



# CERN Agile Infrastructure Road to Production

Steve Traylen,

[steve.traylen@cern.ch](mailto:steve.traylen@cern.ch)

@traylenator

CERN, IT Department

HEPiX Autumn 2012 Workshop

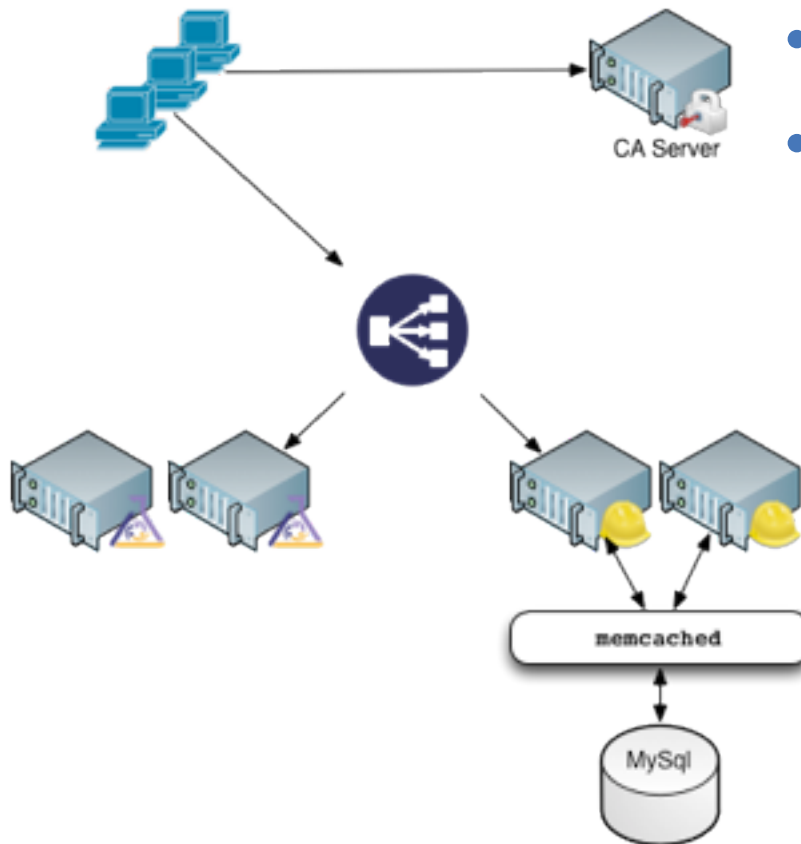
- Motivation
- Component Releases
  - Configuration
    - Puppet, Foreman and Hiera
    - Punch -> Judy
  - Provision
    - OpenStack
  - Other Services
    - koji , git, jira
- Community Interactions.
- AI as a Production Service
  - Expanding user base

- CERN IT is changing strategy for machine provision and configuration.
- Rationale
  - Need to manage twice as many servers as today
  - No increase in staff numbers
  - Our deployment of configuration tools becoming increasingly brittle.
  - New services take far to long to deploy.
- Approach
  - We are no longer a special case for compute.
  - Adopt open source tool chain model
  - Contribute new function back to community.

- Puppet (2.7)
  - Responsible for configuration, an industry standard.
- Foreman (1.0)
  - Groups hosts into hostgroups of similar configuration.
  - Generates kickstart files from where puppet can take over.
- Hiera (1.0)
  - A data store used by puppet.
- Mcollective (2.2)
  - pub sub messaging to control and query hosts.
- CDB legacy (old)
  - Still some items in CDB... e.g warranty information.

- First puppet infrastructure known as “Punch”
  - One 4 core node, set up “by hand”.
  - puppet, foreman running behind passenger (mod\_ruby)
  - In built own puppetca (cert authority)
  - All project members with root access.
    - Secret files uploaded by hand.
    - Secret files being distributed by puppet
  - Node started to struggle once 400 puppet agents attached - CPU limitation on server.
    - This was with reconfigurations every 15 minutes which is excessive.
- Punch ran for 6 months.
  - Punch was never a scalable solution.

