Maybe 3.0, or possibly
4.0. Could even be 5.0…

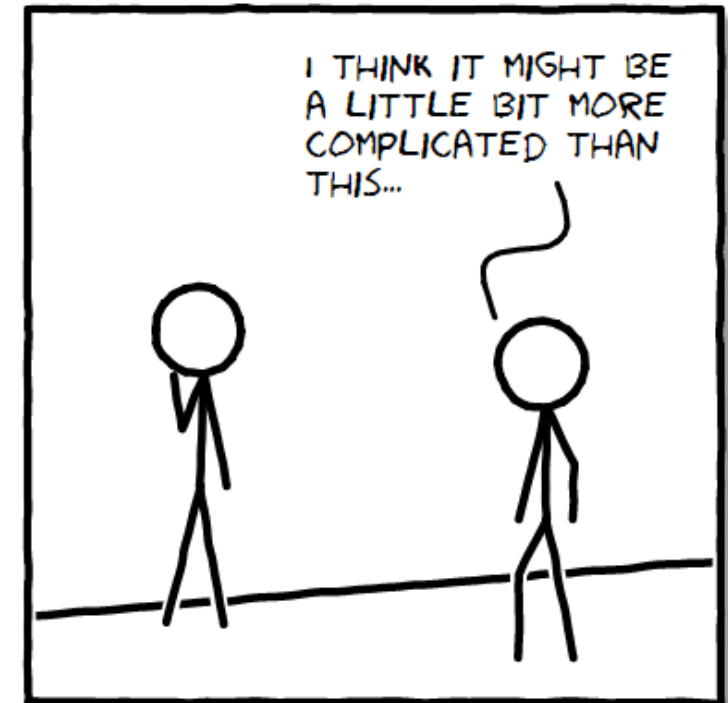Glasgow
# ScotGrid^2.0

ScotGrid Face-to-Face

Edinburgh

October 2018

# Data Centre Migration: Early Planning

1. Build data centre

2. Install kit in data centre

3. Mumble mumble configuration blah blah blah

4. Turn it on

http://cmx.io

INTERNAUL UI FOR LOCAL DEV

- freeipa /accounts
- logging - secure
    " - nats?
- firewalls
- f2b /~~tripwire~~/etc...
    aide
- oscap/lynis
- investigate SLURM
- prometheus /alerting
- monit?
- ansible/git host
- DOCS - BACKUP
    RETIC + CRON
    ?

croquembouche 10.1.10.1

OVIRT
- EXT INTF.

REBUILD
node065
node188

(1) LIST SERVERS TO DECOMMISON
- disk037
- svr001
- svr015
- provision

IPV6 Rollout
- ENABLE PERF/SONAR
    => LOOK AT F/W CONF

VM HOST
SNICKER    10.1.20.1
DOODLE     10.1.20.2
KRUMKAKE   10.1.20.3

VPN
ROCKYROAD  130.209.231.125
           10.1.10.2

node274
node275

GRIPP 41
to DAVE/JEREMY NEW SITE??

PDU/UPS   1
SWITCH    2
RAID      3
GPU       4
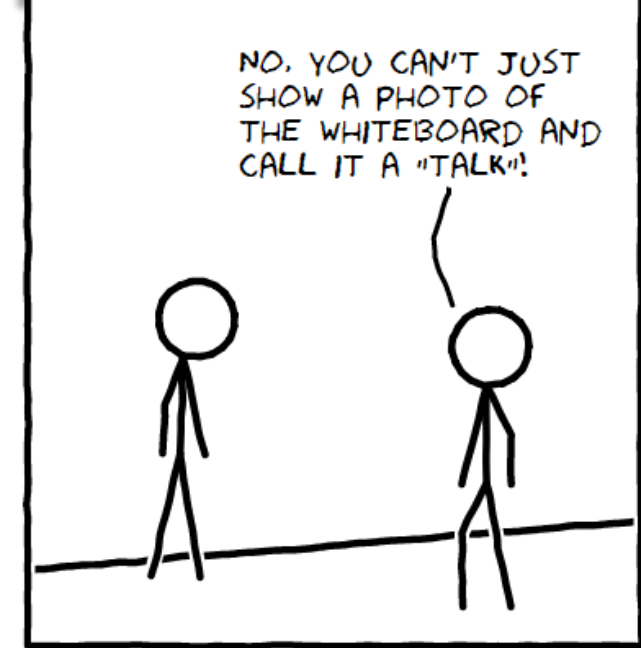10.0.

