# MONITORING HPC SERVICES WITH CHECKMK

Monitoring practice at EPCC

Presenter: Craig Manzi

Authors: Philip Cass, Kieran Leach

THE UNIVERSITY of EDINBURGH      SC19 Denver, hpc CO is now.      |epcc|

---

## Problem statement

- EPCC manages a variety of HPC and research computing services in addition to critical support infrastructure.

- EPCC sysadmins spent a lot of time tracking the state of various systems; problem detection and diagnosis typically requires looking in multiple locations:
  - Time intensive, difficult and requires a constant wide awareness.
  - Difficult to effectively diagnose new systems where team members are typically under pressure to get things up and running in short timeframe.

- We needed a "single pane of glass" approach.

|epcc|

---

## Our solution – CheckMK and locally developed checks

- Originally a Nagios extension, now a Nagios derivative monitoring system.
- Many checks (both Nagios and CheckMK) available already.
  - CPU, Memory, Filesystem, Interface status etc.
- Simple to create new checks
  - DDN controller, Lustre, PBS etc.
- Very simple to add new hosts, and can alter check parameters from the central user interface
- CheckMK server first installed at EPCC in 2015. Now core to our service management for HPC services.
- Since 2015 this has allowed us to provide bespoke integrated monitoring solutions for a variety of HPC technologies.

https://checkmk.com

|epcc|

---

## Panopticon system setup

- CheckMK server hosted on VM – named Panopticon.

- Each service on site has at least one monitoring and management VLAN on our site management network (current network design has little inter-VLAN routing).

- CheckMK VM is presented all monitoring and management VLANs and granted an IP on each.

- Panopticon presently has 18 virtual interfaces for monitoring various systems and services on site.

|epcc|

## Panopticon networking



Office switching

vlan 10

vlans 1,2,3,10

Management Network

Panopticon

vlan 1    vlan 2    vlan 3

ARCHER    Tesseract    RDF

|epcc|

## Day-to-day use

- Monitored during working hours by "on-shift" team member
  - Numerous issues caught this way including partial power failures, system, disk and component failures and networking issues.
- Monitorable over email out of hours
- Used as a "first port of call" for investigating issues reported by users or other team members
  - At-a-glance view of system status.

|epcc|

## Panopticon user interface



|epcc|

## Panopticon user interface

- Delegable access: Partners and vendors can access limited views of only the items we select.
- Controllable email alerts: Email alerting can be managed and pruned on a per user basis.
  - Example: Cray staff supporting ARCHER have beepers which support email. A curated set of alerts will email the relevant address for immediate attention 24/7.
- Customisable views for different people/purposes.
  - Example: Each sysadmin office at EPCC has a screen dedicated to Panopticon using a customized view (large font, less detail, no unnecessary menus etc).

|epcc|

## CheckMK on the Client

- CheckMK agent script runs on the node, conducts most checks when polled by the server but also runs more heavyweight checks in the background.
- Traditionally uses xinetd or systemd socket – can also configure the server to query via SSH or any custom command.
- Packaged install (rpm/deb) available.
- For both ARCHER (4920 node Cray XC30) and Tesseract (1476 node HPE SGI 8600) we have installed on all management and login nodes.
- Installed throughout the RDF (23PB DDN) storage system and within the attached Data Analytics Cluster.
- Pervasively installed on all support systems (authorisation, web, management servers etc.)

https://www.archer.ac.uk

|epcc|

## "Out of the box" checks

- Checks available by default include:
  - CPU, Memory, disk utilisation and load.
  - IPMI checks (fans, temperatures, voltages).
  - Network interface status and statistics.
  - File system mounts.
  - Individual processes can be assigned for monitoring (e.g. pbs_mom or license servers).
  - Number of users logged in
    - We have no need to alert on this, but we installed the check anyway... and 18 months later we were asked to provide statistics for how many concurrent users we had on one of our systems
- Server can also directly query clients using e.g. SNMP
  - We host out-of-band interfaces for switches, tape libraries etc. on management vlans visible to panopticon. Not quite "one click" monitoring, but not far.
  - Can also check SSH servers, HTTP(s), DNS, LDAP etc.

|epcc|

## Imported checks

- Some checks we need have already been developed elsewhere
- We had a specific requirement for monitoring RAID controllers in a storage system.
- Investigation online showed that this had already been achieved and was available in a published repo.
- These checks are now in place and have been extremely useful.
- We have also used (and partially adapted) some lustre server performance checks

|epcc|

## Custom in-house checks

- Developed in whatever language is appropriate.
  - Some bash, some python.
- Method used in checks presented here is to place correctly formatted output into /var/lib/checkmk/spool or to place the check in /var/lib/checkmk and have CheckMK run the check directly.
- Checks must output 1 service per line for monitoring.
- Service line details status, service name, metrics and status details (plaintext).

