams use LHCb terminology but other experiments have equivalent components.)

The Cloud

Infrastructure-as-a-Service
(IaaS)

In LHCb, use the
same TQ as for Grid
and direct DIRAC
execution of jobs.



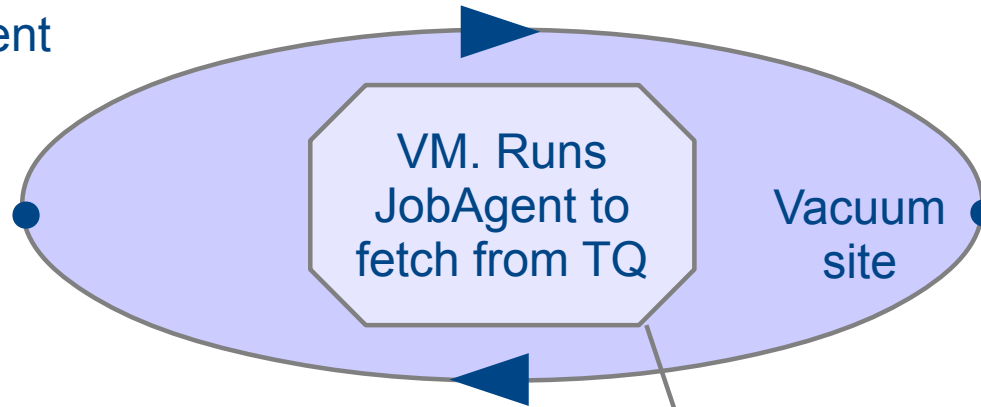
Again uses Push to
start Virtual Machines,
contextualize them
and run JobAgents
which set up the Pull
mechanism used to
fetch real jobs.

“The Vacuum”

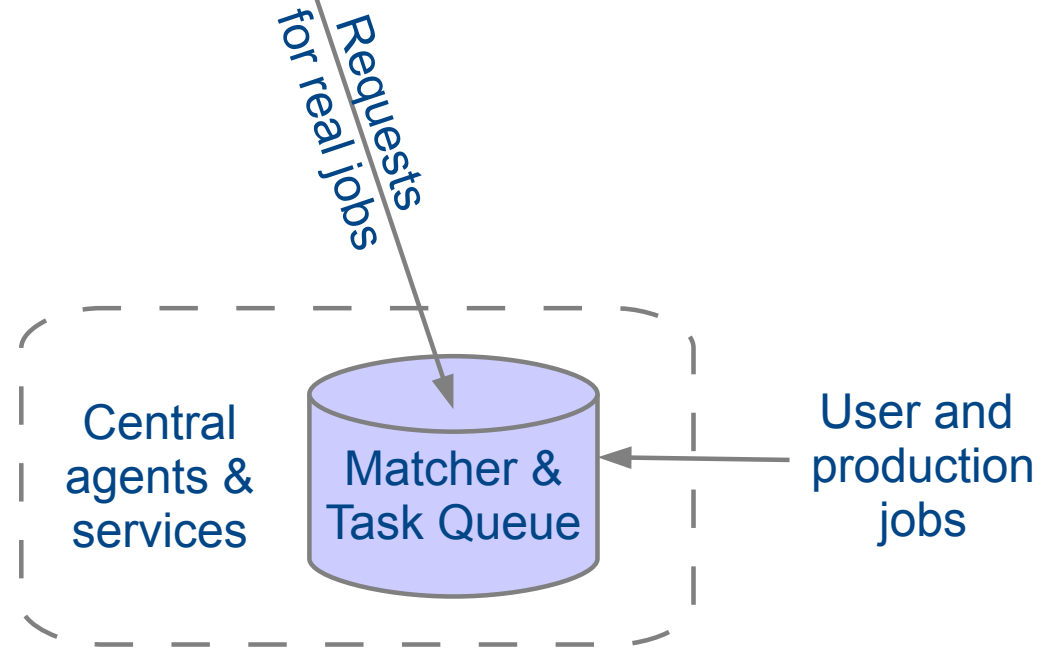
Infrastructure-as-a-Client
(IaaS)

Instead of being created by VOs, the Virtual Machines appear spontaneously “out of the vacuum” at sites.