10.0.X.X IPMI

GRID    VPD        30
        SERVICES   40
        STORAGE    50
        COMPUTE    60
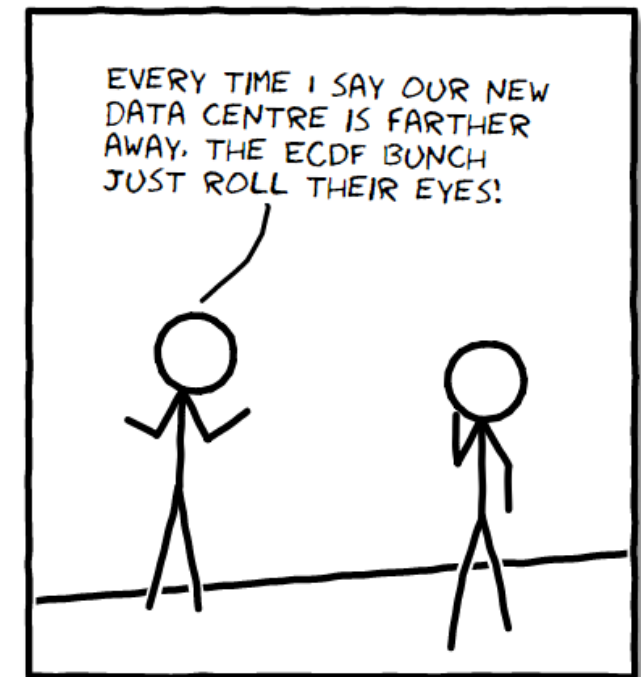INFRA   HYPERVISORS 20
        META       10
                        } 10.1.X.X

TWO WEEKS AGO...

NO. YOU CAN'T JUST SHOW A PHOTO OF THE WHITEBOARD AND CALL IT A "TALK"!

# ScotGrid 2.0

- Data centre move offers perfect opportunity to redesign site from the ground up

- Change in working practices: the cluster will no longer be just down the stairs

- Many important questions to answer:
  - How do we deploy and configure systems?
  - How do we manage jobs?
  - How do we manage storage?
  - How do we check that everything is working?
  - How do we secure it?

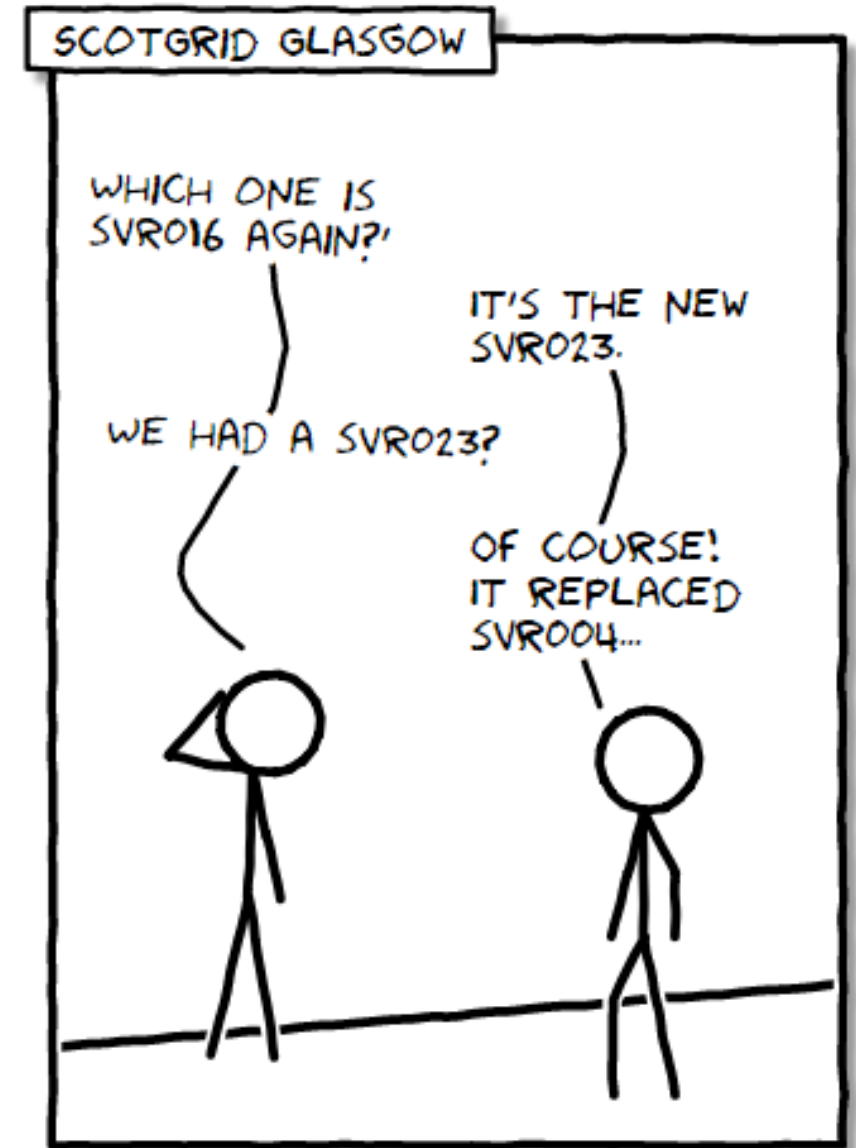- Most important question to answer:
  - What do we call it?



EVERY TIME I SAY OUR NEW DATA CENTRE IS FARTHER AWAY, THE ECDF BUNCH JUST ROLL THEIR EYES!

# Legacy Names

- Currently use an imaginative scheme:
  - `provision`
  - `svr000 – svr031`
  - `disk032 – disk089`
  - `node001 – node282`
  - `nat005 – nat007`
    - *Almost makes sense, except this isn't a NAT, it's a Squid…*

- Some "advantages" (e.g. can obtain certs without knowing what the machine will do) but can we do better?

- Decided to look to market leader for inspiration regarding comprehensive, distinctive, logical, memorable scheme…



SCOTGRID GLASGOW

WHICH ONE IS SVR016 AGAIN?'