|epcc|

## Developed checks



```
0 node_power_usage watts=156 mpower reports 156W
```

1        2        3        4

- **1**. Service status (0 OK, 1 WARN, 3 CRIT, 4 UNKOWN)
- **2**. Service name
- **3**. Metrics
- **4**. Status details – free text passed back with the check.

|epcc|

---

## Developed checks – Tesseract node check



|epcc|

---

## Developed checks – Tesseract node check

- When Tesseract arrived we had an immediate requirement for monitoring of compute node availability.
  - During handover training the "cpower" command showed how a basic health check could be put together.
  - "cpower node status r*i*n*" output was easily mangled into an appropriately formatted status report.
  - Appropriately formatted status report can then be dropped into an appropriate folder and picked up automatically.
  - Also makes use of CheckMK "piggybacking" - script run on one host but provides check output for other hosts

|epcc|

---

## Developed checks – Tesseract node check



```
[root@tesseract-sac ~]# head -n 50 /var/lib/check_mk_agent/spool/420get-node-power
<<<<ts-r1i0n0>>>>
<<<local>>>
0 node_power_usage watts=201 mpower reports 201W
<<<o>>>
<<<<ts-r1i0n1>>>>
<<<local>>>
0 node_power_usage watts=200 mpower reports 200W
<<<o>>>
<<<<ts-r1i0n2>>>>
<<<local>>>
0 node_power_usage watts=192 mpower reports 192W
<<<o>>>
<<<<ts-r1i0n3>>>>
<<<local>>>
0 node_power_usage watts=190 mpower reports 190W
<<<o>>>
<<<<ts-r1i0n4>>>>
<<<local>>>
0 node_power_usage watts=191 mpower reports 191W
<<<o>>>
```

|epcc|

## Developed checks – Tesseract Omnipatch checks

- Tesseract was our first system on site with an Omnipath network.
- Tesseract was not delivered with any particular setup for monitoring the fabric.
- Investigation showed that built-in tools (opareport) could provide a great deal of data but lacked context.

|epcc|

## Developed checks – Tesseract Omnipath checks

- Manually runnable checks were put together fairly quickly.
- The lack of a specification of how the fabric *should* look held things back initially.
  - Remonstrations with the vendor eventually got us this.
- Upon setting up and running a full verification of the fabric we were able to show the vendor that a number of links had been mis-cabled.
- Developed code for these checks was easily adapted into a CheckMK check.
- Bonus outcome: after implementing this we were able to show that OPA speed downgrades could occur on **neighbouring** nodes in the shared backplane whenever a blade was reconnected to the system

|epcc|

## Developed checks – Tesseract Omnipath checks



|epcc|

## Developed checks – Tesseract Omnipath checks



|epcc|

## Additional checks developed at EPCC

- DDN controller monitoring
- GPFS disk status
- SELinux/AppArmour status
- Vulnerability monitoring (Meltdown/Spectre etc.)
- Compute node status, power consumption, stuck processes
- Mainframe cabinet environmental status
- Unplaceable and orphan jobs
- Lustre servers read/write/metadata statistics, LNET statistics

|epcc|

## Facts & figures

- Directly query ~300 hosts
- 40,000 services
- Close to 2,000 hosts in CheckMK – the scripts mentioned previously provide "piggyback" information for compute nodes and OPA switches
- Hosted on dual-socket E7-4820 (2011, 16 cores, 2.00GHz) and close to capacity, especially if an issue results in 1,000+ events
  - Main issue is process forking - the paid version of CheckMK uses a different core which solves this problem
- Using about 4GB of RAM and 40GB of disk space; disk I/O is light

|epcc|

## Approach for new systems

- Initially, we concentrate on administration, service and front-end servers – these are easier to set up and we might get enough information about the compute nodes from them
- Add all applicable checks we know of – don't wait until we have a use-case (as long as they don't affect performance)
- If there's no available check for disk health, prioritise developing this (even if all it returns is a single good/bad status for the entire system)
- Any time we encounter issues, ask if it can be easily detected in the future, either by a custom check or especially using the built-in process or TCP checks

|epcc|

## Next steps

- Our single monitoring server is reaching capacity
  - Multiple CheckMK instances can report to, and be managed by, a single central server.
- Our current version is old (1.2 vs 1.6)
  - Agent format has not changed – we should be able to upgrade the server without issues and roll out agent upgrades over time
- Custom SNMP checks for tape systems, OPA switches etc.
  - Not complicated but more work than simple local checks

|epcc|

## Conclusions

- CheckMK centralised monitoring has been immensely useful.

- As experience has grown we have been able to rapidly deploy monitoring for new systems and to integrate site-developed monitoring for novel or boutique technologies.

- This has had demonstrable benefits for service success.

- Sharing of checks developed for HPC technologies would benefit the community.

|epcc|

## Thanks to...

- UKRI, EPSRC, NERC, DiRAC

- HPC Systems Team at EPCC

- The HPE Support Team

- Neelofer Banglawala

Thanks for listening!

|epcc|