As with the other models, the JobAgent runs and requests real jobs from the Matcher and normal Task Queue.



Hypervisors/hosts can run VMs for particular VOs depending on work available and target shares for each VO.





Vacuum Model

- For the experiments, VMs appear by “spontaneous production in the vacuum”
 - Like virtual particles in the physical vacuum: they appear, potentially interact, and then disappear
- From the conference paper:
 - *“The Vacuum model can be defined as a scenario in which virtual machines are created and contextualized for experiments by the site itself. The contextualization procedures are supplied in advance by the experiments and launch clients within the virtual machines to obtain work from the experiments' central queue of tasks.”*
- At many sites, 90% of the work is done by 2 or 3 experiments
 - This justifies effort to set them up at the site (comparable to site-info.def etc)
 - In return, the site is no longer dependent on all the CREAM/batch or Cloud machinery for these jobs



Vac implementation

- On each physical node, Vac VM factory daemon runs to create and apply contextualization to transient VMs
- Multiple VM flavours (“VM types”) are supported, ~1 per experiment
 - Could also include a PBS worker node in a VM, a PROOF worker VM etc
- Each site or Vac “space” is composed of independent factory nodes
 - All using the same `/etc/vac.conf`, `/etc/vac-targetshares.conf` etc
- Factories communicate with each other via UDP
 - Type of VM to start in a free slot based on what else is running and target shares
 - So no headnode central point of failure; robust against losing individual nodes
- So far supports CernVM image and ISO (AMI) contextualization
 - VO supplies `user_data`, `prolog.sh` and `epilog.sh`; Vac makes the ISO image
 - Vac also makes and NFS exports HEPiX `machinefeatures` etc directories to the VM, including `shutdowntime` and `shutdown_command`

Vac UDP protocol

- Each factory has a list of all the factories in the same Vac space
 - (May offer multicast as an alternative in the future...)
- Sends UDP packet containing a JSON-encoded Python dictionary:
 - {"cookie": "179e6....cd3a5", "method": "status",
"space": "vac01.tier2.hep.manchester.ac.uk"}
- One UDP reply for each VM assigned to a factory:
 - State (shutdown / starting / running), VM type etc
 - Outcomes of last VM instance run here for each VM type
- Use cookies to avoid external denial of service, since UDP
- Can also be used to query nodes (vac command \approx PBS qstat)
- Use UDP port 995 (Roman Numerals: V=5, M=1000. 995=1000-5 !)

“Back off”

- To avoid overloading Matcher/TaskQueue, Vac implements “back off”
- If a VM finishes with “no work” / “banned” / “site misconfigured” outcomes then it counts as an abort
 - If no outcome given, then if a VM finishes after less than fizzle_seconds (600sec?) then it counts as an abort
- For a VM type (~experiment), if an abort has happened on any factory in the last backoff_seconds (600 sec?), then no more VMs of that type will be started
- After that, if an abort happened in the last backoff_seconds + fizzle_seconds and any new VMs have run for less than fizzle_seconds, then no more VMs of that type will be started
 - ie try to run one or two test VMs to see if ok now
- If backoff_seconds + fizzle_seconds have passed without more aborts, then can start VMs again as fast as slots become available



Target shares

- Vac avoids the phrase “fair shares”
- No history recorded: just a targetshares list in the configuration
 - Each time a new VM must be created, the VM type with the least running VMs, weighted by the shares, is tried
 - The back-off mechanism is there to stop trying VM types which have “recently” failed to get any work and failed to stay running
- This approach is very simple, and means the factory nodes can decide themselves what to do
 - Avoids a central management daemon which would be a single point of failure
- But these target shares are instantaneous
 - They are fair, in that if all experiments submit lots of jobs, the site shares out the capacity according to the stated shares
 - But they are unfair in that quiet periods aren't credited and carried forward



Long-term target shares

- The intention is that sites address this by updating the targetshares list and pushing it out to the factory nodes
 - Can use puppet or whatever they use for configurations elsewhere
 - Separate `/etc/vac-targetshares.conf` can be used for convenience
- The plan is to provide a tool to use the local APEL database to calculate targetshares to achieve long term target shares figures
 - For example: “We've not run any jobs from Expt. X this quarter so far. Set X's instantaneous target share so if some jobs do arrive from X, they will all be run.”
- Rather like an offline version of MAUI with a long time constant
- The big advantage is that if this tool fails to run, the system still carries on running jobs with the current share values and doing useful work while you fix it.
- Smaller sites can also just do this by hand every week or so (as some do with short term MAUI shares...)