WE HAD A SVR023?

IT'S THE NEW SVR023.

OF COURSE! IT REPLACED SVR004…

# ScotGrid 2.0 Names: Cake*

**Provisioning / Configuration Management**

croquembouche

**VM Hosts**

snicker

doodle

krumkake

**Remote Access (VPN)**

rockyroad

*\* Including sweet baked goods and confectionery*

# Old Address Scheme

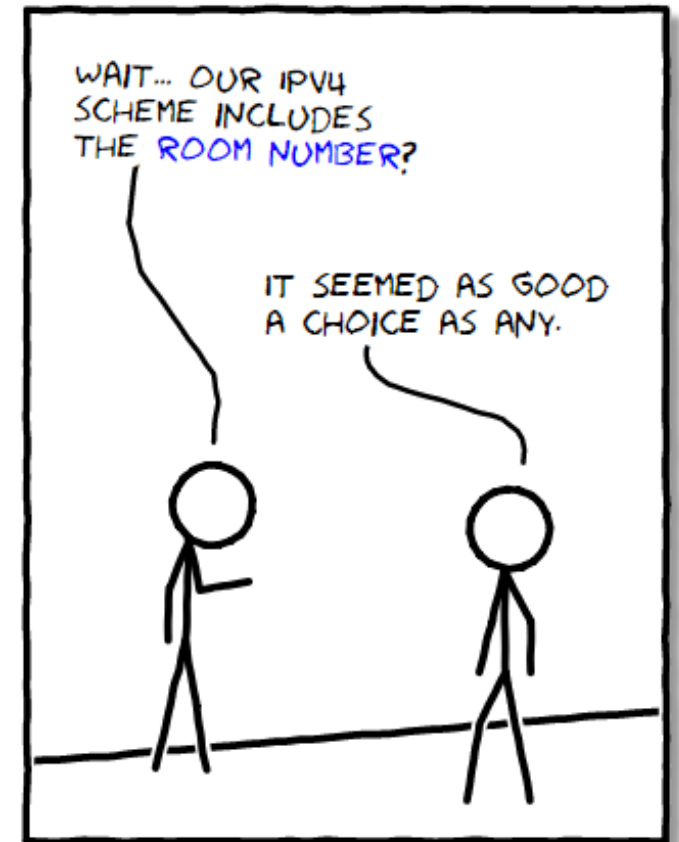- Everything in 10.141.0.0/16

  - `provision → 10.141.100.1`

  - `svrN → 10.141.255.(0 + N)`

  - `diskN → 10.141.245.(0 + N)`

  - `natN → 10.141.246.(0 + N)`

  - `nodeN → 10.141.(floor((N – 1) / 253)).(((N – 1) % 253) + 1)`

*Apart from svr000, which is actually svr016 in disguise*



WAIT... OUR IPV4 SCHEME INCLUDES THE ROOM NUMBER?

IT SEEMED AS GOOD A CHOICE AS ANY.

# New Address Scheme

`10.<NETWORK>.<TYPE>.N`

| | |
|---|---|
| 0 | IPMI / Management |
| 1 | Primary |

| | |
|---|---|
| 1 | PDU / UPS |
| 2 | Network infrastructure |
| 3 | RAID / Storage |
| 4 | Environmental |
| 10..19 | Bare metal servers |
| 20..29 | Hypervisors |
| 30..39 | VPN |
| 40..49 | Services |
| 50..59 | Storage |
| 60.. | Compute |

*Haven't decided on this bit, but my vote is for .pasticceria*

```
DNS (.beowulf.cluster)

10.1.10.1     croquembouche
10.1.20.1     snicker
10.1.20.2     doodle
10.1.20.3     krumkake
10.1.30.1     rockyroad
...
```

# IPv6

- Doesn't this all become irrelevant when we move to IPv6?
    - It's taken over 20 years to get this far – we're not going to drop IPv4 before next summer
    - Some IPv6 addressing schemes incorporate IPv4 address

Still need to plan ahead so we get this...

...and not this

# Networking

- 4 × Lenovo G8272: 48 × 10 Gbps SFP+, 6 × 40 Gbps QSFP+)

- 1 × Lenovo G8332: 32 × 40 Gbps QSFP+

- Total: 192 × 10 Gbps, 56 × 40 Gbps

- Aggregated throughput: 8.32 Tbps

GRATUITOUS USE OF IMPRESSIVE-SOUNDING YET MEANINGLESS STATISTICS: CHECK!

# Provisioning / Configuration Management

- PXE deployment: Cobbler? *But its name fits!*
  - Use at present, but concerns regarding continued support and development
  - PPE using bespoke alternative
    - *Because what you really want to do if you're worried about support is roll your own…!*



