# OpenHPC Automation with Ansible

**Baptiste Gerondeau, Takeharu Kato, Renato Golin**

ISC18, Arm Workshop, 28th June 2018
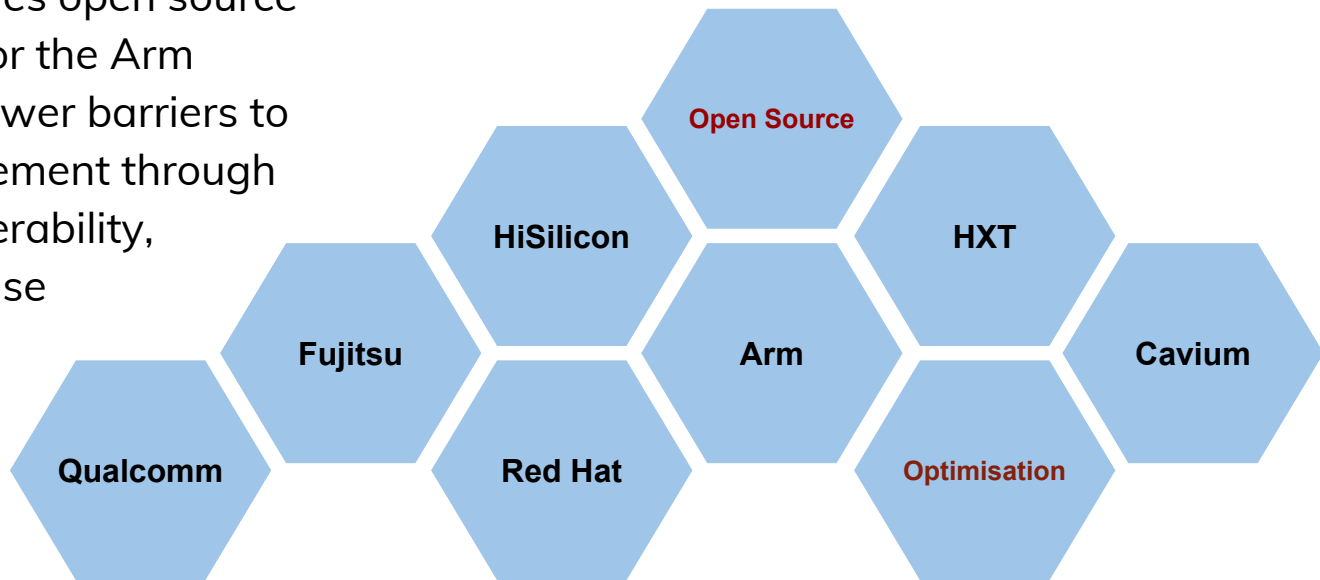
Linaro