- Punch replaced by Judy in August 2012.
  - All components are deployed with puppet.
  - 2 backend puppetmasters, 2 backend foreman.
  - mod\_loadbalance redirecting requests.



- Using CERN CA.
- CertBaby Service
  - Hooks up users kerberos identity, machine ownership and certificate requests.

- Currently 1200 puppet agents.
  - 500 node added in the last week.
  - 100 a day being added right now.
  - Agents are running on
    - Hardware
    - CVI Service (hyper-v)
    - OpenStack Nova (kvm) (all new ones)
  - Organized in 37 hostgroups with 60 subgroups.
- Adding more puppetmasters or foreman backends is easy.
  - Same problem as scaling web pages, e.g.
    - Number of active connections at redirector.
    - Consistency across back end servers.











- Puppet manifests are very (too?) quick to develop.
  - Takes little longer than configuring the service.
  - e.g an apollo module written in two days.
    - while apollo configuration was being learnt.
  - later parametrization of hardcoded values easy.
- Puppet code to be executed on nodes is distributed by puppet first.
  - i.e no need to package any puppet modules.
  - Makes new feature development, deployment very fast.
- We and others will get better at sharing puppet manifests as hiera becomes normal.



- Git used for puppet modules & manifests.
- Git branches map to dynamic environments
  - local development can be ‘puppet apply’d.
  - admins push changes to a (gitolite) repository
  - puppet masters pull branches and translate to environments
  - Production, Testing & Devel branches
  - Topic branches for major changes
  - Some services live in their own branches
    - risk of divergence...
- Atlassian Crucible & Fisheye for module review process ... not really started.



- Groups hosts of similar configuration.
- Top group -> service. e.g lxbatch, cernfts, ...
- Subgroups may be very different e.g
  - cvmfs/stratum0 vs cvmfs/lxcvms.

<input checked="" type="checkbox"/>	Name	Operating System	Environment	Model	Host Group	Last report
<input checked="" type="checkbox"/>	 lxbsp2701.cern.ch	 SLC 6.3	straylen_ai366	SOLAR 820 S4	base/steves/a	4 minutes ago
<input checked="" type="checkbox"/>	 lxfssm4006.cern.ch	 RedHat 6.3	production	X8DT6	base/steves/a	7 minutes ago
<input type="checkbox"/>	 lxfssm4301.cern.ch	 SLC 6.3	production	X8DT6	base/steves	23 minutes ago
<input type="checkbox"/>	 steve01.cern.ch	 SLC 5.8	straylen_ai366	Virtual Mac...	base/steves	16 minutes ago
<input type="checkbox"/>	 steve03.cern.ch	 SLC 6.3	straylen_ai366	Virtual Mac...	base/steves	17 minutes ago

- Quattor separated code and data well:
  - It was one motivation to write Quattor and drop LCFGng in the first place.
- hiera takes the separation to a new level:
  - puppet asks for a value from hiera?
    - `$myNTP = hiera('ntp_servers')`
  - result can be string , array, hash, ....
  - The lookup is based on a nodes properties, e.g.
    - Since I am at CERN answer is ntp1.cern.ch
    - Since I am in Budapest answer is ntp2.cern.ch
  - The schema of results for CERN nodes, Budapest nodes, SLC5 nodes, debian nodes can be arranged and changed as we please.

# Hiera and Hostgroups

- We arrange nodes in to (sub)hostgroups in foreman.
- A tree of YAML files stored in git maps on to these. e.g for castor hostgroups
  - hostgroup/castor/diskserver/atlas.yaml
  - hostgroup/castor/diskserver.yaml
  - hostgroup/castor.yaml
  - os/slc5.yaml
  - common.yaml
- The files above contain increasingly general keyvalues for look up in hiera.
- Schema and can be fully customized to CERN space with no fear of polluting the code.

```
# A YAML file.
---
castorns: ns.cern.ch
```



- Deploy puppetdb
  - Performance improvements - community raving.
  - Repository for configuration data mining.
- Deploy mcollective
  - Pub and Sub system for sending action commands to hosts.
  - Message broker needs ACLs on queues corresponding to full diversity of CERN hosts and actions.
  - Data mine puppetdb.
- Workflow
  - Move to git pull request process for central configuration.