- Configuration management: Ansible
  - PPE switched 2.5 years ago
  - Making increased use within ScotGrid, particularly for ad hoc tasks
  - Plan to switch entire configuration management to Ansible, other than in specific cases where another tool is required (e.g. third-party Puppet modules)
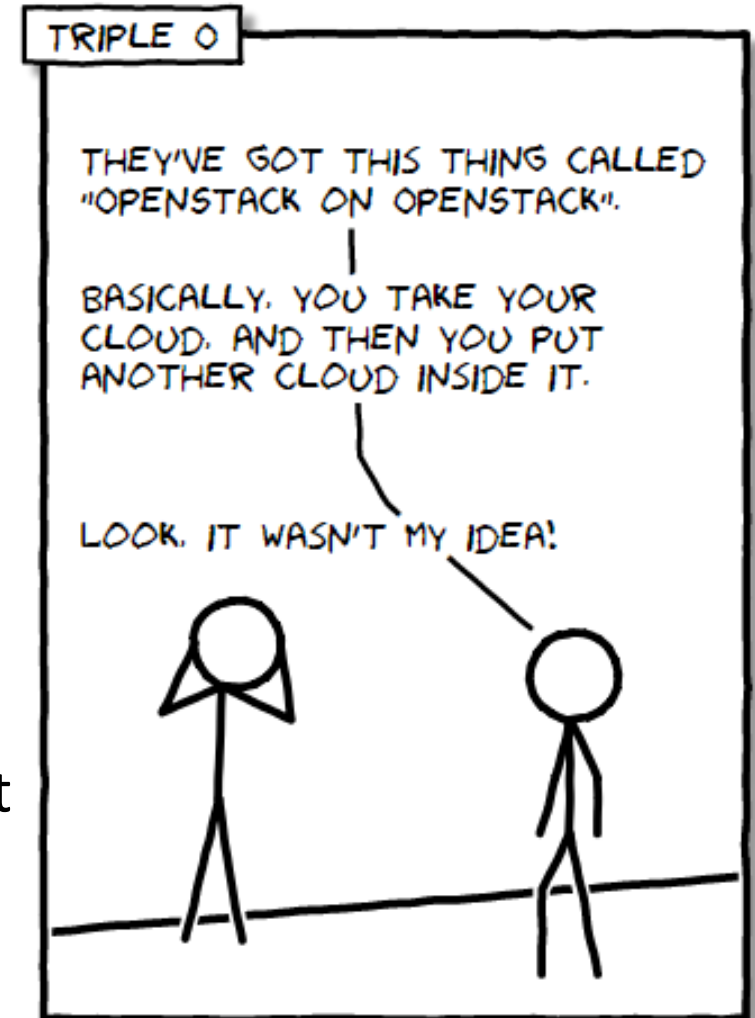
# Virtualisation

- libvirt / KVM
  - Simple but limited features

*Our current approach, also used by PPE*

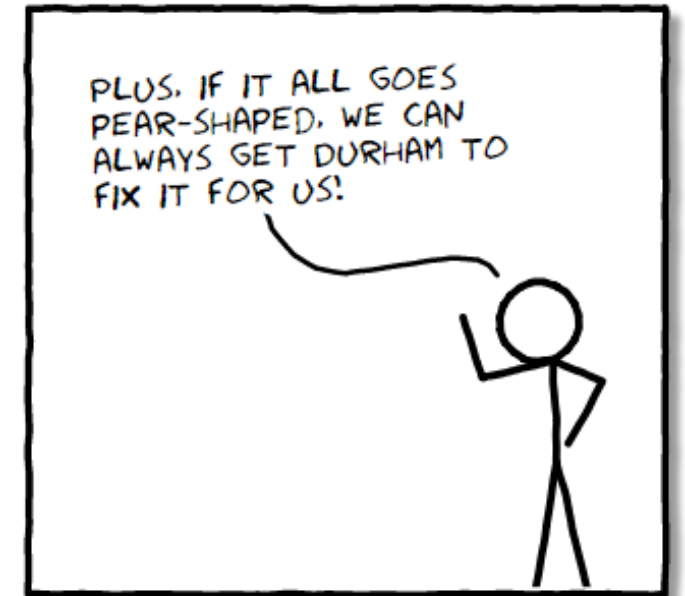- oVirt
  - Complex but feature-rich

- OpenStack
  - Spent about a year (2015) investigating OpenStack
  - Gareth had another look recently
  - Offers greatest flexibility, but is incredibly complicated and would require significant investment of time
  - Unless you truly need a multi-tenant cloud, is it worth the pain?

*But they have cake!*



TRIPLE O

THEY'VE GOT THIS THING CALLED "OPENSTACK ON OPENSTACK".

BASICALLY, YOU TAKE YOUR CLOUD, AND THEN YOU PUT ANOTHER CLOUD INSIDE IT.