Controlling VM lifetimes

- Vac strategy is to use the HEPiX machinefeatures directory
 - NFS exported from factory node into VM, so easy to populate and to update
- So far we set shutdowntime when the VM is created
 - Always kill the VM at that time if still running
 - Will add other VM/job length values in response to new task force effort
- We set `/etc/machinefeatures/shutdown_command`
 - Allows the VM to shut itself down. This may not be appropriate for clouds, but is extremely useful for IaaS systems, like Vac and BOINC
- Vac's default shutdown wrapper saves any command line arguments as the shutdown message
 - These help the site see why VMs are terminating
 - This is a polite thing the VM can do for the site's benefit
 - Uses Vac's writeable NFS directory `/etc/machineoutputs`



Shutdown message codes

- Passed as arguments to shutdown_command:
 - 100 Shutdown as requested by the VM's host/hypervisor
 - 200 Intended work completed ok
 - 300 No more work available from task queue
 - 400 Site/host/VM is currently banned/disabled from receiving more work
 - 500 Problem detected with environment/VM provided by the site
 - 600 Error related to job agent or application within VM
- As with HTTP codes, room to insert more numbers for finer grained information in the future
- Vac uses this information programmatically, but useful to admins too
- More details and rationale:
 - https://www.gridpp.ac.uk/wiki/HEPiX_shutdown_command

Interface with LHCb JobAgent

- Contextualization procedure causes JobAgent to be started in VM
- Vac expects VMs to shutdown if they can't find any work to do
- Use the shutdown_command protocol for site/factory to tell the VM how to shut itself down
 - vac-shutdown-vm is a wrapper around “sudo shutdown -h now”
- Use shutdown code + message saying why it has shutdown (eg “300 Nothing to do” if no more work in TQ)
 - vac-shutdown-vm writes this to /etc/machineoutputs/shutdown_message
 - this directory is NFS-exported from factory into VM; Vac examines it afterwards
- TimeLeft.py is being extended to read HEPiX shutdowntime
 - Vac creates shutdowntime using maximum allowed lifetime of a VM of this type
 - Can also be used to ask VM to stop with, say, 24 hours warning



Vac testing at sites

- During the summer ran routine LHCb production Monte Carlo:
 - Jobs run at Manchester, Lancaster and Imperial Vac sites
 - 3300 LHCb production jobs run across these sites with Vac
- Aim is to make Vac sites look like “normal” WLCG/EGI sites wherever possible
 - `uk.ac.gridpp.vac` service endpoint type registered in GOCDB
- Accounting data successfully published to site APEL database
 - `vacd` writes PBS and BLAHP format accounting files which work with the existing APEL PBS parser
 - Next test with EMI3 site APEL and publish into global APEL database
- Once accounting is demonstrated, next steps are to try larger deployments with more than a handful of machines per site
 - Need accounting done so UK sites don't lose out in work-based funding



