# A Framework for Benchmarking and Testing HPC Applications for the SDP

**PRACE-CERN-GÉANT-SKAO kick-off workshop on High Performance Computing**

John Taylor, Steve Brasier, John Garbutt
firstname.secondname@stackhpc.com

StackHPC

# Contents

StackHPC

- Context & aims
- Preliminary Study
- Proposed implementation

# Context

StackHPC

- ○ Characterise performance differences for Systems Under Test
- ○ Reduce differences if possible via tuning
- ○ Document results
- ○ Maintain reproducibility and repeatability across SUT
- ○ Optimise solution space

# Preliminary Study

StackHPC

- Aims:
  - Characterise performance differences for Infiniband vs. Ethernet (RDMA)
  - Test matrix - choose a range of industry standard benchmarks
  - Reduce differences if possible via profiling and monitoring
  - Document results
- Based on previous tests by John Taylor: [High Performance Ethernet for HPC – Are we there yet?](#)

# Anticipated Problems

StackHPC

- Test matrix complexity
    - IB vs RoCE
    - MPI libraries
    - Number of nodes + processes etc
    - System parameters
    - Tuning parameters
- Changes to test matrix
    - *"MPI x version y has just come out, can we try that?"*
- Correctness & Repeatability
    - Was the right combination actually run?
    - System changes

# Proposed Solution

**StackHPC**

- Testing-as-code: https://github.com/stackhpc/hpc-tests
- Propose using ReFrame - HPC regression tests
  - Tests defined in python, ReFrame handles interaction with system
  - Easy to define/integrate results extraction and processing
  - Some extension to functionality required: presently only done for slurm
  - In production use at CSCS, NERSC, OSC, responsive developers
- Automate build process
- Monitoring:
  - (Software-defined) monitoring for live view of system with context from tests

# Test Hardware
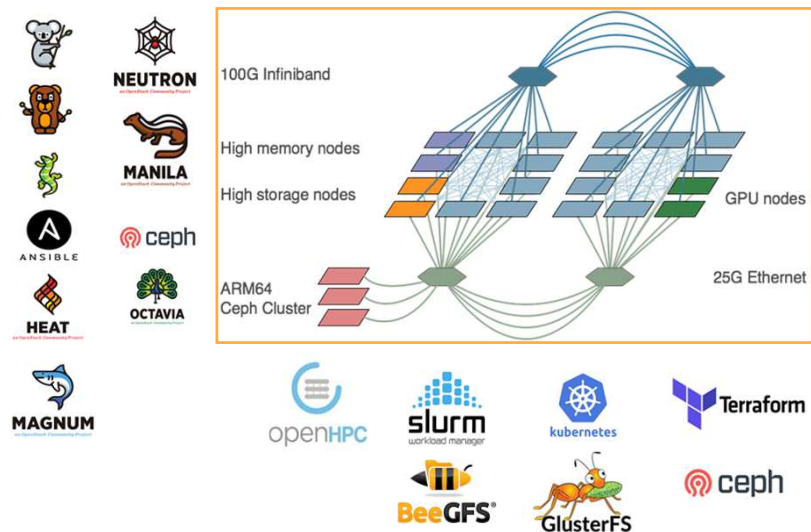
- Two OpenStack Bare Metal Clouds
  - AlaSKA - Baseline System
    - 25GbE RoCE
    - 100 Gbps EDR
    - Broadwell (HT-on)
    - Up to 16 nodes Slurm, K8s
    - OpenHPC Image
  - Cambridge CSD3
    - Cascade lake (HT-off)
    - 50 GbE RoCE
    - 100 Gbps - HDR100
    - Currently up to 56 nodes (Larger by end of the week)
    - Customised Image
    - Large A100 system later in the year!

# ReFrame

Three key aspects:

- [System configuration](#) broken down into:
  - Systems
  - Partitions - logical divisions
  - Environments - software configuration
- [Tests](#) don't need to know about any configuration
- Outputs (from ReFrame's PoV) are:
  - **Test outputs**: e.g. stdout/stderr/files
  - **performance variable logs**
- All under source control

# Results Processing

StackHPC

Plots/tables/reports via jupyter notebooks:

- Web-based interactive python notebook
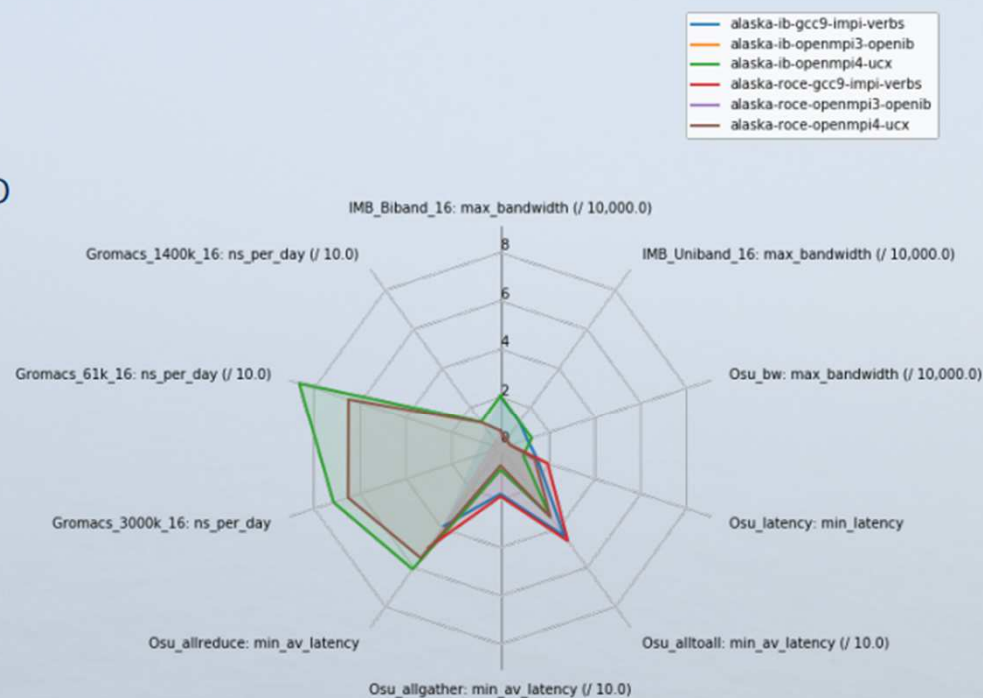- Pages rendered in github

Current demo:
- System info
- Includes automated setup of self-signed key for https:// server
- Separate notebook per-test:
  - IMB - plot of raw results, performance variable history
  - Gromacs - performance variable scaling w/ nodes & history
- Shared code between tests e.g. basic plots of performance variable history

# Optimisation for Key Workloads

- Ensure that infrastructure meets required levels of application performance

- Ensure that performance levels do not regress after reconfiguration

- Find optimal combinations of application and infrastructure configuration

Legend:
- alaska-ib-gcc9-impi-verbs
- alaska-ib-openmpi3-openib
- alaska-ib-openmpi4-ucx
- alaska-roce-gcc9-impi-verbs
- alaska-roce-openmpi3-openib
- alaska-roce-openmpi4-ucx

Radar chart axes:
- IMB_Biband_16: max_bandwidth (/ 10,000.0)
- IMB_Uniband_16: max_bandwidth (/ 10,000.0)
- Osu_bw: max_bandwidth (/ 10,000.0)
- Osu_latency: min_latency
- Osu_alltoall: min_av_latency (/ 10.0)
- Osu_allgather: min_av_latency (/ 10.0)
- Osu_allreduce: min_av_latency
- Gromacs_3000k_16: ns_per_day
- Gromacs_61k_16: ns_per_day (/ 10.0)
- Gromacs_1400k_16: ns_per_day (/ 10.0)

StackHPC

Thanks to
Steve Brasier, John Garbutt,
UIS at the University of Cambridge

# Back Up Slides

# Application Install

**StackHPC**

- Some packages available via *openhpc* repos.
- Also using *spack*: source-based package manager - no root required
- (Multiple) installs define version, build options, compiler, dependencies (mpi)

E.g.:

```
spack install gromacs@2016.4 ^openmpi@4: fabrics=ucx
schedulers=auto
```

- Show `spack info gromacs`
- Integrates with `lmod` and therefore with openhpc & ReFrame (docs for this somewhat lacking)

# Tests & Benchmarks

StackHPC

Synthetics:

- MPI OSU (latency, bandwidth, alltoall, allgather, allreduce):
  - various options here
- IMB: uniband, biband
- HPL
- HPCG

# Tests & Benchmarks

**StackHPC**

Applications

- GROMACS, NAMD (molecular dynamics): HecBioSim 61k/1.5M?/ 3M atoms:
  - One or both codes?
  - Gromacs 1.4M run so far (also used for Archer "small" benchmark)
  - Can only use  2018.x Gromacs
- LS-Dyna (dynamic FEA): Neon (*neon_refined_revised?*), car2car, ODB-10M.
  - Licences? LSTC licence server too.
- Star-CCM+ (CFD): LeMans_100M, TurboCharger, Civil 20M
  - Licences? Flex-LM licence server too
- WRF (CONUS2.5, 12.5 and customer dataset)
  - CONUS require <= WRFV3.8.1. Difficulty getting convergence in previous tests
  - Customer dataset?
- Tensorflow: ResNet50

Potentially also relevant from Archer benchmarks: CASTEP, OpenFOAM

# Test Matrix

StackHPC

- Network
  - IB (100GB?)
  - RoCE (25GB)?
- MPI libraries:
  - OpenMPI4 using UCX
  - Intel MPI (using UCX?): Up to 2019.6 available via yum, .7 in release notes? Early 2019.x known to be problematic.
- Launcher: Only use slurm's `srun` for openmpi (via pmix), impi (via pmi2) at least?
- Number of nodes/cpus/gpus + number of jobs + possibly placement/pinning?
- Number of OpenMP threads (where supported)
- Other MPI tuning parameters
- Any application tuning parameters