LOOK, IT WASN'T MY IDEA!
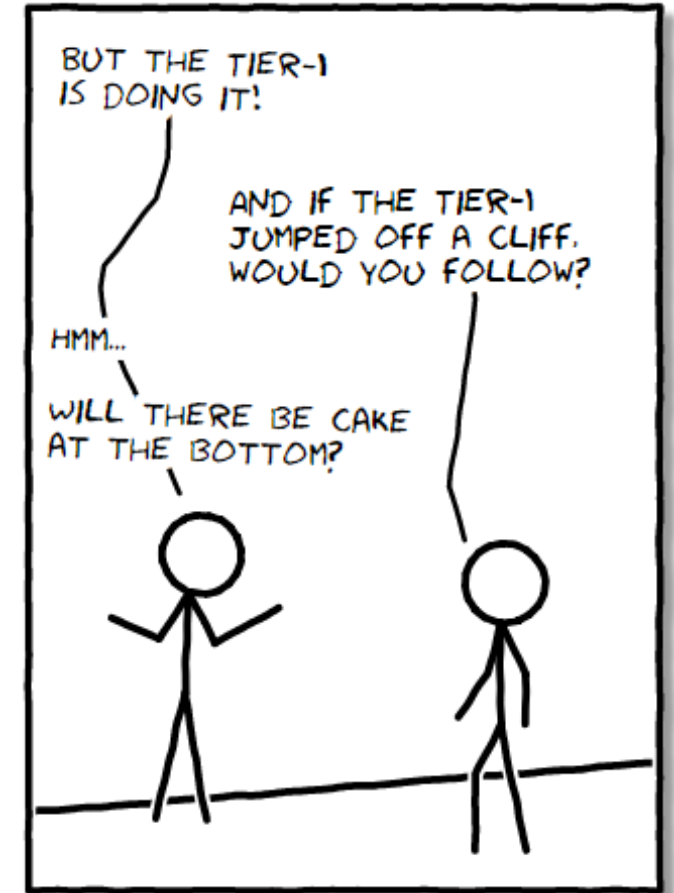
# Identity Management

- Currently use local accounts
  - Compute nodes all have pool accounts (4,663!) defined locally
  - Admins tend to SSH from one place to another as root
  - No central management
  - No audit trail
  - Compute node deployment needlessly lengthy
- FreeIPA
  - Backed by Red Hat
  - Mature toolset, including Web UI

PLUS, IF IT ALL GOES PEAR-SHAPED, WE CAN ALWAYS GET DURHAM TO FIX IT FOR US!
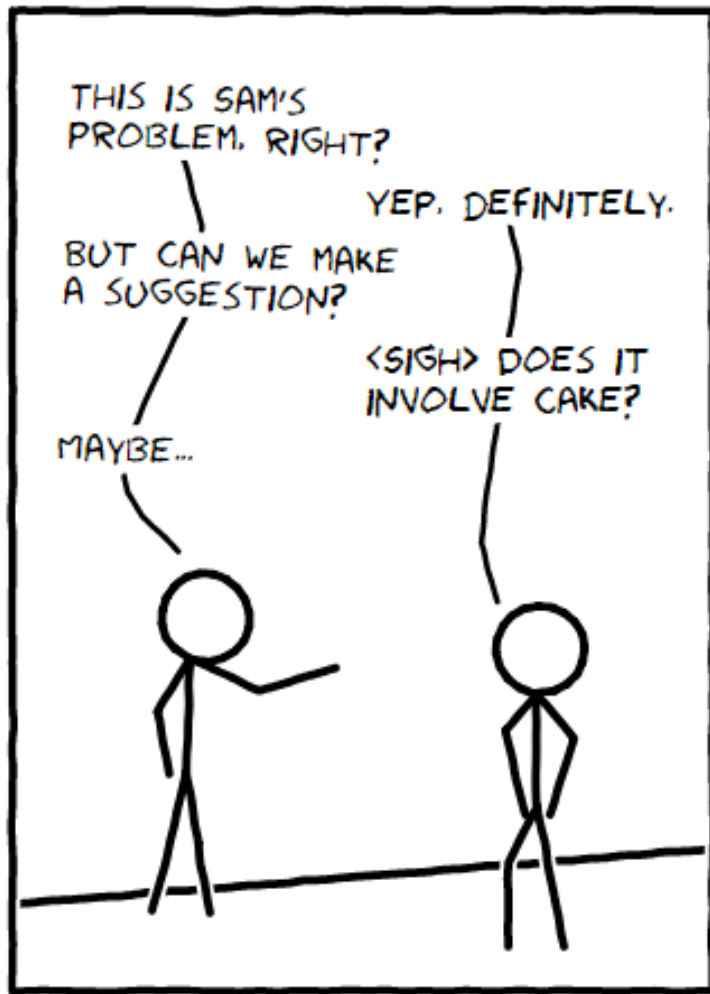
# Batch Systems

- Currently use HTCondor
  - We like it, but our configuration has grown over time, is now overly-complicated and needs to be rewritten from scratch

- SLURM
  - If we're starting again anyway, why not consider alternatives?