- Currently Essex code base from the EPEL repository
- Good experience with the Fedora cloud-sig team
- Cloud-init for contextualisation, oz for images with RHEL/Fedora
- Components
  - Nova on KVM and Hyper-V
  - Keystone integrated with Active Directory
  - Glance with Oz
  - Horizon
- Test bed of 100 Hypervisors, 2000 VMs integrated with CERN infrastructure, Puppet



- CERN's Active Directory
- Unified identity management across the site
  - 44,000 users, 29,000 groups
  - 200 arrivals/departures per month
- Full integration with Active Directory via LDAP
  - Slightly different schema from OpenLDAP
  - Aim to minimise changes to AD Schema
  - 7 patches submitted around hard coded values and additional filtering
- Now in use for our pre-production instance
  - Model project definitions in Active Directory
  - Map roles to groups

- We currently use Hyper-V/System Centre for our server consolidation
  - Over 3,200 VMs, 60% Linux/40% Windows
- Choice of hypervisors should be tactical
  - Performance
  - Compatibility/Support with integration components
  - Image migration
- CERN is working closely with the Hyper-V OpenStack team
  - Puppet to configure hypervisors on Windows
  - Most functions work well but further work on Console, Ceilometer, ...



- Deploy into production
  - Target for production is start of 2013 with Folsom
  - Use current grid model running on top of OpenStack
- Deploy multi-site
  - Extend to 2nd data centre in Hungary and disaster recovery
- Deploy new functionality
  - Ceilometer for accounting
  - Bare metal for non-virtualised use cases such as high I/O servers
  - PKI and X.509 user certificate authentication
  - Load balancing as a service
- Deploy at scale
  - Move towards 15,000 hypervisors over next two years
  - Estimate 100-300,000 virtual machines

- CERN presenting to community/vendors.
  - PuppetConf , San Francisco, Sep 2012
  - Openstack Summit, San Francisco Apr 2012
  - Openstack Summit, San Diego , Oct 2012 (now)
  - PuppetCamp, Geneva, July 2012
- CERN has code contributions to:
  - facter, the foreman, puppet, various puppet modules, mcollective, openstack nova, keystone and swift.
  - This is increasing as new students/fellows are employed for their puppet, ruby, .. skills.
- CERN puppet-users meeting , IT, ATLAS pit, ..
- Share our own <http://github.com/cernops>

- Agile is not just Puppet and Openstack.
- AI created a gitolite ACL'ed GIT service.
  - CERN IT is now provisioning a public GIT service based on this.
  - AI will migrate its projects ASAP.
- AI created a Koji service for RPMs.
  - Creates RPMS and publishes to yum.
  - The service is now being used by others with in IT. e.g castor builds, data management, lemon, ...
- AI ran jira early before a central service was created.
  - AI already migrated to central service.

- Several Services running now on AI.
  - Some CVMFS components.
  - SLC6 batch services
  - SLC build machines
  - GIT gateways.
  - CASTOR (compass VO)
  - Test systems, glusterfs, swift, ..
  - New top level hostgroups every week now.
- From November AI opening up more.
  - Experiment services (voboxes) will start to use AI service.
  - Documentation to be updated/consolidated.

- Agile Infrastructure Project
- We are ready for hardware arriving in Budapest in 2013.
  - Puppet configured VMs on Puppet configured OpenStack.
- Documentation:
  - More user facing documentation needed.
- Configuration with Puppet:
  - Services needing knowledge of everything
  - Inter sysadmin trust.
  - Test facility for AI.
- OpenStack deployment
  - Increase scale.

- AI Project Pages: <http://cern.ch/go/7vFF>
- CERN modules <http://github.com/cernops>
- CERN agile tickets <https://agileinf.its.cern.ch/jira>
- AI Presentations : <http://cern.ch/go/6qRG>