Linaro HPC

# Agenda

- Linaro's HPC-SIG Lab
- OpenHPC Ansible Automation
- Results

# The HPC-SIG Lab

# Linaro High Performance Computing
## Special Interest Group

**The Linaro HPC SIG** drives open source software development for the Arm architecture. It aims to lower barriers to deployment and management through standardisation, interoperability, orchestration and use case development.
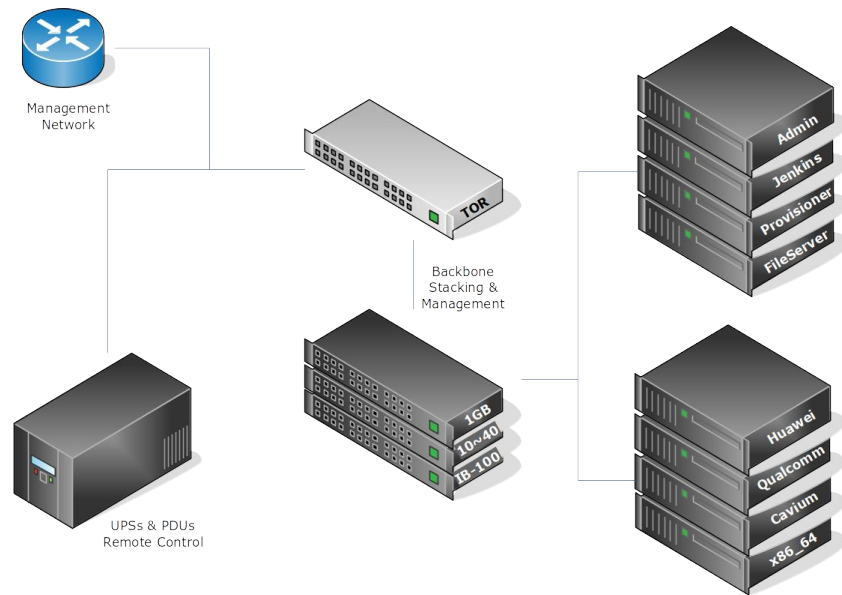
**linaro.org/hpc**

# The HPC-SIG Lab

Goals:

- Cluster Automation & Validation
- Benchmarking & Performance Investigation
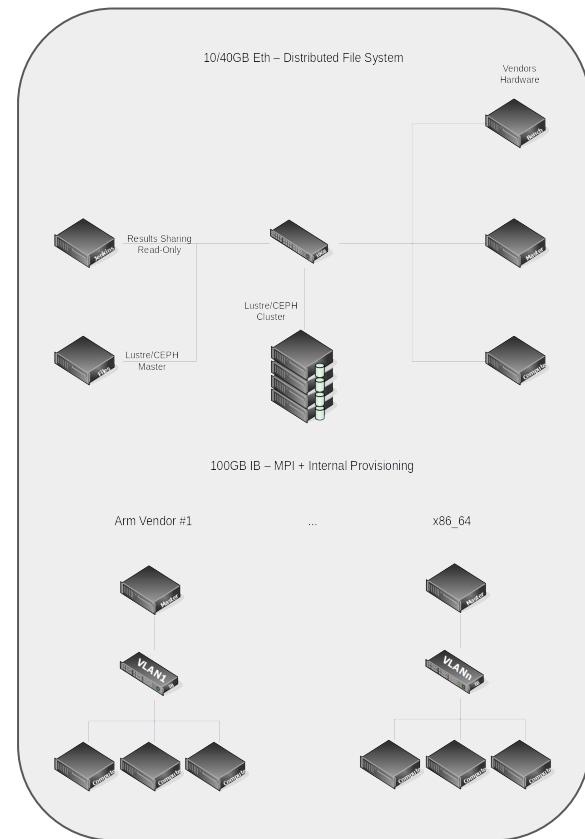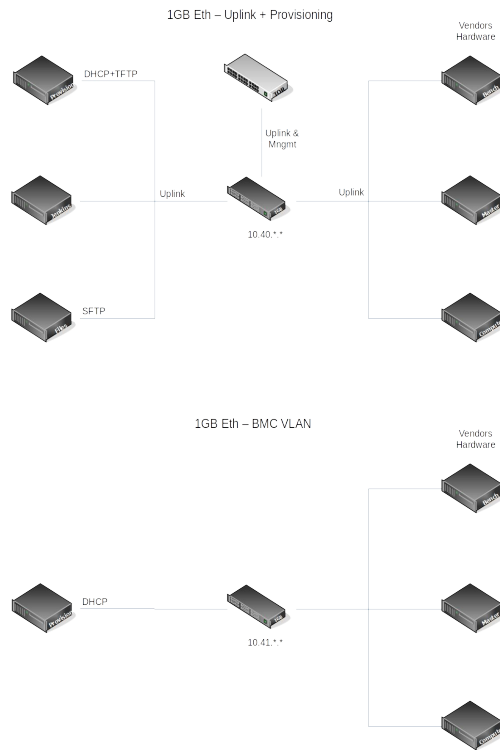
Requirements:

- Stability & Repeatability
- Close-to-production environment
- Upstream technology (reproducibility)
- Vendor isolation (hardware, results)

# The HPC-SIG Lab

Network Layout

- Flat Ethernet
  - Uplink/Provision
- BMC subnet (VLAN)
- File System subnet
  - 10/40GB Ethernet
  - Lustre/Ceph (future)
- MPI subnet
  - 100GB InfiniBand
  - Slave provision

# The HPC-SIG Lab

- Stability & Repeatability
  - Critical external components cached locally
  - Strict migration plans (staging)
- Close-to-production environment
  - Hardware and firmware updated frequently
- Upstream technology (reproducibility)
  - All components are open source / upstream
- Vendor isolation (hardware, results)
  - VPN, Provisioner, Jenkins, SSH control

# Open Source and Upstream

**Lab admin & tools:**

- Jenkins: https://jenkins.io/
- Mr-Provisioner: https://github.com/Linaro/mr-provisioner

**Lab Automation:**

- https://github.com/Linaro/ans_setup_jenkins
- https://github.com/Linaro/mr-provisioner-role
- https://github.com/Linaro/mr-provisioner-kea-dhcp4-role
- https://github.com/Linaro/ansible-role-mr-provisioner
- https://github.com/Linaro/mr-provisioner-client

**HPC-specific automation:**

- https://github.com/Linaro/hpc_lab_jenkins
- https://github.com/Linaro/hpc_deploy_benchmarks
- https://github.com/Linaro/benchmark_harness
- https://github.com/Linaro/ansible-playbook-for-ohpc

# Test Suite

Most tests green, however:

- Intel-specific tests (CILK, TBB, IMB) disabled
- Others need package install (PDF, CDF, HDF) but pass when installed
- TAU fails because LMod defaults to openmpi (needs openmpi3)
- Lustre fails as package depends on kernel 4.2 (which won't work on our machines)
- MiniDFT and PRK had make failures, but we haven't investigated yet
- `--enable-long` doesn't really, need to look into why not

The plan from now on is:

1. Automate package install conditional on enabled tests, fix remaining errors
2. Work with members to prioritise long term ones (like Lustre)
3. Use it for *additional* packages, so we can test them before sending upstream
4. Add a benchmark mode, making sure to use entire cluster

# OpenHPC Ansible Automation

# Existing OpenHPC Automation

- A recipe with all rules described in the official [documents](#)
- LaTex snippets containing shell code

```
% begin_ohpc_run
\begin{lstlisting}[language=bash]
[sms](*\#*) (*\install*) lmod-defaults-gnu7-openmpi3-ohpc
\end{lstlisting}
% end_ohpc_run
```

- Are converted and merged into a (OS-dependent) bash script

```
# ---------------------------------------
# Install Performance Tools (Section 4.4)
# ---------------------------------------
yum -y install ohpc-gnu7-perf-tools
yum -y install lmod-defaults-gnu7-openmpi3-ohpc
```

- Plus a *input.local* file, with some cluster-specific configuration / environment

# Existing OpenHPC Automation

**Shortcomings:**

- The *input.local* file exports shell variables, and don't have enough information

- The *recipe.sh* is **not** idempotent

- Extensibility is impossible without editing the files (downstream work)

# Ansible Playbooks

Ansible is a widely used automation tool which can describe the structure and configuration of IT infrastructure with YAML *"playbooks"*.