- However…
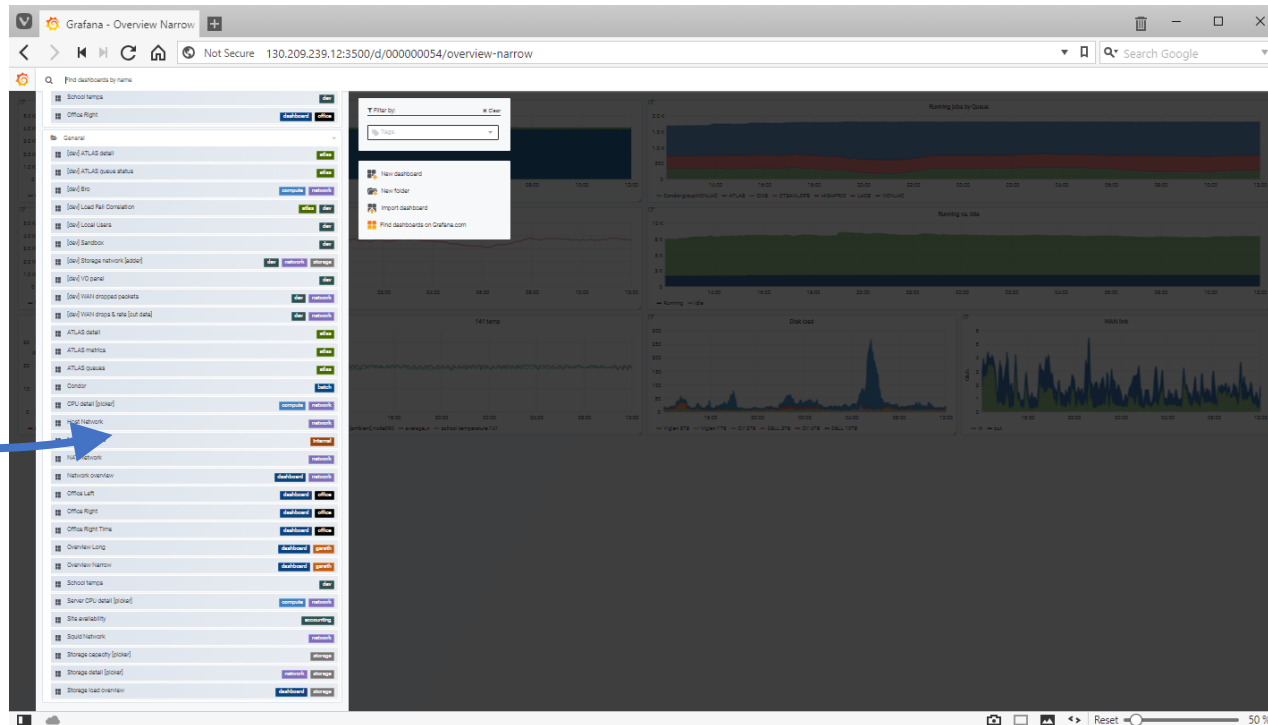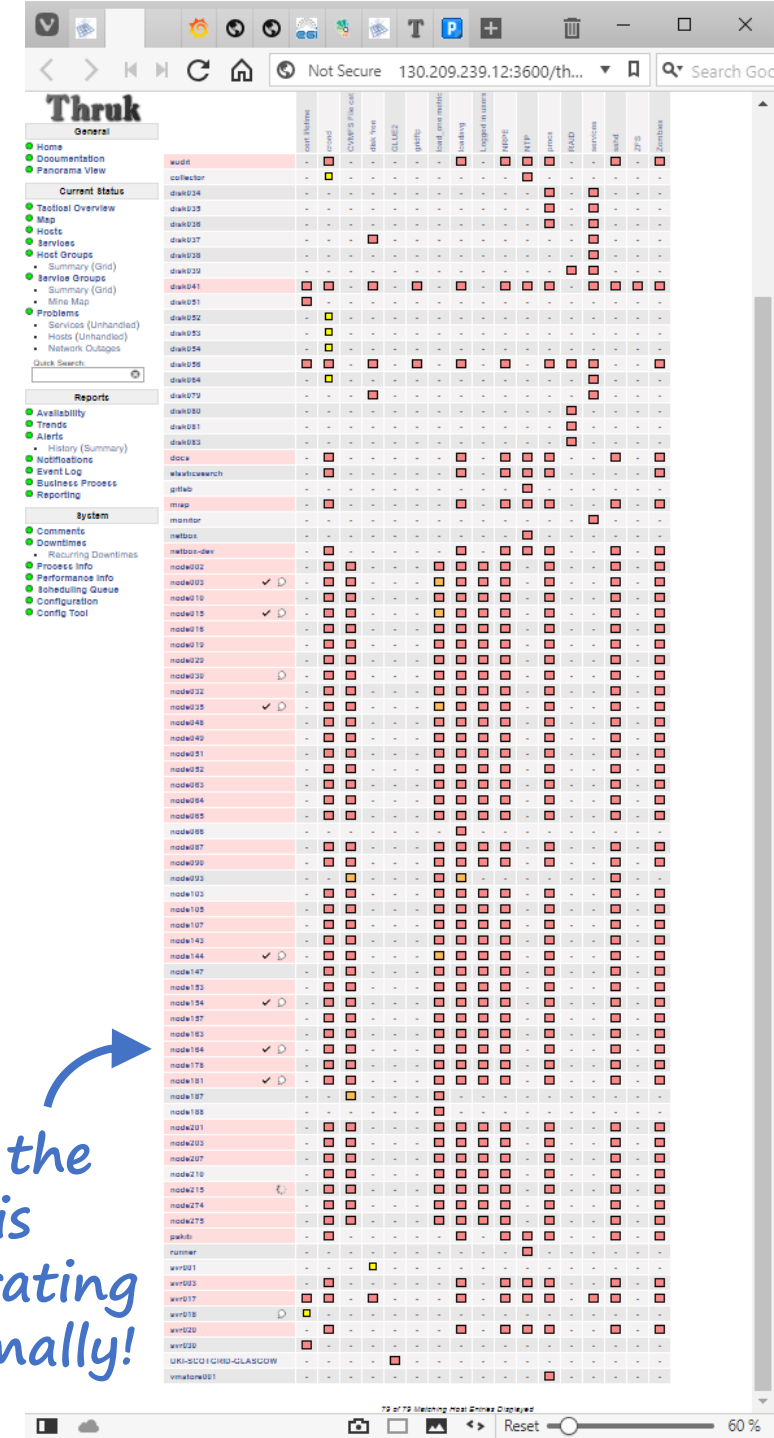  - RAL is fully invested in HTCondor – politically it may be wise to follow