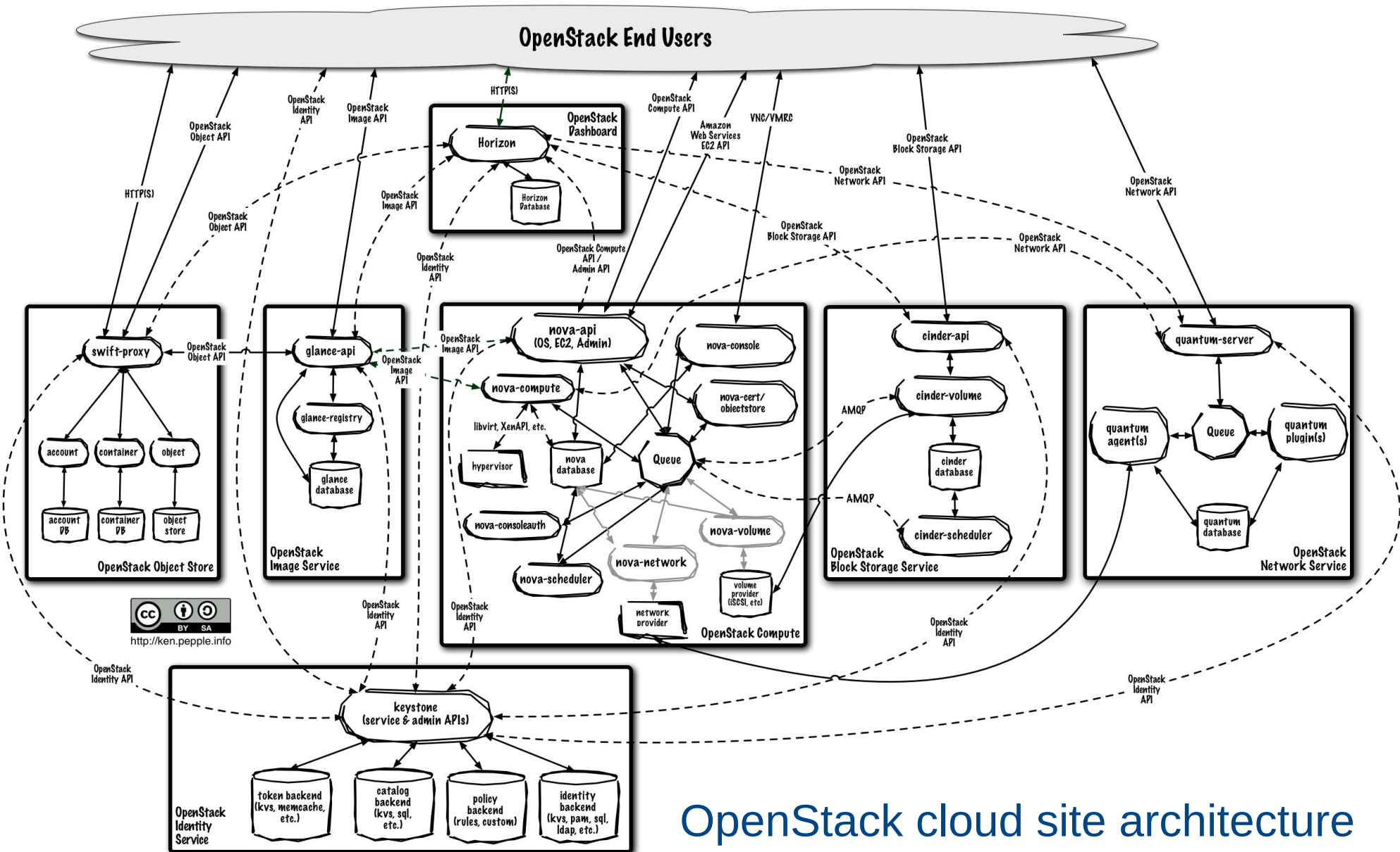
Summary

- Vacuum is an alternative to Grid and Cloud models
 - Quite complementary to Clouds + VMs with the same VMs on both
- Vac is an implementation of this model
 - VM factory on each physical machine, communicating via UDP
 - See <http://www.gridpp.ac.uk/vac/> for RPMs, documentation etc
- Tested at multiple sites; successfully ran production LHCb jobs
- Next step is to try with larger deployments