OpenHPC with Ansible:

- Ansible playbooks can more easily be **idempotent**

- Ansible can manage nodes/tasks according to the the structure of the cluster

- Configuration is passed as a YAML file (no environment handling)

- Composition, using playbooks and roles, building on third-party content

# Ansible OpenHPC Recipes

- Flexible cluster configuration
  - Fine grained / composable
  - Cluster wide / node group wide / node specific
- Works on both x86_64 and AArch64
  - Ansible gathers information about architecture
  - Same playbook runs on both
- OS is directly inferred by Ansible (gather_facts)
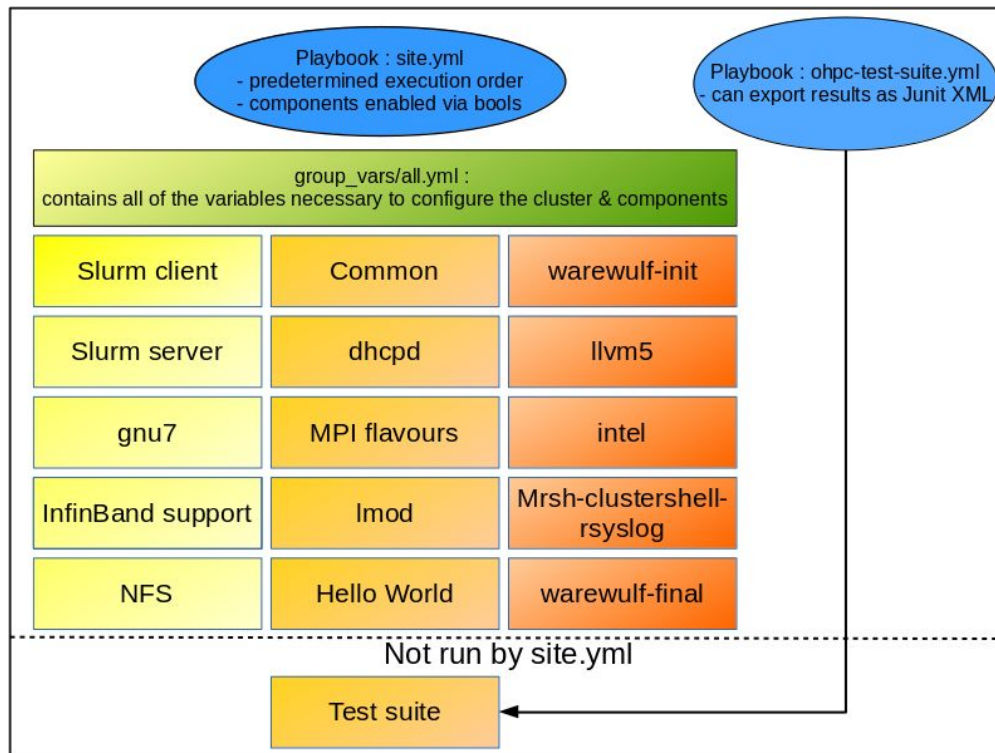  - Yum, apt, zypper... can be switched in the roles' logic

# Ansible OpenHPC Recipes

**The basic structure of the Ansible playbook**

```
playbook
  +---- group_vars/
      +-- all.yml              cluster wide configurations
      +-- group1,group2 … node group(e.g., computing nodes) specific configurations
  +---- host_vars/
      +--- host1,host2         … host specific configurations
  +---- roles/                 package specific tasks, templates, config files, and config variables
      +--- package-name/
              +--- tasks/main.yml … YAML file to describe installation method of package-name
              +--- vars/main.yml … package specific configuration variables
              +--- templates/*.j2 … template files to generate configuration files
```

# Ansible OpenHPC Recipes

# Upstreaming our work

**Option #1:**

- Generate from LaTex sources at the same time as recipe.sh
  - One bundle per recipe (warewulf/xcat, slurm/pbs, centos/suse, arm/x86)
  - Provide as a zip/tar file on the docs package, like the recipes
- Problems:
  - Still need to install package to get recipes that will install OpenHPC
  - Still don't have one recipe to rule them all

**Option #2:**

- Keep as a separate repository, updated in parallel with the doc
  - Easier to integrate to existing automation, update and collaborate
- Problems:
  - Out-of-sync with LaTex sources, can end up meaningless

# Results

# Thank You!

Linaro

Linaro HPC