# Storage

# Monitoring

- We have ~~lots of~~ *too much* monitoring
  - Some has fallen into disrepair since Dave switched focus to security (and moved to STFC!)
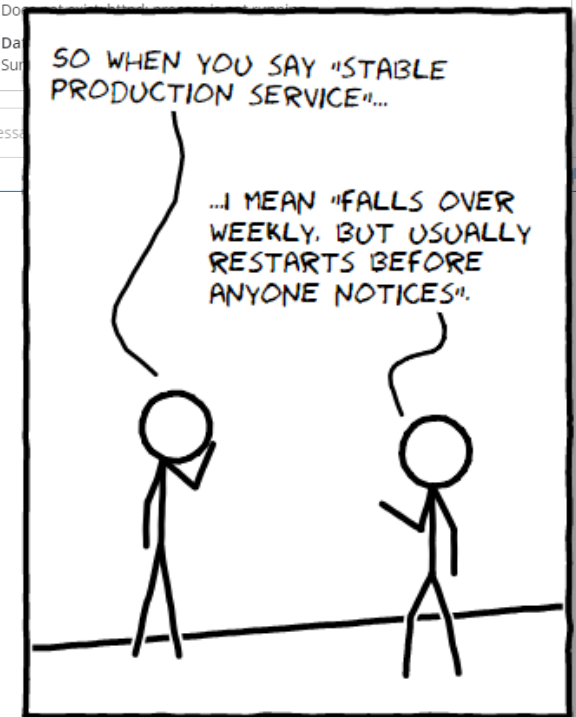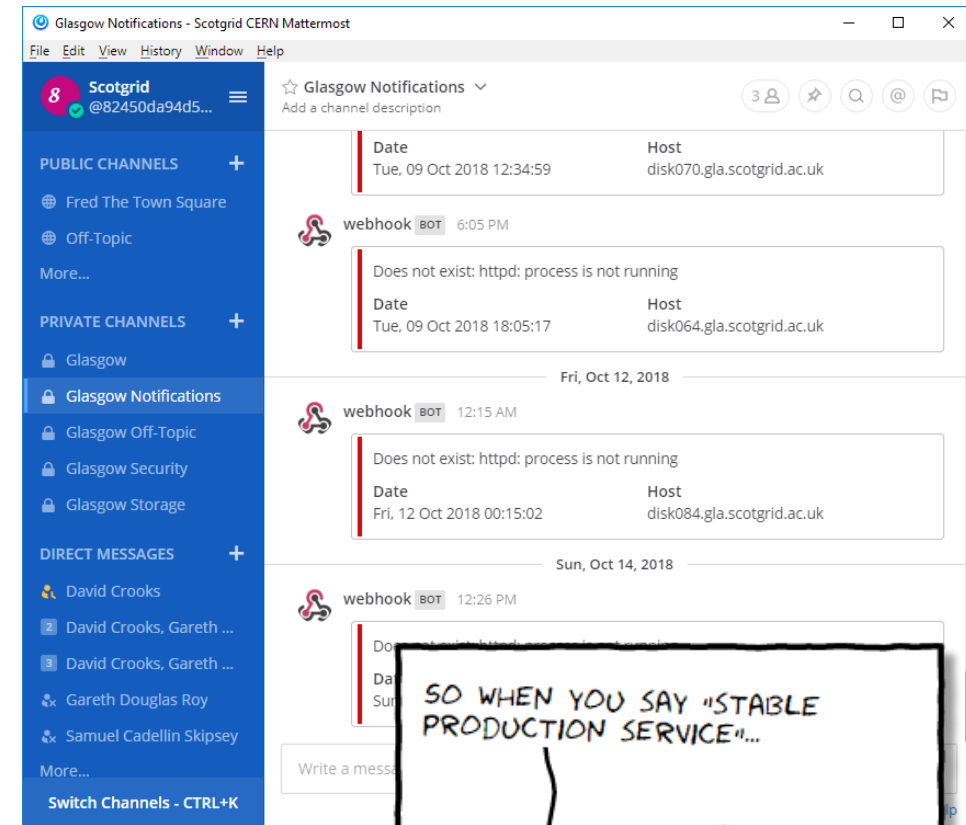
- Identify what we need and remove the rest?