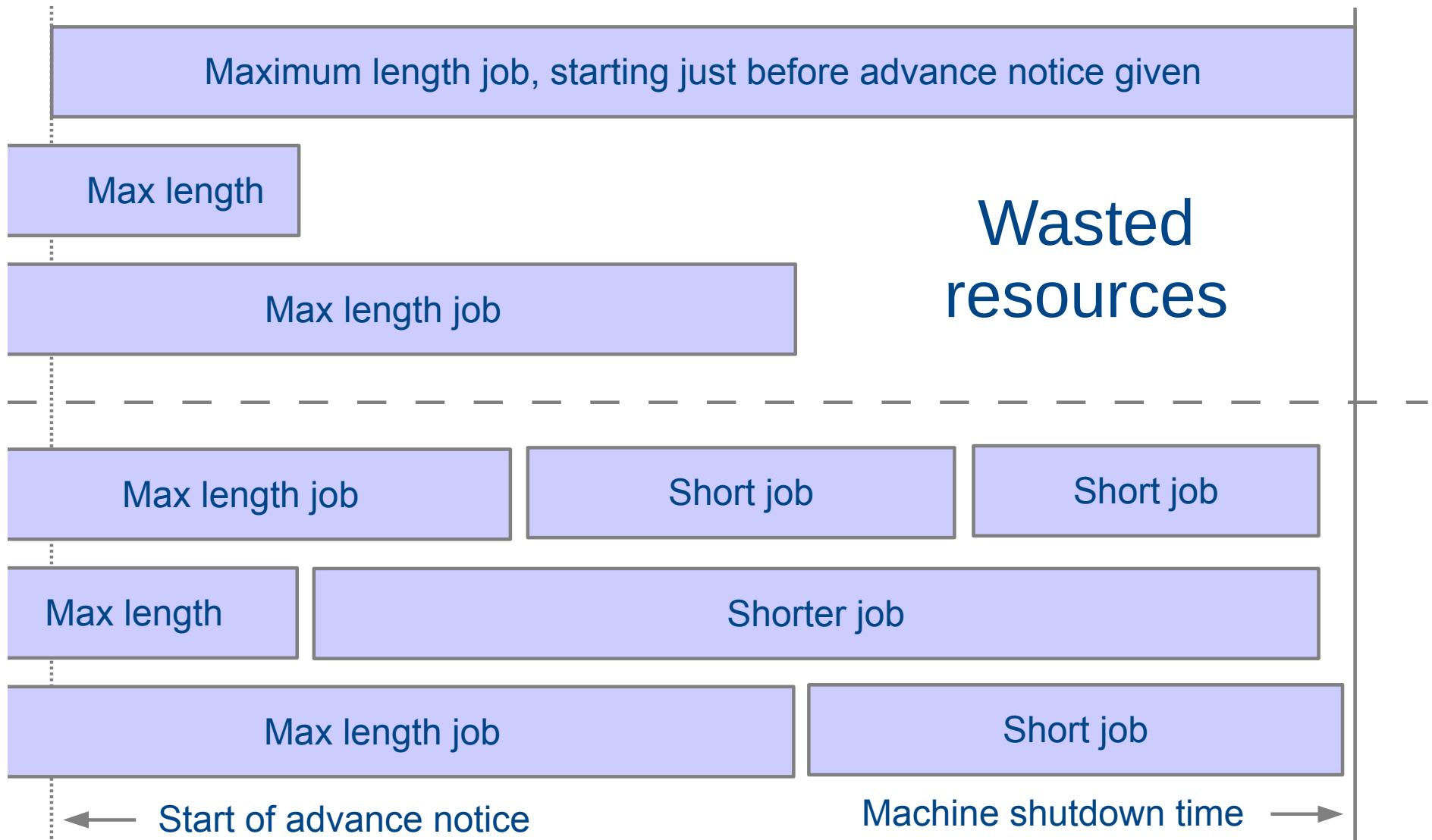
Extra slides

OpenStack End Users



OpenStack cloud site architecture

Graceful mechanisms allow job “masonry” ...



Documentation...

```
work — root@wn2209120:/usr/man/man5 — nc — 93x29
vac.conf(5)                               Vac Manual                               vac.conf(5)

NAME
    vac.conf - Vac configuration file

DESCRIPTION
    vacd is a daemon which implements the Vacuum model on a factory (hypervisor)
    machine. vacd reads its configuration from /etc/vac.conf and this file is
    also read by the vac utility command to find default values.

    vac.conf uses the Python ConfigParser syntax, which is similar to MS Windows
    INI files. The file is divided into sections, with each section name in
    square brackets. For example: [settings]. Each section contains a series of
    option=value pairs.

    For ease of management, three more optional files are read if they exist:
    "/etc/vac-virtualmachines.conf"
    "/etc/vac-factories.conf"
    "/etc/vac-targetshares.conf"

    These files are named after important sections, but sections can be placed
    in any of the four files, or all placed in /etc/vac.conf

[SETTINGS] OPTIONS
    The [settings] section is required and has options which apply to all vir-
    tual machines.

    vac_space is required and gives the name of this Vac space. A single space
:|
```