*Why have one dashboard when you can have 31?*

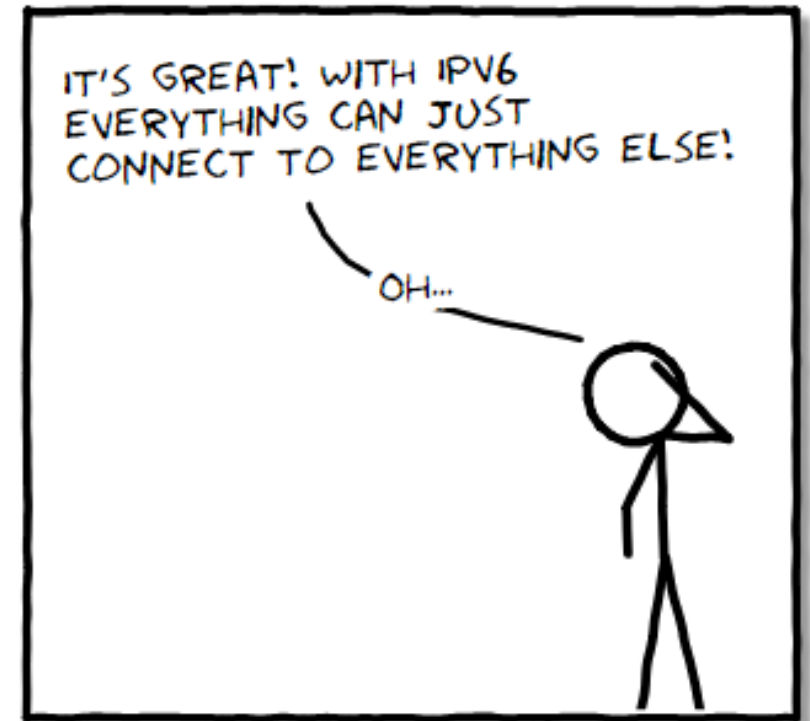*Yes, the site is operating normally!*

# Monitoring

- Prometheus
  - Gareth has been looking into this
  - Running on many systems, both on grid and in PPE
  - Particularly useful for investigating performance issues:
    - ZFS grid storage
    - Jobs on local batch system

- Monit
  - Acts as watchdog
  - Already used to monitor parts of the storage
  - Probably don't need to be notified every time it does something!

- Alerting

# Security

- Firewalls
  - Presently, we don't really have any
  - More important with move to IPv6
- Fail2ban
- AIDE
  - Advanced Intrusion Detection Environment
- OpenSCAP
  - Security Content Automation Protocol
- Lynis
  - Security auditing

# Miscellaneous

- Centralised syslog?
  - Currently dumps to the console, which is incredibly annoying!

- Back-up
  - restic + cron?

- Git
  - How many bells and whistles do we need?

- Documentation
  - We should have some!
  - We always moan about wikis, and then usually decide that anything else involves far too much effort
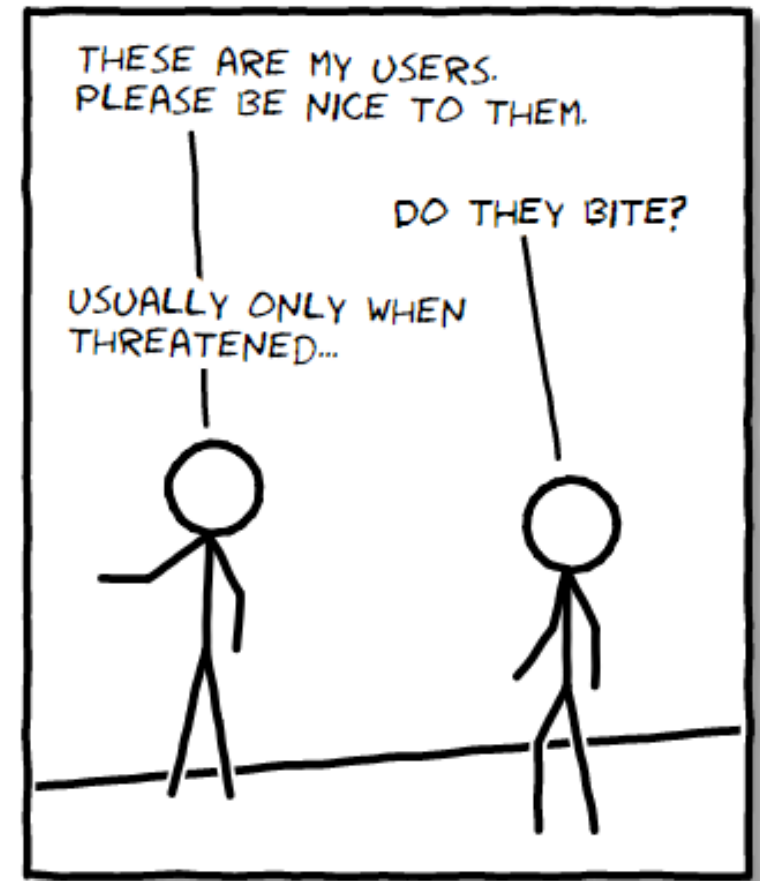
# Internal UI

- Continuation of effort to merge tier-2 and tier-3 resources
  - We're making good progress!
  - "Key influencers" now indoctrinating new RAs and PhD students
- Increase adoption of grid by local users by providing entry point
- Offer some resources for interactive or rapid testing of jobs in grid environment to simplify development work

*I talked about this last year*

# Summary

- Data centre move offers unique opportunity to completely redesign Glasgow ScotGrid site

- Many questions about best solutions remain to be answered

- Time is short!

- You can never have too much cake

# Summary

- Data centre move offers unique opportunity to completely redesign Glasgow ScotGrid site

- Many questions about best solutions remain to be answered

- Time is short!

- You can never have too much cake

PLEASE TELL ME YOU DIDN'T SPEND THE WHOLE WEEK DRAWING CARTOONS AND GOOGLING FOR PICTURES OF CAKE?

ERM...